

MapReduce 服务

组件操作指南（普通版）

文档版本 01
发布日期 2024-12-13



版权所有 © 华为云计算技术有限公司 2024。保留一切权利。

非经本公司书面许可，任何单位和个人不得擅自摘抄、复制本文档内容的部分或全部，并不得以任何形式传播。

商标声明



HUAWEI和其他华为商标均为华为技术有限公司的商标。

本文档提及的其他所有商标或注册商标，由各自的所有人拥有。

注意

您购买的产品、服务或特性等应受华为云计算技术有限公司商业合同和条款的约束，本文档中描述的全部或部分产品、服务或特性可能不在您的购买或使用范围之内。除非合同另有约定，华为云计算技术有限公司对本文档内容不做任何明示或暗示的声明或保证。

由于产品版本升级或其他原因，本文档内容会不定期进行更新。除非另有约定，本文档仅作为使用指导，本文档中的所有陈述、信息和建议不构成任何明示或暗示的担保。

华为云计算技术有限公司

地址：贵州省贵安新区黔中大道交兴功路华为云数据中心 邮编：550029

网址：<https://www.huaweicloud.com/>

目录

1 使用 Alluxio.....	1
1.1 配置底层存储系统.....	1
1.2 通过数据应用访问 Alluxio.....	2
1.3 Alluxio 常用操作.....	6
2 使用 CarbonData（MRS 3.x 之前版本）.....	9
2.1 从零开始使用 CarbonData.....	9
2.2 CarbonData 表简介.....	11
2.3 创建 CarbonData 表.....	12
2.4 删除 CarbonData 表.....	14
3 使用 CarbonData（MRS 3.x 及之后版本）.....	15
3.1 CarbonData 数据类型概述.....	15
3.2 CarbonData 表用户权限说明.....	18
3.3 使用 Spark 客户端创建 CarbonData 表.....	20
3.4 CarbonData 数据分析.....	22
3.4.1 新建 CarbonData Table.....	23
3.4.2 删除 CarbonData Table.....	24
3.4.3 修改 CarbonData Table.....	25
3.4.4 加载 CarbonData 表数据.....	25
3.4.5 删除 CarbonData 表 Segments.....	26
3.4.6 合并 CarbonData 表 Segments.....	27
3.5 CarbonData 性能调优.....	29
3.5.1 CarbonData 调优思路.....	30
3.5.2 CarbonData 性能调优常见配置参数.....	32
3.5.3 创建 CarbonData Table 的建议.....	34
3.6 CarbonData 常见配置参数.....	36
3.7 CarbonData 语法参考.....	48
3.7.1 DDL.....	48
3.7.1.1 CREATE TABLE.....	48
3.7.1.2 CREATE TABLE As SELECT.....	51
3.7.1.3 DROP TABLE.....	52
3.7.1.4 SHOW TABLES.....	52
3.7.1.5 ALTER TABLE COMPACTION.....	53

3.7.1.6 TABLE RENAME.....	55
3.7.1.7 ADD COLUMNS.....	55
3.7.1.8 DROP COLUMNS.....	56
3.7.1.9 CHANGE DATA TYPE.....	57
3.7.1.10 REFRESH TABLE.....	58
3.7.1.11 REGISTER INDEX TABLE.....	59
3.7.2 DML.....	60
3.7.2.1 LOAD DATA.....	60
3.7.2.2 UPDATE CARBON TABLE.....	65
3.7.2.3 DELETE RECORDS from CARBON TABLE.....	66
3.7.2.4 INSERT INTO CARBON TABLE.....	67
3.7.2.5 DELETE SEGMENT by ID.....	68
3.7.2.6 DELETE SEGMENT by DATE.....	69
3.7.2.7 SHOW SEGMENTS.....	70
3.7.2.8 CREATE SECONDARY INDEX.....	71
3.7.2.9 SHOW SECONDARY INDEXES.....	72
3.7.2.10 DROP SECONDARY INDEX.....	73
3.7.2.11 CLEAN FILES.....	73
3.7.2.12 SET/RESET.....	74
3.7.3 CarbonData 表操作并发语法说明.....	77
3.7.4 CarbonData Segment API 语法说明.....	81
3.7.5 CarbonData 表空间索引语法说明.....	82
3.8 CarbonData 故障处理.....	95
3.8.1 当在 Filter 中使用 Big Double 类型数值时，过滤结果与 Hive 不一致.....	95
3.8.2 executor 内存不足导致查询性能下降.....	96
3.9 CarbonData 常见问题.....	96
3.9.1 为什么对 decimal 数据类型进行带过滤条件的查询时会出现异常输出？.....	97
3.9.2 如何避免对历史数据进行 minor compaction？.....	97
3.9.3 如何在 CarbonData 数据加载时修改默认的组名？.....	98
3.9.4 为什么 INSERT INTO CARBON TABLE 失败？.....	98
3.9.5 为什么含转义字符的输入数据记录到 Bad Records 中的值与原始数据不同？.....	99
3.9.6 当初始 Executor 为 0 时，为什么 INSERT INTO/LOAD DATA 任务分配不正确，打开的 task 少于可用的 Executor？.....	99
3.9.7 为什么并行度大于待处理的 block 数目时，CarbonData 仍需要额外的 executor？.....	100
3.9.8 为什么在 off heap 时数据加载失败？.....	100
3.9.9 为什么创建 Hive 表失败？.....	100
3.9.10 如何在不同的 namespaces 上逻辑地分割数据.....	101
3.9.11 为什么在 Spark Shell 中不能执行更新命令？.....	102
3.9.12 如何在 CarbonData 中配置非安全内存？.....	103
3.9.13 设置了 HDFS 存储目录的磁盘空间配额，CarbonData 为什么会发生异常？.....	103
3.9.14 为什么数据查询/加载失败，且发生“org.apache.carbondata.core.memory.MemoryException: Not enough memory”异常？.....	104
3.9.15 开启防误删后为什么 Carbon 表没有执行 drop 命令，回收站中也会存在该表的文件？.....	104

4 使用 ClickHouse.....	106
4.1 ClickHouse 概述.....	106
4.2 ClickHouse 用户权限管理.....	115
4.2.1 ClickHouse 用户及权限管理.....	115
4.2.2 ClickHouse 使用 OpenLDAP 认证.....	120
4.3 使用 ClickHouse 客户端.....	124
4.4 ClickHouse 表创建.....	127
4.5 ClickHouse 数据导入.....	132
4.5.1 配置 ClickHouse 对接 RDS MySQL 数据库.....	132
4.5.2 配置 ClickHouse 对接 OBS 源文件.....	134
4.5.3 同步 Kafka 数据至 ClickHouse.....	136
4.5.4 导入 DWS 表数据至 ClickHouse.....	139
4.5.5 ClickHouse 数据导入导出.....	143
4.6 ClickHouse 企业级能力增强.....	144
4.6.1 通过 ELB 访问 ClickHouse.....	144
4.6.2 ClickHouse 开启 mysql_port 配置.....	149
4.7 ClickHouse 性能调优.....	149
4.7.1 数据表报错 Too many parts 解决方法.....	149
4.7.2 加速 Merge 操作.....	151
4.7.3 加速 TTL 操作.....	152
4.8 ClickHouse 运维管理.....	152
4.8.1 ClickHouse 日志介绍.....	152
4.8.2 ClickHouse 集群管理.....	154
4.8.2.1 ClickHouse 集群配置说明.....	154
4.8.2.2 ClickHouse 增加磁盘容量.....	157
4.8.3 通过数据文件备份恢复 ClickHouse 数据.....	161
4.8.4 配置 ClickHouse 系统表的生命周期.....	163
4.8.5 集群内 ClickHouseServer 节点间数据迁移.....	163
4.9 ClickHouse 常用 SQL 语法.....	166
4.9.1 CREATE DATABASE 创建数据库.....	166
4.9.2 CREATE TABLE 创建表.....	166
4.9.3 INSERT INTO 插入表数据.....	167
4.9.4 SELECT 查询表数据.....	168
4.9.5 ALTER TABLE 修改表结构.....	169
4.9.6 ALTER TABLE 修改表数据.....	169
4.9.7 DESC 查询表结构.....	170
4.9.8 DROP 删除表.....	170
4.9.9 SHOW 显示数据库和表信息.....	171
4.10 ClickHouse 常见问题.....	171
4.10.1 在 System.disks 表中查询到磁盘 status 是 fault 或者 abnormal.....	171
4.10.2 如何迁移 Hive/HDFS 的数据到 ClickHouse.....	171
4.10.3 使用辅助 Zookeeper 或者副本数据同步表数据时，日志报错.....	172

4.10.4 如何为 ClickHouse 用户赋予数据库级别的 Select 权限.....	172
5 使用 DBService.....	174
5.1 DBService 日志介绍.....	174
6 使用 Flink.....	177
6.1 Flink 作业引擎概述.....	177
6.2 Flink 用户权限管理.....	180
6.2.1 Flink 安全认证机制说明.....	180
6.2.2 Flink 用户权限说明.....	182
6.2.3 创建 FlinkServer 权限角色.....	183
6.2.4 配置 Flink 对接 Kafka 安全认证.....	184
6.2.5 配置 Flink 认证和加密.....	186
6.3 Flink 客户端使用实践.....	193
6.4 创建 FlinkServer 作业前准备.....	202
6.4.1 访问 FlinkServer WebUI 界面.....	203
6.4.2 创建 FlinkServer 应用.....	204
6.4.3 创建 FlinkServer 集群连接.....	204
6.4.4 创建 FlinkServer 数据连接.....	206
6.4.5 创建 FlinkServer 流表源.....	207
6.5 创建 FlinkServer 作业.....	209
6.6 管理 FlinkServer 作业.....	212
6.6.1 配置 FlinkServer 重启策略.....	213
6.6.2 配置 FlinkServer 作业中使用 UDF.....	214
6.7 Flink 运维管理.....	216
6.7.1 Flink 常用配置参数.....	216
6.7.2 Flink 日志介绍.....	235
6.8 Flink 性能调优.....	238
6.8.1 优化 Flink 内存 GC 参数.....	239
6.8.2 配置 Flink 任务并行度.....	239
6.8.3 配置 Flink 任务进程参数.....	240
6.8.4 优化 Flink Netty 网络通信参数.....	241
6.9 Flink 客户端常见命令说明.....	241
6.10 Flink 常见问题.....	245
6.11 签发 Flink 证书样例.....	246
7 使用 Flume.....	251
7.1 Flume 日志采集概述.....	251
7.2 Flume 业务模型配置说明.....	254
7.3 安装 Flume 客户端.....	259
7.3.1 安装 MRS 3.x 之前版本 Flume 客户端.....	259
7.3.2 安装 MRS 3.x 及之后版本 Flume 客户端.....	262
7.4 快速使用 Flume 采集节点日志.....	265
7.5 配置 Flume 非加密传输数据采集任务.....	271

7.5.1 生成 Flume 服务端和客户端的配置文件.....	271
7.5.2 使用 Flume 服务端从本地采集静态日志保存到 Kafka.....	275
7.5.3 使用 Flume 服务端从本地采集静态日志保存到 HDFS.....	277
7.5.4 使用 Flume 服务端从本地采集动态日志保存到 HDFS.....	280
7.5.5 使用 Flume 服务端从 Kafka 采集日志保存到 HDFS.....	283
7.5.6 使用 Flume 客户端从 Kafka 采集日志保存到 HDFS.....	286
7.5.7 使用多级 agent 串联从本地采集静态日志保存到 HBase.....	290
7.6 配置 Flume 加密传输数据采集任务.....	297
7.6.1 配置 Flume 加密传输.....	297
7.6.2 使用多级 agent 串联从本地采集静态日志保存到 HDFS.....	303
7.7 Flume 企业级能力增强.....	312
7.7.1 使用 Flume 客户端加密工具.....	312
7.7.2 配置 Flume 对接安全模式 Kafka.....	313
7.8 Flume 运维管理.....	313
7.8.1 Flume 常用配置参数.....	313
7.8.2 Flume 业务配置指南.....	338
7.8.3 Flume 日志介绍.....	361
7.8.4 查看 Flume 客户端日志.....	363
7.8.5 查看 Flume 客户端监控信息.....	364
7.8.6 停止或卸载 Flume 客户端.....	364
7.9 Flume 常见问题.....	365
7.9.1 如何查看 Flume 日志.....	365
7.9.2 如何在 Flume 配置文件中设置环境变量.....	366
7.9.3 如何开发 Flume 第三方插件.....	367
7.9.4 如何配置 Flume 定制脚本.....	368
8 使用 HBase.....	371
8.1 创建 HBase 权限角色.....	371
8.2 HBase 客户端使用实践.....	374
8.3 快速使用 HBase 进行离线数据分析.....	376
8.4 使用 BulkLoad 工具向 HBase 迁移数据.....	381
8.5 HBase 数据操作.....	381
8.5.1 创建 HBase 索引进行数据查询.....	381
8.5.2 配置 HBase 数据压缩格式和编码.....	382
8.6 HBase 企业级能力增强.....	384
8.6.1 配置 HBase 本地二级索引提升查询效率.....	384
8.6.1.1 HBase 本地二级索引介绍.....	384
8.6.1.2 批量加载 HBase 数据并生成本地二级索引.....	394
8.6.1.3 使用 TableIndexer 工具生成 HBase 本地二级索引.....	397
8.6.1.4 迁移 HBase 索引数据.....	399
8.6.2 增强 HBase BulkLoad 工具数据迁移能力.....	401
8.6.2.1 使用 BulkLoad 工具批量导入 HBase 数据.....	401
8.6.2.2 使用 BulkLoad 工具批量更新 HBase 数据.....	405

8.6.2.3 使用 BulkLoad 工具批量删除 HBase 数据.....	406
8.6.2.4 使用 BulkLoad 工具查询 HBase 表的行统计数.....	406
8.6.2.5 BulkLoad 工具配置文件说明.....	407
8.6.3 配置 RSGroup 管理 RegionServer 资源.....	410
8.7 HBase 性能调优.....	412
8.7.1 提升 HBase BulkLoad 工具批量加载效率.....	412
8.7.2 提升 HBase 连续 Put 数据场景性能.....	413
8.7.3 提升 HBase Put 和 Scan 性能综合调优.....	414
8.7.4 提升 HBase 实时写数据效率.....	417
8.7.5 提升 HBase 实时读数据效率.....	424
8.7.6 HBase JVM 参数优化说明.....	431
8.8 HBase 运维管理.....	431
8.8.1 HBase 日志介绍.....	431
8.8.2 HBase 常用参数配置.....	435
8.8.3 配置 Region Transition 恢复线程.....	436
8.8.4 启用集群间拷贝功能备份集群数据.....	437
8.8.5 配置 HBase 主备集群数据自动备份.....	439
8.8.6 HBase 集群容灾高可用.....	440
8.8.6.1 配置 HBase 主备集群容灾.....	440
8.8.6.2 HBase 容灾集群主备倒换.....	447
8.8.6.3 HBase 容灾集群业务切换指导.....	449
8.9 HBase 常见问题.....	450
8.9.1 结束 BulkLoad 客户端程序导致作业执行失败.....	450
8.9.2 如何修复长时间处于 RIT 状态的 Region.....	451
8.9.3 HMaster 等待 NameSpace 表上线时超时退出.....	451
8.9.4 客户端查询 HBase 出现 SocketTimeoutException 异常.....	452
8.9.5 在启动 HBase shell 时，报错“java.lang.UnsatisfiedLinkError: Permission denied”.....	453
8.9.6 停止运行的 RegionServer，在 HMaster WebUI 中显示的“Dead Region Servers”信息什么时候会被清除掉.....	453
8.9.7 访问 HBase Phoenix 提示权限不足如何处理.....	454
8.9.8 租户使用 HBase BulkLoad 功能提示权限不足如何处理.....	454
8.9.9 如何修复 Overlap 状态的 HBase Region.....	455
8.9.10 Phoenix BulkLoad Tool 使用限制说明.....	456
8.9.11 CTBase 对接 Ranger 权限插件提示权限不足.....	457
8.10 HBase 故障排除.....	457
8.10.1 HBase 客户端连接服务端时，长时间无法连接成功.....	458
8.10.2 在 HBase 连续对同一个表名做删除创建操作时出现创建表异常.....	459
8.10.3 HBase 占用网络端口，连接数过大会导致其他服务不稳定.....	460
8.10.4 有 210000 个 map 和 10000 个 reduce 的 HBase BulkLoad 任务运行失败.....	461
8.10.5 使用 scan 命令仍然可以查询到已修改和已删除的数据.....	461
8.10.6 如何处理由于 Region 处于 FAILED_OPEN 状态而造成的建表失败异常.....	462
8.10.7 如何清理由于建表失败残留在 ZooKeeper 中的 table-lock 节点下的表名.....	462
8.10.8 为什么给 HBase 使用的 HDFS 目录设置 quota 会造成 HBase 故障.....	463

8.10.9 使用 OfflineMetaRepair 工具重新构建元数据后 HMaster 启动失败.....	464
8.10.10 HMaster 日志中频繁打印出 FileNotFoundException 信息.....	464
8.10.11 ImportTsv 工具执行失败报“Permission denied”异常.....	466
8.10.12 使用 HBase BulkLoad 导入数据成功，执行相同的查询时却可能返回不同的结果.....	467
8.10.13 HBase 数据恢复任务报错回滚失败.....	467
8.10.14 HBase RegionServer GC 参数 Xms 和 Xmx 的配置为 31GB，导致 RegionServer 启动失败.....	468
8.10.15 在集群内节点使用 LoadIncrementalHFiles 批量导入数据，报错权限不足.....	468
8.10.16 使用 Phoenix Sqlline 脚本报 import argparse 错误.....	469
9 使用 HDFS.....	470
9.1 HDFS 文件系统目录简介.....	470
9.2 HDFS 用户权限管理.....	477
9.2.1 创建 HDFS 权限角色.....	477
9.2.2 配置 HDFS 用户访问 HDFS 文件权限.....	478
9.3 HDFS 客户端使用实践.....	479
9.4 快速使用 Hadoop.....	481
9.5 配置 HDFS 文件回收站机制.....	485
9.6 配置 HDFS DataNode 数据均衡.....	486
9.7 配置 HDFS DiskBalancer 磁盘均衡.....	490
9.8 配置 HDFS Mover 命令迁移数据.....	493
9.9 配置 HDFS 文件目录标签策略（NodeLabel）.....	494
9.10 配置 NameNode 内存参数.....	499
9.11 设置 HBase 和 HDFS 的句柄数限制.....	500
9.12 配置 HDFS 单目录文件数量.....	501
9.13 HDFS 企业级能力增强.....	502
9.13.1 配置 DataNode 节点容量不一致时的副本放置策略.....	502
9.13.2 配置 DataNode 预留磁盘百分比.....	503
9.13.3 配置 NameNode 黑名单功能.....	504
9.13.4 配置 Hadoop 数据传输加密.....	506
9.14 HDFS 性能调优.....	507
9.14.1 提升 HDFS 写数据性能.....	507
9.14.2 配置 HDFS 客户端元数据缓存提高读取性能.....	508
9.14.3 使用活动缓存提升 HDFS 客户端连接性能.....	509
9.14.4 HDFS 网络不稳定场景调优.....	511
9.14.5 优化 HDFS NameNode RPC 的服务质量.....	511
9.14.6 优化 HDFS DataNode RPC 的服务质量.....	513
9.14.7 执行 HDFS 文件并发操作命令.....	514
9.14.8 使用 LZC 压缩算法存储 HDFS 文件.....	516
9.15 HDFS 运维管理.....	517
9.15.1 HDFS 常用配置参数.....	518
9.15.2 HDFS 日志介绍.....	518
9.15.3 查看 HDFS 容量状态.....	522
9.15.4 更改 DataNode 的存储目录.....	526

9.15.5 调整 DataNode 磁盘坏卷信息.....	529
9.15.6 配置 HDFS token 的最大存活时间.....	530
9.15.7 使用 distcp 命令跨集群复制 HDFS 数据.....	530
9.15.8 配置 NFS 服务器存储 NameNode 元数据.....	534
9.16 HDFS 常见问题.....	535
9.16.1 执行 distcp 命令报错如何处理.....	535
9.16.2 HDFS 执行 Balance 时被异常停止如何处理.....	536
9.16.3 访问 HDFS WebUI 时，界面提示无法显示此页.....	536
9.16.4 HDFS WebUI 无法正常刷新损坏数据的信息.....	537
9.16.5 NameNode 节点长时间满负载导致客户端无响应.....	537
9.16.6 为什么主 NameNode 重启后系统出现双备现象.....	538
9.16.7 为什么 DataNode 无法正常上报数据块.....	539
9.16.8 是否可以手动调整 DataNode 数据存储目录.....	540
9.16.9 DataNode 的容量计算出错如何处理.....	541
9.16.10 为什么存储小文件过程中，缓存中的数据会丢失.....	541
9.16.11 当分级存储策略为 LAZY_PERSIST 时为什么文件的副本的存储类型为 DISK.....	542
9.16.12 为什么 NameNode UI 上显示有一些块缺失.....	542
9.17 HDFS 故障排除.....	543
9.17.1 往 HDFS 写数据时报错“java.net.SocketException”.....	543
9.17.2 删除大量文件后重启 NameNode 耗时长.....	544
9.17.3 EditLog 不连续导致 NameNode 启动失败.....	545
9.17.4 当备 NameNode 存储元数据时，断电后备 NameNode 启动失败.....	546
9.17.5 dfs.datanode.data.dir 中定义的磁盘数量等于 dfs.datanode.failed.volumes.tolerated 的值时， DataNode 启动失败.....	547
9.17.6 HDFS 调用 FileInputFormat 的 getsplit 的时候出现数组越界.....	548
10 使用 Hive.....	549
10.1 Hive 用户权限管理.....	549
10.1.1 Hive 用户权限说明.....	549
10.1.2 创建 Hive 角色.....	552
10.1.3 配置 Hive 表、列或数据库的用户权限.....	556
10.1.4 配置 Hive 业务使用其他组件的用户权限.....	560
10.2 Hive 客户端使用实践.....	564
10.3 快速使用 Hive 进行数据分析.....	567
10.4 Hive 数据存储及加密配置.....	572
10.4.1 使用 HDFS Colocation 存储 Hive 表.....	572
10.4.2 配置 Hive 分区元数据冷热存储.....	573
10.4.3 Hive 支持 ZSTD 压缩格式.....	574
10.4.4 配置 Hive 列加密功能.....	575
10.5 Hive on HBase.....	576
10.5.1 配置跨集群互信下 Hive on HBase.....	576
10.5.2 删除 Hive on HBase 表中的单行记录.....	578
10.6 配置 Hive 读取关系型数据库数据.....	578

10.7 Hive 企业级能力增强.....	580
10.7.1 配置 Hive 目录旧数据自动移除至回收站.....	580
10.7.2 配置 Hive 插入数据到不存在的目录中.....	581
10.7.3 配置创建 Hive 内部表时不能指定 Location.....	581
10.7.4 配置用户在具有读和执行权限的目录中创建外表.....	583
10.7.5 配置基于 HTTPS/HTTP 协议的 REST 接口.....	584
10.7.6 配置 Hive Transform 功能开关.....	587
10.7.7 切换 Hive 执行引擎为 Tez.....	588
10.7.8 Hive 负载均衡.....	590
10.7.8.1 配置 Hive 任务的最大 map 数.....	590
10.7.8.2 配置用户租约隔离访问指定节点的 HiveServer.....	591
10.7.9 配置 Hive 单表动态视图的访问控制权限.....	593
10.7.10 配置创建临时函数的用户不需要具有 ADMIN 权限.....	594
10.7.11 配置具备表 select 权限的用户可查看表结构.....	594
10.7.12 配置仅 Hive 管理员用户能创建库和在 default 库建表.....	595
10.7.13 配置 Hive 支持创建超过 32 个角色.....	596
10.7.14 创建 Hive 用户自定义函数.....	597
10.7.15 配置 Hive Beeline 高可靠性.....	600
10.8 Hive 性能调优.....	601
10.8.1 建立 Hive 表分区提升查询效率.....	602
10.8.2 Hive Join 数据优化.....	603
10.8.3 Hive Group By 语句优化.....	605
10.8.4 Hive ORC 数据存储优化.....	605
10.8.5 Hive SQL 逻辑优化.....	606
10.8.6 使用 Hive CBO 功能优化查询效率.....	607
10.9 Hive 运维管理.....	608
10.9.1 Hive 常用配置参数.....	608
10.9.2 Hive 日志介绍.....	610
10.10 Hive 常见 SQL 语法说明.....	613
10.10.1 Hive SQL 扩展语法说明.....	613
10.10.2 自定义 Hive 表行分隔符.....	615
10.10.3 Hive 支持的传统关系型数据库语法说明.....	616
10.11 Hive 常见问题.....	617
10.11.1 如何删除所有 HiveServer 中的永久函数.....	617
10.11.2 为什么已备份的 Hive 表无法执行 drop 操作.....	619
10.11.3 如何在 Hive 自定义函数中操作本地文件.....	619
10.11.4 如何强制停止 Hive 执行的 MapReduce 任务.....	620
10.11.5 Hive 不支持复杂类型字段名称中包含哪些特殊字符.....	620
10.11.6 如何对 Hive 表大小数据进行监控.....	620
10.11.7 如何防止 insert overwrite 语句误操作导致数据丢失.....	621
10.11.8 未安装 HBase 时 Hive on Spark 任务卡顿如何处理.....	621
10.11.9 Hive 使用 WHERE 条件查询超过 3.2 万分区的表报错.....	622

10.11.10 使用 IBM 的 JDK 访问 Beeline 客户端出现连接 HiveServer 失败.....	622
10.11.11 Hive 表的 Location 支持跨 OBS 和 HDFS 路径吗.....	623
10.11.12 MapReduce 引擎无法查询 Tez 引擎执行 union 语句写入的数据.....	623
10.11.13 Hive 是否支持对同一张表或分区进行并发写数据.....	624
10.11.14 Hive 是否支持向量化查询.....	624
10.11.15 Hive 表的 HDFS 数据目录被误删，但是元数据仍然存在，导致执行任务报错.....	624
10.11.16 如何关闭 Hive 客户端日志.....	625
10.11.17 为什么在 Hive 自定义配置中添加 OBS 快删目录后不生效.....	625
10.11.18 Hive 配置类问题.....	626
10.12 Hive 故障排除.....	627
10.12.1 如何对 insert overwrite 自读自写场景进行优化.....	627
10.12.2 Hive SQL 运行变慢阶段如何排查.....	628
11 使用 Hudi.....	629
11.1 Hudi 表概述.....	629
11.2 使用 Spark Shell 创建 Hudi 表.....	630
11.3 使用 Hudi-Cli.sh 操作 Hudi 表.....	633
11.4 Hudi 写操作.....	635
11.4.1 批量写入 Hudi 表.....	635
11.4.2 流式写入 Hudi 表.....	639
11.4.3 将 Hudi 表数据同步到 Hive.....	640
11.5 Hudi 读操作.....	642
11.5.1 读取 Hudi 数据概述.....	642
11.5.2 读取 Hudi cow 表视图.....	643
11.5.3 读取 Hudi mor 表视图.....	644
11.6 数据管理维护.....	644
11.6.1 Hudi Clustering 操作说明.....	644
11.6.2 Hudi Cleaning 操作说明.....	647
11.6.3 Hudi Compaction 操作说明.....	647
11.6.4 Hudi Savepoint 操作说明.....	648
11.7 Hudi 常见配置参数.....	648
11.7.1 写入操作配置.....	648
11.7.2 同步 Hive 表配置.....	650
11.7.3 index 相关配置.....	651
11.7.4 存储配置.....	653
11.7.5 compaction&cleaning 配置.....	654
11.7.6 单表并发控制配置.....	656
11.8 Hudi 性能调优.....	657
11.9 Hudi 常见问题.....	658
11.9.1 数据写入.....	658
11.9.1.1 写入更新数据时报错 Parquet/Avro schema.....	658
11.9.1.2 写入更新数据时报错 UnsupportedOperationException.....	658
11.9.1.3 写入更新数据时报错 SchemaCompatibilityException.....	658

11.9.1.4 Hudi 在 upsert 时占用了临时文件夹中大量空间.....	659
11.9.1.5 Hudi 写入小精度 Decimal 数据失败.....	659
11.9.2 数据采集.....	659
11.9.2.1 使用 kafka 采集数据时报错 IllegalArgumentException.....	660
11.9.2.2 采集数据时报错 HoodieException.....	660
11.9.2.3 采集数据时报错 HoodieKeyException.....	660
11.9.3 Hive 同步.....	660
11.9.3.1 Hive 同步数据报错 SQLException.....	661
11.9.3.2 Hive 同步数据报错 HoodieHiveSyncException.....	661
11.9.3.3 Hive 同步数据报错 SemanticException.....	661
12 使用 Hue（MRS 3.x 之前版本）.....	662
12.1 访问 Hue WebUI 界面.....	662
12.2 使用 Hue WebUI 操作 Hive 表.....	663
12.3 在 Hue WebUI 使用 HiveQL 编辑器.....	664
12.4 在 Hue WebUI 使用元数据浏览器.....	666
12.5 在 Hue WebUI 使用文件浏览器.....	669
12.6 在 Hue WebUI 使用作业浏览器.....	672
12.7 Hue 常用配置参数.....	673
13 使用 Hue（MRS 3.x 及之后版本）.....	675
13.1 访问 Hue WebUI 界面.....	675
13.2 使用 Hue WebUI 操作 Hive 表.....	676
13.3 创建 Hue 操作任务.....	677
13.3.1 在 Hue WebUI 使用 HiveQL 编辑器.....	677
13.3.2 在 Hue WebUI 使用 SparkSql 编辑器.....	679
13.3.3 在 Hue WebUI 使用元数据浏览器.....	681
13.3.4 在 Hue WebUI 使用文件浏览器.....	682
13.3.5 在 Hue WebUI 使用作业浏览器.....	685
13.3.6 在 Hue WebUI 使用 HBase.....	686
13.4 使用 Hue WebUI 典型场景.....	687
13.4.1 HDFS on Hue.....	687
13.4.2 配置 HDFS 冷热数据迁移.....	691
13.4.3 Hive on Hue.....	698
13.4.4 Oozie on Hue.....	700
13.5 Hue 常用配置参数.....	701
13.6 Hue 日志介绍.....	702
13.7 Hue 常见问题.....	705
13.7.1 使用 Hive 输入 use database 语句失效.....	705
13.7.2 使用 Hue WebUI 访问 HDFS 文件失败.....	705
13.7.3 在 Hue 页面上上传大文件失败.....	706
13.7.4 Hue WebUI 中 Oozie 编辑器的时区设置问题.....	707
13.7.5 访问 Hue 原生页面时间长，文件浏览器报错 Read timed out.....	708

14 使用 Impala	709
14.1 Impala 客户端使用实践.....	709
14.2 访问 Impala WebUI 界面.....	711
14.3 使用 Impala 操作 Kudu 表.....	714
14.4 Impala 对接外部 LDAP.....	715
14.5 Impala 启用并配置动态资源池.....	716
14.6 使用 Impala 查询管理界面.....	719
14.7 Impala 常见配置参数.....	720
14.8 Impala 常见问题.....	721
14.8.1 Impala 服务是否支持磁盘热插拔.....	721
15 使用 Kafka	722
15.1 Kafka 数据消费概述.....	722
15.2 Kafka 用户权限管理.....	723
15.2.1 Kafka 用户权限说明.....	723
15.2.2 创建 Kafka 权限角色.....	726
15.2.3 配置 Kafka 用户 Token 认证信息.....	727
15.3 Kafka 客户端使用实践.....	729
15.4 快速使用 Kafka 生产消费数据.....	732
15.5 创建 Kafka Topic.....	736
15.6 在 Kafka Topic 中接入消息.....	738
15.7 管理 Kafka Topic.....	741
15.7.1 查看 Kafka Topic 信息.....	741
15.7.2 修改 Kafka Topic 配置.....	743
15.7.3 增加 Kafka Topic 分区.....	744
15.7.4 管理 Kafka Topic 中的消息.....	745
15.7.5 查看 Kafka 数据生产消费详情.....	746
15.8 Kafka 企业级能力增强.....	748
15.8.1 配置 Kafka 高可用和高可靠.....	749
15.8.2 配置 Kafka 数据安全传输协议.....	752
15.8.3 配置 Kafka 数据均衡工具.....	755
15.9 Kafka 性能调优.....	757
15.10 Kafka 运维管理.....	758
15.10.1 Kafka 常用配置参数.....	758
15.10.2 Kafka 日志介绍.....	762
15.10.3 更改 Broker 的存储目录.....	765
15.10.4 迁移 Kafka 节点内数据.....	767
15.10.5 均衡 Kafka 扩容节点后数据.....	770
15.11 Kafka 常见问题.....	772
15.11.1 Kafka 业务规格说明.....	772
15.11.2 Kafka 相关特性说明.....	773
15.11.3 基于 binlog 的 MySQL 数据同步到 MRS 集群中.....	775
15.11.4 如何解决 Kafka topic 无法删除的问题.....	780

16 使用 KafkaManager	781
16.1 KafkaManager 介绍.....	781
16.2 访问 KafkaManager 的 WebUI.....	781
16.3 管理 Kafka 集群.....	782
16.4 Kafka 集群监控管理.....	785
17 使用 Loader	793
17.1 从零开始使用 Loader.....	793
17.2 Loader 使用简介.....	794
17.3 Loader 常用参数.....	795
17.4 创建 Loader 角色.....	796
17.5 Loader 连接配置说明.....	798
17.6 管理 Loader 连接（MRS 3.x 之前版本）.....	800
17.7 管理 Loader 连接（MRS 3.x 及之后版本）.....	802
17.8 Loader 作业源连接配置说明.....	806
17.9 Loader 作业目的连接配置说明.....	809
17.10 管理 Loader 作业.....	812
17.11 准备 MySQL 数据库连接的驱动.....	814
17.12 数据导入.....	816
17.12.1 概述.....	816
17.12.2 使用 Loader 导入数据.....	818
17.12.3 典型场景：从 SFTP 服务器导入数据到 HDFS/OBS.....	831
17.12.4 典型场景：从 SFTP 服务器导入数据到 HBase.....	836
17.12.5 典型场景：从 SFTP 服务器导入数据到 Hive.....	842
17.12.6 典型场景：从 FTP 服务器导入数据到 HBase.....	847
17.12.7 典型场景：从关系型数据库导入数据到 HDFS/OBS.....	853
17.12.8 典型场景：从关系型数据库导入数据到 HBase.....	858
17.12.9 典型场景：从关系型数据库导入数据到 Hive.....	863
17.12.10 典型场景：从 HDFS/OBS 导入数据到 HBase.....	868
17.12.11 典型场景：从关系型数据库导入数据到 ClickHouse.....	872
17.12.12 典型场景：从 HDFS 导入数据到 ClickHouse.....	876
17.13 数据导出.....	879
17.13.1 概述.....	879
17.13.2 使用 Loader 导出数据.....	881
17.13.3 典型场景：从 HDFS/OBS 导出数据到 SFTP 服务器.....	889
17.13.4 典型场景：从 HBase 导出数据到 SFTP 服务器.....	894
17.13.5 典型场景：从 Hive 导出数据到 SFTP 服务器.....	899
17.13.6 典型场景：从 HDFS/OBS 导出数据到关系型数据库.....	903
17.13.7 典型场景：从 HBase 导出数据到关系型数据库.....	908
17.13.8 典型场景：从 Hive 导出数据到关系型数据库.....	912
17.13.9 典型场景：从 HBase 导出数据到 HDFS/OBS.....	916
17.14 作业管理.....	919
17.14.1 批量迁移 Loader 作业.....	919

17.14.2 批量删除 Loader 作业.....	920
17.14.3 批量导入 Loader 作业.....	921
17.14.4 批量导出 Loader 作业.....	921
17.14.5 查看作业历史信息.....	922
17.15 算子帮助.....	923
17.15.1 概述.....	923
17.15.2 输入算子.....	925
17.15.2.1 CSV 文件输入.....	925
17.15.2.2 固定宽度文件输入.....	927
17.15.2.3 表输入.....	929
17.15.2.4 HBase 输入.....	931
17.15.2.5 HTML 输入.....	932
17.15.2.6 Hive 输入.....	935
17.15.2.7 Spark 输入.....	936
17.15.3 转换算子.....	938
17.15.3.1 长整型时间转换.....	938
17.15.3.2 空值转换.....	940
17.15.3.3 增加常量字段.....	941
17.15.3.4 随机值转换.....	943
17.15.3.5 拼接转换.....	944
17.15.3.6 分隔转换.....	945
17.15.3.7 取模转换.....	947
17.15.3.8 剪切字符串.....	948
17.15.3.9 EL 操作转换.....	949
17.15.3.10 字符串大小写转换.....	951
17.15.3.11 字符串逆序转换.....	952
17.15.3.12 字符串空格清除转换.....	953
17.15.3.13 过滤行转换.....	954
17.15.3.14 更新域.....	956
17.15.4 输出算子.....	957
17.15.4.1 Hive 输出.....	957
17.15.4.2 Spark 输出.....	959
17.15.4.3 表输出.....	961
17.15.4.4 文件输出.....	963
17.15.4.5 HBase 输出.....	965
17.15.4.6 ClickHouse 输出.....	967
17.15.5 关联、编辑、导入、导出算子的字段配置信息.....	969
17.15.6 配置项中使用宏定义.....	972
17.15.7 算子数据处理规则.....	973
17.16 客户端工具说明.....	976
17.16.1 使用命令行运行 Loader 作业.....	976
17.16.2 loader-tool 工具使用指导.....	980

17.16.3 loader-tool 工具使用示例.....	987
17.16.4 schedule-tool 工具使用指导.....	989
17.16.5 schedule-tool 工具使用示例.....	993
17.16.6 使用 loader-backup 工具备份作业数据.....	995
17.16.7 开源 sqoop-shell 工具使用指导.....	998
17.16.8 开源 sqoop-shell 工具使用示例（SFTP - HDFS）.....	1008
17.16.9 开源 sqoop-shell 工具使用示例（Oracle - HBase）.....	1017
17.17 Loader 日志介绍.....	1026
17.18 样例：通过 Loader 将数据从 OBS 导入 HDFS.....	1028
17.19 Loader 常见问题.....	1029
17.19.1 IE 10&IE 11 浏览器无法保存数据.....	1030
17.19.2 将 Oracle 数据库中的数据导入 HDFS 时各连接器的区别.....	1030
18 使用 Kudu.....	1032
18.1 从零开始使用 Kudu.....	1032
18.2 访问 Kudu 的 WebUI.....	1033
19 使用 MapReduce.....	1036
19.1 配置使用分布式缓存执行 MapReduce 任务.....	1036
19.2 配置 MapReduce shuffle address.....	1038
19.3 配置 MapReduce 集群管理员列表.....	1039
19.4 通过 Windows 系统提交 MapReduce 任务.....	1039
19.5 配置 MapReduce 任务日志归档和清理机制.....	1040
19.6 MapReduce 性能调优.....	1042
19.6.1 多 CPU 内核下的 MapReduce 调优配置.....	1042
19.6.2 配置 MapReduce Job 基线.....	1045
19.6.3 MapReduce Shuffle 调优.....	1047
19.6.4 MapReduce 大任务的 AM 调优.....	1051
19.6.5 配置 MapReduce 任务推测执行.....	1051
19.6.6 通过 Slow Start 调优 MapReduce 任务.....	1052
19.6.7 MapReduce 任务 commit 阶段优化.....	1052
19.6.8 降低 MapReduce 客户端运行任务失败率.....	1053
19.7 MapReduce 日志介绍.....	1054
19.8 MapReduce 常见问题.....	1056
19.8.1 ResourceManager 进行主备切换后，任务中断后运行时间过长.....	1056
19.8.2 MapReduce 任务长时间无进展.....	1057
19.8.3 为什么运行任务时客户端不可用.....	1057
19.8.4 在缓存中找不到 HDFS_DELEGATION_TOKEN 如何处理.....	1058
19.8.5 如何在提交 MapReduce 任务时设置任务优先级.....	1058
19.8.6 MapReduce 任务运行失败，ApplicationMaster 出现物理内存溢出异常.....	1059
19.8.7 MapReduce 作业信息无法通过 ResourceManager Web UI 页面的 Tracking URL 打开.....	1060
19.8.8 多个 NameService 环境下运行 MapReduce 任务失败.....	1060
19.8.9 基于分区的任务黑名单异常如何处理.....	1061

20 使用 OpenTSDB.....	1062
20.1 使用 MRS 客户端操作 OpenTSDB 指标数据.....	1062
20.2 使用 curl 命令操作 OpenTSDB.....	1064
21 使用 Oozie.....	1066
21.1 使用 Oozie 客户端提交作业.....	1066
21.1.1 Oozie 客户端配置说明.....	1066
21.1.2 使用 Oozie 客户端提交 Hive 任务.....	1068
21.1.3 使用 Oozie 客户端提交 Spark2x 任务.....	1070
21.1.4 使用 Oozie 客户端提交 Loader 任务.....	1071
21.1.5 使用 Oozie 客户端提交 DistCp 任务.....	1073
21.1.6 使用 Oozie 客户端提交其它任务.....	1075
21.2 使用 Hue 提交 Oozie 作业.....	1078
21.2.1 使用 Hue 创建工作流.....	1078
21.2.2 使用 Hue 提交 Oozie Hive2 作业.....	1079
21.2.3 使用 Hue 提交 Oozie HQL 脚本.....	1081
21.2.4 使用 Hue 提交 Oozie Spark2x 作业.....	1081
21.2.5 使用 Hue 提交 Oozie Java 作业.....	1083
21.2.6 使用 Hue 提交 Oozie Loader 作业.....	1083
21.2.7 使用 Hue 提交 Oozie Mapreduce 作业.....	1084
21.2.8 使用 Hue 提交 Oozie Sub workflow 作业.....	1085
21.2.9 使用 Hue 提交 Oozie Shell 作业.....	1086
21.2.10 使用 Hue 提交 Oozie HDFS 作业.....	1088
21.2.11 使用 Hue 提交 Oozie Streaming 作业.....	1088
21.2.12 使用 Hue 提交 Oozie Distcp 作业.....	1089
21.2.13 使用 Hue 提交 Oozie SSH 作业.....	1090
21.2.14 使用 Hue 提交 Coordinator 定时调度作业.....	1091
21.2.15 使用 Hue 提交提交 Bundle 批处理作业.....	1092
21.2.16 在 Hue 界面中查询 Oozie 作业结果.....	1093
21.2.17 配置 Oozie 节点间用户互信.....	1094
21.3 开启 Oozie HA 机制.....	1095
21.4 Oozie 日志介绍.....	1096
21.5 Oozie 常见问题.....	1099
21.5.1 Oozie 定时任务没有准时运行如何处理.....	1099
21.5.2 HDFS 上更新了 oozie 的 share lib 目录但没有生效.....	1099
21.5.3 Oozie 作业执行失败常用排查手段.....	1099
22 使用 Presto.....	1101
22.1 访问 Presto 的 WebUI.....	1101
22.2 使用 Presto 客户端执行查询语句.....	1103
22.3 Presto 常见问题.....	1104
22.3.1 Presto 配置多 Hive 连接	1105
23 使用 Ranger (MRS 1.9.2)	1106

23.1 创建 Ranger 集群.....	1106
23.2 访问 Ranger WebUI 及同步 Unix 用户到 Ranger WebUI.....	1107
23.3 在 Ranger 中配置 Hive/Impala 的访问权限.....	1109
23.4 在 Ranger 中配置 HBase 的访问权限.....	1114
24 使用 Ranger (MRS 3.x)	1120
24.1 登录 Ranger WebUI 界面.....	1120
24.2 MRS 集群服务启用 Ranger 鉴权.....	1122
24.3 添加 Ranger 权限策略.....	1122
24.4 Ranger 权限策略配置示例.....	1124
24.4.1 添加 HDFS 的 Ranger 访问权限策略.....	1124
24.4.2 添加 HBase 的 Ranger 访问权限策略.....	1128
24.4.3 添加 Hive 的 Ranger 访问权限策略.....	1131
24.4.4 添加 Impala 的 Ranger 访问权限策略.....	1139
24.4.5 添加 Yarn 的 Ranger 访问权限策略.....	1142
24.4.6 添加 Spark2x 的 Ranger 访问权限策略.....	1145
24.4.7 添加 Kafka 的 Ranger 访问权限策略.....	1153
24.4.8 添加 Storm 的 Ranger 访问权限策略.....	1160
24.5 查看 Ranger 审计信息.....	1163
24.6 配置 Ranger 安全区信息.....	1163
24.7 普通集群修改 Ranger 数据源为 Ldap.....	1166
24.8 查看 Ranger 用户权限同步信息.....	1167
24.9 Ranger 日志介绍.....	1169
24.10 Ranger 常见问题.....	1172
24.10.1 安装集群过程中 Ranger 启动失败.....	1172
24.10.2 如何判断某个服务是否使用了 Ranger 鉴权.....	1172
24.10.3 新创建用户修改完密码后无法登录 Ranger.....	1172
24.10.4 Ranger 界面添加或者修改 HBase 策略时，无法使用通配符搜索已存在的 HBase 表.....	1173
24.10.5 在 Ranger 管理界面查看不到创建的 MRS 用户.....	1173
24.10.6 MRS 用户无法同步至 Ranger 管理界面.....	1174
25 使用 Spark (MRS 3.x 之前版本)	1175
25.1 从零开始使用 Spark.....	1175
25.2 从零开始使用 Spark SQL.....	1179
25.3 使用 Spark 客户端.....	1181
25.4 访问 Spark Web UI 界面.....	1182
25.5 Spark 对接 OpenTSDB.....	1185
25.5.1 创建表关联 OpenTSDB.....	1185
25.5.2 插入数据至 OpenTSDB 表.....	1186
25.5.3 查询 OpenTSDB 表.....	1186
25.5.4 默认配置修改.....	1187
26 使用 Spark2x (MRS 3.x 及之后版本)	1188
26.1 Spark 用户权限管理.....	1188

26.1.1 SparkSQL 权限介绍.....	1188
26.1.2 创建 SparkSQL 角色.....	1192
26.1.3 配置 Spark 表、列和数据库的用户权限.....	1195
26.1.4 配置 SparkSQL 业务用户权限.....	1197
26.1.5 配置 Spark2x Web UI ACL.....	1199
26.1.6 Spark 客户端和服务端权限参数配置说明.....	1201
26.2 Spark 客户端使用实践.....	1202
26.3 配置 Spark 读取 HBase 表数据.....	1205
26.4 配置 Spark 任务不获取 HBase Token 信息.....	1208
26.5 Spark Core 企业级能力增强.....	1209
26.5.1 配置 Spark HA 增强高可用.....	1209
26.5.1.1 配置多主实例模式.....	1209
26.5.1.2 配置 Spark 多租户模式.....	1210
26.5.1.3 配置多主实例与多租户模式切换.....	1211
26.5.2 配置 Spark 事件队列大小.....	1212
26.5.3 配置 parquet 表的压缩格式.....	1213
26.5.4 使用 Ranger 时适配第三方 JDK.....	1214
26.5.5 使用 Spark 小文件合并工具说明.....	1215
26.5.6 配置流式读取 Spark Driver 执行结果.....	1217
26.6 Spark SQL 企业级能力增强.....	1218
26.6.1 配置矢量化读取 ORC 数据.....	1218
26.6.2 配置过滤掉分区表中路径不存在的分区.....	1219
26.6.3 配置 Hive 表分区动态覆盖.....	1220
26.6.4 配置 Spark SQL 开启 Adaptive Execution 特性.....	1220
26.6.5 配置 SparkSQL 的分块个数.....	1223
26.7 Spark Streaming 企业级能力增强.....	1223
26.7.1 配置 Kafka 后进先出.....	1224
26.7.2 配置对接 Kafka 可靠性.....	1225
26.8 Spark Core 性能调优.....	1226
26.8.1 Spark Core 数据序列化.....	1226
26.8.2 Spark Core 内存调优.....	1227
26.8.3 Spark Core 内存调优.....	1227
26.8.4 配置 Spark Core 广播变量.....	1228
26.8.5 配置 Spark Executor 堆内存参数.....	1228
26.8.6 使用 External Shuffle Service 提升 Spark Core 性能.....	1229
26.8.7 配置 Yarn 模式下 Spark 动态资源调度.....	1230
26.8.8 调整 Spark Core 进程参数.....	1231
26.8.9 Spark DAG 设计规范说明.....	1232
26.8.10 经验总结.....	1234
26.9 Spark SQL 性能调优.....	1235
26.9.1 Spark SQL join 优化.....	1235
26.9.2 优化数据倾斜场景下的 Spark SQL 性能.....	1237

26.9.3 优化小文件场景下的 Spark SQL 性能.....	1239
26.9.4 Spark INSERT SELECT 语句调优.....	1239
26.9.5 动态分区插入场景内存优化.....	1240
26.9.6 小文件优化.....	1240
26.9.7 聚合算法优化.....	1241
26.9.8 Datasource 表优化.....	1242
26.9.9 合并 CBO 优化.....	1243
26.9.10 多级嵌套子查询以及混合 Join 的 SQL 调优.....	1244
26.10 Spark Streaming 性能调优.....	1246
26.11 Spark 运维管理.....	1248
26.11.1 快速配置参数.....	1248
26.11.2 常用参数.....	1256
26.11.3 Spark2x 日志介绍.....	1273
26.11.4 调整 Spark 日志级别.....	1276
26.11.5 配置 WebUI 上查看 Container 日志.....	1278
26.11.6 获取运行中 Spark 应用的 Container 日志.....	1279
26.11.7 配置 Spark Eventlog 日志回滚.....	1280
26.11.8 配置 WebUI 上显示的 Lost Executor 信息的个数.....	1281
26.11.9 配置 JobHistory 本地磁盘缓存.....	1281
26.11.10 增强有限内存下的稳定性.....	1282
26.11.11 配置 YARN-Client 和 YARN-Cluster 不同模式下的环境变量.....	1283
26.11.12 Hive 分区修剪的谓词下推增强.....	1285
26.11.13 配置列统计值直方图 Histogram 用以增强 CBO 准确度.....	1285
26.11.14 CarbonData 首查优化工具.....	1287
26.12 Spark2x 常见问题.....	1288
26.12.1 Spark Core.....	1288
26.12.1.1 日志聚合下如何查看 Spark 已完成应用日志.....	1288
26.12.1.2 Driver 返回码和 RM WebUI 上应用状态显示不一致.....	1289
26.12.1.3 为什么 Driver 进程不能退出.....	1289
26.12.1.4 网络连接超时导致 FetchFailedException.....	1290
26.12.1.5 当事件队列溢出时如何配置事件队列的大小.....	1291
26.12.1.6 Spark 应用执行过程中，日志中一直打印 getApplicationReport 异常且应用较长时间不退出.....	1292
26.12.1.7 Spark 执行应用时上报“Connection to ip:port has been quiet for xxx ms while there are outstanding requests”并导致应用结束.....	1292
26.12.1.8 NodeManager 关闭导致 Executor(s)未移除.....	1294
26.12.1.9 Password cannot be null if SASL is enabled 异常.....	1294
26.12.1.10 向动态分区表中插入数据时，在重试的 task 中出现"Failed to CREATE_FILE"异常.....	1295
26.12.1.11 使用 Hash shuffle 出现任务失败.....	1295
26.12.1.12 访问 Spark 应用的聚合日志页面报“DNS 查找失败”错误.....	1296
26.12.1.13 由于 Timeout waiting for task 异常导致 Shuffle FetchFailed.....	1297
26.12.1.14 Executor 进程 Crash 导致 Stage 重试.....	1297
26.12.1.15 执行大数据量的 shuffle 过程时 Executor 注册 shuffle service 失败.....	1298
26.12.1.16 在 Spark 应用执行过程中 NodeManager 出现 OOM 异常.....	1299

26.12.1.17 安全集群使用 HiBench 工具运行 sparkbench 获取不到 realm.....	1300
26.12.2 SQL 和 DataFrame.....	1301
26.12.2.1 Spark SQL ROLLUP 和 CUBE 使用的注意事项.....	1301
26.12.2.2 Spark SQL 在不同 DB 都可以显示临时表.....	1302
26.12.2.3 如何在 Spark 命令中指定参数值.....	1303
26.12.2.4 SparkSQL 建表时的目录权限.....	1303
26.12.2.5 为什么不同服务之间互相删除 UDF 失败.....	1304
26.12.2.6 Spark SQL 无法查询到 Parquet 类型的 Hive 表的新插入数据.....	1304
26.12.2.7 cache table 使用指导.....	1305
26.12.2.8 Repartition 时有部分 Partition 没数据.....	1305
26.12.2.9 16T 的文本数据转成 4T Parquet 数据失败.....	1306
26.12.2.10 当表名为 table 时，执行相关操作时出现异常.....	1307
26.12.2.11 执行 analyze table 语句，因资源不足出现任务卡住.....	1307
26.12.2.12 为什么有时访问没有权限的 parquet 表时，在上报“Missing Privileges”错误提示之前，会运行一个 Job?	1308
26.12.2.13 spark-sql 退出时打印 RejectedExecutionException 异常栈.....	1308
26.12.2.14 健康检查时，误将 JDBCServer Kill.....	1309
26.12.2.15 日期类型的字段作为过滤条件时匹配'2016-6-30'时没有查询结果.....	1309
26.12.2.16 为什么在启动 spark-beeline 的命令中指定“--hivevar”选项无效.....	1310
26.12.2.17 执行复杂 SQL 语句时报“Code of method ... grows beyond 64 KB”的错误.....	1310
26.12.2.18 在 Beeline/JDBCServer 模式下连续运行 10T 的 TPCDS 测试套会出现内存不足的现象.....	1311
26.12.2.19 连上不同的 JDBCServer，function 不能正常使用.....	1311
26.12.2.20 用 add jar 方式创建 function，执行 drop function 时出现问题.....	1313
26.12.2.21 Spark2x 无法访问 Spark1.5 创建的 DataSource 表.....	1314
26.12.2.22 Spark SQL 无法查询到 ORC 类型的 Hive 表的新插入数据.....	1315
26.12.3 Spark Streaming.....	1315
26.12.3.1 Streaming 任务打印两次相同 DAG 日志.....	1315
26.12.3.2 Spark Streaming 任务一直阻塞.....	1317
26.12.3.3 运行 Spark Streaming 任务参数调优的注意事项.....	1317
26.12.3.4 为什么提交 Spark Streaming 应用超过 token 有效期，应用失败.....	1318
26.12.3.5 为什么 Spark Streaming 应用创建输入流，但该输入流无输出逻辑时，应用从 checkpoint 恢复启动失败.....	1319
26.12.3.6 Spark Streaming 应用运行过程中重启 Kafka，Web UI 界面部分 batch time 对应 Input Size 为 0 records.....	1320
26.12.4 访问 Spark 应用获取的 restful 接口信息有误.....	1321
26.12.5 为什么从 Yarn Web UI 页面无法跳转到 Spark Web UI 界面.....	1322
26.12.6 HistoryServer 缓存的应用被回收，导致此类应用页面访问时出错.....	1323
26.12.7 加载空的 part 文件时，app 无法显示在 JobHistory 的页面上.....	1324
26.12.8 Spark2x 导出带有相同字段名的表，结果导出失败.....	1324
26.12.9 为什么多次运行 Spark 应用程序会引发致命 JRE 错误.....	1325
26.12.10 IE 浏览器访问 Spark2x 原生 UI 界面失败，无法显示此页或者页面显示错误.....	1325
26.12.11 Spark2x 如何访问外部集群组件.....	1326
26.12.12 对同一目录创建多个外表，可能导致外表查询失败.....	1327

26.12.13 访问 Spark2x JobHistory 中某个应用的原生页面时页面显示错误.....	1328
26.12.14 对接 OBS 场景中，spark-beeline 登录后指定 loaction 到 OBS 建表失败.....	1328
26.12.15 Spark shuffle 异常处理.....	1329
27 使用 Sqoop.....	1331
27.1 Sqoop 客户端使用实践.....	1331
27.2 Sqoop1.4.7 适配 MRS 3.x 集群.....	1335
27.3 Sqoop 常用命令及参数介绍.....	1337
27.4 Sqoop 常见问题.....	1340
27.4.1 报错找不到 QueryProvider 类.....	1340
27.4.2 使用 hcatalog 方式同步数据，报错 getHiveClient 方法不存在.....	1340
27.4.3 连接 postgresql 或者 gaussdb 时报错.....	1341
27.4.4 使用 hive-table 方式同步数据到 obs 上的 hive 表报错.....	1343
27.4.5 使用 hive-table 方式同步数据到 orc 表或者 parquet 表失败.....	1343
27.4.6 使用 hive-table 方式同步数据报错.....	1343
27.4.7 使用 hcatalog 方式同步 hive parquet 表报错.....	1344
27.4.8 使用 Hcatalog 方式同步 Hive 和 MySQL 之间的数据，timestamp 和 data 类型字段会报错.....	1345
28 使用 Storm.....	1346
28.1 从零开始使用 Storm.....	1346
28.2 使用 Storm 客户端.....	1347
28.3 使用客户端提交 Storm 拓扑.....	1348
28.4 访问 Storm 的 WebUI.....	1349
28.5 管理 Storm 拓扑.....	1350
28.6 查看 Storm 拓扑日志.....	1351
28.7 Storm 常用参数.....	1352
28.8 配置 Storm 业务用户密码策略.....	1353
28.9 迁移 Storm 业务至 Flink.....	1355
28.9.1 概述.....	1355
28.9.2 完整迁移 Storm 业务.....	1355
28.9.3 嵌入式迁移 Storm 业务.....	1357
28.9.4 迁移 Storm 对接的外部安全组件业务.....	1357
28.10 Storm 日志介绍.....	1358
28.11 性能调优.....	1362
28.11.1 Storm 性能调优.....	1362
29 使用 Tez.....	1365
29.1 访问 Tez WebUI 查看任务执行结果.....	1365
29.2 Tez 常用配置参数.....	1365
29.3 Tez 日志介绍.....	1366
29.4 Tez 常见问题.....	1368
29.4.1 TezUI 无法展示 Tez 任务执行细节.....	1368
29.4.2 进入 Tez WebUI 界面显示异常.....	1368
29.4.3 TezUI 界面无法查看 Yarn 日志.....	1369

29.4.4 TezUI HiveQueries 界面表格数据为空.....	1369
30 使用 Yarn.....	1371
30.1 Yarn 用户权限管理.....	1371
30.1.1 创建 Yarn 角色.....	1371
30.2 使用 Yarn 客户端提交任务.....	1373
30.3 配置 Container 日志聚合功能.....	1374
30.4 启用 Yarn CGroups 功能限制 Container CPU 使用率.....	1379
30.5 Yarn 企业级能力增强.....	1381
30.5.1 配置 Yarn 权限控制开关.....	1381
30.5.2 手动指定运行 Yarn 任务的用户.....	1382
30.5.3 配置 AM 失败重试次数.....	1383
30.5.4 配置 AM 自动调整分配内存.....	1383
30.5.5 配置 AM 作业自动保留.....	1385
30.5.6 配置 Yarn 数据访问通道协议.....	1386
30.5.7 配置自定义调度器的 WebUI.....	1387
30.5.8 配置 NodeManager 角色实例使用的资源.....	1387
30.5.9 配置 ResourceManager 重启后自动加载 Container 信息.....	1388
30.6 Yarn 性能调优.....	1390
30.6.1 调整 Yarn 任务抢占机制.....	1390
30.6.2 手动配置 Yarn 任务优先级.....	1392
30.6.3 Yarn 节点配置调优.....	1393
30.7 Yarn 运维管理.....	1398
30.7.1 Yarn 常用配置参数.....	1398
30.7.2 Yarn 日志介绍.....	1401
30.7.3 配置 Yarn 本地化日志级别.....	1403
30.7.4 检测 Yarn 内存使用情况.....	1404
30.7.5 更改 NodeManager 的存储目录.....	1405
30.8 Yarn 常见问题.....	1408
30.8.1 任务完成后 Container 挂载的文件目录未清除.....	1408
30.8.2 作业执行失败时会发生 HDFS_DELEGATION_TOKEN 到期的异常.....	1409
30.8.3 重启 YARN，本地日志不被删除.....	1409
30.8.4 执行任务时 AppAttempts 重试次数超过 2 次还没有运行失败.....	1409
30.8.5 在 ResourceManager 重启后，应用程序会移回原来的队列.....	1410
30.8.6 YARN 资源池的所有节点都被加入黑名单，任务一直处于运行状态.....	1410
30.8.7 ResourceManager 持续主备倒换.....	1411
30.8.8 当一个 NodeManager 处于 unhealthy 的状态 10 分钟时，新应用程序失败.....	1411
30.8.9 Superior 通过 REST 接口查看已结束或不存在的 applicationID，页面提示 Error Occurred.....	1412
30.8.10 Superior 调度模式下，单个 NodeManager 故障可能导致 MapReduce 任务失败.....	1412
30.8.11 当应用程序从 lost_and_found 队列移动到其他队列时，应用程序不能继续执行.....	1413
30.8.12 如何限制存储在 ZKstore 中的应用程序诊断消息的大小.....	1413
30.8.13 为什么将非 ViewFS 文件系统配置为 ViewFS 时 MapReduce 作业运行失败.....	1414
30.8.14 开启 Native Task 特性后，Reduce 任务在部分操作系统运行失败.....	1415

31 使用 ZooKeeper	1416
31.1 使用 ZooKeeper 客户端.....	1416
31.2 配置 ZooKeeper ZNode ACL.....	1420
31.3 ZooKeeper 常用配置参数.....	1424
31.4 ZooKeeper 日志介绍.....	1425
31.5 ZooKeeper 常见问题.....	1427
31.5.1 创建大量 ZNode 后 ZooKeeper Server 启动失败.....	1427
31.5.2 为什么 ZooKeeper Server 出现 java.io.IOException: Len 的错误日志.....	1428
31.5.3 为什么 ZooKeeper 节点上 netcat 命令无法正常运行.....	1430
31.5.4 如何查看哪个 ZooKeeper 实例是 Leader.....	1430
31.5.5 使用 IBM JDK 时客户端无法连接 ZooKeeper.....	1431
31.5.6 ZooKeeper 客户端刷新 TGT 失败.....	1431
31.5.7 使用 deleteall 命令删除大量 znode 时偶现报错 “Node does not exist”	1431
32 常见操作	1432
32.1 修改集群服务配置参数.....	1432
32.2 访问集群 Manager.....	1436
32.2.1 访问 MRS Manager（MRS 3.x 之前版本）	1436
32.2.2 访问 FusionInsight Manager（MRS 3.x 及之后版本）	1442
32.3 使用 MRS 客户端.....	1446
32.3.1 安装客户端（3.x 及之后版本）	1446
32.3.2 安装客户端（3.x 之前版本）	1453
32.3.3 更新客户端（3.x 及之后版本）	1458
32.3.4 更新客户端（3.x 之前版本）	1460

1 使用 Alluxio

1.1 配置底层存储系统

用户想要通过统一的客户端API和全局命名空间访问包括HDFS和OBS在内的持久化存储系统，从而实现了计算和存储的分离时，可以在MRS Manager页面中配置Alluxio的底层存储系统来实现。集群创建后，默认的底层存储地址是hdfs://hacluster/，即将HDFS的根目录映射到Alluxio。

前提条件

- 已安装Alluxio服务的集群。
- 获取用户“admin”账号密码。“admin”密码在创建MRS集群时由用户指定。

配置 HDFS 作为 Alluxio 的底层文件系统

📖 说明

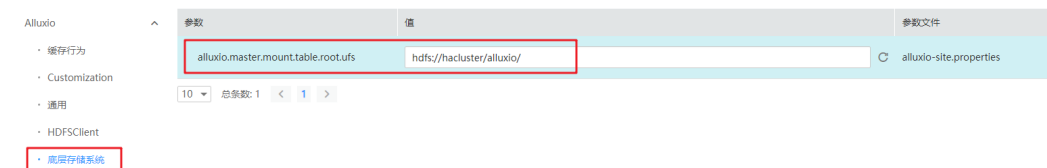
开启Kerberos认证的安全集群不支持该功能。

步骤1 请参考[修改集群服务配置参数](#)，进入Alluxio的“全部配置”页面。

步骤2 在左侧边栏中选择“**Alluxio > 底层存储系统**”，修改参数“alluxio.master.mount.table.root.ufs”的值为“hdfs://hacluster/XXX/”。

例如：若想将“HDFS根目录/alluxio/”作为alluxio的根目录，则修改参数“alluxio.master.mount.table.root.ufs”的值为“hdfs://hacluster/alluxio/”。

图 1-1 HDFS 作为 Alluxio 的底层文件系统



步骤3 单击“保存配置”，并在弹出窗口中勾选“重新启动受影响的服务和实例。”

步骤4 单击“确定”重启Alluxio服务。

----结束

配置 Huawei OBS 作为 Alluxio 的底层文件系统

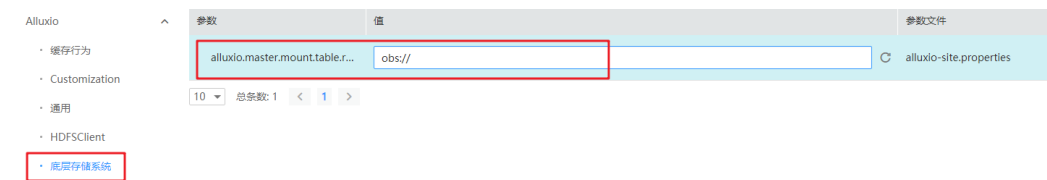
方法一：

步骤1 给集群配置有OBS OperateAccess权限的委托，具体请参见[配置存算分离集群（委托方式）](#)。

步骤2 请参考[修改集群服务配置参数](#)，进入Alluxio的“全部配置”页面。

步骤3 在左侧边栏中选择“**Alluxio > 底层存储系统**”，修改参数“alluxio.master.mount.table.root.ufs”的值为“obs://<OBS_BUCKET>/<OBS_DIRECTORY>/”。OBS_BUCKET为一个已有的OBS文件系统名，OBS_DIRECTORY为该文件系统下的目录。

图 1-2 OBS 作为 Alluxio 的底层文件系统



步骤4 单击“保存配置”，并在弹出窗口中勾选“重新启动受影响的服务和实例。”

步骤5 单击“确定”重启Alluxio服务。

----结束

方法二：

步骤1 给集群配置有OBS OperateAccess权限的委托，具体请参见[配置存算分离集群（委托方式）](#)。

步骤2 登录主Master节点，主节点请参考[如何确认MRS Manager的主备管理节点](#)。

步骤3 执行如下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

说明

/opt/client为举例当前集群客户端的安装目录，请根据实际情况修改。

步骤4 执行如下命令将OBS容器内部的目录挂载到Alluxio的/obs目录。

```
alluxio fs mount /obs obs://<OBS_BUCKET>/<OBS_DIRECTORY>/
```

----结束

1.2 通过数据应用访问 Alluxio

访问Alluxio文件系统的端口号是19998，即地址为alluxio://<alluxio的master节点ip>:19998/<PATH>，本节将通过示例介绍如何通过数据应用（Spark、Hive、Hadoop MapReduce和Presto）访问Alluxio。

使用 Alluxio 作为 Spark 应用程序的输入和输出

步骤1 以root用户登录集群的Master节点，密码为用户创建集群时设置的root密码。

步骤2 执行如下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

步骤3 如果当前集群已启用Kerberos认证，执行如下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit MRS集群用户
```

例如，`kinit admin`

步骤4 准备输入文件，将本地数据复制到Alluxio文件系统中。

如在本地/home目录下准备一个输入文件test_input.txt，然后执行如下命令，将test_input.txt文件放入Alluxio中。

```
alluxio fs copyFromLocal /home/test_input.txt /input
```

步骤5 执行如下命令启动spark-shell。

```
spark-shell
```

步骤6 在spark-shell中运行如下命令。

```
val s = sc.textFile("alluxio://<Alluxio的节点名称>:19998/input")
```

```
val double = s.map(line => line + line)
```

```
double.saveAsTextFile("alluxio://<Alluxio的节点名称>:19998/output")
```

说明

<Alluxio的节点名称>:19998，请根据实际情况替换为AlluxioMaster实例所在所有节点的节点名称与端口号，各个名称与端口之间以英文逗号间隔，例如：`node-ana-corempb.mrs-m0va.com:19998,node-master2kiww.mrs-m0va.com:19998,node-master1cqww.mrs-m0va.com:19998`

步骤7 使用“Ctrl + C”退出spark-shell。

步骤8 通过alluxio命令行**alluxio fs ls /**查看alluxio根目录下存在一个输出目录/output，其中包含了输入文件input的双倍内容。

----结束

在 Alluxio 上创建 Hive 表

步骤1 以root用户登录集群的Master节点，密码为用户创建集群时设置的root密码。

步骤2 执行如下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

步骤3 如果当前集群已启用Kerberos认证，执行如下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit MRS集群用户
```

例如，`kinit admin`

步骤4 准备输入文件，如在本地/home目录下准备一个输入文件hive_load.txt， 内容为

```
1, Alice, company A
2, Bob, company B
```

步骤5 执行如以下命令，将hive_load.txt文件放入Alluxio中。

```
alluxio fs copyFromLocal /home/hive_load.txt /hive_input
```

步骤6 执行如下命令启动hive beeline。

```
beeline
```

步骤7 在beeline中运行如下命令，根据Alluxio中的输入文件进行创表。

```
CREATE TABLE u_user(id INT, name STRING, company STRING) ROW FORMAT
DELIMITED FIELDS TERMINATED BY ',' STORED AS TEXTFILE;
```

```
LOAD DATA INPATH 'alluxio://<Alluxio的节点名称>:19998/hive_input' INTO
TABLE u_user;
```

说明

<Alluxio的节点名称>:19998，请根据实际情况替换为AlluxioMaster实例所在所有节点的节点名称与端口号，各个名称与端口之间以英文逗号间隔，例如：node-ana-coremspb.mrs-m0va.com:19998,node-master2kiww.mrs-m0va.com:19998,node-master1cqww.mrs-m0va.com:19998

步骤8 执行如下命令查看创建的表。

```
select * from u_user;
```

----结束

在 Alluxio 上运行 Hadoop Wordcount

步骤1 以root用户登录集群的Master节点，密码为用户创建集群时设置的root密码。

步骤2 执行如下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

步骤3 如果当前集群已启用Kerberos认证，执行如下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit MRS集群用户
```

例如，`kinit admin`

步骤4 准备输入文件，将本地数据复制到Alluxio文件系统中。

如在本地/home目录下准备一个输入文件test_input.txt，然后执行如下命令，将test_input.txt文件放入Alluxio中。

```
alluxio fs copyFromLocal /home/test_input.txt /input
```

步骤5 通过yarn jar执行wordcount作业。

```
yarn jar /opt/share/hadoop-mapreduce-examples-<hadoop版本号>-mrs-<mrs
集群版本号>/hadoop-mapreduce-examples-<hadoop版本号>-mrs-<mrs集群版本
号>.jar wordcount alluxio://<Alluxio的节点名称>:19998/input alluxio://<Alluxio
的节点名称>:19998/output
```

📖 说明

- <hadoop版本号>请根据实际情况替换。
- <mrs集群版本号>替换为MRS的大版本号，如MRS 1.9.2版本集群此处为**mrs-1.9.0**。
- <Alluxio的节点名称>:19998，请根据实际情况替换为AlluxioMaster实例所在所有节点的节点名称与端口号，各个名称与端口之间以英文逗号间隔，例如：node-ana-coremspb.mrs-m0va.com:19998,node-master2kiww.mrs-m0va.com:19998,node-master1cqww.mrs-m0va.com:19998

步骤6 通过alluxio命令行**alluxio fs ls /**查看alluxio根目录下存在一个输出目录/output，包含了wordcount的结果。

---结束

使用 Presto 在 Alluxio 上查询表

步骤1 以root用户登录集群的Master节点，密码为用户创建集群时设置的root密码。

步骤2 执行如下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

步骤3 如果当前集群已启用Kerberos认证，执行如下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit MRS集群用户
```

例如，`kinit admin`

步骤4 启动hive beeline在alluxio上创建表。

```
beeline
```

```
CREATE TABLE u_user (id int, name string, company string) ROW FORMAT  
DELIMITED FIELDS TERMINATED BY ',' LOCATION 'alluxio://<Alluxio的节点名称>:  
:19998/u_user';
```

```
insert into u_user values(1,'Alice','Company A'),(2, 'Bob', 'Company B');
```

📖 说明

<Alluxio的节点名称>:19998，请根据实际情况替换为AlluxioMaster实例所在所有节点的节点名称与端口号，各个名称与端口之间以英文逗号间隔，例如：node-ana-coremspb.mrs-m0va.com:19998,node-master2kiww.mrs-m0va.com:19998,node-master1cqww.mrs-m0va.com:19998

步骤5 启动Presto客户端，具体请参见[使用Presto客户端执行查询语句](#)的**步骤2~步骤8**。

步骤6 在Presto客户端中执行查询语句**select * from hive.default.u_user;** 查询alluxio上创建表。

图 1-3 Presto 查询 alluxio 上创建的表

```
presto> select * from hive.default.u_user;
id | name | company
---+---+-----
 1 | Alice | Company A
 2 | Bob   | Company B
(2 rows)
```

----结束

1.3 Alluxio 常用操作

前期准备

1. 创建安装Alluxio组件的集群。
2. 以root用户登录集群的主Master节点，密码为用户创建集群时设置的root密码。
3. 执行如下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

使用 Alluxio Shell

Alluxio shell包含多种与Alluxio交互的命令行操作。

- 要查看文件系统操作命令列表。

```
alluxio fs
```

- 使用ls命令列出Alluxio里的文件。例如列出根目录下所有文件。

```
alluxio fs ls /
```

- 使用copyFromLocal命令可以复制本地文件到Alluxio中。

```
alluxio fs copyFromLocal /home/test_input.txt /test_input.txt
```

命令执行后回显：

```
Copied file:///home/test_input.txt to /test_input.txt
```

- 再次使用ls命令列出Alluxio中的文件，可以看到刚刚拷贝的test_input.txt文件。

```
alluxio fs ls /
```

命令执行后回显：

```
12  PERSISTED 11-28-2019 17:10:17:449 100% /test_input.txt
```

输出显示test_input.txt文件在Alluxio中，各参数含义为文件的大小、是否被持久化、创建日期、Alluxio中这个文件的缓存占比、文件名。

- 使用cat命令打印文件的内容。

```
alluxio fs cat /test_input.txt
```

命令执行后回显：

```
Test Alluxio
```

Alluxio 中的挂载功能

Alluxio通过统一命名空间的特性统一了对存储系统的访问。详情请参考：<https://docs.alluxio.io/os/user/2.0/cn/advanced/namespace-management.html>

这个特性允许用户挂载不同的存储系统到Alluxio命名空间中并且通过Alluxio命名空间无缝地跨存储系统访问文件。

1. 在Alluxio中创建一个目录作为挂载点。

```
alluxio fs mkdir /mnt
```

```
Successfully created directory /mnt
```

2. 挂载一个已有的OBS文件系统到Alluxio（前提：给集群配置有OBS OperateAccess权限的委托，具体请参见[配置存算分离集群（委托方式）](#)）。此处以obs-mrstest文件系统为例，请根据实际情况替换文件系统名。

```
alluxio fs mount /mnt/obs obs://obs-mrstest/data
```

```
Mounted obs://obs-mrstest/data at /mnt/obs
```

3. 通过Alluxio命名空间列出OBS文件系统下的文件。使用ls命令列出OBS挂载目录下的文件。

```
alluxio fs ls /mnt/obs
```

```
38  PERSISTED 11-28-2019 17:42:54:554 0% /mnt/obs/hive_load.txt  
12  PERSISTED 11-28-2019 17:43:07:743 0% /mnt/obs/test_input.txt
```

新挂载的文件和目录也可以通过Alluxio WebUI查看。

4. 挂载完成后，通过Alluxio统一命名空间，可以无缝地从不同存储系统中交互数据。例如，使用ls -R命令，递归地列举出一个目录下的所有文件。

```
alluxio fs ls -R /
```

```
0  PERSISTED 11-28-2019 11:15:19:719 DIR /app-logs  
1  PERSISTED 11-28-2019 11:18:36:885 DIR /apps  
1  PERSISTED 11-28-2019 11:18:40:209 DIR /apps/templeton  
239440292 PERSISTED 11-28-2019 11:18:40:209 0% /apps/templeton/hive.tar.gz  
.....  
1  PERSISTED 11-28-2019 19:00:23:879 DIR /mnt  
2  PERSISTED 11-28-2019 19:00:23:879 DIR /mnt/obs  
38  PERSISTED 11-28-2019 17:42:54:554 0% /mnt/obs/hive_load.txt  
12  PERSISTED 11-28-2019 17:43:07:743 0% /mnt/obs/test_input.txt  
.....
```

输出显示了Alluxio文件系统根目录（默认值是HDFS的根目录，即hdfs://hacluster/）中来源于挂载存储系统的所有文件。/app-logs和/apps目录在HDFS文件系统中，/mnt/obs/目录在OBS中。

用 Alluxio 加速数据访问

由于Alluxio利用内存存储数据，它可以加速数据的访问。例如：

1. 上传一个文件test_data.csv（文件是一份记录了食谱的样本）到obs-mrstest文件系统的/data目录下。通过ls命令显示文件状态：

```
alluxio fs ls /mnt/obs/test_data.csv
```

```
294520189 PERSISTED 11-28-2019 19:38:55:000 0% /mnt/obs/test_data.csv
```

输出显示了该文件在Alluxio中缓存占比为0%，即不在Alluxio内存中。

2. 统计该文件中单词"milk"出现的次数，并计算耗时。

```
time alluxio fs cat /mnt/obs/test_data.csv | grep -c milk
```

```
52180
```

```
real 0m10.765s  
user 0m5.540s  
sys 0m0.696s
```

3. 第一次读取数据后会将数据放在内存中，Alluxio再次读取时可以提高访问该数据的速度。例如：在通过cat命令获取文件后，用ls命令再查看文件的状态。

```
alluxio fs ls /mnt/obs/test_data.csv
```

```
294520189 PERSISTED 11-28-2019 19:38:55:000 100% /mnt/obs/test_data.csv
```

输出显示文件已经100%被加载到Alluxio中。

- 再次访问该文件，统计单词“eggs”出现的次数，并计算耗时。

```
time alluxio fs cat /mnt/obs/test_data.csv | grep -c eggs  
59510  
  
real 0m5.777s  
user 0m5.992s  
sys 0m0.592s
```

对比两次耗时可以看出存储在Alluxio内存中的数据，数据访问耗时明显缩短。

2 使用 CarbonData（MRS 3.x 之前版本）

2.1 从零开始使用 CarbonData

MRS 3.x之前版本参考本章节，MRS 3.x及后续版本请参考[使用CarbonData（MRS 3.x及之后版本）](#)。

本章节介绍使用Spark CarbonData的基本流程，所有任务场景基于spark-beeline环境。CarbonData快速入门包含以下任务：

1. 连接到Spark
在对CarbonData进行任何操作之前，需要先连接到Spark。
2. 创建CarbonData表
连接CarbonData之后，需要创建CarbonData Table，用于加载数据和执行查询操作。
3. 加载数据到CarbonData表
用户从HDFS中的CSV文件加载数据到所创建的表中。
4. 在CarbonData中查询数据
在CarbonData表加载数据之后，用户可以执行所需的查询操作，例如groupby或者where等。

前提条件

已安装客户端，具体参见[使用MRS客户端](#)。

操作步骤

步骤1 连接到Spark CarbonData。

1. 根据业务情况，准备好客户端，使用root用户登录安装客户端的节点。
例如在Master2节点更新客户端，则在该节点登录客户端，具体参见[使用MRS客户端](#)。
2. 切换用户与配置环境变量。

```
sudo su - omm
source /opt/client/bigdata_env
```

3. 启用Kerberos认证的集群，执行以下命令认证用户身份。未启用Kerberos认证集群无需执行。

kinit Spark组件用户名

📖 说明

用户需要加入用户组hadoop、hive，主组hadoop。

4. 执行以下命令，连接到Spark运行环境：

spark-beeline

步骤2 执行命令创建CarbonData表。

CarbonData表可用于加载数据和执行查询操作，例如执行以下命令创建CarbonData表：

```
CREATE TABLE x1 (imei string, deviceInformationId int, mac string,
productdate timestamp, updatetime timestamp, gamePointId double,
contractNumber double)
```

```
STORED BY 'org.apache.carbondata.format'
```

```
TBLPROPERTIES
```

```
('DICTIONARY_EXCLUDE'='mac','DICTIONARY_INCLUDE'='deviceInformationId'
);
```

命令执行结果如下：

```
+-----+
| result |
+-----+
+-----+
No rows selected (1.551 seconds)
```

步骤3 从CSV文件加载数据到CarbonData表。

根据所要求的参数运行命令从CSV文件加载数据，且仅支持CSV文件。**LOAD**命令中配置的CSV列名，需要和CarbonData表列名相同，顺序也要对应。CSV文件中的数据的列数，以及数据格式需要和CarbonData表匹配。

文件需要保存在HDFS中。用户可以将文件上传到OBS，并在MRS管理控制台“文件管理”将文件从OBS导入HDFS，具体请参考[导入导出数据](#)。

如果集群启用了Kerberos认证，则需要在工作环境准备CSV文件，然后可以使用开源HDFS命令，参考[5](#)将文件从工作环境导入HDFS，并设置Spark组件用户在HDFS中对文件有读取和执行的权限。

例如，HDFS的“tmp”目录有一个文件“data.csv”，内容如下：

```
x123,111,dd,2017-04-20 08:51:27,2017-04-20 07:56:51,2222,33333
```

执行导入命令：

```
LOAD DATA inpath 'hdfs://hacluster/tmp/data.csv' into table x1
options('DELIMITER'=',','QUOTECHAR'='','FILEHEADER'='imei,
deviceinformationid,mac,productdate,updatetime,gamepointid,contractnumb
er');
```

命令执行结果如下：

```
+-----+
| Result |
```

```
+-----+  
+-----+  
No rows selected (3.039 seconds)
```

步骤4 在CarbonData中查询数据。

- **获取记录数**

为了获取在CarbonData table中的记录数，可以执行以下命令。

```
select count(*) from x1;
```

- **使用Groupby查询**

为了获取不重复的“deviceinformationid”记录数，可以执行以下命令。

```
select deviceinformationid,count (distinct deviceinformationid) from x1  
group by deviceinformationid;
```

- **使用条件查询**

为了获取特定deviceinformationid的记录，可以执行以下命令。

```
select * from x1 where deviceinformationid='111';
```

 **说明**

在执行数据查询操作后，如果查询结果中某一列的结果含有中文字等其他非英文字符，会导致查询结果中的列不能对齐，这是由于不同语言的字符在显示时所占的字宽不尽相同。

步骤5 执行以下命令退出Spark运行环境。

```
!quit
```

```
----结束
```

2.2 CarbonData 表简介

简介

CarbonData表与RDBMS中的表类似，RDBMS数据存储在与行和列构成的表中。CarbonData表存储的也是结构化的数据，具有固定列和数据类型。CarbonData中的数据存储在与表实体文件中。

支持的数据类型

CarbonData表支持以下数据类型：

- Int
- String
- BigInt
- Decimal
- Double
- TimeStamp

[表2-1](#)对所支持的数据类型和对应的范围进行了详细说明。

表 2-1 CarbonData 数据类型

数据类型	描述
Int	4字节有符号整数，从-2,147,483,648到2,147,483,647。 说明 非字典列如果是Int类型，会在内部存储为BigInt类型。
String	最大支持字符长度为100000。
BigInt	使用64-bit存储数据，支持从-9,223,372,036,854,775,808到9,223,372,036,854,775,807。
Decimal	默认值是(10,0)，最大值是(38,38)。 说明 当进行带过滤条件的查询时，为了得到准确的结果，需要在数字后面加上BD。例如， <code>select * from carbon_table where num = 1234567890123456.22BD</code> 。
Double	使用64-bit存储数据，从4.9E-324到1.7976931348623157E308。
TimeStamp	默认格式为“yyyy-MM-dd HH:mm:ss”。

 说明

所有Integer类型度量均以BigInt类型进行处理与显示。

2.3 创建 CarbonData 表

操作场景

使用CarbonData前需先创建表，才可从表中加载数据和查询数据。

使用自定义列创建表

可通过指定各列及其数据类型来创建表。启用Kerberos认证的分析集群创建CarbonData表时，如果用户需要在默认数据库“default”以外的数据库创建新表，则需要Hive角色管理中为用户绑定的角色添加指定数据库的“Create”权限。

命令示例：

```
CREATE TABLE IF NOT EXISTS productdb.productSalesTable (
productNumber Int,
productName String,
storeCity String,
storeProvince String,
revenue Int)
STORED BY 'org.apache.carbondata.format'
TBLPROPERTIES (
```

```
'table_blocksize'='128',
'DICTIONARY_EXCLUDE'='productName',
'DICTIONARY_INCLUDE'='productNumber');
```

上述命令所创建的表的详细信息如下：

表 2-2 表信息定义

参数	描述
productSalesTable	待创建的表的名称。该表用于加载数据进行分析。 表名由字母、数字、下划线组成。
productdb	数据库名称。该数据库将与其中的表保持逻辑连接以便于识别和管理。 数据库名称由字母、数字、下划线组成。
productNumber productName storeCity storeProvince revenue	表中的列，代表执行分析所需的业务实体。 列名（字段名）由字母、数字、下划线组成。 说明 CarbonData暂不支持设置列是否允许为空、默认值以及主键。
table_blocksize	CarbonData表使用的数据文件的block大小，默认值为1024，取值范围为1~2048，单位为MB。 <ul style="list-style-type: none"> 如果“table_blocksize”值太小，数据加载时将生成过多的小数据文件，可能会影响HDFS的使用性能。 如果“table_blocksize”值太大，数据查询时索引匹配的block数据量较大，导致读取并发度不高，从而降低查询性能。 一般情况下，建议根据数据量级别来选择大小。例如：GB级别用256，TB级别用512，PB级别用1024。
DICTIONARY_EXCLUDE	设置指定列不生成字典，适用于数值复杂度高的列。系统默认为String类型的列做字典编码，但是如果字典值过多，会导致字典转换操作增加造成性能下降。 一般情况下，列的数值复杂度高于5万，可以被认定为高复杂度，则需要排除掉字典编码，该参数为可选参数。 说明 在非字典列中，只支持String和Timestamp数据类型。
DICTIONARY_INCLUDE	设置指定列生成字典，适用于数值复杂度低的列，可以提升字典列上的groupby性能，为可选参数。一般情况下，字典列的复杂度不应该高于5万。

2.4 删除 CarbonData 表

操作场景

用户根据业务使用情况，可以删除不再使用的 CarbonData 表。删除表后，其所有的元数据以及表中已加载的数据都会被删除。

操作步骤

步骤1 运行如下命令删除表。

```
DROP TABLE [IF EXISTS] [db_name.]table_name;
```

“db_name”为可选参数。如果没有指定“db_name”，那么将会删除当前数据库下名为“table_name”的表。

例如执行命令，删除数据库“productdb”下的表“productSalesTable”：

```
DROP TABLE productdb.productSalesTable;
```

步骤2 执行以下命令查询表是否被删除：

```
SHOW TABLES;
```

----结束

3 使用 CarbonData（MRS 3.x 及之后版本）

3.1 CarbonData 数据类型概述

简介

CarbonData 中的数据存储在 table 实体中。CarbonData table 与 RDBMS 中的表类似。RDBMS 数据存储在由行和列构成的表中。CarbonData table 存储的也是结构化的数据，拥有固定列和数据类型。

支持数据类型

CarbonData 支持以下数据类型：

- Int
- String
- BigInt
- Smallint
- Char
- Varchar
- Boolean
- Decimal
- Double
- TimeStamp
- Date
- Array
- Struct
- Map

下表对所支持的数据类型及其各自的范围进行了详细说明。

表 3-1 CarbonData 数据类型

数据类型	范围
Int	4字节有符号整数，从-2,147,483,648到2,147,483,647 说明 非字典列如果是Int类型，会在内部存储为BigInt类型。
String	100000字符 说明 如果在CREATE TABLE中使用Char或Varchar数据类型，则这两种数据类型将自动转换为String数据类型。 如果存在字符长度超过32000的列，需要在建表时，将该列加入到tblproperties的LONG_STRING_COLUMNS属性里。
BigInt	64-bit，从-9,223,372,036,854,775,808到9,223,372,036,854,775,807
SmallInt	范围-32,768到32,767
Char	范围A到Z&a到z
Varchar	范围A到Z&a到z&0到9
Boolean	范围true或者false
Decimal	默认值是(10,0)，最大值是(38,38) 说明 当进行带过滤条件的查询时，为了得到准确的结果，需要在数字后面加上BD。例如， <i>select * from carbon_table where num = 1234567890123456.22BD.</i>
Double	64-bit，从4.9E-324到1.7976931348623157E308
TimeStamp	NA，默认格式为“yyyy-MM-dd HH:mm:ss”。
Date	DATE数据类型用于存储日历日期。默认格式为“yyyy-MM-dd”。
Array<data_type>	NA
Struct<col_name: data_type COMMENT col_comment, ...>	说明 现仅支持2层复杂类型的嵌套。
Map<primitive_type, data_type>	

CarbonData 主要规格

表 3-2 CarbonData 主要规格

实体	测试值	测试环境
表数	10000	3个节点，每个executor 4个CPU核，20GB。Driver内存5GB，3个Executor。 总列数：107 String：75 Int：13 BigInt：7 Timestamp：6 Double：6
表的列数	2000	3个节点，每个executor 4个CPU核，20GB。Driver内存5GB，3个Executor。
原始CSV文件大小的最大值	200GB	17个cluster节点，每个executor 150GB，25个CPU核。Driver内存10 GB，17个Executor。
每个文件夹的CSV文件数	100个文件夹，每个文件夹10个文件，每个文件大小50MB。	3个节点，每个executor 4个CPU核，20GB。Driver内存5GB，3个Executor。
加载文件夹数	10000	3个节点，每个executor 4个CPU核，20GB。Driver内存5GB，3个Executor。

数据加载所需的内存取决于以下因素：

- 列数
- 列值大小
- 并发（使用“carbon.number.of.cores.while.loading”进行配置）
- 在内存中排序的大小（使用“carbon.sort.size”进行配置）
- 中间缓存（使用“carbon.graph.rowset.size”进行配置）

加载包含1000万条记录和300列的8 GB CSV文件的数据，每行大小约为0.8KB的8GB CSV文件的数据，需要约为10GB的executor执行内存，也就是说，“carbon.sort.size”配置为“100000”，所有其他前面的配置保留默认值。

二级索引表规格

表 3-3 二级索引表规格

实体	测试值
二级索引表数量	10
二级索引表中的组合列的列数	5
二级索引表中的列名长度（单位：字符）	120
二级索引表名长度（单位：字符）	120
表中所有二级索引表的表名+列名的累积长度*（单位：字符）	3800**

说明

- * Hive允许的上限值或可用资源的上限值。
- ** 二级索引表使用hive注册，并以json格式的值存储在HiveSERDEPROPERTIES中。由hive支持的SERDEPROPERTIES的最大字符数为4000个字符，无法更改。

3.2 CarbonData 表用户权限说明

下表提供了对CarbonData Table执行相应操作所需的Hive ACL特权的信息。

前提条件

已经设置了[表3-21](#)或[表3-22](#)中Carbon相关参数。

Hive ACL 权限

表 3-4 CarbonData 表级操作所需的 Hive ACL 权限

场景	所需权限
DESCRIBE TABLE	SELECT (of table)
SELECT	SELECT (of table)
EXPLAIN	SELECT (of table)
CREATE TABLE	CREATE (of database)
CREATE TABLE As SELECT	CREATE (on database), INSERT (on table), RW on data file, and SELECT (on table)
LOAD	INSERT (of table) RW on data file
DROP TABLE	OWNER (of table)

场景	所需权限
DELETE SEGMENTS	DELETE (of table)
SHOW SEGMENTS	SELECT (of table)
CLEAN FILES	DELETE (of table)
INSERT OVERWRITE / INSERT INTO	INSERT (of table) RW on data file and SELECT (of table)
CREATE INDEX	OWNER (of table)
DROP INDEX	OWNER (of table)
SHOW INDEXES	SELECT (of table)
ALTER TABLE ADD COLUMN	OWNER (of table)
ALTER TABLE DROP COLUMN	OWNER (of table)
ALTER TABLE CHANGE DATATYPE	OWNER (of table)
ALTER TABLE RENAME	OWNER (of table)
ALTER TABLE COMPACTION	INSERT (on table)
FINISH STREAMING	OWNER (of table)
ALTER TABLE SET STREAMING PROPERTIES	OWNER (of table)
ALTER TABLE SET TABLE PROPERTIES	OWNER (of table)
UPDATE CARBON TABLE	UPDATE (of table)
DELETE RECORDS	DELETE (of table)
REFRESH TABLE	OWNER (of main table)
REGISTER INDEX TABLE	OWNER (of table)
SHOW PARTITIONS	SELECT (on table)
ALTER TABLE ADD PARTITION	OWNER (of table)
ALTER TABLE DROP PARTITION	OWNER (of table)

 说明

- 如果数据库下的表由多个用户创建，那么执行Drop database命令会失败，即使执行的用户是数据库的拥有者。
- 在二级索引中，当父表（parent table）触发时，insert和compaction将在索引表上触发。如果选择具有过滤条件匹配索引表列的查询，用户应该为父表和索引表提供选择权限。
- LockFiles文件夹和LockFiles文件夹中创建的锁定文件将具有完全权限，因为LockFiles文件夹不包含任何敏感数据。
- 如果使用ACL，确保不要为DDL或DML配置任何被其他进程使用中的路径，建议创建新路径。
以下配置项需要配置路径：
 - 1) carbon.badRecords.location
 - 2) 创建数据库时Db_Path及其他。
- 对于非安全集群中的Carbon ACL权限，hive-site.xml中的参数hive.server2.enable.doAs必须设置为false。将此属性设置为false，查询将以hiveserver2进程运行的用户身份运行。

3.3 使用 Spark 客户端创建 CarbonData 表

本章节介绍创建CarbonData table、加载数据，以及查询数据的快速入门流程。该快速入门提供基于Spark Beeline客户端的操作。如果使用Spark shell，需将查询命令写在spark.sql()的括号中。

本操作以从CSV文件加载数据到CarbonData Table为例

表 3-5 CarbonData 快速入门

操作	说明
准备CSV文件	准备加载到CarbonData Table的CSV文件。
连接到CarbonData	在对CarbonData进行任何一种操作之前，首先需要连接到CarbonData。
创建CarbonData Table	连接到CarbonData之后，需要创建CarbonData table用于加载数据和执行查询操作。
加载数据到CarbonData Table	创建CarbonData table之后，可以从CSV文件加载数据到所创建的table中。
在CarbonData中查询数据	创建CarbonData table并加载数据之后，可以执行所需的查询操作，例如filters, groupby等。

准备 CSV 文件

1. 在本地准备CSV文件，文件名为：test.csv，样例如下：

```
13418592122,1001,MAC地址,2017-10-23 15:32:30,2017-10-24 15:32:30,62.50,74.56
13418592123,1002,MAC地址,2017-10-23 16:32:30,2017-10-24 16:32:30,17.80,76.28
13418592124,1003,MAC地址,2017-10-23 17:32:30,2017-10-24 17:32:30,20.40,92.94
13418592125,1004,MAC地址,2017-10-23 18:32:30,2017-10-24 18:32:30,73.84,8.58
13418592126,1005,MAC地址,2017-10-23 19:32:30,2017-10-24 19:32:30,80.50,88.02
13418592127,1006,MAC地址,2017-10-23 20:32:30,2017-10-24 20:32:30,65.77,71.24
13418592128,1007,MAC地址,2017-10-23 21:32:30,2017-10-24 21:32:30,75.21,76.04
13418592129,1008,MAC地址,2017-10-23 22:32:30,2017-10-24 22:32:30,63.30,94.40
```

```
13418592130,1009,MAC地址,2017-10-23 23:32:30,2017-10-24 23:32:30,95.51,50.17  
13418592131,1010,MAC地址,2017-10-24 00:32:30,2017-10-25 00:32:30,39.62,99.13
```

2. 使用WinSCP工具将CSV文件导入客户端节点，例如“/opt”目录下。
3. 登录FusionInsight Manager页面，选择“系统 > 权限 > 用户”，添加人机用户sparkuser，用户组（hadoop、hive），主组（hadoop）。
4. 进入客户端目录，加载环境变量并认证用户：

```
cd /客户端安装目录  
source ./bigdata_env  
source ./Spark2x/component_env  
kinit sparkuser
```

5. 上传CSV中的文件到HDFS的“/data”目录：
hdfs dfs -put /opt/test.csv /data/

连接到 CarbonData

- 使用Spark SQL或Spark shell连接到Spark并执行Spark SQL命令。
- 开启JDBCServer并使用JDBC客户端（例如，Spark Beeline）连接。
执行如下命令：

```
cd ./Spark2x/spark/bin  
./spark-beeline
```

创建 CarbonData Table

在Spark Beeline被连接到JDBCServer之后，需要创建一个CarbonData table用于加载数据和执行查询操作。下面是创建一个简单的表的命令。

```
create table x1 (imei string, deviceInformationId int, mac string, productdate  
timestamp, updatetime timestamp, gamePointId double, contractNumber  
double) STORED AS carbondata TBLPROPERTIES  
(SORT_COLUMNS='imei,mac');
```

命令执行结果如下：

```
+-----+  
| Result |  
+-----+  
+-----+  
No rows selected (1.093 seconds)
```

加载数据到 CarbonData Table

创建CarbonData table之后，可以从CSV文件加载数据到所创建的表中。

用所要求的参数运行以下命令从CSV文件加载数据。该表的列名需要与CSV文件的列名匹配。

```
LOAD DATA inpath 'hdfs://hacluster/data/test.csv' into table x1  
options('DELIMITER',';', 'QUOTECHAR=''','FILEHEADER='imei,  
deviceinformationid,mac, productdate,updatetime,  
gamepointid,contractnumber');
```

其中，“test.csv”为[准备CSV文件](#)的CSV文件，“x1”为示例的表名。

CSV样例内容如下：

```
13418592122,1001,MAC地址,2017-10-23 15:32:30,2017-10-24 15:32:30,62.50,74.56
13418592123,1002,MAC地址,2017-10-23 16:32:30,2017-10-24 16:32:30,17.80,76.28
13418592124,1003,MAC地址,2017-10-23 17:32:30,2017-10-24 17:32:30,20.40,92.94
13418592125,1004,MAC地址,2017-10-23 18:32:30,2017-10-24 18:32:30,73.84,8.58
13418592126,1005,MAC地址,2017-10-23 19:32:30,2017-10-24 19:32:30,80.50,88.02
13418592127,1006,MAC地址,2017-10-23 20:32:30,2017-10-24 20:32:30,65.77,71.24
13418592128,1007,MAC地址,2017-10-23 21:32:30,2017-10-24 21:32:30,75.21,76.04
13418592129,1008,MAC地址,2017-10-23 22:32:30,2017-10-24 22:32:30,63.30,94.40
13418592130,1009,MAC地址,2017-10-23 23:32:30,2017-10-24 23:32:30,95.51,50.17
13418592131,1010,MAC地址,2017-10-24 00:32:30,2017-10-25 00:32:30,39.62,99.13
```

命令执行结果如下：

```
+-----+
|Segment ID |
+-----+
|0          |
+-----+
No rows selected (3.039 seconds)
```

在 CarbonData 中查询数据

创建CarbonData table并加载数据之后，可以执行所需的数据查询操作。以下为一些查询操作举例。

- **获取记录数**

为了获取在CarbonData table中的记录数，可以运行以下命令。

```
select count(*) from x1;
```

- **使用Groupby查询**

为了获取不重复的deviceinformationid记录数，可以运行以下命令。

```
select deviceinformationid,count (distinct deviceinformationid) from x1
group by deviceinformationid;
```

- **用Filter查询**

为了获取特定deviceinformationid的记录，可以运行以下命令。

```
select * from x1 where deviceinformationid='1010';
```

说明

在执行数据查询操作后，如果查询结果中某一列的结果含有中文字等非英文字符，会导致查询结果中的列不能对齐，这是由于不同语言的字符在显示时所占的字宽不尽相同。

在 Spark-shell 上使用 CarbonData

用户若需要在Spark-shell上使用CarbonData，需通过如下方式创建CarbonData Table，加载数据到CarbonData Table和在CarbonData中查询数据的操作。

```
spark.sql("CREATE TABLE x2(imei string, deviceInformationId int, mac string, productdate timestamp,
updatetime timestamp, gamePointId double, contractNumber double) STORED AS carbondata")
spark.sql("LOAD DATA inpath 'hdfs://hacluster/data/x1_without_header.csv' into table x2
options('DELIMITER',';', 'QUOTECHAR','\",'FILEHEADER'='imei, deviceinformationid,mac,
productdate,updatetime, gamepointid,contractnumber')")
spark.sql("SELECT * FROM x2").show()
```

3.4 CarbonData 数据分析

3.4.1 新建 CarbonData Table

操作场景

使用CarbonData前需先创建表，才可在其中加载数据和查询数据。可通过**Create Table**命令来创建表。该命令支持使用自定义列创建表。

使用自定义列创建表

可通过指定各列及其数据类型来创建表。

命令示例：

```
CREATE TABLE IF NOT EXISTS productdb.productSalesTable (  
  productNumber Int,  
  productName String,  
  storeCity String,  
  storeProvince String,  
  productCategory String,  
  productBatch String,  
  saleQuantity Int,  
  revenue Int)  
STORED AS carbondata  
TBLPROPERTIES (  
  'table_blocksize'='128');
```

上述命令所创建的表的详细信息如下：

表 3-6 表信息定义

参数	描述
productSalesTable	待创建的表的名称。该表用于加载数据进行分析。 表名由字母、数字、下划线组成。
productdb	数据库名称。该数据库将与其中的表保持逻辑连接以便于识别和管理。 数据库名称由字母、数字、下划线组成。

参数	描述
productName storeCity storeProvince productCategory productBatch saleQuantity revenue	表中的列，代表执行分析所需的业务实体。 列名（字段名）由字母、数字、下划线组成。
table_blocksize	CarbonData表使用的数据文件的block大小，默认值为1024，最小值为1，最大值为2048，单位为MB。 如果“table_blocksize”值太小，数据加载时，生成过多的小数据文件，可能会影响HDFS的使用性能。 如果“table_blocksize”值太大，数据查询时，索引匹配的block数据量较大，某些block会包含较多的blocklet，导致读取并发度不高，从而降低查询性能。 一般情况下，建议根据数据量级别来选择大小。例如：GB级别用256，TB级别用512，PB级别用1024。

📖 说明

- 所有Integer类型度量均以BigInt类型进行处理与显示。
- CarbonData遵循严格解析，因此任何不可解析的数据都会被保存为null。例如，在BigInt列中加载double值（3.14），将会保存为null。
- 在Create Table中使用的Short和Long数据类型在DESCRIBE命令中分别显示为Smallint和Bigint。
- 可以使用DESCRIBE格式化命令查看表数据大小和表索引大小。

操作结果

根据命令创建表。

3.4.2 删除 CarbonData Table

操作场景

可使用**DROP TABLE**命令删除表。删除表后，所有metadata以及表中已加载的数据都会被删除。

操作步骤

运行如下命令删除表。

命令：

```
DROP TABLE [IF EXISTS] [db_name.]table_name;
```

一旦执行该命令，将会从系统中删除表。命令中的“db_name”为可选参数。如果没有指定“db_name”，那么将会删除当前数据库下名为“table_name”的表。

示例：

```
DROP TABLE productdb.productSalesTable;
```

通过上述命令，删除数据库“productdb”下的表“productSalesTable”。

操作结果

从系统中删除命令中指定的表。删除完成后，可通过**SHOW TABLES**命令进行查询，确认所需删除的表是否成功被删除，详见**SHOW TABLES**。

3.4.3 修改 CarbonData Table

SET 和 UNSET

当使用set命令时，所有新set的属性将会覆盖已存在的旧的属性。

- SORT SCOPE

SET SORT SCOPE命令示例：

```
ALTER TABLE tablename SET TBLPROPERTIES('SORT_SCOPE'='no_sort')
```

当UNSET SORT SCOPE后，会使用默认值NO_SORT。

UNSET SORT SCOPE命令示例：

```
ALTER TABLE tablename UNSET TBLPROPERTIES('SORT_SCOPE')
```

- SORT COLUMNS

SET SORT COLUMNS命令示例：

```
ALTER TABLE tablename SET TBLPROPERTIES('SORT_COLUMNS'='column1')
```

在执行该命令后，新的导入会使用新的SORT_COLUMNS配置值。用户可以根据查询的情况来调整SORT_COLUMNS，但是不会直接影响旧的数据。所以对历史的segments的查询性能不会受到影响，因为历史的segments不是按照新的SORT_COLUMNS。

不支持UNSET命令，但是可以使用set SORT_COLUMNS等于空字符串来代替UNSET命令。

```
ALTER TABLE tablename SET TBLPROPERTIES('SORT_COLUMNS'='')
```

说明

- 后续版本会加强自定义合并来对旧的segment重新排序。
- 流式表不支持修改SORT_COLUMNS。
- 如果inverted index的列从SORT_COLUMNS里面移除了，该列不会再创建inverted index。但是旧的INVERTED_INDEX配置值不会变化。

3.4.4 加载 CarbonData 表数据

操作场景

CarbonData table创建成功后，可使用**LOAD DATA**命令在表中加载数据，并可供查询。

触发数据加载后，数据以CarbonData格式进行编码，并将多维列式存储格式文件压缩后复制到存储CarbonData文件的HDFS路径下供快速分析查询使用。

HDFS路径可以配置在carbon.properties文件中。

CarbonData相关配置参数请参考[CarbonData常见配置参数](#)。

3.4.5 删除 CarbonData 表 Segments

操作场景

如果用户将错误数据加载到表中，或者数据加载后出现许多错误记录，用户希望修改并重新加载数据时，可删除对应的segment。可使用segment ID来删除segment，也可以使用加载数据的时间来删除segment。

📖 说明

删除segment操作只能删除未合并的segment，已合并的segment可以通过**CLEAN FILES**命令清除segment。

通过 Segment ID 删除

每个Segment都有与其关联的唯一Segment ID。使用这个Segment ID可以删除该Segment。

步骤1 运行如下命令获取Segment ID。

命令：

```
SHOW SEGMENTS FOR Table dbname.tablename LIMIT number_of_loads;
```

示例：

```
SHOW SEGMENTS FOR TABLE carbonTable;
```

上述命令可显示tablename为carbonTable的表的所有Segment信息。

```
SHOW SEGMENTS FOR TABLE carbonTable LIMIT 2;
```

上述命令可显示*number_of_loads*规定条数的Segment信息。

输出结果如下：

```
+-----+-----+-----+-----+-----+-----+-----+-----+
+
| ID | Status | Load Start Time | Load Time Taken | Partition | Data Size | Index Size | File Format |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 3 | Success | 2020-09-28 22:53:26.336 | 3.726S | {} | 6.47KB | 3.30KB | columnar_v3 |
| 2 | Success | 2020-09-28 22:53:01.702 | 6.688S | {} | 6.47KB | 3.30KB | columnar_v3 |
+-----+-----+-----+-----+-----+-----+-----+-----+
```

📖 说明

SHOW SEGMENTS命令输出包括ID、Status、Load Start Time、Load Time Taken、Partition、Data Size、Index Size、File Format。最新的加载信息在输出中第一行显示。

步骤2 获取到需要删除的Segment的Segment ID后，执行如下命令删除对应Segment：

命令：

```
DELETE FROM TABLE tableName WHERE SEGMENT.ID IN (load_sequence_id1,  
load_sequence_id2, ....);
```

示例：

```
DELETE FROM TABLE carbonTable WHERE SEGMENT.ID IN (1,2,3);
```

详细信息，请参阅[DELETE SEGMENT by ID](#)。

----结束

通过加载数据的时间删除

用户可基于特定的加载时间删除数据。

命令：

```
DELETE FROM TABLE db_name.table_name WHERE SEGMENT.STARTTIME  
BEFORE date_value;
```

示例：

```
DELETE FROM TABLE carbonTable WHERE SEGMENT.STARTTIME BEFORE  
'2017-07-01 12:07:20';
```

上述命令可删除'2017-07-01 12:07:20'之前的所有segment。

有关详细信息，请参阅[DELETE SEGMENT by DATE](#)。

删除结果

数据对应的segment被删除，数据将不能再被访问。可通过**SHOW SEGMENTS**命令显示segment状态，查看是否成功删除。

说明

- 调用**DELETE SEGMENT**命令时，物理上而言，Segment并没有从文件系统中被删除。使用命令**SHOW SEGMENTS**查看Segment信息，可看见被删除的Segment的状态被标识为“Marked for Delete”。但使用**SELECT * FROM tablename**命令查询时，不会显示被删除的Segment的内容。
- 下一次加载数据且达到最大查询执行时间（由“max.query.execution.time”配置，默认为“60分钟”）后，Segment才会从文件系统中真正删除。
- 如果用户想要强制删除物理Segment文件，那么可以使用**CLEAN FILES**命令。

示例：

```
CLEAN FILES FOR TABLE table1;
```

该命令将从物理上删除状态为“Marked for delete”的Segment文件。

如果在“max.query.execution.time”规定的时间到达之前使用该命令，可能会导致查询失败。“max.query.execution.time”可在“carbon.properties”文件中设置，表示一次查询允许花费的最长时间。

3.4.6 合并 CarbonData 表 Segments

操作场景

频繁的数据获取导致在存储目录中产生许多零碎的CarbonData文件。由于数据排序只在每次加载时进行，所以，索引也只在每次加载时执行。这意味着，对于每次加载都

会产生一个索引，随着数据加载数量的增加，索引的数量也随之增加。由于每个索引只在一次加载时工作，索引的性能被降低。CarbonData提供加载压缩。压缩过程通过合并排序各segment中的数据，将多个segment合并为一个大的segment。

前提条件

已经加载了多次数据。

操作描述

有Minor合并、Major合并和Custom合并三种类型。

- Minor合并：**
 在Minor合并中，用户可指定合并数据加载的数量。如果设置了参数“carbon.enable.auto.load.merge”，每次数据加载都可触发Minor合并。如果任意segment均可合并，那么合并将与数据加载时并行进行。
 Minor合并有两个级别。
 - Level 1：合并未合并的segment。
 - Level 2：合并已合并的segment，以形成更大的segment。
- Major合并：**
 在Major合并中，许多segment可以合并为一个大的segment。用户将指定合并尺寸，将对未达到该尺寸的segment进行合并。Major合并通常在非高峰时段进行。
- Custom合并：**
 在Custom合并中，用户可以指定几个segment的id合并为一个大的segment。所有指定的segment的id必须存在并且有效，否则合并将会失败。Custom合并通常在非高峰时段进行。

具体的命令操作，请参考[ALTER TABLE COMPACTION](#)。

表 3-7 合并参数

参数	默认值	应用类型	描述
carbon.enable.auto.load.merge	false	Minor	数据加载时启用合并。 “true”：数据加载时自动触发segment合并。 “false”：数据加载时不触发segment合并。
carbon.compaction.level.threshold	4,3	Minor	对于Minor合并，该属性参数决定合并segment的数量。 例如，如果该参数设置为“2,3”，在Level 1，每2个segment触发一次Minor合并。在Level2，每3个Level 1合并的segment将被再次合并为新的segment。 合并策略根据实际的数据大小和可用资源决定。 有效值为0-100。

参数	默认值	应用类型	描述
carbon.major.compaction.size	1024mb	Major	通过配置该参数可配置Major合并。低于该阈值的segment之和将被合并。 例如，如果该阈值是1024MB，且有5个大小依次为300MB，400MB，500MB，200MB，100MB的segment用于Major合并，那么只有相加的总数小于阈值的segment会被合并，也就是300+400+200+100 = 1000MB的segment会被合并，而500MB的segment将会被跳过。
carbon.numberof.preserve.segments	0	Minor/ Major	如果用户希望从被合并的segment中保留一定数量的segment，可通过该属性参数进行设置。 例如， “carbon.numberof.preserve.segments” = “2”，那么最新的2个segment将不会包含在合并中。 默认不保留任何segment。
carbon.allowed.compaction.days	0	Minor/ Major	合并将合并指定的配置天数中加载的segment。 例如，如果配置为“2”，那么只有在2天的时间框架中被加载的segment可以被合并。在2天以外被加载的segment将不被合并。 默认为禁用。
carbon.number.of.cores.while.compacting	2	Minor/ Major	在合并过程中写入数据时所用的核数。配置的核数越大合并性能越好。如果CPU资源充足可以增加此值。
carbon.merge.index.in.segment	true	SEGMENT_INDEX	如果设置为true，则一个segment中所有Carbon索引文件（.carbonindex）将合并为单个Carbon索引合并文件（.carbonindexmerge）。这增强了首次查询性能。

参考信息

建议避免对历史数据进行minor compaction，请参考[如何避免对历史数据进行minor compaction?](#)

3.5 CarbonData 性能调优

3.5.1 CarbonData 调优思路

查询性能调优

CarbonData可以通过调整各种参数来提高查询性能。大部分参数聚焦于增加并行性处理和更好地使用系统资源。

- Spark Executor数量：Executor是Spark并行性的基础实体。通过增加Executor数量，集群中的并行数量也会增加。关于如何配置Executor数量，请参考Spark资料。
- Executor核：每个Executor内，并行任务数受Executor核的配置控制。通过增加Executor核数，可增加并行任务数，从而提高性能。
- HDFS block容量：CarbonData通过给不同的处理器分配不同的block来分配查询任务。所以一个HDFS block是一个分区单元。另外，CarbonData在Spark驱动器中，支持全局block级索引，这有助于减少需要被扫描的查询block的数量。设置较大的block容量，可提高I/O效率，但是会降低全局索引效率；设置较小的block容量，意味着更多的block数量，会降低I/O效率，但是会提高全局索引效率，同时，对于索引查询会要求更多的内存。
- 扫描线程数量：扫描仪（Scanner）线程控制每个任务中并行处理的数据块的数量。通过增加扫描仪线程数，可增加并行处理的数据块的数量，从而提高性能。可使用“carbon.properties”文件中的“carbon.number.of.cores”属性来配置扫描仪线程数。例如，“carbon.number.of.cores = 4”。
- B-Tree缓存：为了获得更好的查询特性，可以通过B-tree LRU（least recently used，最近最少使用）缓存来优化缓存内存。在driver中，B-Tree LRU缓存配置将有助于通过释放未被访问或未使用的表segments来释放缓存。类似地，在executor中，B-Tree LRU缓存配置将有助于释放未被访问或未使用的表blocks。具体可参考表3-18中的参数“carbon.max.driver.lru.cache.size”和“carbon.max.executor.lru.cache.size”的详细描述。

CarbonData 查询流程

当CarbonData首次收到对某个表（例如表A）的查询任务时，系统会加载表A的索引数据到内存中，执行查询流程。当CarbonData再次收到对表A的查询任务时，系统则不需要再加载其索引数据。

在CarbonData中执行查询时，查询任务会被分成几个扫描任务。即，基于CarbonData数据存储的HDFS block对扫描任务进行分割。扫描任务由集群中的执行器执行。扫描任务可以并行、部分并行，或顺序处理，具体采用的方式取决于执行器的数量以及配置的执行器核数。

查询任务的某些部分可在独立的任务级上处理，例如select和filter。查询任务的某些部分可在独立的任务级上进行部分处理，例如group-by、count、distinct count等。

某些操作无法在任务级上处理，例如Having Clause（分组后的过滤），sort等。这些无法在任务级上处理，或只能在任务级上部分处理的操作需要在集群内跨执行器来传输数据（部分结果）。这个传送操作被称为shuffle。

任务数量越多，需要shuffle的数据就越多，会对查询性能产生不利影响。

由于任务数量取决于HDFS block的数量，而HDFS block的数量取决于每个block的大小，因此合理选择HDFS block的大小很重要，需要在提高并行性，进行shuffle操作的数据量和聚合表的大小之间达到平衡。

分割和 Executors 的关系

如果分割数小于等于Executor数乘以Executor核数，那么任务将以并行方式运行。否则，某些任务只有在其他任务完成之后才能开始。因此，要确保Executor数乘以Executor核数大于等于分割数。同时，还要确保有足够的分割数，这样一个查询任务可被分为足够多的子任务，从而确保并行性。

配置扫描仪线程

扫描仪线程属性决定了每个分割的数据被划分的可并行处理的数据块的数量。如果数量过多，会产生很多小数据块，性能会受到影响。如果数量过少，并行性不佳，性能也会受到影响。因此，决定扫描仪线程数时，需要考虑一个分割内的平均数据大小，选择一个使数据块不会很小的值。经验法则是将单个块大小（MB）除以250得到的值作为扫描仪线程数。

增加并行性还需考虑的重要一点是集群中实际可用的CPU核数，确保并行计算数不超过实际CPU核数的75%至80%。

CPU核数约等于：

并行任务数x扫描仪线程数。其中并行任务数为分割数和执行器数x执行器核数两者之间的较小值。

数据加载性能调优

数据加载性能调优与查询性能调优差异很大。跟查询性能一样，数据加载性能也取决于可达到的并行性。在数据加载情况下，工作线程的数量决定并行的单元。因此，更多的执行器就意味着更多的执行器核数，每个执行器都可以提高数据加载性能。

同时，为了得到更好的性能，可在HDFS中配置如下参数。

表 3-8 HDFS 配置

参数	建议值
dfs.datanode.drop.cache.behind.reads	false
dfs.datanode.drop.cache.behind.writes	false
dfs.datanode.sync.behind.writes	true

压缩调优

CarbonData结合少数轻量级压缩算法和重量级压缩算法来压缩数据。虽然这些算法可处理任何类型的数据，但如果数据经过排序，相似值在一起出现时，就会获得更好的压缩率。

CarbonData数据加载过程中，数据基于Table中的列顺序进行排序，从而确保相似值在一起出现，以获得更好的压缩率。

由于CarbonData按照Table中定义的列顺序将数据进行排序，因此列顺序对于压缩效率起重要作用。如果低cardinality维度位于左边，那么排序后的数据分区范围较小，压缩效率较高。如果高cardinality维度位于左边，那么排序后的数据分区范围较大，压缩效率较低。

内存调优

CarbonData为内存调优提供了一个机制，其中数据加载会依赖于查询中需要的列。不论何时，接收到一个查询命令，将会获取到该查询中的列，并确保内存中这些列有数据加载。在该操作期间，如果达到内存的阈值，为了给查询需要的列提供内存空间，最少使用加载级别的文件将会被删除。

3.5.2 CarbonData 性能调优常见配置参数

操作场景

CarbonData的性能与配置参数相关，本章节提供了能够提升性能的相关配置介绍。

操作步骤

用于CarbonData查询的配置介绍，详情请参见[表3-9](#)和[表3-10](#)。

表 3-9 Shuffle 过程中，启动 Task 的个数

参数	spark.sql.shuffle.partitions
所属配置文件	spark-defaults.conf
适用于	数据查询
场景描述	Spark shuffle时启动的Task个数。
如何调优	一般建议将该参数值设置为执行器核数的1到2倍。例如，在聚合场景中，将task个数从200减少到32，有些查询的性能可提升2倍。

表 3-10 设置用于 CarbonData 查询的 Executor 个数、CPU 核数以及内存大小

参数	spark.executor.cores spark.executor.instances spark.executor.memory
所属配置文件	spark-defaults.conf
适用于	数据查询
场景描述	设置用于CarbonData查询的Executor个数、CPU核数以及内存大小。
如何调优	在银行方案中，为每个执行器提供4个CPU内核和15GB内存，可以获得良好的性能。这2个值并不意味着越多越好，在资源有限的情况下，需要正确配置。例如，在银行方案中，每个节点有足够的32个CPU核，而只有64GB的内存，这个内存是不够的。例如，当每个执行器有4个内核和12GB内存，有时在查询期间发生垃圾收集（GC），会导致查询时间从3秒增加到超过15秒。在这种情况下需要增加内存或减少CPU内核。

用于CarbonData数据加载的配置参数，详情请参见[表3-11](#)、[表3-12](#)和[表3-13](#)。

表 3-11 设置数据加载使用的 CPU core 数量

参数	carbon.number.of.cores.while.loading
所属配置文件	carbon.properties
适用于	数据加载
场景描述	数据加载过程中，设置处理数据使用的CPU core数量。
如何调优	如果有更多的CPU个数，那么可以增加CPU值来提高性能。例如，将该参数值从2增加到4，那么CSV文件读取性能可以增加大约1倍。

表 3-12 是否使用 YARN 本地目录进行多磁盘数据加载

参数	carbon.use.local.dir
所属配置文件	carbon.properties
适用于	数据加载
场景描述	是否使用YARN本地目录进行多磁盘数据加载。
如何调优	如果将该参数值设置为“true”，CarbonData将使用YARN本地目录进行多表加载磁盘负载均衡，以提高数据加载性能。

表 3-13 加载时是否使用多路径

参数	carbon.use.multiple.temp.dir
所属配置文件	carbon.properties
适用于	数据加载
场景描述	是否使用多个临时目录存储sort临时文件。
如何调优	设置为true，则数据加载时使用多个临时目录存储sort临时文件。此配置能提高数据加载性能并避免磁盘单点故障。

用于CarbonData数据加载和数据查询的配置参数，详情请参见[表3-14](#)。

表 3-14 设置数据加载和查询使用的 CPU core 数量

参数	carbon.compaction.level.threshold
所属配置文件	carbon.properties
适用于	数据加载和查询

场景描述	对于minor压缩，在阶段1中要合并的segment数量和阶段2中要合并的已压缩的segment数量。
如何调优	每次CarbonData加载创建一个segment，如果每次加载的数据量较小，将在一段时间内生成许多小文件，影响查询性能。配置该参数将小的segment合并为一个大的segment，然后对数据进行排序，可提高查询性能。 压缩的策略根据实际的数据大小和可用资源决定。如某银行1天加载一次数据，且加载数据选择在晚上无查询时进行，有足够的资源，压缩策略可选择为6、5。

表 3-15 使用索引缓存服务器时是否开启数据预加载

参数	carbon.indexserver.enable.prepriming
所属配置文件	carbon.properties
适用于	数据加载
场景描述	使用索引缓存服务器过程中开启数据预加载可以提升首次查询的性能。
如何调优	用户可以将该参数设置为true来开启预加载。默认情况，该参数为false。

3.5.3 创建 CarbonData Table 的建议

操作场景

本章节根据超过50个测试用例总结出建议，帮助用户创建拥有更高查询性能的CarbonData表。

表 3-16 CarbonData 表中的列

Column name	Data type	Cardinality	Attribution
msname	String	3千万	dimension
BEGIN_TIME	bigint	1万	dimension
host	String	1百万	dimension
dime_1	String	1千	dimension
dime_2	String	500	dimension
dime_3	String	800	dimension
counter_1	numeric(20,0)	NA	measure
...	...	NA	measure

Column name	Data type	Cardinality	Attribution
counter_100	numeric(20,0)	NA	measure

操作步骤

- 如果待创建的表有一个常用于过滤的列，例如80%以上的场景使用此列过滤。

针对此类场景，调优方法如下：

将常用于过滤的列放在sort_columns第一列。

例如，msname作为过滤条件在查询中使用的最多，则将其放在第一列。创建表的命令如下，其中采用msname作为过滤条件的查询性能将会很好。

```
create table carbondata_table(
  msname String,
  ...
)STORED AS carbondata TBLPROPERTIES ('SORT_COLUMNS'='msname');
```

- 如果待创建的表有多个常用于过滤的列。

针对此类场景，调优方法如下：

为常用的过滤列创建索引。

例如，如果msname，host和dime_1是过滤经常使用的列，根据cardinality，sort_columns列的顺序是dime_1-> host-> msname…。创建表命令如下，以下命令可提高dime_1，host和msname上的过滤性能。

```
create table carbondata_table(
  dime_1 String,
  host String,
  msname String,
  dime_2 String,
  dime_3 String,
  ...
)STORED AS carbondata
TBLPROPERTIES ('SORT_COLUMNS'='dime_1,host,msname');
```

- 如果每个用于过滤的列的频率相当。

针对此类场景，调优方法如下：

sort_columns按照cardinality从低到高的顺序排列。

创建表的命令如下：

```
create table carbondata_table(
  Dime_1 String,
  BEGIN_TIME bigint,
  HOST String,
  msname String,
  ...
)STORED AS carbondata
TBLPROPERTIES ('SORT_COLUMNS'='dime_2,dime_3,dime_1, BEGIN_TIME,host,msname');
```

- 按照维度的cardinality从低到高创建表后，再为高Cardinality列创建SECONDARY INDEX。创建索引的语句如下：

```
create index carbondata_table_index_msidx on tablecarbondata_table (
  msname String) as 'carbondata' PROPERTIES ('table_blocksize'='128');
create index carbondata_table_index_host on tablecarbondata_table (
  host String) as 'carbondata' PROPERTIES ('table_blocksize'='128');
```

- 对于不需要高精度的度量，无需使用numeric (20,0)数据类型，建议使用double数据类型来替换numeric (20,0)数据类型，以提高查询性能。

在一个测试用例中，使用double来替换numeric (20, 0)，查询时间从15秒降低到3秒，查询速度提高了5倍。创建表命令如下：

```
create table carbondata_table(  
  Dime_1 String,  
  BEGIN_TIME bigint,  
  HOST String,  
  msname String,  
  counter_1 double,  
  counter_2 double,  
  ...  
  counter_100 double,  
  )STORED AS carbondata  
;
```

- 如果列值总是递增的，如start_time。

例如，每天将数据加载到CarbonData，start_time是每次加载的增量。对于这种情况，建议将start_time列放在sort_columns的最后，因为总是递增的值可以始终使用最小/最大索引。创建表命令如下：

```
create table carbondata_table(  
  Dime_1 String,  
  HOST String,  
  msname String,  
  counter_1 double,  
  counter_2 double,  
  BEGIN_TIME bigint,  
  ...  
  counter_100 double,  
  )STORED AS carbondata  
  TBLPROPERTIES ('SORT_COLUMNS'='dime_2,dime_3,dime_1..BEGIN_TIME');
```

3.6 CarbonData 常见配置参数

本章节介绍CarbonData所有配置的详细信息。

carbon.properties 相关参数

根据用户实际使用场景在服务端或者客户端配置CarbonData相关参数。

- 服务端：登录FusionInsight Manager页面，选择“集群 > 服务 > Spark2x > 配置 > 全部配置 > JDBCServer（角色） > 自定义”，在参数“spark.carbon.customized.configs”中添加CarbonData相关参数配置。
- 客户端：登录客户端节点，在“{客户端安装目录}/Spark/spark/conf/carbon.properties”文件中配置相关参数。

表 3-17 carbon.properties 中的系统配置

参数	默认值	描述
carbon.ddl.base.hdfs.url	hdfs://hacluster/opt/data	此属性用于从HDFS基本路径配置HDFS相对路径，在“fs.defaultFS”中进行配置。在“carbon.ddl.base.hdfs.url”中配置的路径将被追加到在“fs.defaultFS”中配置的HDFS路径中。如果配置了这个路径，则用户不需要通过完整路径加载数据。 例如：如果CSV文件的绝对路径是“hdfs://10.18.101.155:54310/data/cnbc/2016/xyz.csv”，其中，路径“hdfs://10.18.101.155:54310”来源于属性“fs.defaultFS”并且用户可以把“/data/cnbc/”作为“carbon.ddl.base.hdfs.url”配置。 当前，在数据加载时，用户可以指定CSV文件为“/2016/xyz.csv”。
carbon.badRecords.location	-	指定Bad records的存储路径。此路径为HDFS路径。默认值为Null。如果启用了bad records日志记录或者bad records操作重定向，则该路径必须由用户进行配置。
carbon.badRecords.action	fail	以下是bad records的四种行为类型： FORCE：通过将bad records存储为NULL来自动更正数据。 REDIRECT：Bad records被写入carbon.badRecords.location配置路径下的CSV文件而不是被加载。 IGNORE：Bad records既不被加载也不被写入CSV文件。 FAIL：如果找到任何bad records，则数据加载失败。
carbon.update.sync.folder	/tmp/carbondata	modifiedTime.mdt文件路径，可以设置为已有路径或新路径。 说明 如果设置为已有路径，需确保所有用户都可以访问该路径，且该路径具有777权限。
carbon.enable.badrecord.action.redirect	false	是否在数据加载中开启redirect方式来处理bad records。启用该配置后，源文件中的bad records会被记录在指定存储位置生成的CSV文件中。在Windows操作系统中打开此类CSV文件时，可能会发生CSV注入。

表 3-18 carbon.properties 中的性能配置

参数	默认值	描述
数据加载配置		
carbon.sort.file.write.buffer.size	16384	为了限制内存的使用，CarbonData会将数据排序并写入临时文件中。该参数控制读取和写入临时文件过程使用的缓存大小。单位：字节。 取值范围为：10240~10485760。
carbon.graph.rowset.size	100000	数据加载图步骤之间交换的行集大小。 最小值=500，最大值=1000000
carbon.number.of.cores.while.loading	6	数据加载时所使用的核数。配置的核数越大压缩性能越好。如果CPU资源充足可以增加此值。
carbon.sort.size	500000	内存排序的数据大小。
carbon.enableXXHash	true	用于hashkey计算的hashmap算法。
carbon.number.of.cores.block.sort	7	数据加载时块排序所使用的核数。
carbon.max.driver.lru.cache.size	-1	在driver端加载数据所达到的最大LRU缓存大小。以MB为单位，默认值为-1，表示缓存没有内存限制。只允许使用大于0的整数值。
carbon.max.executor.lru.cache.size	-1	在executor端加载数据所达到的最大LRU缓存大小。以MB为单位，默认值为-1，表示缓存没有内存限制。只允许使用大于0的整数值。如果未配置该参数，则将考虑参数“carbon.max.driver.lru.cache.size”的值。
carbon.merge.sort.prefetch	true	在数据加载过程中，从排序的临时文件中读取数据进行合并排序时，启用数据预取。
carbon.update.persist.enable	true	启用此参数将考虑持久化数据，减少UPDATE操作的执行时间。
enable.unsafe.sort	true	指定在数据加载期间是否使用非安全排序。非安全的排序减少了数据加载操作期间的垃圾回收（GC），从而提高了性能。默认值为“true”，表示启用非安全排序功能。
enable.offheap.sort	true	在数据加载期间启用堆排序。
offheap.sort.chunk.size.in.mb	64	指定需要用于排序的数据块的大小。最小值为1MB，最大值为1024MB。

参数	默认值	描述
carbon.unsafe.working.memory.in.mb	512	<p>指定非安全工作内存的大小。这将用于排序数据，存储列页面等。单位是MB。</p> <p>数据加载所需内存： (“ carbon.number.of.cores.while.loading” 的值[默认值 = 6]) x 并行加载数据的表格 x (“ offheap.sort.chunk.size.inmb” 的值[默认值 = 64 MB] + “ carbon.blockletgroup.size.in.mb” 的值[默认值 = 64 MB] + 当前的压缩率[64 MB/3.5]) = ~900 MB 每表格</p> <p>数据查询所需内存： (SPARK_EXECUTOR_INSTANCES. [默认值 = 2]) x (carbon.blockletgroup.size.in.mb [默认值 = 64 MB] + “ carbon.blockletgroup.size.in.mb” 解压内容[默认值 = 64 MB * 3.5]) x (每个执行器核数[默认值 = 1]) = ~ 600 MB</p>
carbon.sort.inmemory.storage.size.in.mb	512	<p>指定要存储在内存中的中间排序数据的大小。达到该指定的值，系统会将数据写入磁盘。单位是MB。</p>
sort.inmemory.size.inmb	1024	<p>指定要保存在内存中的中间排序数据的大小。达到该指定值后，系统会将数据写入磁盘。单位：MB。</p> <p>如果配置了 “ carbon.unsafe.working.memory.in.mb” 和 “ carbon.sort.inmemory.storage.size.in.mb”，则不需要配置该参数。如果此时也配置了该参数，那么这个内存的20%将用于工作内存 “ carbon.unsafe.working.memory.in.mb”，80%将用于排序存储内存 “ carbon.sort.inmemory.storage.size.in.mb”。</p> <p>说明 Spark配置参数 “ spark.yarn.executor.memoryOverhead” 的值应该大于CarbonData配置参数 “ sort.inmemory.size.inmb” 的值，否则如果堆外（off heap）访问超出配置的executor内存，则YARN可能会停止 executor。</p>
carbon.blockletgroup.size.in.mb	64	<p>数据作为blocklet group被系统读入。该参数指定blocklet group的大小。较高的值会有更好的顺序IO访问性能。</p> <p>最小值为16MB，任何小于16MB的值都将重置为默认值（64MB）。</p> <p>单位：MB。</p>
enable.inmemory.merge.sort	false	<p>指定是否启用内存合并排序（inmemorymerge sort）。</p>

参数	默认值	描述
use.offheap.in.query.processing	true	指定是否在查询处理中启用offheap。
carbon.load.sort.scope	local_sort	指定加载操作的排序范围。支持两种类型的排序，batch_sort和local_sort。选择batch_sort将提升加载性能，但会降低查询性能。 说明 local_sort与分区表的DDL操作存在冲突，不能同时使用，且对分区表性能提升不明显，不建议在分区表上启用该特性。
carbon.batch.sort.size.inmb	-	指定在数据加载期间为批处理排序而考虑的数据大小。推荐值为小于总排序数据的45%。该值以MB为单位。 说明 如果没有设置参数值，那么默认情况下其大约等于“sort.inmemory.size.inmb”参数值的45%。
enable.unsafe.columnpage	true	指定在数据加载或查询期间，是否将页面数据保留在堆内存中，以避免垃圾回收受阻。
carbon.use.local.dir	false	是否使用YARN本地目录加载多个磁盘的数据。设置为true，则使用YARN本地目录加载多个磁盘的数据，以提高数据加载性能。
carbon.use.multiple.temp.dir	false	是否使用多个临时目录存储临时文件以提高数据加载性能。
carbon.load.datamaps.parallel.db_name.table_name	NA	值为true或者false。可以设置数据库名和表名，使得该表的首次查询性能得到提升。
压缩配置		
carbon.number.of.cores.while.compacting	2	在压缩过程中用于写入数据所使用的核数。配置的核数越大压缩性能越好。如果CPU资源充足可以增加此值。
carbon.compaction.level.threshold	4,3	该属性用于Minor压缩，决定合并segment的数量。例如：如果被设置为“2,3”，则将每2个segment触发一次Minor压缩。“3”是Level 1压缩的segment个数，这些segment将进一步被压缩为新的segment。有效值为0-100。
carbon.major.compaction.size	1024	使用该参数配置Major压缩的大小。总数低于该阈值的segment将被合并。 单位为MB。

参数	默认值	描述
carbon.horizontal.compaction.enable	true	该参数用于配置打开/关闭水平压缩。在每个DELETE和UPDATE语句之后，如果增量（DELETE / UPDATE）文件超过指定的阈值，则可能发生水平压缩。默认情况下，该参数值设置为“true”，打开水平压缩功能，可将参数值设置为“false”来关闭水平压缩功能。
carbon.horizontal.update.compaction.threshold	1	该参数指定segment内的UPDATE增量文件数的阈值限制。在增量文件数量超过阈值的情况下，segment内的UPDATE增量文件变得适合水平压缩，并压缩为单个UPDATE增量文件。默认情况下，该参数值设置为1。可以设置为1到10000之间的值。
carbon.horizontal.delete.compaction.threshold	1	该参数指定segment的block中的DELETE增量文件数量的阈值限制。在增量文件数量超过阈值的情况下，segment特定block的DELETE增量文件变得适合水平压缩，并压缩为单个DELETE增量文件。默认情况下，该参数值设置为1。可以设置为1到10000之间的值。
查询配置		
carbon.number.of.cores	4	查询时所使用的核数。
carbon.limit.block.distribution.enable	false	当查询语句中包含关键字limit时，启用或禁用CarbonData块分布。默认值为“false”，将对包含关键字limit的查询语句禁用块分布。此参数调优请参考 CarbonData性能调优常见配置参数 。
carbon.custom.block.distribution	false	指定是使用Spark还是CarbonData的块分配功能。默认情况下，其配置值为“false”，表明启用Spark块分配。若要使用CarbonData块分配，请将配置值更改为“true”。
carbon.infilter.subquery.pushdown.enable	false	如果启用此参数，并且用户在具有subquery的过滤器中触发Select查询，则执行子查询，并将输出作为IN过滤器广播到左表，否则将执行SortMergeSemiJoin。建议在IN过滤器子查询未返回太多记录时启用此参数。例如，IN子查询返回10k或更少的记录时，启用此参数将更快地给出查询结果。 示例： <i>select * from flow_carbon_256b where cus_no in (select cus_no from flow_carbon_256b where dt>='20260101' and dt<='20260701' and txn_bk='tk_1' and txn_br='tr_1') limit 1000;</i>
carbon.scheduler.minRegisteredResourcesRatio	0.8	启动块分布所需的最小资源（executor）比率。默认值为“0.8”，表示所请求资源的80%被分配用于启动块分布。
carbon.dynamicAllocation.schedulerTimeout	5	此参数值指示调度器等待executors处于活动状态的最长时间。默认值为“5”秒，允许的最大值为“15”秒。

参数	默认值	描述
enable.unsafe.in.query.processing	true	指定在查询操作期间是否使用非安全排序。非安全排序减少查询操作期间的垃圾回收（GC），从而提高性能。默认值为“true”，表示启用非安全排序功能。
carbon.enable.vector.reader	true	为结果收集（result collection）启用向量处理，以增强查询性能。
carbon.query.show.datamaps	true	SHOW TABLES 会展示所有的表包含主表和datamap。如果需要过滤掉datamap，将该配置设置为false。
二级索引配置		
carbon.secondary.index.creation.threads	1	该参数用于配置启动二级索引创建期间并行处理segments的线程数。当表的segments数较多时，该参数有助于微调系统生成二级索引的速度。该参数值范围为1到50。
carbon.silookup.partialstring	true	<ul style="list-style-type: none"> 当配置为true时，它包括开始，结尾和包含。 当配置为false时，它只包括从二级索引开始。
carbon.sisegment.merge	true	<p>开启这个配置后会合并二级索引表segment内的carbondata文件。合并发生在导入操作后，在二级索引表导入操作的最后，会检查并合并小文件。</p> <p>说明 Table Block Size会用作合并小文件的大小阈值。</p>

表 3-19 carbon.properties 中的其它配置

参数	默认值	描述
数据加载配置		
carbon.lock.type	HDFSLOCK	<p>该配置指定了表上并发操作过程中所要求的锁的类型。</p> <p>有以下几种类型锁实现方式：</p> <ul style="list-style-type: none"> LOCALLOCK：基于本地文件系统的文件来创建的锁。该锁只适用于一台机器上只运行一个Spark Driver（或者JDBCServer）的情况。 HDFSLOCK：基于HDFS文件系统上的文件来创建的锁。该锁适用于集群上有多个运行的Spark应用而且没有可用的ZooKeeper的情况。
carbon.sort.intermediate.files.limit	20	中间文件的最小数量。生成中间文件后开始排序合并。此参数调优请参考 CarbonData性能调优常见配置参数 。

参数	默认值	描述
carbon.csv.read.buffer.size.byte	1048576	CSV读缓冲区大小。
carbon.merge.sort.reader.thread	3	用于读取中间文件进行最终合并的最大线程数。
carbon.concurrent.lock.retries	100	指定获取并发操作锁的最大重试次数。该参数用于并发加载。
carbon.concurrent.lock.retry.timeout.sec	1	指定获取并发操作的锁重试之间的间隔。
carbon.lock.retries	3	指定除导入操作外其他所有操作尝试获取锁的次数。
carbon.lock.retry.timeout.sec	5	指定除导入操作外其他所有操作尝试获取锁的时间间隔。
carbon.tempstore.location	/opt/Carbon/TempStoreLocation	临时存储位置。默认情况下，采用“System.getProperty("java.io.tmpdir")”方法获取。此参数调优请参考 CarbonData性能调优常见配置参数 中关于“carbon.use.local.dir”的描述。
carbon.load.log.counter	500000	数据加载记录计数日志。
SERIALIZATION_NULL_FORMAT	\N	指定需要替换为NULL的值。
carbon.skip.empty.line	false	设置此属性将在数据加载期间忽略CSV文件中的空行。
carbon.load.datamaps.parallel	false	该配置项将会开启对所有会话所有表的datamap并行加载。该配置项通过将导入datamap到内存的工作分发给所有的executor来缩短时间，进而提升查询性能。
合并配置		
carbon.numberof.preserve.segments	0	若用户希望从被合并的segment中保留一定数量的segment，可设置该属性参数。 例如：“carbon.numberof.preserve.segments” = “2”，那么合并的segment中将不包含最新的2个segment。 默认保留No segment的状态。

参数	默认值	描述
carbon.allow.ed.compaction.days	0	合并将合并并在配置的指定天数中加载的 segment。 例如：如果配置值为“2”，那么只有在2天时间框架中加载的segment被合并。2天以外被加载的segment不会被合并。 该参数默认为禁用。
carbon.enable.auto.load.merge	false	在数据加载时启用压缩。
carbon.merge.index.in.segment	true	如果设置，则Segment内的所有Carbon索引文件（.carbonindex）将合并为单个Carbon索引合并文件（.carbonindexmerge）。这增强了首次查询性能
查询配置		
max.query.execution.time	60	单次查询允许的最大时间。 单位为分钟。
carbon.enableMinMax	true	MinMax用于提高查询性能。设置为false可禁用该功能。
carbon.lease.recovery.retry.count	5	需要为恢复文件租约所需的最大尝试次数。 最小值：1 最大值：50
carbon.lease.recovery.retry.interval	1000 (ms)	尝试在文件上进行租约恢复之后的间隔（Interval）或暂停（Pause）时间。 最小值：1000（ms） 最大值：10000（ms）

spark-defaults.conf 相关参数

- 登录客户端节点，在“{客户端安装目录}/Spark/spark/conf/spark-defaults.conf”文件中配置表3-20相关参数。

表 3-20 spark-defaults.conf 中的 Spark 配置参考

参数	默认值	描述
spark.driver.memory	4G	指定用于driver端进程的内存，其中 SparkContext已初始化。 说明 在客户端模式下，不要使用SparkConf在应用程序中设置该参数，因为驱动程序JVM已经启动。要配置该参数，请在--driver-memory命令行选项或默认属性文件中进行配置。

参数	默认值	描述
spark.executor.memory	4GB	指定每个执行程序进程使用的内存。
spark.sql.crossJoin.enabled	true	如果查询包含交叉连接，请启用此属性，以便不会发生错误，此时使用交叉连接而不是连接，可实现更好的性能。

- 在Spark Driver端的“spark-defaults.conf”文件中配置以下参数。
 - 在spark-sql模式下配置：登录Spark客户端节点，在“{客户端安装目录}/Spark/spark/conf/spark-defaults.conf”文件中配置表3-21相关参数。

表 3-21 spark-sql 模式下的配置参数

参数	配置值	描述
spark.driver.extraJavaOptions	- Dlog4j.configuration=file:/opt/client/Spark2x/spark/conf/log4j.properties - Djetty.version=x.y.z - Dzookeeper.server.principal=zookeeper/hadoop.<系统域名> - Djava.security.krb5.conf=/opt/client/KrbClient/kerberos/var/krb5kdc/krb5.conf - Djava.security.auth.login.config=/opt/client/Spark2x/spark/conf/jaas.conf - Dorg.xerial.snappy.tmpdir=/opt/client/Spark2x/tmp - Dcarbon.properties.filepath=/opt/client/Spark2x/spark/conf/carbon.properties - Djava.io.tmpdir=/opt/client/Spark2x/tmp	默认值中“/opt/client/Spark2x/spark”为客户端的CLIENT_HOME，且该默认值是追加到参数“spark.driver.extraJavaOptions”其他值之后的，此参数用于指定Driver端的“carbon.properties”文件路径。 说明 请注意“=”两边不要有空格。
spark.sql.session.state.builder	org.apache.spark.sql.hive.FIHiveACLSessionStateBuilder	指定会话状态构造器。
spark.carbon.sqlastbuilder.class.name	org.apache.spark.sql.hive.CarbonInternalSqlAstBuilder	指定AST构造器。
spark.sql.catalog.classes	org.apache.spark.sql.hive.HiveACLExternalCatalog	指定Hive的外部目录实现。启用Spark ACL时必须提供。

参数	配置值	描述
spark.sql.hive.implementation	org.apache.spark.sql.hive.HiveACLClientImpl	指定Hive客户端调用的实现。启用Spark ACL时必须提供。
spark.sql.hiveClient.isolation.enabled	false	启用Spark ACL时必须提供。

- 在JDBCServer服务中配置：登录JDBCServer安装节点，在“{BIGDATA_HOME}/FusionInsight_Spark_*/*_JDBCServer/etc/spark-defaults.conf”文件中配置[表3-22](#)相关参数。

表 3-22 JDBCServer 服务中的配置参数

参数	配置值	描述
spark.driver.extraJavaOptions	-Xloggc:\${SPARK_LOG_DIR}/indexserver-omm-%p-gc.log - XX:+PrintGCDetails -XX:-OmitStackTracenFastThrow - XX:+PrintGCTimeStamps - XX:+PrintGCDateStamps - XX:MaxDirectMemorySize=512M - XX:MaxMetaspaceSize=512M - XX:+UseGCLogFileRotation - XX:NumberOfGCLogFiles=20 - XX:GCLogFileSize=10M - XX:OnOutOfMemoryError='kill -9 %p' - Djetty.version=x.y.z - Dorg.xerial.snappy.tmpdir=\${BIGDATA_HOME}/tmp/spark2x/JDBCServer/snappy_tmp - Djava.io.tmpdir=\${BIGDATA_HOME}/tmp/spark2x/JDBCServer/io_tmp - Dcarbon.properties.filepath=\${SPARK_CONF_DIR}/carbon.properties -	默认值中\${SPARK_CONF_DIR}需视具体的集群而定，且该默认值是追加到参数“spark.driver.extraJavaOptions”其他值之后的，此参数用于指定Driver端的“carbon.properties”文件路径。 说明 请注意“=”两边不要有空格。

参数	配置值	描述
	Djdk.tls.ephemeralDHKeySize=2048 - Dspark.ssl.keyStore=\${SPARK_CONF_DIR}/child.keystore #{java_stack_prefer}	
spark.sql.session.state.builder	org.apache.spark.sql.hive.FIHiveACLSessionStateBuilder	指定会话状态构造器。
spark.carbon.sqlastbuilder.classname	org.apache.spark.sql.hive.CarbonInternalSqlAstBuilder	指定AST构造器。
spark.sql.catalog.classes	org.apache.spark.sql.hive.HiveACLExternalCatalog	指定Hive的外部目录实现。启用Spark ACL时必须提供。
spark.sql.hive.implementation	org.apache.spark.sql.hive.HiveACLClientImpl	指定Hive客户端调用的实现。启用Spark ACL时必须提供。
spark.sql.hiveClient.isolation.enabled	false	启用Spark ACL时必须提供。

3.7 CarbonData 语法参考

3.7.1 DDL

3.7.1.1 CREATE TABLE

命令功能

CREATE TABLE命令通过指定带有表属性的字段列表来创建CarbonData Table。

命令格式

CREATE TABLE *[IF NOT EXISTS]* *[db_name.]table_name*

[(col_name data_type, ...)]

STORED AS *carbodata*

[TBLPROPERTIES (property_name=property_value, ...)];

所有表的附加属性都会放到TBLPROPERTIES中来定义。

参数描述

表 3-23 CREATE TABLE 参数描述

参数	描述
db_name	Database名称，由字母、数字和下划线（_）组成。
col_name data_type	以逗号分隔的带数据类型的列表。列名由字母、数字和下划线（_）组成。 说明 在CarbonData表创建过程中，不允许使用tupleId，PositionId和PositionReference为列命名，因为具有这些名称的列由二级索引命令在内部使用。
table_name	Database中的表名，由字母、数字和下划线（_）组成。
STORED AS	参数carbodata，定义和创建CarbonData table。
TBLPROPERTIES	CarbonData table属性列表。

注意事项

以下是表格属性的使用。

- **Block大小**

单个表的数据文件block大小可以通过TBLPROPERTIES进行定义，系统会选择数据文件实际大小和设置的blocksize大小中的较大值，作为该数据文件在HDFS上存储的实际blocksize大小。单位为MB，默认值为1024MB，范围为1MB~2048MB。若设置值不在[1, 2048]之间，系统将会报错。

一旦block大小达到配置值，写入程序将启动新的CarbonData数据的block。数据以页面大小（32000个记录）的倍数写入，因此边界在字节级别上不严格。如果新页面跨越配置block的边界，则不会将其写入当前block，而是写入新的block。

TBLPROPERTIES('table_blocksize'='128')

说明

- 当在CarbonData表中配置了较小的blocksize，而加载的数据生成的数据文件比较大时，在HDFS上显示的blocksize会与设置值不同。这是因为，对于每一个本地block文件的首次写入，即使待写入数据的大小大于blocksize的配置值，也直接将待写入数据写入此block。所以，HDFS上blocksize的实际值为待写入数据大小与blocksize配置值中的较大值。
- 当CarbonData表中的数据文件block.num小于任务并行度（parallelism）时，CarbonData数据文件的block会被切为新的block，使得blocks.num大于parallelism，这样所有core均可被使用。这种优化称为block distribution。
- **SORT_SCOPE**：指定表创建时的排序范围。如下为四种排序范围。

- GLOBAL_SORT: 它提高了查询性能，特别是点查询。
TBLPROPERTIES('SORT_SCOPE'='GLOBAL_SORT')
- LOCAL_SORT: 数据会本地排序（任务级别排序）。

📖 说明

LOCAL_SORT与分区表的DDL操作存在冲突，不能同时使用，且对分区表性能提升不明显，不建议在分区表上启用该特性。

- NO_SORT: 默认排序。它将以不排序的方式加载数据，这将显著提升加载性能。
- SORT_COLUMNS
此表属性指定排序列的顺序。

TBLPROPERTIES('SORT_COLUMNS'='column1, column3')

📖 说明

- 如果未指定此属性，则默认情况下，没有列会被排序。
- 如果指定了此属性，但具有空参数，则表将被加载而不进行排序。例如，*('SORT_COLUMNS='')*。
- SORT_COLUMNS将接受string, date, timestamp, short, int, long, byte和boolean数据类型。
- RANGE_COLUMN
此表属性指定一列，该列将会按照一个范围值来对输入的数据进行分区。仅可配置一列。在数据导入过程中，可以使用“global_sort_partitions”或者“scale_factor”来避免生成小文件。

TBLPROPERTIES('RANGE_COLUMN'='column1')

- LONG_STRING_COLUMNS
普通String类型的长度不能超过32000字符，如果需要存储超过32000字符的字符串，指定LONG_STRING_COLUMNS配置为该列。

TBLPROPERTIES('LONG_STRING_COLUMNS'='column1, column3')

📖 说明

LONG_STRING_COLUMNS仅可以设置string/char/varchar类型的列，并且不能为SORT_COLUMNS和复杂列。

使用场景

通过指定列创建表

CREATE TABLE命令与Hive DDL相同。CarbonData的额外配置将作为表格属性给出。

CREATE TABLE *[IF NOT EXISTS] [db_name.]table_name*

[(col_name data_type, ...)]

STORED AS *carbondata*

[TBLPROPERTIES (property_name=property_value, ...)];

示例

CREATE TABLE IF NOT EXISTS *productdb.productSalesTable (*

```
productNumber Int,  
productName String,  
storeCity String,  
storeProvince String,  
productCategory String,  
productBatch String,  
saleQuantity Int,  
revenue Int)  
  
STORED AS carbondata  
  
TBLPROPERTIES (  
  'table_blocksize'='128',  
  'SORT_COLUMNS'='productBatch, productName')
```

系统响应

Table创建成功，创建成功的消息将被记录在系统日志中。

3.7.1.2 CREATE TABLE As SELECT

命令功能

CREATE TABLE As SELECT命令通过指定带有表属性的字段列表来创建CarbonData Table。

命令格式

```
CREATE TABLE [IF NOT EXISTS] [db_name.]table_name STORED AS carbondata  
[TBLPROPERTIES (key1=val1, key2=val2, ...)] AS select_statement;
```

参数描述

表 3-24 CREATE TABLE 参数描述

参数	描述
db_name	Database名称，由字母、数字和下划线（_）组成。
table_name	Database中的表名，由字母、数字和下划线（_）组成。
STORED AS	使用CarbonData数据格式存储数据。
TBLPROPERTIES	CarbonData table属性列表。详细信息，见 注意事项 。

注意事项

NA

示例

```
CREATE TABLE ctas_select_parquet STORED AS carbondata as select * from  
parquet_ctas_test;
```

系统响应

该命令会从Parquet表上创建一个Carbon表，同时导入所有Parquet表的数据。

3.7.1.3 DROP TABLE

命令功能

DROP TABLE的功能是用来删除已存在的Table。

命令格式

```
DROP TABLE [IF EXISTS] [db_name.]table_name;
```

参数描述

表 3-25 DROP TABLE 参数描述

参数	描述
db_name	Database名称。如果未指定，将选择当前database。
table_name	需要删除的Table名称。

注意事项

在该命令中，IF EXISTS和db_name是可选配置。

示例

```
DROP TABLE IF EXISTS productDatabase.productSalesTable;
```

系统响应

Table将被删除。

3.7.1.4 SHOW TABLES

命令功能

SHOW TABLES命令用于显示所有在当前database中的table，或所有指定database的table。

命令格式

```
SHOW TABLES [IN db_name];
```

参数描述

表 3-26 SHOW TABLES 参数描述

参数	描述
IN db_name	Database名称，仅当需要显示指定Database的所有Table时配置。

注意事项

IN db_Name为可选配置。

示例

```
SHOW TABLES IN ProductDatabase;
```

系统响应

显示所有Table。

3.7.1.5 ALTER TABLE COMPACTION

命令功能

ALTER TABLE COMPACTION命令将合并指定数量的segment为一个segment。这将提高该表的查询性能。

命令格式

```
ALTER TABLE [db_name.]table_name COMPACT 'MINOR/MAJOR/  
SEGMENT_INDEX';
```

```
ALTER TABLE [db_name.]table_name COMPACT 'CUSTOM' WHERE SEGMENT.ID IN  
(id1, id2, ...);
```

参数描述

表 3-27 ALTER TABLE COMPACTION 参数描述

Parameter	Description
db_name	数据库名。若未指定，则选择当前数据库。
table_name	表名。
MINOR	Minor合并，详见 合并Segments 。

Parameter	Description
MAJOR	Major合并，详见 合并Segments 。
SEGMENT_INDEX	这会将一个segment内的所有Carbon索引文件（.carbonindex）合并为一个Carbon索引合并文件（.carbonindexmerge）。这增强了首次查询性能。详见 表3-7 。
CUSTOM	Custom合并，详见 合并Segments 。

注意事项

NA

示例

```
ALTER TABLE ProductDatabase COMPACT 'MINOR';
```

```
ALTER TABLE ProductDatabase COMPACT 'MAJOR';
```

```
ALTER TABLE ProductDatabase COMPACT 'SEGMENT_INDEX';
```

```
ALTER TABLE ProductDatabase COMPACT 'CUSTOM' WHERE SEGMENT.ID IN (0, 1);
```

系统响应

由于为后台运行，**ALTER TABLE COMPACTION**命令不会显示压缩响应。

如果想要查看MINOR合并和MAJOR合并的响应结果，用户可以检查日志或运行**SHOW SEGMENTS**命令查看。

示例：

```
+-----+-----+-----+-----+-----+-----+-----+-----+
+--+
| ID | Status | Load Start Time | Load Time Taken | Partition | Data Size | Index Size | File
Format |
+-----+-----+-----+-----+-----+-----+-----+-----+
+--+
| 3 | Success | 2020-09-28 22:53:26.336 | 3.726S | {} | 6.47KB | 3.30KB | columnar_v3 |
| 2 | Success | 2020-09-28 22:53:01.702 | 6.688S | {} | 6.47KB | 3.30KB | columnar_v3 |
| 1 | Compacted | 2020-09-28 22:51:15.242 | 5.82S | {} | 6.50KB | 3.43KB |
columnar_v3 |
| 0.1 | Success | 2020-10-30 20:49:24.561 | 16.66S | {} | 12.87KB | 6.91KB | columnar_v3
|
| 0 | Compacted | 2020-09-28 22:51:02.6 | 6.819S | {} | 6.50KB | 3.43KB | columnar_v3
|
+-----+-----+-----+-----+-----+-----+-----+-----+
+--+
```

其中，

- Compacted表示该数据已被合并。
- 0.1表示segment0与segment1合并之后的结果。

数据合并前后的其他操作没有差别。

被合并的segments（例如segment0和segment1）即成为无用的segments，会占用空间，因此建议合并之后使用***CLEAN FILES***命令进行彻底删除，再进行其他操作。**CLEAN FILES**命令的使用方法可参考**CLEAN FILES**。

3.7.1.6 TABLE RENAME

命令功能

RENAME命令用于重命名现有表。

命令语法

```
ALTER TABLE [db_name].table_name RENAME TO new_table_name;
```

参数描述

表 3-28 RENAME 参数描述

参数	描述
<i>db_name</i>	数据库名。若未指定，则选择当前数据库。
<i>table_name</i>	现有表名。
<i>new_table_name</i>	现有表名的新表名。

注意事项

- 并行运行的查询（需要使用表名获取路径，以读取CarbonData存储文件）可能会在此操作期间失败。
- 不允许二级索引表重命名。

示例

```
ALTER TABLE carbon RENAME TO carbondata;
```

```
ALTER TABLE test_db.carbon RENAME TO test_db.carbondata;
```

系统响应

CarbonData库中的文件夹将显示新表名称，可以通过运行SHOW TABLES显示新表名称。

3.7.1.7 ADD COLUMNS

命令功能

ADD COLUMNS命令用于为现有表添加新列。

命令语法

```
ALTER TABLE [db_name.]table_name ADD COLUMNS (col_name data_type,...)  
TBLPROPERTIES ('COLUMNPROPERTIES.columnName.shared_column'='sharedFolder.sharedColumnName,...', 'DEFAULT.VALUE.COLUMN_NAME'='default_value');
```

参数描述

表 3-29 ADD COLUMNS 参数描述

参数	描述
db_name	数据库名。若未指定，则选择当前数据库。
table_name	表名。
col_name data_type	带数据类型且用逗号分隔的列的名称。列名称包含字母，数字和下划线（_）。 说明 创建CarbonData表时，不要将列名命名为tupleId，PositionId和PositionReference，因为将在UPDATE，DELETE和二级索引命令内部使用这些名称。

注意事项

- 除了shared_column和default_value之外，将不会读取其他属性。如果指定了任何其他属性名称，则不会发生错误，其他属性将被忽略。
- 如果未指定默认值，则新列的默认值将被视为null。
- 如果在该列上应用filter，则在排序期间不会考虑新增列，新增列可能会影响查询性能。

示例

- ALTER TABLE carbon ADD COLUMNS (a1 INT, b1 STRING);**
- ALTER TABLE carbon ADD COLUMNS (a1 INT, b1 STRING)
TBLPROPERTIES ('COLUMNPROPERTIES.b1.shared_column'='sharedFolder.b1');**
- ALTER TABLE carbon ADD COLUMNS (a1 INT, b1 STRING)
TBLPROPERTIES ('DEFAULT.VALUE.a1'='10');**

系统响应

通过运行DESCRIBE命令，可显示新添加的列。

3.7.1.8 DROP COLUMNS

命令功能

DROP COLUMNS命令用于删除表中现有的列或多个列。

命令语法

```
ALTER TABLE [db_name.]table_name DROP COLUMNS (col_name, ...);
```

参数描述

表 3-30 DROP COLUMNS 参数描述

参数	描述
db_name	数据库名。若未指定，则选择当前数据库。
table_name	表名。
col_name	表中的列名称。支持多列。列名称包含字母，数字和下划线（_）。

注意事项

对于删除列操作，至少要有有一个key列在删除操作后存在于schema中，否则将显示出错信息，删除列操作将失败。

示例

假设表包含4个列，分别命名为a1，b1，c1和d1。

- 删除单个列：

```
ALTER TABLE carbon DROP COLUMNS (b1);
```

```
ALTER TABLE test_db.carbon DROP COLUMNS (b1);
```
- 删除多个列：

```
ALTER TABLE carbon DROP COLUMNS (b1,c1);
```

```
ALTER TABLE test_db.carbon DROP COLUMNS (b1,c1);
```

系统响应

运行DESCRIBE命令，将不会显示已删除的列。

3.7.1.9 CHANGE DATA TYPE

命令功能

CHANGE命令用于将数据类型从INT更改为BIGINT或将Decimal精度从低精度改为高精度。

命令语法

```
ALTER TABLE [db_name.]table_name CHANGE col_name col_name  
changed_column_type;
```

参数描述

表 3-31 CHANGE DATA TYPE 参数描述

参数	描述
db_name	数据库名。若未指定，则选择当前数据库。
table_name	表名。
col_name	表中的列名称。列名称包含字母，数字和下划线（_）。
changed_column_type	所要更改为的新数据类型。

注意事项

- 仅在没有数据丢失的情况下支持将Decimal数据类型从较低精度更改为较高精度
例如：
 - 无效场景：将Decimal数据精度从（10,2）更改为（10,5）无效，因为在这种情况下，只有scale增加，但总位数保持不变。
 - 有效场景：将Decimal数据精度从（10,2）更改为（12,3）有效，因为总位数增加2，但是scale仅增加1，这不会导致任何数据丢失。
- 将Decimal数据类型从较低精度更改为较高精度，其允许的最大精度（precision,scale）范围为（38,38），并且只适用于不会导致数据丢失的有效提升精度的场景。

示例

- 将列a1的数据类型从INT更改为BIGINT。
ALTER TABLE test_db.carbon CHANGE a1 a1 BIGINT;
- 将列a1的精度从10更改为18。
ALTER TABLE test_db.carbon CHANGE a1 a1 DECIMAL(18,2);

系统响应

通过运行DESCRIBE命令，将显示被修改列变更后的数据类型。

3.7.1.10 REFRESH TABLE

命令功能

REFRESH TABLE命令用于将已有的Carbon表数据注册到Hive元数据库中。

命令语法

REFRESH TABLE db_name.table_name;

参数描述

表 3-32 REFRESH TABLE 参数描述

参数	描述
db_name	数据库名。若未指定，则选择当前数据库。
table_name	表名。

注意事项

- 在执行此命令之前，应将旧表的表结构定义schema和数据复制到新数据库位置。
- 对于旧版本仓库，源集群和目的集群的时区应该相同。
- 新的数据库和旧数据库的名字应该相同。
- 执行命令前，旧表的表结构定义schema和数据应该复制到新的数据库位置。
- 如果表是聚合表，则应将所有聚合表复制到新的数据库位置。
- 如果旧集群使用HIVE元数据库来存储表结构，则刷新将不起作用，因为文件系统中不存在表结构定义schema文件。

示例

```
REFRESH TABLE dbcarbon.productSalesTable;
```

系统响应

通过运行该命令，已有的Carbon表数据会被注册到Hive元数据库中。

3.7.1.11 REGISTER INDEX TABLE

命令功能

REGISTER INDEX TABLE命令用于将索引表注册到主表。

命令语法

```
REGISTER INDEX TABLE indextable_name ON db_name.maintable_name;
```

参数描述

表 3-33 REFRESH INDEX TABLE 参数描述

参数	描述
db_name	数据库名。若未指定，则选择当前数据库。
indextable_name	索引表名。
maintable_name	主表名。

注意事项

在执行此命令之前，使用REFRESH TABLE将主表和二级索引表都注册到Hive元数据中。

示例

```
create database productdb;
use productdb;
CREATE TABLE productSalesTable(a int,b string,c string) stored as carbondata;
create index productNameIndexTable on table productSalesTable(c) as
'carbondata';
insert into table productSalesTable select 1,'a','aaa';
```

```
create database productdb2;
```

使用hdfs命令将productdb数据库下的productSalesTable和productNameIndexTable拷贝到productdb2。

```
refresh table productdb2.productSalesTable ;
```

```
refresh table productdb2.productNameIndexTable ;
```

```
explain select * from productdb2.productSalesTable where c = 'aaa'; //可以发现该查询命令没有使用索引表
```

```
REGISTER INDEX TABLE productNameIndexTable ON
productdb2.productSalesTable;
```

```
explain select * from productdb2.productSalesTable where c = 'aaa'; //可以发现该查询命令使用了索引表
```

系统响应

通过运行该命令，索引表会被注册到主表。

3.7.2 DML

3.7.2.1 LOAD DATA

命令功能

LOAD DATA命令以CarbonData特定的数据存储类型加载原始的用户数据，这样，CarbonData可以在查询数据时提供良好的性能。

说明

仅支持加载位于HDFS上的原始数据。

命令格式

```
LOAD DATA INPATH 'folder_path' INTO TABLE [db_name.]table_name
OPTIONS(property_name=property_value, ...);
```

参数描述

表 3-34 LOAD DATA 参数描述

参数	描述
folder_path	原始CSV数据文件夹或者文件的路径。
db_name	Database名称。若未指定，则使用当前database。
table_name	所提供的database中的表的名称。

注意事项

以下是可以在加载数据时使用的配置选项：

- DELIMITER：可以在加载命令中提供分隔符和引号字符。默认值为,。
`OPTIONS('DELIMITER'=',' , 'QUOTECHAR'='')`
可使用'DELIMITER'='\t'来表示用制表符tab对CSV数据进行分隔。
`OPTIONS('DELIMITER'='\t')`
CarbonData也支持\001和\017作为分隔符。

📖 说明

对于CSV数据，分隔符为单引号（'）时，单引号必须在双引号（" "）内。例如：
`'DELIMITER'='''`。

- QUOTECHAR：可以在加载命令中提供分隔符和引号字符。默认值为"。
`OPTIONS('DELIMITER'=',' , 'QUOTECHAR'='')`
- COMMENTCHAR：可以在加载命令中提供注释字符。在加载操作期间，如果在行的开头遇到注释字符，那么该行将被视为注释，并且不会被加载。默认值为#。
`OPTIONS('COMMENTCHAR'='#')`
- FILEHEADER：如果源文件中没有表头，可在LOAD DATA命令中提供表头。
`OPTIONS('FILEHEADER'='column1,column2')`
- ESCAPECHAR：如果用户想在CSV上对Escape字符进行严格验证，可以提供Escape字符。默认值为\。
`OPTIONS('ESCAPECHAR'='\')`

📖 说明

如果在CSV数据中输入ESCAPECHAR，该ESCAPECHAR必须在双引号（" "）内。例如：`"a \b"`。

- Bad Records处理：
为了使数据处理应用程序为用户增值，不可避免地需要对数据进行某种程度的集成。在大多数情况下，数据质量问题源于生成源数据的上游（主要）系统。
有两种完全不同的方式处理Bad Data：
 - 按照原始数据加载所有数据，之后进行除错处理。
 - 在进入数据源的过程中，可以清理或擦除Bad Data，或者在发现Bad Data时让数据加载失败。

有多个选项可用于在CarbonData数据加载过程中清除源数据。对于CarbonData数据中的Bad Records管理，请参见表3-35。

表 3-35 Bad Records Logger

配置项	默认值	描述
BAD_RECORDS_LOGGER_ENABLE	false	若设置为true，则将创建Bad Records日志文件，其中包含Bad Records的详细信息。
BAD_RECORDS_ACTION	FAIL	<p>以下为Bad Records的四种操作类型：</p> <ul style="list-style-type: none"> ● FORCE：通过将Bad Records存储为NULL来自动校正数据。 ● REDIRECT：无法加载Bad Records，并将其写入BAD_RECORD_PATH下的CSV文件中，默认不开启该类型，如需使用该类型，需要设置参数carbon.enable.badrecord.action.redirect为true。 ● IGNORE：既不加载Bad Records也不将其写入CSV文件。 ● FAIL：如果发现存在Bad Records，数据加载将会失败。 <p>说明 在加载数据时，如果所有记录都是Bad Records，则参数BAD_RECORDS_ACTION将不起作用，加载数据操作将会失败。</p>
IS_EMPTY_DATA_BAD_RECORD	false	如果设置为“false”，则空（""或,,）数据将不被视为Bad Records，如果设置为“true”，则空数据将被视为Bad Records。
BAD_RECORD_PATH	-	指定存储Bad Records的HDFS路径。默认值为Null。如果启用了Bad Records日志记录或者Bad Records操作重定向，则该路径必须由用户进行配置。

示例：

```
LOAD DATA INPATH 'filepath.csv' INTO TABLE tablename
OPTIONS('BAD_RECORDS_LOGGER_ENABLE'='true',
'BAD_RECORD_PATH'='hdfs://hacluster/tmp/carbon',
'BAD_RECORDS_ACTION'='REDIRECT',
'IS_EMPTY_DATA_BAD_RECORD'='false');
```

说明

使用“REDIRECT”选项，CarbonData会将所有的Bad Records添加到单独的CSV文件中，但是该文件内容不能用于后续的数据加载，因为其内容可能无法与源记录完全匹配。用户必须清理原始源记录以便于进一步的数据提取。该选项的目的只是让用户知道哪些记录被视为Bad Records。

- MAXCOLUMNS: 该可选参数指定了在一行中，由CSV解析器解析的最大列数。
OPTIONS('MAXCOLUMNS'='400')

表 3-36 MAXCOLUMNS

可选参数名称	默认值	最大值
MAXCOLUMNS	2000	20000

表 3-37 MAXCOLUMNS 可选参数的行为图

MAXCOLUMNS值	在文件Header选项中的列数	考虑的最终值
在加载项中未指定	5	2000
在加载项中未指定	6000	6000
40	7	文件header列数与MAXCOLUMNS值，两者中的最大值
22000	40	20000
60	在加载项中未指定	CSV文件第一行的列数与MAXCOLUMNS值，两者中的最大值

说明

对于设置MAXCOLUMNS Option的最大值，要求executor具有足够的内存，否则，数据加载会由于内存不足的错误而失败。

- 如果在创建表期间将SORT_SCOPE定义为GLOBAL_SORT，则可以指定在对数据进行排序时要使用的分区数。如果未配置或配置小于1，则将使用map任务的数量作为reduce任务的数量。建议每个reduce任务处理512MB - 1GB数据。
OPTIONS('GLOBAL_SORT_PARTITIONS'='2')

说明

增加分区数可能需要增加“spark.driver.maxResultSize”，因为在driver中收集的采样数据随着分区的增加而增加。

- DATEFORMAT: 此选项用于指定表的日期格式。
OPTIONS('DATEFORMAT'='dateFormat')

📖 说明

日期格式由日期模式字符串指定。Carbon中的日期模式字母与JAVA中的日期模式字母相同。

- **TIMESTAMPFORMAT**: 此选项用于指定表的时间戳格式。
- `OPTIONS('TIMESTAMPFORMAT'='timestampFormat')`
- **SKIP_EMPTY_LINE**: 数据加载期间，此选项将忽略CSV文件中的空行。
`OPTIONS('SKIP_EMPTY_LINE'='TRUE/FALSE')`
- **可选: SCALE_FACTOR**: 针对RANGE_COLUMN，SCALE_FACTOR用来控制分区的数量，根据如下公式：
splitSize = max(blocklet_size, (block_size - blocklet_size)) * scale_factor
numPartitions = total size of input data / splitSize
默认值为3，range的范围为[1, 300]。
`OPTIONS('SCALE_FACTOR'='10')`

📖 说明

- 如果GLOBAL_SORT_PARTITIONS和SCALE_FACTOR同时使用，只有GLOBAL_SORT_PARTITIONS生效。
- RANGE_COLUMN合并默认使用LOCAL_SORT。
LOCAL_SORT与分区表的DDL操作存在冲突，不能同时使用，且对分区表性能提升不明显，不建议在分区表上启用该特性。

使用场景

可使用下列语句从CSV文件加载CarbonData table。

```
LOAD DATA INPATH 'folder path' INTO TABLE tablename  
OPTIONS(property_name=property_value, ...);
```

示例

data.csv源文件数据如下所示：

```
ID,date,country,name,phonetype,serialname,salary  
4,2014-01-21 00:00:00,xxx,aaa4,phone2435,ASD66902,15003  
5,2014-01-22 00:00:00,xxx,aaa5,phone2441,ASD90633,15004  
6,2014-03-07 00:00:00,xxx,aaa6,phone294,ASD59961,15005
```

```
CREATE TABLE carbontable(ID int, date Timestamp, country String, name String,  
phonetype String, serialname String,salary int) STORED AS carbondata;
```

```
LOAD DATA inpath 'hdfs://hacluster/tmp/data.csv' INTO table carbontable  
options('DELIMITER'=',');
```

系统响应

可在driver日志中查看命令运行成功或失败。

3.7.2.2 UPDATE CARBON TABLE

命令功能

UPDATE命令根据列表表达式和可选的过滤条件更新CarbonData表。

命令格式

- 格式1：
`UPDATE <CARBON TABLE> SET (column_name1, column_name2, ... column_name n) = (column1_expression , column2_expression , column3_expression ... column n_expression) [WHERE { <filter_condition> }];`
- 格式2：
`UPDATE <CARBON TABLE> SET (column_name1, column_name2,) = (select sourceColumn1, sourceColumn2 from sourceTable [WHERE { <filter_condition> }]) [WHERE { <filter_condition> }];`

参数描述

表 3-38 UPDATE 参数

参数	描述
CARBON TABLE	在其中执行更新操作的CarbonData表的名称。
column_name	待更新的目标列。
sourceColumn	需在目标表中更新的源表的列值。
sourceTable	将其记录更新到目标CarbonData表中的表。

注意事项

以下是使用UPDATE命令的条件：

- 如果源表中的多个输入行与目标表中的单行匹配，则UPDATE命令失败。
- 如果源表生成空记录，则UPDATE操作将在不更新表的情况下完成。
- 如果源表的行与目标表中任何已有的行不对应，则UPDATE操作将完成，不更新表。
- 具有二级索引的表不支持UPDATE命令。
- 在子查询中，如果源表和目标表相同，则UPDATE操作失败。
- 如果在UPDATE命令中使用的子查询包含聚合函数或group by子句，则UPDATE操作失败。

例如，`update t_carbn01 a set (a.item_type_code, a.profit) = (select b.item_type_cd, sum(b.profit) from t_carbn01b b where item_type_cd =2 group by item_type_code);`

其中，在子查询中使用聚合函数sum(b.profit)和group by子句，因此UPDATE操作失败。

- 如果查询的表设置了carbon.input.segments属性，则UPDATE操作失败。要解决这个问题，在查询前执行以下语句。

语法：

```
SET carbon.input.segments. <database_name>. <table_name>=*;
```

示例

- 示例1：

```
update carbonTable1 d set (d.column3,d.column5 ) = (select s.c33 ,s.c55  
from sourceTable1 s where d.column1 = s.c11) where d.column1 =  
'country' exists( select * from table3 o where o.c2 > 1);
```
- 示例2：

```
update carbonTable1 d set (c3) = (select s.c33 from sourceTable1 s where  
d.column1 = s.c11) where exists( select * from iud.other o where o.c2 >  
1);
```
- 示例3：

```
update carbonTable1 set (c2, c5 ) = (c2 + 1, concat(c5 , "y" ));
```
- 示例4：

```
update carbonTable1 d set (c2, c5 ) = (c2 + 1, "xyx") where d.column1 =  
'india';
```
- 示例5：

```
update carbonTable1 d set (c2, c5 ) = (c2 + 1, "xyx") where d.column1 =  
'india' and exists( select * from table3 o where o.column2 > 1);
```

系统响应

可在driver日志和客户端中查看命令运行成功或失败。

3.7.2.3 DELETE RECORDS from CARBON TABLE

命令功能

DELETE RECORDS命令从CarbonData表中删除记录。

命令格式

```
DELETE FROM CARBON_TABLE [WHERE expression];
```

参数描述

表 3-39 DELETE RECORDS 参数

参数	描述
CARBON TABLE	在其中执行删除操作的CarbonData表的名称。

注意事项

- 删除segment将删除相应segment的所有二级索引。
- 如果查询的表设置了carbon.input.segments属性，则DELETE操作失败。要解决该问题，在查询前执行以下语句。

语法：

```
SET carbon.input.segments. <database_name>.<table_name>=*;
```

示例

- 示例1：
`delete from columncarbonTable1 d where d.column1 = 'country';`
- 示例2：
`delete from dest where column1 IN ('country1', 'country2');`
- 示例3：
`delete from columncarbonTable1 where column1 IN (select column11 from sourceTable2);`
- 示例4：
`delete from columncarbonTable1 where column1 IN (select column11 from sourceTable2 where column1 = 'xxx');`
- 示例5：
`delete from columncarbonTable1 where column2 >= 4;`

系统响应

可在driver日志和客户端中查看命令运行成功或失败。

3.7.2.4 INSERT INTO CARBON TABLE

命令功能

INSERT命令用于将SELECT查询结果加载到CarbonData表中。

命令格式

```
INSERT INTO [CARBON TABLE] [select query];
```

参数描述

表 3-40 INSERT INTO 参数

参数	描述
CARBON TABLE	需要执行INSERT命令的CarbonData表的名称。
select query	Source表上的SELECT查询（支持CarbonData、Hive和Parquet表）。

注意事项

- 表必须已经存在。
- 用户应属于数据加载组以执行数据加载操作。默认情况下，数据加载组被命名为“ficommon”。
- CarbonData表不支持Overwrite。
- 源表和目标表的数据类型应该相同，否则源表中的数据将被视为Bad Records。
- **INSERT INTO**命令不支持部分成功（partial success），如果存在Bad Records，该命令会失败。
- 在从源表插入数据到目标表的过程中，无法在源表中加载或更新数据。若要在INSERT操作期间启用数据加载或更新，请将以下参数配置为“true”。
“carbon.insert.persist.enable” = “true”
默认上述参数配置为“false”。

📖 说明

启用该参数将降低INSERT操作的性能。

示例

```
create table carbon01(a int,b string,c string) stored as carbondata;  
insert into table carbon01 values(1,'a','aa'),(2,'b','bb'),(3,'c','cc');  
create table carbon02(a int,b string,c string) stored as carbondata;  
INSERT INTO carbon02 select * from carbon01 where a > 1;
```

系统响应

可在driver日志中查看命令运行成功或失败。

3.7.2.5 DELETE SEGMENT by ID

命令功能

DELETE SEGMENT by ID命令是使用Segment ID来删除segment。

命令格式

```
DELETE FROM TABLE db_name.table_name WHERE SEGMENT.ID IN  
(segment_id1,segment_id2);
```

参数描述

表 3-41 DELETE LOAD 参数描述

参数	描述
segment_id	将要删除的Segment的ID。

参数	描述
db_name	Database名称，若未指定，则使用当前database。
table_name	在给定的database中的表名。

注意事项

流式表不支持删除segment。

示例

```
DELETE FROM TABLE CarbonDatabase.CarbonTable WHERE SEGMENT.ID IN (0);
```

```
DELETE FROM TABLE CarbonDatabase.CarbonTable WHERE SEGMENT.ID IN (0,5,8);
```

系统响应

操作成功或失败会在CarbonData日志中被记录。

3.7.2.6 DELETE SEGMENT by DATE

命令功能

DELETE SEGMENT by DATE命令用于通过加载日期删除CarbonData segment，在特定日期之前创建的segment将被删除。

命令格式

```
DELETE FROM TABLE db_name.table_name WHERE SEGMENT.STARTTIME BEFORE date_value;
```

参数描述

表 3-42 DELETE SEGMENT by DATE 参数描述

参数	描述
db_name	Database名称，若未指定，则使用当前database。
table_name	给定database中的表名。
date_value	有效Segment加载启动时间。在这个指定日期前的Segment将被删除。

注意事项

流式表不支持删除segment。

示例

```
DELETE FROM TABLE db_name.table_name WHERE SEGMENT.STARTTIME  
BEFORE '2017-07-01 12:07:20';
```

其中，STARTTIME是不同负载的加载启动时间。

系统响应

操作成功或失败会在CarbonData日志中被记录。

3.7.2.7 SHOW SEGMENTS

命令功能

SHOW SEGMENTS命令是用来向用户展示CarbonData table的Segment。

命令格式

```
SHOW SEGMENTS FOR TABLE [db_name.]table_name LIMIT number_of_loads;
```

参数描述

表 3-43 SHOW SEGMENTS FOR TABLE 参数描述

参数	描述
db_name	Database名，若未指定，则使用当前database。
table_name	在给定database中的表名。
number_of_loads	加载数的限制。

注意事项

无。

示例

```
create table carbon01(a int,b string,c string) stored as carbondata;  
insert into table carbon01 select 1,'a','aa';  
insert into table carbon01 select 2,'b','bb';  
insert into table carbon01 select 3,'c','cc';  
SHOW SEGMENTS FOR TABLE carbon01 LIMIT 2;
```

系统响应

```
+-----+-----+-----+-----+-----+-----+-----+-----+  
+  
| ID | Status | Load Start Time | Load Time Taken | Partition | Data Size | Index Size | File Format |
```

id	status	start_time	end_time	size	format	block_size	row_size	columnar_v3
3	Success	2020-09-28 22:53:26.336	2020-09-28 22:53:01.702	3.726S	{}	6.47KB	3.30KB	columnar_v3
2	Success	2020-09-28 22:53:01.702	2020-09-28 22:53:01.702	6.688S	{}	6.47KB	3.30KB	columnar_v3

3.7.2.8 CREATE SECONDARY INDEX

命令功能

该命令用于在CarbonData表中创建二级索引表。

命令格式

```
CREATE INDEX index_name
ON TABLE [db_name.]table_name (col_name1, col_name2)
AS 'carbodata'
PROPERTIES ('table_blocksize'='256');
```

参数描述

表 3-44 CREATE SECONDARY INDEX 参数

参数	描述
index_name	索引表的名称。表名称应由字母数字字符和下划线（_）特殊字符组成。
db_name	数据库的名称。数据库名称应由字母数字字符和下划线（_）特殊字符组成。
table_name	数据库中的表名称。表名称应由字母数字字符和下划线（_）特殊字符组成。
col_name	表中的列名称。支持多列。列名称应由字母数字字符和下划线（_）特殊字符组成。
table_blocksize	数据文件的block大小。更多详细信息，请参考 Block大小 。

注意事项

db_name为可选项。

示例

```
create table productdb.productSalesTable(id int,price int,productName
string,city string) stored as carbodata;

CREATE INDEX productNameIndexTable on table productdb.productSalesTable
(productName,city) as 'carbodata' ;
```


上述示例将创建名为“productdb.productNameIndexTable”的二级表并加载所提供的索引信息。

系统响应

将创建二级索引表，加载与所提供的列相关的索引信息到二级索引表中，并将成功消息记录在系统日志中。

3.7.2.9 SHOW SECONDARY INDEXES

命令功能

该命令用于在所提供的CarbonData表中显示所有的二级索引表。

命令格式

```
SHOW INDEXES ON db_name.table_name;
```

参数描述

表 3-45 SHOW SECONDARY INDEXES 参数

参数	描述
db_name	数据库的名称。数据库名称应由字母数字字符和下划线（_）特殊字符组成
table_name	数据库中的表名称。表名称应由字母数字字符和下划线（_）特殊字符组成。

注意事项

db_name为可选项。

示例

```
create table productdb.productSalesTable(id int,price int,productName  
string,city string) stored as carbondata;
```

```
CREATE INDEX productNameIndexTable on table productdb.productSalesTable  
(productName,city) as 'carbondata' ;
```

```
SHOW INDEXES ON productdb.productSalesTable;
```

系统响应

显示列出给定CarbonData表中的所有索引表和相应的索引列。

3.7.2.10 DROP SECONDARY INDEX

命令功能

该命令用于删除给定表中存在的二级索引表。

命令格式

```
DROP INDEX [IF EXISTS] index_name ON [db_name.]table_name;
```

参数描述

表 3-46 DROP SECONDARY INDEX 参数

参数	描述
index_name	索引表的名称。表名称应由字母数字字符和下划线（_）特殊字符组成。
db_name	数据库的名称。若未指定，选择当前默认数据库。
table_name	需要删除的表的名称。

注意事项

该命令中IF EXISTS和db_name为可选项。

示例

```
DROP INDEX if exists productNameIndexTable ON  
productdb.productSalesTable;
```

系统响应

二级索引表将被删除，索引信息将在CarbonData表中被清除，删除成功的消息将记录在系统日志中。

3.7.2.11 CLEAN FILES

命令功能

DELETE SEGMENT命令会将删除的segments标识为delete状态；segment合并后，旧的segments状态会变为compacted。这些segments的数据文件不会从物理上删除。如果用户希望强制删除这些文件，可以使用**CLEAN FILES**命令。

但是，使用该命令可能会导致查询命令执行失败。

命令格式

```
CLEAN FILES FOR TABLE [db_name.]table_name ;
```

参数描述

表 3-47 CLEAN FILES FOR TABLE 参数描述

参数	描述
db_name	数据库名称。数据库名称由字母，数字和下划线组成。
table_name	数据库中的表的名称。表名由字母，数字和下划线组成。

注意事项

无。

示例

添加carbon配置参数

```
carbon.clean.file.force.allowed = true
```

```
create table carbon01(a int,b string,c string) stored as carbondata;
```

```
insert into table carbon01 select 1,'a','aa';
```

```
insert into table carbon01 select 2,'b','bb';
```

```
delete from table carbon01 where segment.id in (0);
```

```
show segments for table carbon01;
```

```
CLEAN FILES FOR TABLE carbon01 options('force'='true');
```

```
show segments for table carbon01;
```

上述命令将从物理上删除所有DELETE SEGMENT命令删除的segment和合并后的旧的segment。

系统响应

可在driver日志中查看命令运行成功或失败。

3.7.2.12 SET/RESET

命令功能

此命令用于动态Add，Update，Display或Reset CarbonData参数，而无需重新启动driver。

命令格式

- Add或Update参数值：
SET *parameter_name=parameter_value*
此命令用于添加或更新“parameter_name”的值。

- Display参数值：
SET parameter_name
此命令用于显示指定的“parameter_name”的值。
- Display会话参数：
SET
此命令显示所有支持的会话参数。
- Display会话参数以及使用细节：
SET -v
此命令显示所有支持的会话参数及其使用细节。
- Reset参数值：
RESET
此命令清除所有会话参数。

参数描述

表 3-48 SET 参数描述

参数	描述
parameter_name	其值需要被动态添加（add），更新（update）或显示（display）的参数名称。
parameter_value	将要设置的“parameter_name”的新值。

注意事项

以下为分别使用SET和RESET命令进行动态设置或清除操作的属性：

表 3-49 属性描述

属性	描述
carbon.options.bad.records.logger.enable	启用或禁用bad record日志记录。
carbon.options.bad.records.action	指定bad record操作，例如，强制（force），重定向（redirect），失败（fail）或忽略（ignore）。有关详细信息，请参阅 Bad Records处理 ：
carbon.options.is.empty.data.bad.record	指定空数据是否被视为bad record。有关详细信息，请参阅 Bad Records处理 ：
carbon.options.sort.scope	指定数据加载期间排序的范围。
carbon.options.bad.record.path	指定需要存储bad record的HDFS路径。

属性	描述
carbon.custom.block.distribution	指定是否使用Spark或CarbonData的块分布功能。
enable.unsafe.sort	指定在数据加载期间是否使用不安全的排序。不安全的排序可减少数据加载操作期间的垃圾回收，从而实现更好的性能。
carbon.si.lookup.partialstring	当参数设置为TRUE时，二级索引采用 starts-with、ends-with、contains和 LIKE分区条件字符串。 当参数设置为FALSE时，二级索引只采用 starts-with分区条件字符串。
carbon.input.segments	<p>指定要查询的段ID。此属性允许您查询指定表的指定段。CarbonScan将仅从指定的段ID读取数据。</p> <p>语法： “carbon.input.segments. <database_name>. <table_name> = <list of segment ids >”</p> <p>如果用户想在多线程模式下查询指定段，可使用CarbonSession.threadSet代替SET语句。</p> <p>语法： “CarbonSession.threadSet ("carbon.input.segments. <database_name>. <table_name>","<list of segment ids >");”</p> <p>说明 不建议在carbon.properties文件中设置该属性，因为所有会话都包含段列表，除非发生会话级或线程级覆盖。</p>

示例

- 添加（Add）或更新（Update）：
SET enable.unsafe.sort=true
- 显示（Display）属性值：
SET enable.unsafe.sort
- 显示段ID列表，段状态和其他所需详细信息的示例，然后指定要读取的段列表：
SHOW SEGMENTS FOR TABLE carbontable1;
SET carbon.input.segments.db.carbontable1 = 1, 3, 9;
- 多线程模式查询指定段示例如下：
CarbonSession.threadSet
("carbon.input.segments.default.carbon_table_Multi_Thread", "1,3");

- 在多线程环境中使用 `CarbonSession.threadSet` 查询段示例如下（以 Scala 代码为例）：

```
def main(args: Array[String]) {
  Future
  {
    CarbonSession.threadSet("carbon.input.segments.default.carbon_table_MuTI_THread", "1")
    spark.sql("select count(empno) from carbon_table_MuTI_THread").show()
  }
}
```
- 重置（Reset）：

RESET

系统响应

- 若运行成功，将记录在 driver 日志中。
- 若出现故障，将显示在用户界面（UI）中。

3.7.3 CarbonData 表操作并发语法说明

DDL 和 DML 中的操作，执行前，需要获取对应的锁，各操作需要获取锁的情况见 [表1 操作获取锁一览表](#)，√ 表示需要获取该锁，一个操作仅在获取到所有需要获取的锁后，才能继续执行。

任意两个操作是否可以并发执行，可以通过如下方法确定：[表3-50](#) 两行代表两个操作，这两行没有任意一列都标记√，即不存在某一列两行全为√。

表 3-50 操作获取锁一览表

操作	MET ADA TA_L OCK	COM PAC TIO N_L OCK	DRO P_TA BLE_ LOC K	DELE TE_S EGM ENT_ LOC K	CLE AN_ FILE S_LO CK	ALTE R_PA RTIT ION_ LOC K	UPD ATE_ LOC K	STRE AMI NG_ LOC K	CON CUR REN T_LO AD_ LOC K	SEG MEN T_LO CK
CRE ATE TABL E	-	-	-	-	-	-	-	-	-	-
CRE ATE TABL E As SELE CT	-	-	-	-	-	-	-	-	-	-
DRO P TABL E	√	-	√	-	-	-	-	√	-	-

操作	MET ADA TA_L OCK	COM PAC TION N_L OCK	DRO P_TA BLE_ LOC K	DELE TE_S EGM ENT_ LOC K	CLE AN_ FILE S_LO CK	ALTE R_PA RTIT ION_ LOC K	UPD ATE_ LOC K	STRE AMI NG_ LOC K	CON CUR REN T_LO AD_ LOC K	SEG MEN T_LO CK
ALTE R TABL E COM PACT ION	-	√	-	-	-	-	√	-	-	-
TABL E REN AME	-	-	-	-	-	-	-	-	-	-
ADD COL UM NS	√	√	-	-	-	-	-	-	-	-
DRO P COL UM NS	√	√	-	-	-	-	-	-	-	-
CHA NGE DAT A TYPE	√	√	-	-	-	-	-	-	-	-
REFR ESH TABL E	-	-	-	-	-	-	-	-	-	-
REGI STER INDE X TABL E	√	-	-	-	-	-	-	-	-	-
REFR ESH INDE X	-	√	-	-	-	-	-	-	-	-

操作	MET ADA TA_L OCK	COM PAC TION N_L OCK	DRO P_TA BLE_ LOC K	DELE TE_S EGM ENT_ LOC K	CLE AN_ FILE S_LO CK	ALTE R_PA RTIT ION_ LOC K	UPD ATE_ LOC K	STRE AMI NG_ LOC K	CON CURRE NT_LO AD_ LOC K	SEG MEN T_LO CK
LOA D DAT A/ INSE RT INTO	-	-	-	-	-	-	-	-	√	√
UPD ATE CAR BON TABL E	√	√	-	-	-	-	√	-	-	-
DELE TE REC ORD S from CAR BON TABL E	√	√	-	-	-	-	√	-	-	-
DELE TE SEG MEN T by ID	-	-	-	√	√	-	-	-	-	-
DELE TE SEG MEN T by DAT E	-	-	-	√	√	-	-	-	-	-
SHO W SEG MEN TS	-	-	-	-	-	-	-	-	-	-

操作	MET ADA TA_L OCK	COM PAC TION N_L OCK	DRO P_TA BLE_ LOC K	DELE TE_S EGM ENT_ LOC K	CLE AN_ FILE S_LO CK	ALTE R_PA RTIT ION_ LOC K	UPD ATE_ LOC K	STRE AMI NG_ LOC K	CON CUR REN T_LO AD_ LOC K	SEG MEN T_LO CK
CRE ATE SEC OND ARY INDE X	√	√	-	√	-	-	-	-	-	-
SHO W SEC OND ARY INDE XES	-	-	-	-	-	-	-	-	-	-
DRO P SEC OND ARY INDE X	√	-	√	-	-	-	-	-	-	-
CLEA N FILES	-	-	-	-	-	-	-	-	-	-
SET/ RESE T	-	-	-	-	-	-	-	-	-	-
Add Hive Parti tion	-	-	-	-	-	-	-	-	-	-
Drop Hive Parti tion	√	√	√	√	√	√	-	-	-	-
Drop Parti tion	√	√	√	√	√	√	-	-	-	-

操作	MET ADA TA_L OCK	COM PAC TIO N_L OCK	DRO P_TA BLE_ LOC K	DELE TE_S EGM ENT_ LOC K	CLE AN_ FILE S_LO CK	ALTE R_PA RTIT ION_ LOC K	UPD ATE_ LOC K	STRE AMI NG_ LOC K	CON CUR REN T_LO AD_ LOC K	SEG MEN T_LO CK
Alter table set	√	√	-	-	-	-	-	-	-	-

3.7.4 CarbonData Segment API 语法说明

本章节描述Segment的API以及使用方法，所有方法在org.apache.spark.util.CarbonSegmentUtil类中。

如下方法已废弃：

```
/**
 * Returns the valid segments for the query based on the filter condition
 * present in carbonScanRdd.
 *
 * @param carbonScanRdd
 * @return Array of valid segments
 */
@deprecated def getFilteredSegments(carbonScanRdd: CarbonScanRDD[InternalRow]): Array[String];
```

使用方法

使用如下方法从查询语句中获得CarbonScanRDD：

```
val df=carbon.sql("select * from table where age='12'")
val myscan=df.queryExecution.sparkPlan.collect {
case scan: CarbonDataSourceScan if scan.rdd.isInstanceOf[CarbonScanRDD[InternalRow]] => scan.rdd
case scan: RowDataSourceScanExec if scan.rdd.isInstanceOf[CarbonScanRDD[InternalRow]] => scan.rdd
}.head
val carbonrdd=myscan.asInstanceOf[CarbonScanRDD[InternalRow]]
```

例子：

```
CarbonSegmentUtil.getFilteredSegments(carbonrdd)
```

可以通过传入sql语句来获取过滤后的segment：

```
/**
 * Returns an array of valid segment numbers based on the filter condition provided in the sql
 * NOTE: This API is supported only for SELECT Sql (insert into,ctas,... is not supported)
 *
 * @param sql
 * @param sparkSession
 * @return Array of valid segments
 * @throws UnsupportedOperationException because Get Filter Segments API supports if and only
 * if only one carbon main table is present in query.
 */
def getFilteredSegments(sql: String, sparkSession: SparkSession): Array[String];
```

例子：

```
CarbonSegmentUtil.getFilteredSegments("select * from table where age='12'", sparkSession)
```

传入数据库名和表名，获取会被合并的segment列表，得到的segment列表可以当做getMergedLoadName函数的参数传入：

```
/**
 * Identifies all segments which can be merged with MAJOR compaction type.
 * NOTE: This result can be passed to getMergedLoadName API to get the merged load name.
 *
 * @param sparkSession
 * @param tableName
 * @param dbName
 * @return list of LoadMetadataDetails
 */
def identifySegmentsToBeMerged(sparkSession: SparkSession,
                               tableName: String,
                               dbName: String) : util.List[LoadMetadataDetails];
```

例子:

```
CarbonSegmentUtil.identifySegmentsToBeMerged(sparkSession, "table_test", "default")
```

传入数据库名、表名和自定义的segment列表，获取自定义合并操作会被合并的segment列表，得到的segment列表可以当做getMergedLoadName函数的参数传入：

```
/**
 * Identifies all segments which can be merged with CUSTOM compaction type.
 * NOTE: This result can be passed to getMergedLoadName API to get the merged load name.
 *
 * @param sparkSession
 * @param tableName
 * @param dbName
 * @param customSegments
 * @return list of LoadMetadataDetails
 * @throws UnsupportedOperationException if customSegments is null or empty.
 * @throws MalformedCarbonCommandException if segment does not exist or is not valid
 */
def identifySegmentsToBeMergedCustom(sparkSession: SparkSession,
                                      tableName: String,
                                      dbName: String,
                                      customSegments: util.List[String]): util.List[LoadMetadataDetails];
```

例子:

```
val customSegments = new util.ArrayList[String]()
customSegments.add("1")
customSegments.add("2")
CarbonSegmentUtil.identifySegmentsToBeMergedCustom(sparkSession, "table_test", "default",
                                                    customSegments)
```

给定segment列表，返回合并后新的导入名称:

```
/**
 * Returns the Merged Load Name for given list of segments
 *
 * @param list of segments
 * @return Merged Load Name
 * @throws UnsupportedOperationException if list of segments is less than 1
 */
def getMergedLoadName(list: util.List[LoadMetadataDetails]): String;
```

例子:

```
val carbonTable = CarbonEnv.getCarbonTable(Option(databaseName), tableName)(sparkSession)
val loadMetadataDetails = SegmentStatusManager.readLoadMetadata(carbonTable.getMetadataPath)
CarbonSegmentUtil.getMergedLoadName(loadMetadataDetails.toList.asJava)
```

3.7.5 CarbonData 表空间索引语法说明

快速示例

```
create table IF NOT EXISTS carbonTable
(
```

```
COLUMN1 BIGINT,  
LONGITUDE BIGINT,  
LATITUDE BIGINT,  
COLUMN2 BIGINT,  
COLUMN3 BIGINT  
)  
STORED AS carbondata  
TBLPROPERTIES  
(  
  'SPATIAL_INDEX.mygeohash.type'='geohash',  
  'SPATIAL_INDEX.mygeohash.sourcecolumns'='longitude,  
  latitude',  
  'SPATIAL_INDEX.mygeohash.originLatitude'='39.850713',  
  'SPATIAL_INDEX.mygeohash.gridSize'='50',  
  'SPATIAL_INDEX.mygeohash.minLongitude'='115.828503',  
  'SPATIAL_INDEX.mygeohash.maxLongitude'='720.000  
  000',  
  'SPATIAL_INDEX.mygeohash.minLatitude'='39.850713',  
  'SPATIAL_INDEX.mygeohash.maxLatitude'='720.0  
  0000',  
  'SPATIAL_INDEX'='mygeohash',  
  'SPATIAL_INDEX.mygeohash.conversionRatio'='1000000',  
  'SORT_COLUMNS'='column1,column2,column3,latitude,longitude');
```

空间索引介绍

空间数据包括多维点、线、矩形、立方体、多边形和其他几何对象。空间数据对象占据空间的某一区域，称为空间范围，通过其位置和边界描述。空间数据可以是点数据，也可以是区域数据。

- 点数据：一个点具有一个空间范围，仅通过其位置描述。它不占用空间，没有相关的边界。点数据由二维空间中的点的集合组成。点可以存储为一对经纬度。
- 区域数据：一个区域有空间范围，有位置和边界。位置可以看作是一个定点在区域内的位置，例如它的质心。在二维中，边界可以可视化为一组线（有限区域，闭环）。区域数据包含一系列区域。

目前仅限于支持点数据，存储点数据。

经纬度可以编码为唯一的 GeoID。Geohash 是 Gustavo Niemeyer 发明的公共域地理编码系统，它将地理位置编码为一串由字母和数字组成的短字符串。它是一种分层的空间数据结构，把空间细分为网格形状的桶，是被称为 Z 阶曲线和通常称为空间填充曲线的许多应用之一。

点在多维中的 Z 值是简单地通过交织其坐标值的二进制表示来计算的，如下图所示。使用 Geohash 创建 GeoID 时，数据按照 GeoID 排序，而不是按照经纬度排序，数据按照空间就近性排序存储。

	x: 0	1	2	3	4	5	6	7
	000	001	010	011	100	101	110	111
y: 0	000000	000001	000100	000101	010000	010001	010100	010101
1	000010	000011	000110	000111	010010	010011	010110	010111
2	001000	001001	001100	001101	011000	011001	011100	011101
3	001010	001011	001110	001111	011010	011011	011110	011111
4	100000	100001	100100	100101	110000	110001	110100	110101
5	100010	100011	100110	100111	110010	110011	110110	110111
6	101000	101001	101100	101101	111000	111001	111100	111101
7	101010	101011	101110	101111	111010	111011	111110	111111

建表

GeoHash编码:

```
create table IF NOT EXISTS carbonTable
(
...
`LONGITUDE` BIGINT,
`LATITUDE` BIGINT,
...
)
STORED AS carbondata
TBLPROPERTIES
('SPATIAL_INDEX.mygeohash.type='geohash','SPATIAL_INDEX.mygeohash.sourcecolumns='longitude,
latitude','SPATIAL_INDEX.mygeohash.originLatitude='xx.xxxxxx','SPATIAL_INDEX.mygeohash.gridSize='xx','SP
ATIAL_INDEX.mygeohash.minLongitude='xxx.xxxxxx','SPATIAL_INDEX.mygeohash.maxLongitude='xxx.xxxxxx'
','SPATIAL_INDEX.mygeohash.minLatitude='xx.xxxxxx','SPATIAL_INDEX.mygeohash.maxLatitude='xxx.xxxxxx','
SPATIAL_INDEX='mygeohash','SPATIAL_INDEX.mygeohash.conversionRatio='1000000','SORT_COLUMNS='co
lumn1,column2,column3,latitude,longitude');
```

SPATIAL_INDEX: 自定义索引处理器。此处理程序允许用户从表结构列集中创建新的列。新创建的列名与处理程序名相同。处理程序的type和sourcecolumns属性是必需的属性。目前，type属性只支持“geohash”。Carbon提供一个简单的默认实现类。用户可以通过扩展默认实现类来挂载geohash的自定义实现类。该默认处理程序还需提供以下的表属性:

- SPATIAL_INDEX.xxx.originLatitude: Double类型，坐标原点纬度
- SPATIAL_INDEX.xxx.gridSize: Int类型，栅格长度（米）
- SPATIAL_INDEX.xxx.minLongitude: Double类型，最小经度
- SPATIAL_INDEX.xxx.maxLongitude: Double类型，最大经度
- SPATIAL_INDEX.xxx.minLatitude: Double类型，最小纬度

- SPATIAL_INDEX.xxx.maxLatitude: Double类型，最大纬度
- SPATIAL_INDEX.xxx.conversionRatio: Int类型，将经纬度小数值转换为整型值

用户可以按照上述格式为处理程序添加自己的表属性，并在自定义实现类中访问它们。originLatitude, gridSize及conversionRatio是必选参数，其余属性在Carbon中都是可选的。可以使用“SPATIAL_INDEX.xxx.class”属性指定它们的实现类。

默认实现类可以为每一行的sourcecolumns生成handler列值，并且支持基于sourcecolumns的过滤条件查询。生成的handler列对用户不可见。除SORT_COLUMNS表属性外，任何DDL命令和属性都不允许包含handler列。

说明

- 生成的handler列默认被视为排序列。如果SORT_COLUMNS不包含任何sourcecolumns，则将handler列追加到现有的SORT_COLUMNS最后。如果在SORT_COLUMNS中已经指定了该handler列，则它在SORT_COLUMNS的顺序将保持不变。
- 如果SORT_COLUMNS包含任意的sourcecolumns，但是没有包含handler列，则handler列将自动插入到SORT_COLUMNS中的sourcecolumns之前。
- 如果SORT_COLUMNS需要包含任意的sourcecolumns，那么需要保证handler列出现在sourcecolumns之前，这样handler列才能在排序中生效。

GeoSOT编码:

```
CREATE TABLE carbontable(
...
longitude DOUBLE,
latitude DOUBLE,
...)
STORED AS carbondata
TBLPROPERTIES ('SPATIAL_INDEX'='xxx',
'SPATIAL_INDEX.xxx.type'='geosot',
'SPATIAL_INDEX.xxx.sourcecolumns'='longitude, latitude',
'SPATIAL_INDEX.xxx.level'='21',
'SPATIAL_INDEX.xxx.class'='org.apache.carbondata.geo.GeoSOTIndex')
```

表 3-51 参数说明

参数	说明
SPATIAL_INDEX	指定表属性“SPATIAL_INDEX”，空间索引列，列名与该属性的值相同。
SPATIAL_INDEX.xxx.type	必填参数，值为geosot。
SPATIAL_INDEX.xxx.sourcecolumns	必填参数，空间索引列属性，指定计算空间索引的源数据列，需为2个存在的列，且类型为double。
SPATIAL_INDEX.xxx.level	可选参数，用于计算空间索引列。默认值为17，因为该值可以计算出足够精确的结果，同时拥有良好的性能。
SPATIAL_INDEX.xxx.class	可选参数，用于指定geo的实现类，默认为“org.apache.carbondata.geo.GeoSOTIndex”。

使用示例:

```
create table geosot(
timevalue bigint,
```

```
longitude double,  
latitude double)  
stored as carbondata  
TBLPROPERTIES ('SPATIAL_INDEX'='mygeosot',  
'SPATIAL_INDEX.mygeosot.type'='geosot',  
'SPATIAL_INDEX.mygeosot.level'='21', 'SPATIAL_INDEX.mygeosot.sourcecolumns'='longitude, latitude');
```

准备数据

- 准备数据文件1: geosotdata.csv

```
timevalue,longitude,latitude  
1575428400000,116.285807,40.084087  
1575428400000,116.372142,40.129503  
1575428400000,116.187332,39.979316  
1575428400000,116.337069,39.951887  
1575428400000,116.359102,40.154684  
1575428400000,116.736367,39.970323  
1575428400000,116.720179,40.009893  
1575428400000,116.346961,40.13355  
1575428400000,116.302895,39.930753  
1575428400000,116.288955,39.999101  
1575428400000,116.17609,40.129953  
1575428400000,116.725575,39.981115  
1575428400000,116.266922,40.179415  
1575428400000,116.353706,40.156483  
1575428400000,116.362699,39.942444  
1575428400000,116.325378,39.963129
```

- 准备数据文件2: geosotdata2.csv

```
timevalue,longitude,latitude  
1575428400000,120.17708,30.326882  
1575428400000,120.180685,30.326327  
1575428400000,120.184976,30.327105  
1575428400000,120.189311,30.327549  
1575428400000,120.19446,30.329698  
1575428400000,120.186965,30.329133  
1575428400000,120.177481,30.328911  
1575428400000,120.169713,30.325614  
1575428400000,120.164563,30.322243  
1575428400000,120.171558,30.319613  
1575428400000,120.176365,30.320687  
1575428400000,120.179669,30.323688  
1575428400000,120.181001,30.320761  
1575428400000,120.187094,30.32354  
1575428400000,120.193574,30.323651  
1575428400000,120.186192,30.320132  
1575428400000,120.190055,30.317464  
1575428400000,120.195376,30.318094  
1575428400000,120.160786,30.317094  
1575428400000,120.168211,30.318057  
1575428400000,120.173618,30.316612  
1575428400000,120.181001,30.317316  
1575428400000,120.185162,30.315908  
1575428400000,120.192415,30.315871  
1575428400000,120.161902,30.325614  
1575428400000,120.164306,30.328096  
1575428400000,120.197093,30.325985  
1575428400000,120.19602,30.321651  
1575428400000,120.198638,30.32354  
1575428400000,120.165421,30.314834
```

导入数据

GeoHash默认实现类扩展自定义索引抽象类。如果没有配置handler属性为自定义的实现类，则使用默认的实现类。用户可以通过扩展默认实现类来挂载geohash的自定义实现类。自定义索引抽象类方法包括：

- Init方法，用来提取、验证和存储handler属性。在失败时发生异常，并显示错误信息。
- Generate方法，用来生成索引。它为每行数据生成一个索引数据。
- Query方法，用来对给定输入生成索引值范围列表。

导入命令同普通Carbon表：

```
LOAD DATA inpath '/tmp/geosotdata.csv' INTO TABLE geosot OPTIONS ('DELIMITER'= ',');
```

```
LOAD DATA inpath '/tmp/geosotdata2.csv' INTO TABLE geosot OPTIONS ('DELIMITER'= ',');
```

📖 说明

geosotdata.csv和geosotdata2.csv表请参考[准备数据](#)。

不规则空间集合的聚合查询

查询语句及Filter UDF

- 根据polygon过滤数据

IN_POLYGON(pointList)

UDF输入参数：

参数	类型	说明
pointList	String	将多个点输入为一个字符串，每个点以 longitude latitude 表示。经纬度间用空格分隔，每对经纬度用逗号分隔，字符串首尾经纬度一致。

UDF输出参数：

参数	类型	说明
inOrNot	Boolean	判断数据是否在指定的polygon_list之内。

使用示例：

```
select longitude, latitude from geosot where IN_POLYGON('116.321011 40.123503, 116.137676 39.947911, 116.560993 39.935276, 116.321011 40.123503');
```

- 根据polygon列表过滤数据。

IN_POLYGON_LIST(polygonList, opType)

UDF输入参数：

参数	类型	说明
polygonList	String	<p>将多个polygon输入为一个字符串，每个polygon以POLYGON ((longitude1 latitude1, longitude2 latitude2, ...))表示。注意“POLYGON”后有空格，经纬度间用空格分隔，每对经纬度用逗号分隔，一个polygon的首尾经纬度一致。IN_POLYGON_LIST必须输入2个以上polygon。</p> <p>一个polygon示例： POLYGON ((116.137676 40.163503, 116.137676 39.935276, 116.560993 39.935276, 116.137676 40.163503))</p>
opType	String	<p>对多个polygon进行并交集操作。</p> <p>目前支持的操作类型：</p> <ul style="list-style-type: none"> • OR: A U B U C (假设输入了三个POLYGON, A、B、C) • AND: A ∩ B ∩ C

UDF输出参数：

参数	类型	说明
inOrNot	Boolean	判断数据是否在指定的polygon_list之内。

使用示例：

```
select longitude, latitude from geosot where IN_POLYGON_LIST('POLYGON ((120.176433 30.327431,120.171283 30.322245,120.181411 30.314540, 120.190509 30.321653,120.185188 30.329358,120.176433 30.327431)), POLYGON ((120.191603 30.328946,120.184179 30.327465,120.181819 30.321464, 120.190359 30.315388,120.199242 30.324464,120.191603 30.328946))', 'OR');
```

- 根据polyline列表过滤数据。

IN_POLYLINE_LIST(polylineList, bufferInMeter)

UDF输入参数：

参数	类型	说明
polylineList	String	<p>将多个polyline输入为一个字符串，每个polyline以LINestring (longitude1 latitude1, longitude2 latitude2, ...)表示。注意“LINestring”后有空格，经纬度间用空格分隔，每组经纬度用逗号分隔。</p> <p>对多个polyline区域内的数据会输出并集结果。</p> <p>一个polyline示例： <code>LINestring (116.137676 40.163503, 116.137676 39.935276, 116.260993 39.935276)</code></p>
bufferInMeter	Float	<p>polyline的buffer距离，单位为米。末端使用直角创建缓冲区。</p>

UDF输出参数：

参数	类型	说明
inOrNot	Boolean	判断数据是否在指定的polyline_list之内。

使用示例：

```
select longitude, latitude from geosot where IN_POLYLINE_LIST('LINestring (120.184179 30.327465, 120.191603 30.328946, 120.199242 30.324464, 120.190359 30.315388)', 65);
```

- 根据Geold区间列表过滤数据。

IN_POLYGON_RANGE_LIST(polygonRangeList, opType)

UDF输入参数：

参数	类型	说明
polygonRangeList	String	<p>将多个rangeList输入为一个字符串，每个rangeList以RANGELIST (startGeold1 endGeold1, startGeold2 endGeold2, ...)表示。注意“RANGELIST”后有空格，首尾Geold间用空格分隔，每组Geold range用逗号分隔。</p> <p>一个rangeList示例： <code>RANGELIST (855279368848 855279368850, 855280799610 855280799612, 855282156300 855282157400)</code></p>

参数	类型	说明
opType	String	对多个rangeList进行并交集操作。 目前支持的操作类型： <ul style="list-style-type: none"> OR: A U B U C (假设输入了三个 RANGELIST, A、B、C) AND: A ∩ B ∩ C

UDF输出参数:

参数	类型	说明
inOrNot	Boolean	判断数据是否在指定的polyRange_list之内。

使用示例:

```
select mygeosot, longitude, latitude from geosot where IN_POLYGON_RANGE_LIST('RANGELIST (526549722865860608 526549722865860618, 532555655580483584 532555655580483594)', 'OR');
```

- polygon连接查询

IN_POLYGON_JOIN(GEO_HASH_INDEX_COLUMN, POLYGON_COLUMN)

两张表做join查询，一张表为空间数据表（有经纬度列和GeoHashIndex列），另一张表为维度表，保存polygon数据。

查询使用IN_POLYGON_JOIN UDF，参数GEO_HASH_INDEX_COLUMN和polygon表的POLYGON_COLUMN。Polygon_column列是一系列的点（经纬度列）。Polygon表的每一行的第一个点和最后一个点必须是相同的。Polygon表的每一行的所有点连接起来形成一个封闭的几何对象。

UDF输入参数:

参数	类型	说明
GEO_HASH_INDEX_COLUMN	Long	空间数据表的GeoHashIndex列。
POLYGON_COLUMN	String	Polygon表的polygon列，数据为polygon的字符串表示。例如，一个polygon是POLYGON ((longitude1 latitude1, longitude2 latitude2, ...))

使用示例:

```
CREATE TABLE polygonTable(
polygon string,
poiType string,
poiId String)
STORED AS carbondata;
```

```
insert into polygonTable select 'POLYGON ((120.176433 30.327431,120.171283 30.322245, 120.181411 30.314540,120.190509 30.321653,120.185188 30.329358,120.176433 30.327431))','abc','1';
```

```
insert into polygonTable select 'POLYGON ((120.191603 30.328946,120.184179 30.327465,
120.181819 30.321464,120.190359 30.315388,120.199242 30.324464,120.191603 30.328946))','abc','2';

select t1.longitude,t1.latitude from geosot t1
inner join
(select polygon,poild from polygonTable where poitype='abc') t2
on in_polygon_join(t1.mygeosot,t2.polygon) group by t1.longitude,t1.latitude;
```

- range_list连接查询

IN_POLYGON_JOIN_RANGE_LIST(GEO_HASH_INDEX_COLUMN, POLYGON_COLUMN)

同IN_POLYGON_JOIN，使用IN_POLYGON_JOIN_RANGE_LIST UDF关联空间数据表和polygon维度表，关联基于Polygon_RangeList。直接使用range list可以避免polygon到range list的转换。

UDF输入参数：

参数	类型	说明
GEO_HASH_INDEX_COLUMN	Long	空间数据表的GeoHashIndex列。
POLYGON_COLUMN	String	Polygon表的rangelist列，数据为rangeList的字符串。例如，一个rangelist是RANGELIST (startGeold1 endGeold1, startGeold2 endGeold2, ...)

使用示例：

```
CREATE TABLE polygonTable(
polygon string,
poiType string,
poild String)
STORED AS carbondata;

insert into polygonTable select 'RANGELIST (526546455897309184 526546455897309284,
526549831217315840 526549831217315850, 532555655580483534 532555655580483584)','xyz','2';

select t1.*
from geosot t1
inner join
(select polygon,poild from polygonTable where poitype='xyz') t2
on in_polygon_join_range_list(t1.mygeosot,t2.polygon);
```

空间索引工具类UDF

- Geold转栅格行列号。

GeoldToGridXy(geold)

UDF输入参数：

参数	类型	说明
geold	Long	根据Geold计算栅格行列号。

UDF输出参数：

参数	类型	说明
gridArray	Array[Int]	返回该geoid所包含的栅格行列号，以数组的方式返回，第一位为行，第二位为列。

使用示例：

```
select longitude, latitude, mygeohash, GeoidToGridXy(mygeohash) as GridXY from geoTable;
```

- 经纬度转Geoid。

LatLngToGeoid(latitude, longitude oriLatitude, gridSize)

UDF输入参数：

参数	类型	说明
longitude	Long	经度，注：转换后的整数类型。
latitude	Long	纬度，注：转换后的整数类型。
oriLatitude	Double	原点纬度，计算Geoid需要参数。
gridSize	Int	栅格大小，计算Geoid需要参数。

UDF输出参数：

参数	类型	说明
geoid	Long	通过编码获得一个表示经纬度的数。

使用示例：

```
select longitude, latitude, mygeohash, LatLngToGeoid(latitude, longitude, 39.832277, 50) as geoid from geoTable;
```

- Geoid转经纬度。

GeoidToLatLng(geoid, oriLatitude, gridSize)

UDF输入参数：

参数	类型	说明
geoid	Long	根据Geoid计算经纬度。
oriLatitude	Double	原点纬度，计算经纬度需要参数。
gridSize	Int	栅格大小，计算经纬度需要参数。

 说明

由于Geoid由栅格坐标生成，坐标为栅格中心点，则计算出的经纬度是栅格中心点经纬度，与生成该Geoid的经纬度可能有[0度~半个栅格度数]的误差。

UDF输出参数：

参数	类型	说明
latitudeAndLongitude	Array[Double]	返回该geoid所表示的栅格的中心点的经纬度坐标，以数组的方式返回，第一位为latitude，第二位为longitude。

使用示例：

```
select longitude, latitude, mygeohash, GeoidToLatLng(mygeohash, 39.832277, 50) as LatitudeAndLongitude from geoTable;
```

- 计算金字塔模型向上汇聚一层的Geoid。

ToUpperLayerGeoid(geoid)

UDF输入参数：

参数	类型	说明
geoid	Long	根据输入Geoid计算金字塔模型上一层Geoid。

UDF输出参数：

参数	类型	说明
geoid	Long	金字塔模型上一层Geoid。

使用示例：

```
select longitude, latitude, mygeohash, ToUpperLayerGeoid(mygeohash) as upperLayerGeoid from geoTable;
```

- 输入polygon获得Geoid范围列表。

ToRangeList(polygon, oriLatitude, gridSize)

UDF输入参数：

参数	类型	说明
polygon	String	输入polygon字符串，用一组经纬度表示。 经纬度间用空格分隔，每对经纬度间用逗号分隔，首尾经纬度一致。
oriLatitude	Double	原点纬度，计算Geoid需要参数。
gridSize	Int	栅格大小，计算Geoid需要参数。

UDF输出参数:

参数	类型	说明
geoidList	Buffer[Array[Long]]	将polygon转换为一串geoid的范围列表。

使用示例:

```
select ToRangeList('116.321011 40.123503, 116.137676 39.947911, 116.560993 39.935276, 116.321011 40.123503', 39.832277, 50) as rangeList from geoTable;
```

- 计算金字塔模型向上汇聚一层的longitude。

ToUpperLongitude (longitude, gridSize, oriLat)

UDF输入参数:

参数	类型	说明
longitude	Long	输入longitude，用一个长整型表示。
gridSize	Int	栅格大小，计算longitude需要参数。
oriLatitude	Double	原点纬度，计算longitude需要参数。

UDF输出参数:

参数	类型	说明
longitude	Long	返回上一层的longitude。

使用示例:

```
select ToUpperLongitude (-23575161504L, 50, 39.832277) as upperLongitude from geoTable;
```

- 计算金字塔模型向上汇聚一层的Latitude。

ToUpperLatitude(Latitude, gridSize, oriLat)

UDF输入参数:

参数	类型	说明
latitude	Long	输入latitude，用一个长整型表示。
gridSize	Int	栅格大小，计算latitude需要参数。
oriLatitude	Double	原点纬度，计算latitude需要参数。

UDF输出参数:

参数	类型	说明
Latitude	Long	返回上一层的latitude。

使用示例：

```
select ToUpperLatitude (-23575161504L, 50, 39.832277) as upperLatitude from geoTable;
```

- 经纬度转GeoSOT

LatLngToGridCode(latitude, longitude, level)

UDF输入参数：

参数	类型	说明
latitude	Double	输入latitude。
longitude	Double	输入longitude。
level	Int	输入level，值区间[0-32]。

UDF输出参数：

参数	类型	说明
geold	Long	通过GeoSOT编码获得一个表示经纬度的数。

使用示例：

```
select LatLngToGridCode(39.930753, 116.302895, 21) as geold;
```

3.8 CarbonData 故障处理

3.8.1 当在 Filter 中使用 Big Double 类型数值时，过滤结果与 Hive 不一致

现象描述

当在filter中使用更高精度的double数据类型的数值时，过滤结果没有按照所使用的filter的要求返回正确的值。

可能原因

如果filter使用更高精度的double数据类型的数值，系统将会对该值四舍五入进行比较，因此在这种情况下，即使小数部分不同，系统仍然会认为double数据类型的值是相同的。

定位思路

无。

处理步骤

当需要高精度的数据比较时，可以使用Decimal数据类型的数值，例如，在财务应用程序中，equality和inequality检查，以及取整运算，均可使用Decimal数据类型的数值。

参考信息

无。

3.8.2 executor 内存不足导致查询性能下降

现象描述

在不同的查询周期内运行查询功能，查询性能会有起伏。

可能原因

在处理数据加载时，为每个executor程序实例配置的内存不足，可能会产生更多的Java GC（垃圾收集）。当GC发生时，会发现查询性能下降。

定位思路

在Spark UI上，会发现某些executors的GC时间明显比其他executors高，或者所有的executors都表现出高GC时间。

处理步骤

登录Manager页面，选择“集群 > 服务 > Spark2x > 配置 > 全部配置”，在搜索框搜索“spark.executor.memory”，通过参数“spark.executor.memory”配置更高的内存值。

spark.executor.memory

4G

参考信息

无。

3.9 CarbonData 常见问题

3.9.1 为什么对 decimal 数据类型进行带过滤条件的查询时会出现异常输出？

问题

当对 decimal 数据类型进行带过滤条件的查询时，输出结果不正确。

例如，

```
select * from carbon_table where num = 1234567890123456.22;
```

输出结果：

```
+-----+-----+-----+-----+
| name |      num      |
+-----+-----+-----+-----+
| IAA | 1234567890123456.22 |
| IAA | 1234567890123456.21 |
+-----+-----+-----+-----+
```

回答

为了得到准确的输出结果，需在数字后面加上“BD”。

例如，

```
select * from carbon_table where num = 1234567890123456.22BD;
```

输出结果：

```
+-----+-----+-----+-----+
| name |      num      |
+-----+-----+-----+-----+
| IAA | 1234567890123456.22 |
+-----+-----+-----+-----+
```

3.9.2 如何避免对历史数据进行 minor compaction？

问题

如何避免对历史数据进行 minor compaction？

回答

如果要先加载历史数据，后加载增量数据，则以下步骤可避免对历史数据进行 minor compaction：

1. 加载所有历史数据。
2. 将 major compaction 大小配置为小于历史数据 segment 大小的值。
3. 对历史数据进行一次 major compaction，之后将不会考虑这些 segments 进行 minor compaction。
4. 加载增量数据。
5. 用户可以根据自己的需要配置 minor compaction 阈值。

配置示例和预期输出：

1. 用户将所有历史数据加载到 CarbonData，此数据的一个 segment 的大小假定为 500GB。

2. 用户设置major compaction参数的阈值：“carbon.major.compaction.size” = “491520（480gb * 1024）”。其中，491520可配置。
3. 运行major compaction。由于每个segment的大小超过配置值的大小，因此这些segments将会被压缩。
4. 加载增量负载。
5. 配置minor compaction参数的阈值：“compaction.level.threshold” = “6,6”。
6. 运行minor compaction。此时只考虑增量负载。

3.9.3 如何在 CarbonData 数据加载时修改默认的组名？

问题

如何在CarbonData数据加载时修改默认的组名？

回答

CarbonData数据加载时，默认的组名为“ficommon”。可以根据需要修改默认的组名。

1. 编辑“carbon.properties”文件。
2. 根据需要修改关键字“carbon.dataload.group.name”的值。其默认值为“ficommon”。

3.9.4 为什么 INSERT INTO CARBON TABLE 失败？

问题

为什么 *INSERT INTO CARBON TABLE* 命令无法在日志文件中记录以下信息？

```
Data load failed due to bad record
```

回答

在以下场景中，*INSERT INTO CARBON TABLE* 命令会失败：

- 当源表和目标表的列数据类型不同时，源表中的数据将被视为Bad Records，则 *INSERT INTO* 命令会失败。
- 源列上的aggregation函数的结果超过目标列的最大范围，则 *INSERT INTO* 命令会失败。

解决方法：

在进行插入操作时，可在对应的列上使用cast函数。

示例：

- a. 使用DESCRIBE命令查询目标表和源表。

```
DESCRIBE newcarbontable;
```

结果：

```
col1 int  
col2 bigint
```

```
DESCRIBE sourcetable;
```

结果：

```
col1 int  
col2 int
```

- b. 添加cast函数以将BigInt类型数据转换为Integer类型数据。

```
INSERT INTO newcarbontable select col1, cast(col2 as integer) from  
sourcetable;
```

3.9.5 为什么含转义字符的输入数据记录到 Bad Records 中的值与原始数据不同？

问题

为什么含转义字符的输入数据记录到Bad Records中的值与原始数据不同？

回答

转义字符以反斜线“\”开头，后跟一个或几个字符。如果输入记录包含类似\t, \b, \n, \r, \f, \', \", \\的转义字符，Java将把转义符\和它后面的字符一起处理得到转义后的值。

例如：如果CSV数据类似“2010\\10,test”，将这两列插入“String,int”类型时，因为“test”无法转换为int类型，表会将这条记录重定向到Bad Records中。但记录到Bad Records中的值为“2010\10”，Java会将原始数据中的“\\”转义为“\”。

3.9.6 当初始 Executor 为 0 时，为什么 INSERT INTO/LOAD DATA 任务分配不正确，打开的 task 少于可用的 Executor？

问题

当初始Executor为0时，为什么INSERT INTO/LOAD DATA任务分配不正确，打开的task少于可用的Executor？

回答

在这种场景下，CarbonData会给每个节点分配一个INSERT INTO或LOAD DATA任务。如果Executor不是不同的节点分配的，CarbonData将会启动较少的task。

解决措施：

您可以适当增大Executor内存和Executor核数，以便YARN可以在每个节点上启动一个Executor。具体的配置方法如下：

1. 配置Executor核数。
 - 将“spark-defaults.conf”中的“spark.executor.cores”配置项或者“spark-env.sh”中的“SPARK_EXECUTOR_CORES”配置项设置为合适大小。
 - 在使用spark-submit命令时，添加“--executor-cores NUM”参数设置核数。
2. 配置Executor内存。
 - 将“spark-defaults.conf”中的“spark.executor.memory”配置项或者“spark-env.sh”中的“SPARK_EXECUTOR_MEMORY”配置项设置为合适大小。

- 在使用spark-submit命令时，添加“--executor-memory MEM”参数设置内存。

3.9.7 为什么并行度大于待处理的 block 数目时，CarbonData 仍需要额外的 executor？

问题

为什么并行度大于待处理的block数目时，CarbonData仍需要额外的executor？

回答

CarbonData块分布对于数据处理进行了如下优化：

1. 优化数据处理并行度。
2. 优化了读取块数据的并行性。

为了优化并行数据处理及并行读取块数据，CarbonData根据块的局域性申请 executor，因此CarbonData可获得所有节点上的executor。

为了优化并行数据处理及并行读取块数据，运用动态分配的用户需配置以下特性。

1. 使用参数“spark.dynamicAllocation.executorIdleTimeout”并将此参数值设置为15min（或平均查询时间）。
2. 正确配置参数“spark.dynamicAllocation.maxExecutors”，不推荐使用默认值（2048），否则CarbonData将申请最大数量的executor。
3. 对于更大的集群，配置参数“carbon.dynamicAllocation.schedulerTimeout”为10~15sec，默认值为5sec。
4. 配置参数“carbon.scheduler.minRegisteredResourcesRatio”为0.1~1.0，默认值为0.8。只要达到此参数值，块分布可启动。

3.9.8 为什么在 off heap 时数据加载失败？

问题

为什么在off heap时数据加载失败？

回答

YARN Resource Manager将（Java堆内存 + “spark.yarn.am.memoryOverhead”）作为内存限制。

因此在off heap时，内存可能会超出此限制。

您需配置参数“spark.yarn.am.memoryOverhead”以增加memory。

3.9.9 为什么创建 Hive 表失败？

问题

为什么创建Hive表失败？

回答

当源表或子查询具有大数据量的Partition时，创建Hive表失败。执行查询需要很多的task，此时输出的文件数就会很多，从而导致driver OOM。

可以在创建Hive表的语句中增加***distribute by***子句来解决这个问题，其中***distribute by***的字段要选取合适的cardinality（即distinct值的个数）。

distribute by子句限制了Hive表的Partition数量。增加***distribute by***子句后，最终的输出文件数取决于指定列的cardinality和“spark.sql.shuffle.partitions”参数值。但如果***distribute by***的字段的cardinality值很小，例如，“spark.sql.shuffle.partitions”参数值为200，但***distribute by***字段的cardinality只有100，则输出的200个文件中，只有其中100个文件有数据，剩下的100个文件为空文件。也就是说，如果选取的字段的cardinality过低，如1，则会造成严重的数据倾斜，从而严重影响查询性能。

因此，建议选取的***distribute by***字段的cardinality个数要大于“spark.sql.shuffle.partitions”参数，可大于2~3倍。

示例：

```
create table hivetable1 as select * from sourcetable1 distribute by col_age;
```

3.9.10 如何在不同的 namespaces 上逻辑地分割数据

问题

如何在不同的namespaces上逻辑地分割数据？

回答

- 配置：
要在不同namespaces之间逻辑地分割数据，必须更新HDFS，Hive和Spark的“core-site.xml”文件中的以下配置。

说明

改变Hive组件将改变carbonstore的位置和warehouse的位置。

- HDFS中的配置

- fs.defaultFS - 默认文件系统的名称。URI模式必须设置为“viewfs”。当使用“viewfs”模式时，权限部分必须是“ClusterX”。
- fs.viewfs.mountable.ClusterX.homedir - 主目录基本路径。每个用户都可以使用在“FileSystem/FileContext”中定义的getHomeDirectory()方法访问其主目录。
- fs.viewfs.mountable.default.link.<dir_name> - ViewFS安装表。

示例：

```
<property>  
<name>fs.defaultFS</name>  
<value>viewfs://ClusterX</value>  
</property>  
<property>  
<name>fs.viewfs.mountable.ClusterX.link./folder1</name>  
<value>hdfs://NS1/folder1</value>  
</property>  
</property>
```

```
<name>fs.viewfs.mounttable.ClusterX.link./folder2</name>  
<value>hdfs://NS2/folder2</value>  
</property>
```

- Hive和Spark中的配置

fs.defaultFS - 默认文件系统的名称。URI模式必须设置为“viewfs”。当使用“viewfs”模式时，权限部分必须是“ClusterX”。

- 命令格式:

```
LOAD DATA INPATH 'path to data' INTO TABLE table_name OPTIONS ('...');
```

📖 说明

每当Spark配置有viewFS文件系统时，当尝试从HDFS加载数据时，用户必须在LOAD语句中指定如“viewfs://”这样的路径或相对路径作为文件路径。

- 示例:

- viewFS路径举例:

```
LOAD DATA INPATH 'viewfs://ClusterX/dir/data.csv' INTO TABLE  
table_name OPTIONS ('...');
```

- 相对路径举例:

```
LOAD DATA INPATH '/apps/input_data1.txt' INTO TABLE table_name;
```

3.9.11 为什么在 Spark Shell 中不能执行更新命令？

问题

为什么在Spark Shell中不能执行更新命令？

回答

本文档中给出的语法和示例是关于Beeline的命令，而不是Spark Shell中的命令。

若要在Spark Shell中使用更新命令，可以使用以下语法。

- 语法1

```
<carbon_context>.sql("UPDATE <CARBON TABLE> SET (column_name1,  
column_name2, ... column_name n) = (column1_expression ,  
column2_expression , column3_expression ... column n_expression)  
[ WHERE { <filter_condition> } ]");.show
```

- 语法2

```
<carbon_context>.sql("UPDATE <CARBON TABLE> SET (column_name1,  
column_name2,) = (select sourceColumn1, sourceColumn2 from  
sourceTable [ WHERE { <filter_condition> } ] ) [ WHERE  
{ <filter_condition> } ]");.show
```

示例:

如果CarbonData的context是“carbon”，那么更新命令如下:

```
carbon.sql("update carbonTable1 d set (d.column3,d.column5) = (select  
s.c33 ,s.c55 from sourceTable1 s where d.column1 = s.c11) where d.column1 =  
'country' exists( select * from table3 o where o.c2 > 1);".show
```

3.9.12 如何在 CarbonData 中配置非安全内存？

问题

如何在CarbonData中配置非安全内存？

回答

在Spark配置中，“spark.yarn.executor.memoryOverhead”参数的值应大于CarbonData配置参数“sort.inmemory.size.inmb”与“Netty offheapmemory required”参数值的总和，或者“carbon.unsafe.working.memory.in.mb”、“carbon.sort.inmemory.storage.size.in.mb”与“Netty offheapmemory required”参数值的总和。否则，如果堆外（off heap）访问超出配置的executor内存，则YARN可能会停止executor。

“Netty offheapmemory required”说明：当“spark.shuffle.io.preferDirectBufs”设为true时，Spark中netty 传输服务从"spark.yarn.executor.memoryOverhead"中拿掉部分堆内存[~ 384 MB or 0.1 x 执行器内存]。

详细信息请参考常见[配置Spark Executor堆内存参数](#)。

3.9.13 设置了 HDFS 存储目录的磁盘空间配额，CarbonData 为什么会发生异常？

问题

设置了HDFS存储目录的磁盘空间配额，CarbonData为什么会发生异常。

回答

创建、加载、更新表或进行其他操作时，数据会被写入HDFS。若HDFS目录的磁盘空间配额不足，则操作失败并发生以下异常。

```
org.apache.hadoop.hdfs.protocol.DSQuotaExceededException: The DiskSpace quota of /user/tenant is exceeded: quota = 314572800 B = 300 MB but diskspace consumed = 402653184 B = 384 MB at org.apache.hadoop.hdfs.server.namenode.DirectoryWithQuotaFeature.verifyStorageSpaceQuota(DirectoryWithQuotaFeature.java:211) at org.apache.hadoop.hdfs.server.namenode.DirectoryWithQuotaFeature.verifyQuota(DirectoryWithQuotaFeature.java:239) at org.apache.hadoop.hdfs.server.namenode.FSDirectory.verifyQuota(FSDirectory.java:941) at org.apache.hadoop.hdfs.server.namenode.FSDirectory.updateCount(FSDirectory.java:745)
```

若发生此异常，请为租户配置足够的磁盘空间配额。

例如：

需要的磁盘空间配置可以按照如下方法计算：

如果HDFS的副本数为3，HDFS默认的块大小为128MB，则最小需要384MB的磁盘空间用于写表的schema文件到HDFS上。计算公式： $\text{no. of block} \times \text{block_size} \times \text{replication_factor of the schema file} = 1 \times 128 \times 3 = 384 \text{ MB}$

说明

数据加载时，由于默认块大小为1024MB，每个fact文件需要的最小空间为3072MB。

3.9.14 为什么数据查询/加载失败，且发生“org.apache.carbondata.core.memory.MemoryException: Not enough memory”异常？

问题

为什么数据查询/加载失败，且发生“org.apache.carbondata.core.memory.MemoryException: Not enough memory”异常？

回答

当执行器中此次数据查询和加载所需要的堆外内存不足时，便会发生此异常。

在这种情况下，请增大“carbon.unsafe.working.memory.in.mb”和“spark.yarn.executor.memoryOverhead”的值。

详细信息请参考[如何在CarbonData中配置非安全内存？](#)

该内存被数据查询和加载共享。所以如果加载和查询需要同时进行，建议将“carbon.unsafe.working.memory.in.mb”和“spark.yarn.executor.memoryOverhead”的值配置为2048 MB以上。

可以使用以下公式进行估算：

数据加载所需内存：

$$\begin{aligned} & (\text{“carbon.number.of.cores.while.loading”的值[默认值 = 6]}) \times \text{并行加载数据的表格} \\ & \times (\text{“offheap.sort.chunk.size.inmb”的值[默认值 = 64 MB]} + \\ & \text{“carbon.blockletgroup.size.in.mb”的值[默认值 = 64 MB]} + \text{当前的压缩率}[64 MB / \\ & 3.5]) \end{aligned}$$

= ~900 MB 每表格

数据查询所需内存：

$$\begin{aligned} & (\text{SPARK_EXECUTOR_INSTANCES. [默认值 = 2]}) \times (\text{carbon.blockletgroup.size.in.mb} \\ & \text{[默认值 = 64 MB]} + \text{“carbon.blockletgroup.size.in.mb”解压内容[默认值 = 64 MB *} \\ & 3.5]) \times (\text{每个执行器核数[默认值 = 1]}) \end{aligned}$$

= ~ 600 MB

3.9.15 开启防误删后为什么 Carbon 表没有执行 drop 命令，回收站中也会存在该表的文件？

问题

开启防误删下，为什么Carbon表没有执行drop table命令，回收站中也会存在该表的文件？

回答

在Carbon适配防误删后，调用文件删除命令，会将删除的文件放入回收站中。

在insert、load等命令中会有中间文件.carbonindex文件的删除，所以在未执行drop table命令的时候，回收站中也可能会存在该表的文件。

如果这个时候再执行drop table命令，那么按照回收站机制，会生成一个带时间戳的该表目录，该目录中的文件是完整的。

4 使用 ClickHouse

4.1 ClickHouse 概述

ClickHouse 表引擎介绍

表引擎在ClickHouse中的作用十分关键，不同的表引擎决定了：

- 数据存储和读取的位置
- 支持哪些查询方式
- 能否并发式访问数据
- 能不能使用索引
- 是否可以执行多线程请求
- 数据复制使用的参数

其中MergeTree和Distributed是ClickHouse表引擎中最重要，也是最常使用的两个引擎，本文将重点进行介绍。

其他表引擎详细可以参考官网链接：<https://clickhouse.tech/docs/en/engines/table-engines>。

- **MergeTree系列引擎**

MergeTree用于高负载任务的最通用和功能最强大的表引擎，其主要有以下关键特征：

- 基于分区键（partitioning key）的数据分区分块存储
- 数据索引排序（基于primary key和order by）
- 支持数据复制（带Replicated前缀的表引擎）
- 支持数据抽样

在写入数据时，该系列引擎表会按照分区键将数据分成不同的文件夹，文件夹内每列数据为不同的独立文件，以及创建数据的序列化索引排序记录文件。该结构使得数据读取时能够减少数据检索时的数据量，极大的提高查询效率。

- MergeTree

建表语法：

```
CREATE TABLE [IF NOT EXISTS] [db.]table_name [ON CLUSTER cluster]  
(
```

```

name1 [type1] [DEFAULT|MATERIALIZED|ALIAS expr1] [TTL expr1],
name2 [type2] [DEFAULT|MATERIALIZED|ALIAS expr2] [TTL expr2],
...
INDEX index_name1 expr1 TYPE type1(...) GRANULARITY value1,
INDEX index_name2 expr2 TYPE type2(...) GRANULARITY value2
) ENGINE = MergeTree()
ORDER BY expr
[PARTITION BY expr]
[PRIMARY KEY expr]
[SAMPLE BY expr]
[TTL expr [DELETE|TO DISK 'xxx'|TO VOLUME 'xxx'], ...]
[SETTINGS name=value, ...]

```

使用示例:

```

CREATE TABLE default.test (
  name1 DateTime,
  name2 String,
  name3 String,
  name4 String,
  name5 Date,
  ...
) ENGINE = MergeTree()
PARTITION BY toYYYYMM(name5)
ORDER BY (name1, name2)
SETTINGS index_granularity = 8192

```

示例参数说明如下:

- **ENGINE = MergeTree():** MergeTree表引擎。
- **PARTITION BY toYYYYMM(name4):** 分区，示例数据将以月份为分区，每个月份一个文件夹。
- **ORDER BY:** 排序字段，支持多字段的索引排序，第一个相同的时候按照第二个排序依次类推。
- **index_granularity = 8192:** 排序索引的颗粒度，每8192条数据记录一个排序索引值。

如果被查询的数据存在于分区或排序字段中，能极大降低数据查找时间。

- ReplacingMergeTree

该引擎和MergeTree的不同之处在于它会删除排序键值相同的重复项。ReplacingMergeTree适合于清除重复数据节省存储空间，但是它不保证重复数据不出现，一般不建议使用。

建表语法:

```

CREATE TABLE [IF NOT EXISTS] [db.]table_name [ON CLUSTER cluster]
(
  name1 [type1] [DEFAULT|MATERIALIZED|ALIAS expr1],
  name2 [type2] [DEFAULT|MATERIALIZED|ALIAS expr2],
  ...
) ENGINE = ReplacingMergeTree([ver])
[PARTITION BY expr]
[ORDER BY expr]
[SAMPLE BY expr]
[SETTINGS name=value, ...]

```

- SummingMergeTree

当合并SummingMergeTree表的数据片段时，ClickHouse会把所有具有相同主键的行合并为一行，该行包含了被合并的行中具有数值数据类型的列的汇总值。如果主键的组合方式使得单个键值对应于大量的行，则可以显著减少存储空间并加快数据查询的速度。

建表语法:

```
CREATE TABLE [IF NOT EXISTS] [db.]table_name [ON CLUSTER cluster]
(
  name1 [type1] [DEFAULT|MATERIALIZED|ALIAS expr1],
  name2 [type2] [DEFAULT|MATERIALIZED|ALIAS expr2],
  ...
) ENGINE = SummingMergeTree([columns])
[PARTITION BY expr]
[ORDER BY expr]
[SAMPLE BY expr]
[SETTINGS name=value, ...]
```

使用示例：

创建一个SummingMergeTree表testTable：

```
CREATE TABLE testTable
(
  id UInt32,
  value UInt32
)
ENGINE = SummingMergeTree()
ORDER BY id
```

插入表数据：

```
INSERT INTO testTable Values(5,9),(5,3),(4,6),(1,2),(2,5),(1,4),(3,8);
INSERT INTO testTable Values(88,5),(5,5),(3,7),(3,5),(1,6),(2,6),(4,7),(4,6),(43,5),(5,9),(3,6);
```

在未合并parts查询所有数据：

```
SELECT * FROM testTable
```

id	value
1	6
2	5
3	8
4	6
5	12

id	value
1	6
2	6
3	18
4	13
5	14
43	5
88	5

ClickHouse还没有汇总所有行，如果需要通过ID进行汇总聚合，需要用到sum和GROUP BY子句：

```
SELECT id, sum(value) FROM testTable GROUP BY id
```

id	sum(value)
4	19
3	26
88	5
2	11
5	26
1	12
43	5

手工执行合并操作：

```
OPTIMIZE TABLE testTable
```

此时再查询testTable表数据：

```
SELECT * FROM testTable
```

id	value
1	12
2	11
3	26
4	19
5	26

43	5
88	5

SummingMergeTree根据ORDER BY排序键作为聚合数据的条件Key。即如果排序key是相同的，则会合并成一条数据，并对指定的合并字段进行聚合。

后台执行合并操作时才会进行数据的预先聚合，而合并操作的执行时机无法预测，所以可能存在部分数据已经被预先聚合、部分数据尚未被聚合的情况。因此，在执行聚合计算时，SQL中仍需要使用GROUP BY子句。

- AggregatingMergeTree

AggregatingMergeTree是预先聚合引擎的一种，用于提升聚合计算的性能。AggregatingMergeTree引擎能够在合并分区时，按照预先定义的条件聚合数据，同时根据预先定义的聚合函数计算数据并通过二进制的格式存入表内。

建表语法：

```
CREATE TABLE [IF NOT EXISTS] [db.]table_name [ON CLUSTER cluster]
(
    name1 [type1] [DEFAULT|MATERIALIZED|ALIAS expr1],
    name2 [type2] [DEFAULT|MATERIALIZED|ALIAS expr2],
    ...
) ENGINE = AggregatingMergeTree()
[PARTITION BY expr]
[ORDER BY expr]
[SAMPLE BY expr]
[TTL expr]
[SETTINGS name=value, ...]
```

使用示例：

AggregatingMergeTree无单独参数设置，在分区合并时，在每个数据分区内，会按照ORDER BY聚合，使用何种聚合函数，对哪些列字段计算，则是通过定义AggregateFunction函数类型实现，例如：

```
create table test_table (
    name1 String,
    name2 String,
    name3 AggregateFunction(uniq,String),
    name4 AggregateFunction(sum,Int),
    name5 DateTime
) ENGINE = AggregatingMergeTree()
PARTITION BY toYYYYMM(name5)
ORDER BY (name1,name2)
PRIMARY KEY name1;
```

AggregateFunction类型的数据在写入和查询时需要分别调用*state、*merge函数，*表示定义字段类型时使用的聚合函数。如上示例表test_table定义的名称3、名称4字段分别使用了uniq、sum函数，那么在写入数据时需要调用uniqState、sumState函数，并使用INSERT SELECT语法。

```
insert into test_table select '8','test1',uniqState('name1'),sumState(toInt32(100)),2021-04-30 17:18:00';
insert into test_table select '8','test1',uniqState('name1'),sumState(toInt32(200)),2021-04-30 17:18:00';
```

在查询数据时也需要调用对应的函数uniqMerge、sumMerge：

```
select name1,name2,uniqMerge(name3),sumMerge(name4) from test_table group by name1,name2;
```

name1	name2	uniqMerge(name3)	sumMerge(name4)
8	test1	1	300

AggregatingMergeTree更常用的方式是结合物化视图使用，物化视图即其它数据表上层的一种查询视图。详细可以参考：<https://clickhouse.tech/docs/en/engines/table-engines/mergetree-family/aggregatingmergetree/>

- CollapsingMergeTree

CollapsingMergeTree它通过定义一个sign标记位字段记录数据行的状态。如果sign标记为1，则表示这是一行有效的数据；如果sign标记为-1，则表示这行数据需要被删除。

建表语法：

```
CREATE TABLE [IF NOT EXISTS] [db.]table_name [ON CLUSTER cluster]
(
    name1 [type1] [DEFAULT|MATERIALIZED|ALIAS expr1],
    name2 [type2] [DEFAULT|MATERIALIZED|ALIAS expr2],
    ...
) ENGINE = CollapsingMergeTree(sign)
[PARTITION BY expr]
[ORDER BY expr]
[SAMPLE BY expr]
[SETTINGS name=value, ...]
```

使用示例：

具体的使用示例可以参考：<https://clickhouse.tech/docs/en/engines/table-engines/mergetree-family/collapsingmergetree/>。

- VersionedCollapsingMergeTree

VersionedCollapsingMergeTree表引擎在建表语句中新增了一列version，用于在乱序情况下记录状态行与取消行的对应关系。主键相同，且Version相同、Sign相反的行，在Compaction时会被删除。

建表语法：

```
CREATE TABLE [IF NOT EXISTS] [db.]table_name [ON CLUSTER cluster]
(
    name1 [type1] [DEFAULT|MATERIALIZED|ALIAS expr1],
    name2 [type2] [DEFAULT|MATERIALIZED|ALIAS expr2],
    ...
) ENGINE = VersionedCollapsingMergeTree(sign, version)
[PARTITION BY expr]
[ORDER BY expr]
[SAMPLE BY expr]
[SETTINGS name=value, ...]
```

使用示例：

具体的使用示例可以参考：<https://clickhouse.tech/docs/en/engines/table-engines/mergetree-family/versionedcollapsingmergetree/>。

- GraphiteMergeTree

GraphiteMergeTree引擎用来存储时序数据库Graphite的数据。

建表语法：

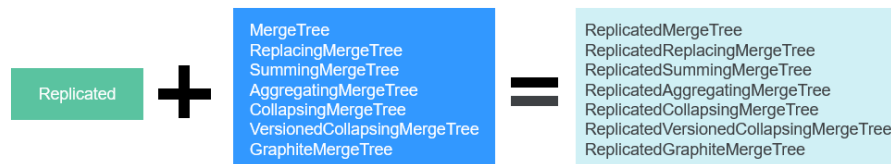
```
CREATE TABLE [IF NOT EXISTS] [db.]table_name [ON CLUSTER cluster]
(
    Path String,
    Time DateTime,
    Value <Numeric_type>,
    Version <Numeric_type>
    ...
) ENGINE = GraphiteMergeTree(config_section)
[PARTITION BY expr]
[ORDER BY expr]
[SAMPLE BY expr]
[SETTINGS name=value, ...]
```

使用示例：

具体的使用示例可以参考：<https://clickhouse.tech/docs/en/engines/table-engines/mergetree-family/graphitemergetree/>。

- Replicated*MergeTree引擎

ClickHouse中的所有MergeTree家族引擎前面加上Replicated就成了支持副本的合并树引擎。



Replicated系列引擎借助ZooKeeper实现数据的同步，创建Replicated复制表时通过注册到ZooKeeper上的信息实现同一个分片的所有副本数据进行同步。

Replicated表引擎的创建模板：

```
ENGINE = Replicated*MergeTree('ZooKeeper存储路径',副本名称, ...)
```

Replicated表引擎需指定两个参数：

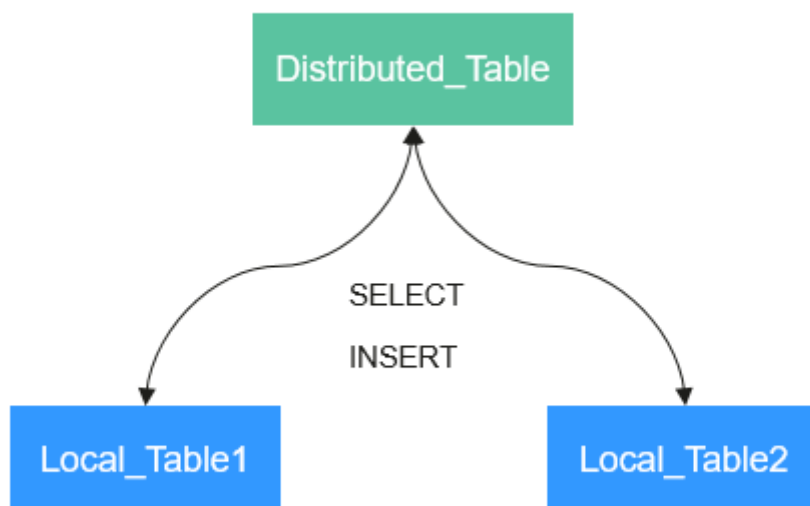
- ZooKeeper存储路径：ZooKeeper中该表相关数据的存储路径，建议规范化，如：`/clickhouse/tables/{shard}/数据库名/表名`。
- 副本名称，一般用`{replica}`即可。

Replicated表引擎使用示例可以参考：[ClickHouse表创建](#)。

- **Distributed表引擎**

Distributed表引擎本身不存储任何数据，而是作为数据分片的透明代理，能够自动路由数据到集群中的各个节点，分布式表需要和其他本地数据表一起协同工作。分布式表会将接收到的读写任务分发到各个本地表，而实际上数据的存储在各个节点的本地表中。

图 4-1 Distributed



Distributed表引擎的创建模板：

```
ENGINE = Distributed(cluster_name, database_name, table_name, [sharding_key])
```


Distributed表参数解析如下：

- cluster_name：集群名称，在对分布式表执行读写的过程中，使用集群的配置信息查找对应的ClickHouse实例节点。
- database_name：数据库名称。
- table_name：数据库下对应的本地表名称，用于将分布式表映射到本地表上。
- sharding_key：分片键（可选参数），分布式表会按照这个规则，将数据分发到各个本地表中。

Distributed表引擎使用示例：

```
--先创建一个表名为test的ReplicatedMergeTree本地表
CREATE TABLE default.test ON CLUSTER default_cluster_1
(
    `EventDate` DateTime,
    `id` UInt64
)
ENGINE = ReplicatedMergeTree('/clickhouse/tables/{shard}/default/test', '{replica}')
PARTITION BY toYYYYMM(EventDate)
ORDER BY id

--基于本地表test创建表名为test_all的Distributed表
CREATE TABLE default.test_all ON CLUSTER default_cluster_1
(
    `EventDate` DateTime,
    `id` UInt64
)
ENGINE = Distributed(default_cluster_1, default, test, rand())
```

分布式表创建规则：

- 创建Distributed表时需加上**on cluster cluster_name**，这样建表语句在某一个ClickHouse实例上执行一次即可分发到集群中所有实例上执行。
- 分布式表通常以本地表加“_all”命名。它与本地表形成一对多的映射关系，之后可以通过分布式表代理操作多张本地表。
- 分布式表的表结构尽量和本地表的结构一致。如果不一致，在建表时不会报错，但在查询或者插入时可能会发生异常。

ClickHouse 数据类型说明

MRS的ClickHouse服务数据类型如[表4-1](#)所示。

ClickHouse完整数据类型介绍，请参考[开源官方数据类型介绍](#)。

表 4-1 ClickHouse 数据类型

分类	关键字	数据类型	描述
数据类型	Int8	Int8	取值范围：【 -128, 127 】
	Int16	Int16	取值范围：【 -32768, 32767 】
	Int32	Int32	取值范围：【 -2147483648, 2147483647 】
	Int64	Int64	取值范围：【 -9223372036854775808, 9223372036854775807 】

分类	关键字	数据类型	描述
浮点类型	Float32	单精度浮点数	同C语言Float类型，单精度浮点数在机内占4个字节，用32位二进制描述。
	Float64	双精度浮点数	同C语言Double类型，双精度浮点数在机内占8个字节，用64位二进制描述。
Decimal类型	Decimal	Decimal	<p>有符号的定点数，可在加、减和乘法运算过程中保持精度。支持几种写法：</p> <ul style="list-style-type: none"> • Decimal(P, S) • Decimal32(S) • Decimal64(S) • Decimal128(S) <p>说明</p> <ul style="list-style-type: none"> • P: 精度，有效范围：[1:38]，决定可以有多少个十进制数字（包括分数）。 • S: 规模，有效范围：[0: P]，决定数字的小数部分中包含的小数位。
字符串类型	String	字符串	字符串可以是任意长度的。它可以包含任意的字节集，包含空字节。因此，字符串类型可以代替其他 DBMSs 中的VARCHAR、BLOB、CLOB 等类型。
	FixedString	固定字符串	<p>当数据的长度恰好为N个字节时，FixedString类型是高效的。在其他情况下，这可能会降低效率。可以有效存储在FixedString类型的列中的值的示例：</p> <ul style="list-style-type: none"> • 二进制表示的IP地址 • 语言代码（ru_RU, en_US ...） • 货币代码（RUB ...） • 二进制表示的哈希值（MD5使用FixedString（16），SHA256使用FixedString（32））
时间日期类型	Date	日期	用两个字节存储，表示从1970-01-01（无符号）到当前的日期值。日期中没有存储时区信息。

分类	关键字	数据类型	描述
	DateTime	时间戳	用四个字节（无符号的）存储 Unix 时间戳。允许存储与日期类型相同的范围内的值。最小值为 1970-01-01 00:00:00。时间戳类型值精确到秒（没有闰秒）。时区使用启动客户端或服务时的系统时区。
	DateTime64	DateTime64	此类型允许以日期（date）加时间（time）的形式来存储一个时刻的时间值。
布尔型	Boolean	Boolean	ClickHouse没有单独的类型来存储布尔值。可以使用UInt8 类型，取值限制为0或1。
数组类型	Array	Array	Array(T)，由 T 类型元素组成的数组。T 可以是任意类型，包含数组类型。但不推荐使用多维数组，ClickHouse对多维数组的支持有限。例如，不能在 MergeTree表中存储多维数组。
元组类型	Tuple	Tuple	Tuple(T1, T2, ...)，元组，其中每个元素都有单独的类型，不能在表中存储元组（除了内存表）。它们可以用于临时列分组。在查询中，IN表达式和带特定参数的 lambda 函数可以用来对临时列进行分组。
Domains数据类型	Domains	Domains	Domains类型是特定实现的类型： IPv4是与UInt32类型保持二进制兼容的Domains类型，用于存储IPv4地址的值。它提供了更为紧凑的二进制存储的同时支持识别可读性更加友好的输入输出格式。
枚举类型	Enum8	Enum8	取值范围：【-128, 127】 Enum 保存 'string'= integer 的对应关系，例如： Enum8('hello' = 1, 'world' = 2)
	Enum16	Enum16	取值范围：【-32768, 32767】

分类	关键字	数据类型	描述
可为空	Nullable	Nullable	除非在ClickHouse服务器配置中另有说明，否则NULL是任何Nullable类型的默认值。Nullable类型字段不能包含在表索引中。 可以与TypeName的正常值存放一起。例如，Nullable(Int8) 类型的列可以存储 Int8 类型值，而没有值的行将存储 NULL。
嵌套类型	nested	nested	嵌套的数据结构就像单元格内的表格。嵌套数据结构的参数（列名和类型）的指定方式与CREATE TABLE查询中的指定方式相同。每个表行都可以对应于嵌套数据结构中的任意数量的行。 示例: Nested(Name1 Type1, Name2 Type2, ...)

4.2 ClickHouse 用户权限管理

4.2.1 ClickHouse 用户及权限管理

用户权限模型

ClickHouse用户权限管理实现了对集群中各个ClickHouse实例上用户、角色、权限的统一管理。通过Manager UI的权限管理模块进行创建用户、创建角色、绑定ClickHouse访问权限配置等操作，通过用户绑定角色的方式，实现用户权限控制。

管理资源：Clickhouse权限管理支持的资源如表4-2所示。

资源权限：ClickHouse支持的资源权限如表4-3所示。

说明

本章节不支持MRS 3.1.0普通模式（未开启Kerberos认证）集群。

表 4-2 ClickHouse 支持的权限管理对象

资源列表	是否集成	备注
数据库	是（一级）	-
表	是（二级）	-
视图	是（二级）	与表一致

表 4-3 资源权限列表

资源对象	可选权限	备注
数据库（DATABASE）	CREATE	CREATE DATABASE/TABLE/ VIEW/DICTIONARY权限
表/视图（TABLE/ VIEW）	SELECT/INSERT	-

前提条件

- ClickHouse服务运行正常，ZooKeeper服务运行正常。
- 用户在集群中创建数据库或者表时需使用ON CLUSTER语句，保证各个ClickHouse节点上数据库、表的元信息相同。

说明

ClickHouse赋权成功后，权限生效时间大约为1分钟。

添加 ClickHouse 角色

步骤1 登录Manager，选择“系统 > 权限 > 角色”，在“角色”界面单击“添加角色”按钮，进入添加角色页面。

步骤2 在添加角色界面输入“角色名称”，在配置资源权限处单击集群名称，进入服务列表页面，单击ClickHouse服务，进入ClickHouse权限资源页面。

根据业务需求确定是否要创建具有ClickHouse管理员权限的角色。

说明

- ClickHouse管理员权限为：除去对user/role的创建、删除和修改之外的所有数据库操作权限。
- 对于用户和角色的管理，仅有ClickHouse的内置用户clickhouse具有权限。
- 是，执行[步骤3](#)。
- 否，执行[步骤4](#)。

角色 > 添加角色

角色名称:

配置资源权限: 所有资源 > B0118 > ClickHouse

视图名称

Clickhouse管理员权限

[Clickhouse Scope](#)

描述:

步骤3 勾选“ClickHouse管理员权限”，单击“确定”操作结束。

步骤4 单击“ClickHouse Scope”，进入ClickHouse数据库资源列表。勾选“创建”权限，则该角色将拥有该数据库下的创建（CREATE）权限。

角色名称:

配置资源权限: 所有资源 > B0118 > ClickHouse > Clickhouse Scope

资源名称	资源类型	权限
_temporary_and_external_tables	数据库	<input checked="" type="checkbox"/> 创建 <input type="checkbox"/>
db1	数据库	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/>
db10	数据库	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/>
db2	数据库	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/>
db3	数据库	<input type="checkbox"/>
db4	数据库	<input type="checkbox"/>
db5	数据库	<input type="checkbox"/>
db6	数据库	<input type="checkbox"/>
db7	数据库	<input type="checkbox"/>
db8	数据库	<input type="checkbox"/>

根据业务需求确定是否赋权。

- 是，单击“确定”操作结束。
- 否，执行**步骤5**。

步骤5 单击“资源名称 > 待操作的数据库资源名称”，进入表、视图页面，根据业务需要，勾选“读”（SELECT权限）或者“写”（INSERT权限）权限，单击“确定”。

角色名称:

配置资源权限: 所有资源 > B0118 > ClickHouse > Clickhouse Scope > db2

资源名称	资源类型	权限
tb3	表	<input type="checkbox"/> 读 <input checked="" type="checkbox"/> 写
tb4	表	<input type="checkbox"/> 读 <input checked="" type="checkbox"/> 写

描述:

----结束

添加用户并将 ClickHouse 对应角色绑定到该用户

步骤1 登录Manager，选择“系统 > 权限 > 用户”，单击“添加用户”，进入添加用户页面。

步骤2 “用户类型”选择“人机”，在“密码”和“确认密码”参数设置该用户对应的密码。

说明

- 用户名：添加的用户名不能包含字符“-”，否则会导致认证失败。
- 密码：设置的密码不能携带“\$”、“.”、“#”特殊字符，否则会导致认证失败。

步骤3 在“角色”处单击“添加”，在弹框中选择具有ClickHouse权限的角色，单击“确定”添加到角色，单击“确定”完成操作。

The screenshot shows a web form for adding a user. The fields are: Username (testuser), User Type (Human selected), Password (masked), Confirm Password (masked), User Groups (Add, Remove All, Create New User Group), Main Group (dropdown), Roles (Add, Remove All, Create New Role) with 'testrole' selected, and Description (text area). At the bottom are '确定' (Confirm) and '取消' (Cancel) buttons.

步骤4 登录ClickHouse客户端安装节点，使用新添加的用户及设置的密码连接ClickHouse服务。

- 执行以下命令，切换到客户端安装目录。
cd /opt/客户端安装目录
- 执行以下命令配置环境变量。
source bigdata_env
- 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建ClickHouse表的权限，具体请参见[添加ClickHouse角色](#)，为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行本步骤。
如果是MRS 3.1.0版本集群，则需要先执行：**export CLICKHOUSE_SECURITY_ENABLED=true**
kinit 步骤1中添加的用户
- 使用新添加的用户登录验证。
当前集群未启用Kerberos认证：
clickhouse client --host ClickHouse的实例IP --multiline --port ClickHouse的端口号 --secure

当前集群已启用Kerberos认证:

```
clickhouse client --host ClickHouse的实例IP --user 用户名 --password --port 9440 --secure
```

输入用户密码

📖 说明

普通模式的用户为默认的default用户，或者使用 ClickHouse社区开源能力添加管理用户。不能使用在FusionInsight Manager页面创建的用户。

---结束

异常场景下登录客户端操作赋权

ClickHouse集群默认每个节点上的表元信息是相同的，因此在Manager的权限管理页面上默认采集的是任意ClickHouse节点的表信息，如果有个别节点上创建DATABASE/TABLE时未使用ON CLUSTER语句，则权限操作可能无法展示该资源，不保证可以对其赋权。对于这样单个ClickHouse节点中的本地表，如果需要赋权，可以通过后台客户端进行操作。

📖 说明

以下操作，需要提前获取到需要赋权的角色、数据库或表名称、对应的ClickHouseServer实例所在的节点IP和系统域名。

- ClickHouseServer的实例IP地址可登录集群FusionInsight Manager，然后选择“集群 > 服务 > ClickHouse > 实例”，获取ClickHouseServer实例对应的业务IP地址。
- 系统域名：默认为hadoop.com。可登录集群FusionInsight Manager，单击“系统 > 权限 > 域和互信”，“本端域”参数值即为系统域名。在执行命令时改为小写。

步骤1 以root用户登录ClickHouseServer实例所在的节点。

步骤2 执行以下命令获取“clickhouse.keytab”文件路径。

```
ls ${BIGDATA_HOME}/FusionInsight_ClickHouse_*/install/FusionInsight-ClickHouse-*/clickhouse/keytab/clickhouse.keytab
```

步骤3 以客户端安装用户，登录安装客户端的节点。

步骤4 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤5 执行以下命令配置环境变量。

```
source bigdata_env
```

如果是MRS 3.1.0版本集群并且集群已启用Kerberos认证，则还需要执行：

```
export CLICKHOUSE_SECURITY_ENABLED=true
```

步骤6 执行如下命令使用客户端命令连接ClickHouseServer实例。

如果当前集群已启用Kerberos认证，执行以下命令：

```
clickhouse client --host ClickHouseServer实例所在节点IP --user clickhouse/hadoop.<系统域名> --password 步骤2中获取的clickhouse.keytab路径 --port ClickHouse的端口号 --secure
```

如果当前集群未启用Kerberos认证，执行以下命令：


```
clickhouse client --host ClickHouseServer实例所在节点IP --user clickhouse --port  
ClickHouse的端口号
```

步骤7 对某DATABASE进行赋权操作，执行如下命令。

授权操作语法，其中DATABASE为要操作的数据库名称，role为需要操作的角色。

```
GRANT [ON CLUSTER cluster_name] privilege ON {DATABASE/TABLE} TO {user /  
role}
```

例如，给用户testuser授予数据库t2的CREATE权限：

```
GRANT CREATE ON t2 to testuser;
```

步骤8 对TABLE/VIEW进行赋权操作，执行如下命令，其中TABLE为要操作的表或视图名称，user为需要操作的角色。

对某数据库下的表赋予查询权限：

```
GRANT SELECT ON TABLE TO user;
```

对某数据库下的表赋予写入权限：

```
GRANT INSERT ON TABLE TO user;
```

📖 说明

更多ClickHouse授权操作及详细权限说明可参考<https://clickhouse.tech/docs/zh/sql-reference/statements/grant/>。

步骤9 执行如下命令，退出客户端。

```
quit;
```

----结束

4.2.2 ClickHouse 使用 OpenLDAP 认证

ClickHouse支持和OpenLDAP进行对接，通过在ClickHouse上添加OpenLDAP服务器配置和创建用户，实现账号和权限的统一集中管理和权限控制等操作。此方案适合从OpenLDAP服务器中批量向ClickHouse中导入用户。

本章节操作仅支持MRS 3.1.0及以上集群版本。

前提条件

- MRS集群及ClickHouse实例运行正常，已安装ClickHouse客户端。
- OpenLDAP已安装且状态正常。

对接 OpenLDAP 服务器创建 ClickHouse 用户

步骤1 登录集群Manager页面，选择“集群 > 服务 > ClickHouse > 配置 > 全部配置”。

步骤2 参考下图图4-2，选择“ClickHouseServer（角色）> 自定义”，在“clickhouse-config-customize”配置项中添加如下OpenLDAP配置参数。

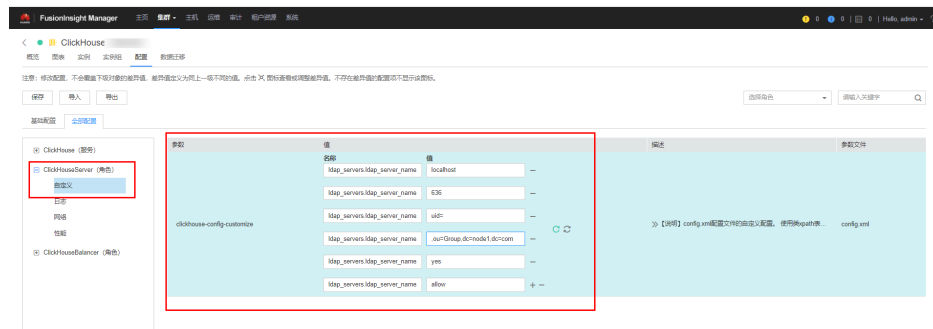
表 4-4 OpenLDAP 参数说明

参数名	参数值说明	参数取值参考
ldap_servers.ldap_server_name.host	OpenLDAP服务器主机名或IP，不能为空。	localhost
ldap_servers.ldap_server_name.port	OpenLDAP服务器端口。 如果enable_tls参数设置为true，则默认端口号为636，否则为389。	636
ldap_servers.ldap_server_name.auth_dn_prefix	用于构造要绑定到的DN的前缀和后缀。	uid=
ldap_servers.ldap_server_name.auth_dn_suffix	生成的DN将被构造为 auth_dn_prefix + escape(user_name) + auth_dn_suffix字符串。 auth_dn_suffix通常应将逗号“,”作为其第一个非空格字符。	,ou=Group,dc=node1,dc=com
ldap_servers.ldap_server_name.enable_tls	触发使用OpenLDAP服务器安全连接的标志。 <ul style="list-style-type: none">纯文本（ldap://）协议指定“no”（不推荐）。LDAP over SSL/TLS（ldaps://）协议指定“yes”。	yes
ldap_servers.ldap_server_name.tls_require_cert	SSL/TLS对端证书校验行为。 取值范围为：'never'、'allow'、'try'、'require'。	allow

说明

其他参数说明详细可以参考[<ldap_servers>配置参数详解](#)。

图 4-2 OpenLDAP 配置



步骤3 添加完配置后，单击“保存”，在弹出对话框中单击“确定”，配置保存成功后，单击“完成”。

步骤4 Manager页面，单击“实例”，选择ClickHouseServer实例，单击“更多 > 重启实例”，弹出对话框输入密码，单击“确定”。重启实例对话框，单击“确定”，根据界面提示信息确认实例重启成功，单击“完成”重启操作完成。

步骤5 登录ClickHouseServer实例所在主机节点，进入“`${BIGDATA_HOME}/FusionInsight_ClickHouse_版本号/x_x_ClickHouseServer/etc`”目录。

```
cd ${BIGDATA_HOME}/FusionInsight_ClickHouse_*/x_x_ClickHouseServer/etc
```

步骤6 执行以下命令，查看配置文件config.xml，确认OpenLDAP参数是否配置成功。

```
cat config.xml
```

```
[root@k 3 etc]# cat config.xml
```

```
<vandex>
  <ldap_servers>
    <ldap_server_name>
      <auth_dn_prefix>uid=</auth_dn_prefix>
      <port>636</port>
      <host>localhost</host>
      <enable_tls>yes</enable_tls>
      <tls_require_cert>allow</tls_require_cert>
      <auth_dn_suffix>,ou=Group,dc=node1,dc=com</auth_dn_suffix>
    </ldap_server_name>
  </ldap_servers>
</vandex>
```

步骤7 以root用户登录ClickHouseServer实例所在的节点。

步骤8 执行以下命令获取“clickhouse.keytab”文件路径。

```
ls ${BIGDATA_HOME}/FusionInsight_ClickHouse_*/install/FusionInsight-ClickHouse-*/clickhouse/keytab/clickhouse.keytab
```

步骤9 以客户端安装用户，登录安装客户端的节点。

步骤10 执行以下命令，切换到ClickHouse客户端安装目录。

```
cd /opt/client
```

步骤11 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤12 执行如下命令使用客户端命令连接ClickHouseServer实例。

- 如果当前集群已启用Kerberos认证，使用clickhouse.keytab连接ClickHouseServer实例：

```
clickhouse client --host ClickHouseServer实例所在节点IP --user clickhouse/hadoop.<系统域名> --password 步骤8中获取的clickhouse.keytab路径 --port ClickHouse的端口号
```

📖 说明

系统域名：默认为hadoop.com。具体可登录集群FusionInsight Manager，单击“系统 > 权限 > 域和互信”，“本端域”参数值即为系统域名。在执行命令时改为小写。

- 如果当前集群未启用Kerberos认证，使用ClickHouse管理员用户连接ClickHouseServer实例：

```
clickhouse client --host ClickHouseServer实例所在节点IP --user clickhouse --port ClickHouse的端口号
```

步骤13 创建OpenLDAP中的普通用户。

如以下语句，在集群default_cluster上创建testUser用户，设置ldap_server为步骤6中<ldap_servers>标签下的OpenLDAP服务名，本示例为ldap_server_name。

```
CREATE USER testUser ON CLUSTER default_cluster IDENTIFIED WITH  
ldap_server BY 'ldap_server_name';
```

testUser用户为OpenLDAP中已有的用户名，请根据实际情况修改。

步骤14 退出客户端，使用新建的用户登录验证配置是否成功。

```
exit;
```

```
clickhouse client --host ClickHouseServer实例IP --user testUser --password --  
port ClickHouse的端口号
```

输入testUser对应的密码

----结束

<ldap_servers>配置参数详解

- host
OpenLDAP服务器主机名或IP，必选参数，不能为空。
- port
OpenLDAP服务器端口，如果enable_tls参数设置为true，则默认为636，否则为389。
- auth_dn_prefix, auth_dn_suffix
用于构造要绑定到的DN的前缀和后缀。
实际上，生成的DN将被构造为auth_dn_prefix + escape(user_name) + auth_dn_suffix字符串。
注意，这意味着auth_dn_suffix通常应将逗号“，”作为其第一个非空格字符。
- enable_tls
触发使用OpenLDAP服务器安全连接的标志。
为纯文本（ldap://）协议指定“no”（不推荐）。
为LDAP over SSL/TLS (ldaps://)协议指定“yes”（建议为默认值）。
- tls_minimum_protocol_version
SSL/TLS的最小协议版本。
接受的值是：'ssl2'、'ssl3'、'tls1.0'、'tls1.1'、'tls1.2'（默认值）。
- tls_require_cert
SSL/TLS对端证书校验行为。
接受的值是：'never'、'allow'、'try'、'require'（默认值）。
- tls_cert_file
证书文件。
- tls_key_file
证书密钥文件。
- tls_ca_cert_file
CA证书文件。

- `tls_ca_cert_dir`
CA证书所在的目录。
- `tls_cipher_suite`
允许加密套件。

4.3 使用 ClickHouse 客户端

ClickHouse是面向联机分析处理的列式数据库，支持SQL查询，且查询性能好，特别是基于大宽表的聚合分析查询性能非常优异，比其他分析型数据库速度快一个数量级。

前提条件

已安装客户端，例如安装目录为“`/opt/client`”。以下操作的客户端目录只是举例，请根据实际安装目录修改。在使用客户端前，需要先下载并更新客户端配置文件，确认Manager的主管理节点后才能使用客户端。

操作步骤

步骤1 安装客户端，具体请参考[安装客户端](#)章节。

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤4 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤5 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建ClickHouse表的权限，具体请参见[ClickHouse用户及权限管理](#)章节，为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行本步骤。

如果是MRS 3.1.0版本集群，则需要先执行：**export CLICKHOUSE_SECURITY_ENABLED=true**

```
kinit 组件业务用户
```

例如，**kinit clickhouseuser**。

步骤6 执行ClickHouse组件的客户端命令。

执行**clickhouse -h**，查看ClickHouse组件命令帮助。

回显信息如下：

```
Use one of the following commands:
clickhouse local [args]
clickhouse client [args]
clickhouse benchmark [args]
clickhouse server [args]
clickhouse performance-test [args]
clickhouse extract-from-config [args]
clickhouse compressor [args]
clickhouse format [args]
clickhouse copier [args]
```

```
clickhouse obfuscator [args]
...
```

MRS 3.1.0版本，使用**clickhouse client**命令连接ClickHouse服务端：

- 例如，当前集群未启用Kerberos认证，使用非ssl方式登录：
clickhouse client --host ClickHouse的实例IP --port 9000 --user 用户名 --password
输入用户密码
- 例如，当前集群已启用Kerberos认证，使用ssl安全方式登录。
Kerberos集群场景下没有默认用户，必须在Manager上创建用户，详细参考[ClickHouse用户及权限管理](#)。
使用kinit认证成功后，客户端登录时可以不携带--user和--password参数，即使用kinit认证的用户登录。

clickhouse client --host ClickHouse的实例IP --port 9440 --secure

MRS 3.1.2及之后版本，使用**clickhouse client**命令连接ClickHouse服务端：

- 例如，当前集群未启用Kerberos认证，使用非ssl方式登录：
clickhouse client --host ClickHouse的实例IP --port 9000 --user 用户名 --password
输入用户密码
- 例如，当前集群已启用Kerberos认证，使用ssl安全方式登录。
Kerberos集群场景下没有默认用户，必须在Manager上创建用户，详细参考[ClickHouse用户及权限管理](#)。
clickhouse client --host ClickHouse的实例IP --port 9440 --user 用户名 --password --secure
输入用户密码

执行**quit;**命令，退出ClickHouse服务端连接。

相关参数使用说明如[表4-5](#)：

表 4-5 clickhouse client 命令行参数说明

参数名	参数说明
--host	服务端的host名称，默认是localhost。您可以选择使用ClickHouse实例所在节点主机名或者IP地址。 说明 ClickHouse的实例IP地址可登录集群FusionInsight Manager，然后选择“集群 > 服务 > ClickHouse > 实例”，获取ClickHouseServer实例对应的业务IP地址。
--port	连接的端口。 <ul style="list-style-type: none"> • 如果使用ssl安全连接则默认端口为9440，并且需要携带参数--secure。具体的端口值可通过ClickHouseServer实例配置搜索“tcp_port_secure”参数获取。 • 如果使用非ssl安全连接则默认端口为9000，不需要携带参数--secure。具体的端口值可通过ClickHouseServer实例配置搜索“tcp_port”参数获取。

参数名	参数说明
--user	<p>用户名。</p> <p>可以在Manager上创建该用户名并绑定对应的角色权限，具体请参见ClickHouse用户及权限管理章节。</p> <ul style="list-style-type: none"> 如果当前集群已启用Kerberos认证（集群为安全模式），使用kinit认证成功后，客户端登录时可以不携带--user和--password参数，即使用kinit认证的用户登录。Kerberos集群场景下没有默认用户，必须在Manager上创建该用户名。 如果当前集群未启用Kerberos认证（集群为普通模式），客户端登录时如果需要指定用户名和密码，不能使用FusionInsight Manager页面创建的ClickHouse用户，需要使用客户端命令行执行create user SQL语句创建ClickHouse用户。客户端登录时如果不需要指定用户名和密码参数时，默认使用default用户登录。
--password	密码。默认值：空字符串。该参数和--user参数配套使用，可以在Manager上创建用户名时设置该密码。
--query	使用非交互模式查询。
--database	默认当前操作的数据库。默认值：服务端默认的配置（默认是default）。
--multiline	如果指定，允许多行语句查询（Enter仅代表换行，不代表查询语句完结）。
--multiquery	如果指定，允许处理用;号分隔的多个查询，只在非交互模式下生效。
--format	使用指定的默认格式输出结果。
--vertical	如果指定，默认情况下使用垂直格式输出结果。在这种格式中，每个值都在单独的行上打印，适用显示宽表的场景。
--time	如果指定，非交互模式下会打印查询执行的时间到stderr中。
--stacktrace	如果指定，如果出现异常，会打印堆栈跟踪信息。
--config-file	配置文件的名称。
--secure	如果指定，将通过ssl安全模式连接到服务器。
--history_file	存放命令历史的文件的路径。
--param_<name>	带有参数的查询，并将值从客户端传递给服务器。具体用法详见 https://clickhouse.tech/docs/zh/interfaces/cli/#cli-queries-with-parameters 。

----结束

4.4 ClickHouse 表创建

ClickHouse依靠ReplicatedMergeTree引擎与ZooKeeper实现了复制表机制，用户在创建表时可以通过指定引擎选择该表是否高可用，每张表的分片与副本都是互相独立的。

同时ClickHouse依靠Distributed引擎实现了分布式表机制，在所有分片（本地表）上建立视图进行分布式查询，使用很方便。ClickHouse有数据分片（shard）的概念，这也是分布式存储的特点之一，即通过并行读写提高效率。

CPU架构为鲲鹏计算的ClickHouse集群表引擎不支持使用HDFS和Kafka。

查看 ClickHouse 服务 cluster 等环境参数信息

步骤1 使用ClickHouse客户端连接到ClickHouse服务端，具体请参考[使用ClickHouse客户端](#)。

步骤2 查询集群标识符cluster等其他环境参数信息。

```
select cluster,shard_num,replica_num,host_name from system.clusters;
```

```
SELECT
  cluster,
  shard_num,
  replica_num,
  host_name
FROM system.clusters
```

cluster	shard_num	replica_num	host_name
default_cluster_1	1	1	node-master1dOnG
default_cluster_1	1	2	node-group-1tXED0001
default_cluster_1	2	1	node-master2OXQS
default_cluster_1	2	2	node-group-1tXED0002
default_cluster_1	3	1	node-master3QsRI
default_cluster_1	3	2	node-group-1tXED0003

6 rows in set. Elapsed: 0.001 sec.

步骤3 查询分片标识符shard和副本标识符replica。

```
select * from system.macros;
```

```
SELECT *
FROM system.macros
```

macro	substitution
id	76
replica	2
shard	3

3 rows in set. Elapsed: 0.001 sec.

----结束

创建本地复制表和分布式表

步骤1 客户端登录ClickHouse节点，例如：`clickhouse client --host node-master3QsRI --multiline --port 9440 --secure;`

 说明

node-master3QsRI 参数为[查看ClickHouse服务cluster等环境参数信息](#)中步骤2对应的 host_name 参数的值。

步骤2 使用ReplicatedMergeTree引擎创建复制表。

详细的语法说明请参考：<https://clickhouse.tech/docs/zh/engines/table-engines/mergetree-family/replication/#creating-replicated-tables>。

例如，如下在default_cluster_1集群节点上和default数据库下创建表名为test的ReplicatedMergeTree表：

```
CREATE TABLE default.test ON CLUSTER default_cluster_1
(
  `EventDate` DateTime,
  `id` UInt64
)
ENGINE = ReplicatedMergeTree('/clickhouse/tables/{shard}/default/test',
'{replica}')
PARTITION BY toYYYYMM(EventDate)
ORDER BY id;
```

参数说明如下：

- ON CLUSTER语法表示分布式DDL，即执行一次就可在集群所有实例上创建同样的本地表。
- default_cluster_1为[查看ClickHouse服务cluster等环境参数信息](#)中步骤2查询到的cluster集群标识符。

 注意

ReplicatedMergeTree引擎接收两个参数：

- ZooKeeper中该表相关数据的存储路径。

该路径必须在/clickhouse目录下，否则后续可能因为ZooKeeper配额不够导致数据插入失败。

为了避免不同表在ZooKeeper上数据冲突，目录格式必须按照如下规范填写：

/clickhouse/tables/{shard} default/test，其中/clickhouse/tables/{shard}为固定值，default为数据库名，test为创建的表名。

- 副本名称，一般用{replica}即可。

```
CREATE TABLE default.test ON CLUSTER default_cluster_1
(
  `EventDate` DateTime,
  `id` UInt64
)
ENGINE = ReplicatedMergeTree('/clickhouse/tables/{shard}/default/test', '{replica}')
PARTITION BY toYYYYMM(EventDate)
ORDER BY id
```

```
┌ host ───────────────────┬ port ┤ status┤ error┤ num_hosts_remaining┤
num_hosts_active┐
node-group-1tXED0002    │ 9000 │ 0    │    │ 5    │ 3    │
node-group-1tXED0003    │ 9000 │ 0    │    │ 4    │ 3    │
node-master1dOnG        │ 9000 │ 0    │    │ 3    │ 3    │
└────────────────────────┴───┘

┌ host ───────────────────┬ port ┤ status┤ error┤ num_hosts_remaining┤
num_hosts_active┐
node-master3QsRI        │ 9000 │ 0    │    │ 2    │ 0    │
node-group-1tXED0001    │ 9000 │ 0    │    │ 1    │ 0    │
node-master2OXQS        │ 9000 │ 0    │    │ 0    │ 0    │
└────────────────────────┴───┘
```

6 rows in set. Elapsed: 0.189 sec.

步骤3 使用Distributed引擎创建分布式表。

例如，以下将在default_cluster_1集群节点上和default数据库下创建名为test_all 的Distributed表：

```
CREATE TABLE default.test_all ON CLUSTER default_cluster_1
(
  `EventDate` DateTime,
  `id` UInt64
)
ENGINE = Distributed(default_cluster_1, default, test, rand());
```

```
CREATE TABLE default.test_all ON CLUSTER default_cluster_1
(
  `EventDate` DateTime,
  `id` UInt64
)
ENGINE = Distributed(default_cluster_1, default, test, rand())
```

```
┌ host ───────────────────┬ port ┤ status┤ error┤ num_hosts_remaining┤
num_hosts_active┐
node-group-1tXED0002    │ 9000 │ 0    │    │ 5    │ 0    │
node-master3QsRI        │ 9000 │ 0    │    │ 4    │ 0    │
node-group-1tXED0003    │ 9000 │ 0    │    │ 3    │ 0    │
node-group-1tXED0001    │ 9000 │ 0    │    │ 2    │ 0    │
node-master1dOnG        │ 9000 │ 0    │    │ 1    │ 0    │
node-master2OXQS        │ 9000 │ 0    │    │ 0    │ 0    │
└────────────────────────┴───┘
```

6 rows in set. Elapsed: 0.115 sec.

📖 说明

Distributed引擎需要以下几个参数：

- default_cluster_1为[查看ClickHouse服务cluster等环境参数信息](#)中[步骤2](#)查询到的cluster集群标识符。
- default本地表所在的数据库名称。
- test为本地表名称，该例中为[步骤2](#)中创建的表名。
- （可选的）分片键（sharding key）

该键与config.xml中配置的分片权重（weight）一同决定写入分布式表时的路由，即数据最终落到哪个物理表上。它可以是表中一列的原始数据（如site_id），也可以是函数调用的结果，如上面的SQL语句采用了随机值rand()。注意该键要尽量保证数据均匀分布，另外一个常用的操作是采用区分度较高的列的哈希值，如intHash64(user_id)。

----结束

ClickHouse 表数据操作

步骤1 客户端登录ClickHouse节点。例如：

```
clickhouse client --host node-master3QsRI --multiline --port 9440 --secure;
```

📖 说明

node-master3QsRI 参数为[查看ClickHouse服务cluster等环境参数信息](#)中[步骤2](#)对应的host_name参数的值。

步骤2 参考[创建本地复制表和分布式表](#)创建表后，可以插入数据到本地表。

例如插入数据到本地表：test

```
insert into test values(toDateTime(now()), rand());
```

步骤3 查询本地表信息。

例如查询[步骤2](#)中的表test数据信息：

```
select * from test;
```

```
SELECT *  
FROM test
```

EventDate	id
2020-11-05 21:10:42	1596238076

1 rows in set. Elapsed: 0.002 sec.

步骤4 查询Distributed分布式表。

例如[步骤3](#)中因为分布式表test_all基于test创建，所以test_all表也能查询到和test相同的数据。

```
select * from test_all;
```

```
SELECT *  
FROM test_all
```

EventDate	id
2020-11-05 21:10:42	1596238076

```
1 rows in set. Elapsed: 0.004 sec.
```

步骤5 切换登录节点为相同shard_num的shard节点，并且查询当前表信息，能查询到相同的表数据。

例如，退出原有登录节点：**exit;**

切换到节点node-group-1tXED0003:

```
clickhouse client --host node-group-1tXED0003 --multiline --port 9440 --secure;
```

📖 说明

通过**步骤2**可以看到node-group-1tXED0003和node-master3QsRI的shard_num值相同。

```
show tables;
```

```
SHOW TABLES
```

name
test
test_all

步骤6 查询本地表数据。例如在节点node-group-1tXED0003查询test表数据。

```
select * from test;
```

```
SELECT *  
FROM test
```

EventDate	id
2020-11-05 21:10:42	1596238076

```
1 rows in set. Elapsed: 0.005 sec.
```

步骤7 切换到不同shard_num的shard节点，并且查询之前创建的表数据信息。

例如退出之前的登录节点node-group-1tXED0003:

```
exit;
```

切换到node-group-1tXED0001节点。通过**步骤2**可以看到node-group-1tXED0001和node-master3QsRI的shard_num值不相同。

```
clickhouse client --host node-group-1tXED0001 --multiline --port 9440 --secure;
```

查询test本地表数据，因为test是本地表所以在不同分片节点上查询不到数据。

```
select * from test;
```

```
SELECT *  
FROM test
```

```
Ok.
```

查询test_all分布式表数据，能正常查询到数据信息。

```
select * from test_all;
```

```
SELECT *  
FROM test
```

```
EventDate | id |
2020-11-05 21:12:19 | 3686805070 |

1 rows in set. Elapsed: 0.002 sec.
```

---结束

4.5 ClickHouse 数据导入

4.5.1 配置 ClickHouse 对接 RDS MySQL 数据库

ClickHouse面向OLAP场景提供高效的数据分析能力，支持通过MySQL等数据库引擎将远程数据库服务器中的表映射到ClickHouse集群中，后续可以在ClickHouse中进行数据分析。以下操作通过ClickHouse集群和RDS服务下的MySQL数据库实例对接进行举例说明。

前提条件

- 已提前准备好对接的RDS数据库实例及数据库用户名、密码。详细操作可以参考[创建和连接RDS数据库实例](#)。
- 已成功创建ClickHouse集群且集群和实例状态正常。

约束限制

- RDS数据库实例和ClickHouse集群在相同的VPC和子网内。
- 在进行数据同步操作时需要评估对源数据库和目标数据库性能的影响，同时建议您在业务低峰期执行数据同步。
- 当前ClickHouse支持和RDS服务下的MySQL、PostgreSQL实例进行对接，不支持对接SQL Server实例。

ClickHouse 通过 MySQL 引擎对接 RDS 服务

MySQL引擎用于将远程的MySQL服务器中的表映射到ClickHouse中，并允许您对表进行INSERT和SELECT查询，以方便您在ClickHouse与MySQL之间进行数据交换。

MySQL引擎使用语法：

```
CREATE DATABASE [IF NOT EXISTS] db_name [ON CLUSTER cluster]
ENGINE = MySQL('host:port', ['database' | database], 'user', 'password')
```

MySQL数据库引擎参数说明：

- host:port：RDS服务MySQL数据库实例IP地址和端口。
- database：RDS服务MySQL数据库名。
- user：RDS服务MySQL数据库用户名。
- password：RDS服务MySQL数据库用户密码，命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。

MySQL引擎使用示例：

步骤1 连接到RDS服务的MySQL数据库。详细操作可以参考[RDS服务MySQL实例连接](#)。

步骤2 在MySQL数据库上创建表，并插入数据。

创建表mysql_table：

```
CREATE TABLE `mysql_table` (  
  `int_id` INT NOT NULL AUTO_INCREMENT,  
  `float` FLOAT NOT NULL,  
  PRIMARY KEY (`int_id`));
```

插入表数据：

```
insert into mysql_table (`int_id`, `float`) VALUES (1,2);
```

步骤3 登录ClickHouse客户端安装节点。执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤4 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤5 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建ClickHouse表的权限，具体请参见[ClickHouse用户及权限管理](#)章节，为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行本步骤。

1. 如果是MRS 3.1.0版本集群，则需要先执行：**export CLICKHOUSE_SECURITY_ENABLED=true**
2. **kinit 组件业务用户**
例如，**kinit clickhouseuser**。

步骤6 使用客户端命令连接ClickHouse。

```
clickhouse client --host clickhouse实例IP --user 用户名 --password --port 端口号  
输入用户密码
```

步骤7 在ClickHouse中创建MySQL引擎的数据库，创建成功后自动与MySQL服务器交换数据。

```
CREATE DATABASE mysql_db ENGINE = MySQL('RDS服务MySQL数据库实例IP地址:MySQL数据库实例端口', 'MySQL数据库名', 'MySQL数据库用户名', 'MySQL数据库用户名密码');
```

步骤8 切换到新建的数据库mysql_db，并查询表数据。

```
USE mysql_db;
```

在ClickHouse中查询MySQL数据库表数据。

```
SELECT * FROM mysql_table;
```

```
┌──int_id──┬──float──┐  
└── 1 ───┴── 2 ───┘
```

新增插入数据后也可以正常进行查询。

```
INSERT INTO mysql_table VALUES (3,4);
```

```
SELECT * FROM mysql_table;
```

int_id	float
1	2
3	4

----结束

4.5.2 配置 ClickHouse 对接 OBS 源文件

使用 S3 表函数

步骤1 登录主OMS节点。

步骤2 执行以下命令获取OBS内的存储数据。

```
select * from S3(path, [ak, sk,] format, structure, [compression])
```

说明

- path: 访问域名/OBS文件路径, 登录OBS管理控制台, 在左侧导航栏单击“并行文件系统”, 在“并行文件系统”页面单击对应的文件系统名称, 在“文件”页面单击文件名称, 文件“链接”即path路径, 如图4-3所示。

图 4-3 文件路径



- ak: 参数可选, 具备访问OBS权限的ak。
- sk: 参数可选, 具备访问OBS权限的sk。
- format: 文件的格式。
- structure: 表的结构。
- compression: 参数可选, 压缩类型。

```
node-group-1sWT06081:~$ select * from s3('https://obs.obs.cn-east-3.amazonaws.com/clickhouse/S3_engine_test/*', 'CSV', 'name String,age int')
SELECT *
FROM s3('https://obs.obs.cn-east-3.amazonaws.com/clickhouse/S3_engine_test/*', 'CSV', 'name String,age int')
Query id: 999bb342-c799-4cd4-9296-fd8db99c972a
+----+----+
|name|age|
+----+----+
| 4  | 4  |
+----+----+
|name|age|
+----+----+
|xx2 | 3  |
+----+----+
2 rows in set. Elapsed: 0.266 sec.
```

----结束

使用 S3 表引擎

步骤1 登录主OMS节点。

步骤2 执行以下命令创建表。

```
CREATE TABLE test1_s3 ('name' String, 'age' int)
```

ENGINE = S3(path, [ak, sk,] format, [compression])

```
node-group-1sWT0001 :) CREATE TABLE test1_s3 ('name' String, 'age' int) ENGINE = S3('https://obs.obs.cn-east-3.amazonaws.com/clickhouse/s3_engine_test/*', 'ak', 'sk', 'CSV');
CREATE TABLE test1_s3
(
  'name' String,
  'age' int
)
ENGINE = S3('https://obs.obs.cn-east-3.amazonaws.com/clickhouse/s3_engine_test/*', 'ak', 'sk', 'CSV')
Query id: b0586eb4-a95f-4543-9868-b0f3b9ef3bb6
ok.
0 rows in set. Elapsed: 0.006 sec.
```

步骤3 执行以下命令查询表。

select * from test1_s3;

```
node-group-1sWT0001 :) select * from test1_s3;
SELECT *
FROM test1_s3
Query id: 079fe21d-54c1-4cc9-be8a-0f20a198f9dd
+----+----+
| name | age |
+----+----+
| xx2  | 3   |
+----+----+
| name | age |
+----+----+
| 4    | 4   |
+----+----+
2 rows in set. Elapsed: 0.277 sec.
```

----结束

修改 Manager 配置

登录FusionInsight Manager，选择“集群 > 服务 > ClickHouse > 配置 > 全部配置”。搜索参数项“clickhouse-config-customize”并添加参数值。参数值的添加参考下表。

参数	值
s3.endpoint-name.endpoint	OBS桶地址
s3.endpoint-name.access_key_id	OBS ak, 获取方法请参考 如何获取访问密钥AK/SK
s3.endpoint-name.secret_access_key	OBS sk, 获取方法请参考 如何获取访问密钥AK/SK

对于OBS直接分享出来的URL，一般是带HTTPS的，如果不能直接访问，请按如下步骤修改配置。

登录FusionInsight Manager，选择“集群 > 服务 > ClickHouse > 配置 > 全部配置”。搜索参数项“clickhouse-config-customize”并添加参数值。参数值的添加参考下表。

参数	值
openSSL.client.loadDefaultCAFile	true
openSSL.client.cacheSessions	true
openSSL.client.disableProtocols	ssl2,ssl3
openSSL.client.preferServerCiphers	true

参数	值
openssl.client.invalidCertificateHandler.name	AcceptCertificateHandler

修改完成后，单击“保存”。

4.5.3 同步 Kafka 数据至 ClickHouse

您可以通过创建Kafka引擎表将Kafka数据自动同步至ClickHouse集群，具体操作详见本章节描述。

前提条件

- 已创建Kafka集群。已安装Kafka客户端，详细可以参考[安装客户端](#)。
- 已创建ClickHouse集群，并且ClickHouse集群和Kafka集群在同一VPC下，网络可以互通，并安装ClickHouse客户端。

约束限制

当前ClickHouse不支持和开启安全模式的Kafka集群进行对接。

Kafka 引擎表使用语法说明

- **语法**

```
CREATE TABLE [IF NOT EXISTS] [db.]table_name [ON CLUSTER cluster]
(
  name1 [type1] [DEFAULT|MATERIALIZED|ALIAS expr1],
  name2 [type2] [DEFAULT|MATERIALIZED|ALIAS expr2],
  ...
) ENGINE = Kafka()
SETTINGS
  kafka_broker_list = 'host1:port1,host2:port2',
  kafka_topic_list = 'topic1,topic2,...',
  kafka_group_name = 'group_name',
  kafka_format = 'data_format';
[kafka_row_delimiter = 'delimiter_symbol',]
[kafka_schema = ',']
[kafka_num_consumers = N]
```
- **参数说明**

表 4-6 Kafka 引擎表参数说明

参数名	是否必选	参数说明
kafka_broker_list	是	Kafka集群broker实例的IP和端口列表。例如： <i>kafka集群broker实例IP1:9092,kafka集群broker实例IP2:9092,kafka集群broker实例IP3:9092</i> 。 说明 启用Kerberos认证下，使用21005端口需要“allow.everyone.if.no.acl.found”参数值设置为true；若不设置此参数，操作会报错。 Kafka集群broker实例IP获取方法如下： <ul style="list-style-type: none">• MRS 3.x及后续版本，登录FusionInsight Manager，然后选择“集群 > 待操作的集群名称 > 服务 > Kafka”。单击“实例”，查看Kafka角色实例的IP地址。
kafka_topic_list	是	Kafka的topic列表。
kafka_group_name	是	Kafka的Consumer Group名称，可以自己指定。
kafka_format	是	Kafka消息体格式。例如JSONEachRow、CSV、XML等。
kafka_row_delimiter	否	每个消息体（记录）之间的分隔符。
kafka_schema	否	如果解析格式需要一个schema时，此参数必填。
kafka_num_consumers	否	单个表的消费者数量。默认值是：1，如果一个消费者的吞吐量不足，则指定更多的消费者。消费者的总数不应该超过topic中分区的数量，因为每个分区只能分配一个消费者。

Kafka 数据同步至 ClickHouse 操作示例

步骤1 参考[Kafka客户端使用实践](#)，切换到Kafka客户端安装目录。

1. 以Kafka客户端安装用户，登录Kafka安装客户端的节点。
2. 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

3. 执行以下命令配置环境变量。

```
source bigdata_env
```

4. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
如果是MRS 3.1.0版本集群，则需要先执行：export  
CLICKHOUSE_SECURITY_ENABLED=true
```

kinit 组件业务用户

步骤2 执行以下命令，创建Kafka的Topic。详细的命令使用可以参考[创建Kafka Topic](#)。

```
kafka-topics.sh --topic kafkacktest2 --create --zookeeper ZooKeeper角色实例  
IP:ZooKeeper侦听客户端连接的端口/kafka --partitions 2 --replication-factor 1
```

📖 说明

- **--topic**参数值为要创建的Topic名称，本示例创建的名称为kafkacktest2。
- **--zookeeper**: ZooKeeper角色实例所在节点IP地址，填写三个角色实例其中任意一个的IP地址即可。ZooKeeper角色实例所在节点IP获取参考如下。
 - MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > ZooKeeper > 实例”。查看ZooKeeper角色实例的IP地址。
 - MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > ZooKeeper > 实例”。查看ZooKeeper角色实例的IP地址。
- **--partitions**主题分区数和**--replication-factor**主题备份个数不能大于Kafka角色实例数量。
- **ZooKeeper侦听客户端连接的端口**获取方式：登录FusionInsight Manager，选择“集群 > 服务 > ZooKeeper”，在“配置”页签查看“clientPort”的值。默认为24002。

步骤3 参考[使用ClickHouse客户端](#)登录ClickHouse客户端。

1. 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

2. 执行以下命令配置环境变量。

```
source bigdata_env
```

3. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建ClickHouse表的权限，具体请参见[ClickHouse用户及权限管理](#)章节，为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行本步骤。

kinit 组件业务用户

例如，**kinit clickhouseuser**。

4. 执行以下命令连接到要导入数据的ClickHouse实例节点。

```
clickhouse client --host ClickHouse的实例IP --user 登录名 --password --port  
ClickHouse的端口号 --database 数据库名 --multiline
```

输入用户密码

步骤4 参考[Kafka引擎表使用语法说明](#)，在ClickHouse中创建Kafka引擎表。例如，如下建表语句在default数据库下，创建表名为kafka_src_tbl3，Topic名为kafkacktest2、消息格式为JSONEachRow的Kafka引擎表。

```
create table kafka_src_tbl3 on cluster default_cluster  
(id UInt32, age UInt32, msg String)  
ENGINE=Kafka()  
SETTINGS  
kafka_broker_list='kafka集群broker实例IP1:9092,kafka集群broker实例IP2:9092,kafka集群broker实例  
IP3:9092',  
kafka_topic_list='kafkacktest2',  
kafka_group_name='cg12',  
kafka_format='JSONEachRow';
```

步骤5 创建ClickHouse本地复制表。例如，如下创建表名为kafka_dest_tbl3的ReplicatedMergeTree表。

```
create table kafka_dest_tbl3 on cluster default_cluster  
( id UInt32, age UInt32, msg String )  
engine = ReplicatedMergeTree('/clickhouse/tables/{shard}/default/kafka_dest_tbl3', '{replica}')
```

```
partition by age  
order by id;
```

步骤6 创建MATERIALIZED VIEW，该视图会在后台转换Kafka引擎中的数据并将其放入创建的ClickHouse表中。

```
create materialized view consumer3 on cluster default_cluster to kafka_dest_tbl3 as select * from  
kafka_src_tbl3;
```

步骤7 再次执行**步骤1**，进入Kafka客户端安装目录。

步骤8 执行以下命令，在Kafka的Topic中产生消息。例如，如下命令向**步骤2**中创建的Topic发送消息。

```
kafka-console-producer.sh --broker-list kafka集群broker实例IP1:9092,kafka集群  
broker实例IP2:9092,kafka集群broker实例IP3:9092 --topic kafkacktest2  
>{"id":31,"age":30,"msg":"31 years old"}  
>{"id":32,"age":30,"msg":"31 years old"}  
>{"id":33,"age":30,"msg":"31 years old"}  
>{"id":35,"age":30,"msg":"31 years old"}
```

步骤9 使用ClickHouse客户端登录**步骤3**中ClickHouse实例节点，查询ClickHouse表数据。例如，查询kafka_dest_tbl3本地复制表，Kafka消息中的数据已经同步到该表。

```
select * from kafka_dest_tbl3;
```

```
ClickHouseXkya.mrs-2xxk.com :) select * from kafka_dest_tbl3;  
  
SELECT *  
FROM kafka_dest_tbl3  
  
Query id: 386ed9db-26b6-4823-8ce5-5b2de14384bb  
  
  id  age  msg  
  ---  ---  ---  
  31   30  31 years old  
  
  id  age  msg  
  ---  ---  ---  
  32   30  31 years old  
  
  id  age  msg  
  ---  ---  ---  
  33   30  31 years old  
  
  id  age  msg  
  ---  ---  ---  
  35   30  31 years old  
  
4 rows in set. Elapsed: 0.003 sec.
```

----结束

4.5.4 导入 DWS 表数据至 ClickHouse

ClickHouse支持CSV、JSON等格式文件的数据导入导出操作。本章节主要介绍怎么把DWS数据仓库服务中的表数据导出到CSV文件，再把CSV文件数据导入到ClickHouse表中。

前提条件

- ClickHouse集群和实例状态正常。
- DWS集群已创建，已获取到相关表所在的数据库用户名和密码。
- 已安装MRS客户端，例如安装目录为“/opt/client”。以下操作的客户端目录只是举例，请根据实际安装目录修改。在使用客户端前，需要先下载并更新客户端配置文件，确认Manager的主管理节点后才能使用客户端。

DWS 服务数据导入到 ClickHouse

- 步骤1** 参考[下载Data Studio图形界面客户端](#)中的“Data Studio图形界面客户端”下载Data Studio工具。
- 步骤2** 使用已创建好的DWS集群中的数据库用户名、密码等信息，参考[使用Data Studio工具连接](#)章节连接DWS数据库。
- 步骤3** 将DWS数据库中的表数据导出到CSV格式文件。

1. （可选）如果DWS数据库对应的表和数据已经存在，该步骤请忽略。本文通过演示在DWS创建测试表，并插入测试数据进行演示。

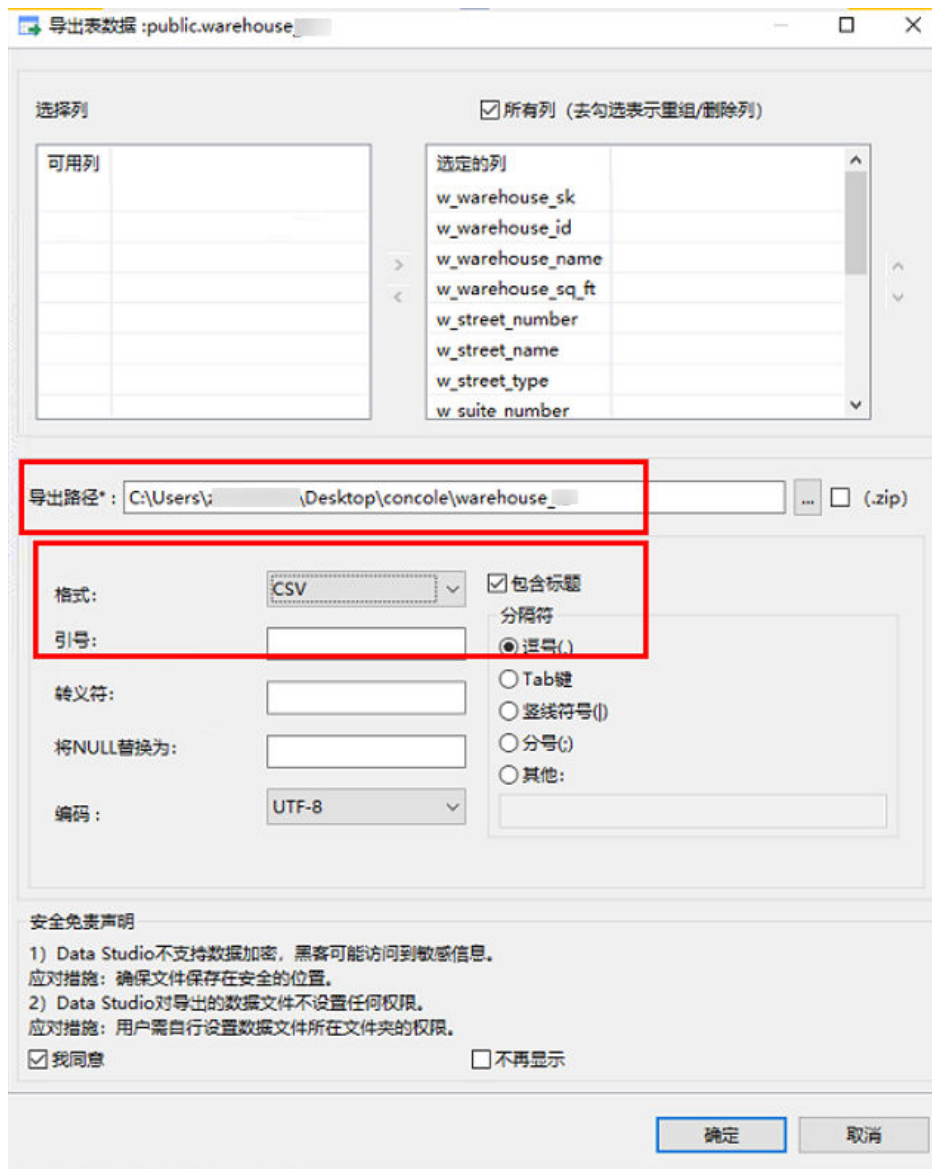
使用Data Studio创建测试表warehouse_t1，并插入测试数据。

```
CREATE TABLE warehouse_t1
(
  W_WAREHOUSE_SK INTEGER NOT NULL,
  W_WAREHOUSE_ID CHAR ( 16 ) NOT NULL,
  W_WAREHOUSE_NAME VARCHAR ( 20 ),
  W_WAREHOUSE_SQ_FT INTEGER,
  W_STREET_NUMBER CHAR ( 10 ),
  W_STREET_NAME VARCHAR ( 60 ),
  W_STREET_TYPE CHAR ( 15 ),
  W_SUITE_NUMBER CHAR ( 10 ),
  W_CITY VARCHAR ( 60 ),
  W_COUNTY VARCHAR ( 30 ),
  W_STATE CHAR ( 2 ),
  W_ZIP CHAR ( 10 ),
  W_COUNTRY VARCHAR ( 20 ),
  W_GMT_OFFSET DECIMAL ( 5,2 ),
  W_DATE DATE
);

INSERT INTO warehouse_t1 VALUES(1314, 123, 'name1', 2324, 123, 'STREET_NAME1', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:07');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name2', 2324, 123, 'STREET_NAME2', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:08');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name3', 2324, 123, 'STREET_NAME3', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:09');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name4', 2324, 123, 'STREET_NAME4', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:00');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name5', 2324, 123, 'STREET_NAME5', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:01');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name6', 2324, 123, 'STREET_NAME6', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:02');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name7', 2324, 123, 'STREET_NAME7', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:03');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name8', 2324, 123, 'STREET_NAME8', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:04');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name9', 2324, 123, 'STREET_NAME9', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:05');
INSERT INTO warehouse_t1 VALUES(1314, 123, 'name0', 2324, 123, 'STREET_NAME0', '12', '12',
'guangzhou', 'zhongguo', '1', '12', 'zn', 50.2, '2021-07-05 17:45:06');
INSERT INTO warehouse_t1(W_WAREHOUSE_SK, W_WAREHOUSE_ID, W_WAREHOUSE_NAME,
W_DATE) VALUES(1314, 123, 'name0', '2021-07-05 17:45:06');
```

2. 导出DWS表数据为CSV格式文件。

在Data Studio左侧的“对象浏览器”中，右键要导出的表，选择“导出表数据”。在导出界面选择具体的导出路径，格式选择CSV、分隔符选择逗号，在安全免责声明下选择“我同意”，单击“确定”完成数据导出。例如，本文导出表warehouse_t1数据文件为“warehouse_t1.csv”。



步骤4 使用WinSCP工具将导出的CSV文件上传到ClickHouse实例节点主机目录下。比如，当前上传“warehouse_t1.csv”文件到/opt目录下。

步骤5 以客户端安装用户，登录安装ClickHouse客户端的节点。

步骤6 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤7 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤8 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建ClickHouse表的权限，具体请参见[ClickHouse用户及权限管理](#)章节，为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行本步骤。

1. 如果是MRS 3.1.0版本集群，则需要先执行：**export CLICKHOUSE_SECURITY_ENABLED=true**
2. **kinit 组件业务用户**

例如，`kinit clickhouseuser`。

步骤9 执行以下命令连接到要导入数据的ClickHouse实例节点。

```
clickhouse client --host ClickHouse的实例IP --user 登录名 --password --port  
ClickHouse的端口号 --database 数据库名
```

输入用户密码

步骤10 在ClickHouse实例节点上创建和DWS表结构相同的表。

例如，当前执行以下建表语句，在ClickHouse实例上的默认数据库和用户下创建和**步骤3**中相同表结构的ReplicatedMergeTree表warehouse_t1。

```
CREATE TABLE warehouse_t1  
(  
    `W_WAREHOUSE_SK` Int32 NOT NULL,  
    `W_WAREHOUSE_ID` String NOT NULL,  
    `W_WAREHOUSE_NAME` String,  
    `W_WAREHOUSE_SQ_FT` Int32,  
    `W_STREET_NUMBER` String,  
    `W_STREET_NAME` String,  
    `W_STREET_TYPE` String,  
    `W_SUITE_NUMBER` String,  
    `W_CITY` String,  
    `W_COUNTY` String,  
    `W_STATE` String,  
    `W_ZIP` String,  
    `W_COUNTRY` String,  
    `W_GMT_OFFSET` Decimal(5, 2),  
    `W_DATE` DateTime  
)  
ENGINE = ReplicatedMergeTree('/clickhouse/tables/{shard}/default/warehouse_t1', '{replica}')  
PARTITION BY toYear(W_DATE)  
ORDER BY (W_DATE, W_WAREHOUSE_ID);
```

步骤11 退出ClickHouse客户端。

```
exit;
```

步骤12 执行以下命令，将导出的CSV文件数据导入到ClickHouse表中。

```
clickhouse client --host ClickHouse实例IP地址 --database 数据库名 --port 端口号  
--format_csv_delimiter="csv文件数据分隔符" --query="INSERT INTO 数据表名  
FORMAT CSV" < csv文件所在主机路径
```

例如，导入以逗号分隔的CSV文件“warehouse_t1.csv”数据到默认数据库和用户下的表warehouse_t1。

```
clickhouse client --host 10.248.12.10 --format_csv_delimiter="," --  
query="INSERT INTO warehouse_t1 FORMAT CSV" < /opt/warehouse_t1.csv
```

步骤13 导入完成后，登录ClickHouse客户端连接导入数据的ClickHouse实例节点，执行查询命令查看导入的结果。

例如，导入完成后查询表warehouse_t1数据，结果如下：

```
clickhouse client --host ClickHouse的实例IP --user 登录名 --password --port  
ClickHouse的端口号 --database 数据库名
```

输入用户密码

```
select * from warehouse_t1;
```

```
SELECT *
FROM warehouse_t1
```

Query id: 0af421a5-200f-4005-b254-c6013a5c5541

M_WAREHOUSE_SK	M_WAREHOUSE_ID	M_WAREHOUSE_NAME	M_WAREHOUSE_SQ_FT	M_STREET_NUMBER	M_STREET_NAME	M_STREET_TYPE	M_SUITE_NUMBER	M_CITY	M_COUNTY	M_STATE	M_ZIP	M_COUNTRY	M_GMT_OFFSET	M_DATE
1314	123	name4	2324	123	STREET_NAME4	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:00
1314	123	name5	2324	123	STREET_NAME5	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:01
1314	123	name6	2324	123	STREET_NAME6	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:02
1314	123	name7	2324	123	STREET_NAME7	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:03
1314	123	name8	2324	123	STREET_NAME8	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:04
1314	123	name9	2324	123	STREET_NAME9	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:05
1314	123	name0	2324	123	STREET_NAME0	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:06
1314	123	name1	2324	123	STREET_NAME1	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:07
1314	123	name2	2324	123	STREET_NAME2	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:08
1314	123	name3	2324	123	STREET_NAME3	12	12	guangzhou	zhongguo	1	50	zn	50:00	2021-07-05 17:45:09

----结束

4.5.5 ClickHouse 数据导入导出

使用 ClickHouse 客户端导入导出数据

本章节主要介绍使用ClickHouse客户端导入导出文件数据的基本语法和使用说明。

- CSV格式数据导入

```
clickhouse client --host 主机名/ClickHouse实例IP地址 --database 数据库名 --port 端口号 --secure --format_csv_delimiter="csv文件数据分隔符" --query="INSERT INTO 数据表名 FORMAT CSV" < csv文件所在主机路径
```

使用示例:

```
clickhouse client --host 10.5.208.5 --database testdb --port 9440 --secure --format_csv_delimiter="," --query="INSERT INTO testdb.csv_table FORMAT CSV" < /opt/data
```

数据表需提前创建好。

- CSV格式数据导出



注意

导出数据为CSV格式的文件，可能存在CSV注入的安全风险，请谨慎使用。

```
clickhouse client --host 主机名/ClickHouse实例IP地址 --database 数据库名 --port 端口号 -m --secure --query="SELECT * FROM 表名" > csv文件导出路径
```

使用示例:

```
clickhouse client --host 10.5.208.5 --database testdb --port 9440 -m --secure --query="SELECT * FROM test_table" > /opt/test
```

- parquet格式数据导入

```
cat parquet格式文件 | clickhouse client --host 主机名/ClickHouse实例IP --database 数据库名 --port 端口号 -m --secure --query="INSERT INTO 表名 FORMAT Parquet"
```

使用示例:

```
cat /opt/student.parquet | clickhouse client --host 10.5.208.5 --database testdb --port 9440 -m --secure --query="INSERT INTO parquet_tab001 FORMAT Parquet"
```

- parquet格式数据导出

```
clickhouse client --host 主机名/ClickHouse实例IP --database 数据库名 --port 端口号 -m --secure --query="select * from 表名 FORMAT Parquet" > parquet格式文件输出路径
```

使用示例:

```
clickhouse client --host 10.5.208.5 --database testdb --port 9440 -m --secure --query="select * from test_table FORMAT Parquet" > /opt/student.parquet
```

- ORC格式数据导入


```
cat orc格式文件路径 | clickhouse client --host 主机名/ClickHouse实例IP --  
database 数据库名 --port 端口号 -m --secure --query="INSERT INTO 表名  
FORMAT ORC"
```

使用示例：

```
cat /opt/student.orc | clickhouse client --host 10.5.208.5 --database testdb --port 9440 -m --secure --  
query="INSERT INTO orc_tab001 FORMAT ORC"  
#orc格式文件格式文件数据可以从HDFS中导出，例如：  
hdfs dfs -cat /user/hive/warehouse/hivedb.db/emp_orc/000000_0_copy_1 | clickhouse client --host  
10.5.208.5 --database testdb --port 9440 -m --secure --query="INSERT INTO orc_tab001 FORMAT  
ORC"
```

- ORC格式数据导出

```
clickhouse client --host 主机名/ClickHouse实例IP --database 数据库名 --port  
端口 -m --secure --query="select * from 表名 FORMAT ORC" > 输出的ORC格  
式文件路径
```

使用示例：

```
clickhouse client --host 10.5.208.5 --database testdb --port 9440 -m --secure --query="select * from  
csv_tab001 FORMAT ORC" > /opt/student.orc
```

- JSON格式数据导入

```
INSERT INTO 表名 FORMAT JSONEachRow JSON格式字符串1 JSON格式字  
符串2
```

使用示例：

```
INSERT INTO test_table001 FORMAT JSONEachRow {"PageViews":5,  
"UserID":"4324182021466249494", "Duration":146,"Sign":-1}  
{"UserID":"4324182021466249494","PageViews":6,"Duration":185,"Sign":1}
```

- JSON格式数据导出

```
clickhouse client --host 主机名/ClickHouse实例IP --database 数据库名 --port  
端口号 -m --secure --query="SELECT * FROM 表名 FORMAT JSON|  
JSONEachRow|JSONCompact|..." > json文件输出路径
```

使用示例

```
#导出json  
clickhouse client --host 10.5.208.5 --database testdb --port 9440 -m --secure --query="SELECT *  
FROM test_table FORMAT JSON" > /opt/test.json
```

```
#导出json(JSONEachRow)  
clickhouse client --host 10.5.208.5 --database testdb --port 9440 -m --secure --query="SELECT *  
FROM test_table FORMAT JSONEachRow" > /opt/test_jsoneachrow.json
```

```
#导出json(JSONCompact)  
clickhouse client --host 10.5.208.5 --database testdb --port 9440 -m --secure --query="SELECT *  
FROM test_table FORMAT JSONCompact" > /opt/test_jsoncompact.json
```

4.6 ClickHouse 企业级能力增强

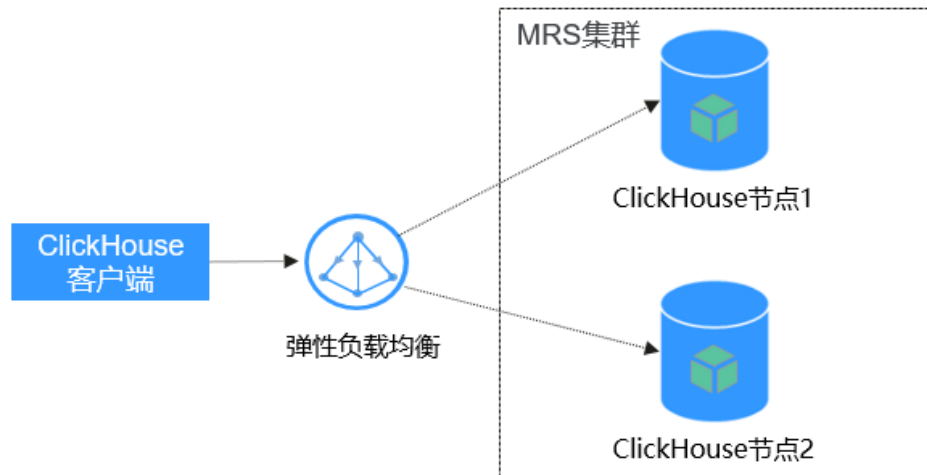
4.6.1 通过 ELB 访问 ClickHouse

当前ClickHouse不管是多分片还是多副本都是以集群方式部署，如果对外直接提供服务，将暴露多个节点服务，没有统一的访问入口。ClickHouse官方虽然提供了BalancedClickhouseDataSource的驱动方案，可以支持多节点的随机分配，提供了一定程度的负载均衡能力，但其故障检测能力不足，而且在扩缩容时，需要客户端感知集群节点变化，易用性不佳。

针对上述风险，MRS服务提供了基于弹性负载均衡ELB的部署架构图4-4。基于ELB的部署架构，可以将用户访问流量自动均匀分发到多台后端节点，扩展系统对外的服务

能力，实现更高水平的应用容错。当其中一台ClickHouse后端节点发生故障时，ELB通过故障转移方式正常对外提供服务。

图 4-4 通过弹性负载均衡访问 ClickHouse



ELB的部署架构对比BalancedClickhouseDataSource的优势可以参考表4-7说明。

表 4-7 ELB 和 BalancedClickhouseDataSource 两种负载均衡方案对比

负载均衡方案	方案对比
ELB	<ul style="list-style-type: none">• 支持多种请求策略• 故障自动检测转移• 后端ClickHouse扩容新增节点只需要修改ELB上的配置即可
BalancedClickhouseDataSource	<ul style="list-style-type: none">• 内部随机方式分发请求，可能会导致负载不均匀• 故障检测能力不足

当前通过ELB访问ClickHouse支持的协议和端口请参考表4-8，请根据实际使用场景选择配置。

表 4-8 通过 ELB 访问 ClickHouse 支持的协议和端口列表

协议	端口	场景描述
TCP	9000	通过客户端请求到ELB连接ClickHouse场景时配置。例如使用clickhouse client命令连接，host参数为ELB的私有IP地址。

协议	端口	场景描述
HTTP	8123	发送http请求到ELB连接ClickHouse场景时配置。

本章节演示如何实现客户端通过ELB访问ClickHouse。具体操作分为以下几个步骤：

- 步骤一：购买ELB并获取其私有IP地址。
- 步骤二：添加ELB监听器，配置协议端口。
- 步骤三：在ELB上添加ClickHouse后端服务器。
- 步骤四：使用客户端通过ELB访问ClickHouse。

前提条件

- MRS集群已创建，ClickHouse实例状态正常。
- 已安装MRS客户端，例如安装目录为“/opt/client”。以下操作的客户端目录只是举例，请根据实际安装目录修改。

购买 ELB 并配置对接 ClickHouse

购买ELB并获取其私有IP地址

详细操作步骤请参考[创建共享型负载均衡器](#)。

- 步骤1** 登录“弹性负载均衡器”控制台，在“负载均衡器”界面单击“购买弹性负载均衡”。
- 步骤2** 在“购买弹性负载均衡”界面，“实例规格类型”选择“共享型”，“所属VPC”和“子网”参数需要和MRS集群保持一致，其他参数保持默认即可。
- 步骤3** 单击“立即购买”，确认配置信息，并单击“提交”。
- 步骤4** 创建完成后，在“负载均衡器”界面，选择对应的区域即可看到新建的负载均衡器。查看并获取该负载均衡器的私有IP地址。



----结束

添加ELB监听器

详细操作步骤请参考[添加监听器](#)。

- 步骤1** 在“负载均衡器”界面，单击需要添加监听器的负载均衡名称。
- 步骤2** 选择“监听器 > 添加监听器”。



步骤3 在“添加监听器”界面，根据界面提示完成具体配置。

1. 配置监听器。

“前端协议/端口”选择“TCP”、端口填写“9000”，其他参数保持默认。配置完成单击“下一步”。

说明

如果是通过HTTP请求访问，则“前端协议/端口”选择“HTTP”、端口填写“8123”。

2. 配置后端服务器组。

“分配策略类型”参数选择“加权轮询算法”。单击“完成”，添加成功后，单击“确定”完成配置。

----结束

添加ClickHouse后端服务器

详细操作步骤请参考[添加后端服务器](#)。

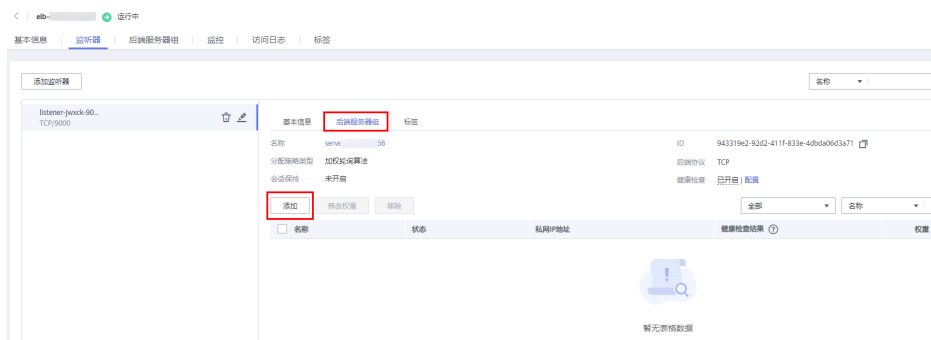
步骤1 登录MRS控制台，单击要对接的MRS集群名称。

步骤2 在MRS集群页面，单击“节点管理”，在ClickHouse节点组名称下，获取ClickHouse实例节点名称和IP地址。



步骤3 登录“弹性负载均衡器”控制台，单击已创建的负载均衡器名称。

步骤4 单击“监听器”，在“监听器”界面选择“后端服务器组”页签，单击“添加”。



步骤5 在“添加后端服务器”界面，根据**步骤2**中获取到的ClickHouse实例节点名称和IP地址勾选后端服务器。单击“下一步”。

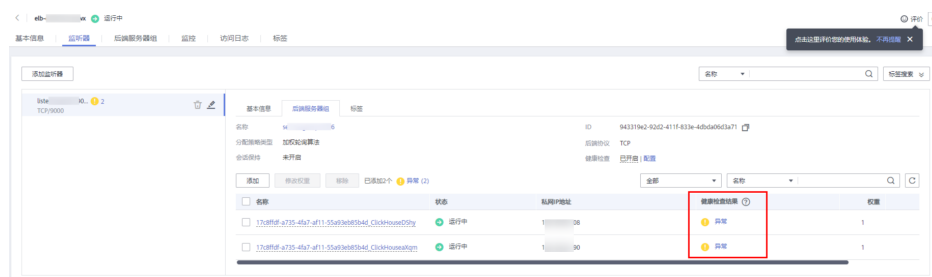
步骤6 “批量添加端口”参数填写为“9000”，单击“确定”。确认后端口配置无误后，单击“完成”。

📖 说明

如果是通过HTTP请求访问，端口填写“8123”。

步骤7 后端服务器配置安全组。

配置完成后，在“监听器”界面的“后端服务器组”页签下，对应的后端服务器显示“健康检查结果”状态为“异常”。

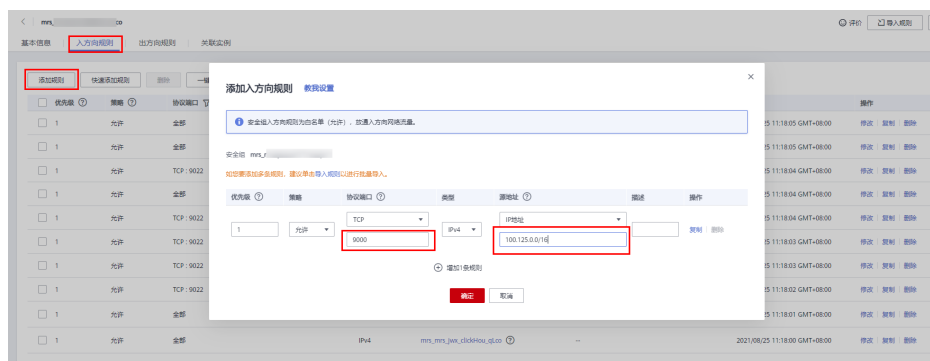


解决如上问题需要在ClickHouse后端服务器对应的安全组下放通“100.125.0.0/16”网段，具体操作如下：

1. 在“监听器”界面的“后端服务器”页签下，单击任意一个服务器名称。
2. 单击“安全组 > 配置规则”，选择“入方向规则 > 添加规则”。
3. 在“添加入方向规则”界面添加协议为TCP，端口为9000，IP地址配置“100.125.0.0/16”。单击“确定”完成配置。

📖 说明

如果是通过HTTP请求访问，端口填写“8123”。



4. 重新进入到创建的负载均衡器，刷新浏览器页面，单击“监听器”界面中的“后端服务器组”页签，对应的后端服务器“健康检查结果”状态显示为“正常”。

----结束

通过ELB访问ClickHouse

步骤1 登录Manager页面选择“集群 > 服务 > ClickHouse > 配置 > 全部配置”，修改参数“SSL_NONESSL_BOTH_ENABLE”值为“true”。

步骤2 参考[使用ClickHouse客户端](#)使用客户端登录ClickHouse服务实例节点。**注意：**客户端命令clickhouse client中的host参数填写[步骤4](#)中获取的ELB私有IP地址。

步骤3 在客户端界面查看通过ELB可以正常连接到ClickHouse实例节点。

📖 说明

手工通过客户端命令连接时，因为并发请求数较少，ELB可能始终将请求发送给一个后端ClickHouse节点，属于正常现象。

如果并发请求数多时，ELB会把请求轮询分配给多个后端ClickHouse节点。

----结束

4.6.2 ClickHouse 开启 mysql_port 配置

本章节指导用户使用MySQL客户端连接ClickHouse。

操作步骤

步骤1 登录FusionInsight Manager，选择“集群 > 服务 > ClickHouse > 配置 > 全部配置”。搜索参数项“clickhouse-config-customize”添加名称为“mysql_port”，值为“9004”的参数值。

📖 说明

参数值可以自行设置。

修改完成后，单击“保存”。

步骤2 单击“概览”页签，选择“更多 > 重启实例”或者“更多 > 滚动重启实例”。

----结束

4.7 ClickHouse 性能调优

4.7.1 数据表报错 Too many parts 解决方法

问题排查步骤

1. 磁盘或其他存储介质问题导致merge过慢或者中止。
登录Manager页面，检查是否存在磁盘容量不足或其他磁盘告警，如果存在，请按照告警指导处理。
如果是磁盘容量不足，也可以联系客服删除部分过期数据，释放空间，快速恢复业务。
2. Zookeeper异常导致merge无法正常执行。
登录Manager页面，检查ZooKeeper是否存在服务不可用、ClickHouse服务在ZooKeeper的数量配额使用率超过阈值等相关告警，如果存在，请按照告警指导处理。
3. 执行如下SQL排查是否存在副本同步队列任务积压：

```
select FQDN() as node,type,count() from
clusterAllReplicas(default_cluster, system.replication_queue) group by
node,type;
```

如果存在积压，请查看副本队列中的任务是否报错，并根据报错信息处理。

4. 执行如下SQL排查是否存在节点间表结构不一致。

```
select FQDN(), create_table_query from
clusterAllReplicas(default_cluster,system.tables) where name = '$
{table_name}' group by FQDN(),create_table_query;
```

如果存在，请将不一致的表结构修改一致。

5. 执行如下SQL排查是否存在mutation任务异常：

```
select FQDN(), database, table, mutation_id, create_time, command from
clusterAllReplicas(default_cluster, system.mutations) where is_done = '0'
order by create_time asc;
```

如果mutation任务正常，等待mutation任务完成，如果mutation任务异常，清理异常的mutation任务。

6. 业务写入压力过大导致merge速度小于insert速度。

可以用以下SQL语句检查报错节点最近一小时的写入条数和频次：

```
select tables,written_rows,count() from system.query_log where
type='QueryFinish' and query_start_time between
(toUnixTimestamp(now()) - 3600) AND toUnixTimestamp(now()) and
query_kind = 'Insert' group by tables,written_rows order by written_rows
limit 10;
```

业务上建议一次写入一个分区，写入频率不要太快，不要小批量数据的插入，适当增大每次插入的时间间隔。

7. 如果没有触发Merge，或者Merge较慢，需要调整参数加快Merge。

加速Merge，需要调整如下参数，请参考[加速Merge操作](#)：

配置项	参考值
max_threads	CPU核数*2
background_pool_size	CPU核数
merge_max_block_size	8192的整数倍，根据CPU内存资源大小调整
cleanup_delay_period	适当小于默认值 30

修改 parts_to_throw_insert 值

注意

增大Too many parts的触发阈值，除非特殊场景，不建议修改此配置。此配置在一定程度上起到潜在问题预警的作用，如果集群硬件资源不足，此配置调整不合理，会导致服务潜在问题不能及时发现，可能进一步引起其他故障，恢复难度增加。

- MRS 3.2.0之前版本：登录FusionInsight Manager界面，选择“集群 > ClickHouse > 配置 > 全部配置 > ClickHouseServer > 自定义 > clickhouse-config-customize”，添加[表4-9](#)中参数，保存配置，重启服务。

- MRS 3.2.0及之后版本：登录FusionInsight Manager界面，选择“集群 > ClickHouse > 配置 > 全部配置”，搜索并修改参数“merge_tree.parts_to_throw_insert”的值，保存配置，重启服务。

表 4-9 参数说明

名称	值
merge_tree.parts_to_throw_insert	clickhouse实例内存 / 32GB * 300（保守估计值）

验证修改结果：

登录ClickHouse客户端，执行命令 `select * from system.merge_tree_settings where name = 'parts_to_throw_insert';`

4.7.2 加速 Merge 操作

加速后台任务，需要优先调整Zookeeper服务配置，否则Zookeeper会因为znode等资源不足，导致ClickHouse服务异常，后台任务异常。

1. 调整Zookeeper配置：登录FusionInsight Manager界面，选择“集群 > Zookeeper > 配置 > 全部配置 > quorumpeer > 系统”，修改参数“GC_OPTS”的值，保存配置，滚动重启Zookeeper服务，如下表所示

配置项	参考值	描述
GC_OPTS	Xmx最大内存数 参考值：（ Master 节点内存 - 16GB） * 0.65（保守估计 值）	JVM用于gc的参数。仅当GC_PROFILE 设置为custom时该配置才会生效。需 确保GC_OPT参数设置正确，否则进程 启动会失败。 注意 请谨慎修改该项。如果配置不当，将造成 服务不可用。

2. 调整ClickHouse配置：在FusionInsight Manager界面，选择“集群 > ClickHouse > 配置 > 全部配置 > ClickHouseServer > Zookeeper”，修改如下参数，保存配置，无需重启服务。

配置项	参考值	描述
clickhouse.zookeeper.quota.node.count	Xmx最大内存 数/4GB * 1500000	ClickHouse在ZooKeeper上的顶层目录 的节点数量配额。 数量配额的单位是个，最小值是-1 （无限制），不能等于0。 注意 设置的数量配额值，如果小于当前 ZooKeeper目录的实际值，保存配置可成 功，但是配置值不会生效，并且界面上 报告警。

配置项	参考值	描述
clickhouse.zookeeper.quota.size	Xmx最大内存数/4GB * 1G	ClickHouse在ZooKeeper上的顶层目录的容量配额。 注意 设置的数量配额值，如果小于当前ZooKeeper目录的实际值，保存配置可成功，但是配置值不会生效，并且界面上报告警。

4.7.3 加速 TTL 操作

ClickHouse触发TTL的时候，对CPU和内存会存在较大消耗和占用。

登录FusionInsight Manager界面，选择“集群 > ClickHouse > 配置 > 全部配置 > ClickHouseServer > 自定义 > clickhouse-config-customize”，添加如下配置，保存配置，重启服务。

配置项	参考值	作用
merge_tree.max_replicated_merges_with_ttl_in_queue	CPU核数一半	在ReplicatedMergeTree队列中允许同时使用TTL合并部件的任务数。
merge_tree.max_number_of_merges_with_ttl_in_pool	CPU核数	在ReplicatedMergeTree队列中允许TTL合并部件的线程池。

📖 说明

- 当集群写入压力较大，不建议修改此配置。需要给常规Merge留出空闲线程，避免“Too many parts parts”。
- 对于已存在的大表（数据量超亿级别），不要使用修改或新设置TTL的方式来实现数据老化能力，推荐使用定时任务“alter table table_name on cluster default_cluster drop partition partition_name”来实现。

4.8 ClickHouse 运维管理

4.8.1 ClickHouse 日志介绍

日志描述

日志路径：ClickHouse相关日志的默认存储路径为“\${BIGDATA_LOG_HOME}/clickhouse”。

日志归档规则：ClickHouse日志启动了自动压缩归档功能，缺省情况下，当日志大小超过100MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>.[编号].gz”。默认最多保留最近的10个压缩文件，压缩文件保留个数可以在Manager界面中配置。

表 4-10 ClickHouse 日志列表

日志类型	日志文件名	描述
运行日志	/var/log/Bigdata/clickhouse/clickhouseServer/clickhouse-server.err.log	ClickHouseServer服务运行错误日志文件路径。
	/var/log/Bigdata/clickhouse/clickhouseServer/checkService.log	ClickHouseServer服务运行关键日志文件路径。
	/var/log/Bigdata/clickhouse/clickhouseServer/clickhouse-server.log	
	/var/log/Bigdata/clickhouse/balance/start.log	ClickHouseBalancer服务启动日志文件路径。
	/var/log/Bigdata/clickhouse/balance/error.log	ClickHouseBalancer服务运行错误日志文件路径。
	/var/log/Bigdata/clickhouse/balance/access_http.log	ClickHouseBalancer服务运行日志文件路径。
数据迁移日志	/var/log/Bigdata/clickhouse/migration/数据迁移任务名/clickhouse-copier_{timestamp}_{processId}/copier.log	参考 集群内ClickHouseServer节点间数据迁移 使用迁移工具时产生的运行日志。
	/var/log/Bigdata/clickhouse/migration/数据迁移任务名/clickhouse-copier_{timestamp}_{processId}/copier.err.log	参考 集群内ClickHouseServer节点间数据迁移 使用迁移工具时产生的错误日志。

日志级别

ClickHouse提供了如[表4-11](#)所示的日志级别。

运行日志的级别优先级从高到低分别是error、warning、trace、information、debug，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 4-11 日志级别

级别	描述
error	error表示系统运行的错误信息。
warning	warning表示当前事件处理存在异常信息。
trace	trace表示当前事件处理跟踪信息。
information	information表示记录系统及各事件正常运行状态信息。
debug	debug表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1 登录FusionInsight Manager系统。
- 步骤2 选择“集群 > 服务 > ClickHouse > 配置”。
- 步骤3 单击“全部配置”。
- 步骤4 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤5 选择所需修改的日志级别。
- 步骤6 单击“保存”，然后单击“确定”，成功后配置生效。

----结束

说明

配置完成后即生效，不需要重启服务。

日志格式

ClickHouse的日志格式如下所示：

表 4-12 日志格式

日志类型	格式	示例
ClickHouse运行日志	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level><产生该日志的线程名字> <log中的message> <日志事件的发生位置>	2021.02.23 15:26:30.691301 [6085] {} <Error> DynamicQueryHandler: Code: 516, e.displayText() = DB::Exception: default: Authentication failed: password is incorrect or there is no user with such name, Stack trace (when copying this message, always include the lines below): 0. Poco::Exception::Exception(std::__1::basic_string<char, std::__1::char_traits<char>, std::__1::allocator<char> > const&, int) @ 0x1250e59c

4.8.2 ClickHouse 集群管理

4.8.2.1 ClickHouse 集群配置说明

背景介绍

ClickHouse通过多分片多副本的部署架构实现了集群的高可用，每个集群定义多个分片，每个分片具有2个或2个以上副本。当某节点故障时，分片内其他主机节点上的副本可替代工作，保证服务能正常运行，提高集群的稳定性。

 说明

本章节仅适用于MRS 3.1.0版本。

集群配置

步骤1 登录集群Manager页面，选择“集群 > 服务 > ClickHouse > 配置 > 全部配置”。

步骤2 在“clickhouse-metrika-customize”参数中添加表4-13中自定义配置项。

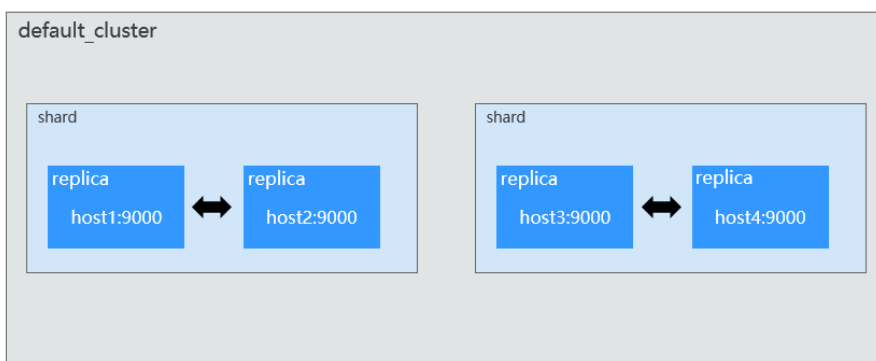
表 4-13 自定义参数

参数	值
clickhouse_remote_servers.example_cluster.shard[1].replica[1].host	host1.9bf17e66-e7ed-4f21-9dfc-34575f955ae6.com
clickhouse_remote_servers.example_cluster.shard[1].replica[1].port	9000
clickhouse_remote_servers.example_cluster.shard[1].replica[2].host	host2.9bf17e66-e7ed-4f21-9dfc-34575f955ae6.com
clickhouse_remote_servers.example_cluster.shard[1].replica[2].port	9000
clickhouse_remote_servers.example_cluster.shard[1].internal_replication	true
clickhouse_remote_servers.example_cluster.shard[2].replica[1].host	host3.9bf17e66-e7ed-4f21-9dfc-34575f955ae6.com
clickhouse_remote_servers.example_cluster.shard[2].replica[1].port	9000
clickhouse_remote_servers.example_cluster.shard[2].replica[2].host	host4.9bf17e66-e7ed-4f21-9dfc-34575f955ae6.com
clickhouse_remote_servers.example_cluster.shard[2].replica[2].port	9000
clickhouse_remote_servers.example_cluster.shard[2].internal_replication	true

步骤3 单击“保存”，保存配置。

----结束

上述集群架构如下图所示：



配置参数说明如下：

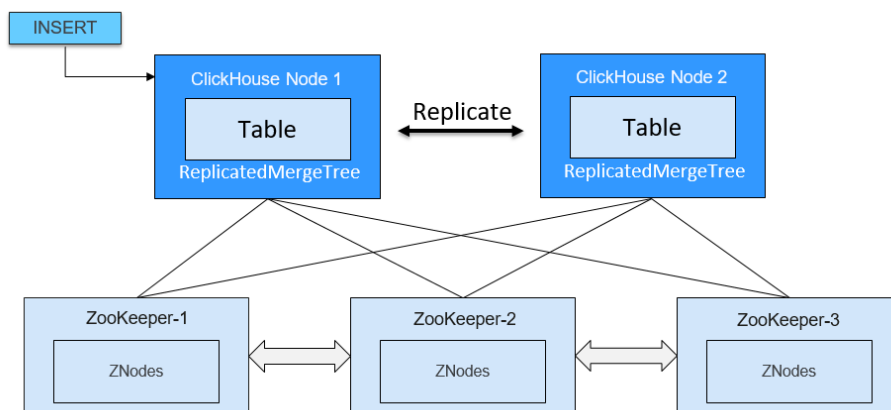
- default_cluster 标签
 - default_cluster 表示当前集群的名称。
 - 当前集群有两个分片 shard，每个 shard 下面有两个副本 replica，每个副本 replica 对应了一个 ClickHouse 实例节点。
 - internal_replication 表示副本间是否为内部复制，当通过集群向分片插入数据时会起作用。
默认配置为 true，表示只向其中的一个副本写入数据（副本间通过复制表来完成同步，能保证数据的一致性）。
如果配置为 false（**不建议配置**），表示向该分片的所有副本中写入相同的数据（副本间数据一致性不强，无法保证完全同步）。
- macros 标签
当前实例节点所在的分片和副本编号，可以用于区别不同的副本。
例如，上述配置对应 host3 节点实例，该实例所在分片编号 shard 为 2，副本编号 replica 为 1。

本章节详细描述了分片和副本信息的配置说明，具体 ClickHouse 集群副本之间如何进行数据同步，详见[副本机制](#)详细说明。

副本机制

ClickHouse 利用 ZooKeeper，通过 ReplicatedMergeTree 引擎（Replicated 系列引擎）实现了副本机制。副本机制是多主架构，可以将 INSERT 语句发送给任意一个副本，分片内其余副本会进行数据的异步复制。

下图中的 Node1 和 Node2 对应为[集群配置](#)中的 host1 和 host2 主机节点。



ClickHouse集群创建成功后，默认会创建3个Zookeeper节点，Zookeeper中存储了ClickHouse的表在复制过程中的元数据信息。

Zookeeper节点信息可以参考config.xml文件内容，具体路径在“\${BIGDATA_HOME}/FusionInsight_ClickHouse_版本号/x_x_clickhouse实例名/etc”目录下。

```
<yandex>
...
<zookeeper>
  <node index="1">
    <host>node-master1lrj.9bf17e66-e7ed-4f21-9dfc-34575f955ae6.com</host>
    <port>2181</port>
  </node>
  <node index="2">
    <port>2181</port>
    <host>node-master2vocd.9bf17e66-e7ed-4f21-9dfc-34575f955ae6.com</host>
  </node>
  <node index="3">
    <host>node-master3xwmu.9bf17e66-e7ed-4f21-9dfc-34575f955ae6.com</host>
    <port>2181</port>
  </node>
</zookeeper>
...
```

集群配置完成后，具体后续怎么使用可以参考[ClickHouse表创建](#)章节描述。

4.8.2.2 ClickHouse 增加磁盘容量

随着业务量的增长，ClickHouse节点数据盘的磁盘容量已不能满足业务需求，需要扩容数据盘磁盘容量。

注意

如果购买MRS集群的计费模式为按需计费，扩容磁盘容量后MRS集群不支持转包周期。

本章节仅适用于MRS 3.1.0版本。

前提条件

- ClickHouse集群和实例状态正常。

- 已评估好要扩容的ClickHouse节点数据盘磁盘容量大小。

扩容数据盘磁盘容量

步骤1 登录MRS控制台，在左侧导航栏选择“现有集群”，单击集群名称。

步骤2 单击“节点管理”，在对应的ClickHouse节点组下，单击要扩容的节点名称，进入到“云硬盘”界面。



步骤3 在对应的数据盘单击“扩容”，进入到扩容磁盘界面。

说明

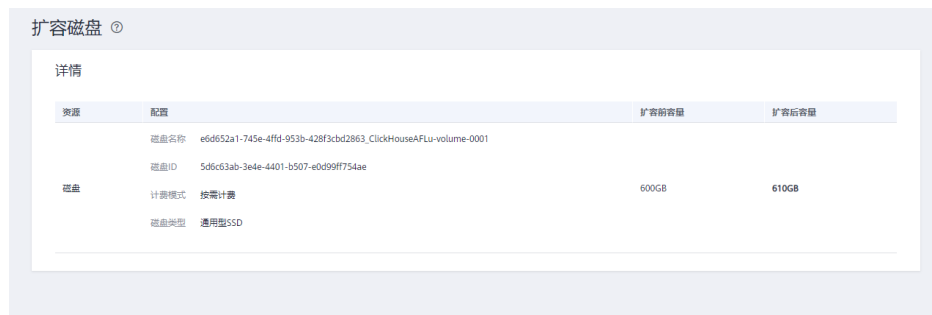
如果当前界面只能看到系统盘，没有数据盘则表示当前ClickHouse节点数据盘暂不支持通过该操作进行扩容。



步骤4 在“新增容量(GB)”参数下修改需要增加的磁盘容量，修改完成后单击“下一步”。



步骤5 按照提示仔细阅读扩容须知，单击“我已阅读，继续扩容”，确认扩容的磁盘容量信息无误后，单击“提交订单”。



步骤6 以root用户登录到ClickHouse的扩容节点上，执行命令：`df -hl`，查看当前已有的数据目录和磁盘分区信息。

```
[root@ClickHouseAFLu ~]# df -hl
Filesystem      Size  Used Avail Use% Mounted on
/dev/vda1       217G   38G  170G  19% /
devtmpfs        32G     0   32G   0% /dev
tmpfs           32G     0   32G   0% /dev/shm
tmpfs           32G    73M   32G   1% /run
tmpfs           32G     0   32G   0% /sys/fs/cgroup
/dev/vda5        9.8G   37M   9.3G   1% /tmp
/dev/vda7        59G   147M   56G   1% /srv/BigData
/dev/vda6        9.8G   583M   8.7G   7% /var
/dev/vda8       177G   154M  168G   1% /var/log
/dev/vdb1       590G   75M  590G   1% /srv/BigData/data1
tmpfs           6.3G     0   6.3G   0% /run/user/2000
```

ClickHouse默认数据目录格式为：“/srv/BigData/dataN”。如上图举例所示，当前ClickHouse数据目录为：“/srv/BigData/data1”，对应分区为：“/dev/vdb1”。

步骤7 执行以下操作使得新扩容的磁盘容量生效。

- 如果是新增分区操作，请执行**步骤8**。新增分区操作是指把扩容的磁盘容量分配给新的分区，并挂载新的ClickHouse数据目录到新增分区下，该操作不会有中断业务的影响。
- 如果是扩大已有分区操作，请执行**步骤15**。扩大已有分区是指把扩容的磁盘容量分配给已存在分区下，操作期间会有中断业务的影响，请谨慎操作，建议操作前先停止业务。

步骤8 新增分区操作请参考**扩容云硬盘分区和文件系统（Linux）**中的“新增MBR分区”或“新增GPT分区”章节进行操作。

步骤9 以root用户登录到ClickHouse的扩容节点上，执行以下命令创建ClickHouse数据目录，为新增分区创建挂载点。目录建议按照当前编号递增。

如当前数据目录为“/srv/BigData/data1”，则新增目录“/srv/BigData/data2”。

```
cd /srv/BigData/  
mkdir data2  
cd data2  
mkdir clickhouse  
cd /srv/BigData/  
chmod 750 -R data2  
chown omm:wheel -R data2
```

步骤10 执行以下命令，挂载新建分区。

```
mount 磁盘分区 挂载目录
```

比如当前新增分区为：“/dev/vdb2”，挂载目录为：“/srv/BigData/data2”，则执行以下命令：

```
mount /dev/vdb2 /srv/BigData/data2
```

📖 说明

弹性云服务器重启后，挂载会失效。您可以修改“/etc/fstab”文件，将新建磁盘分区设置为开机自动挂载，具体请参见[设置开机自动挂载磁盘分区](#)。

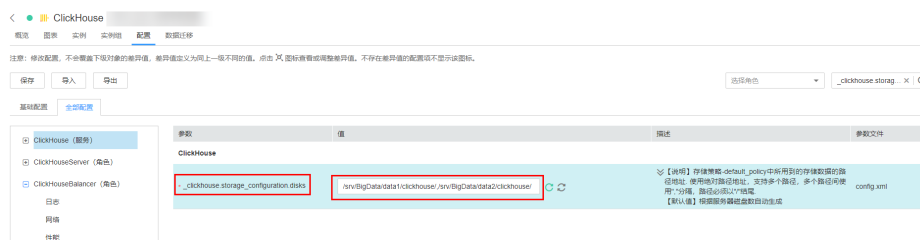
步骤11 参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)，登录FusionInsight Manager。选择“集群 > ClickHouse > 配置 > 全部配置”。

步骤12 搜索“_clickhouse.storage_configuration.disks”，在该配置项下，添加新增的ClickHouse数据目录。

📖 说明

多个目录之间需用“，”分隔，添加的目录以“/”结尾。

例如：在“/srv/BigData/data1/clickhouse/”基础上，添加新增的“/srv/BigData/data2/clickhouse/”目录。添加之后为“/srv/BigData/data1/clickhouse/,/srv/BigData/data2/clickhouse/”。



步骤13 添加完新增目录后，单击“保存”保存配置。单击“概览”，选择“更多 > 同步配置”，单击“确认”完成配置同步。

步骤14 登录到ClickHouse的扩容节点上，进入到以下目录，查看新增的数据目录是否已更新到配置文件中。确认无误后新增分区操作完成。

```
cd ${BIGDATA_HOME}/FusionInsight_ClickHouse_*/x_x_clickhouse实例名/etc  
cat config.xml
```

举例如下，新增的“/srv/BigData/data2/clickhouse/”目录已添加到config.xml中。

```
</trace_log>
<storage_configuration>
  <policies>
    <default>
      <volumes>
        <volume1>
          <disk>disk1</disk>
          <disk>disk2</disk>
        </volume1>
      </volumes>
    </default>
  </policies>
  <disks>
    <disk2>
      <keep_free_space_bytes>104857600</keep_free_space_bytes>
      <path>/srv/BigData/data2/clickhouse/</path>
    </disk2>
    <disk1>
      <keep_free_space_bytes>104857600</keep_free_space_bytes>
      <path>/srv/BigData/data1/clickhouse/</path>
    </disk1>
  </disks>
</storage_configuration>
<access_control_path>/srv/BigData/data1/clickhouse_path/access/</access_control_path>
```

步骤15 如果是扩大已有分区操作，请提前确认ClickHouse业务已停止，否则操作期间会有中断业务的影响。

步骤16 根据**步骤6**确认要扩大的分区，参考[扩容云硬盘分区和文件系统（Linux）](#)中的“扩大已有分区”章节进行操作。

步骤17 扩大已有分区操作完成后，重新执行ClickHouse业务。

----结束

4.8.3 通过数据文件备份恢复 ClickHouse 数据

操作场景

本章节主要介绍通过把ClickHouse中的表数据导出到CSV文件进行备份，后续可以通过备份的CSV文件数据再进行恢复操作。

前提条件

- 已安装ClickHouse客户端。
- 在Manager已创建具有ClickHouse相关表权限的用户。
- 已准备好备份服务器。

备份数据

步骤1 以客户端安装用户，登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建ClickHouse表的权限。如果当前集群未启用Kerberos认证，则无需执行本步骤。

1. 如果是MRS 3.1.0版本集群，则需要先执行：

```
export CLICKHOUSE_SECURITY_ENABLED=true
```

2. **kinit 组件业务用户**

例如，**kinit clickhouseuser**。

步骤5 执行ClickHouse组件的客户端命令，将要备份ClickHouse表数据导出到指定目录下。

```
clickhouse client --host 主机名/实例IP --secure --port 9440 --query="表查询语句" > 输出的csv格式文件路径
```

例如，如下是在ClickHouse实例10.244.225.167下备份test表数据到default_test.csv文件中。

```
clickhouse client --host 10.244.225.167 --secure --port 9440 --query="select * from default.test FORMAT CSV" > /opt/clickhouse/default_test.csv
```

步骤6 将导出的csv数据文件上传至备份服务器。

----结束

恢复数据

步骤1 将备份服务器上的备份数据文件上传到ClickHouse客户端所在目录。

例如，上传default_test.csv备份文件到：/opt/clickhouse目录下。

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤4 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤5 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建ClickHouse表的权限。如果当前集群未启用Kerberos认证，则无需执行本步骤。

1. 如果是MRS 3.1.0版本集群，则需要先执行：

```
export CLICKHOUSE_SECURITY_ENABLED=true
```

2. **kinit 组件业务用户**

例如，**kinit clickhouseuser**。

步骤6 执行ClickHouse组件的客户端命令，登录ClickHouse集群。

```
clickhouse client --host 主机名/实例IP --secure --port 9440
```

步骤7 创建与CSV备份数据文件格式对应的表。

```
CREATE TABLE [IF NOT EXISTS] [database_name.]table_name [ON CLUSTER Cluster名]
```

```
(
```

```
name1 [type1] [DEFAULT|materialized|ALIAS expr1],
```

```
name2 [type2] [DEFAULT|materialized|ALIAS expr2],
```

```
...
```

```
) ENGINE = engine
```

步骤8 将备份数据文件中的内容导入到**步骤7**创建的表中进行数据恢复。

```
clickhouse client --host 主机名/实例IP --secure --port 9440 --query="insert into  
表信息 FORMAT CSV" < csv文件路径
```

例如，如下在ClickHouse实例10.244.225.167中，恢复default_test.csv备份文件数据到test_cpy表中。

```
clickhouse client --host 10.244.225.167 --secure --port 9440 --query="insert  
into default.test_cpy FORMAT CSV" < /opt/clickhouse/default_test.csv
```

----结束

4.8.4 配置 ClickHouse 系统表的生命周期

操作场景

ClickHouse没有默认配置系统表的TTL，长期使用可能会占用大量磁盘空间。

本章节指导用户配置系统表的生命周期，仅保留近30天的数据，减少系统表的磁盘使用。

说明

本章节仅适用于MRS 3.2.0之前版本。

操作步骤

步骤1 使用具有ClickHouse系统表修改权限的用户登录客户端节点，连接到ClickHouse服务端，具体请参考[使用ClickHouse客户端](#)。

步骤2 清理存量系统表数据，执行如下命令：

```
truncate table system.query_log on cluster default_cluster;  
truncate table system.query_thread_log on cluster default_cluster;  
truncate table system.trace_log on cluster default_cluster;
```

步骤3 系统表配置TTL，保留近30天的数据，执行如下命令：

```
alter table system.query_log on cluster default_cluster modify TTL event_time  
+ INTERVAL 30 day;  
alter table system.query_thread_log on cluster default_cluster modify TTL  
event_time + INTERVAL 30 day;  
alter table system.trace_log on cluster default_cluster modify TTL event_time  
+ INTERVAL 30 day;
```

----结束

4.8.5 集群内 ClickHouseServer 节点间数据迁移

ClickHouse数据迁移工具可以将某几个ClickHouseServer实例节点上的一个或多个MergeTree引擎分区表的部分分区迁移至其他ClickHouseServer节点上相同的表中。在

扩容场景中，可以使用该工具将原节点上的部分数据迁移至新增节点上，从而达到扩容后的数据均衡。

前提条件

- ClickHouse服务运行正常，Zookeeper服务运行正常，迁入、迁出节点的ClickHouseServer实例状态正常。
- 请确保迁入节点已有待迁移数据表，且确保该表是MergeTree系列引擎的分区表。
- 创建迁移任务前请确保所有对待迁移数据表的写入任务已停止，且任务启动后，只允许对待迁移数据表进行查询操作，禁止对该表进行写入、删除等操作，否则可能会造成迁移前后数据不一致。
- 迁入节点的ClickHouse数据目录有足够的空间。

操作步骤

步骤1 登录Manager，选择“集群 > 服务 > ClickHouse”，在ClickHouse服务界面单击“数据迁移”页签，进入数据迁移界面。



步骤2 单击“创建迁移任务”。



步骤3 在创建迁移任务界面，填写迁移任务的相关参数，具体参考如下表4-14。

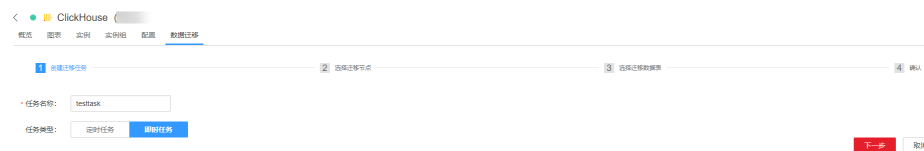
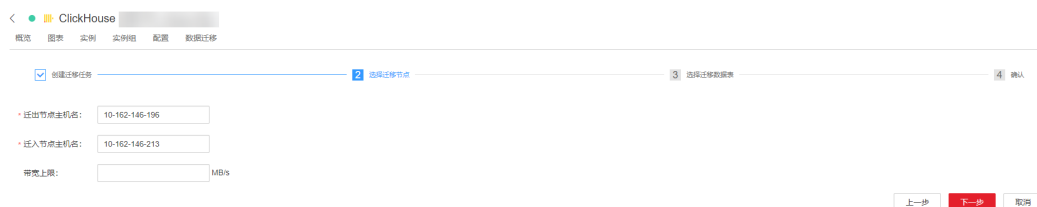


表 4-14 迁移任务参数说明

参数名	参数取值说明
任务名称	填写具体的任务名称。可由字母、数组及下划线组成，长度为1~50位，且不能与已有的迁移任务相同。
任务类型	<ul style="list-style-type: none"> 定时任务：选择定时任务时，可以设置“开始时间”参数，设定任务在当前时间以后的某个时间点执行。 即时任务：任务启动后立即开始执行。
开始时间	在“任务类型”参数选择“定时任务”时填写，有效值为当前时间以后的某个时间（最长为90天以后）。

步骤4 在选择迁移节点界面，填写“迁入节点主机名”、“迁出节点主机名”，单击“下一步”。



说明

- “迁入节点主机名”与“迁出节点主机名”只能各填写一个主机名，不支持多节点迁移。
具体的参数值可以在ClickHouse服务界面单击“实例”页签，查看当前ClickHouseServer实例所在“主机名称”列获取。
- “带宽上限”为可选参数，若不填写则为无上限，最大可设置为10000MB/s。

步骤5 在选择迁移数据表界面，单击“数据库”后的▼，选择待迁出节点上存在的数据库，在“数据表”处选择待迁移的数据表，数据表下拉列表中展示的是所选数据库中的MergeTree系列引擎的分区表。“节点信息”中展示的为当前迁入节点、迁出节点上ClickHouse服务数据目录的空间使用情况，单击“下一步”。



步骤6 确认任务信息，确认无误后可以单击“提交”提交任务。

数据迁移工具将根据待迁移数据表的大小自动计算需要迁移的分区，数据迁移量则是计算出的需要迁移的分区总大小。

步骤7 提交迁移任务成功后，单击操作列的“启动”。如果任务类型是即时任务则开始执行任务，如果是定时任务则开始倒计时。



任务名称	任务类型	数据源	数据库	创建日期	状态	操作
test	即时任务	default	hbs_v1	2021/03/23 16:22:34 GMT+08:00	0%	取消 更多

步骤8 迁移任务执行过程中，可单击“取消”取消正在执行的迁移任务，若取消任务，则会回退掉迁入节点上已迁移的数据。

可以单击“更多 > 详情”查看迁移过程中的日志信息。

步骤9 迁移完成后，选择“更多 > 结果”查看迁移结果；选择“更多 > 删除”清理 ZooKeeper 以及迁出节点上该迁移任务相关的目录。

----结束

4.9 ClickHouse 常用 SQL 语法

4.9.1 CREATE DATABASE 创建数据库

本章节主要介绍 ClickHouse 创建数据库的 SQL 基本语法和使用说明。

基本语法

```
CREATE DATABASE [IF NOT EXISTS] database_name [ON CLUSTER ClickHouse 集群名]
```

说明

ON CLUSTER *ClickHouse 集群名*的语法，使得该 DDL 语句执行一次即可在集群中所有实例上都执行。集群名信息可以使用以下语句的 **cluster** 字段获取：

```
select cluster, shard_num, replica_num, host_name from system.clusters;
```

使用示例

```
--创建数据库名为test的数据库  
CREATE DATABASE test ON CLUSTER default_cluster;  
--创建成功后，通过查询命令验证  
show databases;
```

```
name  
default  
system  
test
```

4.9.2 CREATE TABLE 创建表

本章节主要介绍 ClickHouse 创建表的 SQL 基本语法和使用说明。

基本语法

- 方法一：在指定的“*database_name*”数据库中创建一个名为“*table_name*”的表。

如果建表语句中没有包含“*database_name*”，则默认使用客户端登录时选择的数据库作为数据库名称。

```
CREATE TABLE [IF NOT EXISTS] [database_name.]table_name [ON CLUSTER ClickHouse 集群名]
```

```
(  
  name1 [type1] [DEFAULT|MATERIALIZED|ALIAS expr1],  
  name2 [type2] [DEFAULT|MATERIALIZED|ALIAS expr2],  
  ...  
) ENGINE = engine_name()  
[PARTITION BY expr_list]  
[ORDER BY expr_list]
```

⚠ 注意

ClickHouse在创建表时建议携带**PARTITION BY**创建表分区。因为ClickHouse数据迁移工具是基于表的分区进行数据迁移，在创建表时如果不携带**PARTITION BY**创建表分区，则在[集群内ClickHouseServer节点间数据迁移](#)界面无法对该表进行数据迁移。

- 方法二：创建一个与database_name2.table_name2具有相同结构的表，同时可以对其指定不同的表引擎声明。

如果没有表引擎声明，则创建的表将与database_name2.table_name2使用相同的表引擎。

```
CREATE TABLE [IF NOT EXISTS] [database_name.]table_name AS  
[database_name2.]table_name2 [ENGINE = engine_name]
```

- 方法三：使用指定的引擎创建一个与SELECT子句的结果具有相同结构的表，并使用SELECT子句的结果填充它。

```
CREATE TABLE [IF NOT EXISTS] [database_name.]table_name ENGINE =  
engine_name AS SELECT ...
```

使用示例

```
--在default数据库和default_cluster集群下创建名为test表  
CREATE TABLE default.test ON CLUSTER default_cluster  
(  
  `EventDate` DateTime,  
  `id` UInt64  
)  
ENGINE = ReplicatedMergeTree('/clickhouse/tables/{shard}/default/test', '{replica}')  
PARTITION BY toYYYYMM(EventDate)  
ORDER BY id
```

4.9.3 INSERT INTO 插入表数据

本章节主要介绍ClickHouse插入表数据的SQL基本语法和使用说明。

基本语法

- 方法一：标准格式插入数据。

```
INSERT INTO [database_name.]table [(c1, c2, c3)] VALUES (v11, v12, v13),  
(v21, v22, v23), ...
```
- 方法二：使用SELECT的结果写入。

```
INSERT INTO [database_name.]table [(c1, c2, c3)] SELECT ...
```


使用示例

```
--给test2表插入数据
insert into test2 (id, name) values (1, 'abc'), (2, 'bbbb');
--查询test2表数据
select * from test2;
```

id	name
1	abc
2	bbbb

4.9.4 SELECT 查询表数据

本章节主要介绍ClickHouse查询表数据的SQL基本语法和使用说明。

基本语法

```
SELECT [DISTINCT] expr_list
[FROM [database_name.]table | (subquery) | table_function] [FINAL]
[SAMPLE sample_coeff]
[ARRAY JOIN ...]
[GLOBAL] [ANY|ALL|ASOF] [INNER|LEFT|RIGHT|FULL|CROSS] [OUTER|SEMI|
ANTI] JOIN (subquery)|table (ON <expr_list>)|(USING <column_list>)
[PREWHERE expr]
[WHERE expr]
[GROUP BY expr_list] [WITH TOTALS]
[HAVING expr]
[ORDER BY expr_list] [WITH FILL] [FROM expr] [TO expr] [STEP expr]
[LIMIT [offset_value, ]n BY columns]
[LIMIT [n, ]m] [WITH TIES]
[UNION ALL ...]
[INTO OUTFILE filename]
[FORMAT format]
```

使用示例

```
--查看ClickHouse集群信息
select * from system.clusters;
--显示当前节点设置的宏
select * from system.macros;
--查看数据库容量
select
sum(rows) as "总行数",
formatReadableSize(sum(data_uncompressed_bytes)) as "原始大小",
formatReadableSize(sum(data_compressed_bytes)) as "压缩大小",
round(sum(data_compressed_bytes) / sum(data_uncompressed_bytes) * 100,
0) "压缩率"
from system.parts;
--查询test表容量。where条件根据实际情况添加修改
select
```

```
sum(rows) as "总行数",
formatReadableSize(sum(data_uncompressed_bytes)) as "原始大小",
formatReadableSize(sum(data_compressed_bytes)) as "压缩大小",
round(sum(data_compressed_bytes) / sum(data_uncompressed_bytes) * 100,
0) "压缩率"
from system.parts
where table in ('test')
and partition like '2020-11-%'
group by table;
```

4.9.5 ALTER TABLE 修改表结构

本章节主要介绍ClickHouse修改表结构的SQL基本语法和使用说明。

基本语法

```
ALTER TABLE [database_name].name [ON CLUSTER cluster] ADD|DROP|CLEAR|COMMENT|MODIFY COLUMN ...
```

说明

ALTER仅支持 *MergeTree ， MergeI以及Distributed等引擎表。

使用示例

```
--给表t1增加列test01
ALTER TABLE t1 ADD COLUMN test01 String DEFAULT 'defaultvalue';
--查询修改后的表t1
desc t1
+----+-----+-----+-----+-----+-----+
| name | type | default_type | default_expression | comment | codec_expression |
+----+-----+-----+-----+-----+-----+
| id | UInt8 | | | | |
| name | String | | | | |
| address | String | | | | |
| test01 | String | DEFAULT | 'defaultvalue' | | |
+----+-----+-----+-----+-----+-----+

--修改表t1列name类型为UInt8
ALTER TABLE t1 MODIFY COLUMN name UInt8;
--查询修改后的表t1
desc t1
+----+-----+-----+-----+-----+-----+
| name | type | default_type | default_expression | comment | codec_expression |
+----+-----+-----+-----+-----+-----+
| id | UInt8 | | | | |
| name | UInt8 | | | | |
| address | String | | | | |
| test01 | String | DEFAULT | 'defaultvalue' | | |
+----+-----+-----+-----+-----+-----+

--删除表t1的列test01
ALTER TABLE t1 DROP COLUMN test01;
--查询修改后的表t1
desc t1
+----+-----+-----+-----+-----+-----+
| name | type | default_type | default_expression | comment | codec_expression |
+----+-----+-----+-----+-----+-----+
| id | UInt8 | | | | |
| name | UInt8 | | | | |
| address | String | | | | |
+----+-----+-----+-----+-----+-----+
```

4.9.6 ALTER TABLE 修改表数据

- 建议慎用delete、update的mutation操作
标准SQL的更新、删除操作是同步的，即客户端要等服务端返回执行结果（通常是int值）；而ClickHouse的update、delete是通过异步方式实现的，当执行

update语句时，服务端立即返回执行成功还是失败结果，但是实际上此时数据还没有修改完成，而是在后台排队等着进行真正的修改，可能会出现操作覆盖的情况，也无法保证操作的原子性。

业务场景要求有update、delete等操作，建议使用ReplacingMergeTree、CollapsingMergeTree、VersionedCollapsingMergeTree引擎，使用方式参见：<https://clickhouse.tech/docs/zh/engines/table-engines/mergetree-family/collapsingmergetree/>。

- 建议少或不增删数据列

业务提前规划列个数，如果将来有更多列要使用，可以规划预留多列，避免在生产系统跑业务过程中进行大量的alter table modify列操作，导致不可以预知的性能、数据一致性问题。

4.9.7 DESC 查询表结构

本章节主要介绍ClickHouse查询表结构的SQL基本语法和使用说明。

基本语法

```
DESC|DESCRIBE TABLE [database_name.]table [INTO OUTFILE filename]  
[FORMAT format]
```

使用示例

例如查询表t1的表结构：

```
desc t1;
```

name	type	default_type	default_expression	comment	codec_expression	t1_expression
id	UInt8					
name	UInt8					
address	String					

4.9.8 DROP 删除表

本章节主要介绍ClickHouse删除表的SQL基本语法和使用说明。

基本语法

```
DROP [TEMPORARY] TABLE [IF EXISTS] [database_name.]name [ON CLUSTER  
cluster] [SYNC]
```

使用示例

```
--删除表t1  
drop table t1 SYNC;
```

📖 说明

在删除复制表时，因为复制表需要在Zookeeper上建立一个路径，存放相关数据。ClickHouse默认的库引擎是原子数据库引擎，删除Atomic数据库中的表后，它不会立即删除，而是会在480秒后删除。在删除表时，加上SYNC字段，即可解决该问题，例如：**drop table t1 SYNC;**

删除本地表和分布式表，则不会出现该问题，可不带SYNC字段，例如：**drop table t1;**

4.9.9 SHOW 显示数据库和表信息

本章节主要介绍ClickHouse显示数据库和表信息的SQL基本语法和使用说明。

基本语法

```
show databases
```

```
show tables
```

使用示例

```
--查询数据库  
show databases;
```

```
name  
default  
system  
test
```

```
--查询表信息  
show tables;
```

```
name  
t1  
test  
test2  
test5
```

4.10 ClickHouse 常见问题

4.10.1 在 System.disks 表中查询到磁盘 status 是 fault 或者 abnormal

问题

在System.disks表中查询到磁盘status是fault或者abnormal。

回答

这种情况是由于磁盘存在IO错误，处理方法如下：

- 方法一：登录FusionInsight Manager页面，检查Manager界面上是否磁盘IO异常的告警，如果有，可参考对应的告警帮助文档，通过更换硬盘恢复。
- 方法二：登录FusionInsight Manager页面，重启ClickHouse实例，恢复磁盘状态。

📖 说明

此时磁盘未更换，有IO错误发生时，磁盘状态还会被置为fault或者abnormal。

4.10.2 如何迁移 Hive/HDFS 的数据到 ClickHouse

问题

如何迁移Hive/HDFS的数据到ClickHouse。

回答

可以将Hive中的数据导出为CSV文件，再将CSV文件导入到 ClickHouse。

1. 从Hive中导出数据为 CSV:

```
hive -e "select * from db_hive.student limit 1000"| tr "\t" "," > /data/  
bigdata/hive/student.csv;
```

2. 导入到ClickHouse的default数据库中的student_hive表中，命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。

```
clickhouse --client --port 9002 --password xxx -m --query='INSERT INTO  
default.student_hive FORMAT CSV' < /data/bigdata/hive/student.csv
```

4.10.3 使用辅助 Zookeeper 或者副本数据同步表数据时，日志报错

问题

使用辅助Zookeeper或者副本数据同步表数据时，日志报错:

```
DB::Exception: Cannot parse input: expected 'quorum:' before: 'merge_type: 2'...  
Too many parts (315). Merges are processing significantly slower than inserts...
```

回答

复制表副本版本不一致存在兼容性问题，表结构中有TTL语句，ClickHouse 20.9之后版本新加了TTL_DELETE，之前的版本不识别，高版本复制表副本被选作leader时会出现该问题。

可修改高版本ClickHouse 配置文件config.xml文件做规避，需尽可能保证复制表副本见ClickHouse 版本一致。

4.10.4 如何为 ClickHouse 用户赋予数据库级别的 Select 权限

操作步骤

- 步骤1** 登录到MRS集群装有ClickHouse客户端的节点，执行如下命令：

```
su - omm
```

```
source {客户端安装目录}/bigdata_env
```

```
kinit 组件用户（普通集群无需执行kinit命令）
```

```
clickhouse client --host clickhouse实例节点IP --port 9000 -m --user clickhouse -  
password 'clickhouse用户密码'
```

📖 说明

- 查看ClickHouse用户密码：
登录FusionInsight Manager界面，选择“集群 > 服务 > ClickHouse > 实例”，单击任意ClickHouseServer角色名称。进入ClickHouseServer“概览”页面，单击“配置文件”中的users.xml文件，查看ClickHouse用户密码。
- 命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。

步骤2 创建指定数据库只读权限角色，有如下两种方案：

方案一：

1. 创建指定数据库只读权限角色（以default数据库为例，下同）：
create role ck_role on cluster default_cluster;
GRANT SELECT ON default.* TO ck_role on cluster default_cluster;
2. 创建普通用户
CREATE USER user_01 on cluster default_cluster IDENTIFIED WITH PLAINTEXT_PASSWORD BY 'password';
3. 将只读权限角色赋予普通用户
GRANT ck_role to user_01 on cluster default_cluster;
4. 查看用户权限
show grants for user_01;
select * from system.grants where role_name = 'ck_role';

方案二：

创建指定数据库只读权限用户：

1. 创建用户：
CREATE USER user_01 on cluster default_cluster IDENTIFIED WITH PLAINTEXT_PASSWORD BY 'password';
2. 给用户赋予指定数据库的查询权限：
grant select on default.* to user_01 on cluster default_cluster;
3. 查询用户权限：
select * from system.grants where user_name = 'user_01';

----结束

5 使用 DBService

5.1 DBService 日志介绍

日志描述

日志存储路径：DBService相关日志的默认存储路径为“/var/log/Bigdata/dbservice”。

- gaussDB：“/var/log/Bigdata/dbservice/DB”（gaussDB运行日志目录），“/var/log/Bigdata/dbservice/scriptlog/gaussdbinstall.log”（gaussDB安装日志），“/var/log/gaussdbuninstall.log”（gaussDB卸载日志）。
- HA：“/var/log/Bigdata/dbservice/ha/runlog”（HA运行日志目录），“/var/log/Bigdata/dbservice/ha/scriptlog”（HA脚本日志目录）。
- DBServer：“/var/log/Bigdata/dbservice/healthCheck”（服务进程健康状态检查日志目录）。
“/var/log/Bigdata/dbservice/scriptlog”（运行日志目录），“/var/log/Bigdata/audit/dbservice/”（审计日志目录）。

日志归档规则：DBService的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过1MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>_<编号>.gz”。最多保留最近的20个压缩文件。

说明

日志归档规则用户不能修改。

表 5-1 DBService 日志列表

日志类型	日志文件名	描述
DBServer运行相关日志	dbservice_serviceCheck.log	服务检查脚本运行日志
	dbservice_processCheck.log	进程检查脚本运行日志

日志类型	日志文件名	描述
	backup.log	备份恢复操作运行日志 (需执行DBService备份恢复操作)
	checkHaStatus.log	HA检查日志
	cleanupDBService.log	卸载日志(需执行DBService卸载日志操作)
	componentUserManager.log	数据库用户添加删除操作日志 (需添加依赖DBService的服务)
	install.log	安装日志
	preStartDBService.log	预启动日志
	start_dbserver.log	DBServer启动操作日志 (需执行启动DBService服务的操作)
	stop_dbserver.log	DBServer停止操作日志 (需执行停止DBService服务的操作)
	status_dbserver.log	DBServer状态检查日志 (需执行 \$DBSERVICE_HOME/ sbin/status- dbserver.sh)
	modifyPassword.log	DBService修改密码脚本运行日志
	modifyDBPwd_yyyy-mm-dd.log	修改密码工具运行日志
	dbserver_switchover.log	DBServer执行主备倒换脚本的日志(需执行主备倒换操作)
GAUSSDB运行日志	gaussdb.log	记录数据库运行信息
	gs_ctl-current.log	记录gs_ctl工具的操作
	gs_guc-current.log	记录gs_guc工具的操作,主要是参数修改
	gaussdbinstall.log	gaussDB安装日志
	gaussdbuninstall.log	gaussDB卸载日志
HA脚本相关运行日志	floatip_ha.log	Floatip资源脚本日志
	gaussDB_ha.log	gaussDB资源脚本日志

日志类型	日志文件名	描述
	ha_monitor.log	HA进程监控日志
	send_alarm.log	告警发送日志
	ha.log	HA运行日志
DBService审计日志	dbservice_audit.log	dbservice操作审计日志（例如：备份恢复操作）

日志级别

DBService中提供了如表5-2所示的日志级别。日志级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG。程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 5-2 日志级别

级别	描述
ERROR	ERROR表示当前时间处理存在错误信息。
WARN	WARN表示当前事件处理存在异常信息。
INFO	INFO表示记录系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及系统的调试信息。

日志格式

DBService的日志格式如下所示：

表 5-3 日志格式

日志类型	格式	示例
运行日志	[<yyyy-MM-dd HH:mm:ss> <Log Level>: [<产生该日志的脚本名称: 行号>]: <log中的message>	[2020-12-19 15:56:42] INFO [postinstall.sh:653] Is cloud flag is false. (main)
审计日志	[<yyyy-MM-dd HH:mm:ss,SSS>] UserName:<用户名称> UserIP:<用户IP> Operation:<操作内容> Result:<操作结果> Detail:<具体信息>	[2020-05-26 22:00:23] UserName:omm UserIP:192.168.10.21 Operation:DBService data backup Result: SUCCESS Detail: DBService data backup is successful.

6 使用 Flink

6.1 Flink 作业引擎概述

Flink WebUI提供基于Web的可视化开发平台，用户只需要编写SQL即可开发作业，极大降低作业开发门槛。同时通过作业平台能力开放，支持业务人员自行编写SQL开发作业来快速应对需求，大大减少Flink作业开发工作量。

📖 说明

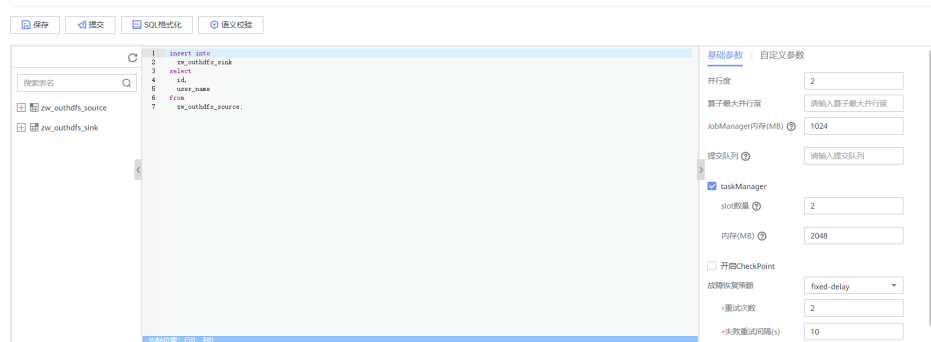
Flink WebUI功能仅支持MRS 3.1.0及之后版本。

Flink WebUI 特点

Flink WebUI主要有以下特点：

- 企业级可视化运维：运维管理界面化、作业监控、作业开发Flink SQL标准化等。

作业管理 > zw_outdfs_3 > 编辑



- 快速建立集群连接：通过集群连接功能配置访问一个集群，需要客户端配置、用户认证密钥文件。
- 快速建立数据连接：通过数据连接功能配置访问一个组件。创建“数据连接类型”为“HDFS”类型时需创建集群连接，其他数据连接类型的“认证类型”为“KERBEROS”需创建集群连接，“认证类型”为“SIMPLE”不需创建集群连接。

📖 说明

“数据连接类型”为“Kafka”时，认证类型不支持“KERBEROS”。

- 可视化开发平台：支持自定义输入/输出映射表，满足不同输入来源、不同输出目标端的需求。
- 图形化作业管理：简单易用。



Flink WebUI 关键能力

Flink WebUI关键能力如表6-1：

表 6-1 Flink WebUI 关键能力

关键能力分类	描述
批流一体	<ul style="list-style-type: none"> ● 支持一套FlinkSQL定义批作业和流作业。
Flink SQL内核能力	<ul style="list-style-type: none"> ● Flink SQL支持自定义大小窗、24小时以内流计算、超出24小时批处理。 ● FlinkSQL支持Kafka、HDFS读取；支持写入Kafka和HDFS。 ● 支持同一个作业定义多个FlinkSQL，多个指标合并在一个作业计算。当一个作业是相同主键、相同的输入和输出时，该作业支持多个窗口的计算。 ● 支持AVG、SUM、COUNT、MAX和MIN统计方法。
Flink SQL可视化定义	<ul style="list-style-type: none"> ● 集群连接管理，配置Kafka、HDFS等服务所属的集群信息。 ● 数据连接管理，配置Kafka、HDFS等服务信息。 ● 数据表管理，定义Sql访问的数据表信息，用于生成DDL语句。 ● FlinkSQL作业定义，根据用户输入的Sql，校验、解析、优化、转换成Flink作业并提交运行。
Flink作业可视化管理	<ul style="list-style-type: none"> ● 支持可视化定义流作业和批作业。 ● 支持作业资源、故障恢复策略、Checkpoint策略可视化配置。 ● 流作业和批作业的状态监控。 ● Flink作业运维能力增强，包括原生监控页面跳转。
性能&可靠性	<ul style="list-style-type: none"> ● 流处理支持24小时窗口聚合计算，毫秒级性能。 ● 批处理支持90天窗口聚合计算，分钟级计算完成。 ● 支持对流处理和批处理的数据进行过滤配置，过滤无效数据。 ● 读取HDFS数据时，提前根据计算周期过滤。 ● 作业定义平台故障、服务降级，不支持再定义作业，但是不影响已有作业计算。 ● 作业故障有自动重启机制，重启策略可配置。

Flink WebUI 应用流程

Flink WebUI应用流程参考如下步骤：

图 6-1 Flink WebUI 应用流程

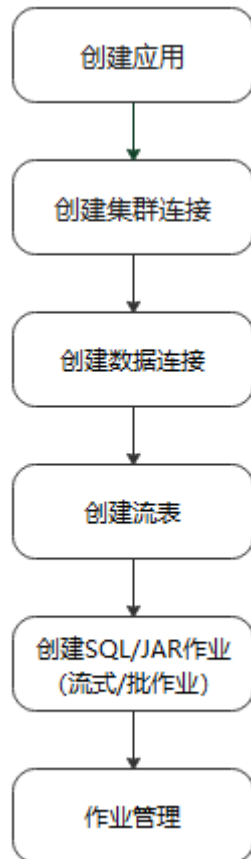


表 6-2 Flink WebUI 应用流程说明

阶段	说明	参考章节
创建应用	通过应用来隔离不同的上层业务。	创建FlinkServer应用
创建集群连接	通过集群连接配置访问不同的集群。	创建FlinkServer集群连接
创建数据连接	通过数据连接，访问不同的数据服务，包括HDFS、Kafka等。	创建FlinkServer数据连接
创建流表	通过数据表，定义源表、维表、输出表的基本属性和字段信息。	创建FlinkServer流表源
创建SQL/JAR作业 (流式/批作业)	定义Flink作业的API，包括Flink SQL和Flink Jar作业。	创建FlinkServer作业
作业管理	管理创建的作业，包括作业启动、开发、停止、删除和编辑等。	创建FlinkServer作业

6.2 Flink 用户权限管理

6.2.1 Flink 安全认证机制说明

Flink 认证和加密

- Flink集群中，各部件支持认证。
 - Flink集群内部各部件和外部部件之间，支持和外部部件如YARN、HDFS、ZooKeeper进行Kerberos认证。
 - Flink集群内部各部件之间，如Flink client和JobManager、JobManager和TaskManager、TaskManager和TaskManager之间支持security cookie认证。
- Flink集群中，各部件支持SSL加密传输；集群内部各部件之间，如Flink client和JobManager、JobManager和TaskManager、TaskManager和TaskManager之间支持SSL加密传输。
详情可参考[配置Flink认证和加密](#)。

ACL 控制

在HA模式下，支持ACL控制。

Flink在HA模式下，支持用ZooKeeper来管理集群和发现服务。ZooKeeper支持SASL ACL控制，即只有通过SASL（kerberos）认证的用户，才有往ZK上操作文件的权限。如果要在Flink上使用SASL ACL控制，需要在Flink配置文件中设置如下配置：

```
high-availability.zookeeper.client.acl: creator  
zookeeper.sasl.disable: false
```

具体配置项介绍请参考[HA](#)。

Web 安全

Flink Web安全加固，支持白名单过滤，Flink Web只能通过YARN代理访问，支持安全头域增强。在Flink集群中，各部件的监测端口支持范围可配置。

- 编码规范：
 - 说明：Web Service客户端和服务端间使用相同的编码方式，是为了防止出现乱码现象，也是实施输入校验的基础。
 - 安全加固：web server响应消息统一采用UTF-8字符编码。
- 支持IP白名单过滤：
 - 说明：防止非法用户登录，需在web server侧添加IP Filter过滤源IP非法的请求。
 - 安全加固：支持IP Filter实现Web白名单配置，配置项是“jobmanager.web.allow-access-address”，默认情况下只支持YARN用户接入。

📖 说明

安装客户端之后需要将客户端节点IP追加到jobmanager.web.allow-access-address配置项中。

- 禁止将文件绝对路径发送到客户端：
 - 说明：文件绝对路径发送到客户端会暴露服务端的目录结构信息，有助于攻击者遍历了解系统，为攻击者攻击提供帮助。
 - 安全加固：Flink配置文件中所有配置项中如果包含以/开头的，则删掉第一级目录。
- 同源策略：

适用于MRS 3.x及之后版本。

 - 说明：如果两个URL的协议，主机和端口均相同，则它们同源；如果不同源，默认不能相互访问；除非被访问者在其服务端显示指定访问者的来源。
 - 安全加固：响应头“Access-Control-Allow-Origin”头域默认配置为YARN集群ResourceManager的IP地址，如果源不是来自YARN的，则不能互相访问。
- 防范敏感信息泄露：

适用于MRS 3.x及之后版本。

 - 说明：带有敏感数据的Web页面都应该禁止缓存，以防止敏感信息泄漏或通过代理服务器上网的用户数据互窜现象。
 - 安全加固：添加“Cache-control”、“Pragma”、“Expires”安全头域，默认值为：“Cache-Control: no-store”，“Pragma : no-cache”，“Expires : 0”。实现了安全加固，Flink和web server交互的内容将不会被缓存。
- 防止劫持：

适用于MRS 3.x及之后版本。

 - 说明：由于点击劫持（ClickJacking）和框架盗链都利用到框架技术，所以需要采用安全措施。
 - 安全加固：添加“X-Frame-Options”安全头域，给浏览器提供允许一个页面可否在“iframe”、“frame”或“object”网站中的展现页面的指示，如果默认配置为“X-Frame-Options: DENY”，则确保任何页面都不能被嵌入到别的“iframe”、“frame”或“object”网站中，从而避免了点击劫持（clickjacking）的攻击。
- 对Web Service接口调用记录日志：

适用于MRS 3.x及之后版本。

 - 说明：对“Flink webmonitor restful”接口调用进行日志记录。
 - 安全加固：“access log”支持配置：“jobmanager.web.accesslog.enable”，默认为“true”。且日志保存在单独的“webaccess.log”文件中。
- 跨站请求（CSRF）伪造防范：

适用于MRS 3.x及之后版本。

 - 说明：在B/S应用中，对于涉及服务器端数据改动（如增加、修改、删除）的操作必须进行跨站请求伪造的防范。跨站请求伪造是一种挟制终端用户在当前已登录的Web应用程序上执行非本意的操作的攻击方法。
 - 安全加固：现有请求修改的接口有2个post，1个delete，其余均是get请求，非get请求的接口均已删除。

- 异常处理：
适用于MRS 3.x及之后版本。
 - 说明：应用程序出现异常时，捕获异常，过滤返回给客户端的信息，并在日志中记录详细的错误信息。
 - 安全加固：默认的错误提示页面，进行信息过滤，并在日志中记录详细的错误信息。新加四个配置项，默认配置为FusionInsight提供的跳转URL，错误提示页面跳转到固定配置的URL中，防止暴露不必要的信息。

表 6-3 四个配置项参数介绍

参数	描述	默认值	是否必选配置
jobmanager.web.403-redirect-url	web403页面，访问若遇到403错误，则会重定向到配置的页面。	-	是
jobmanager.web.404-redirect-url	web404页面，访问若遇到404错误，则会重定向到配置的页面。	-	是
jobmanager.web.415-redirect-url	web415页面，访问若遇到415错误，则会重定向到配置的页面。	-	是
jobmanager.web.500-redirect-url	web500页面，访问若遇到500错误，则会重定向到配置的页面。	-	是

- HTML5安全：
适用于MRS 3.x及之后版本。
 - 说明：HTML5是下一代的Web开发规范，为开发者提供了许多新的功能并扩展了标签。这些新的标签及功能增加了攻击面，存在被攻击的风险（例如跨域资源共享、客户端存储、WebWorker、WebRTC、WebSocket等）。
 - 安全加固：添加“Access-Control-Allow-Origin”配置，如运用到跨域资源共享功能，可对HTTP响应头的“Access-Control-Allow-Origin”属性进行控制。

📖 说明

Flink不涉及如客户端存储、WebWorker、WebRTC、WebSocket等安全风险。

安全声明

- Flink的安全都为开源社区提供和自身研发。有些是需要用户自行配置的安全特性，如认证、SSL传输加密等，这些特性可能对性能和使用方便性造成一定影响。
- Flink作为大数据计算和分析平台，对客户输入的数据是否包含敏感信息无法感知，因此需要客户保证输入数据是脱敏的。
- 客户可以根据应用环境，权衡配置安全与否。
- 任何与安全有关的问题，请联系运维人员。

6.2.2 Flink 用户权限说明

访问并使用Flink WebUI进行业务操作需为用户赋予FlinkServer相关权限，Manager的admin用户没有FlinkServer的业务操作权限。

FlinkServer中应用（租户）是最大管理范围，包含集群连接管理、数据连接管理、应用管理、流表和作业管理等。

FlinkServer中有如表6-4所示三种资源权限：

表 6-4 FlinkServer 资源权限

权限名称	权限描述	备注
FlinkServer 管理员权限	具有所有应用的编辑、查看权限。	是FlinkServer的最高权限。如果已经具有FlinkServer管理员权限，则会自动具备所有应用的权限。
应用编辑权限	具有当前应用编辑权限的用户，可以执行创建、编辑和删除集群连接、数据连接，创建流表、创建作业及运行作业等操作。	同时具有当前应用查看权限。
应用查看权限	具有当前应用查看权限的用户，可以查看应用。	-

6.2.3 创建 FlinkServer 权限角色

该任务指导MRS集群管理员在Manager创建并设置FlinkServer的角色。FlinkServer角色可设置FlinkServer管理员权限以及应用的编辑和查看权限。

用户需要在FlinkServer中对指定的用户设置权限，才能够更新数据、查询数据和删除数据等。

前提条件

集群管理员已根据业务需要规划权限。

操作步骤

步骤1 登录Manager。

步骤2 选择“系统 > 权限 > 角色”。

步骤3 单击“添加角色”，然后在“角色名称”和“描述”输入角色名字与描述。

步骤4 设置角色“配置资源权限”。

FlinkServer权限类型：

- FlinkServer管理员权限：是最高权限，具有FlinkServer所有应用的业务操作权限。
- FlinkServer应用权限：可设置对应用的“应用查看”、“应用编辑”权限。

表 6-5 设置角色

任务场景	角色授权操作
设置FlinkServer管理员权限	在“配置资源权限”的表格中选择“待操作集群的名称 > Flink”，勾选“FlinkServer管理操作权限”。
设置FlinkServer应用权限	在“配置资源权限”的表格中选择“待操作集群的名称 > Flink > FlinkServer应用”，在“权限”列，根据需要勾选“应用查看”或“应用编辑”。

步骤5 单击“确定”完成，返回角色管理。

步骤6 （可选）创建具有FlinkServer相关权限的用户。

FlinkServer角色创建成功后，可创建一个FlinkServer用户，并绑定设置的FlinkServer角色和用户组。

----结束

6.2.4 配置 Flink 对接 Kafka 安全认证

Flink样例工程的数据存储在Kafka组件中。向Kafka组件发送数据（需要有Kafka权限用户），并从Kafka组件接收数据。

步骤1 确保集群安装完成，包括HDFS、Yarn、Flink和Kafka。

步骤2 创建Topic。

- 用户使用Linux命令行创建topic，执行命令前需要使用kinit命令进行人机认证，如 `kinit flinkuser`。

📖 说明

flinkuser需要用户自己创建，并拥有创建Kafka的topic权限。

具体操作请参考[准备开发用户](#)章节。

创建topic的命令格式：`{zkQuorum}`表示ZooKeeper集群信息，格式为IP:port。
`{Topic}`表示Topic名称。

`bin/kafka-topics.sh --create --zookeeper {zkQuorum}/kafka --replication-factor 1 --partitions 5 --topic {Topic}`

例如此处以topic1的数据为例：

```
/opt/client/Kafka/kafka/bin/kafka-topics.sh --create --zookeeper
10.96.101.32:2181,10.96.101.251:2181,10.96.101.177:2181,10.91.8.160:2181/kafka --replication-factor
1 --partitions 5 --topic topic1
```

📖 说明

ZooKeeper集群信息如下：

- ZooKeeper的quorumpeer实例业务IP

ZooKeeper服务所有quorumpeer实例业务IP。登录FusionInsight Manager，选择“集群 > 服务 > ZooKeeper > 实例”，可查看所有quorumpeer实例所在主机业务IP地址。

- ZooKeeper客户端端口号

登录FusionInsight Manager，选择“集群 > 服务 > ZooKeeper”，在“配置”页签查看“clientPort”的值。

- 服务端topic权限配置。

将Kafka的Broker配置参数“allow.everyone.if.no.acl.found”的值修改为“true”。

步骤3 安全认证。

安全认证的方式有三种：Kerberos认证、SSL加密认证和Kerberos+SSL模式认证，用户在使用的时候可任选其中一种方式进行认证。

📖 说明

针对MRS 3.x之前版本，安全认证的方式只支持Kerberos认证。

● Kerberos认证配置

- 客户端配置。

在Flink配置文件“flink-conf.yaml”中，增加kerberos认证相关配置（主要在“contexts”项中增加“KafkaClient”），示例如下：

```
security.kerberos.login.keytab: /home/demo/keytab/flinkuser.keytab
security.kerberos.login.principal: flinkuser
security.kerberos.login.contexts: Client,KafkaClient
security.kerberos.login.use-ticket-cache: false
```

📖 说明

针对MRS 3.x之前版本，配置security.kerberos.login.keytab示例为：/home/demo/flink/release/keytab/flinkuser.keytab。

- 运行参数。

关于“SASL_PLAINTEXT”协议的运行参数示例如下：

```
--topic topic1 --bootstrap.servers 10.96.101.32:21007 --security.protocol SASL_PLAINTEXT --
sasl.kerberos.service.name kafka --kerberos.domain.name hadoop.系统域名.com //
10.96.101.32:21007表示kafka服务器的IP:port
```

● SSL加密配置

- 服务端配置。

登录FusionInsight Manager页面，选择“集群 > 服务 > Kafka > 配置”，参数类别设置为“全部配置”，搜索“ssl.mode.enable”并配置为“true”。

- 客户端配置。

- i. 登录集群的FusionInsight Manager，选择“集群 > 待操作的集群名称 > 服务 > Kafka > 更多 > 下载客户端”，下载客户端压缩文件到本地机器。
- ii. 使用客户端根目录中的“ca.crt”证书文件生成客户端的“truststore”。

执行命令如下：

```
keytool -noprompt -import -alias myservercert -file ca.crt -keystore truststore.jks
```

命令执行结果查看：

```
drwx-----, 5 zgd users 4096 Feb 4 16:22 .
drwxr-xr-x, 10 zgd users 4096 Jan 22 17:38 ..
-rwx-----, 1 zgd users 135 Jan 22 17:31 application.properties
-rwx-----, 1 zgd users 790 Jan 22 17:31 bigdata_env.sample
-rw-----, 1 zgd users 1322 Jan 22 17:31 ca.crt
-rwx-----, 1 zgd users 4508 Jan 22 17:31 conf.py
-rw-----, 1 zgd users 120 Jan 22 17:31 hosts
-rwx-----, 1 zgd users 745 Jan 22 17:31 install.bat
-rwx-----, 1 zgd users 15082 Jan 22 17:31 install.sh
drwx-----, 2 zgd users 4096 Jan 22 17:38 JDK
-rwx-----, 1 zgd users 37021723 Jan 22 17:31 jython-standalone-2.7.0.jar
drwx-----, 5 zgd users 4096 Jan 22 17:38 Kafka
drwx-----, 3 zgd users 4096 Jan 22 17:38 KrbClient
-rwx-----, 1 zgd users 473 Jan 22 17:31 log4j.properties
-rwx-----, 1 zgd users 2107 Jan 22 17:31 README
-rwx-----, 1 zgd users 6949 Jan 22 17:31 refreshConfig.sh
-rwx-----, 1 zgd users 1736 Jan 22 17:31 switchuser.py
-rw-r--r--, 1 root root 1004 Feb 4 16:22 truststore.jks
```

iii. 运行参数。

“ssl.truststore.password”参数内容需要跟创建“truststore”时输入的密码保持一致，执行以下命令运行参数。命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。

```
--topic topic1 --bootstrap.servers 10.96.101.32:9093 --security.protocol SSL --  
ssl.truststore.location /home/zgd/software/FusionInsight_Kafka_ClientConfig/truststore.jks  
--ssl.truststore.password XXX
```

• **Kerberos+SSL模式配置**

完成上文中Kerberos和SSL各自的服务端和客户端配置后，只需要修改运行参数中的端口号和协议类型即可启动Kerberos+SSL模式。

```
--topic topic1 --bootstrap.servers 10.96.101.32:21009 --security.protocol SASL_SSL --  
sasl.kerberos.service.name kafka --ssl.truststore.location --kerberos.domain.name hadoop.系统域  
名.com /home/zgd/software/FusionInsight_Kafka_ClientConfig/truststore.jks --ssl.truststore.password  
XXX
```

---结束

6.2.5 配置 Flink 认证和加密

安全认证

Flink整个系统存在三种认证方式：

- 使用kerberos认证：Flink yarn client与Yarn Resource Manager、JobManager与Zookeeper、JobManager与HDFS、TaskManager与HDFS、Kafka与TaskManager、TaskManager和Zookeeper。
- 使用security cookie进行认证：Flink yarn client与Job Manager、JobManager与TaskManager、TaskManager与TaskManager。
- 使用YARN内部的认证机制：Yarn Resource Manager与Application Master（简称AM）。

说明

- Flink的JobManager与YARN的AM是在同一个进程下。
- 如果用户集群开启Kerberos认证需要使用kerberos认证。
- 针对MRS 3.x之前版本，Flink不支持使用security cookie方式进行认证。

表 6-6 安全认证方式

安全认证方式	说明	配置方法
Kerberos 认证	当前只支持 keytab 认证方式。	<ol style="list-style-type: none">从 FusionInsight Manager 下载用户 keytab，并将 keytab 放到 Flink 客户端所在主机的某个文件夹下。在 “flink-conf.yaml” 上配置：<ol style="list-style-type: none">keytab 路径。 security.kerberos.login.keytab: /home/flinkuser/keytab/abc222.keytab 说明： “/home/flinkuser/keytab/abc222.keytab” 表示的是用户目录。principal 名。 security.kerberos.login.principal: abc222对于 HA 模式，如果配置了 ZooKeeper，还需要设置 ZK kerberos 认证相关的配置。配置如下： zookeeper.sasl.disable: false security.kerberos.login.contexts: Client如果用户对于 Kafka client 和 Kafka broker 之间也需要做 kerberos 认证，配置如下： security.kerberos.login.contexts: Client,KafkaClient

安全认证方式	说明	配置方法
Security Cookie 认证	-	<p>1. 参考签发Flink证书样例章节生成“generate_keystore.sh”脚本并放置在Flink客户端的“bin”目录下，调用“generate_keystore.sh”脚本，生成“Security Cookie”、“flink.keystore”文件和“flink.truststore”文件。 执行sh generate_keystore.sh，输入用户自定义密码。密码不允许包含#。</p> <p>说明 执行脚本后，在Flink客户端的“conf”目录下生成“flink.keystore”和“flink.truststore”文件，并且在客户端配置文件“flink-conf.yaml”中将以下配置项进行了默认赋值。</p> <ul style="list-style-type: none"> • 将配置项“security.ssl.keystore”设置为“flink.keystore”文件所在绝对路径。 • 将配置项“security.ssl.truststore”设置为“flink.truststore”文件所在的绝对路径。 • 将配置项“security.cookie”设置为“generate_keystore.sh”脚本自动生成的一串随机规则密码。 • 默认“flink-conf.yaml”中“security.ssl.encrypt.enabled: false”，“generate_keystore.sh”脚本将配置项“security.ssl.key-password”、“security.ssl.keystore-password”和“security.ssl.truststore-password”的值设置为调用“generate_keystore.sh”脚本时输入的密码。配置文件中包含认证密码信息可能存在安全风险，建议当前场景执行完毕后删除相关配置文件或加强安全管理。 • MRS 3.x及之后版本，如果需要使用密文时，设置“flink-conf.yaml”中“security.ssl.encrypt.enabled: true”，“generate_keystore.sh”脚本不会配置“security.ssl.key-password”、“security.ssl.keystore-password”和“security.ssl.truststore-password”的值，需要使用Manager明文加密API进行获取，执行curl -k -i -u user name:password -X POST -HContent-type:application/json -d '{"plainText":"password"}' 'https://x.x.x.x:28443/web/api/v2/tools/encrypt' 其中user name:password分别为当前系统登录用户名和密码；"plainText"的password为调用“generate_keystore.sh”脚本时的密码；x.x.x.x为集群Manager的浮动IP。命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。 <p>2. 查看是否打开“Security Cookie”开关，即查看配置“flink-conf.yaml”文件中的“security.enable: true”，查看“security cookie”是否已配置成功，例如：</p> <pre>security.cookie: ae70acc9-9795-4c48-ad35-8b5adc8071744f605d1d-2726-432e-88ae-dd39bfec40a9</pre>

安全认证方式	说明	配置方法
		<p>说明</p> <p>用户需要获取SSL证书，放置到Flink客户端中。具体操作可参考签发Flink证书样例。</p> <p>使用MRS客户端预制“generate_keystore.sh”脚本获取SSL证书有效期为5年。参考签发Flink证书样例获取的SSL证书有效期为10年。</p> <p>若要关闭默认的SSL认证方式，需在“flink-conf.yaml”文件中配置“security.ssl.enabled”的值为“false”，并且注释如下参数：security.ssl.key-password、security.ssl.keystore-password、security.ssl.keystore、security.ssl.truststore-password、security.ssl.truststore。</p>
YARN内部认证方式	该方式是YARN内部的认证方式，不需要用户配置。	-

📖 说明

当前一个Flink集群只支持一个用户，一个用户可以创建多个Flink集群。

加密传输

Flink整个系统存在三种加密传输方式：

- 使用Yarn内部的加密传输方式：Flink yarn client与Yarn Resource Manager、Yarn Resource Manager与Job Manager。
- SSL：Flink yarn client与JobManager、JobManager与TaskManager、TaskManager与TaskManager。
- 使用Hadoop内部的加密传输方式：JobManager和HDFS、TaskManager和HDFS、JobManager与ZooKeeper、TaskManager与ZooKeeper。

📖 说明

Yarn内部和Hadoop内部都不需要用户配置加密，用户只需要配置SSL加密传输方式。

配置SSL传输，用户主要在客户端的“flink-conf.yaml”文件中做如下配置：

1. 打开SSL开关和设置SSL加密算法，针对MRS 3.x及之后版本，配置参数如表6-7所示，请根据实际情况修改对应参数值。

表 6-7 参数描述

参数	参数值示例	描述
security.ssl.enabled	true	打开SSL总开关。

参数	参数值示例	描述
akka.ssl.enabled	true	打开akka SSL开关。
blob.service.ssl.enabled	true	打开blob通道SSL开关。
taskmanager.data.ssl.enabled	true	打开taskmanager之间通信的SSL开关。
security.ssl.algorithms	TLS_DHE_RSA_WITH_AES_128_GCM_SHA256,TLS_ECDHE_RSA_WITH_AES_128_GCM_SHA256,TLS_DHE_RSA_WITH_AES_256_GCM_SHA384,TLS_ECDHE_RSA_WITH_AES_256_GCM_SHA384	设置SSL加密的算法。

针对MRS 3.x之前版本，配置参数如表6-8所示。

表 6-8 参数描述

参数	参数值示例	描述
security.ssl.internal.enabled	true	打开内部SSL开关。
akka.ssl.enabled	true	打开akka SSL开关。
blob.service.ssl.enabled	true	打开blob通道SSL开关。
taskmanager.data.ssl.enabled	true	打开taskmanager之间通信的SSL开关。
security.ssl.algorithms	TLS_RSA_WITH_AES128_CBC_SHA256	设置SSL加密的算法。

针对MRS 3.x之前版本，如下参数见表6-9，在MRS的Flink默认配置中不存在，用户如果开启外部连接SSL，则需要添加以下参数。开启外部连接SSL后，因为YARN目前的开源版本无法代理HTTPS请求，所以无法通过YARN代理访问Flink的原生页面，用户可以在集群的同一个VPC下，创建windows虚拟机，在该虚拟机中访问Flink 原生页面。

表 6-9 参数描述

参数	参数值示例	描述
security.ssl.rest.enabled	true	打开外部SSL开关，若该参数配置为“true”，请参考表6-9配置相关参数。

参数	参数值示例	描述
security.ssl.rest.keystore	\${path}/ flink.keystore	keystore的存放路径。
security.ssl.rest.keystore -password	-	keystore的password，-表示需要用户输入自定义设置的密码值。
security.ssl.rest.key- password	-	ssl key的password，-表示需要用户输入自定义设置的密码值。
security.ssl.rest.truststor e	\${path}/ flink.truststore	truststore存放路径。
security.ssl.rest.truststor e-password	-	truststore的password，-表示需要用户输入自定义设置的密码值。

📖 说明

- 如果打开Task Manager之间data传输通道的SSL，对性能会有较大影响，需要用户从安全和性能综合考虑。
 - 配置文件中包含认证密码信息可能存在安全风险，建议当前场景执行完毕后删除相关配置文件或加强安全管理。
2. 参考[签发Flink证书样例](#)章节生成“generate_keystore.sh”脚本并放置在Flink客户端的bin目录下，执行命令**sh generate_keystore.sh <password>**，请参考[配置Flink认证和加密](#)，针对MRS 3.x及之后版本，[表6-10](#)中的配置项会被默认赋值，用户也可以手动配置。

表 6-10 参数描述

参数	参数值示例	描述
security.ssl.keystore	\${path}/ flink.keystore	keystore的存放路径，“flink.keystore”表示用户通过generate_keystore.sh*工具生成的keystore文件名称。
security.ssl.keystore- password	-	keystore的password，-表示需要用户输入自定义设置的密码值。
security.ssl.key- password	-	ssl key的password，-表示需要用户输入自定义设置的密码值。
security.ssl.truststore	\${path}/ flink.truststore	truststore存放路径，“flink.truststore”表示用户通过generate_keystore.sh*工具生成的truststore文件名称。
security.ssl.truststore- password	-	truststore的password，-表示需要用户输入自定义设置的密码值。

针对MRS 3.x之前版本，*generate_keystore.sh*不需手动生成，[表6-11](#)中的配置项会被默认赋值，用户也可以手动配置。

表 6-11 参数描述

参数	参数值示例	描述
security.ssl.internal.keystore	\${path}/flink.keystore	keystore的存放路径，“flink.keystore”表示用户通过generate_keystore.sh*工具生成的keystore文件名称。
security.ssl.internal.keystore-password	-	keystore的password，表示需要用户输入自定义设置的密码值。
security.ssl.internal.keystore-key-password	-	ssl key的password，表示需要用户输入自定义设置的密码值。
security.ssl.internal.truststore	\${path}/flink.truststore	truststore存放路径，“flink.truststore”表示用户通过generate_keystore.sh*工具生成的truststore文件名称。
security.ssl.internal.truststore-password	-	truststore的password，表示需要用户输入自定义设置的密码值。

针对MRS 3.x之前版本，如果开启外部连接SSL，即 security.ssl.rest.enabled 配置为 true，则如下参数见[表6-12](#)，用户需要配置。

表 6-12 参数说明

参数	参数值示例	描述
security.ssl.rest.enabled	true	打开外部SSL开关，若该参数配置为“true”，请参考 表6-12 配置相关参数。
security.ssl.rest.keystore	\${path}/flink.keystore	keystore的存放路径
security.ssl.rest.keystore-password	-	keystore的password，表示需要用户输入自定义设置的密码值。
security.ssl.rest.key-password	-	ssl key的password，表示需要用户输入自定义设置的密码值。
security.ssl.rest.truststore	\${path}/flink.truststore	truststore存放路径

参数	参数值示例	描述
security.ssl.rest.truststore-password	-	truststore的password，表示需要用户输入自定义设置的密码值。

📖 说明

“path”目录是用来存放SSL keystore、truststore相关配置文件，该目录是由用户自定义创建。

3. 配置客户端访问keystore或truststore文件路径。

- 相对路径（推荐）

请执行如下步骤配置“flink.keystore”和“flink.truststore”文件路径为相对路径，并确保Flink客户端执行命令的目录可以直接访问该相对路径。

i. 在Flink客户端“conf”目录下新建目录，例如ssl。

```
cd /Flink客户端目录/Flink/flink/conf/
```

```
mkdir ssl
```

ii. 移动“flink.keystore”和“flink.truststore”文件到新建目录中。

```
mv flink.keystore ssl/
```

```
mv flink.truststore ssl/
```

iii. 修改“flink-conf.yaml”文件中如下两个参数为相对路径。

```
vi /Flink客户端目录/Flink/flink/conf/flink-conf.yaml
```

```
security.ssl.keystore: ssl/flink.keystore
```

```
security.ssl.truststore: ssl/flink.truststore
```

- 绝对路径

执行“generate_keystore.sh”脚本后，默认在“flink-conf.yaml”文件中将“flink.keystore”和“flink.truststore”文件路径自动配置为绝对路径，此时需要将“conf”目录中的“flink.keystore”和“flink.truststore”文件分别放置在Flink客户端以及Yarn各个节点的该绝对路径上。

6.3 Flink 客户端使用实践

本节提供使用Flink运行wordcount作业的操作指导。

前提条件

- MRS集群中已安装Flink组件。
- 集群正常运行，已安装集群客户端，例如安装目录为“/opt/hadoopclient”。以下操作的客户端目录只是举例，请根据实际安装目录修改。

使用 Flink 客户端（MRS 3.x 及之后版本）

步骤1 安装客户端。

以在集群内节点安装客户端为例：

1. 登录Manager，在“集群”下拉列表中单击需要操作的集群名称，选择“更多 > 下载客户端”，弹出“下载集群客户端”信息提示框。
2. 选择“完整客户端”，选择与待安装节点架构相匹配的平台类型，勾选“仅保存到如下路径”，单击“确定”开始生成客户端文件。
 - 文件生成后默认保存在主管理节点“/tmp/FusionInsight-Client”。
 - 客户端软件包名称格式为：“FusionInsight_Cluster_集群ID_Services_Client.tar”。本章节仅以集群ID为1进行介绍，请以实际集群ID为准。
3. 以客户端安装用户登录将要安装客户端的服务器。
4. 进入安装包所在目录，执行如下命令解压软件包。

```
cd /tmp/FusionInsight-Client
```

```
tar -xvf FusionInsight_Cluster_1_Services_Client.tar
```

5. 执行如下命令校验解压得到的文件，检查回显信息与sha256文件里面的内容是否一致，例如：

```
sha256sum -c FusionInsight_Cluster_1_Services_ClientConfig.tar.sha256
```

```
FusionInsight_Cluster_1_Services_ClientConfig.tar: OK
```

6. 解压获取的安装文件。

```
tar -xvf FusionInsight_Cluster_1_Services_ClientConfig.tar
```

7. 进入安装包所在目录，执行如下命令安装客户端到指定目录（绝对路径），例如安装到“/opt/hadoopclient”目录。

```
cd /tmp/FusionInsight-Client/
```

```
FusionInsight_Cluster_1_Services_ClientConfig
```

```
./install.sh /opt/hadoopclient
```

等待客户端安装完成（以下只显示部分屏显结果）。

```
The component client is installed successfully
```

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，切换到客户端安装目录。

```
cd /opt/hadoopclient
```

步骤4 执行如下命令初始化环境变量。

```
source /opt/hadoopclient/bigdata_env
```

步骤5 若集群开启Kerberos认证，需要执行以下步骤，若集群未开启Kerberos认证请跳过该步骤。

1. 准备一个提交Flink作业的用户。

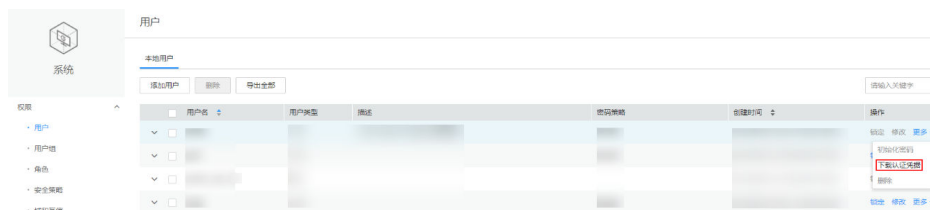
登录Manager，选择“系统 > 权限 > 角色”，单击“添加角色”，输入角色名称与描述。在“配置资源权限”的表格中选择“待操作集群的名称 > Flink”，勾选“FlinkServer管理操作权限”，单击“确定”，返回角色管理。

选择“系统 > 权限 > 用户”，单击“添加用户”，输入用户名、密码等，用户类型选择“人机”，用户组根据需求添加“hadoop”、“yarnviewgroup”和“hadooppmanager”，并添加“System_administrator”、“default”和创建的角色，单击“确定”完成Flink作业用户创建（首次创建的用户需使用该用户登录Manager修改密码）。

2. 登录Manager，下载认证凭据。

登录集群的Manager界面，具体请参见[访问FusionInsight Manager（MRS 3.x 及之后版本）](#)，选择“系统 > 权限 > 用户”，在已增加用户所在行的“操作”列，选择“更多 > 下载认证凭据”。

图 6-2 下载认证凭据



3. 将下载认证凭据压缩包解压，并将得到的文件拷贝到客户端节点中，例如客户端节点的“/opt/hadoopclient/Flink/flink/conf”目录下。如果是在集群外节点安装的客户端，需要将得到的文件拷贝到该节点的“/etc/”目录下。
4. 将客户端安装节点的业务IP和所有Master节点IP添加到配置文件“/opt/hadoopclient/Flink/flink/conf/flink-conf.yaml”中的“jobmanager.web.access-control-allow-origin”和“jobmanager.web.allow-access-address”配置项中，IP地址之间使用英文逗号分隔。

```
jobmanager.web.access-control-allow-origin: xx.xx.xxx.xxx,xx.xx.xxx.xxx,xx.xx.xxx.xxx  
jobmanager.web.allow-access-address: xx.xx.xxx.xxx,xx.xx.xxx.xxx,xx.xx.xxx.xxx
```

说明

客户端安装节点的业务IP获取方法：

- 集群内节点：

登录MapReduce服务管理控制台，选择“现有集群”，选中当前的集群并单击集群名，进入集群信息页面。

在“节点管理”中查看安装客户端所在的节点IP。

- 集群外节点：安装客户端所在的弹性云服务器的IP。

5. 配置安全认证，在“/opt/hadoopclient/Flink/flink/conf/flink-conf.yaml”配置文件中的对应配置添加keytab路径以及用户名。

```
security.kerberos.login.keytab: <user.keytab文件路径>  
security.kerberos.login.principal: <用户名>
```

例如：

```
security.kerberos.login.keytab: /opt/hadoopclient/Flink/flink/conf/user.keytab  
security.kerberos.login.principal: test
```

6. 在Flink的客户端bin目录下，执行如下命令进行安全加固，并设置一个用于提交作业密码。

```
cd /opt/hadoopclient/Flink/flink/bin
```

```
sh generate_keystore.sh
```

该脚本会自动替换“/opt/hadoopclient/Flink/flink/conf/flink-conf.yaml”中关于SSL的值。

📖 说明

执行认证和加密后会在Flink客户端的“conf”目录下生成“flink.keystore”和“flink.truststore”文件，并且在客户端配置文件“flink-conf.yaml”中将以下配置项进行了默认赋值：

- 将配置项“security.ssl.keystore”设置为“flink.keystore”文件所在绝对路径。
- 将配置项“security.ssl.truststore”设置为“flink.truststore”文件所在的绝对路径。
- 将配置项“security.cookie”设置为“generate_keystore.sh”脚本自动生成的一串随机规则密码。
- 默认“flink-conf.yaml”中“security.ssl.encrypt.enabled: false”，“generate_keystore.sh”脚本将配置项“security.ssl.key-password”、“security.ssl.keystore-password”和“security.ssl.truststore-password”的值设置为调用“generate_keystore.sh”脚本时输入的密码。配置文件中包含认证密码信息可能存在安全风险，建议当前场景执行完毕后删除相关配置文件或加强安全管理。
- MRS 3.x及之后版本，如果需要使用密文时，设置“flink-conf.yaml”中“security.ssl.encrypt.enabled: true”，“generate_keystore.sh”脚本不会配置“security.ssl.key-password”、“security.ssl.keystore-password”和“security.ssl.truststore-password”的值，需要使用Manager明文加密API进行获取，执行`curl -k -i -u user name:password -X POST -HContent-type:application/json -d '{"plainText":"password"}' 'https://x.x.x.x:28443/web/api/v2/tools/encrypt'`

其中`user name:password`分别为当前系统登录用户名和密码；“plainText”的password为调用“generate_keystore.sh”脚本时的密码；x.x.x.x为集群Manager的浮动IP。命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。

7. 配置客户端访问flink.keystore和flink.truststore文件的路径。

- 相对路径（推荐）：

执行如下步骤配置flink.keystore和flink.truststore文件路径为相对路径，并确保Flink Client执行命令的目录可以直接访问该相对路径。

- i. 在“/opt/hadoopclient/Flink/flink/conf/”目录下新建目录，例如ssl。

```
cd /opt/hadoopclient/Flink/flink/conf/  
mkdir ssl
```

- ii. 移动flink.keystore和flink.truststore文件到“/opt/hadoopclient/Flink/flink/conf/ssl/”中。

```
mv flink.keystore ssl/  
mv flink.truststore ssl/
```

- iii. 修改flink-conf.yaml文件中如下两个参数为相对路径。

```
security.ssl.keystore: ssl/flink.keystore  
security.ssl.truststore: ssl/flink.truststore
```

- 绝对路径：

执行“generate_keystore.sh”脚本后，在flink-conf.yaml文件中将flink.keystore和flink.truststore文件路径自动配置为绝对路径“/opt/hadoopclient/Flink/flink/conf/”，此时需要将conf目录中的flink.keystore和flink.truststore文件分别放置在Flink Client以及Yarn各个节点的该绝对路径上。

步骤6 运行wordcount作业。

须知

用户在Flink提交作业或者运行作业时，需根据涉及的相关服务（如HDFS、Kafka等）是否启用Ranger鉴权，使该用户应具有如下权限：

- 如果启用Ranger鉴权，当前用户必须属于hadoop组或者已在Ranger中为该用户添加“/flink”的读写权限。
- 如果停用Ranger鉴权，当前用户必须属于hadoop组。

- 普通集群（未开启Kerberos认证）可通过如下两种方式提交作业：

- 执行如下命令启动session，并在session中提交作业。

```
yarn-session.sh -nm "session-name" -d
```

```
flink run /opt/hadoopclient/Flink/flink/examples/streaming/  
WordCount.jar
```

- 执行如下命令在Yarn上提交单个作业。

```
flink run -m yarn-cluster /opt/hadoopclient/Flink/flink/examples/  
streaming/WordCount.jar
```

- 安全集群（开启Kerberos认证）根据flink.keystore和flink.truststore文件的路径有如下两种方式提交作业：

- flink.keystore和flink.truststore文件路径为相对路径时：

- 在“ssl”的同级目录下执行如下命令启动session，并在session中提交作业。

其中“ssl”是相对路径，如“ssl”所在目录是“opt/hadoopclient/Flink/flink/conf/”，则在“opt/hadoopclient/Flink/flink/conf/”目录下执行命令。

```
cd /opt/hadoopclient/Flink/flink/conf
```

```
yarn-session.sh -t ssl/ -nm "session-name" -d
```

```
flink run /opt/hadoopclient/Flink/flink/examples/streaming/  
WordCount.jar
```

- 执行如下命令在Yarn上提交单个作业。

```
cd /opt/hadoopclient/Flink/flink/conf
```

```
flink run -m yarn-cluster -yt ssl/ /opt/hadoopclient/Flink/flink/  
examples/streaming/WordCount.jar
```

- flink.keystore和flink.truststore文件路径为绝对路径时：

- 执行如下命令启动session，并在session中提交作业。

```
cd /opt/hadoopclient/Flink/flink/conf
```

```
yarn-session.sh -nm "session-name" -d
```

```
flink run /opt/hadoopclient/Flink/flink/examples/streaming/  
WordCount.jar
```

- 执行如下命令在Yarn上提交单个作业。

```
flink run -m yarn-cluster /opt/hadoopclient/Flink/flink/  
examples/streaming/WordCount.jar
```

步骤7 作业提交成功后，客户端界面显示如下。

图 6-3 在 Yarn 上提交作业成功

```
[root@node-master1kz2p ~]# flink run -a yarn-cluster /opt/client/flink/flink/examples/streaming/WordCount.jar
2019-07-10 16:20:11,090 | WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
2019-07-10 16:20:11,098 | WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
Starting execution of program
Executing WordCount example with default input data set.
Use --input to specify file input.
Printing result to stdout. Use --output to specify output path.
Program execution finished.
Job with JobID c043b1921e80afe2bb24b51a5beid has finished.
Job Runtime: 7953 ms
```

图 6-4 启动 session 成功

```
[root@node-master1kz2p hive]# yarn session.sh -m "test@hadoop" -d
2019-07-26 09:17:58,919 | WARN | [main] | Unable to load native-hadoop library for your platform... using builtin-java classes where applicable | org.apache.hadoop.util.NativeCodeLoader (NativeCodeLoader.java:62)
2019-07-26 09:17:58,988 | WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
Flink JobManager is now running on node-ana-corehdp:32586 with leader id b9b5a8-1983-435f-bb90-ad128fd1d46b.
JobManager Web Interfaces: https://192.168.2.01:47097
[root@node-master1kz2p hive]#
```

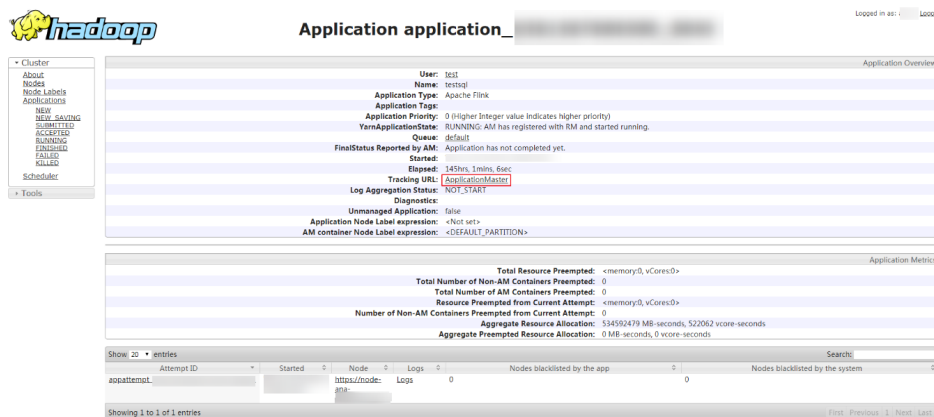
图 6-5 在 session 中提交作业成功

```
[root@node-master1kz2p hive]# flink run /opt/client/flink/flink/examples/streaming/WordCount.jar
WARN properties set default parallelism to 2
2019-07-26 09:19:20,548 | WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
2019-07-26 09:19:20,548 | WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
Starting execution of program
Executing WordCount example with default input data set.
Use --input to specify file input.
Printing result to stdout. Use --output to specify output path.
Program execution finished.
Job with JobID 5bdbc18d6563f3d792a19163c2e7c3c3 has finished.
Job Runtime: 5906 ms
[root@node-master1kz2p hive]#
```

步骤8 使用运行用户进入Yarn服务的原生页面，具体操作参考[查看Flink作业信息](#)，找到对应作业的application，单击application名称，进入到作业详情页面

- 若作业尚未结束，可单击“Tracking URL”链接进入到Flink的原生页面，查看作业的运行信息。
- 若作业已运行结束，对于在session中提交的作业，可以单击“Tracking URL”链接登录Flink原生页面查看作业信息。

图 6-6 application



----结束

使用 Flink 客户端（MRS 3.x 之前版本）

步骤1 安装客户端。

以在Core节点安装客户端为例：

1. 登录MRS Manager页面，选择“服务管理 > 下载客户端”下载客户端安装包至主管理节点。

2. 使用IP地址搜索主管理节点并使用VNC登录主管理节点。
3. 在主管理节点，执行以下命令切换用户。

```
sudo su - omm
```

4. 在MRS管理控制台，查看指定集群“节点管理”页面的“IP”地址。
记录需使用客户端的Core节点IP地址。
5. 在主管理节点，执行以下命令，将客户端安装包从主管理节点文件拷贝到当前Core节点：

```
scp -p /tmp/MRS-client/MRS_Services_Client.tar Core节点的IP地址:/opt/client
```

6. 使用“root”登录Core节点。
Master节点支持Cloud-Init特性，Cloud-init预配置的用户名“root”，密码为创建集群时设置的密码。
7. 执行以下命令，安装客户端：

```
cd /opt/client
tar -xvf MRS_Services_Client.tar
tar -xvf MRS_Services_ClientConfig.tar
cd /opt/client/MRS_Services_ClientConfig
./install.sh 客户端安装目录
例如，执行命令：
./install.sh /opt/hadoopclient
```

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，切换到客户端安装目录。

```
cd /opt/hadoopclient
```

步骤4 执行如下命令初始化环境变量。

```
source /opt/hadoopclient/bigdata_env
```

步骤5 若集群开启Kerberos认证，需要执行以下步骤，若集群未开启Kerberos认证请跳过该步骤。

1. 准备一个提交Flink作业的用户。

登录MRS Manager，选择“系统设置 > 角色管理 > 添加角色”，添加角色例如 **flinkrole**。在“权限”的表格中选择“HDFS > File System > hdfs://hacluster/”，勾选“Read”、“Write”和“Execute”，单击“权限”表格中“服务”返回。选择“Yarn > Scheduler Queue > root”，勾选default的“Submit”，单击“确定”保存。

选择“系统设置 > 用户组管理 > 添加用户组”，为样例工程创建一个用户组，例如 **flinkgroup**。选择“系统设置 > 用户管理 > 添加用户”，为样例工程创建一个用户。填写用户名例如 **flinkuser**，用户类型为“人机”用户，加入用户组 **flinkgroup**和**hadoop**，并绑定角色 **flinkrole**取得权限，单击“确定”（首次创建的用户需使用该用户登录MRS Manager修改密码）。

2. 登录Manager，下载认证凭据。

登录集群的Manager界面，具体请参见[访问MRS Manager（MRS 3.x之前版本）](#)，选择“系统设置 > 用户管理”，在已增加用户所在行的“操作”列，选择“更多 > 下载认证凭据”。

图 6-7 下载认证凭据



3. 将下载的认证凭据压缩包解压缩，并将得到的文件拷贝到客户端节点中，例如客户端节点的“/opt/hadoopclient/Flink/flink/conf”目录下。如果是在集群外节点安装的客户端，需要将得到的文件拷贝到该节点的“/etc/”目录下。
4. 配置安全认证，在“/opt/hadoopclient/Flink/flink/conf/flink-conf.yaml”配置文件中的对应配置添加keytab路径以及用户名。

```
security.kerberos.login.keytab: <user.keytab文件路径>  
security.kerberos.login.principal: <用户名>
```

例如：

```
security.kerberos.login.keytab: /opt/hadoopclient/Flink/flink/conf/user.keytab  
security.kerberos.login.principal: test
```

5. 参考[签发证书样例](#)章节生成“generate_keystore.sh”脚本并放置在Flink的客户端bin目录下，执行如下命令进行安全加固，并设置一个用于提交作业密码。

```
cd /opt/hadoopclient/Flink/flink/bin
```

```
sh generate_keystore.sh
```

该脚本会自动替换“/opt/hadoopclient/Flink/flink/conf/flink-conf.yaml”中关于SSL的值，针对MRS2.x及之前版本，安全集群默认没有开启外部SSL，用户如果需要启用外部SSL，请参考“认证和加密”章节进行配置后再次运行该脚本即可。

📖 说明

- generate_keystore.sh脚本无需手动生成。
 - 执行认证和加密后会将生成的flink.keystore、flink.truststore、security.cookie自动填充到“flink-conf.yaml”对应配置项中。
6. 配置客户端访问flink.keystore和flink.truststore文件的路径。
 - 相对路径（推荐）：

执行如下步骤配置flink.keystore和flink.truststore文件路径为相对路径，并确保Flink Client执行命令的目录可以直接访问该相对路径。

 - i. 在“/opt/hadoopclient/Flink/flink/conf/”目录下新建目录，例如ssl。

```
cd /opt/hadoopclient/Flink/flink/conf/  
mkdir ssl
```
 - ii. 移动flink.keystore和flink.truststore文件到“/opt/hadoopclient/Flink/flink/conf/ssl/”中。

```
mv flink.keystore ssl/  
mv flink.truststore ssl/
```
 - iii. 修改flink-conf.yaml文件中如下两个参数为相对路径。

```
security.ssl.internal.keystore: ssl/flink.keystore  
security.ssl.internal.truststore: ssl/flink.truststore
```
 - 绝对路径：

执行“generate_keystore.sh”脚本后，在flink-conf.yaml文件中将flink.keystore和flink.truststore文件路径自动配置为绝对路径“/opt/

hadoopclient/Flink/flink/conf/”，此时需要将conf目录中的flink.keystore和flink.truststore文件分别放置在Flink Client以及Yarn各个节点的该绝对路径上。

步骤6 运行wordcount作业。

须知

用户在Flink提交作业或者运行作业时，需根据涉及的相关服务（如HDFS、Kafka等）是否启用Ranger鉴权，使该用户应具有如下权限：

- 如果启用Ranger鉴权，当前用户必须属于hadoop组或者已在Ranger中为该用户添加“/flink”的读写权限。
 - 如果停用Ranger鉴权，当前用户必须属于hadoop组。
-
- 普通集群（未开启Kerberos认证）可通过如下两种方式提交作业：
 - 执行如下命令启动session，并在session中提交作业。

```
yarn-session.sh -nm "session-name" -d
flink run /opt/hadoopclient/Flink/flink/examples/streaming/WordCount.jar
```
 - 执行如下命令在Yarn上提交单个作业。

```
flink run -m yarn-cluster /opt/hadoopclient/Flink/flink/examples/streaming/WordCount.jar
```
 - 安全集群（开启Kerberos认证）根据flink.keystore和flink.truststore文件的路径有如下两种方式提交作业：
 - flink.keystore和flink.truststore文件路径为相对路径时：
 - 在“ssl”的同级目录下执行如下命令启动session，并在session中提交作业。
其中“ssl”是相对路径，如“ssl”所在目录是“opt/hadoopclient/Flink/flink/conf/”，则在“opt/hadoopclient/Flink/flink/conf/”目录下执行命令。

```
cd /opt/hadoopclient/Flink/flink/conf
yarn-session.sh -t ssl/ -nm "session-name" -d
flink run /opt/hadoopclient/Flink/flink/examples/streaming/WordCount.jar
```
 - 执行如下命令在Yarn上提交单个作业。

```
cd /opt/hadoopclient/Flink/flink/conf
flink run -m yarn-cluster -yt ssl/ /opt/hadoopclient/Flink/flink/examples/streaming/WordCount.jar
```
 - flink.keystore和flink.truststore文件路径为绝对路径时：
 - 执行如下命令启动session，并在session中提交作业。

```
cd /opt/hadoopclient/Flink/flink/conf
yarn-session.sh -nm "session-name" -d
flink run /opt/hadoopclient/Flink/flink/examples/streaming/WordCount.jar
```

- 执行如下命令在Yarn上提交单个作业。
**flink run -m yarn-cluster /opt/hadoopclient/Flink/flink/
examples/streaming/WordCount.jar**

步骤7 作业提交成功后，客户端界面显示如下。

图 6-8 在 Yarn 上提交作业成功

```
[root@node-master1ks2p ~]# flink run -m yarn-cluster /opt/client/Flink/flink/examples/streaming/WordCount.jar
2019-07-10 16:30:11,090 WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
2019-07-10 16:30:11,090 WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
Starting execution of program
Executing WordCount example with default input data set.
Use --input to specify file input.
Printing result to stdout. Use --output to specify output path.
Program execution finished
Job with JobID c9c3b192ee0a1efc2bb24b51a5beid has finished.
Job Runtime: 7953 ms
```

图 6-9 启动 session 成功

```
[root@node-master1ks2p ~]# yarn-session.sh -nm "test4doc" -d
2019-07-26 09:17:08,919 WARN | [main] | Unable to load native-hadoop library for your platform... using builtin-java classes where applicable | org.apache.hadoop.util.NativeCodeLoader (NativeCodeLoader.java:102)
2019-07-26 09:17:08,986 WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
2019-07-26 09:19:20,648 WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
flink JobManager is now running on node-ana-corehdp:32586 with leader id b9b5a88-1983-435f-bb90-ad128fd1d46b.
JobManager Web Interface: http://192.168.2.61:47897
[root@node-master1ks2p ~]#
```

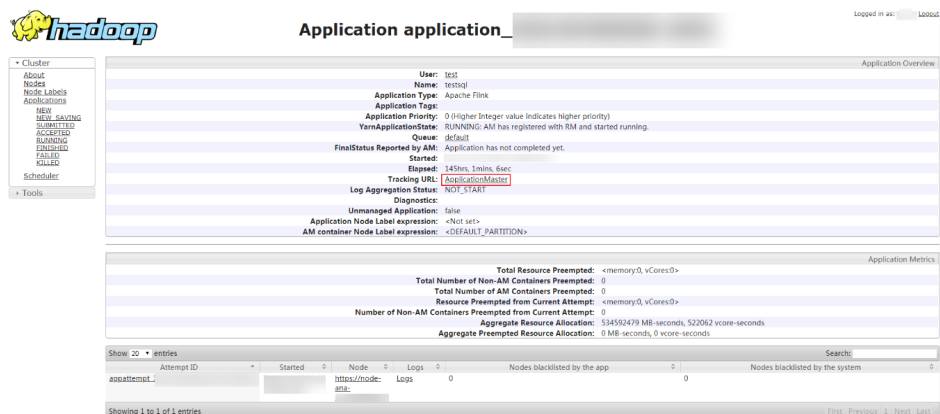
图 6-10 在 session 中提交作业成功

```
[root@node-master1ks2p ~]# flink run /opt/client/Flink/flink/examples/streaming/WordCount.jar
YARN properties set default parallelism to 3
2019-07-26 09:19:20,548 WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
2019-07-26 09:19:20,648 WARN | [main] | The short-circuit local reads feature cannot be used because libhadoop cannot be loaded. | org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory (DomainSocketFactory.java:118)
Starting execution of program
Executing WordCount example with default input data set.
Use --input to specify file input.
Printing result to stdout. Use --output to specify output path.
Program execution finished
Job with JobID Sb0cc1086563fd792a19163c2e7c3c3 has finished.
Job Runtime: 5908 ms
[root@node-master1ks2p ~]#
```

步骤8 使用运行用户进入Yarn服务的原生页面，具体操作参考[查看Flink作业信息](#)，找到对应作业的application，单击application名称，进入到作业详情页面。

- 若作业尚未结束，可单击“Tracking URL”链接进入到Flink的原生页面，查看作业的运行信息。
- 若作业已运行结束，对于在session中提交的作业，可以单击“Tracking URL”链接登录Flink原生页面查看作业信息。

图 6-11 application



----结束

6.4 创建 FlinkServer 作业前准备

6.4.1 访问 FlinkServer WebUI 界面

操作场景

MRS集群安装Flink组件后，用户可以通过Flink的WebUI，在图形化界面进行集群连接、数据连接、流表管理和作业管理等。

该任务指导用户在MRS集群中访问Flink WebUI。

对系统的影响

第一次访问Manager和Flink WebUI，需要在浏览器中添加站点信任以继续访问Flink WebUI。

操作步骤

步骤1 使用具有FlinkServer管理员权限的用户登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)，选择“集群 > 服务 > Flink”。

📖 说明

对于开启了Kerberos认证的MRS集群，访问Flink WebUI，需提前创建具有FlinkServer管理员权限或应用查看、应用编辑权限的角色，并为用户绑定该角色，角色创建可参考[创建FlinkServer 权限角色](#)。

步骤2 在“Flink WebUI”右侧，单击链接，访问Flink的WebUI。

图 6-12 访问 Flink 的 WebUI



Flink WebUI支持以下功能：

- 使用系统管理可以支持以下功能：
 - 使用集群连接管理可以创建、查看、编辑、测试和删除集群连接。
 - 使用数据连接管理可以创建、查看、编辑、测试和删除数据连接。数据连接类型包含HDFS、Kafka等。

- 使用应用管理可以创建、查看、删除应用。
- 使用流表管理可以新建、查看、编辑和删除流表。
- 使用作业管理可以新建、查看、启动、开发、编辑、停止和删除作业等。

----结束

6.4.2 创建 FlinkServer 应用

操作场景

通过应用来隔离不同的上层业务。

创建应用

步骤1 使用具有FlinkServer管理员权限的用户访问Flink WebUI，请参考[访问FlinkServer WebUI界面](#)。

步骤2 选择“系统管理 > 应用管理”，进入应用管理页面。

步骤3 单击“创建应用”，在弹出的页面中填写应用信息，单击“确定”，完成应用创建。

应用创建成功后，在Flink WebUI左上角即可切换待操作的应用，然后进行相关的作业开发。

----结束

6.4.3 创建 FlinkServer 集群连接

操作场景

通过集群连接配置访问不同的集群。

创建集群连接

步骤1 访问Flink WebUI，请参考[访问FlinkServer WebUI界面](#)。

步骤2 选择“系统管理 > 集群连接管理”，进入集群连接管理页面。

步骤3 单击“创建集群连接”，在弹出的页面中参考[表6-13](#)填写信息，单击“确定”，完成集群连接创建。创建完成后，可在对应集群连接的“操作”列对集群连接进行编辑、测试、删除等操作。

图 6-13 创建集群连接

创建集群连接

*集群连接名称	<input type="text" value="请输入集群连接名称"/>
描述	<input type="text" value="请输入描述信息"/>
*版本	<input type="text" value="MRS 3"/>
*是否安全版本	<input checked="" type="radio"/> 是 <input type="radio"/> 否
*访问用户名 	<input type="text" value="请输入访问用户名"/>
*客户端配置文件 	<input type="button" value="集群配置上传"/>
*用户凭据 	<input type="button" value="认证凭据上传"/>

表 6-13 创建集群连接信息

参数名称	参数描述
集群连接名称	集群连接的名称。
描述	集群连接名称描述信息。
版本	选择集群版本。
是否安全版本	<ul style="list-style-type: none">是，安全集群选择是。需要输入访问用户名和上传用户凭证；否，非安全集群选择否。
访问用户名	访问用户需要包含访问集群中服务所需要的最小权限。 “是否安全版本”选择“是”时存在此参数。

参数名称	参数描述
客户端配置文件	集群客户端配置文件，格式为tar。
用户凭据	FusionInsight Manager中用户的认证凭据，格式为tar。 “是否安全版本”选择“是”时存在此参数。 输入访问用户名后才可上传文件。

📖 说明

集群客户端配置文件获取方法：

1. 登录FusionInsight Manager，选择“集群 > 概览”。
2. 选择“更多 > 下载客户端 > 仅配置文件”，选择平台类型后单击“确定”。

用户凭据获取方法：

1. 登录FusionInsight Manager，单击“系统”。
2. 在对应用户的“操作”列，选择“更多 > 下载认证凭据”，选择集群后单击“确定”。

----结束

6.4.4 创建 FlinkServer 数据连接

操作场景

通过数据连接，访问不同的数据服务，当前FlinkServer支持HDFS、Kafka、Redis类型的数据连接。

创建数据连接

- 步骤1** 访问Flink WebUI，请参考[访问FlinkServer WebUI界面](#)。
- 步骤2** 选择“系统管理 > 数据连接管理”，进入数据连接管理页面。
- 步骤3** 单击“创建数据连接”，在弹出的页面中选择数据连接类型，参考[表6-14](#)填写信息，单击“确定”，完成数据连接创建。创建完成后，可在对应数据连接的“操作”列对数据连接进行编辑、测试、删除等操作。

表 6-14 创建数据连接信息

参数名称	参数描述	示例
数据连接类型	选择数据连接的类型，包含HDFS、Kafka、Redis。 选择Redis数据连接类型时，需提前准备“分布式缓存服务 Redis版”实例，并确保其“实例类型”为“Cluster集群”、“访问方式”为“免密访问”、同时“区域”和“虚拟私有云”需与Flink所在集群相同。	-
数据连接名称	数据连接的名称。	-

参数名称	参数描述	示例
集群连接	配置管理里的集群连接名称。 HDFS类型数据连接需配置该参数。	-
Kafka broker	Kafka Broker实例的连接信息，格式为“IP地址:端口”，多个实例之间通过逗号分隔。 Kafka类型数据连接需配置该参数。	192.168.0.1:21005,192.168.0.2:21005
Redis部署方式	Redis部署方式，当前仅支持“Cluster”。 Redis类型数据连接需配置该参数。	Cluster
Redis服务器列表	Redis实例的连接信息，格式为“IP地址:端口”，多个实例之间通过逗号分隔。 Redis类型数据连接需配置该参数。	192.168.0.1:6379,192.168.0.2:6379
认证类型	<ul style="list-style-type: none"> ● SIMPLE：表示对接的服务是非安全模式，无需认证。 ● KERBEROS：表示对接的服务是安全模式，安全模式的服务统一使用Kerberos认证协议进行安全认证。 	-

----结束

6.4.5 创建 FlinkServer 流表源

操作场景

通过数据表，定义源表、维表、输出表的基本属性和字段信息。

新建流表

步骤1 访问Flink WebUI，请参考[访问FlinkServer WebUI界面](#)。

步骤2 单击“流表管理”进入流表管理页面。

步骤3 单击“新建流表”，在新建流表页面参考[表6-15](#)填写信息，单击“确定”，完成流表创建。创建完成后，可在对应流表的“操作”列对流表进行编辑、删除等操作。

图 6-14 新建流表

*流/表名称 ?

描述 ?

*映射表类型 ? Kafka HDFS Redis

*类型 ? Source Sink

*数据连接

*Topic ?

*编码 ? JSON CSV

*前缀 ?

*流/表结构 ? Kafka流/表结构 ⊕

名称	类型	操作
<input type="text" value="请输入名称"/>	<input type="text" value="请选择类型"/>	<input type="text" value=""/>

*Proctime ?

Event Time ?

表 6-15 新建流表信息

参数名称	参数描述	备注
流/表名称	流/表的名称。	例如：flink_sink
描述	流/表的描述信息。	-
映射表类型	Flink SQL本身不带有数据存储功能，所有涉及表创建的操作，实际上均是对外部数据表、存储的引用映射。 类型包含Kafka、HDFS。	-
类型	包含数据源表Source，数据结果表Sink。不同映射表类型包含的表如下所示。 <ul style="list-style-type: none"> • Kafka：Source、Sink • HDFS：Source、Sink 	-
数据连接	选择数据连接。	-

参数名称	参数描述	备注
Topic	读取的Kafka的topic，支持从多个Kafka topic中读取，topic之间使用英文分隔符进行分隔。 “映射表类型”选择“Kafka”时存在此参数。	-
文件路径	要传输的HDFS目录或单个文件路径。 “映射表类型”选择“HDFS”时存在此参数。	例如： “/user/sqoop/” 或“/user/sqoop/ example.csv”
编码	选择不同“映射表类型”对应的编码如下： <ul style="list-style-type: none"> • Kafka：CSV、JSON • HDFS：CSV 	-
前缀	“映射表类型”选择“Kafka”，且“类型”选择“Source”，“编码”选择“JSON”时含义为：多层嵌套json的层级前缀，使用英文逗号(,)进行分隔。	例如：data,info表示取嵌套json中data, info下的内容，作为json格式数据输入
分隔符	选择不同“映射表类型”对应的含义为：用于指定CSV字段分隔符。当数据“编码”为“CSV”时存在此参数。	例如：“,”
行分隔符	文件中的换行符，包含“\r”、“\n”、“\r\n”。 “映射表类型”选择“HDFS”时存在此参数。	-
列分隔符	文件中的字段分隔符。 “映射表类型”选择“HDFS”时存在此参数。	例如：“,”
流/表结构	填写流/表结构，包含名称，类型。	-
Proctime	指系统时间，与数据本身的时间戳无关，即在Flink算子内计算完成的时间。 “类型”选择“Source”时存在此参数。	-
Event Time	指事件产生的时间，即数据产生时自带时间戳。 “类型”选择“Source”时存在此参数。	-

----结束

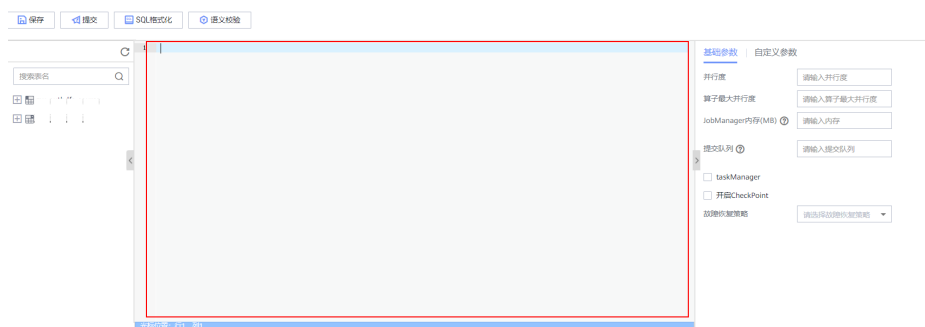
6.5 创建 FlinkServer 作业

操作场景

定义Flink的作业，包括Flink SQL和Flink Jar作业。

新建作业

- 步骤1** 访问Flink WebUI，请参考[访问FlinkServer WebUI界面](#)。
- 步骤2** 单击“作业管理”进入作业管理页面。
- 步骤3** 单击“新建作业”，在新建作业页面可选择新建Flink SQL作业或Flink Jar作业，然后填写作业信息，单击“确定”，创建作业成功并进入作业开发界面。
- 步骤4** （可选）如果需要立即进行作业开发，可以在作业开发界面进行作业配置。
- 新建Flink SQL作业
 - a. 在作业开发界面进行作业开发。



- b. 可以单击上方“语义校验”对输入内容校验，单击“SQL格式化”对SQL语句进行格式化。
- c. 作业SQL开发完成后，请参考[表6-16](#)设置基础参数，还可根据需要设置自定义参数，然后单击“保存”。

表 6-16 基础参数

参数名称	参数描述
并行度	并行数量。
算子最大并行度	算子最大的并行度。
JobManager内存 (MB)	JobManager的内存。输入值最小为512。
提交队列	作业提交队列。不填默认提交到default。
taskManager	taskManager运行参数。该参数需配置以下内容： <ul style="list-style-type: none"> ▪ slot数量：不填默认是1； ▪ 内存 (MB)：输入值最小为512。

参数名称	参数描述
开启CheckPoint	<p>是否开启CheckPoint。开启后，需配置以下内容：</p> <ul style="list-style-type: none"> ▪ 时间间隔（ms）：必填； ▪ 模式：必填； EXACTLY_ONCE：数据或事件仅会被算子处理一次； AT_LEAST_ONCE：数据或事件会被算子至少处理一次； ▪ 最小间隔（ms）：输入值最小为10； ▪ 超时时间：输入值最小为10； ▪ 最大并发量：正整数，且不能超过64个字符； ▪ 是否清理：是/否； ▪ 是否开启增量Checkpoint：是/否。
故障恢复策略	<p>作业的故障恢复策略，包含以下三种，详情请参考配置FlinkServer重启策略。</p> <ul style="list-style-type: none"> ▪ fixed-delay：需配置“重试次数”和“失败重试间隔（s）”； ▪ failure-rate：需配置“最大重试次数”、“时间间隔（min）”和“失败重试间隔（s）”； ▪ none：无。

- d. 单击左上角“提交”提交作业。
- 新建Flink Jar作业
 - a. 单击“选择”，上传本地Jar文件，并参考[表6-17](#)配置参数或添加自定义参数。

表 6-17 参数配置

参数名称	参数描述
本地jar文件	上传jar文件。直接上传本地文件，大小不能超过10M。

参数名称	参数描述
Main Class	Main-Class类型。 <ul style="list-style-type: none"> 默认：默认根据Jar包文件的Mainfest文件指定类名。 指定：手动指定类名。
类名	类名。 “Main Class”选择“指定”时存在该参数。
类参数	类参数，为Main-Class的参数（参数间用空格分隔）。
并行度	并行数量。
JobManager内存（MB）	JobManager的内存。输入值最小为512。
提交队列	作业提交队列。不填默认提交到default。
taskManager	taskManager运行参数。该参数需配置以下内容： <ul style="list-style-type: none"> slot数量：不填默认是1； 内存（MB）：输入值最小为512。

b. 单击“保存”保存配置，单击“提交”提交作业。

步骤5 返回作业管理页面，可以查看到已创建的作业名称、类型、状态、作业种类和描述等信息。

作业创建完成后，可在对应作业的“操作”列对作业进行启动、开发、停止、编辑、删除、查看作业详情和Checkpoint故障恢复等操作。

📖 说明

- 若要使用其他用户在节点上读取已提交的作业相关文件，需确保该用户与提交作业的用户具有相同的用户组和具有对应的FlinkServer应用管理权限角色，如参考[创建FlinkServer权限角色](#)勾选“应用查看”。
- 作业状态为“运行中”的作业可以查看作业详情。
- 作业状态为“运行失败”、“运行成功”和“停止”的作业可以进行Checkpoint故障恢复。

----结束

6.6 管理 FlinkServer 作业

6.6.1 配置 FlinkServer 重启策略

概述

Flink支持不同的重启策略，以在发生故障时控制作业是否重启以及如何重启。若不指定重启策略，集群会使用默认的重启策略。用户也可以在提交作业时指定一个重启策略，可参考[创建FlinkServer作业](#)在作业开发界面配置（MRS 3.1.0及以后版本）。

重启策略也可以通过Flink的配置文件“*客户端安装目录*/Flink/flink/conf/flink-conf.yaml”中的参数“restart-strategy”指定，为全局配置，还可以在应用代码中动态指定，会覆盖全局配置，重启策略包括失败率（failure-rate）和两种默认策略，默认策略为如下：

- 无重启（No restart）：若没有启用CheckPoint，默认使用该策略。
- 固定间隔（fixed-delay）：若启用了CheckPoint，但没有配置重启策略，默认使用该策略。

No restart 策略

发生故障时作业会直接失败，不会尝试重启。

参数配置为：

```
restart-strategy: none
```

fixed-delay 策略

发生故障时会尝试重启作业固定次数，如果超过了最大的尝试次数，作业最终会失败。并且在两次连续重启尝试之间，重启策略会等待固定的时间。

以配置若重启失败了3次则认为该Job失败，重试时间间隔为10s为例，参数配置为：

```
restart-strategy: fixed-delay  
restart-strategy.fixed-delay.attempts: 3  
restart-strategy.fixed-delay.delay: 10 s
```

failure-rate 策略

在作业失败后会直接重启，但超过设置的失败率后，作业会被认定为失败。在两个连续的重启尝试之间，重启策略会等待一个固定的时间。

以配置10分钟内若重启失败了3次则认为该作业失败，重试时间间隔为10s为例，参数配置为：

```
restart-strategy: failure-rate  
restart-strategy.failure-rate.max-failures-per-interval: 3  
restart-strategy.failure-rate.failure-rate-interval: 10 min  
restart-strategy.failure-rate.delay: 10 s
```

重启策略选择

- 如果用户在作业失败后，不希望重试，则推荐使用No restart策略。
- 如果用户在作业失败后，希望对作业进行重试，推荐使用failure-rate策略。因为fixed-delay策略可能会因为网络、内存等硬件故障导致用户作业失败次数达到最大重试次数，从而导致作业失败。

为了防止在failure-rate策略下的无限重启，推荐如下参数配置：

```
restart-strategy: failure-rate
restart-strategy.failure-rate.max-failures-per-interval: 3
restart-strategy.failure-rate.failure-rate-interval: 10 min
restart-strategy.failure-rate.delay: 10 s
```

6.6.2 配置 FlinkServer 作业中使用 UDF

本章节适用于MRS 3.1.2及之后的版本。

用户可以自定义一些函数，用于扩展SQL以满足个性化的需求，这类函数称为UDF。用户可以在Flink WebUI界面中上传并管理UDF jar包，然后在运行作业时调用相关UDF函数。

Flink支持以下3类自定义函数，如表6-18。

表 6-18 函数分类

分类	描述
UDF (User Defined Scalar Function)	自定义函数，支持一个或多个输入参数，返回一个结果值。详情请参考 UDF java代码及SQL样例 。
UDAF (User Defined Aggregation Function)	自定义聚合函数，将多条记录聚合成一个值。详情请参考 UDAF java代码及SQL样例 。
UDTF (User Defined Table-valued Function)	自定义表值函数，支持一个或多个输入参数，可返回多行多列。详情请参考 UDTF java代码及SQL样例 。

前提条件

准备UDF jar文件，大小不能超过200MB。

上传 UDF

步骤1 访问Flink WebUI，请参考[访问FlinkServer WebUI界面](#)。

步骤2 单击“UDF管理”进入UDF管理页面。

步骤3 单击“添加UDF”，在“本地Jar文件”参数后选择并上传本地已准备好的UDF jar文件。

步骤4 填写UDF名称以及描述信息后，单击“确定”。

📖 说明

- “UDF名称”最多可添加10项，“名称”可自定义，“类名”需与上传的UDF jar文件中UDF函数全限定类名一一对应。
- 上传UDF jar文件后，服务器默认保留5分钟，5分钟内单击确定则完成UDF创建，超时时单击确定则创建UDF失败并弹出错误提示：本地UDF文件路径有误。

步骤5 在UDF列表中，可查看当前应用内所有的UDF信息。可在对应UDF信息的“操作”列编辑或删除UDF信息（只能删除未被使用的UDF项）。

步骤6（可选）如果需要立即运行或开发作业，可在“作业管理”进行相关作业配置，可参考[创建FlinkServer作业](#)。

----结束

UDF java 代码及 SQL 样例

- UDF java使用样例

```
package com.xxx.udf;
import org.apache.flink.table.functions.ScalarFunction;
public class UdfClass_UDF extends ScalarFunction {
    public int eval(String s) {
        return s.length();
    }
}
```

- UDF SQL使用样例

```
CREATE TEMPORARY FUNCTION udf as 'com.xxx.udf.UdfClass_UDF';
CREATE TABLE udfSource (a VARCHAR) WITH ('connector' = 'datagen','rows-per-second'=1);
CREATE TABLE udfSink (a VARCHAR,b int) WITH ('connector' = 'print');
INSERT INTO
    udfSink
SELECT
    a,
    udf(a)
FROM
    udfSource;
```

UDAF java 代码及 SQL 样例

- UDAF java使用样例

```
package com.xxx.udf;
import org.apache.flink.table.functions.AggregateFunction;
public class UdfClass_UDAF {
    public static class AverageAccumulator {
        public int sum;
    }
    public static class Average extends AggregateFunction<Integer, AverageAccumulator> {
        public void accumulate(AverageAccumulator acc, Integer value) {
            acc.sum += value;
        }
        @Override
        public Integer getValue(AverageAccumulator acc) {
            return acc.sum;
        }
        @Override
        public AverageAccumulator createAccumulator() {
            return new AverageAccumulator();
        }
    }
}
```

- UDAF SQL使用样例

```
CREATE TEMPORARY FUNCTION udaf as 'com.xxx.udf.UdfClass_UDAF$Average';
CREATE TABLE udfSource (a int) WITH ('connector' = 'datagen','rows-per-second'=1,'fields.a.min'=1,'fields.a.max'=3);
CREATE TABLE udfSink (b int,c int) WITH ('connector' = 'print');
INSERT INTO
    udfSink
SELECT
    a,
    udaf(a)
FROM
    udfSource group by a;
```


UDTF java 代码及 SQL 样例

- UDTF java使用样例

```
package com.xxx.udf;
import org.apache.flink.api.java.tuple.Tuple2;
import org.apache.flink.table.functions.TableFunction;
public class UdfClass_UDTF extends TableFunction<Tuple2<String, Integer>> {
    public void eval(String str) {
        Tuple2<String, Integer> tuple2 = Tuple2.of(str, str.length());
        collect(tuple2);
    }
}
```

- UDTF SQL使用样例

```
CREATE TEMPORARY FUNCTION udtf as 'com.xxx.udf.UdfClass_UDTF';
CREATE TABLE udfSource (a VARCHAR) WITH ('connector' = 'datagen','rows-per-second'=1);
CREATE TABLE udfSink (b VARCHAR,c int) WITH ('connector' = 'print');
INSERT INTO
    udfSink
SELECT
    str,
    strLength
FROM
    udfSource,lateral table(udtf(udfSource.a)) as T(str,strLength);
```

6.7 Flink 运维管理

6.7.1 Flink 常用配置参数

配置说明

Flink所有的配置参数都可以在客户端侧进行配置，建议用户直接修改客户端的“flink-conf.yaml”配置文件进行配置，如果通过Manager界面修改Flink服务参数，配置完成之后需要重新下载安装客户端：

- 配置文件路径：*客户端安装路径*/Flink/flink/conf/flink-conf.yaml。
- 文件的配置格式为 *key: value*。

例：**taskmanager.heap.size: 1024mb**

注意配置项key:与value之间需有空格分隔。

配置详情

本章节介绍如下参数配置：

- JobManager & TaskManager：**

JobManager和TaskManager是Flink的主要组件，针对各种安全场景和性能场景，配置项包括通信端口，内存管理，连接重试等。

- Blob服务端：**

JobManager节点上的Blob服务端是用于接收用户在客户端上传的Jar包，或将Jar包发送给TaskManager，传输log文件等，配置项包括端口，SSL，重试次数，并发等。

- Distributed Coordination (via Akka)：**

Flink客户端与JobManager的通信，JobManager与TaskManager的通信和TaskManager与TaskManager的通信都基于Akka actor模型。相关参数可以根据

网络环境或调优策略进行配置，配置项包括消息发送和等待的超时设置，以及 Akka DeathWatch 检测机制参数等。

- **SSL:**
当需要配置安全 Flink 集群时，需要配置 SSL 相关配置项，配置项包括 SSL 开关，证书，密码，加密算法等。
- **Network communication (via Netty):**
Flink 运行 Job 时，Task 之间的数据传输和反压检测都依赖 Netty，某些环境下可能需要对 Netty 参数进行配置。对于高级调优，可调整部分 Netty 配置项，默认配置已可满足大规模集群并发高吞吐量的任务。
- **JobManager Web Frontend:**
JobManager 启动时，会在同一进程内启动 Web 服务器，访问 Web 服务器可以获得当前 Flink 集群的信息，包括 JobManager，TaskManager 及集群内运行的 Job。Web 服务器参数的配置项包括端口，临时目录，显示项目，错误重定向，安全相关等。
- **File Systems:**
Task 运行中会创建结果文件，支持对文件创建行为进行配置，配置项包括文件覆盖策略，目录创建等。
- **State Backend:**
Flink 提供了 HA 和作业的异常恢复，并且提供版本升级时作业的暂停恢复。对于作业状态的存储，Flink 依赖于 state backend，作业的重启依赖于重启策略，用户可以对这两部分进行配置。配置项包括 state backend 类型，存储路径，重启策略等。
- **Kerberos-based Security:**
Flink 安全模式下必须配置 Kerberos 相关配置项，配置项包括 kerberos 的 keytab、principal 等。
- **HA:**
Flink 的 HA 模式依赖于 ZooKeeper，所以必须配置 ZooKeeper 相关配置，配置项包括 ZooKeeper 地址，路径，安全认证等。
- **Environment:**
对于 JVM 配置有特定要求的场景，可以通过配置项传递 JVM 参数到客户端，JobManager，TaskManager 等。
- **Yarn:**
Flink 运行在 Yarn 集群上时，JobManager 运行在 Application Master 上。JobManager 的一些配置参数依赖于 Yarn，通过配置 YARN 相关的配置，使 Flink 更好的运行在 Yarn 上，配置项包括 yarn container 的内存，虚拟内核，端口等。
- **Pipeline:**
为适应某些场景对降低时延的需求，设计多个 Job 间采用 Netty 直接相连的方式传递数据，即分别使用 NettySink 用于 Server 端、NettySource 用于 Client 端进行数据传输。配置项包括 NettySink 的信息存放路径、NettySink 的端口监测范围、连接是否通过 SSL 加密以及 NettySink 监测所使用的网络所在域等。

JobManager & TaskManager

表 6-19 JobManager & TaskManager 参数说明

参数	描述	默认值	是否必选	备注
taskmanager.memory.size	TaskManager在JVM堆内存中保留空间的大小，此内存用于排序，哈希表和中间状态的缓存。如果未指定，则会使用JVM堆内存乘以比例taskmanager.memory.fraction。单位：MB。	0	否	仅MR S 3.x之前版本
taskmanager.registration.initial-backoff	两次连续注册的初始间隔时间。单位：ms/s/m/h/d。 时间数值和单位之间有半角字符空格。ms/s/m/h/d表示毫秒、秒、分钟、小时、天。	500 ms	否	
taskmanager.registration.refused-backoff	JobManager拒绝注册后到允许再次注册的间隔时间。	5 min	否	
taskmanager.rpc.port	TaskManager的IPC端口范围。	32326-32390	否	适用于所有版本
taskmanager.memory.segment-size	内存管理器和网络堆栈使用的内存缓冲区大小。单位：bytes。	32768	否	
taskmanager.data.port	TaskManager数据交换端口范围。	32391-32455	否	
taskmanager.data.ssl.enabled	TaskManager之间数据传输是否使用SSL加密，仅在全局开关security.ssl开启时有效。	false	否	
taskmanager.numberOfTaskSlots	TaskManager占用的slot数，一般配置成物理机的核数，yarn-session模式下只能使用-s参数传递，yarn-cluster模式下只能使用-ys参数传递。	1	否	
parallelism.default	默认并行度，用于未指定并行度的作业。	1	否	
taskmanager.memory.fraction	TaskManager在JVM堆内存中保留空间的比例，此内存用于排序，哈希表和中间状态的缓存。	0.7	否	
taskmanager.memory.off-heap	TaskManager是否使用堆外内存，此内存用于排序，哈希表和中间状态的缓存。建议对于大内存，开启此配置提高内存操作的效率。	false	是	

参数	描述	默认值	是否必选	备注
taskmanager.memory.preallocate	TaskManager是否在启动时分配保留内存空间。当开启堆外内存时，建议开启此配置项。	false	否	仅 MR S 3.x 及之后版本
task.cancellation.interval	两次连续任务取消操作的间隔时间。单位：ms。	30000	否	
client.rpc.port	Flink client端Akka system监测端口。	32651-32720	否	
jobmanager.heap.size	JobManager堆内存大小，yarn-session模式下只能使用-jm参数传递，yarn-cluster模式下只能使用-yjm参数传递，如果小于YARN配置文件中yarn.scheduler.minimum-allocation-mb大小，则使用YARN配置中的值。单位：B/KB/MB/GB/TB。	1024mb	否	
taskmanager.heap.size	TaskManager堆内存大小，yarn-session模式下只能使用-tm参数传递，yarn-cluster模式下只能使用-ytm参数传递，如果小于YARN配置文件中yarn.scheduler.minimum-allocation-mb大小，则使用YARN配置中的值。单位：B/KB/MB/GB/TB。	1024mb	否	
taskmanager.network.numberofbuffers	TaskManager网络传输缓冲栈数量，如果作业运行中出错提示系统中可用缓冲不足，可以增加这个配置项的值。	2048	否	
taskmanager.debug.memory.startLogThread	调试Flink内存和GC相关问题时可开启，TaskManager会定时采集内存和GC的统计信息，包括当前堆内，堆外，内存池的使用率和GC时间。	false	否	
taskmanager.debug.memory.logIntervalMs	TaskManager定时采集内存和GC的统计信息的采集间隔。	0	否	
taskmanager.maxRegistrationDuration	TaskManager向JobManager注册自己的最长时间，如果超过时间，TaskManager会关闭。	5 min	否	

参数	描述	默认值	是否必选	备注
taskmanager.initial-registration-pause	两次连续注册的初始间隔时间。该值需带一个时间单位（ms/s/min/h/d）（比如5秒）。 时间数值和单位之间有半角字符空格。ms/s/m/h/d表示毫秒、秒、分钟、小时、天。	500 ms	否	
taskmanager.max-registration-pause	TaskManager注册失败最大重试间隔。单位：ms/s/m/h/d。	30 s	否	
taskmanager.refused-registration-pause	TaskManager注册连接被JobManager拒绝后的重试间隔。单位：ms/s/m/h/d。	10 s	否	
classloader.resolve-order	从用户代码加载类时定义类解析策略，这意味着是首先检查用户代码jar（“child-first”）还是应用程序类路径（“parent-first”）。默认设置指示首先从用户代码jar加载类，这意味着用户代码jar可以包含和加载不同于Flink使用的（依赖）依赖项。	child-first	否	
slot.idle.timeout	Slot Pool中空闲Slot的超时时间（以ms为单位）。	50000	否	
slot.request.timeout	从Slot Pool请求Slot的超时（以ms为单位）。	300000	否	
task.cancellation.timeout	取消任务超时时间（以ms为单位），超时会触发TaskManager致命错误。设置为0，取消任务卡住则不会报错。	180000	否	
taskmanager.network.detailed-metrics	启用网络队列长度的详细指标监控。	false	否	
taskmanager.network.memory.buffer-per-channel	每个传出/传入通道（子分区/输入通道）使用的最大网络缓冲区数。在基于信用的流量控制模式下，这表示每个输入通道中有多少信用。它应配置至少2以获得良好的性能。1个缓冲区用于接收子分区中的飞行中数据，1个缓冲区用于并行序列化。	2	否	

参数	描述	默认值	是否必选	备注
taskmanager.network.memory.floating-buffers-per-gate	每个输出/输入门（结果分区/输入门）使用的额外网络缓冲区数。在基于信用的流量控制模式中，这表示在所有输入通道之间共享多少浮动信用。浮动缓冲区基于积压（子分区中的实时输出缓冲区）反馈来分布，并且可以帮助减轻由子分区之间的不平衡数据分布引起的背压。如果节点之间的往返时间较长和/或群集中的机器数量较多，则应增加此值。	8	否	
taskmanager.network.memory.fraction	用于网络缓冲区的JVM内存的占比。这决定了TaskManager可以同时拥有多少流数据交换通道以及通道缓冲的程度。如果作业被拒绝或者收到系统没有足够缓冲区的警告，请增加此值或 “taskmanager.network.memory.min”和 “taskmanager.network.memory.max”。另请注意， “taskmanager.network.memory.min”和 “taskmanager.network.memory.max”可能会覆盖此占比。	0.1	否	
taskmanager.network.memory.max	网络缓冲区的最大内存大小。该值需带一个大小单位（B/KB/MB/GB/TB）。	1 GB	否	
taskmanager.network.memory.min	网络缓冲区的最小内存大小。该值需带一个大小单位（B/KB/MB/GB/TB）。	64 MB	否	
taskmanager.network.request-backoff.initial	输入通道的分区请求的最小退避（以ms为单位）。	100	否	
taskmanager.network.request-backoff.max	输入通道的分区请求的最大退避（以ms为单位）。	10000	否	
taskmanager.registration.timeout	TaskManager注册的超时时间，在该时间内未成功注册，TaskManager将终止。该值需带一个时间单位（ms/s/min/h/d）。	5 min	否	

参数	描述	默认值	是否必选	备注
resourcemanager.taskmanager.timeout	释放空闲TaskManager的超时（以ms为单位）。	30000	否	

Blob 服务端

表 6-20 Blob 服务端参数说明

参数	描述	默认值	是否必选
blob.server.port	blob服务器端口。	32456-32520	否
blob.service.ssl.enabled	blob传输通道是否加密传输，仅在全局开关security.ssl开启时有。	true	是
blob.fetch.retries	TaskManager从JobManager下载blob文件的重试次数。	50	否
blob.fetch.num-concurrent	JobManager支持的下载blob的并发数。	50	否
blob.fetch.backlog	JobManager支持的blob下载队列大小，比如下载Jar包等。单位：个。	1000	否
library-cache-manager.cleanup.interval	当用户取消flink job后，jobmanager删除HDFS上存放用户jar包的时间，单位为s。 仅适用于MRS 3.x及之后版本。	3600	否

Distributed Coordination (via Akka)

表 6-21 Distributed Coordination 参数说明

参数	描述	默认值	是否必选	备注
akka.ask.timeout	akka所有异步请求和阻塞请求的超时时间。如果Flink发生超时失败，可以增大这个值。当机器处理速度慢或者网络阻塞时会发生超时。单位：ms/s/m/h/d。	10s	否	适用于所有版本
akka.lookup.timeout	查找JobManager actor对象的超时时间。单位：ms/s/m/h/d。	10s	否	

参数	描述	默认值	是否必选	备注
akka.framesize	JobManager和TaskManager间最大消息传输大小。当Flink出现消息大小超过限制的错误时，可以增大这个值。单位：b/B/KB/MB。	10485760b	否	
akka.watch.heartbeat.interval	Akka DeathWatch机制检测失联TaskManager的心跳间隔。如果TaskManager经常发生由于心跳消息丢失或延误而被错误标记为失联的情况，可以增大这个值。单位：ms/s/m/h/d。	10s	否	
akka.watch.heartbeat.pause	Akka DeathWatch可接受的心跳暂停时间，较小的数值表示不允许不规律的心跳。单位：ms/s/m/h/d。	60s	否	
akka.watch.threshold	DeathWath失败检测阈值，较小的数值容易把正常TaskManager标记为失败，较大的值增加了失败检测的时间。	12	否	
akka.tcp.timeout	发送连接TCP超时时间，如果经常发生满网络环境下连接TaskManager超时，可以增大这个值。单位：ms/s/m/h/d。	20s	否	
akka.throughput	Akka批量处理消息的数量，一次操作完后把处理线程归还线程池。较小的数值代表actor消息处理的公平调度，较大的值以牺牲调度公平的代价提高整体性能。	15	否	
akka.log.lifecycle.events	Akka远程时间日志开关，当需要调试时可打开此开关。	false	否	
akka.startup-timeout	远程组件启动失败前的超时时间。该值需带一个时间单位（ms/s/min/h/d）	与 akka.ask.timeout 的值一致	否	
akka.ssl.enabled	Akka通信SSL开关，仅在全局开关 security.ssl开启时有。	true	是	
akka.client-socket-worker-pool.pool-size-factor	计算线程池大小的因子，计算公式： $\text{ceil}(\text{可用处理器} \times \text{因子})$ ，计算结果限制在pool-size-min和pool-size-max之间。	1.0	否	仅适用于 MR S 3.x 及之后版本

参数	描述	默认值	是否必选	备注
akka.client-socket-worker-pool.pool-size-max	基于因子计算的线程数上限。	2	否	
akka.client-socket-worker-pool.pool-size-min	基于因子计算的线程数下限。	1	否	
akka.client.timeout	【说明】客户端超时时间。该值需带一个时间单位（ms/s/min/h/d）。	60s	否	
akka.server-socket-worker-pool.pool-size-factor	【说明】计算线程池大小的因子，计算公式： $\text{ceil}(\text{可用处理器} \times \text{因子})$ ，计算结果限制在pool-size-min和pool-size-max之间。	1.0	否	
akka.server-socket-worker-pool.pool-size-max	基于因子计算的线程数上限。	2	否	
akka.server-socket-worker-pool.pool-size-min	基于因子计算的线程数下限。	1	否	

SSL

表 6-22 SSL 参数说明

参数	描述	默认值	是否必选	备注
security.ssl.internal.enabled	内部通信SSL总开关，按照集群的安全模式自动配置。	<ul style="list-style-type: none"> 安全模式：true 普通模式：false 	是	仅MRS 3.x之前版本
security.ssl.internal.keystore	Java keystore文件。	-	是	
security.ssl.internal.keystore-password	keystore文件解密密码。	-	是	

参数	描述	默认值	是否必选	备注
security.ssl.internal.key-password	keystore文件中服务端key的解密密码。	-	是	
security.ssl.internal.truststore	truststore文件包含公共CA证书。	-	是	
security.ssl.internal.truststore-password	truststore文件解密密码。	-	是	
security.ssl.rest.enabled	外部通信SSL总开关，按照集群的安全模式自动配置。	<ul style="list-style-type: none"> 安全模式: true 普通模式: false 	是	
security.ssl.rest.keystore	Java keystore文件。	-	是	
security.ssl.rest.keystore-password	keystore文件解密密码。	-	是	
security.ssl.rest.key-password	keystore文件中服务端key的解密密码。	-	是	
security.ssl.rest.truststore	truststore文件包含公共CA证书。	-	是	
security.ssl.rest.truststore-password	truststore文件解密密码。	-	是	
security.ssl.protocol	SSL传输的协议版本。	TLSv1.2	是	适用于所有版本
security.ssl.algorithms	支持的SSL标准算法，具体可参考java官网： http://docs.oracle.com/javase/8/docs/technotes/guides/security/StandardNames.html#ciphersuites 。	TLS_DHE_RSA_WITH_AES_128_GCM_SHA256, TLS_ECDHE_RSA_WITH_AES_128_GCM_SHA256, TLS_DHE_RSA_WITH_AES_256_GCM_SHA384, TLS_ECDHE_RSA_WITH_AES_256_GCM_SHA384	是	

参数	描述	默认值	是否必选	备注
security.ssl.enabled	内部通信SSL总开关，按照集群的安装模式自动配置。	<ul style="list-style-type: none"> 安全模式: true 普通模式: false 	是	仅 MRS 3.x 及之后版本
security.ssl.keystore	Java keystore文件。	-	是	
security.ssl.keystore-password	keystore文件解密密码。	-	是	
security.ssl.key-password	keystore文件中服务端key的解密密码。	-	是	
security.ssl.truststore	truststore文件包含公共CA证书。	-	是	
security.ssl.truststore-password	truststore文件解密密码。	-	是	

Network communication (via Netty)

表 6-23 Network communication 参数说明

参数	描述	默认值	是否必选
taskmanager.network.netty.num-arenas	Netty内存块数。	1	否
taskmanager.network.netty.server.numThreads	Netty服务器线程的数量。	1	否
taskmanager.network.netty.client.numThreads	Netty客户端线程数。	1	否
taskmanager.network.netty.client.connectTimeoutSec	Netty客户端连接超时。单位: s。	120	否
taskmanager.network.netty.sendReceiveBufferSize	Netty发送和接收缓冲区大小。默认为系统缓冲区大小 (cat / proc / sys / net / ipv4 / tcp_ [rw] mem) , 在现代Linux中为4MB。单位: bytes。	4096	否

参数	描述	默认值	是否必选
taskmanager.net work.netty.transp ort	Netty传输类型，“nio”或 “epoll”。	nio	否

JobManager Web Frontend

表 6-24 JobManager Web Frontend 参数说明

参数	描述	默认值	是否必选	备注
jobmanager. web.port	web端口，支持范围： 32261-32325。	32261-32 325	否	仅 MR S 3.x 之前 版本
jobmanager. web.allow- access- address	web访问白名单，ip以逗号隔开。只 有在白名单中的ip才能访问web。	*	是	适用 于所 有版 本
flink.security .enable	用户安装Flink集群时，需要选择“安 全模式”或“普通模式”。 <ul style="list-style-type: none"> 当选择“安全模式”，自动配置 为“true”。 当选择“普通模式”，自动配置 为“false”。 对于已经安装好的Flink集群，用户可 以通过查看配置的值来区分当前安装 的是安全模式还是普通模式。	自动配置	否	仅 MR S 3.x 及之 后版 本
rest.bind- port	web端口，支持范围： 32261-32325。	32261-32 325	否	
jobmanager. web.history	显示“flink.security.enable”最近的 job数目。	5	否	
jobmanager. web.checkp oints.disable	禁用checkpoint统计。	false	否	
jobmanager. web.checkp oints.history	Checkpoint统计记录数。	10	否	

参数	描述	默认值	是否必选	备注
jobmanager.web.backpressure.cleanup-interval	未访问反压记录清理周期。单位：ms。	600000	否	
jobmanager.web.backpressure.refresh-interval	反压记录刷新周期。单位：ms。	60000	否	
jobmanager.web.backpressure.num-samples	计算反压使用的堆栈跟踪记录数。	100	否	
jobmanager.web.backpressure.delay-between-samples	计算反压的采样间隔。单位：ms	50	否	
jobmanager.web.ssl.enabled	web是否使用SSL加密传输，仅在全局开关security.ssl开启时有。	false	是	
jobmanager.web.accesslog.enable	web操作日志使能开关，日志会存放在webaccess.log中。	true	是	
jobmanager.web.x-frame-options	http安全头X-Frame-Options的值，可选范围为：SAMEORIGIN、DENY、ALLOW-FROM uri。	DENY	是	
jobmanager.web.cache-directive	web页面是否支持缓存。	no-store: 所有内容都不会被保存到缓存	是	
jobmanager.web.expires-time	web页面缓存过期时长。单位：ms。	0	是	
jobmanager.web.access-control-allow-origin	网页同源策略，防止跨域攻击。*表示允许任意网站跨域访问该服务端口，可配置为指定网址。	*(非安全集群)	是	
jobmanager.web.refresh-interval	web网页刷新时间。单位：ms。	3000	是	

参数	描述	默认值	是否必选	备注
jobmanager.web.logout-timer	配置无操作情况下自动登出时间间隔。单位：ms。	600000	是	
jobmanager.web.403-redirect-url	web403页面，访问若遇到403错误，则会重定向到配置的页面。	自动配置	是	
jobmanager.web.404-redirect-url	web404页面，访问若遇到404错误，则会重定向到配置的页面。	自动配置	是	
jobmanager.web.415-redirect-url	web415页面，访问若遇到415错误，则会重定向到配置的页面。	自动配置	是	
jobmanager.web.500-redirect-url	web500页面，访问若遇到500错误，则会重定向到配置的页面。	自动配置	是	
rest.await-leader-timeout	客户端等待Leader地址的时间（以ms为单位）。	30000	否	
rest.client.max-content-length	客户端处理的最大内容长度（以字节为单位）。	104857600	否	
rest.connection-timeout	客户端建立TCP连接的最长时间（以ms为单位）。	15000	否	
rest.idleness-timeout	连接保持空闲状态的最长时间（以ms为单位）。	300000	否	
rest.retry.delay	客户端在连续重试之间等待的时间（以ms为单位）。	3000	否	
rest.retry.max-attempts	如果可重试算子操作失败，客户端将尝试重试的次数。	20	否	
rest.server.max-content-length	服务端处理的最大内容长度（以字节为单位）。	104857600	否	
rest.server.numThreads	异步处理请求的最大线程数。	4	否	
web.timeout	web监控超时时间（以ms为单位）。	10000	否	

File Systems

表 6-25 File Systems 参数说明

参数	描述	默认值	是否必选
fs.overwrite-files	文件输出写操作是否默认覆盖已有文件。	false	否
fs.output.always-create-directory	<p>当文件写入程序的并行度大于1时，输出文件的路径下会创建一个目录，并将不同的结果文件（每个并行写程序任务）放入该目录。</p> <ul style="list-style-type: none"> • 设置为true，那么并行度为1的写入程序也将创建一个目录并将一个结果文件放入其中。 • 设置为false，则并行度为1的写入程序将直接在输出路径中创建文件，而不再创建目录。 	false	否

State Backend

表 6-26 State Backend 参数说明

参数	描述	默认值	是否必选
state.backend.fs.checkpointdir	当backend为filesystem时的路径，路径必须能够被JobManager访问到，本地路径只支持local模式，集群模式下请使用HDFS路径。	hdfs:///flink/checkpoints	否
state.savepoints.dir	Flink用于恢复和更新作业的保存点存储目录。当触发保存点的时候，保存点元数据信息将会保存到该目录中。	hdfs:///flink/savepoint	安全模式下必配
restart-strategy	<p>默认重启策略，用于未指定重启策略的作业：</p> <ul style="list-style-type: none"> • fixed-delay • failure-rate • none 	none	否

参数	描述	默认值	是否必选
restart-strategy.fixed-delay.attempts	fixed-delay策略重试次数。	<ul style="list-style-type: none"> 作业中开启了checkpoint, 默认值为Integer.MAX_VALUE。 作业中未开启checkpoint, 默认值为3。 	否
restart-strategy.fixed-delay.delay	fixed-delay策略重试间隔时间。单位：ms/s/m/h/d。	<ul style="list-style-type: none"> 作业中开启了checkpoint, 默认值是10s。 作业中未开启checkpoint, 默认值和配置项akka.ask.timeout的值一致。 	否
restart-strategy.failure-rate.max-failures-per-interval	故障率策略下作业失败前给定时间段内的最大重启次数。	1	否
restart-strategy.failure-rate.failure-rate-interval	failure-rate策略重试时间。单位：ms/s/m/h/d。	60 s	否

参数	描述	默认值	是否必选
restart-strategy.failure-rate.delay	failure-rate策略重试间隔时间。单位：ms/s/m/h/d。	默认值和akka.ask.timeout配置值一样。可参考 Distributed Coordination (via Akka) 。	否

Kerberos-based Security

表 6-27 Kerberos-based Security 参数说明

参数	描述	默认值	是否必选
security.kerberos.login.keytab	该参数为客户端参数，keytab路径。	根据实际业务配置	是
security.kerberos.login.principal	该参数为客户端参数，如果keytab和principal都设置，默认会使用keytab认证。	根据实际业务配置	否
security.kerberos.login.contexts	该参数为服务器端参数，flink生成jass文件的contexts。	Client、KafkaClient	是

HA

表 6-28 HA 参数说明

参数	描述	默认值	是否必选
high-availability	HA模式，是启用HA还是非HA模式。当前支持两种模式： <ul style="list-style-type: none"> • none，只运行单个jobManager，jobManager的状态不进行Checkpoint。 • ZooKeeper。 <ul style="list-style-type: none"> - 非YARN模式下，支持多个jobManager，通过选举产生leader。 - YARN模式下只存在一个jobManager。 	zookeeper	否

参数	描述	默认值	是否必选
high-availability.zookeeper.quorum	ZooKeeper quorum地址。	自动配置	否
high-availability.zookeeper.path.root	Flink在ZooKeeper上创建的根目录，存放HA模式必须的元数据。	/flink	否
high-availability.storageDir	存放state backend中JobManager元数据，ZooKeeper只保存实际数据的指针。	hdfs:///flink/recovery	否
high-availability.zookeeper.client.session-timeout	ZooKeeper客户端会话超时时间。单位：ms。	60000	否
high-availability.zookeeper.client.connection-timeout	ZooKeeper客户端连接超时时间。单位：ms。	15000	否
high-availability.zookeeper.client.retry-wait	ZooKeeper客户端重试等待时间。单位：ms。	5000	否
high-availability.zookeeper.client.max-retry-attempts	ZooKeeper客户端最大重试次数。	3	否
high-availability.job.delay	当jobManager恢复后重启job的延迟时间。 仅适用于MRS 3.x及之后版本。	默认值和akka.ask.timeout配置值保持一致	否
high-availability.zookeeper.client.acl	设置ZooKeeper节点的ACL (open creator)，按照集群的安全模式自动配置。设置ACL选项请参考： https://zookeeper.apache.org/doc/r3.5.1-alpha/zookeeperProgrammers.html#sc_BuiltinACLschemes 。	<ul style="list-style-type: none"> 安全模式：creator 非安全模式：open 	是
zookeeper.sasl.disable	基于SASL认证的使能开关，按照集群的安全模式自动配置：。	<ul style="list-style-type: none"> 安全模式：false 非安全模式：true 	是

参数	描述	默认值	是否必选
zookeeper.sasl.service-name	<ul style="list-style-type: none"> 如果ZooKeeper服务端配置了不同于“ZooKeeper”的服务名，可以设置此配置项。 如果客户端和服务端的服务名不一致，认证会失败。 	zookeeper	是

Environment

表 6-29 Environment 参数说明

参数	描述	默认值	是否必选
env.java.opts	JVM参数，会传递到启动脚本，JobManager，TaskManager，Yarn客户端。比如传递远程调试的参数等。	-Xloggc:<LOG_DIR>/gc.log -XX:+PrintGCDetails -XX:-OmitStackTraceInFastThrow -XX:+PrintGCTimeStamps -XX:+PrintGCDateStamps -XX:+UseGCLogFileRotation -XX:NumberOfGCLogFiles=20 -XX:GCLogFileSize=20M -Djdk.tls.ephemeralDHKeySize=2048 -Djava.library.path=\$ {HADOOP_COMMON_HOME}/lib/native -Djava.net.preferIPv4Stack=true -Djava.net.preferIPv6Addresses=false -Dbeetle.application.home.path=/opt/xxx/Bigdata/common/runtime/security/config	否

Yarn

表 6-30 Yarn 参数说明

参数	描述	默认值	是否必选
yarn.maximum-failed-containers	当TaskManager所属容器出错后，重新申请container次数。默认值为Flink集群启动时TaskManager的数量。	5	否
yarn.application-attempts	Application master重启次数，次数是算在一个validity interval的最大次数，validity interval在flink中设置为akka的timeout。重启后AM的地址和端口会变化，client需要手动连接。	2	否

参数	描述	默认值	是否必选
yarn.heartbeat-delay	Application Master和YARN Resource Manager心跳的时间间隔。单位：seconds	5	否
yarn.containers.vcores	每个Yarn容器的虚拟核数。	TaskManager的slot数	否
yarn.application-master.port	Application Master端口号设置，支持端口范围。	32586-32650	否

Pipeline

适用于MRS 3.x及之后版本。

表 6-31 Pipeline 参数说明

参数	描述	默认值	是否必选
nettyconnector.registerserver.topic.storage	设置NettySink的IP、端口及并发度信息在第三方注册服务器上的路径。建议用户使用ZooKeeper进行存储。	/flink/nettyconnector	否，当使用pipeline特性为必选
nettyconnector.sinkserver.port.range	设置NettySink的端口范围。	28444-28843	否，当使用pipeline特性为必选
nettyconnector.ssl.enabled	设置NettySink与NettySource之间通信是否配置SSL加密。其中加密密钥以及加密协议等请参见SSL。	false	否，当使用pipeline特性为必选
nettyconnector.message.delimiter	用来配置nettysink发送给nettysource消息的分隔符，长度为2-4个字节，不可包含“\n”，“ ”，“#”。	默认使用“\$ _”	否，当使用pipeline特性为必选

6.7.2 Flink 日志介绍

日志描述

日志存储路径：

- Flink作业运行日志：“\${BIGDATA_DATA_HOME}/hadoop/data\${i}/nm/containerlogs/application_\${appid}/container_\${scontid}”。

说明

运行中的任务日志存储在以上路径中，运行结束后会基于Yarn的配置确定是否汇聚到HDFS目录中。

- FlinkResource运行日志：“/var/log/Bigdata/flink/flinkResource”。

日志归档规则：

1. FlinkResource运行日志：

- 服务日志默认20MB滚动存储一次，最多保留20个文件，不压缩。

说明

针对MRS 3.x之前版本，Executor日志默认30MB滚动存储一次，最多保留20个文件，不压缩。

- 日志大小和压缩文件保留个数可以在Manager界面中配置或者修改客户端“客户端安装目录/Flink/flink/conf/”中的log4j-cli.properties、log4j.properties、log4j-session.properties中对应的配置项。

表 6-32 FlinkResource 日志列表

日志类型	日志文件名	描述
FlinkResource运行日志	checkService.log	健康检查日志。
	kinit.log	初始化日志。
	postinstall.log	服务安装日志。
	prestart.log	prestart脚本日志。
	start.log	启动日志。

2. FlinkServer服务日志、审计日志。

- FlinkServer服务日志、审计日志默认100MB滚动存储一次，服务日志最多保留30天，审计日志最多保留90天。
- 日志大小和压缩文件保留个数可以在Manager界面中配置或者修改客户端“客户端安装目录/Flink/flink/conf/”中的log4j-cli.properties、log4j.properties、log4j-session.properties中对应的配置项。

表 6-33 FlinkServer 日志列表

日志类型	日志文件名	描述
FlinkServer运行日志	checkService.log	健康检查日志。
	checkFlinkServer.log	FlinkServer健康检查日志。
	localhost_access_log.yyyy-mm-dd.txt	FlinkServer访问URL日志。

日志类型	日志文件名	描述
	start_thrift_server.out	thrift server启动日志。
	thrift_server_thriftServer_XXX.log.last	
	cleanup.log	安装卸载实例时的清理日志。
	flink-omm-client-IP.log	作业启动日志。
	flinkserver_yyyymmdd-x.log.gz	业务归档日志。
	flinkserver.log	业务日志。
	flinkserver---pidxxx-gc.log.x.current	GC日志。
	kinit.log	初始化日志。
	postinstall.log	服务安装日志。
	prestart.log	prestart脚本日志。
	start.log	启动日志。
	stop.log	停止日志。
	catalina.yyyy-mm-dd.log	tomcat运行日志。
	catalina.out	
	host-manager.yyyy-mm-dd.log	
localhost.yyyy-mm-dd.log		
manager.yyyy-mm-dd.log		
FlinkServer审计日志	flinkserver_audit_yyyymmdd-x.log.gz	审计归档日志。
	flinkserver_audit.log	审计日志。

日志级别

Flink中提供了如表6-34所示的日志级别。日志级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG。程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 6-34 日志级别

级别	描述
ERROR	ERROR表示当前时间处理存在错误信息。

级别	描述
WARN	WARN表示当前事件处理存在异常信息。
INFO	INFO表示记录系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

步骤1 请参考[修改集群服务配置参数](#)，进入Flink的“全部配置”页面。

步骤2 左边菜单栏中选择所需修改的角色所对应的日志菜单。

步骤3 选择所需修改的日志级别。

步骤4 保存配置，在弹出窗口中单击“确定”使配置生效。

----结束

📖 说明

- 配置完成后不需要重启服务，重新下载客户端使配置生效。
- 也可以直接修改客户端“客户端安装目录/Flink/flink/conf/”中log4j-cli.properties、log4j.properties、log4j-session.properties文件中对应的日志级别配置项。
- 通过客户端提交作业时会在客户端log文件夹中生成相应日志文件，由于系统默认umask值是0022，所以日志默认权限为644；如果需要修改文件权限，需要修改umask值；例如修改omm用户umask值：
 - 在“/home/omm/.baskrc”文件末尾添加“umask 0026”；
 - 执行命令`source /home/omm/.baskrc`使文件权限生效。

日志格式

表 6-35 日志格式

日志类型	格式	示例
运行日志	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置>	2019-06-27 21:30:31,778 INFO [flink-akka.actor.default-dispatcher-3] TaskManager container_e10_1498290698388_0004_02_0000 07 has started. org.apache.flink.yarn.YarnFlinkResourceManager (FlinkResourceManager.java:368)

6.8 Flink 性能调优

6.8.1 优化 Flink 内存 GC 参数

操作场景

Flink是依赖内存计算，计算过程中内存不够对Flink的执行效率影响很大。可以通过监控GC（Garbage Collection），评估内存使用及剩余情况来判断内存是否变成性能瓶颈，并根据情况优化。

监控节点进程的YARN的Container GC日志，如果频繁出现Full GC，需要优化GC。

📖 说明

GC的配置：在客户端的“conf/flink-conf.yaml”配置文件中，在“env.java.opts”配置项中添加参数：“-Xloggc:<LOG_DIR>/gc.log -XX:+PrintGCDetails -XX:-OmitStackTraceInFastThrow -XX:+PrintGCTimeStamps -XX:+PrintGCDateStamps -XX:+UseGCLogFileRotation -XX:NumberOfGCLogFiles=20 -XX:GCLogFileSize=20M”。此处默认已经添加GC日志。

操作步骤

- 优化GC。
调整老年代和新生代的比值。在客户端的“conf/flink-conf.yaml”配置文件中，在“env.java.opts”配置项中添加参数：“-XX:NewRatio”。如“-XX:NewRatio=2”，则表示老年代与新生代的比值为2:1，新生代占整个堆空间的1/3，老年代占2/3。
- 开发Flink应用程序时，优化DataStream的数据分区或分组操作。
 - 当分区导致数据倾斜时，需要考虑优化分区。
 - 避免非并行度操作，有些对DataStream的操作会导致无法并行，例如WindowAll。
 - keyBy尽量不要使用String。

6.8.2 配置 Flink 任务并行度

操作场景

并行度控制任务的数量，影响操作后数据被切分成的块数。调整并行度让任务的数量和每个任务处理的数据与机器的处理能力达到更优。

查看CPU使用情况和内存占用情况，当任务和数据不是平均分布在各节点，而是集中在个别节点时，可以增大并行度使任务和数据更均匀的分布在各个节点。增加任务的并行度，充分利用集群机器的计算能力。

操作步骤

任务的并行度可以通过以下四种层次（按优先级从高到低排列）指定，用户可以根据实际的内存、CPU、数据以及应用程序逻辑的情况调整并行度参数。

- 算子层次
一个算子、数据源和sink的并行度可以通过调用setParallelism()方法来指定，例如
- ```
final StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();

DataStream<String> text = [...]
DataStream<Tuple2<String, Integer>> wordCounts = text
 .flatMap(new LineSplitter())
```



```
.keyBy(0)
.timeWindow(Time.seconds(5))
.sum(1).setParallelism(5);

wordCounts.print();

env.execute("Word Count Example");
```

- 执行环境层次

Flink程序运行在执行环境中。执行环境为所有执行的算子、数据源、data sink定义了一个默认的并行度。

执行环境的默认并行度可以通过调用setParallelism()方法指定。例如：

```
final StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();
env.setParallelism(3);
DataStream<String> text = [...];
DataStream<Tuple2<String, Integer>> wordCounts = [...];
wordCounts.print();
env.execute("Word Count Example");
```

- 客户端层次

并行度可以在客户端将job提交到Flink时设定。对于CLI客户端，可以通过“-p”参数指定并行度。例如：

```
./bin/flink run -p 10 ../examples/*WordCount-java*.jar
```

- 系统层次

在系统级可以通过修改Flink客户端conf目录下的“flink-conf.yaml”文件中的“parallelism.default”配置选项来指定所有执行环境的默认并行度。

## 6.8.3 配置 Flink 任务进程参数

### 操作场景

Flink on YARN模式下，有JobManager和TaskManager两种进程。在任务调度和运行的过程中，JobManager和TaskManager承担了很大的责任。

因而JobManager和TaskManager的参数配置对Flink应用的执行有着很大的影响意义。用户可通过如下操作对Flink集群性能做优化。

### 操作步骤

#### 步骤1 配置JobManager内存。

JobManager负责任务的调度，以及TaskManager、RM之间的消息通信。当任务数变多，任务平行度增大时，JobManager内存都需要相应增大。

您可以根据实际任务数量的多少，为JobManager设置一个合适的内存。

- 在使用yarn-session命令时，添加“-jm MEM”参数设置内存。
- 在使用yarn-cluster命令时，添加“-yjm MEM”参数设置内存。

#### 步骤2 配置TaskManager个数。

每个TaskManager每个核同时能跑一个task，所以增加了TaskManager的个数相当于增大了任务的并发度。在资源充足的情况下，可以相应增加TaskManager的个数，以提高运行效率。

#### 步骤3 配置TaskManager Slot数。

每个TaskManager多个核同时能跑多个task，相当于增大了任务的并发度。但是由于所有核共用TaskManager的内存，所以要在内存和核数之间做好平衡。

- 在使用yarn-session命令时，添加“-s NUM”参数设置SLOT数。
- 在使用yarn-cluster命令时，添加“-ys NUM”参数设置SLOT数。

#### 步骤4 配置TaskManager内存。

TaskManager的内存主要用于任务执行、通信等。当一个任务很大的时候，可能需要较多资源，因而内存也可以做相应的增加。

- 将在使用yarn-session命令时，添加“-tm MEM”参数设置内存。
- 将在使用yarn-cluster命令时，添加“-ytm MEM”参数设置内存。

----结束

## 6.8.4 优化 Flink Netty 网络通信参数

### 操作场景

Flink通信主要依赖netty网络，所以在Flink应用执行过程中，netty的设置尤为重要，网络通信的好坏决定着数据交换的速度以及任务执行的效率。

### 操作步骤

以下配置均可在客户端的“conf/flink-conf.yaml”配置文件中进行修改适配，默认已经是相对较优解，请谨慎修改，防止性能下降。

- “taskmanager.network.netty.num-arenas”：默认是“taskmanager.numberOfWorkers”，表示netty的域的数量。
- “taskmanager.network.netty.server.numThreads”和“taskmanager.network.netty.client.numThreads”：默认是“taskmanager.numberOfWorkers”，表示netty的客户端和服务端的线程数目设置。
- “taskmanager.network.netty.client.connectTimeoutSec”：默认是120s，表示taskmanager的客户端连接超时的时间。
- “taskmanager.network.netty.sendReceiveBufferSize”：默认是系统缓冲区大小（cat /proc/sys/net/ipv4/tcp\_[rw]mem），一般为4MB，表示netty的发送和接收的缓冲区大小。
- “taskmanager.network.netty.transport”：默认为“nio”方式，表示netty的传输方式，有“nio”和“epoll”两种方式。

## 6.9 Flink 客户端常见命令说明

本章节适用于MRS 3.x及之后版本。

在使用Flink的Shell脚本前，首先需要执行以下操作，详细使用场景可参考[Flink客户端使用实践](#)运行wordcount作业：

**步骤1** 安装Flink客户端，例如安装目录为“/opt/client”。

**步骤2** 初始化环境变量。

**source /opt/client/bigdata\_env**

**步骤3** 如果当前集群已启用Kerberos认证，需先配置客户端认证，可参考**步骤5**。如果当前集群未启用Kerberos认证，则无需执行该步骤。

**步骤4** 参考**表6-36**运行相关命令。

**表 6-36** Flink Shell 命令参考

| 命令              | 参数说明                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                | 描述                              |
|-----------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------|
| yarn-session.sh | -at,--applicationType <arg>: 为Yarn application自定义类型。<br>-D <property=value>: 动态参数配置。<br>-d,--detached: 关闭交互模式，启动一个分离的Flink YARN session。<br>-h,--help: 显示Yarn session CLI的帮助。<br>-id,--applicationId <arg>: 绑定到一个已经运行的Yarn session。<br>-j,--jar <arg>: 设置用户jar包路径。<br>-jm,--jobManagerMemory <arg>: 为JobManager设置内存。<br>-m,--jobmanager <arg>: 要连接的JobManager的地址，使用该参数可以连接特定的JobManager。<br>-nl,--nodeLabel <arg>: 指定YARN application的nodeLabel。<br>-nm,--name <arg>: 为Yarn application自定义名称。<br>-q,--query: 查询可用的Yarn 资源。<br>-qu,--queue <arg>: 指定YARN 队列。<br>-s,--slots <arg>: 设置每个Taskmanager的SLOT个数。<br>-t,--ship <arg>: 指定待发送文件的目录。<br>-tm,--taskManagerMemory <arg>: 为TaskManager设置内存。<br>-yd,--yarndetached: 以分离模式启动。<br>-z,--zookeeperNamespace <args>: 指定zookeeper的namespace。<br>-h: 获取帮助。 | 启动一个常驻的Flink集群，接受来自Flink客户端的任务。 |

| 命令        | 参数说明                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         | 描述                                                                                                                                                               |
|-----------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| flink run | <p>-c,--class &lt;classname&gt;: 指定一个类作为程序运行的入口点。</p> <p>-C,--classpath &lt;url&gt;: 指定classpath。</p> <p>-d,--detached: 以分离方式运行job。</p> <p>-n,--allowNonRestoredState: 从快照点恢复时允许跳过不能恢复的状态。比如删除了程序中某个操作符，那么在恢复快照点时需要增加该参数。</p> <p>-m,--jobmanager &lt;host:port&gt;: 指定JobManager。</p> <p>-p,--parallelism &lt;parallelism&gt;: 指定job并行度，会覆盖配置文件中配置的并行度参数。</p> <p>-q,--sysoutLogging: 禁止flink日志输出至控制台。</p> <p>-s,--fromSavepoint &lt;savepointPath&gt;: 指定用于恢复job的savepoint路径。</p> <p>-z,--zookeeperNamespace &lt;zookeeperNamespace&gt;: 指定zookeeper的namespace。</p> <p>-yat,--yarnapplicationType &lt;arg&gt;: 为Yarn application自定义类型。</p> <p>-yD &lt;arg&gt;: 动态参数配置。</p> <p>-yd,--yarn detached: 以分离模式启动。</p> <p>-yh,--yarnhelp: 获取yarn帮助。</p> <p>-yid,--yarnapplicationId &lt;arg&gt;: 绑定到yarn session运行job。</p> <p>-yj,--yarnjar &lt;arg&gt;: 设置Flink jar文件路径。</p> <p>-yjm,--yarnjobManagerMemory &lt;arg&gt;: 为JobManager设置内存（MB）。</p> <p>-ynm,--yarnname &lt;arg&gt;: 为Yarn application自定义名称。</p> <p>-yq,--yarnquery: 查询可用的YARN资源（内存、CPU）。</p> <p>-yqu,--yarnqueue &lt;arg&gt;: 指定YARN队列。</p> <p>-ys,--yarnslots: 设置每个TaskManager的SLOT个数。</p> <p>-yt,--yarnship &lt;arg&gt;: 指定待发送文件的路径。</p> <p>-ytm,--yarn taskManagerMemory &lt;arg&gt;: 为TaskManager设置内存（MB）。</p> <p>-yz,--yarnzookeeperNamespace &lt;arg&gt;: 指定zookeeper的namespace，需与yarn-session.sh -z 保持一致。</p> | <p>Flink提交作业。</p> <ol style="list-style-type: none"> <li>1. "-y*"参数是指 yarn-cluster模式下使用。</li> <li>2. 非"-y*"参数用户在用该命令提交任务前需要先用 yarn-session启动 Flink集群。</li> </ol> |

| 命令           | 参数说明                                                                                                                                                                                                                                                                                                                   | 描述                                                                            |
|--------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------|
|              | -h: 获取帮助。                                                                                                                                                                                                                                                                                                              |                                                                               |
| flink info   | -c,--class <classname>: 指定一个类作为程序运行的入口点。<br>-p,--parallelism <parallelism>: 指定程序运行的并行度。<br>-h: 获取帮助。                                                                                                                                                                                                                   | 显示所运行程序的执行计划（JSON）                                                            |
| flink list   | -a,--all: 显示所有的Job。<br>-m,--jobmanager <host:port>: 指定JobManager。<br>-r,--running: 仅显示running状态的Job。<br>-s,--scheduled: 仅显示scheduled状态的Job。<br>-z,--zookeeperNamespace <zookeeperNamespace>: 指定zookeeper的namespace。<br>-yid,--yarnapplicationId <arg>: 绑定YARN session。<br>-h: 获取帮助。                                    | 查询集群中运行的程序。                                                                   |
| flink stop   | -d,--drain: 在触发savepoint和停止作业之前，发送MAX_WATERMARK。<br>-p,--savepointPath <savepointPath>: savepoint的储存路径，默认目录state.savepoints.dir。<br>-m,--jobmanager <host:port>: 指定JobManager。<br>-z,--zookeeperNamespace <zookeeperNamespace>: 指定zookeeper的namespace。<br>-yid,--yarnapplicationId <arg>: 绑定YARN session。<br>-h: 获取帮助。 | 强制停止一个运行中的Job（仅支持streaming jobs、业务代码 source 端需要 implements StoppableFunction） |
| flink cancel | -m,--jobmanager <host:port>: 指定JobManager。<br>-s,--withSavepoint <targetDirectory>: 取消Job时触发savepoint，默认目录state.savepoints.dir<br>-z,--zookeeperNamespace <zookeeperNamespace>: 指定zookeeper的namespace。<br>-yid,--yarnapplicationId <arg>: 绑定YARN session。<br>-h: 获取帮助。                                                 | 取消一个运行中Job                                                                    |

| 命令                                | 参数说明                                                                                                                                                                                                                        | 描述                                                                                                                                             |
|-----------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------|
| flink savepoint                   | -d,--dispose <arg>: 指定savepoint的保存目录。<br>-m,--jobmanager <host:port>: 指定JobManager。<br>-z,--zookeeperNamespace <zookeeperNamespace>: 指定zookeeper的namespace。<br>-yid,--yarnapplicationId <arg>: 绑定YARN session。<br>-h: 获取帮助。 | 触发一个savepoint                                                                                                                                  |
| <b>source</b> 客户端安装目录/bigdata_env | 无                                                                                                                                                                                                                           | 导入客户端环境变量。<br>使用限制：如果用户使用自定义脚本（例如A.sh）并在脚本中调用该命令，则脚本A.sh不能传入参数。如果确实需要给A.sh传入参数，则需采用二次调用方式。<br>例如A.sh中调用B.sh，在B.sh中调用该命令。A.sh可以传入参数，B.sh不能传入参数。 |
| start-scala-shell.sh              | local   remote <host> <port>   yarn: 运行模式                                                                                                                                                                                   | scala shell启动脚本                                                                                                                                |
| sh generate_keystore.sh           | -                                                                                                                                                                                                                           | 用户调用“generate_keystore.sh”脚本工具生成“Security Cookie”、“flink.keystore”和“flink.truststore”。需要输入自定义密码（不能包含#）。                                        |

----结束

## 6.10 Flink 常见问题

### 数据倾斜

当数据发生倾斜（某一部分数据量特别大），虽然没有GC（Gabbage Collection，垃圾回收），但是task执行时间严重不一致。

- 需要重新设计key，以更小粒度的key使得task大小合理化。
- 修改并行度。
- 调用rebalance操作，使数据分区均匀。

## 缓冲区超时设置

- 由于task在执行过程中存在数据通过网络进行交换，数据在不同服务器之间传递的缓冲区超时时间可以通过setBufferTimeout进行设置。
- 当设置“setBufferTimeout(-1)”，会等待缓冲区满之后才会刷新，使其达到最大吞吐量；当设置“setBufferTimeout(0)”时，可以最小化延迟，数据一旦接收到就会刷新；当设置“setBufferTimeout”大于0时，缓冲区会在该时间之后超时，然后进行缓冲区的刷新。

示例可以参考如下：

```
env.setBufferTimeout(timeoutMillis);

env.generateSequence(1,10).map(new MyMapper()).setBufferTimeout(timeoutMillis);
```

## 资源冗余量

Flink任务运行时，建议整个集群的Yarn资源留有一定的余量。比如当前Yarn总体的资源有100Vcore，200GB，则建议Yarn的任务使用90vcore，180GB，保留10%的资源用于当部分节点故障时，任务可以自动重试恢复。

## 6.11 签发 Flink 证书样例

将该样例代码生成generate\_keystore.sh脚本，放置在Flink客户端的bin目录下。

```
#!/bin/bash

KEYTOOL=${JAVA_HOME}/bin/keytool
KEYSTOREPATH="$FLINK_HOME/conf/"
CA_ALIAS="ca"
CA_KEYSTORE_NAME="ca.keystore"
CA_DNAME="CN=Flink_CA"
CA_KEYALG="RSA"
CLIENT_CONF_YAML="$FLINK_HOME/conf/flink-conf.yaml"
KEYTABPRINCEPAL=""

function getConf()
{
 if [$# -ne 2]; then
 echo "invalid parameters for getConf"
 exit 1
 fi

 confName="$1"
 if [-z "$confName"]; then
 echo "conf name is empty."
 exit 2
 fi

 configFile=$FLINK_HOME/conf/client.properties
 if [! -f $configFile]; then
 echo "$configFile" is not exist."
 exit 3
 fi

 defaultValue="$2"
 cnt=$(grep $1 $configFile | wc -l)
 if [$cnt -gt 1]; then
```

```
 echo "$confName" has multi values in "$configFile"
 exit 4
 elif [$cnt -lt 1]; then
 echo $defaultValue
 else
 line=$(grep $1 $configFile)
 confValue=$(echo "${line#*=}")
 echo "$confValue"
 fi
}

function createSelfSignedCA()
{
 #variable from user input
 keystorePath=$1
 storepassValue=$2
 keypassValue=$3

 #generate ca keystore
 rm -rf $keystorePath/$CA_KEYSTORE_NAME
 $KEYTOOL -genkeypair -alias $CA_ALIAS -keystore $keystorePath/$CA_KEYSTORE_NAME -dname
$CA_DNAME -storepass $storepassValue -keypass $keypassValue -validity 3650 -keyalg $CA_KEYALG -
keysize 3072 -ext bc=ca:true
 if [$? -ne 0]; then
 echo "generate ca.keystore failed."
 exit 1
 fi

 #generate ca.cer
 rm -rf "$keystorePath/ca.cer"
 $KEYTOOL -keystore "$keystorePath/$CA_KEYSTORE_NAME" -storepass "$storepassValue" -alias
$CA_ALIAS -validity 3650 -exportcert > "$keystorePath/ca.cer"
 if [$? -ne 0]; then
 echo "generate ca.cer failed."
 exit 1
 fi

 #generate ca.truststore
 rm -rf "$keystorePath/flink.truststore"
 $KEYTOOL -importcert -keystore "$keystorePath/flink.truststore" -alias $CA_ALIAS -storepass
"$storepassValue" -noprompt -file "$keystorePath/ca.cer"
 if [$? -ne 0]; then
 echo "generate ca.truststore failed."
 exit 1
 fi
}

function generateKeystore()
{
 #get path/pass from input
 keystorePath=$1
 storepassValue=$2
 keypassValue=$3

 #get value from conf
 aliasValue=$(getConf "flink.keystore.rsa.alias" "flink")
 validityValue=$(getConf "flink.keystore.rsa.validity" "3650")
 keyalgValue=$(getConf "flink.keystore.rsa.keyalg" "RSA")
 dnameValue=$(getConf "flink.keystore.rsa.dname" "CN=flink.huawei.com")
 SANValue=$(getConf "flink.keystore.rsa.ext" "ip:127.0.0.1")
 SANValue=$(echo "$SANValue" | xargs)
 SANValue="ip:$(echo "$SANValue" | sed 's/,/;/g')"

 #generate keystore
 rm -rf $keystorePath/flink.keystore
 $KEYTOOL -genkeypair -alias $aliasValue -keystore $keystorePath/flink.keystore -dname $dnameValue -
ext SAN=$SANValue -storepass $storepassValue -keypass $keypassValue -keyalg $keyalgValue -keysize
3072 -validity 3650
 if [$? -ne 0]; then
```



```
 echo "generate flink.keystore failed."
 exit 1
 fi

 #generate cer
 rm -rf $keystorePath/flink.csr
 $KEYTOOL -certreq -keystore $keystorePath/flink.keystore -storepass $storepassValue -alias $aliasValue -
file $keystorePath/flink.csr
 if [$? -ne 0]; then
 echo "generate flink.csr failed."
 exit 1
 fi

 #generate flink.cer
 rm -rf $keystorePath/flink.cer
 $KEYTOOL -gencert -keystore $keystorePath/ca.keystore -storepass $storepassValue -alias $CA_ALIAS -
ext SAN=$SANValue -infile $keystorePath/flink.csr -outfile $keystorePath/flink.cer -validity 3650
 if [$? -ne 0]; then
 echo "generate flink.cer failed."
 exit 1
 fi

 #import cer into keystore
 $KEYTOOL -importcert -keystore $keystorePath/flink.keystore -storepass $storepassValue -file
$keystorePath/ca.cer -alias $CA_ALIAS -noprompt
 if [$? -ne 0]; then
 echo "importcert ca."
 exit 1
 fi

 $KEYTOOL -importcert -keystore $keystorePath/flink.keystore -storepass $storepassValue -file
$keystorePath/flink.cer -alias $aliasValue -noprompt;
 if [$? -ne 0]; then
 echo "generate flink.truststore failed."
 exit 1
 fi
}

function configureFlinkConf()
{
 # set config
 if [-f "$CLIENT_CONF_YAML"]; then
 SSL_ENCRYPT_ENABLED=$(grep "security.ssl.encrypt.enabled" "$CLIENT_CONF_YAML" | awk '{print
$2}')
 if ["$SSL_ENCRYPT_ENABLED" = "false"];then
 sed -i s/"security.ssl.key-password:.*"/"security.ssl.key-password:"\ "${keyPass}"/g
"$CLIENT_CONF_YAML"
 if [$? -ne 0]; then
 echo "set security.ssl.key-password failed."
 return 1
 fi

 sed -i s/"security.ssl.keystore-password:.*"/"security.ssl.keystore-password:"\ "${storePass}"/g
"$CLIENT_CONF_YAML"
 if [$? -ne 0]; then
 echo "set security.ssl.keystore-password failed."
 return 1
 fi

 sed -i s/"security.ssl.truststore-password:.*"/"security.ssl.truststore-password:"\ "${storePass}"/g
"$CLIENT_CONF_YAML"
 if [$? -ne 0]; then
 echo "set security.ssl.keystore-password failed."
 return 1
 fi

 echo "security.ssl.encrypt.enabled is false, set security.ssl.key-password security.ssl.keystore-
password security.ssl.truststore-password success."
 fi
 fi
}
```

```
else
 echo "security.ssl.encrypt.enabled is true, please enter security.ssl.key-password security.ssl.keystore-
password security.ssl.truststore-password encrypted value in flink-conf.yaml."
fi

keystoreFilePath="{keystorePath}/flink.keystore
sed -i 's#"security.ssl.keystore:".*#"security.ssl.keystore:"\ "$keystoreFilePath"#g'
"$CLIENT_CONF_YAML"
if [$? -ne 0]; then
 echo "set security.ssl.keystore failed."
 return 1
fi

truststoreFilePath="{keystorePath}/flink.truststore"
sed -i 's#"security.ssl.truststore:".*#"security.ssl.truststore:"\ "$truststoreFilePath"#g'
"$CLIENT_CONF_YAML"
if [$? -ne 0]; then
 echo "set security.ssl.truststore failed."
 return 1
fi

command -v sha256sum >/dev/null
if [$? -ne 0];then
 echo "sha256sum is not exist, it will produce security.cookie with date +%F-%H-%M-%s-%N."
 cookie=$(date +%F-%H-%M-%s-%N)
else
 cookie=$(echo "${KEYTABPRINCEPAL}" | sha256sum | awk '{print $1}')
fi

sed -i s/"security.cookie:".*"/"security.cookie:"\ "${cookie}"/g "$CLIENT_CONF_YAML"
if [$? -ne 0]; then
 echo "set security.cookie failed."
 return 1
fi
fi
return 0;
}

main()
{
 #check environment variable is set or not
 if [-z ${FLINK_HOME+x}]; then
 echo "erro: environment variables are not set."
 exit 1
 fi
 stty -echo
 read -rp "Enter password:" password
 stty echo
 echo

 KEYTABPRINCEPAL=$(grep "security.kerberos.login.principal" "$CLIENT_CONF_YAML" | awk '{print $2}')
 if [-z "$KEYTABPRINCEPAL"];then
 echo "please config security.kerberos.login.principal info first."
 exit 1
 fi

 #get input
 keystorePath="$KEYSTOREPATH"
 storePass="$password"
 keyPass="$password"

 #generate self signed CA
 createSelfSignedCA "$keystorePath" "$storePass" "$keyPass"
 if [$? -ne 0]; then
 echo "create self signed ca failed."
 exit 1
 fi
}
```

```
fi

#generate keystore
generateKeystore "$keystorePath" "$storePass" "$keyPass"
if [$? -ne 0]; then
 echo "create keystore failed."
 exit 1
fi

echo "generate keystore/truststore success."

set flink config
configureFlinkConf "$keystorePath" "$storePass" "$keyPass"
if [$? -ne 0]; then
 echo "configure Flink failed."
 exit 1
fi

return 0;
}

#the start main
main "$@"

exit 0
```

### 📖 说明

执行命令 “sh generate\_keystore.sh <password>” 即可，<password>由用户自定义输入

- 若<password>中包含特殊字符"\$"，应使用如下方式，以防止被转义，“sh generate\_keystore.sh '*password*'”。命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。
- 密码不允许包含“#”。
- 使用该generate\_keystore.sh脚本前需要在客户端目录下执行source bigdata\_env。
- 使用该generate\_keystore.sh脚本会自动将security.ssl.keystore、security.ssl.truststore的绝对路径填写到flink-conf.yaml中，所以需要用户根据实际情况手动修改为相对路径。例如：
  - 将security.ssl.keystore: /opt/client/Flink/flink/conf//flink.keystore修改为 security.ssl.keystore: ssl/flink.keystore;
  - 将security.ssl.truststore: /opt/client/Flink/flink/conf//flink.truststore修改为 security.ssl.truststore: ssl/flink.truststore;
  - 需要在Flink客户端环境中任意目录下创建ssl文件夹，如在“/opt/client/Flink/flink/conf/”目录下新建目录ssl，将flink.keystore、flink.truststore文件放入ssl文件夹中；
  - 执行yarn-session或者flink run -m yarn-cluster命令时需要在ssl文件夹同级目录下执行：yarn-session.sh -t ssl -d 或者 flink run -m yarn-cluster -yt ssl -d WordCount.jar。

# 7 使用 Flume

## 7.1 Flume 日志采集概述

Flume是一个分布式、可靠和高可用的海量日志聚合的系统。它能够将不同数据源的海量日志数据进行高效收集、聚合、移动，最后存储到一个中心化数据存储系统中。支持在系统中定制各类数据发送方，用于收集数据。同时，提供对数据进行简单处理，并写到各种数据接受方（可定制）的能力。

Flume分为客户端和服务端，两者都是FlumeAgent。服务端对应着FlumeServer实例，直接部署在集群内部。而客户端部署更灵活，可以部署在集群内部，也可以部署在集群外。它们之间没有必然联系，都可以独立工作，并且提供的功能是一样的。

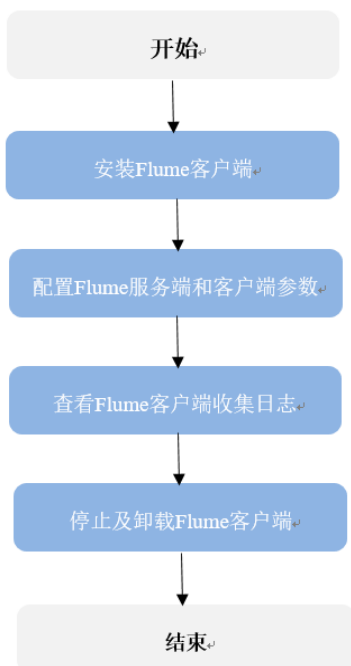
Flume客户端需要单独安装，支持将数据直接导出到集群中的HDFS和Kafka等组件上，也可以结合Flume服务端一起使用。

### 使用流程

通过同时利用Flume服务端和客户端，构成Flume的级联任务，采集日志的流程如下所示。

1. 安装Flume客户端。
2. 配置Flume服务端和客户端参数。
3. 查看Flume客户端收集日志。
4. 停止及卸载Flume客户端。

图 7-1 Flume 使用流程



## Flume 模块介绍

Flume客户端/服务端由一个或多个Agent组成，而每个Agent是由Source、Channel、Sink三个模块组成，数据先进入Source然后传递到Channel，最后由Sink发送到下一个Agent或目的地（客户端外部）。各模块说明见表7-1。

表 7-1 模块说明

| 名称     | 说明                                                                                                                                                                                                                                                                     |
|--------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Source | Source负责接收数据或产生数据，并将数据批量放到一个或多个Channel。Source有两种类型：数据驱动和轮询。<br>典型的Source样例如下： <ul style="list-style-type: none"><li>和系统集成并接收数据的Sources：Syslog、Netcat。</li><li>自动生成事件数据的Sources：Exec、SEQ。</li><li>用于Agent和Agent之间通信的IPC Sources：Avro。</li></ul> Source必须至少和一个Channel关联。 |

| 名称      | 说明                                                                                                                                                                                                                                                                                                                                                                                   |
|---------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Channel | <p>Channel位于Source和Sink之间，用于缓存Source传递的数据，当Sink成功将数据发送到下一跳的Channel或最终数据处理端，缓存数据将自动从Channel移除。</p> <p>不同类型的Channel提供的持久化水平也是不一样的：</p> <ul style="list-style-type: none"> <li>• Memory Channel：非持久化</li> <li>• File Channel：基于预写式日志（Write-Ahead Logging，简称WAL）的持久化实现</li> <li>• JDBC Channel：基于嵌入Database的持久化实现</li> </ul> <p>Channel支持事务特性，可保证简易的顺序操作，同时可以配合任意数量的Source和Sink共同工作。</p> |
| Sink    | <p>Sink负责将数据传输到下一跳或最终目的，成功完成后将数据从Channel移除。</p> <p>典型的Sink样例如下：</p> <ul style="list-style-type: none"> <li>• 存储数据到最终目的终端Sink，比如：HDFS、Kafka</li> <li>• 自动消耗的Sinks，比如：Null Sink</li> <li>• 用于Agent和Agent之间通信的IPC sink：Avro</li> </ul> <p>Sink必须关联到一个Channel。</p>                                                                                                                       |

每个Flume的Agent可以配置多个Source、Channel、Sink模块，即一个Source将数据发送给多个Channel，再由多个Sink发送到下一个Agent或目的地。

Flume支持多个Flume配置级联，即上一个Agent的Sink将数据再发送给另一个Agent的Source。

## 补充说明

### 1. Flume可靠性保障措施。

- Source与Channel、Channel与Sink之间支持事务机制。
- Sink Processor支持配置failover、load\_balance机制。

例如load\_balance示例如下：

```
server.sinkgroups=g1
server.sinkgroups.g1.sinks=k1 k2
server.sinkgroups.g1.processor.type=load_balance
server.sinkgroups.g1.processor.backoff=true
server.sinkgroups.g1.processor.selector=random
```

### 2. Flume多客户端聚合级联时的注意事项。

- 级联时需要走Avro或者Thrift协议进行级联。
- 聚合端存在多个节点时，连接配置尽量配置均衡，不要聚合到单节点上。

### 3. Flume客户端可以包含多个独立的数据流，即在一个配置文件properties.properties中配置多个Source、Channel、Sink。这些组件可以链接以形成多个流。

例如在一个配置中配置两个数据流，示例如下：

```
server.sources = source1 source2
server.sinks = sink1 sink2
```

```
server.channels = channel1 channel2

#dataflow1
server.sources.source1.channels = channel1
server.sinks.sink1.channel = channel1

#dataflow2
server.sources.source2.channels = channel2
server.sinks.sink2.channel = channel2
```

## 7.2 Flume 业务模型配置说明

### 业务模型配置指导

本章节适用于MRS 3.x及之后版本。

本任务旨在提供Flume常用模块的性能差异，用于指导用户进行合理的Flume业务配置，避免出现前端Source和后端Sink性能不匹配进而导致整体业务性能不达标的场景。

本任务只针对于单通道的场景进行比较说明。

Flume业务配置及模块选择过程中，一般要求Sink的极限吞吐量需要大于Source的极限吞吐量，否则在极限负载的场景下，Source往Channel的写入速度大于Sink从Channel取出的速度，从而导致Channel频繁被写满，进而影响性能表现。

Avro Source和Avro Sink一般都是成对出现，用于多个Flume Agent间进行数据中转，因此一般场景下Avro Source和Avro Sink都不会成为性能瓶颈。

### 模块间性能

根据模块间性能对比，可以看到对于前端是SpoolDir Source的场景下，Kafka Sink和HDFS Sink都能满足吞吐量要求，但是HBase Sink由于自身写入性能较低的原因，会成为性能瓶颈，会导致数据都积压在Channel中。但是如果必须使用HBase Sink或者其他性能容易成为瓶颈的Sink的场景时，可以选择使用**Channel Selector**或者**Sink Group**来满足性能要求。

### Channel Selector

Channel Selector可以允许一个Source对接多个Channel，通过选择不同的Selector类型来将Source的数据进行分流或者复制，目前Flume提供的Channel Selector有两种：Replicating和Multiplexing。

Replicating：表示Source的数据同步发送给所有Channel。

Multiplexing：表示根据Event中的Header的指定字段的值来进行判断，从而选择相应的Channel进行发送，从而起到根据业务类型进行分流的目的。

- Replicating配置样例：

```
client.sources = kafkasource
client.channels = channel1 channel2
client.sources.kafkasource.type = org.apache.flume.source.kafka.KafkaSource
client.sources.kafkasource.kafka.topics = topic1,topic2
client.sources.kafkasource.kafka.consumer.group.id = flume
client.sources.kafkasource.kafka.bootstrap.servers = 10.69.112.108:21007
client.sources.kafkasource.kafka.security.protocol = SASL_PLAINTEXT
client.sources.kafkasource.batchDurationMillis = 1000
client.sources.kafkasource.batchSize = 800
client.sources.kafkasource.channels = channel1 channel2
```

```
client.sources.kafkasource.selector.type = replicating
client.sources.kafkasource.selector.optional = channel2
```

表 7-2 Replicating 配置样例参数说明

| 选项名称              | 默认值         | 描述                          |
|-------------------|-------------|-----------------------------|
| Selector.type     | replicating | Selector类型，应配置为 replicating |
| Selector.optional | -           | 可选Channel，可以配置为列表           |

- Multiplexing配置样例：

```
client.sources = kafkasource
client.channels = channel1 channel2
client.sources.kafkasource.type = org.apache.flume.source.kafka.KafkaSource
client.sources.kafkasource.kafka.topics = topic1,topic2
client.sources.kafkasource.kafka.consumer.group.id = flume
client.sources.kafkasource.kafka.bootstrap.servers = 10.69.112.108:21007
client.sources.kafkasource.kafka.security.protocol = SASL_PLAINTEXT
client.sources.kafkasource.batchDurationMillis = 1000
client.sources.kafkasource.batchSize = 800
client.sources.kafkasource.channels = channel1 channel2

client.sources.kafkasource.selector.type = multiplexing
client.sources.kafkasource.selector.header = myheader
client.sources.kafkasource.selector.mapping.topic1 = channel1
client.sources.kafkasource.selector.mapping.topic2 = channel2
client.sources.kafkasource.selector.default = channel1
```

表 7-3 Multiplexing 配置样例参数说明

| 选项名称               | 默认值                   | 描述                           |
|--------------------|-----------------------|------------------------------|
| Selector.type      | replicating           | Selector类型，应配置为 multiplexing |
| Selector.header    | Flume.selector.header | -                            |
| Selector.default   | -                     | -                            |
| Selector.mapping.* | -                     | -                            |

Multiplexing类型的Selector的样例中，选择Event中Header名称为topic的字段来进行判断，当Header中topic字段的值为topic1时，向channel1发送该Event，当Header中topic字段的值为topic2时，向channel2发送该Event。

这种Selector需要借助Source中Event的特定Header来进行Channel的选择，需要根据业务场景选择合理的Header来进行数据分流。

## SinkGroup

当后端单Sink性能不足、需要高可靠性保证或者异构输出时可以使用Sink Group来将指定的Channel和多个Sink对接，从而满足相应的使用场景。目前Flume提供了两种Sink Processor用于对Sink Group中的Sink进行管理：Load Balancing和Failover。



**Failover:** 表示在Sink Group中同一时间只有一个Sink处于活跃状态，其他Sink作为备份处于非活跃状态，当活跃状态的Sink故障时，根据优先级从非活跃状态的Sink中选择一个来接管业务，保证数据不会丢失，多用于高可靠性场景。

**Load Balancing:** 表示在Sink Group中所有Sink都处于活跃状态，每个Sink都会从Channel中去获取数据并进行处理，并且保证在运行过程中该Sink Group的所有Sink的负载是均衡的，多用于性能提升场景。

- **Load Balancing配置样例:**

```
client.sources = source1
client.sinks = sink1 sink2
client.channels = channel1

client.sinkgroups = g1
client.sinkgroups.g1.sinks = sink1 sink2
client.sinkgroups.g1.processor.type = load_balance
client.sinkgroups.g1.processor.backoff = true
client.sinkgroups.g1.processor.selector = random

client.sinks.sink1.type = logger
client.sinks.sink1.channel = channel1

client.sinks.sink2.type = logger
client.sinks.sink2.channel = channel1
```

**表 7-4** Load Balancing 配置样例参数说明

| 选项名称                          | 默认值         | 描述                                                            |
|-------------------------------|-------------|---------------------------------------------------------------|
| sinks                         | -           | Sink Group的sink列表，多个以空格分隔                                     |
| processor.type                | default     | Processor的类型，应配置为load_balance                                 |
| processor.backoff             | false       | 是否以指数的形式退避失败的Sinks                                            |
| processor.selector            | round_robin | 选择机制。必须是round_robin, random或者自定义的类，且该类继承了AbstractSinkSelector |
| processor.selector.maxTimeOut | 30000       | 屏蔽故障sink的时间，默认是30000毫秒                                        |

- **Failover配置样例:**

```
client.sources = source1
client.sinks = sink1 sink2
client.channels = channel1

client.sinkgroups = g1
client.sinkgroups.g1.sinks = sink1 sink2
client.sinkgroups.g1.processor.type = failover
client.sinkgroups.g1.processor.priority.sink1 = 10
client.sinkgroups.g1.processor.priority.sink2 = 5
client.sinkgroups.g1.processor.maxpenalty = 10000

client.sinks.sink1.type = logger
client.sinks.sink1.channel = channel1
```

```
client.sinks.sink2.type = logger
client.sinks.sink2.channel = channel1
```

表 7-5 Failover 配置样例参数说明

| 选项名称                          | 默认值     | 描述                                                                                                        |
|-------------------------------|---------|-----------------------------------------------------------------------------------------------------------|
| sinks                         | -       | Sink Group的sink列表，多个以空格分隔                                                                                 |
| processor.type                | default | Processor的类型，应配置为failover                                                                                 |
| processor.priority.<sinkName> | -       | 优先级值。<sinkName> 必须是sinks中有定义的。优先级值高Sink会更早被激活。值越大，优先级越高。 <b>注：</b> 多个sinks的话，优先级的值不要相同，如果优先级相同的话，只会有一个生效。 |
| processor.maxpenalty          | 30000   | 失败的Sink最大的退避时间(单位：毫秒)                                                                                     |

## Interceptors

Flume的拦截器（Interceptor）支持在数据传输过程中修改或丢弃传输的基本单元Event。用户可以通过在配置中指定Flume内建拦截器的类名列表，也可以开发自定义的拦截器来实现Event的修改或丢弃。Flume内建支持的拦截器如下表所示，本章节选取一个较为复杂的作为示例。其余的用户可以根据需要自行配置使用。官网参考：<http://flume.apache.org/releases/content/1.9.0/FlumeUserGuide.html>

### 说明

1. 拦截器用在Flume的Source、Channel之间，大部分的Source都带有Interceptor参数。用户可以依据需要配置。
2. Flume支持一个Source配置多个拦截器，各拦截器名称用空格分开。
3. 指定拦截器的顺序就是它们被调用的顺序。
4. 使用拦截器在Header中插入的内容，都可以在Sink中读取并使用。

表 7-6 Flume 内建支持的拦截器类型

| 拦截器类型                     | 简要描述                                     |
|---------------------------|------------------------------------------|
| Timestamp Interceptor     | 该拦截器会在Event的Header中插入一个时间戳。              |
| Host Interceptor          | 该拦截器会在Event的Header中插入当前Agent所在节点的IP或主机名。 |
| Remove Header Interceptor | 该拦截器会依据Header中包含的符合正则匹配的字符串，丢弃掉对应的Event。 |
| UUID Interceptor          | 该拦截器会为每个Event的Header生成一个UUID字符串。         |

| 拦截器类型                          | 简要描述                                                                    |
|--------------------------------|-------------------------------------------------------------------------|
| Search and Replace Interceptor | 该拦截器基于Java正则表达式提供简单的基于字符串的搜索和替换功能。与Java Matcher.replaceAll() 的规则相同。     |
| Regex Filtering Interceptor    | 该拦截器通过将Event的Body体解释为文本文件，与配置的正则表达式进行匹配来选择性的过滤Event。提供的正则表达式可用于排除或包含事件。 |
| Regex Extractor Interceptor    | 该拦截器使用正则表达式抽取原始events中的内容，并将该内容加入events的header中。                        |

下面以Regex Filtering Interceptor 为例说明Interceptor使用（其余的可参考官网配置）：

表 7-7 Regex Filtering Interceptor 配置参数说明

| 选项名称          | 默认值   | 描述                                         |
|---------------|-------|--------------------------------------------|
| type          | -     | 组件类型名称，必须写为regex_filter。                   |
| regex         | -     | 用于匹配事件的正则表达式。                              |
| excludeEvents | false | 默认收集匹配到的Event。设置为true，则会删除匹配的Event，保留不匹配的。 |

配置示例（为了方便观察，此模型使用了netcat tcp作为Source源，logger作为Sink）。配置好如下参数后，在Linux的配置的主机节点上执行Linux命令“telnet 主机名或IP 44444”，并任意敲入符合正则和不符合正则的字符串。会在日志中观察到，只有匹配到的字符串被传输了。

```
#define the source、channel、sink
server.sources = r1

server.channels = c1
server.sinks = k1

#config the source
server.sources.r1.type = netcat
server.sources.r1.bind = ${主机IP}
server.sources.r1.port = 44444
server.sources.r1.interceptors= i1
server.sources.r1.interceptors.i1.type= regex_filter
server.sources.r1.interceptors.i1.regex= (flume)|(myflume)
server.sources.r1.interceptors.i1.excludeEvents= false
server.sources.r1.channels = c1

#config the channel
server.channels.c1.type = memory
server.channels.c1.capacity = 1000
server.channels.c1.transactionCapacity = 100

#config the sink
server.sinks.k1.type = logger
server.sinks.k1.channel = c1
```

## 7.3 安装 Flume 客户端

### 7.3.1 安装 MRS 3.x 之前版本 Flume 客户端

#### 操作场景

使用Flume搜集日志时，需要在日志主机上安装Flume客户端。用户可以创建一个新的ECS并安装Flume客户端。

本章节适用于MRS 3.x之前版本。

#### 前提条件

- 已创建包含Flume组件的流集群。
- 日志主机需要与MRS集群在相同的VPC和子网。
- 已获取日志主机的登录方式。

#### 操作步骤

**步骤1** 根据前提条件，创建一个满足要求的弹性云服务器。

**步骤2** 登录集群详情页面，选择“组件管理”。

#### 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

**步骤3** 单击“下载客户端”。

1. 在“客户端类型”选择“完整客户端”。
2. 在“下载路径”选择“远端主机”。
3. 将“主机IP”设置为ECS的IP地址，设置“主机端口”为“22”，并将“保存路径”设置为“/tmp”。
  - 如果使用SSH登录ECS的默认端口“22”被修改，请将“主机端口”设置为新端口。
  - “保存路径”最多可以包含256个字符。
4. “登录用户”设置为“root”。

如果使用其他用户，请确保该用户对保存目录拥有读取、写入和执行权限。
5. 在“登录方式”选择“密码”或“SSH私钥”。
  - 密码：输入创建集群时设置的root用户密码。
  - SSH私钥：选择并上传创建集群时使用的密钥文件。
6. 单击“确定”开始生成客户端文件。

若界面显示以下提示信息表示客户端包已经成功保存。

下载客户端文件到远端主机成功。

若界面显示以下提示信息，请检查用户名密码及远端主机的安全组配置，确保用户名密码正确，及远端主机的安全组已增加SSH(22)端口的入方向规则。然后从**步骤3**执行重新下载客户端。

连接到服务器失败，请检查网络连接或参数设置。

图 7-2 下载客户端

### 下载客户端

**警告：**生成客户端会占用大量的磁盘IO，不建议在集群处于安装中、启动中、打补丁中等非稳态场景进行“下载客户端”操作。

\* 客户端类型  完整客户端  仅配置文件

\* 下载路径  服务器端  远端主机

仅保存到服务器如下路径，如果存在客户端文件，会覆盖路径下已有的客户端文件。

\* 主机IP

\* 主机端口

\* 保存路径

\* 登录用户

\* 登录方式  密码  SSH私钥

\* 密码

**确定**

**步骤4** 选择“Flume”服务，单击“实例”，查看任意一个Flume实例和两个MonitorServer实例的“业务IP”。



**步骤5** 使用VNC方式，登录弹性云服务器。参见[远程登录（VNC方式）](#)。

所有镜像均支持Cloud-init特性。Cloud-init预配置的用户名“root”，密码为创建集群时设置的密码。首次登录建议修改。

**步骤6** 在弹性云服务器，切换到root用户，并将安装包复制到目录“/opt”。

```
sudo su - root
```

```
cp /tmp/MRS_Flume_Client.tar /opt
```

**步骤7** 在“/opt”目录执行以下命令，解压压缩包获取校验文件与客户端配置包。

```
tar -xvf MRS_Flume_Client.tar
```

**步骤8** 执行以下命令，校验文件包。

```
sha256sum -c MRS_Flume_ClientConfig.tar.sha256
```

界面显示如下信息，表明文件包校验成功：

```
MRS_Flume_ClientConfig.tar: OK
```

**步骤9** 执行以下命令，解压“MRS\_Flume\_ClientConfig.tar”。

```
tar -xvf MRS_Flume_ClientConfig.tar
```

**步骤10** 执行以下命令，安装客户端运行环境到新的目录，例如“/opt/Flumeenv”。安装时自动生成目录。

```
sh /opt/MRS_Flume_ClientConfig/install.sh /opt/Flumeenv
```

查看安装输出信息，如有以下结果表示客户端运行环境安装成功：

```
Components client installation is complete.
```

**步骤11** 执行以下命令，配置环境变量。

```
source /opt/Flumeenv/bigdata_env
```

**步骤12** 执行以下命令，解压Flume客户端。

```
cd /opt/MRS_Flume_ClientConfig/Flume
```

```
tar -xvf FusionInsight-Flume-1.6.0.tar.gz
```

**步骤13** 执行以下命令，查看当前用户密码是否过期。

```
chage -l root
```

“Password expires”时间早于当前则表示过期。此时需要修改密码，或执行**chage -M -1 root**设置密码为未过期状态。

**步骤14** 执行以下命令，安装Flume客户端到新目录，例如“/opt/FlumeClient”。安装时自动生成目录。

```
sh /opt/MRS_Flume_ClientConfig/Flume/install.sh -d /opt/FlumeClient -f
MonitorServer实例的业务IP地址 -c Flume配置文件路径 -l /var/log/ -e Flume的业
务IP地址 -n Flume客户端名称
```

各参数说明如下：

- “-d”：表示Flume客户端安装路径。
- “-f”：可选参数，表示两个MonitorServer角色的业务IP地址，中间用英文逗号分隔，若不设置则Flume客户端将不向MonitorServer发送告警信息，同时在MRS Manager界面上看不到该客户端的相关信息。

- “-c”：可选参数，表示Flume客户端在安装后默认加载的配置文件“properties.properties”。如不添加参数，默认使用客户端安装目录的“fusioninsight-flume-1.6.0/conf/properties.properties”。客户端中配置文件为空白模板，根据业务需要修改后Flume客户端将自动加载。
- “-l”：可选参数，表示日志目录，默认值为“/var/log/Bigdata”。
- “-e”：可选参数，表示Flume实例的业务IP地址，主要用于接收客户端上报的监控指标信息。
- “-n”：可选参数，表示自定义的Flume客户端的名称。
- IBM的JDK不支持“-Xloggc”，需要修改“flume/conf/flume-env.sh”，将“-Xloggc”修改为“-Xverbosegclog”，若JDK为32位，“-Xmx”不能大于3.25GB。
- “flume/conf/flume-env.sh”中，“-Xmx”默认为4GB。若客户端机器内存过小，可调整为512M甚至1GB。

例如执行：`sh install.sh -d /opt/FlumeClient`

系统显示以下结果表示客户端运行环境安装成功：

```
install flume client successfully.
```

----结束

## 7.3.2 安装 MRS 3.x 及之后版本 Flume 客户端

### 操作场景

使用Flume搜集日志时，需要在日志主机上安装Flume客户端。用户可以创建一个新的ECS并安装Flume客户端。

本章节适用于MRS 3.x及之后版本。

### 前提条件

- 已创建包含Flume组件的集群。
- 日志主机需要与MRS集群在相同的VPC和子网。
- 已获取日志主机的登录方式。
- 安装目录可以不存在，会自动创建。但如果存在，则必须为空。目录路径不能包含空格。

### 操作步骤

#### 步骤1 获取软件包。

登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume”进入Flume服务界面，在右上角选择“更多 > 下载客户端”，选择“选择客户端类型”为“完整客户端”，下载Flume服务客户端文件。

客户端文件名称为“FusionInsight\_Cluster\_<集群ID>\_Flume\_Client.tar”，本章节以“FusionInsight\_Cluster\_1\_Flume\_Client.tar”为例进行描述。



## 步骤2 上传软件包。

以user用户将软件包上传到将要安装Flume服务客户端的节点目录上，例如“/opt/client”。

### 说明

user用户为安装和运行Flume客户端的用户。

## 步骤3 解压软件包。

以user用户登录将要安装Flume服务客户端的节点。进入安装包所在目录，例如“/opt/client”，执行如下命令解压安装包到当前目录。

```
cd /opt/client
```

```
tar -xvf FusionInsight_Cluster_1_Flume_Client.tar
```

## 步骤4 校验软件包。

执行sha256sum -c命令校验解压得到的文件，返回“OK”表示校验通过。例如：

```
sha256sum -c FusionInsight_Cluster_1_Flume_ClientConfig.tar.sha256
```

```
FusionInsight_Cluster_1_Flume_ClientConfig.tar: OK
```

## 步骤5 解压文件。

```
tar -xvf FusionInsight_Cluster_1_Flume_ClientConfig.tar
```

## 步骤6 若在集群外节点安装Flume客户端，需执行如下步骤配置安装环境。在集群内节点安装可不执行该步骤。

1. 执行以下命令，安装客户端运行环境到新的目录，例如“/opt/Flumeenv”。安装时自动生成目录。

```
sh /opt/client/FusionInsight_Cluster_1_Flume_ClientConfig/install.sh /opt/Flumeenv
```

查看安装输出信息，如有以下结果表示客户端运行环境安装成功：

```
Components client installation is complete.
```

2. 执行以下命令，配置环境变量。

```
source /opt/Flumeenv/bigdata_env
```

## 步骤7 在Flume客户端安装目录下执行以下命令，安装客户端到指定目录（绝对路径），例如安装到“/opt/FlumeClient”目录。客户端安装成功后安装结束。

```
cd /opt/client/FusionInsight_Cluster_1_Flume_ClientConfig/Flume/FlumeClient
```

```
./install.sh -d /opt/FlumeClient -f MonitorServer角色的业务IP或主机名 -c 用户业务配置文件properties.properties放置路径 -s cpu阈值 -l /var/log/Bigdata -e FlumeServer的业务IP或主机名 -n Flume
```



## 📖 说明

- “-d”：Flume客户端安装路径。
- “-f”（可选）：两个MonitorServer角色的业务IP或主机名，中间用逗号分隔，若不设置则Flume客户端将不向MonitorServer发送告警信息，同时在FusionInsight Manager界面上看不到该客户端的相关信息。
- “-c”（可选）：指定业务配置文件，该文件需要用户根据自己业务生成，具体操作可在Flume服务端中“配置工具”页面参考[Flume业务配置指南](#)章节生成，并上传到待安装客户端节点上的任一目录下。若安装时未指定（即不配置该参数），可在安装后上传已经生成的业务配置文件properties.properties到“/opt/FlumeClient/fusioninsight-flume-1.9.0/conf”目录下。
- “-s”（可选）：Cgroup阈值，阈值取值范围为1~100\*N之间的整数，N表示机器cpu核数。默认阈值为“-1”，表示加入到Cgroup的进程不受cpu使用率限制。
- “-l”（可选）：日志路径，默认值为“/var/log/Bigdata”（“user”用户需要对此目录有写权限）。首次安装客户端会生成名为flume-client的子目录，之后安装会依次生成名为“flume-client-n”的子目录，n代表一个序号，从1依次递增。在Flume客户端安装目录下的conf目录中，编辑ENV\_VARS文件，搜索FLUME\_LOG\_DIR属性，可查看客户端日志路径。
- “-e”（可选）：FlumeServer的业务IP地址或主机名，主要用于接收客户端上报的监控指标信息。
- “-n”（可选）：Flume客户端的名称，可以通过在FusionInsight Manager上选择“集群 > 待操作集群名称 > 服务 > Flume > Flume管理”查看对应节点上客户端的名称。
- 若产生以下错误提示，可执行命令**export JAVA\_HOME=JDK路径**进行处理。可使用**echo \$JAVA\_HOME**查找JDK路径。  
JAVA\_HOME is null in current user,please install the JDK and set the JAVA\_HOME
- 集群混搭时，安装跨平台客户端时，请进入/opt/client/FusionInsight\_Cluster\_1\_Flume\_ClientConfig/Flume/FusionInsight-Flume-1.9.0.tar.gz路径下进行Flume客户端安装。

----结束

## Flume 客户端 Cgroup 使用指导

### • 加入Cgroup

执行以下命令，加入Cgroup，假设Flume客户端安装路径为“/opt/FlumeClient”，Cgroup cpu阈值设置为50%：

```
cd /opt/FlumeClient/fusioninsight-flume-1.9.0/bin
./flume-manage.sh cgroup join 50
```

### 📖 说明

- 该命令不仅可以加入Cgroup，同时也可以更改Cgroup cpu阈值。
- Cgroup cpu阈值取值范围为1~100\*N之间的整数，N表示机器cpu核数。
- **查询Cgroup状态**  
执行以下命令，查询Cgroup状态，假设Flume客户端安装路径为“/opt/FlumeClient”：  

```
cd /opt/FlumeClient/fusioninsight-flume-1.9.0/bin
./flume-manage.sh cgroup status
```
- **退出Cgroup**  
执行以下命令，退出Cgroup，假设Flume客户端安装路径为“/opt/FlumeClient”：  

```
cd /opt/FlumeClient/fusioninsight-flume-1.9.0/bin
```

### ./flume-manage.sh cgroup exit

#### 📖 说明

- 客户端安装完成后，会自动创建默认Cgroup。若安装客户端时未配置“-s”参数，则默认值为“-1”，表示agent进程不受cpu使用率限制。
- 加入、退出Cgroup时，agent进程不受影响。若agent进程未启动，加入、退出Cgroup仍然可以成功执行，待下一次agent启动时生效。
- 客户端卸载完成后，安装时期创建的Cgroup会自动删除。

## 7.4 快速使用 Flume 采集节点日志

### 操作场景

Flume支持将采集的日志信息导入到Kafka。

### 前提条件

- 已创建开启Kerberos认证的包含Flume、Kafka等组件的流式集群。可参考[购买自定义集群](#)。
- 已配置网络，使日志生成节点与流集群互通。

### 使用 Flume 客户端（MRS 3.x 之前版本）

#### 📖 说明

普通集群不需要执行[步骤2-步骤6](#)。

#### 步骤1 安装Flume客户端。

可参考[安装MRS 3.x之前版本Flume客户端](#)在日志生成节点安装Flume客户端，例如安装目录为“/opt/Flumeclient”。以下操作的客户端目录只是举例，请根据实际安装目录修改。

#### 步骤2 将Master1节点上的认证服务器配置文件，复制到安装Flume客户端的节点，保存到Flume客户端中“Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/conf”目录下。

文件完整路径为\${BIGDATA\_HOME}/MRS\_Current/1\_X\_KerberosClient/etc/kdc.conf。

其中“X”为随机生成的数字，请根据实际情况修改。同时文件需要以Flume客户端安装用户身份保存，例如root用户。

#### 步骤3 查看任一部署Flume角色节点的“业务IP”。

登录集群详情页面，选择“集群 > 组件管理 > Flume > 实例”，查看任一部署Flume角色节点的“业务IP”。

#### 📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。



The screenshot shows a management interface with a navigation bar at the top containing: 概览, 节点管理, 组件管理 (selected), 告警管理, 文件管理, 作业管理, 租户管理, 标签管理. Below the navigation bar, there is a sub-header '服务 Flume / 实例' and tabs for '服务状态' (selected) and '服务配置'. A table lists the installed services:

| <input type="checkbox"/> | 角色            | 主名称              | 管理IP          | 业务IP          |
|--------------------------|---------------|------------------|---------------|---------------|
| <input type="checkbox"/> | Flume         | node-master3this | 192.168.0.101 | 192.168.0.101 |
| <input type="checkbox"/> | MonitorServer | node-master3this | 192.168.0.101 | 192.168.0.101 |

**步骤4** 将此节点上的用户认证文件，复制到安装Flume客户端的节点，保存到Flume客户端中“Flume客户端安装目录/fusioninsight-flume-*Flume组件版本号*/conf”目录下。

文件完整路径为 $\${BIGDATA\_HOME}/MRS\_XXX/install/FusionInsight-Flume-Flume组件版本号/flume/conf/flume.keytab$ 。

其中“XXX”为产品版本号，请根据实际情况修改。同时文件需要以Flume客户端安装用户身份保存，例如root用户。

**步骤5** 将此节点上的配置文件“jaas.conf”，复制到安装Flume客户端的节点，保存到Flume客户端中“conf”目录。

文件完整路径为 $\${BIGDATA\_HOME}/MRS\_Current/1\_X\_Flume/etc/jaas.conf$ 。

其中“X”为随机生成的数字，请根据实际情况修改。同时文件需要以Flume客户端安装用户身份保存，例如root用户。

**步骤6** 登录安装Flume客户端节点，切换到客户端安装目录，执行以下命令修改文件：

```
vi conf/jaas.conf
```

修改参数“keyTab”定义的用户认证文件完整路径即**步骤4**中保存用户认证文件的目录：“Flume客户端安装目录/fusioninsight-flume-*Flume组件版本号*/conf”，然后保存并退出。

**步骤7** 执行以下命令，修改Flume客户端配置文件“flume-env.sh”：

```
vi Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/conf/flume-env.sh
```

在“-XX:+UseCMSCompactAtFullCollection”后面，增加以下内容：

```
-Djava.security.krb5.conf=Flume客户端安装目录/fusioninsight-flume-1.9.0/conf/kdc.conf -
Djava.security.auth.login.config=Flume客户端安装目录/fusioninsight-flume-1.9.0/conf/jaas.conf -
Dzookeeper.request.timeout=120000
```

例如：“-XX:+UseCMSCompactAtFullCollection -Djava.security.krb5.conf=/opt/FlumeClient/fusioninsight-flume-*Flume组件版本号*/conf/kdc.conf -Djava.security.auth.login.config=/opt/FlumeClient/fusioninsight-flume-*Flume组件版本号*/conf/jaas.conf -Dzookeeper.request.timeout=120000”

请根据实际情况，修改“Flume客户端安装目录”，然后保存并退出。

**步骤8** 执行以下命令，重启Flume客户端：

```
cd Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/bin
./flume-manage.sh restart
```

例如：

```
cd /opt/FlumeClient/fusioninsight-flume-Flume组件版本号/bin
```

## ./flume-manage.sh restart

### 📖 说明

Flume客户端停止后会自动重启，如果不需自动重启，请执行以下命令：

```
./flume-manage.sh stop force
```

需要启动时，可执行以下命令：

```
./flume-manage.sh start force
```

**步骤9** 执行以下命令，根据实际业务需求，可参考[Flume配置参数说明](#)在Flume客户端配置文件“properties.properties”中配置并保存作业。

### vi *Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/conf/properties.properties*

以配置SpoolDir Source+File Channel+Kafka Sink为例：

```


client.sources = static_log_source
client.channels = static_log_channel
client.sinks = kafka_sink

#LOG_TO_HDFS_ONLINE_1

client.sources.static_log_source.type = spooldir
client.sources.static_log_source.spoolDir = 监控目录
client.sources.static_log_source.fileSuffix = .COMPLETED
client.sources.static_log_source.ignorePattern = ^$
client.sources.static_log_source.trackerDir = 传输过程中元数据存储路径
client.sources.static_log_source.maxBlobLength = 16384
client.sources.static_log_source.batchSize = 51200
client.sources.static_log_source.inputCharset = UTF-8
client.sources.static_log_source.deserializer = LINE
client.sources.static_log_source.selector.type = replicating
client.sources.static_log_source.fileHeaderKey = file
client.sources.static_log_source.fileHeader = false
client.sources.static_log_source.basenameHeader = true
client.sources.static_log_source.basenameHeaderKey = basename
client.sources.static_log_source.deletePolicy = never

client.channels.static_log_channel.type = file
client.channels.static_log_channel.dataDirs = 数据缓存路径，设置多个路径可提升性能，中间用逗号分开
client.channels.static_log_channel.checkpointDir = 检查点存放路径
client.channels.static_log_channel.maxFileSize = 2146435071
client.channels.static_log_channel.capacity = 1000000
client.channels.static_log_channel.transactionCapacity = 612000
client.channels.static_log_channel.minimumRequiredSpace = 524288000

client.sinks.kafka_sink.type = org.apache.flume.sink.kafka.KafkaSink
client.sinks.kafka_sink.kafka.topic = 数据写入的topic，如flume_test
client.sinks.kafka_sink.kafka.bootstrap.servers = XXX.XXX.XXX.XXX:kafka端口号,XXX.XXX.XXX.XXX:kafka端口号,XXX.XXX.XXX.XXX:kafka端口号
client.sinks.kafka_sink.flumeBatchSize = 1000
client.sinks.kafka_sink.kafka.producer.type = sync
client.sinks.kafka_sink.kafka.security.protocol = SASL_PLAINTEXT
client.sinks.kafka_sink.kafka.kerberos.domain.name = Kafka Domain名称，安全集群必填，如hadoop.example.com
client.sinks.kafka_sink.requiredAcks = 0

client.sources.static_log_source.channels = static_log_channel
client.sinks.kafka_sink.channel = static_log_channel
```

### 📖 说明

- `client.sinks.kafka_sink.kafka.topic`: 数据写入的topic。若kafka中该topic不存在，默认情况下会自动创建该topic。
- `client.sinks.kafka_sink.kafka.bootstrap.servers`: Kafka brokers列表，多个用英文逗号分隔。默认情况下，安全集群端口21007，普通集群对应端口9092。
- `client.sinks.kafka_sink.kafka.security.protocol`: 安全集群为SASL\_PLAINTEXT，普通集群为PLAINTEXT。
- `client.sinks.kafka_sink.kafka.kerberos.domain.name`:  
普通集群无需配置此参数。安全集群对应此参数的值为Kafka集群中“`kerberos.domain.name`”对应的值。  
具体可到Broker实例所在节点上查看`${BIGDATA_HOME}/MRS_Current/1_X_Broker/etc/server.properties`。  
其中X为随机生成的数字，请根据实际情况修改。同时文件需要以Flume客户端安装用户身份保存，例如root用户。

**步骤10** 参数配置并保存后，Flume客户端将自动加载“`properties.properties`”中配置的内容。当`spoolDir`生成新的日志文件，文件内容将发送到Kafka生产者，并支持Kafka消费者消费。

---结束

## 使用 Flume 客户端（MRS 3.x 及之后版本）

### 📖 说明

普通集群不需要执行**步骤2-步骤6**。

**步骤1** 安装Flume客户端。

可参考[安装MRS 3.x及之后版本Flume客户端](#)在日志生成节点安装Flume客户端，例如安装目录为“`/opt/Flumeclient`”。以下操作的客户端目录只是举例，请根据实际安装目录修改。

**步骤2** 将Master1节点上的认证服务器配置文件，复制到安装Flume客户端的节点，保存到Flume客户端中`Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/conf`目录下。

文件完整路径为“`${BIGDATA_HOME}/FusionInsight_BASE_XXX/1_X KerberosClient/etc/kdc.conf`”。其中“XXX”为产品版本号，“X”为随机生成的数字，请根据实际情况修改。同时文件需要以Flume客户端安装用户身份保存，例如root用户。

**步骤3** 查看任一部署Flume角色节点的“业务IP”。

登录FusionInsight Manager页面，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)，选择“集群 > 服务 > Flume > 实例”。查看任一部署Flume角色节点的“业务IP”。

### 📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

**步骤4** 将此节点上的用户认证文件，复制到安装Flume客户端的节点，保存到Flume客户端中“`Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/conf`”目录下。

文件完整路径为\${BIGDATA\_HOME}/FusionInsight\_Porter\_XXX/install/FusionInsight-Flume-*Flume组件版本号*/flume/conf/flume.keytab。

其中“XXX”为产品版本号，请根据实际情况修改。同时文件需要以Flume客户端安装用户身份保存，例如root用户。

**步骤5** 将此节点上的配置文件“jaas.conf”，复制到安装Flume客户端的节点，保存到Flume客户端中“conf”目录。

文件完整路径为\${BIGDATA\_HOME}/FusionInsight\_Current/1\_X\_Flume/etc/jaas.conf。

其中“X”为随机生成的数字，请根据实际情况修改。同时文件需要以Flume客户端安装用户身份保存，例如root用户。

**步骤6** 登录安装Flume客户端节点，切换到客户端安装目录，执行以下命令修改文件：

```
vi conf/jaas.conf
```

修改参数“keyTab”定义的用户认证文件完整路径即**步骤4**中保存用户认证文件的目录：“Flume客户端安装目录/fusioninsight-flume-*Flume组件版本号*/conf”，然后保存并退出。

**步骤7** 执行以下命令，修改Flume客户端配置文件“flume-env.sh”：

```
vi Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/conf/flume-env.sh
```

在“-XX:+UseCMSCompactAtFullCollection”后面，增加以下内容：

```
-Djava.security.krb5.conf=Flume客户端安装目录/fusioninsight-flume-1.9.0/conf/kdc.conf -
Djava.security.auth.login.config=Flume客户端安装目录/fusioninsight-flume-1.9.0/conf/jaas.conf -
Dzookeeper.request.timeout=120000
```

例如：“-XX:+UseCMSCompactAtFullCollection -Djava.security.krb5.conf=/opt/FlumeClient/fusioninsight-flume-*Flume组件版本号*/conf/kdc.conf -Djava.security.auth.login.config=/opt/FlumeClient/fusioninsight-flume-*Flume组件版本号*/conf/jaas.conf -Dzookeeper.request.timeout=120000”

请根据实际情况，修改“Flume客户端安装目录”，然后保存并退出。

**步骤8** 执行以下命令，重启Flume客户端：

```
cd Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/bin
./flume-manage.sh restart
```

例如：

```
cd /opt/FlumeClient/fusioninsight-flume-Flume组件版本号/bin
./flume-manage.sh restart
```

**步骤9** 根据实际业务场景配置作业。

- MRS 3.x及之后版本部分参数可直接在Manager界面配置，可参考[非加密传输](#)或[加密传输](#)。
- 在“properties.properties”文件中配置，以配置SpoolDir Source+File Channel+Kafka Sink为例。

在安装Flume客户端的节点执行以下命令，根据实际业务需求，可参考[Flume业务配置指南](#)在Flume客户端配置文件“properties.properties”中配置并保存作业。

## vi *Flume*客户端安装目录/fusioninsight-flume-*Flume*组件版本号/conf/properties.properties

```


client.sources = static_log_source
client.channels = static_log_channel
client.sinks = kafka_sink

#LOG_TO_HDFS_ONLINE_1

client.sources.static_log_source.type = spooldir
client.sources.static_log_source.spoolDir = 监控目录
client.sources.static_log_source.fileSuffix = .COMPLETED
client.sources.static_log_source.ignorePattern = ^$
client.sources.static_log_source.trackerDir = 传输过程中元数据存储路径
client.sources.static_log_source.maxBlobLength = 16384
client.sources.static_log_source.batchSize = 51200
client.sources.static_log_source.inputCharset = UTF-8
client.sources.static_log_source.deserializer = LINE
client.sources.static_log_source.selector.type = replicating
client.sources.static_log_source.fileHeaderKey = file
client.sources.static_log_source.fileHeader = false
client.sources.static_log_source.basenameHeader = true
client.sources.static_log_source.basenameHeaderKey = basename
client.sources.static_log_source.deletePolicy = never

client.channels.static_log_channel.type = file
client.channels.static_log_channel.dataDirs = 数据缓存路径, 设置多个路径可提升性能, 中间用逗号分开
client.channels.static_log_channel.checkpointDir = 检查点存放路径
client.channels.static_log_channel.maxFileSize = 2146435071
client.channels.static_log_channel.capacity = 1000000
client.channels.static_log_channel.transactionCapacity = 612000
client.channels.static_log_channel.minimumRequiredSpace = 524288000

client.sinks.kafka_sink.type = org.apache.flume.sink.kafka.KafkaSink
client.sinks.kafka_sink.kafka.topic = 数据写入的topic, 如flume_test
client.sinks.kafka_sink.kafka.bootstrap.servers = XXX.XXX.XXX.XXX:kafka端口号,XXX.XXX.XXX.XXX:kafka
端口号,XXX.XXX.XXX.XXX:kafka端口号
client.sinks.kafka_sink.flumeBatchSize = 1000
client.sinks.kafka_sink.kafka.producer.type = sync
client.sinks.kafka_sink.kafka.security.protocol = SASL_PLAINTEXT
client.sinks.kafka_sink.kafka.kerberos.domain.name = Kafka Domain名称, 安全集群必填
client.sinks.kafka_sink.requiredAcks = 0

client.sources.static_log_source.channels = static_log_channel
client.sinks.kafka_sink.channel = static_log_channel
```

### 📖 说明

- client.sinks.kafka\_sink.kafka.topic: 数据写入的topic。若kafka中该topic不存在, 默认情况下会自动创建该topic。
- client.sinks.kafka\_sink.kafka.bootstrap.servers: Kafkabrokers列表, 多个用英文逗号分隔。默认情况下, 安全集群端口21007, 普通集群对应端口9092。
- client.sinks.kafka\_sink.kafka.security.protocol: 安全集群为SASL\_PLAINTEXT, 普通集群为PLAINTEXT。
- client.sinks.kafka\_sink.kafka.kerberos.domain.name:  
普通集群无需配置此参数。安全集群对应此参数的值为Kafka集群中“kerberos.domain.name”对应的值。  
具体可到Broker实例所在节点上查看\${BIGDATA\_HOME}/MRS\_Current/1\_X\_Broker/etc/server.properties。  
其中X为随机生成的数字, 请根据实际情况修改。同时文件需要以Flume客户端安装用户身份保存, 例如root用户。

**步骤10** 参数配置并保存后，Flume客户端将自动加载“properties.properties”中配置的内容。当spoolDir生成新的日志文件，文件内容将发送到Kafka生产者，并支持Kafka消费者消费。可参考[管理Kafka主题中的消息](#)查看数据消费情况。

----结束

## 7.5 配置 Flume 非加密传输数据采集任务

### 7.5.1 生成 Flume 服务端和客户端的配置文件

#### 操作场景

该操作指导安装工程师在集群及Flume服务安装完成后，分别配置Flume服务的服务端和客户端参数，使其可以正常工作。

本章节适用于MRS 3.x及之后版本。

#### 📖 说明

本配置默认集群网络环境是安全的，数据传输过程不需要启用SSL认证。如需使用加密方式，请参考[配置Flume加密传输数据采集任务](#)。

#### 前提条件

- 已安装Flume客户端。
- 已成功安装集群及Flume服务。
- 确保集群网络环境安全。

#### 操作步骤

**步骤1** 配置Flume角色客户端参数。

1. 使用FusionInsight Manager界面中的Flume配置工具来配置Flume角色客户端参数并生成配置文件。
  - a. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

图 7-3 选择配置工具



- b. “Agent名”选择“client”，然后选择要使用的Source、Channel以及Sink，将其拖到右侧的操作界面中并将其连接。

例如采用SpoolDir Source、File Channel和Avro Sink，如[图7-4](#)所示。



图 7-4 Flume 配置工具示例



- c. 双击对应的Source、Channel以及Sink，根据实际环境并参考表7-8设置对应的配置参数。

**说明**

- 如果对应的Flume角色之前已经配置过客户端参数，为保证与之前的配置保持一致，可以到“客户端安装目录/fusioninsight-flume-1.9.0/conf/properties.properties”获取已有的客户端参数配置文件。然后登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置Source/Channel/Sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。

表 7-8 Flume 角色客户端所需修改的参数列表

| 参数名称 | 参数值填写规则                                                                                                                                    | 参数样例  |
|------|--------------------------------------------------------------------------------------------------------------------------------------------|-------|
| ssl  | 是否启用SSL认证（基于安全要求，建议启用此功能）<br>只有“Avro”类型的Source才有此配置项<br><ul style="list-style-type: none"> <li>▪ true表示启用</li> <li>▪ false表示不启用</li> </ul> | false |

- d. 单击“导出”，将配置文件“properties.properties”保存到本地。
2. 将“properties.properties”文件上传到Flume客户端安装目录下的“flume/conf/”下。

**步骤2** 配置Flume角色的服务端参数，并将配置文件上传到集群。

1. 使用FusionInsight Manager界面中的Flume配置工具来配置服务端参数并生成配置文件。
  - a. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

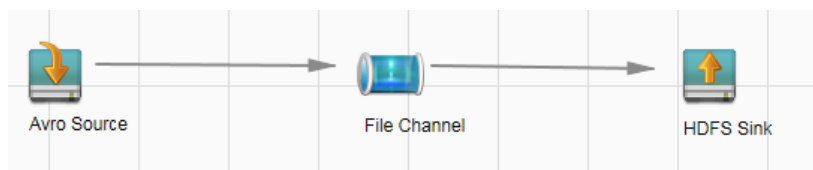
图 7-5 选择配置工具



- b. “Agent名”选择“server”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

例如采用Avro Source、File Channel和HDFS Sink，如图7-6所示。

图 7-6 Flume 配置工具示例



- c. 双击对应的source、channel以及sink，根据实际环境并参考表7-9设置对应的配置参数。

**说明**

- 如果对应的Flume角色之前已经配置过服务端参数，为保证与之前的配置保持一致，在 FusionInsight Manager界面选择“集群 > 服务 > Flume > 实例”，选择相应的Flume角色实例，单击“实例配置”页面“flume.config.file”参数后的“下载文件”，可获取已有的服务端参数配置文件。然后选择“集群 > 服务 > Flume > 配置 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。
- 不同的File Channel均需要配置一个不同的checkpoint目录。

表 7-9 Flume 角色服务端所需修改的参数列表

| 参数名称 | 参数值填写规则                                                                                                                                 | 参数样例  |
|------|-----------------------------------------------------------------------------------------------------------------------------------------|-------|
| ssl  | 是否启用SSL认证（基于安全要求，建议启用此功能）<br>只有“Avro”类型的Source才有此配置项 <ul style="list-style-type: none"> <li>▪ true表示启用</li> <li>▪ false表示不启用</li> </ul> | false |

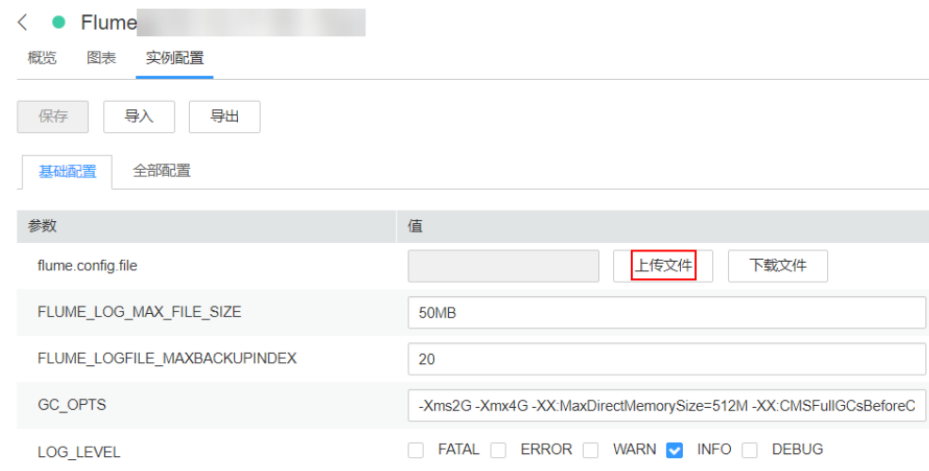
- d. 单击“导出”，将配置文件“properties.properties”保存到本地。
2. 登录FusionInsight Manager，选择“集群 > 服务 > Flume”，在“实例”下单击“Flume”角色。

图 7-7 单击 Flume 角色



3. 选择准备上传配置文件的节点行的“Flume”角色，单击“实例配置”页面“flume.config.file”参数后的“上传文件”，选择“properties.properties”文件完成操作。

图 7-8 上传文件



### 说明

- 每个Flume实例均可以上传单独的服务端配置文件。
  - 更新配置文件需要按照此步骤操作，后台修改配置文件是不规范操作，同步配置时后台做的修改将会被覆盖。
4. 单击“保存”，单击“确定”。

5. 单击“完成”完成操作。

----结束

## 7.5.2 使用 Flume 服务端从本地采集静态日志保存到 Kafka

### 操作场景

该任务指导用户使用Flume服务端从本地采集静态日志保存到Kafka的Topic列表（test1）。

本章节适用于MRS 3.x及之后版本。

#### 说明

本配置默认集群网络环境是安全的，数据传输过程不需要启用SSL认证。如需使用加密方式，请参考[配置Flume加密传输数据采集任务](#)。该配置为只用一个Flume场景，例如：SpoolDir Source +Memory Channel+Kafka Sink。

### 前提条件

- 已成功安装集群，包含Kafka及Flume服务。
- 确保集群网络环境安全。
- MRS集群管理员已明确业务需求，并准备一个Kafka管理员用户flume\_kafka。

### 操作步骤

#### 步骤1 配置Flume的参数。

使用Manager界面中的Flume配置工具来配置Flume角色服务端参数并生成配置文件。

1. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

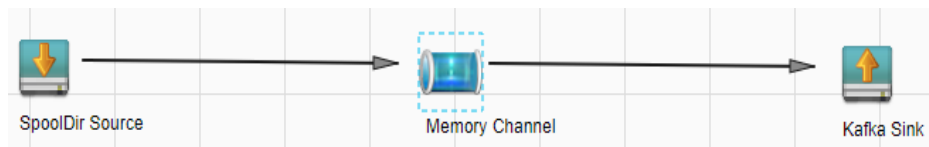
图 7-9 选择配置工具



2. “Agent名”选择“server”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

采用SpoolDir Source、Memory Channel和Kafka Sink，如[图7-10](#)所示。

图 7-10 Flume 配置工具示例



3. 双击对应的source、channel以及sink，根据实际环境并参考[表7-10](#)设置对应的配置参数。

 说明

- 如果想在之前的“properties.properties”文件上进行修改后继续使用，则登录 Manager，选择“集群 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置Source/Channel/Sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。

表 7-10 Flume 角色服务端所需修改的参数列表

| 参数名称                    | 参数值填写规则                                            | 参数样例                              |
|-------------------------|----------------------------------------------------|-----------------------------------|
| 名称                      | 不能为空，必须唯一                                          | test                              |
| spoolDir                | 待采集的文件所在的目录路径，此参数不能为空。该路径需存在，且对flume运行用户有读写执行权限。   | /srv/BigData/hadoop/data1/zb      |
| trackerDir              | flume采集文件信息元数据保存路径。                                | /srv/BigData/hadoop/data1/tracker |
| batchSize               | Flume一次发送的事件个数（数据条数）。增大会提升性能，降低实时性；反之降低性能，提升实时性。   | 61200                             |
| kafka.topics            | 订阅的Kafka topic列表，多个topic用逗号分隔，此参数不能为空。             | test1                             |
| kafka.bootstrap.servers | Kafka的bootstrap地址端口列表，默认值为Kafka集群中所有的Kafkabrokers。 | 192.168.101.10:21007              |

4. 单击“导出”，将配置文件“properties.properties”保存到本地。

**步骤2** 上传配置文件。

登录FusionInsight Manager，选择“集群 > 服务 > Flume”，在“实例”下单击准备上传配置文件的“Flume”角色，单击“实例配置”页面“flume.config.file”参数后的“上传文件”，选择**步骤1.4**导出的“properties.properties”文件完成操作。

**步骤3** 验证日志是否传输成功。

1. 登录Kafka客户端：

```
cd Kafka客户端安装目录/Kafka/kafka
```

```
kinit flume_kafka (输入密码)
```

2. 读取KafkaTopic中的数据（修改命令中的中文为实际参数）。

```
bin/kafka-console-consumer.sh --topic 主题名称 --bootstrap-server Kafka角色实例所在节点的业务IP地址:21007 --consumer.config config/consumer.properties --from-beginning
```

系统显示待采集文件目录下的内容：

```
[root@host1 kafka]# bin/kafka-console-consumer.sh --topic test1 --bootstrap-server
192.168.101.10:21007 --consumer.config config/consumer.properties --from-beginning
Welcome to flume
```

----结束

## 7.5.3 使用 Flume 服务端从本地采集静态日志保存到 HDFS

### 操作场景

该任务指导用户使用Flume服务端从本地采集静态日志保存到HDFS上“/flume/test”目录下。

本章节适用于MRS 3.x及之后版本。

#### 说明

本配置默认集群网络环境是安全的，数据传输过程不需要启用SSL认证。如需使用加密方式，请参考[配置Flume加密传输数据采集任务](#)。该配置为只用一个Flume场景，例如：Spooldir Source +Memory Channel+HDFS Sink。

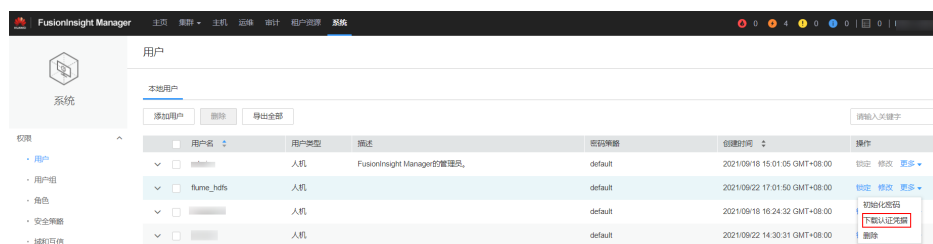
### 前提条件

- 已成功安装集群，包含HDFS及Flume服务。
- 确保集群网络环境安全。
- 已创建用户flume\_hdfs并授权验证日志时操作的HDFS目录和数据。

### 操作步骤

**步骤1** 在FusionInsight Manager管理界面，选择“系统 > 权限 > 用户”，选择用户flume\_hdfs，选择“更多 > 下载认证凭据”下载Kerberos证书文件并保存在本地。

图 7-11 下载认证凭据



**步骤2** 配置Flume参数。

使用FusionInsight Manager界面中的Flume来配置Flume角色服务端参数并生成配置文件。

1. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

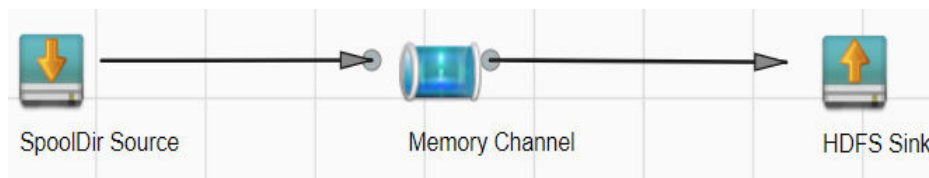
图 7-12 选择配置工具



- “Agent名”选择“server”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

采用SpoolDir Source、Memory Channel和HDFS Sink，如图7-13所示。

图 7-13 Flume 配置工具示例



- 双击对应的source、channel以及sink，根据实际环境并参考表7-11设置对应的配置参数。

#### 说明

- 如果想在之前的“properties.propertites”文件上进行修改后继续使用，则登录 FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。

表 7-11 Flume 角色服务端所需修改的参数列表

| 参数名称                   | 参数值填写规则                                          | 参数样例                              |
|------------------------|--------------------------------------------------|-----------------------------------|
| 名称                     | 不能为空，必须唯一。                                       | test                              |
| spoolDir               | 待采集的文件所在的目录路径，此参数不能为空。该路径需存在，且对flume运行用户有读写执行权限。 | /srv/BigData/hadoop/data1/zb      |
| trackerDir             | flume采集文件信息元数据保存路径。                              | /srv/BigData/hadoop/data1/tracker |
| batchSize              | Flume一次发送数据的最大事件数。                               | 61200                             |
| hdfs.path              | 写入HDFS的目录，此参数不能为空。                               | hdfs://hacluster/flume/test       |
| hdfs.filePrefix        | 数据写入HDFS后文件名的前缀。                                 | TMP_                              |
| hdfs.batchSize         | 一次写入HDFS的最大事件数目。                                 | 61200                             |
| hdfs.kerberosPrincipal | kerberos认证时用户，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。  | flume_hdfs                        |

| 参数名称                   | 参数值填写规则                                                 | 参数样例                                                                                                                                |
|------------------------|---------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------|
| hdfs.kerberosKeytab    | kerberos认证时keytab文件路径，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。 | /opt/test/conf/user.keytab<br><b>说明</b><br>user.keytab文件从下载用户flume_hdfs的kerberos证书文件中获取，另外，确保用于安装和运行Flume客户端的用户对user.keytab文件有读写权限。 |
| hdfs.useLocalTimeStamp | 是否使用本地时间，取值为"true"或者"false"。                            | true                                                                                                                                |

4. 单击“导出”，将配置文件“properties.properties”保存到本地。

### 步骤3 上传配置文件。

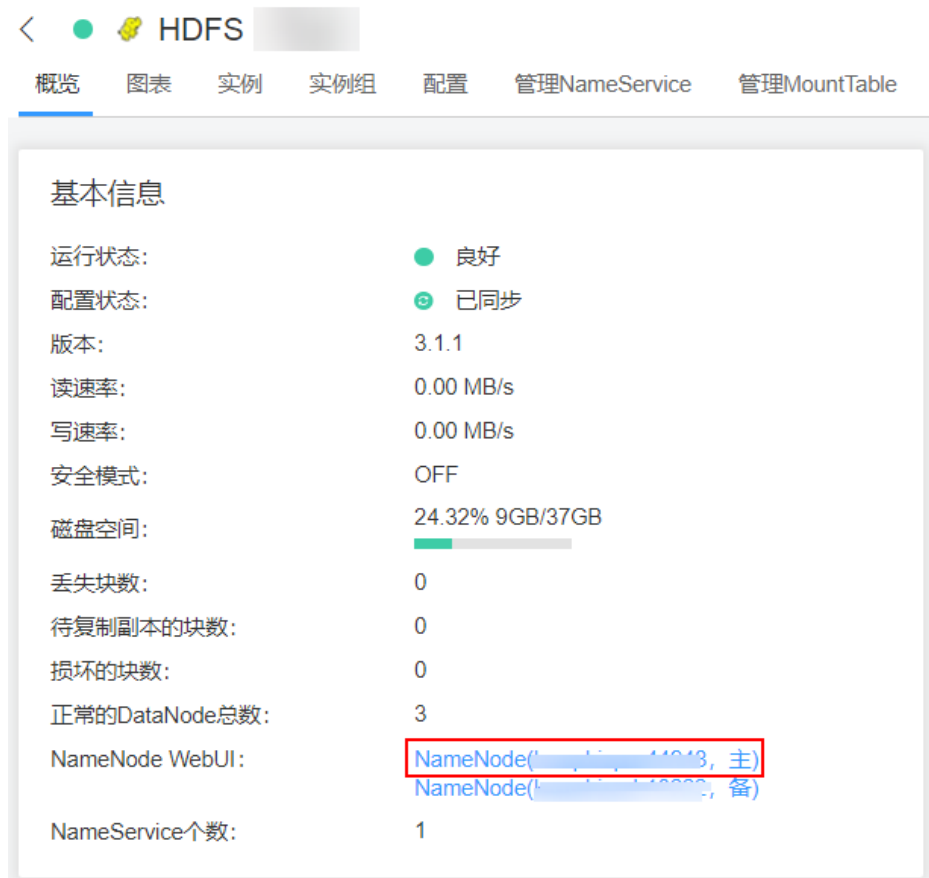
登录FusionInsight Manager，选择“集群 > 服务 > Flume”，在“实例”下单击准备上传配置文件的“Flume”角色，单击“实例配置”页面“flume.config.file”参数后的“上传文件”，选择[步骤2.4](#)导出的“properties.properties”文件完成操作。

### 步骤4 验证日志是否传输成功。

1. 以具有HDFS组件管理权限的用户登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。在FusionInsight Manager界面选择“集群 > 服务 > HDFS”，单击“NameNode(主)”对应的链接，打开HDFS WebUI，然后选择“Utilities > Browse the file system”。

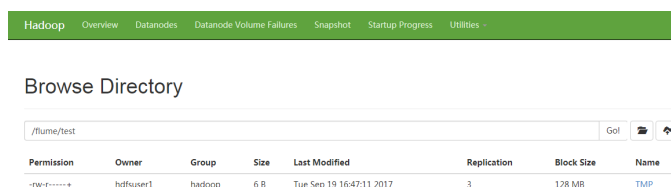


图 7-14 进入 HDFS WebUI



2. 观察HDFS上“/flume/test”目录下是否有产生数据。

图 7-15 查看 HDFS 目录和文件



---结束

## 7.5.4 使用 Flume 服务端从本地采集动态日志保存到 HDFS

### 操作场景

该任务指导用户使用Flume服务端从本地采集动态日志保存到HDFS上“/flume/test”目录下。

本章节适用于MRS 3.x及之后版本。

## 说明

本配置默认集群网络环境是安全的，数据传输过程不需要启用SSL认证。如需使用加密方式，请参考[配置Flume加密传输数据采集任务](#)。该配置为只用一个Flume场景，例如：Taildir Source +Memory Channel+HDFS Sink。

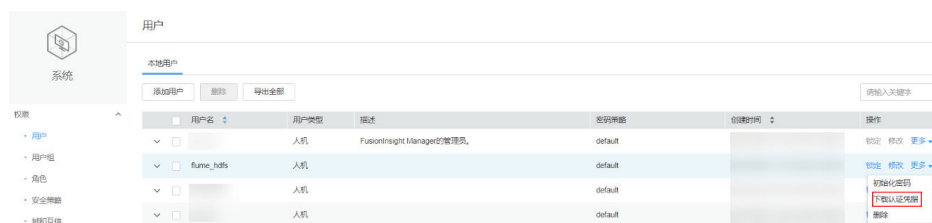
## 前提条件

- 已成功安装集群，包含HDFS及Flume服务。
- 确保集群网络环境安全。
- 已创建用户flume\_hdfs并授权验证日志时操作的HDFS目录和数据。

## 操作步骤

**步骤1** 在FusionInsight Manager管理界面，选择“系统 > 权限 > 用户”，选择“更多 > 下载认证凭据”下载用户flume\_hdfs的kerberos证书文件并保存在本地。

图 7-16 下载认证凭据



**步骤2** 配置Flume参数。

使用FusionInsight Manager界面中的Flume配置工具来配置Flume角色服务端参数并生成配置文件。

1. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

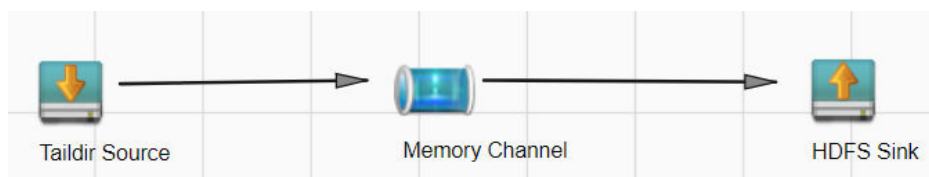
图 7-17 选择配置工具



2. “Agent名”选择“server”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

采用Taildir Source、Memory Channel和HDFS Sink，如[图7-18](#)所示。

图 7-18 Flume 配置工具示例



3. 双击对应的Source、Channel以及Sink，根据实际环境并参考[表7-12](#)设置对应的配置参数。

 说明

- 如果在之前的“properties.propretites”文件上进行修改后继续使用，则登录 FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置Source/Channel/Sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。

表 7-12 Flume 角色服务端所需修改的参数列表

| 参数名称                   | 参数值填写规则                                                                  | 参数样例                                                                                                                                |
|------------------------|--------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------|
| 名称                     | 不能为空，必须唯一。                                                               | test                                                                                                                                |
| filegroups             | 文件分组列表名，此参数不能为空。该值包含如下两项参数：<br>- 名称：文件分组列表名。<br>- filegroups：动态日志文件绝对路径。 | -                                                                                                                                   |
| positionFile           | 保存当前采集文件信息（文件名和已经采集的位置），此参数不能为空。该文件不需要手工创建，但其上层目录需对flume运行用户可写。          | /home/omm/flume/positionfile                                                                                                        |
| batchSize              | Flume一次发送数据的最大事件数。                                                       | 61200                                                                                                                               |
| hdfs.path              | 写入HDFS的目录，此参数不能为空。                                                       | hdfs://hacluster/flume/test                                                                                                         |
| hdfs.filePrefix        | 数据写入HDFS后文件名的前缀。                                                         | TMP_                                                                                                                                |
| hdfs.batchSize         | 一次写入HDFS的最大事件数目。                                                         | 61200                                                                                                                               |
| hdfs.kerberosPrincipal | kerberos认证时用户，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。                          | flume_hdfs                                                                                                                          |
| hdfs.kerberosKeytab    | kerberos认证时keytab文件路径，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。                  | /opt/test/conf/user.keytab<br><b>说明</b><br>user.keytab文件从下载用户flume_hdfs的kerberos证书文件中获取，另外，确保用于安装和运行Flume客户端的用户对user.keytab文件有读写权限。 |
| hdfs.useLocalTimeStamp | 是否使用本地时间，取值为"true"或者"false"。                                             | true                                                                                                                                |

- 单击“导出”，将配置文件“properties.properties”保存到本地。

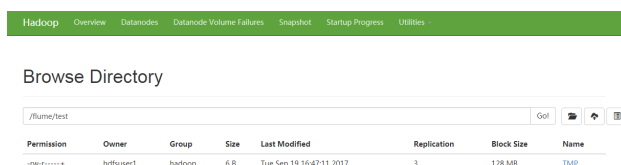
### 步骤3 上传配置文件。

登录FusionInsight Manager，选择“集群 > 服务 > Flume”，在“实例”下单击准备上传配置文件的“Flume”角色，单击“实例配置”页面“flume.config.file”参数后的“上传文件”，选择[步骤2.4](#)导出的“properties.properties”文件完成操作。

### 步骤4 验证日志是否传输成功。

- 以具有HDFS组件管理权限的用户登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。在FusionInsight Manager界面选择“集群 > 服务 > HDFS”，单击“NameNode(节点名称, 主)”对应的链接，打开HDFS WebUI，然后选择“Utilities > Browse the file system”。
- 观察HDFS上“/flume/test”目录下是否有产生数据。

图 7-19 查看 HDFS 目录和文件



----结束

## 7.5.5 使用 Flume 服务端从 Kafka 采集日志保存到 HDFS

### 操作场景

该任务指导用户使用Flume服务端从Kafka的Topic列表(test1)采集日志保存到HDFS上“/flume/test”目录下。

本章节适用于MRS 3.x及之后版本。

### 说明

本配置默认集群网络环境是安全的，数据传输过程不需要启用SSL认证。如需使用加密方式，请参考[配置Flume加密传输数据采集任务](#)。该配置为只用一个Flume场景，例如：Kafka Source +Memory Channel+HDFS Sink。

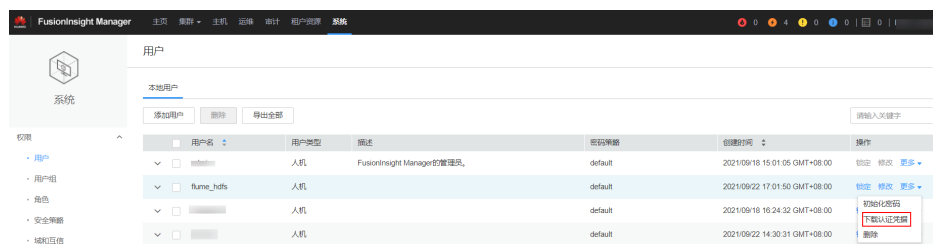
### 前提条件

- 已成功安装集群，包含HDFS、Kafka及Flume服务。
- 确保集群网络环境安全。
- 已创建用户flume\_hdfs并授权验证日志时操作的HDFS目录和数据。

### 操作步骤

- 步骤1** 在FusionInsight Manager管理界面，选择“系统 > 权限 > 用户”，选择“更多 > 下载认证凭据”下载用户flume\_hdfs的kerberos证书文件并保存在本地。

图 7-20 下载认证凭据



**步骤2 配置Flume角色服务端参数。**

使用FusionInsight Manager界面中的Flume配置工具来配置Flume角色服务端参数并生成配置文件。

1. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

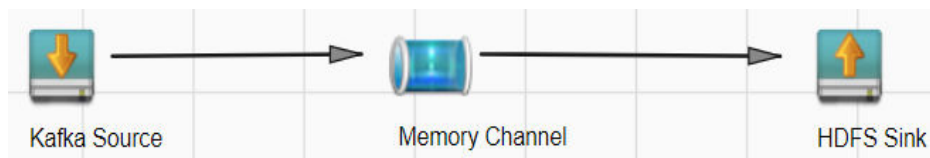
图 7-21 选择配置工具



2. “Agent名”选择“server”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

例如采用Kafka Source、Memory Channel和HDFS Sink，如图7-22所示。

图 7-22 Flume 配置工具示例



3. 双击对应的source、channel以及sink，根据实际环境并参考表7-13设置对应的配置参数。

**说明**

- 如果想在之前的“properties.propertites”文件上进行修改后继续使用，则登录 FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。

表 7-13 Flume 角色服务端所需修改的参数列表

| 参数名称         | 参数值填写规则                         | 参数样例  |
|--------------|---------------------------------|-------|
| 名称           | 不能为空，必须唯一。                      | test  |
| kafka.topics | 订阅的Kafka topic列表，用逗号分隔，此参数不能为空。 | test1 |

| 参数名称                    | 参数值填写规则                                                                    | 参数样例                                                                                                                                |
|-------------------------|----------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------|
| kafka.consumer.group.id | 从Kafka中获取数据的组标识，此参数不能为空。                                                   | flume                                                                                                                               |
| kafka.bootstrap.servers | Kafka的bootstrap地址端口列表,默认值为Kafka集群中所有的Kafka列表。如果集群安装有Kafka并且配置已经同步，可以不配置此项。 | 192.168.101.10:9092                                                                                                                 |
| batchSize               | Flume一次发送的事件个数（数据条数）。                                                      | 61200                                                                                                                               |
| hdfs.path               | 写入HDFS的目录，此参数不能为空。                                                         | hdfs://hacluster/flume/test                                                                                                         |
| hdfs.filePrefix         | 数据写入HDFS后文件名的前缀。                                                           | TMP_                                                                                                                                |
| hdfs.batchSize          | 一次写入HDFS的最大事件数目。                                                           | 61200                                                                                                                               |
| hdfs.kerberosPrincipal  | kerberos认证时用户，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。                            | flume_hdfs                                                                                                                          |
| hdfs.kerberosKeytab     | kerberos认证时keytab文件路径，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。                    | /opt/test/conf/user.keytab<br><b>说明</b><br>user.keytab文件从下载用户flume_hdfs的kerberos证书文件中获取，另外，确保用于安装和运行Flume客户端的用户对user.keytab文件有读写权限。 |
| hdfs.useLocalTimeStamp  | 是否使用本地时间，取值为"true"或者"false"。                                               | true                                                                                                                                |

4. 单击“导出”，将配置文件“properties.properties”保存到本地。

### 步骤3 上传配置文件。

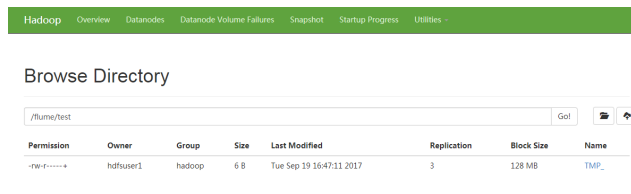
登录FusionInsight Manager，选择“集群 > 服务 > Flume”，在“实例”下单击准备上传配置文件的“Flume”角色，单击“实例配置”页面“flume.config.file”参数后的“上传文件”，选择[步骤2.4](#)导出的“properties.properties”文件完成操作。

### 步骤4 验证日志是否传输成功。

1. 以具有HDFS组件管理权限的用户登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。在FusionInsight Manager界面选择“集群 > 服务 > HDFS”，单击“NameNode(节点名称, 主)”对应的链接，打开HDFS WebUI，然后选择“Utilities > Browse the file system”。

2. 观察HDFS上“/flume/test”目录下是否有产生数据。

图 7-23 查看 HDFS 目录和文件



----结束

## 7.5.6 使用 Flume 客户端从 Kafka 采集日志保存到 HDFS

### 操作场景

该任务指导用户使用Flume客户端从Kafka客户端的Topic列表(test1)采集日志保存到HDFS上“/flume/test”目录下。

本章节适用于MRS 3.x及之后版本。

### 说明

本配置默认集群网络环境是安全的，数据传输过程不需要启用SSL认证。如需使用加密方式，请参考[配置Flume加密传输数据采集任务](#)。

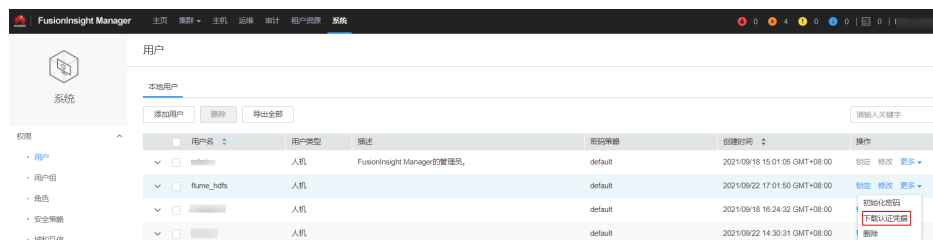
### 前提条件

- 已安装Flume客户端。
- 已成功安装集群，包含HDFS、Kafka及Flume服务。
- 已创建用户flume\_hdfs并授权验证日志时操作的HDFS目录和数据。
- 确保集群网络环境安全。

### 操作步骤

- 步骤1** 在FusionInsight Manager管理界面，选择“系统 > 权限 > 用户”，选择“更多 > 下载认证凭据”下载用户flume\_hdfs的kerberos证书文件并保存在本地。

图 7-24 下载认证凭据



- 步骤2** 配置Flume角色客户端参数。

1. 使用FusionInsight Manager界面中的Flume配置工具来配置Flume角色服务端参数并生成配置文件。

- a. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

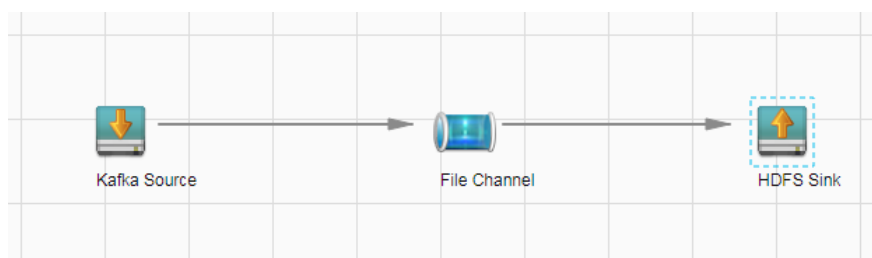
图 7-25 选择配置工具



- b. “Agent名”选择“client”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

例如采用Kafka Source、File Channel和HDFS Sink，如图7-26所示。

图 7-26 Flume 配置工具示例



- c. 双击对应的source、channel以及sink，根据实际环境并参考表7-14设置对应的配置参数。

#### 说明

- 如果想在之前的“properties.propretites”文件上进行修改后继续使用，则登录 FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。

表 7-14 Flume 角色客户端所需修改的参数列表

| 参数名称                    | 参数值填写规则                         | 参数样例  |
|-------------------------|---------------------------------|-------|
| 名称                      | 不能为空，必须唯一。                      | test  |
| kafka.topics            | 订阅的Kafka topic列表，用逗号分隔，此参数不能为空。 | test1 |
| kafka.consumer.group.id | 从Kafka中获取数据的组标识，此参数不能为空。        | flume |



| 参数名称                    | 参数值填写规则                                                                                                                                 | 参数样例                                       |
|-------------------------|-----------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------|
| kafka.bootstrap.servers | Kafka的bootstrap地址端口列表,默认值为Kafka集群中所有的Kafka列表。如果集群安装有Kafka并且配置已经同步,可以不配置此项。当使用Flume客户端时,必须配置此项。                                          | 192.168.101.10:21007                       |
| batchSize               | Flume一次发送的事件个数(数据条数)。                                                                                                                   | 61200                                      |
| dataDirs                | 缓冲区数据保存目录,默认为运行目录。配置多个盘上的目录可以提升传输效率,多个目录使用逗号分隔。如果为集群内,则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/data, dataX为data1~dataN。如果为集群外,则需要单独规划。 | /srv/BigData/hadoop/data1/flume/data       |
| checkpointDir           | checkpoint信息保存目录,默认在运行目录下。如果为集群内,则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/checkpoint, dataX为data1~dataN。如果为集群外,则需要单独规划。                | /srv/BigData/hadoop/data1/flume/checkpoint |
| transactionCapacity     | 事务大小:即当前channel支持事务处理的事件个数,建议和Source的batchSize设置为同样大小,不能小于batchSize。                                                                    | 61200                                      |
| hdfs.path               | 写入HDFS的目录,此参数不能为空。                                                                                                                      | hdfs://hacluster/flume/test                |
| hdfs.filePrefix         | 数据写入HDFS后文件名的前缀。                                                                                                                        | TMP_                                       |
| hdfs.batchSize          | 一次写入HDFS的最大事件数目。                                                                                                                        | 61200                                      |

| 参数名称                   | 参数值填写规则                                                 | 参数样例                                                                                                                                    |
|------------------------|---------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------|
| hdfs.kerberosPrincipal | kerberos认证时用户，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。         | flume_hdfs                                                                                                                              |
| hdfs.kerberosKeytab    | kerberos认证时keytab文件路径，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。 | /opt/test/conf/<br>user.keytab<br><b>说明</b><br>user.keytab文件从下载用户flume_hdfs的kerberos证书文件中获取，另外，确保用于安装和运行Flume客户端的用户对user.keytab文件有读写权限。 |
| hdfs.useLocalTimeStamp | 是否使用本地时间，取值为"true"或者"false"                             | true                                                                                                                                    |

- d. 单击“导出”，将配置文件“properties.properties”保存到本地。
2. 将“properties.properties”文件上传到Flume客户端安装目录下的“flume/conf/”下。
3. Flume客户端连接到HDFS，还需要补充如下配置：
  - a. 通过“用户”下载用户flume\_hdfs的kerberos证书文件获取krb5.conf配置文件，并上传至客户端所在节点安装目录的“fusioninsight-flume-1.9.0/conf/”下。
  - b. 新建jaas.conf配置文件到客户端所在节点安装目录的“fusioninsight-flume-1.9.0/conf/”下。

#### vi jaas.conf

```
KafkaClient {
com.sun.security.auth.module.Krb5LoginModule required
useKeyTab=true
keyTab="/opt/test/conf/user.keytab"
principal="flume_hdfs@<系统域名>"
useTicketCache=false
storeKey=true
debug=true;
};
```

参数keyTab和principal根据实际情况修改。

- c. 从/opt/FusionInsight\_Cluster\_<集群ID>\_Flume\_ClientConfig/Flume/config目录下获取core-site.xml和hdfs-site.xml配置文件，并上传至客户端所在节点安装目录的“fusioninsight-flume-1.9.0/conf/”下。
4. 进入客户端所在节点安装目录的“fusioninsight-flume-1.9.0/bin”下，执行以下命令重启Flume进程。

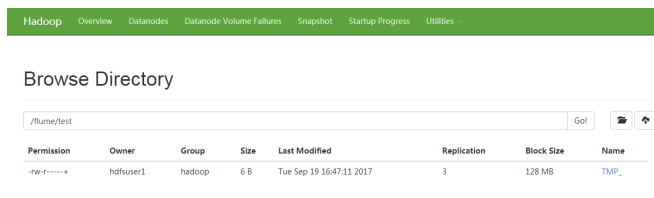
**./flume-manage.sh restart**

#### 步骤3 验证日志是否传输成功。

1. 以具有HDFS组件管理权限的用户登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。在FusionInsight Manager界

- 面选择“集群 > 服务 > HDFS”，单击“NameNode(节点名称, 主)”对应的链接，打开HDFS WebUI，然后选择“Utilities > Browse the file system”。
2. 观察HDFS上“/flume/test”目录下是否有产生数据。

图 7-27 查看 HDFS 目录和文件



----结束

## 7.5.7 使用多级 agent 串联从本地采集静态日志保存到 HBase

### 操作场景

该任务指导用户使用Flume客户端从本地采集静态日志保存到HBase表：flume\_test。该场景介绍的是多级agent串联操作。

本章节适用于MRS 3.x及之后版本。

#### 说明

本配置默认集群网络环境是安全的，数据传输过程不需要启用SSL认证。如需使用加密方式，请参考[配置Flume加密传输数据采集任务](#)。该配置可以只用一个Flume场景，例如Server: Spooldir Source+File Channel+HBase Sink。

### 前提条件

- 已成功安装集群，包含HBase及Flume服务。
- 已安装Flume客户端。
- 确保集群网络环境安全。
- 已创建HBase表：**create 'flume\_test', 'cf'**。
- MRS集群管理员已明确业务需求，并准备一个HBase管理员用户**flume\_hbase**。

### 操作步骤

- 步骤1** 在FusionInsight Manager管理界面，选择“系统 > 权限 > 用户”，选择“更多 > 下载认证凭据”下载用户**flume\_hbase**的kerberos证书文件并保存在本地。

图 7-28 下载认证凭据



**步骤2** 配置Flume角色客户端参数。

1. 使用FusionInsight Manager界面中的Flume配置工具来配置Flume角色客户端参数并生成配置文件。
  - a. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

**图 7-29** 选择配置工具



- b. “Agent名”选择“client”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。  
采用SpoolDir Source、File Channel和Avro Sink，如图7-30所示。

**图 7-30** Flume 配置工具示例



- c. 双击对应的source、channel以及sink，根据实际环境并参考表7-15设置对应的配置参数。

**说明**

- 如果想在之前的“properties.propretites”文件上进行修改后继续使用，则登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。

**表 7-15** Flume 角色客户端所需修改的参数列表

| 参数名称       | 参数值填写规则                                          | 参数样例                              |
|------------|--------------------------------------------------|-----------------------------------|
| 名称         | 不能为空，必须唯一。                                       | test                              |
| spoolDir   | 待采集的文件所在的目录路径，此参数不能为空。该路径需存在，且对flume运行用户有读写执行权限。 | /srv/BigData/hadoop/data1/zb      |
| trackerDir | flume采集文件信息元数据保存路径。                              | /srv/BigData/hadoop/data1/tracker |

| 参数名称                | 参数值填写规则                                                                                                                                | 参数样例                                       |
|---------------------|----------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------|
| batchSize           | Flume一次发送的事件个数（数据条数）。增大会提升性能，降低实时性；反之降低性能，提升实时性。                                                                                       | 61200                                      |
| dataDirs            | 缓冲区数据保存目录，默认为运行目录。配置多个盘上的目录可以提升传输效率，多个目录使用逗号分隔。如果为集群内，则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/data，dataX为data1~dataN。如果为集群外，则需要单独规划。 | /srv/BigData/hadoop/data1/flume/data       |
| checkpointDir       | checkpoint信息保存目录，默认在运行目录下。如果为集群内，则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/checkpoint，dataX为data1~dataN。如果为集群外，则需要单独规划。                | /srv/BigData/hadoop/data1/flume/checkpoint |
| transactionCapacity | 事务大小：即当前channel支持事务处理的事件个数，建议和Source的batchSize设置为同样大小，不能小于batchSize。                                                                   | 61200                                      |
| hostname            | 要发送数据的主机名或者IP，此参数不能为空。须配置为与之相连的avro source所在的主机名或IP。                                                                                   | 192.168.108.11                             |
| port                | 要发送数据的端口，此参数不能为空。须配置为与之相连的avro source监测的端口。                                                                                            | 21154                                      |

| 参数名称 | 参数值填写规则                                                                                                                                   | 参数样例  |
|------|-------------------------------------------------------------------------------------------------------------------------------------------|-------|
| ssl  | 是否启用SSL认证（基于安全要求，建议启用此功能）。<br>只有“Avro”类型的Source才有此配置项。 <ul style="list-style-type: none"> <li>▪ true表示启用</li> <li>▪ false表示不启用</li> </ul> | false |

- d. 单击“导出”，将配置文件“properties.properties”保存到本地。
2. 将“properties.properties”文件上传到Flume客户端安装目录下的“flume/conf/”下。

**步骤3** 配置Flume角色的服务端参数，并将配置文件上传到集群。

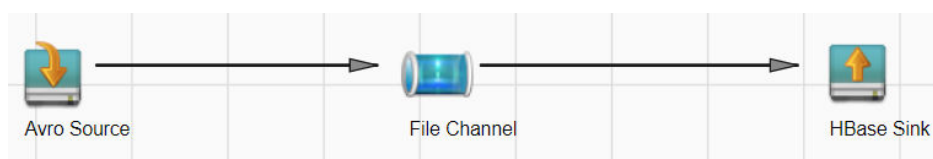
1. 使用FusionInsight Manager界面中的Flume配置工具来配置服务端参数并生成配置文件。
  - a. 登录FusionInsight Manager，选择“集群 > 服务 > Flume > 配置工具”。

**图 7-31** 选择配置工具



- b. “Agent名”选择“server”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。  
 采用Avro Source、File Channel和HBase Sink，如图7-32所示。

**图 7-32** Flume 配置工具示例



- c. 双击对应的source、channel以及sink，根据实际环境并参考表7-16设置对应的配置参数。

 说明

- 如果对应的Flume角色之前已经配置过服务端参数，为保证与之前的配置保持一致，在FusionInsight Manager界面选择“集群 > 待操作集群的名称 > 服务 > Flume > 实例”，选择相应的Flume角色实例，单击“实例配置”页面“flume.config.file”参数后的“下载文件”，可获取已有的服务端参数配置文件。然后选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改非加密传输的相关配置项即可。
- 导入配置文件时，建议配置source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。
- 不同的File Channel均需要配置一个不同的checkpoint目录。

表 7-16 Flume 角色服务端所需修改的参数列表

| 参数名称          | 参数值填写规则                                                                                                                                      | 参数样例                                             |
|---------------|----------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------|
| 名称            | 不能为空，必须唯一。                                                                                                                                   | test                                             |
| bind          | avro source绑定的ip地址，此参数不能为空。须配置为服务端配置文件即将要上传的主机IP。                                                                                            | 192.168.108.11                                   |
| port          | avro source监测的端口，此参数不能为空。须配置为未被使用的端口。                                                                                                        | 21154                                            |
| ssl           | 是否启用SSL认证（基于安全要求，建议启用此功能）。<br>只有“Avro”类型的Source才有此配置项。<br><ul style="list-style-type: none"> <li>■ true表示启用</li> <li>■ false表示不启用</li> </ul> | false                                            |
| dataDirs      | 缓冲区数据保存目录，默认为运行目录。配置多个盘上的目录可以提升传输效率，多个目录使用逗号分隔。如果为集群内，则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/data，dataX为data1~dataN。如果为集群外，则需要单独规划。       | /srv/BigData/hadoop/data1/flumeserver/data       |
| checkpointDir | checkpoint 信息保存目录，默认在运行目录下。如果为集群内，则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/checkpoint，dataX为data1~dataN。如果为集群外，则需要单独规划。                     | /srv/BigData/hadoop/data1/flumeserver/checkpoint |

| 参数名称                | 参数值填写规则                                                              | 参数样例                                                                                                                                     |
|---------------------|----------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------|
| transactionCapacity | 事务大小：即当前channel支持事务处理的事件个数。建议和Source的batchSize设置为同样大小，不能小于batchSize。 | 61200                                                                                                                                    |
| table               | HBase表名，此参数不能为空。                                                     | flume_test                                                                                                                               |
| columnFamily        | HBase列族名，此参数不能为空。                                                    | cf                                                                                                                                       |
| batchSize           | Flume一次写入HBase中的最大事件数。                                               | 61200                                                                                                                                    |
| kerberosPrincipal   | kerberos认证时用户,在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。                      | flume_hbase                                                                                                                              |
| kerberosKeytab      | kerberos认证时文件路径，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。                    | /opt/test/conf/<br>user.keytab<br><b>说明</b><br>user.keytab文件从下载用户flume_hbase的kerberos证书文件中获取，另外，确保用于安装和运行Flume客户端的用户对user.keytab文件有读写权限。 |

- d. 单击“导出”，将配置文件“properties.properties”保存到本地。
2. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume”，在“实例”下单击“Flume”角色。

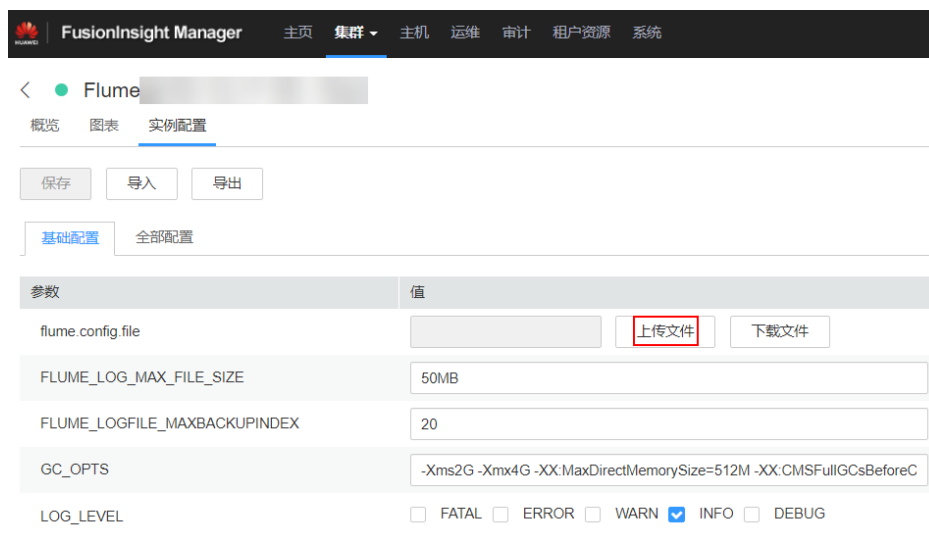


图 7-33 单击 Flume 角色



3. 选择准备上传配置文件的节点行的“Flume”角色，单击“实例配置”页面“flume.config.file”参数后的“上传文件”，选择“properties.properties”文件完成操作。

图 7-34 上传文件



### 说明

- 每个Flume实例均可以上传单独的服务端配置文件。
  - 更新配置文件需要按照此步骤操作，后台修改配置文件是不规范操作，同步配置时后台做的修改将会被覆盖。
4. 单击“保存”，单击“确定”。

5. 单击“完成”完成操作。

#### 步骤4 验证日志是否传输成功。

1. 进入HBase客户端目录：

```
cd /客户端安装目录/HBase/hbase
```

```
kinit flume_hbase (输入密码)
```

2. 执行hbase shell进入HBase客户端。
3. 执行语句：`scan 'flume_test'`，可以看到日志按行写入HBase列族里。

```
hbase(main):001:0> scan 'flume_test'
ROW COLUMN
+CELL

2017-09-18 16:05:36,394 INFO [hconnection-0x415a3f6a-shared--pool2-t1] ipc.AbstractRpcClient:
RPC Server Kerberos principal name for service=ClientService is hbase/hadoop.<系统域名>@<系统域名>
default4021ff4a-9339-4151-a4d0-00f20807e76d column=cf:pCol,
timestamp=1505721909388, value=Welcome to
flume
incRow column=cf:iCol, timestamp=1505721909461, value=
\x00\x00\x00\x00\x00\x00\x00\x01
2 row(s) in 0.3660 seconds
```

----结束

## 7.6 配置 Flume 加密传输数据采集任务

### 7.6.1 配置 Flume 加密传输

#### 操作场景

该操作指导安装工程师在集群安装完成后，设置Flume服务（Flume角色）的服务端和客户端参数，使其可以正常工作。

本章节适用于MRS 3.x及之后版本。

#### 前提条件

已成功安装集群及Flume服务。

#### 操作步骤

- 步骤1 分别生成Flume角色服务端和客户端的证书和信任列表。

1. 使用ECM远程以omm用户登录将要安装Flume服务端的节点。进入“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/bin`”目录。

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/bin
```

#### 📖 说明

此处版本号8.1.0.1为示例，具体以实际环境的版本号为准。

2. 执行以下命令，生成并导出Flume角色服务端、客户端证书。

```
sh geneJKS.sh -f xxx -g xxx
```

生成的证书在“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf`”路径下。其中：

- “`flume_sChat.jks`”是Flume角色服务端的证书库，“`flume_sChat.crt`”是“`flume_sChat.jks`”证书的导出文件，“`-f`”配置项是证书和证书库的密码；
- “`flume_cChat.jks`”是Flume角色客户端的证书库，“`flume_cChat.crt`”是“`flume_cChat.jks`”证书的导出文件，“`-g`”配置项是证书和证书库的密码；
- “`flume_sChatt.jks`”和“`flume_cChatt.jks`”分别为Flume服务端、客户端SSL证书信任列表。

### 📖 说明

本章节涉及到所有的用户自定义密码，需满足以下复杂度要求：

- 至少包含大写字母、小写字母、数字、特殊符号4种类型字符。
- 至少8位，最多64位。
- 出于安全考虑，建议用户定期更换自定义密码（例如三个月更换一次），并重新生成各项证书和信任列表。

**步骤2** 配置Flume角色的服务端参数，并将配置文件上传到集群。

1. 使用ECM远程，以`omm`用户登录任意一个Flume角色所在的节点。执行以下命令进入“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/bin`”。

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/bin
```

2. 执行以下命令，生成并得到Flume服务端密钥库密码、信任列表密码和`keystore-password`加密的私钥信息。连续输入两次密码并确认，该密码是`flume_sChat.jks`证书库的密码。

```
./genPwFile.sh
```

```
cat password.property
```

3. 使用FusionInsight Manager界面中的Flume配置工具来配置服务端参数并生成配置文件。
  - a. 登录FusionInsight Manager，选择“服务 > Flume > 配置工具”。

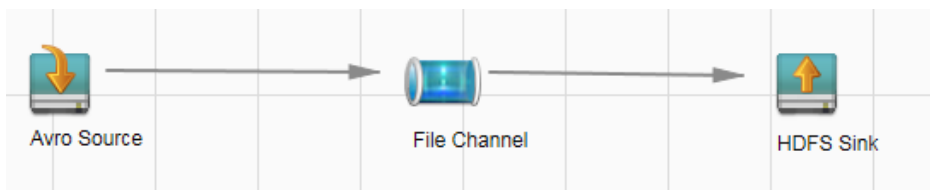
图 7-35 选择配置工具



- b. “Agent名”选择“`server`”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

例如采用Avro Source、File Channel和HDFS Sink，如图7-36所示。

图 7-36 Flume 配置工具示例



- c. 双击对应的source、channel以及sink，根据实际环境并参考表7-17设置对应的配置参数。

**说明**

- 如果对应的Flume角色之前已经配置过服务端参数，为保证与之前的配置保持一致，在FusionInsight Manager界面选择“服务 > Flume > 实例”，选择相应的Flume角色实例，单击“实例配置”页面“flume.config.file”参数后的“下载文件”，可获取已有的服务端参数配置文件。然后选择“服务 > Flume > 导入”，将该文件导入后再修改加密传输的相关配置项即可。
- 导入配置文件时，建议配置Source/Channel/Sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。
- d. 单击“导出”，将配置文件“properties.properties”保存到本地。配置文件中包含认证密码信息可能存在安全风险，建议当前场景执行完毕后删除相关配置文件或加强安全管理。

表 7-17 Flume 角色服务端所需修改的参数列表

| 参数名称              | 参数值填写规则                                                                                                         | 参数样例                                                                                                                     |
|-------------------|-----------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------|
| ssl               | 是否启用SSL认证（基于安全要求，建议启用此功能）。 <ul style="list-style-type: none"> <li>▪ true表示启用。</li> <li>▪ false表示不启用。</li> </ul> | true                                                                                                                     |
| keystore          | 服务端证书。                                                                                                          | <code>\${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf/flume_sChat.jks</code>  |
| keystore-password | 密钥库密码，获取keystore信息所需密码。<br>输入步骤2.2中获取的“password”值。                                                              | -                                                                                                                        |
| truststore        | 服务端的SSL证书信任列表。                                                                                                  | <code>\${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf/flume_sChatt.jks</code> |

| 参数名称                | 参数值填写规则                                                                 | 参数样例 |
|---------------------|-------------------------------------------------------------------------|------|
| truststore-password | 信任列表密码，获取truststore信息所需密码。<br>输入 <a href="#">步骤2.2</a> 中获取的“password”值。 | -    |

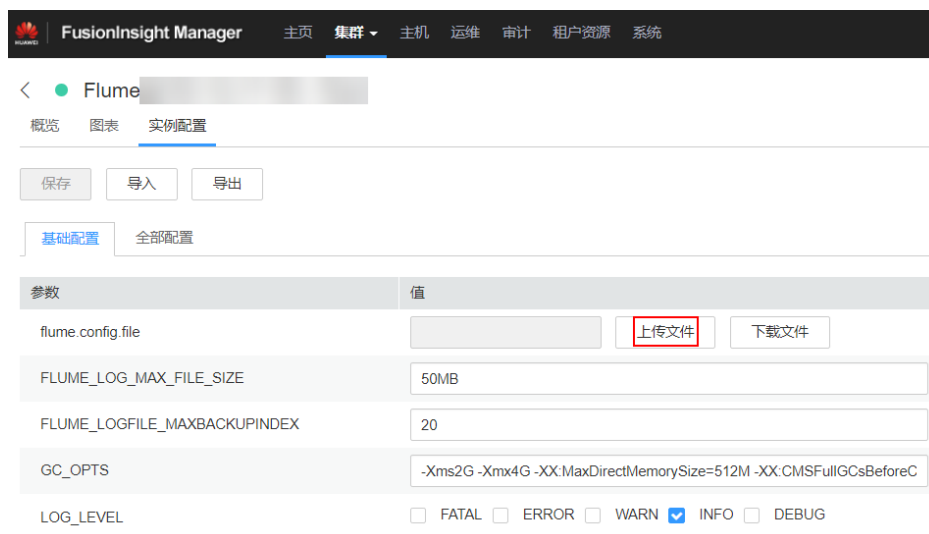
4. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume”服务，在“角色”下单击“Flume”角色。

图 7-37 单击 Flume 角色



5. 选择准备上传配置文件的节点行的“Flume”角色，单击“实例配置”页面“flume.config.file”参数后的“上传文件”，选择“properties.properties”文件完成操作。

图 7-38 上传文件



### 说明

- 每个Flume实例均可以上传单独的服务端配置文件。
- 更新配置文件需要按照此步骤操作，后台修改配置文件是不规范操作，同步配置时后台做的修改将会被覆盖。

6. 单击“保存”，单击“确定”。单击“完成”完成操作。

### 步骤3 设置Flume角色客户端参数。

1. 执行以下命令将生成的客户端证书（flume\_cChat.jks）和客户端信任列表（flume\_cChatt.jks）复制到客户端目录下，如“/opt/flume-client/fusionInsight-flume-1.9.0/conf/”（要求已安装Flume客户端），其中10.196.26.1为客户端所在节点业务平面的IP地址。

```
scp ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf/flume_cChat.jks user@10.196.26.1:/opt/flume-client/fusionInsight-flume-1.9.0/conf/
```

```
scp ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf/flume_cChatt.jks user@10.196.26.1:/opt/flume-client/fusionInsight-flume-1.9.0/conf/
```

### 说明

复制过程中需要输入客户端所在主机（如10.196.26.1）user用户的密码。

2. 以user用户登录解压Flume客户端的节点。执行以下命令进入客户端目录“opt/flume-client/fusionInsight-flume-1.9.0/bin”。

```
cd opt/flume-client/fusionInsight-flume-1.9.0/bin
```

3. 执行以下命令，生成并得到Flume客户端密钥库密码、信任列表密码和keystore-password加密的私钥信息。连续输入两次密码并确认，该密码是别名为flumechatclient的证书和flume\_cChat.jks证书库的密码。

```
./genPwFile.sh
```

```
cat password.property
```

### 说明

若产生以下错误提示，可执行命令`export JAVA_HOME=JDK路径`进行处理。

```
JAVA_HOME is null in current user,please install the JDK and set the JAVA_HOME
```

4. 执行`echo $SCC_PROFILE_DIR`检查SCC\_PROFILE\_DIR环境变量是否为空。
  - 是，执行`source .sccfile`。
  - 否，执行[步骤3.5](#)。
5. 使用FusionInsight Manager界面中的Flume配置工具来配置Flume角色客户端参数并生成配置文件。
  - a. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具”。

图 7-39 选择配置工具



- b. “Agent名”选择“client”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

例如采用SpoolDir Source、File Channel和Avro Sink，如[图7-40](#)所示。

图 7-40 Flume 配置工具示例



- c. 双击对应的source、channel以及sink，根据实际环境并参考[表7-18](#)设置对应的配置参数。

### 说明

- 如果对应的Flume角色之前已经配置过客户端参数，为保证与之前的配置保持一致，可以到“客户端安装目录/fusioninsight-flume-1.9.0/conf/properties.properties”获取已有的客户端参数配置文件。然后登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改加密传输的相关配置项即可。
  - 导入配置文件时，建议配置source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。
  - 不同的File Channel均需要配置一个不同的checkpoint目录。
- d. 单击“导出”，将配置文件“properties.properties”保存到本地。

表 7-18 Flume 角色客户端所需修改的参数列表

| 参数名称                | 参数值填写规则                                                                                                                  | 参数样例                                                                          |
|---------------------|--------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------|
| ssl                 | 是否启用SSL认证（基于安全要求，建议用户启用此功能）。<br><br><ul style="list-style-type: none"> <li>▪ true表示启用。</li> <li>▪ false表示不启用。</li> </ul> | true                                                                          |
| keystore            | 客户端证书。                                                                                                                   | /opt/flume-client/<br>fusionInsight-<br>flume-1.9.0/conf/<br>flume_cChat.jks  |
| keystore-password   | 密钥库密码，获取keystore信息所需密码。<br>输入步骤3.3中获取的“password”值。                                                                       | -                                                                             |
| truststore          | 客户端的SSL证书信任列表。                                                                                                           | /opt/flume-client/<br>fusionInsight-<br>flume-1.9.0/conf/<br>flume_cChatt.jks |
| truststore-password | 信任列表密码，获取truststore信息所需密码。<br>输入步骤3.3中获取的“password”值。                                                                    | -                                                                             |

6. 将“properties.properties”文件上传到Flume客户端安装目录下的“flume/conf/”下。

----结束

## 7.6.2 使用多级 agent 串联从本地采集静态日志保存到 HDFS

### 操作场景

该任务指导用户使用Flume从本地采集静态日志保存到HDFS上如下目录“/flume/test”。

本章节适用于MRS 3.x及之后版本。

### 前提条件

- 已成功安装集群、HDFS及Flume服务、Flume客户端。
- 已创建用户flume\_hdfs并授权验证日志时操作的HDFS目录和数据。



## 操作步骤

**步骤1** 分别生成Flume角色服务端和客户端的证书和信任列表。

1. 以omm用户登录Flume服务端所在节点。进入“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/bin`”目录。

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/bin
```

2. 执行以下命令，生成并导出Flume角色服务端、客户端证书。

```
sh geneJKS.sh -f 密码 -g 密码
```

生成的证书在“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf`”路径下。其中：

- “flume\_sChat.jks”是Flume角色服务端的证书库，“flume\_sChat.crt”是“flume\_sChat.jks”证书的导出文件，“-f”配置项是证书和证书库的密码；
- “flume\_cChat.jks”是Flume角色客户端的证书库，“flume\_cChat.crt”是“flume\_cChat.jks”证书的导出文件，“-g”配置项是证书和证书库的密码；
- “flume\_sChatt.jks”和“flume\_cChatt.jks”分别为Flume服务端、客户端SSL证书信任列表。

### 📖 说明

本章节涉及到所有的用户自定义密码，需满足以下复杂度要求：

- 至少包含大写字母、小写字母、数字、特殊符号4种类型字符
- 至少8位，最多64位
- 出于安全考虑，建议用户定期更换自定义密码（例如三个月更换一次），并重新生成各项证书和信任列表。

**步骤2** 在FusionInsight Manager管理界面，选择“系统 > 权限 > 用户”，选择“更多 > 下载认证凭据”下载用户flume\_hdfs的kerberos证书文件并保存在本地。

**步骤3** 配置Flume角色的服务端参数，并将配置文件上传到集群。

1. 以omm用户登录任意一个Flume角色所在的节点。执行以下命令进入“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/bin`”。

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/bin
```

2. 执行以下命令，生成并得到Flume服务端密钥库密码、信任列表密码和keystore-password加密的私钥信息。连续输入两次密码并确认，该密码是flume\_sChat.jks证书库的密码。

```
./genPwFile.sh
```

```
cat password.property
```

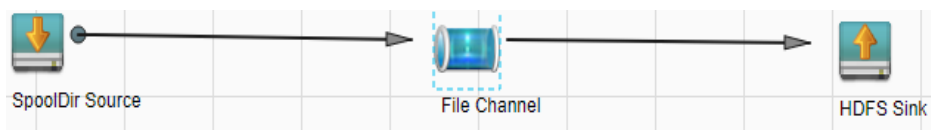
3. 使用FusionInsight Manager界面中的Flume配置工具来配置服务端参数并生成配置文件。
  - a. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具”。

图 7-41 选择配置工具



- b. “Agent名”选择“server”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。  
采用SpoolDir Source、File Channel和HDFS Sink，如图7-42所示。

图 7-42 Flume 配置工具示例



- c. 双击对应的source、channel以及sink，根据实际环境并参考表7-19设置对应的配置参数。

说明

- 如果对应的Flume角色之前已经配置过服务端参数，为保证与之前的配置保持一致，在FusionInsight Manager界面选择“集群 > 待操作集群的名称 > 服务 > Flume > 实例”，选择相应的Flume角色实例，单击“实例配置”页面“flume.config.file”参数后的“下载文件”按钮，可获取已有的服务端参数配置文件。然后选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改加密传输的相关配置项即可。
  - 导入配置文件时，建议配置source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。
  - 不同的File Channel均需要配置一个不同的checkpoint目录。
- d. 单击“导出”，将配置文件“properties.properties”保存到本地。配置文件中包含认证密码信息可能存在安全风险，建议当前场景执行完毕后删除相关配置文件或加强安全管理。

表 7-19 Flume 角色服务端所需修改的参数列表

| 参数名称 | 参数值填写规则                                           | 参数样例           |
|------|---------------------------------------------------|----------------|
| 名称   | 不能为空，必须唯一。                                        | test           |
| bind | avro source绑定的ip地址，此参数不能为空。须配置为服务端配置文件即将要上传的主机IP。 | 192.168.108.11 |
| port | avro source监测的端口,此参数不能为空。须配置为未被使用的端口。             | 21154          |

| 参数名称                | 参数值填写规则                                                                                                                                          | 参数样例                                                                                                                     |
|---------------------|--------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------|
| ssl                 | <p>是否启用SSL认证（基于安全要求，建议用户启用此功能）。只有“Avro”类型的Source才有此配置项。</p> <ul style="list-style-type: none"> <li>▪ true表示启用。</li> <li>▪ false表示不启用。</li> </ul> | true                                                                                                                     |
| keystore            | 服务端证书。                                                                                                                                           | <code>\${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf/flume_sChat.jks</code>  |
| keystore-password   | <p>密钥库密码，获取keystore信息所需密码。</p> <p>输入<a href="#">步骤3.2</a>中获取的“password”值。</p>                                                                    | -                                                                                                                        |
| truststore          | 服务端的SSL证书信任列表。                                                                                                                                   | <code>\${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf/flume_sChatt.jks</code> |
| truststore-password | <p>信任列表密码，获取truststore信息所需密码。</p> <p>输入<a href="#">步骤3.2</a>中获取的“password”值。</p>                                                                 | -                                                                                                                        |
| dataDirs            | 缓冲区数据保存目录，默认为运行目录。配置多个盘上的目录可以提升传输效率，多个目录使用逗号分隔。如果为集群内，则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/data，dataX为data1~dataN。如果为集群外，则需要单独规划。           | <code>/srv/BigData/hadoop/data1/flumeserver/data</code>                                                                  |
| checkpointDir       | checkpoint 信息保存目录，默认在运行目录下。如果为集群内，则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/checkpoint，dataX为data1~dataN。如果为集群外，则需要单独规划。                         | <code>/srv/BigData/hadoop/data1/flumeserver/checkpoint</code>                                                            |

| 参数名称                   | 参数值填写规则                                                              | 参数样例                                                                                                                                |
|------------------------|----------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------|
| transactionCapacity    | 事务大小：即当前channel支持事务处理的事件个数。建议和Source的batchSize设置为同样大小，不能小于batchSize。 | 61200                                                                                                                               |
| hdfs.path              | 写入HDFS的目录，此参数不能为空。                                                   | hdfs://hacluster/flume/test                                                                                                         |
| hdfs.inUsePrefix       | 正在写入HDFS的文件的前缀。                                                      | TMP_                                                                                                                                |
| hdfs.batchSize         | 一次写入HDFS的最大事件数目。                                                     | 61200                                                                                                                               |
| hdfs.kerberosPrincipal | kerberos认证时用户，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。                      | flume_hdfs                                                                                                                          |
| hdfs.kerberosKeytab    | kerberos认证时keytab文件路径，在安全版本下必须填写。安全集群需要配置此项，普通模式集群无需配置。              | /opt/test/conf/user.keytab<br><b>说明</b><br>user.keytab文件从下载用户flume_hdfs的kerberos证书文件中获取，另外，确保用于安装和运行Flume客户端的用户对user.keytab文件有读写权限。 |
| hdfs.useLocalTimeStamp | 是否使用本地时间，取值为"true"或者"false"。                                         | true                                                                                                                                |

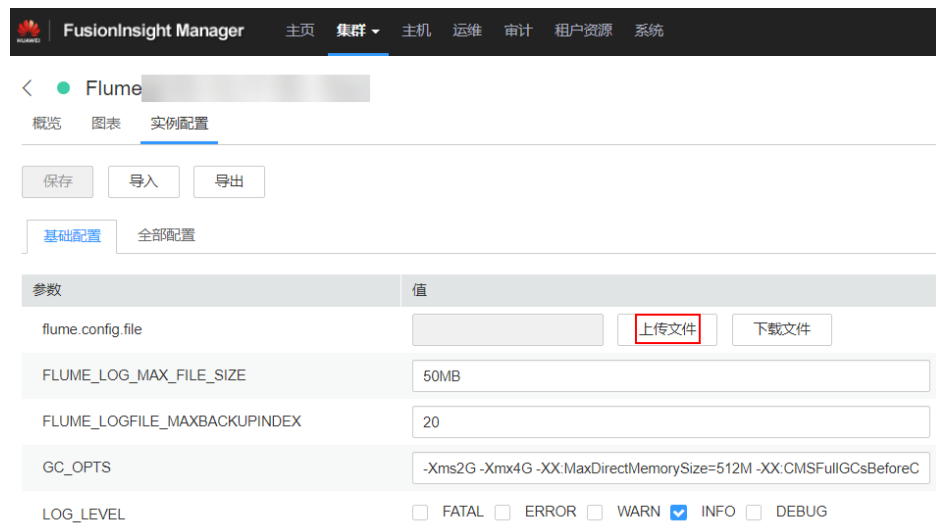
4. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume”，在“角色”下单击“Flume”角色。

图 7-43 单击 Flume 角色



5. 选择准备上传配置文件的节点行的“Flume”角色，单击“实例配置”页面“flume.config.file”参数后的“上传文件”，选择“properties.properties”文件完成操作。

图 7-44 上传文件



### 说明

- 每个Flume实例均可以上传单独的服务端配置文件。
  - 更新配置文件需要按照此步骤操作，后台修改配置文件是不规范操作，同步配置时后台做的修改将会被覆盖。
6. 单击“保存”，单击“确定”。

7. 单击“完成”完成操作。

#### 步骤4 配置Flume角色客户端参数。

1. 执行以下命令将生成的客户端证书（flume\_cChat.jks）和客户端信任列表（flume\_cChatt.jks）复制到客户端目录下，如“/opt/flume-client/fusionInsight-flume-1.9.0/conf/”（要求已安装Flume客户端），其中10.196.26.1为客户端所在节点业务平面的IP地址。

```
scp ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf/flume_cChat.jks user@10.196.26.1:/opt/flume-client/fusionInsight-flume-1.9.0/conf/
```

```
scp ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Flume-1.9.0/flume/conf/flume_cChatt.jks user@10.196.26.1:/opt/flume-client/fusionInsight-flume-1.9.0/conf/
```

#### 📖 说明

复制过程中需要输入客户端所在主机（如10.196.26.1）user用户的密码。

2. 以user用户登录解压Flume客户端的节点。执行以下命令进入客户端目录“/opt/flume-client/fusionInsight-flume-1.9.0/bin”。

```
cd opt/flume-client/fusionInsight-flume-1.9.0/bin
```

3. 执行以下命令，生成并得到Flume客户端密钥库密码、信任列表密码和keystore-password加密的私钥信息。连续输入两次密码并确认，该密码是别名为flumechatclient的证书和flume\_cChat.jks证书库的密码。

```
./genPwFile.sh
```

```
cat password.property
```

#### 📖 说明

若产生以下错误提示，可执行命令export JAVA\_HOME=JDK路径进行处理。

```
JAVA_HOME is null in current user,please install the JDK and set the JAVA_HOME
```

4. 执行echo \$SCC\_PROFILE\_DIR检查SCC\_PROFILE\_DIR环境变量是否为空。
  - 是，执行source .sccfile。
  - 否，执行[步骤4.5](#)。
5. 使用FusionInsight Manager界面中的Flume配置工具来配置Flume角色客户端参数并生成配置文件。
  - a. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具”。

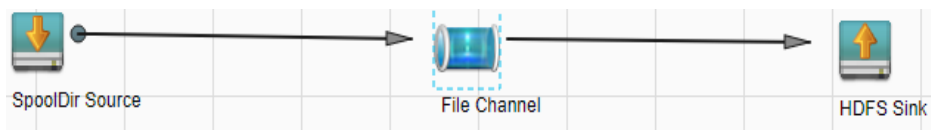
图 7-45 选择配置工具



- b. “Agent名”选择“client”，然后选择要使用的source、channel以及sink，将其拖到右侧的操作界面中并将其连接。

采用SpoolDir Source、File Channel和HDFS Sink，如[图7-46](#)所示。

图 7-46 Flume 配置工具示例



- c. 双击对应的source、channel以及sink，根据实际环境并参考表7-20设置对应的配置参数。

**说明**

- 如果对应的Flume角色之前已经配置过客户端参数，为保证与之前的配置保持一致，可以到“客户端安装目录/fusioninsight-flume-1.9.0/conf/properties.properties”获取已有的客户端参数配置文件。然后登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Flume > 配置工具 > 导入”，将该文件导入后再修改加密传输的相关配置项即可。
  - 导入配置文件时，建议配置中source/channel/sink的各自的个数都不要超过40个，否则可能导致界面响应时间过长。
- d. 单击“导出”，将配置文件“properties.properties”保存到本地。

表 7-20 Flume 角色客户端所需修改的参数列表

| 参数名称       | 参数值填写规则                                                                                                                                | 参数样例                                 |
|------------|----------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------|
| 名称         | 不能为空，必须唯一。                                                                                                                             | test                                 |
| spoolDir   | 待采集的文件所在的目录路径，此参数不能为空。该路径需存在，且对flume运行用户有读写执行权限。                                                                                       | /srv/BigData/hadoop/data1/zb         |
| trackerDir | flume采集文件信息元数据保存路径。                                                                                                                    | /srv/BigData/hadoop/data1/tracker    |
| batch-size | Flume一次发送数据的最大事件数。                                                                                                                     | 61200                                |
| dataDirs   | 缓冲区数据保存目录，默认为运行目录。配置多个盘上的目录可以提升传输效率，多个目录使用逗号分隔。如果为集群内，则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/data，dataX为data1~dataN。如果为集群外，则需要单独规划。 | /srv/BigData/hadoop/data1/flume/data |

| 参数名称                | 参数值填写规则                                                                                                                                          | 参数样例                                                             |
|---------------------|--------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------|
| checkpointDir       | checkpoint信息保存目录，默认在运行目录下。如果为集群内，则可以指定在如下目录/srv/BigData/hadoop/dataX/flume/checkpoint，dataX为data1~dataN。如果为集群外，则需要单独规划。                          | /srv/BigData/hadoop/data1/flume/checkpoint                       |
| transactionCapacity | 事务大小：即当前channel支持事务处理的事件个数，建议和Source的batchSize设置为同样大小，不能小于batchSize。                                                                             | 61200                                                            |
| hostname            | 要发送数据的主机名或者IP，此参数不能为空。须配置为与之相连的avro source所在的主机名或IP。                                                                                             | 192.168.108.11                                                   |
| port                | avro sink监测的端口，此参数不能为空。须配置为与之相连的avro source监测的端口。                                                                                                | 21154                                                            |
| ssl                 | 是否启用SSL认证（基于安全要求，建议用户启用此功能）。<br>只有“Avro”类型的Source才有此配置项。<br><ul style="list-style-type: none"> <li>▪ true表示启用。</li> <li>▪ false表示不启用。</li> </ul> | true                                                             |
| keystore            | 服务端生成的flume_cChat.jks证书。                                                                                                                         | /opt/flume-client/fusionInsight-flume-1.9.0/conf/flume_cChat.jks |
| keystore-password   | 密钥库密码，获取keystore信息所需密码。<br>输入 <a href="#">步骤4.3</a> 中获取的“password”值。                                                                             | -                                                                |



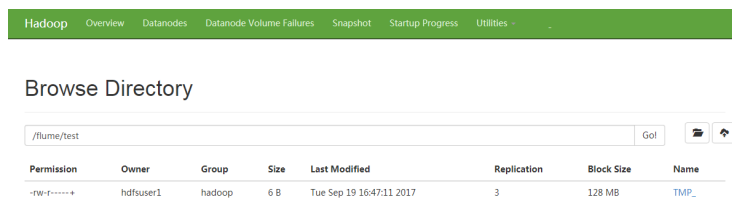
| 参数名称                | 参数值填写规则                                                        | 参数样例                                                              |
|---------------------|----------------------------------------------------------------|-------------------------------------------------------------------|
| truststore          | 服务端的SSL证书信任列表。                                                 | /opt/flume-client/fusionInsight-flume-1.9.0/conf/flume_cChart.jks |
| truststore-password | 信任列表密码，获取truststore信息所需密码。<br>输入 <b>步骤4.3</b> 中获取的“password”值。 | -                                                                 |

6. 将“properties.properties”文件上传到Flume客户端安装目录下的“flume/conf/”下。

#### 步骤5 验证日志是否传输成功。

1. 以具有HDFS组件管理权限的用户登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。在FusionInsight Manager界面选择“集群 > 待操作集群的名称 > 服务 > HDFS”，单击“NameNode(节点名称, 主)”对应的链接，打开HDFS WebUI，然后选择“Utilities > Browse the file system”
2. 观察HDFS上“/flume/test”目录下是否有产生数据。

图 7-47 查看 HDFS 目录和文件



----结束

## 7.7 Flume 企业级能力增强

### 7.7.1 使用 Flume 客户端加密工具

安装Flume客户端后，配置文件的部分参数可能需要填写加密的字符，Flume客户端中提供了加密工具。

**步骤1** 安装Flume客户端。

**步骤2** 登录安装Flume客户端的节点，并切换到客户端安装目录。例如“/opt/FlumeClient”。

**步骤3** 切换到以下目录

```
cd fusioninsight-flume-Flume组件版本号/bin
```

**步骤4** 执行以下命令，加密原始信息：

```
./genPwFile.sh
```

输入两次待加密信息。

**步骤5** 执行以下命令，查看加密后的信息：

```
cat password.property
```

#### 📖 说明

如果加密参数是用于Flume Server，那么需要到相应的Flume Server所在节点执行加密。需要使用omm用户执行加密脚本进行加密。

- 针对MRS 3.x之前版本加密路径为“/opt/Bigdata/MRS\_XXX/install/FusionInsight-Flume-Flume组件版本号/flume/bin/genPwFile.sh”。
- 针对MRS 3.x及之后版本加密路径为“/opt/Bigdata/FusionInsight\_Porter\_XXX/install/FusionInsight-Flume-Flume组件版本号/flume/bin/genPwFile.sh”。其中XXX为产品的版本号。

----结束

## 7.7.2 配置 Flume 对接安全模式 Kafka

使用Flume客户端对接安全kafka。

**步骤1** 新增jaas.conf文件，并保存到“\${Flume客户端安装目录}/conf”下，jaas.conf文件内容如下：

```
KafkaClient {
 com.sun.security.auth.module.Krb5LoginModule required
 useKeyTab=true
 keyTab="/opt/test/conf/user.keytab"
 principal="flume_hdfs@<系统域名>"
 useTicketCache=false
 storeKey=true
 debug=true;
};
```

其中keyTab和principal的值请按照实际情况配置，所配置的principal需要有相应的kafka的权限。

**步骤2** 配置业务，其中kafka.bootstrap.servers的端口号使用21007，kafka.security.protocol使用SASL\_PLAINTEXT。

**步骤3** 如果Kafka所在集群的域名发生了更改，需要对\${Flume客户端安装目录}/conf/flume-env.sh文件中的-Dkerberos.domain.name项的值做修改，具体请根据实际域名进行配置。

**步骤4** 上传所配置的properties.properties文件到\${Flume客户端安装目录}/conf目录下。

----结束

## 7.8 Flume 运维管理

### 7.8.1 Flume 常用配置参数

MRS 3.x之前版本需在“properties.properties”文件中配置。

MRS 3.x及之后版本，部分参数可在Manager界面配置。

## 基本介绍

使用Flume需要配置Source、Channel和Sink，各模块配置参数说明可通过本节内容了解。

MRS 3.x及之后版本部分参数可通过Manager界面配置，选择“集群 > 服务 > Flume > 配置工具”，选择要使用的Source、Channel以及Sink，将其拖到右侧的操作界面中，双击对应的Source、Channel以及Sink，根据实际环境可配置Source、Channel和Sink参数。“channels”、“type”等参数仅在客户端配置文件“properties.properties”中进行配置，配置文件路径为“*Flume客户端安装目录/ fusioninsight-flume-Flume组件版本号/conf/properties.properties*”。

### 说明

部分配置可能需要填写加密后的信息，请参见[使用Flume客户端加密工具](#)。

## 常用 Source 配置

- **Avro Source**

Avro Source监测Avro端口，接收外部Avro客户端数据并放入配置的Channel中。常用配置如[表7-21](#)所示：

图 7-48 Avro Source

### Avro Source-Avro Source

|                     |                                                                         |
|---------------------|-------------------------------------------------------------------------|
| * 名称                | <input type="text"/>                                                    |
| * bind              | <input type="text"/>                                                    |
| * port              | <input type="text"/>                                                    |
| threads             | <input type="text"/>                                                    |
| compression-type    | <input type="text" value="none"/>                                       |
| ssl                 | <input checked="" type="checkbox"/> true <input type="checkbox"/> false |
| keystore-type       | <input type="text" value="JKS"/>                                        |
| keystore            | <input type="text"/>                                                    |
| keystore-password   | <input type="text"/>                                                    |
| truststore-type     | <input type="text" value="JKS"/>                                        |
| truststore          | <input type="text"/>                                                    |
| truststore-password | <input type="text"/>                                                    |
| additional-items    | <input type="text"/>                                                    |

—

表 7-21 Avro Source 常用配置

| 参数                  | 默认值   | 描述                                                                                                                                                                                                                                                                                                                                                                    |
|---------------------|-------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| channels            | -     | 与之相连的Channel，可以配置多个。用空格隔开。<br>在单个代理流程中，是通过channel连接sources和sinks。一个source实例对应多个channels，但一个sink实例只能对应一个channel。<br>格式如下：<br><b>&lt;Agent<br/>&gt;.sources.&lt;Source&gt;.channels =<br/>&lt;channel1&gt; &lt;channel2&gt;<br/>&lt;channel3&gt;...</b><br><b>&lt;Agent &gt;.sinks.&lt;Sink&gt;.channels =<br/>&lt;channel1&gt;</b><br>仅可在“properties.properties”文件中配置。 |
| type                | avro  | 类型，需设置为“avro”。每一种source的类型都为相应的固定值。<br>仅可在“properties.properties”文件中配置。                                                                                                                                                                                                                                                                                               |
| bind                | -     | 绑定和source关联的主机名或IP地址。                                                                                                                                                                                                                                                                                                                                                 |
| port                | -     | 绑定端口号。                                                                                                                                                                                                                                                                                                                                                                |
| ssl                 | false | 是否使用SSL加密。 <ul style="list-style-type: none"><li>• true</li><li>• false</li></ul>                                                                                                                                                                                                                                                                                     |
| truststore-type     | JKS   | Java信任库类型。填写JKS或其他java支持的truststore类型。                                                                                                                                                                                                                                                                                                                                |
| truststore          | -     | Java信任库文件。                                                                                                                                                                                                                                                                                                                                                            |
| truststore-password | -     | Java信任库密码。                                                                                                                                                                                                                                                                                                                                                            |
| keystore-type       | JKS   | 密钥存储类型。填写JKS或其他java支持的truststore类型。                                                                                                                                                                                                                                                                                                                                   |
| keystore            | -     | 密钥存储文件。                                                                                                                                                                                                                                                                                                                                                               |
| keystore-password   | -     | 密钥存储密码。                                                                                                                                                                                                                                                                                                                                                               |

- **SpoolDir Source**

SpoolDir Source监控并传输目录下新增的文件，可实现准实时数据传输。常用配置如表 2 [Spooling Source常用配置](#)所示：

图 7-49 SpoolDir Source

### SpoolDir Source-SpoolDir Source

|                            |                                                                         |
|----------------------------|-------------------------------------------------------------------------|
| * 名称                       | <input type="text"/>                                                    |
| * spoolDir                 | <input type="text"/>                                                    |
| montime                    | <input type="text"/>                                                    |
| fileSuffix                 | <input type="text" value=".COMPLETED"/>                                 |
| deletePolicy               | <input type="text" value="never"/>                                      |
| trackerDir                 | <input type="text" value=".flumespool"/>                                |
| ignorePattern              | <input type="text" value="^\$"/>                                        |
| batchSize                  | <input type="text" value="1000"/>                                       |
| inputCharset               | <input type="text" value="UTF-8"/>                                      |
| selector.type              | <input type="text" value="replicating"/>                                |
| fileHeader                 | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |
| basenameHeader             | <input checked="" type="checkbox"/> true <input type="checkbox"/> false |
| basenameHeaderKey          | <input type="text" value="basename"/>                                   |
| deserializer               | <input type="text" value="LINE"/>                                       |
| deserializer.maxBatchLine  | <input type="text" value="1"/>                                          |
| deserializer.maxLineLength | <input type="text" value="2048"/>                                       |
| additional-items           | <input type="text"/>                                                    |

-

表 7-22 SpoolDir Source 常用配置

| 参数                         | 默认值         | 描述                                                                                                                                                                                                                                                                                  |
|----------------------------|-------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| channels                   | -           | 与之相连的Channel，可以配置多个。仅可在“properties.properties”文件中配置。                                                                                                                                                                                                                                |
| type                       | spooldir    | 类型，需设置为“spooldir”。仅可在“properties.properties”文件中配置。                                                                                                                                                                                                                                  |
| monTime                    | 0（不开启）      | 线程监控阈值，更新时间大于阈值时会重新启动该Source，单位：秒。                                                                                                                                                                                                                                                  |
| spoolDir                   | -           | 监控目录。                                                                                                                                                                                                                                                                               |
| fileSuffix                 | .COMPLETED  | 文件传输完成后添加的后缀。                                                                                                                                                                                                                                                                       |
| deletePolicy               | never       | 文件传输完成后源文件删除策略，支持“never”或“immediate”。分别是从不删除和立即删除。                                                                                                                                                                                                                                  |
| ignorePattern              | ^\$         | 忽略文件的正则表达式表示。                                                                                                                                                                                                                                                                       |
| trackerDir                 | .flumespool | 传输过程中元数据存储路径。                                                                                                                                                                                                                                                                       |
| batchSize                  | 1000        | Source传输粒度。                                                                                                                                                                                                                                                                         |
| decodeErrorPolicy          | FAIL        | <p>编码错误策略。仅可在“properties.properties”文件中配置。</p> <p>可选FAIL、REPLACE、IGNORE。</p> <p>FAIL：发生异常并让解析失败。</p> <p>REPLACE：将不能识别的字符用其它字符代替，通常是字符U+FFFD。</p> <p>IGNORE：直接丢弃不能解析的字符串。</p> <p><b>说明</b><br/>如果文件中有编码错误，请配置“decodeErrorPolicy”为“REPLACE”或“IGNORE”，Flume遇到编码错误将跳过编码错误，继续采集后续日志。</p> |
| deserializer               | LINE        | <p>文件解析器，值为“LINE”或“BufferedLine”。</p> <ul style="list-style-type: none"> <li>配置为“LINE”时，对从文件读取的字符逐个转码。</li> <li>配置为“BufferedLine”时，对文件读取的一行或多行的字符进行批量转码，性能更优。</li> </ul>                                                                                                            |
| deserializer.maxLineLength | 2048        | 按行解析最大长度。0到2,147,483,647。                                                                                                                                                                                                                                                           |

| 参数                        | 默认值         | 描述                                                                                                                                                                                                          |
|---------------------------|-------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| deserializer.maxBatchLine | 1           | 按行解析最多行数，如果行数设置为多行，“maxLineLength”也应该设置为相应的倍数。例如maxBatchLine设置为2，“maxLineLength”相应的设置为2048*2为4096。                                                                                                          |
| selector.type             | replicating | 选择器类型，支持“replicating”或“multiplexing”。 <ul style="list-style-type: none"> <li>“replicating”表示同样的内容会发给每一个channel。</li> <li>“multiplexing”表示根据分发规则，有选择地发给某些channel。</li> </ul>                                 |
| interceptors              | -           | 拦截器配置。详细配置可参考Flume官方文档： <a href="https://flume.apache.org/FlumeUserGuide.html#flume-interceptors">https://flume.apache.org/FlumeUserGuide.html#flume-interceptors</a> 。<br>仅可在“properties.properties”文件中配置。 |

### 📖 说明

Spooling Source在按行读取过程中，会忽略掉每一个Event的最后一个换行符，该换行符所占用的数据量指标不会被Flume统计。

- **Kafka Source**

Kafka Source从Kafka的topic中消费数据，可以设置多个Source消费同一个topic的数据，每个Source会消费topic的不同partitions。常用配置如表 3 [Kafka Source常用配置](#)所示：



图 7-50 Kafka Source

### Kafka Source-Kafka Source

|                           |                                                      |
|---------------------------|------------------------------------------------------|
| * 名称                      | <input type="text"/>                                 |
| * kafka.topics            | <input type="text"/>                                 |
| montime                   | <input type="text"/>                                 |
| nodatime                  | <input type="text" value="0"/>                       |
| kafka.topics.regex        | <input type="text"/>                                 |
| * kafka.consumer.group.id | <input type="text"/>                                 |
| kafka.bootstrap.servers   | <input type="text" value="例: 192.168.1.100:21007;"/> |
| kafka.security.protocol   | <input type="text" value="SASL_PLAINTEXT"/>          |
| batchDurationMillis       | <input type="text" value="1000"/>                    |
| batchSize                 | <input type="text" value="1000"/>                    |
| additional-items          | <input type="text"/>                                 |

确认 取消

表 7-23 Kafka Source 常用配置

| 参数       | 默认值                                       | 描述                                                                                  |
|----------|-------------------------------------------|-------------------------------------------------------------------------------------|
| channels | -                                         | 与之相连的Channel，可以配置多个。仅可在“properties.properties”文件中配置。                                |
| type     | org.apache.flume.source.kafka.KafkaSource | 类型，需设置为“org.apache.flume.source.kafka.KafkaSource”。仅可在“properties.properties”文件中配置。 |

| 参数                         | 默认值            | 描述                                                                                                                                                                          |
|----------------------------|----------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| monTime                    | 0（不开启）         | 线程监控阈值，更新时间大于阈值时重新启动该Source，单位：秒。                                                                                                                                           |
| nodatotime                 | 0（不开启）         | 告警阈值，从Kafka中订阅不到数据的时长大于阈值时发送告警，单位：秒。                                                                                                                                        |
| batchSize                  | 1000           | 每次写入Channel的Event数量。                                                                                                                                                        |
| batchDurationMillis        | 1000           | 每次消费topic数据的最大时长，单位：毫秒。                                                                                                                                                     |
| keepTopicInHeader          | false          | 是否在Event Header中保存topic，如果保存，Kafka Sink配置的topic将无效。<br><ul style="list-style-type: none"> <li>• true</li> <li>• false</li> </ul> 仅可在“properties.properties”文件中配置。           |
| keepPartitionInHeader      | false          | 是否在Event Header中保存partitionID，如果保存，Kafka Sink将写入对应的Partition。<br><ul style="list-style-type: none"> <li>• true</li> <li>• false</li> </ul> 仅可在“properties.properties”文件中配置。 |
| kafka.bootstrap.servers    | -              | brokers地址列表，多个地址用英文逗号分隔。                                                                                                                                                    |
| kafka.consumer.group.id    | -              | Kafka消费者组ID。                                                                                                                                                                |
| kafka.topics               | -              | 订阅的kafka topic列表，用英文逗号分隔。                                                                                                                                                   |
| kafka.topics.regex         | -              | 符合正则表达式的topic会被订阅，优先级高于“kafka.topics”，如果配置将覆盖“kafka.topics”。                                                                                                                |
| kafka.security.protocol    | SASL_PLAINTEXT | Kafka安全协议，未启用Kerberos集群中须配置为“PLAINTEXT”。                                                                                                                                    |
| kafka.kerberos.domain.name | -              | 此参数的值为Kafka集群中kerberos的“default_realm”，仅安全集群需要配置。<br>仅可在“properties.properties”文件中配置。                                                                                       |

| 参数                              | 默认值 | 描述                                                                                  |
|---------------------------------|-----|-------------------------------------------------------------------------------------|
| Other Kafka Consumer Properties | -   | 其他Kafka配置，可以接受任意Kafka支持的消费参数配置，配置需要加前缀“.kafka”。<br>仅可在“properties.properties”文件中配置。 |

- **Taildir Source**

Taildir Source监控目录下文件的变化并自动读取文件内容，可实现实时数据传输，常用配置如表7-24所示：

图 7-51 Taildir Source

### Taildir Source-Taildir Source

|                  |                                                                         |
|------------------|-------------------------------------------------------------------------|
| * 名称             | <input type="text"/>                                                    |
| * filegroups     | <input type="text"/>                                                    |
| * positionFile   | <input type="text"/>                                                    |
| montime          | <input type="text"/>                                                    |
| byteOffsetHeader | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |
| skipToEnd        | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |
| idleTimeout      | <input type="text" value="12000"/>                                      |
| writePosInterval | <input type="text" value="3000"/>                                       |
| batchSize        | <input type="text" value="1000"/>                                       |
| additional-items | <input type="text"/>                                                    |
| fileHeader       | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |

—

表 7-24 Taildir Source 常用配置

| 参数                                      | 默认值     | 描述                                                                       |
|-----------------------------------------|---------|--------------------------------------------------------------------------|
| channels                                | -       | 与之相连的Channel，可以配置多个。<br>仅可在“properties.properties”文件中配置。                 |
| type                                    | taildir | 类型，需配置为“taildir”。<br>仅可在“properties.properties”文件中配置。                    |
| filegroups                              | -       | 设置采集文件目录分组名字，分组名字中间使用空格间隔。                                               |
| filegroups.<filegroup Name>.parentDir   | -       | 父目录，需要配置为绝对路径。<br>仅可在“properties.properties”文件中配置。                       |
| filegroups.<filegroup Name>.filePattern | -       | 相对父目录的文件路径，可以包含目录，支持正则表达式，须与父目录联合使用。<br>仅可在“properties.properties”文件中配置。 |
| positionFile                            | -       | 传输过程中元数据存储路径。                                                            |
| headers.<filegroup Name>.<headerKey>    | -       | 设置某一个分组采集数据时Event中的key-value值。<br>仅可在“properties.properties”文件中配置。       |
| byteOffsetHeader                        | false   | 是否在每一个Event头中携带该Event在源文件中的位置信息，该信息保存在“byteoffset”变量中。                   |
| skipToEnd                               | false   | Flume在重启后是否直接定位到文件最新的位置处，以读取最新的数据。                                       |
| idleTimeout                             | 120000  | 设置读取文件的空闲时间，单位：毫秒。如果在该时间内文件内容没有变更，关闭掉该文件，关闭后如果该文件有数据写入，重新打开并读取数据。        |
| writePosInterval                        | 3000    | 设置将元数据写入到文件的周期，单位：毫秒。                                                    |
| batchSize                               | 1000    | 批次写入Channel的Event数量。                                                     |
| monTime                                 | 0（不开启）  | 线程监控阈值，更新时间大于阈值时重新启动该Source，单位：秒。                                        |

- **Http Source**

Http Source接收外部HTTP客户端发送过来的数据，并放入配置的Channel中，常用配置如表7-25所示：

图 7-52 Http Source

### Http Source-Http Source

|                  |                                                                         |
|------------------|-------------------------------------------------------------------------|
| * 名称             | <input type="text"/>                                                    |
| * bind           | <input type="text"/>                                                    |
| * port           | <input type="text"/>                                                    |
| handler          | <input type="text" value="org.apache.flume.source.ht"/>                 |
| handler.*        | <input type="text"/>                                                    |
| enableSSL        | <input checked="" type="checkbox"/> true <input type="checkbox"/> false |
| keystore         | <input type="text"/>                                                    |
| keystorePassword | <input type="text"/>                                                    |
| additional-items | <input type="text"/>                                                    |

表 7-25 Http Source 常用配置

| 参数       | 默认值  | 描述                                                   |
|----------|------|------------------------------------------------------|
| channels | -    | 与之相连的Channel，可以配置多个。仅可在“properties.properties”文件中配置。 |
| type     | http | 类型，需配置为“http”。仅可在“properties.properties”文件中配置。       |
| bind     | -    | 绑定关联的主机名或IP地址。                                       |
| port     | -    | 绑定端口。                                                |

| 参数               | 默认值                                      | 描述                                                                                                                                                                                                    |
|------------------|------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| handler          | org.apache.flume.source.http.JSONHandler | http请求的消息解析方式，支持以下两种： <ul style="list-style-type: none"><li>“org.apache.flume.source.http.JSONHandler”：表示Json格式解析。</li><li>“org.apache.flume.sink.solr.morphline.BlobHandler”：表示二进制Blob块解析。</li></ul> |
| handler.*        | -                                        | 设置handler的参数。                                                                                                                                                                                         |
| enableSSL        | false                                    | http协议是否启用SSL。                                                                                                                                                                                        |
| keystore         | -                                        | http启用SSL后设置keystore的路径。                                                                                                                                                                              |
| keystorePassword | -                                        | http启用SSL后设置keystore的密码。                                                                                                                                                                              |

## 常用 Channel 配置

- **Memory Channel**

Memory Channel使用内存作为缓存区，Events存放在内存队列中。常用配置如[表7-26](#)所示：

图 7-53 Memory Channel

### Memory Channel-Memory Channel

|                              |                                    |
|------------------------------|------------------------------------|
| * 名称                         | <input type="text"/>               |
| capacity                     | <input type="text" value="10000"/> |
| transactionCapacity          | <input type="text" value="1000"/>  |
| channelFullcount             | <input type="text" value="10"/>    |
| keep-alive                   | <input type="text" value="3"/>     |
| byteCapacity                 | <input type="text"/>               |
| byteCapacityBufferPercentage | <input type="text" value="20"/>    |
| additional-items             | <input type="text"/>               |

-

表 7-26 Memory Channel 常用配置

| 参数                  | 默认值   | 描述                                               |
|---------------------|-------|--------------------------------------------------|
| type                | -     | 类型，需配置为“memory”。仅可在“properties.properties”文件中配置。 |
| capacity            | 10000 | 缓存在Channel中的最大Event数。                            |
| transactionCapacity | 1000  | 每次存取的最大Event数。                                   |
| channelFullcount    | 10    | Channel full次数，达到该次数后发送告警。                       |

- **File Channel**

File Channel使用本地磁盘作为缓存区，Events存放在设置的“dataDirs”配置项文件夹中。常用配置如[表7-27](#)所示：

图 7-54 File Channel

### File Channel-File Channel

|                      |                                                                         |
|----------------------|-------------------------------------------------------------------------|
| * 名称                 | <input type="text"/>                                                    |
| * dataDirs           | <input type="text" value="/srv/BigData/hadoop/data1,"/>                 |
| * checkpointDir      | <input type="text" value="/srv/BigData/hadoop/data1,"/>                 |
| capacity             | <input type="text" value="1000000"/>                                    |
| channelfullcount     | <input type="text" value="10"/>                                         |
| useDualCheckpoints   | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |
| transactionCapacity  | <input type="text" value="10000"/>                                      |
| checkpointInterval   | <input type="text" value="30000"/>                                      |
| maxFileSize          | <input type="text" value="2146435071"/>                                 |
| minimumRequiredSpace | <input type="text" value="524288000"/>                                  |

表 7-27 File Channel 常用配置

| 参数            | 默认值                                    | 描述                                             |
|---------------|----------------------------------------|------------------------------------------------|
| type          | -                                      | 类型，需配置为“file”。仅可在“properties.properties”文件中配置。 |
| checkpointDir | \${BIGDATA_DATA_HOME}/flume/checkpoint | 检查点存放路径。                                       |
| dataDirs      | \${BIGDATA_DATA_HOME}/flume/data       | 数据缓存路径，设置多个路径可提升性能，中间用逗号分开。                    |
| maxFileSize   | 2146435071                             | 单个缓存文件的最大值，单位：字节。                              |



| 参数                   | 默认值       | 描述                         |
|----------------------|-----------|----------------------------|
| minimumRequiredSpace | 524288000 | 缓冲区空闲空间最小值，单位：字节。          |
| capacity             | 1000000   | 缓存在Channel中的最大Event数。      |
| transactionCapacity  | 10000     | 每次存取的最大Event数。             |
| channelFullcount     | 10        | Channel full次数，达到该次数后发送告警。 |

- **Kafka Channel**

Kafka Channel使用kafka集群缓存数据，Kafka提供高可用、多副本，以防Flume或Kafka Broker崩溃，Channel中的数据会立即被Sink消费。常用配置如[表 10 Kafka Channel 常用配置](#)所示：

图 7-55 Kafka Channel

## Kafka Channel-Kafka Channel

|                                  |                                                                         |
|----------------------------------|-------------------------------------------------------------------------|
| * 名称                             | <input type="text"/>                                                    |
| * kafka.bootstrap.servers        | <input type="text"/>                                                    |
| kafka.topic                      | <input type="text" value="flume-channel"/>                              |
| kafka.consumer.group.id          | <input type="text" value="flume"/>                                      |
| parseAsFlumeEvent                | <input checked="" type="checkbox"/> true <input type="checkbox"/> false |
| migrateZookeeperOffsets          | <input checked="" type="checkbox"/> true <input type="checkbox"/> false |
| kafka.consumer.auto.offset.reset | <input type="text" value="latest"/>                                     |
| kafka.producer.security.protocol | <input type="text" value="SASL_PLAINTEXT"/>                             |
| kafka.consumer.security.protocol | <input type="text" value="SASL_PLAINTEXT"/>                             |
| ignoreLongMessage                | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |

表 7-28 Kafka Channel 常用配置

| 参数                      | 默认值 | 描述                                                                                        |
|-------------------------|-----|-------------------------------------------------------------------------------------------|
| type                    | -   | 类型，需配置为“org.apache.flume.channel.kafka.KafkaChannel”。<br>仅可在“properties.properties”文件中配置。 |
| kafka.bootstrap.servers | -   | kafka broker列表。                                                                           |

| 参数                               | 默认值            | 描述                                           |
|----------------------------------|----------------|----------------------------------------------|
| kafka.topic                      | flume-channel  | Channel用来缓存数据的topic。                         |
| kafka.consumer.group.id          | flume          | Kafka消费者组ID。                                 |
| parseAsFlumeEvent                | true           | 是否解析为Flume event。                            |
| migrateZookeeperOffsets          | true           | 当Kafka没有存储offset时，是否从ZooKeeper中查找，并提交到Kafka。 |
| kafka.consumer.auto.offset.reset | latest         | 当没有offset记录时，从指定的位置消费数据。                     |
| kafka.producer.security.protocol | SASL_PLAINTEXT | Kafka生产者安全协议。                                |
| kafka.consumer.security.protocol | SASL_PLAINTEXT | Kafka消费者安全协议。                                |

## 常用 Sink 配置

- **HDFS Sink**  
HDFS Sink将数据写入HDFS。常用配置如[表7-29](#)所示：

图 7-56 HDFS Sink

HDFS Sink-HDFS Sink

|                          |                                                                         |
|--------------------------|-------------------------------------------------------------------------|
| * 名称                     | <input type="text"/>                                                    |
| * hdfs.path              | <input type="text" value="hdfs://hacluster"/>                           |
| montime                  | <input type="text"/>                                                    |
| hdfs.filePrefix          | <input type="text" value="over_{basename}"/>                            |
| hdfs.fileSuffix          | <input type="text"/>                                                    |
| hdfs.inUsePrefix         | <input type="text"/>                                                    |
| hdfs.inUseSuffix         | <input type="text" value=".tmp"/>                                       |
| hdfs.idleTimeout         | <input type="text" value="0"/>                                          |
| hdfs.batchSize           | <input type="text" value="1000"/>                                       |
| hdfs.codeC               | <input type="text"/>                                                    |
| hdfs.fileType            | <input type="text" value="DataStream"/>                                 |
| hdfs.maxOpenFiles        | <input type="text" value="5000"/>                                       |
| hdfs.writeFormat         | <input type="text" value="Writable"/>                                   |
| hdfs.callTimeout         | <input type="text" value="10000"/>                                      |
| hdfs.threadsPoolSize     | <input type="text" value="10"/>                                         |
| hdfs.rollTimerPoolSize   | <input type="text" value="1"/>                                          |
| hdfs.kerberosPrincipal   | <input type="text"/>                                                    |
| hdfs.kerberosKeytab      | <input type="text"/>                                                    |
| hdfs.round               | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |
| hdfs.roundUnit           | <input type="text" value="second"/>                                     |
| hdfs.useLocalTimeStamp   | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |
| hdfs.failcount           | <input type="text" value="10"/>                                         |
| hdfs.fileCloseByEndEvent | <input checked="" type="checkbox"/> true <input type="checkbox"/> false |
| hdfs.batchCallTimeout    | <input type="text" value="0"/>                                          |
| serializer.appendNewline | <input checked="" type="checkbox"/> true <input type="checkbox"/> false |
| additional-items         | <input type="text"/>                                                    |

-

表 7-29 HDFS Sink 常用配置

| 参数                       | 默认值    | 描述                                                      |
|--------------------------|--------|---------------------------------------------------------|
| channel                  | -      | 与之相连的Channel。仅可在“properties.properties”文件中配置。           |
| type                     | hdfs   | 类型，需配置为“hdfs”。仅可在“properties.properties”文件中配置。          |
| monTime                  | 0（不开启） | 线程监控阈值，更新时间大于阈值时重新启动该Sink，单位：秒。                         |
| hdfs.path                | -      | HDFS路径。                                                 |
| hdfs.inUseSuffix         | .tmp   | 正在写入的HDFS文件后缀。                                          |
| hdfs.rollInterval        | 30     | 按时间滚动文件，单位：秒，同时需将“hdfs.fileCloseByEndEvent”设置为“false”。  |
| hdfs.rollSize            | 1024   | 按大小滚动文件，单位：字节，同时需将“hdfs.fileCloseByEndEvent”设置为“false”。 |
| hdfs.rollCount           | 10     | 按Event个数滚动文件，同时需将“hdfs.fileCloseByEndEvent”设置为“false”。  |
| hdfs.idleTimeout         | 0      | 自动关闭空闲文件超时时间，单位：秒。                                      |
| hdfs.batchSize           | 1000   | 每次写入HDFS的Event个数。                                       |
| hdfs.kerberosPrincipal   | -      | 认证HDFS的Kerberos用户名，未启用Kerberos认证集群不配置。                  |
| hdfs.kerberosKeytab      | -      | 认证HDFS的Kerberos keytab路径，未启用Kerberos认证集群不配置             |
| hdfs.fileCloseByEndEvent | true   | 收到最后一个Event时是否关闭文件。                                     |

| 参数                       | 默认值  | 描述                                                                                                                                                                                                                                            |
|--------------------------|------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| hdfs.batchCallTimeout    | -    | 每次写入HDFS超时控制时间，单位：毫秒。<br>当不配置此参数时，对每个Event写入HDFS进行超时控制。当“hdfs.batchSize”大于0时，配置此参数可以提升写入HDFS性能。<br><b>说明</b><br>“hdfs.batchCallTimeout”设置多长时间需要考虑“hdfs.batchSize”的大小，“hdfs.batchSize”越大，“hdfs.batchCallTimeout”也要调整更长时间，设置过短时间容易导致数据写入HDFS失败。 |
| serializer.appendNewline | true | 将一个Event写入HDFS后是否追加换行符（'\n'），如果追加该换行符，该换行符所占用的数据量指标不会被HDFS Sink统计。                                                                                                                                                                            |

- **Avro Sink**

Avro Sink把events转化为Avro events并发送到配置的主机的监测端口。常用配置如表7-30所示：

图 7-57 Avro Sink

### Avro Sink-Avro Sink

|                           |                                                                         |
|---------------------------|-------------------------------------------------------------------------|
| * 名称                      | <input type="text"/>                                                    |
| * hostname                | <input type="text"/>                                                    |
| * port                    | <input type="text"/>                                                    |
| batch-size                | <input type="text" value="1000"/>                                       |
| connect-timeout           | <input type="text" value="20000"/>                                      |
| request-timeout           | <input type="text" value="20000"/>                                      |
| reset-connection-interval | <input type="text" value="0"/>                                          |
| compression-type          | <input type="text" value="none"/>                                       |
| maxIoWorkers              | <input type="text" value="0"/>                                          |
| ssl                       | <input checked="" type="checkbox"/> true <input type="checkbox"/> false |
| keystore-type             | <input type="text" value="JKS"/>                                        |
| keystore                  | <input type="text"/>                                                    |
| keystore-password         | <input type="text"/>                                                    |
| truststore-type           | <input type="text" value="JKS"/>                                        |
| truststore                | <input type="text"/>                                                    |
| truststore-password       | <input type="text"/>                                                    |
| additional-items          | <input type="text"/>                                                    |

-

表 7-30 Avro Sink 常用配置

| 参数                  | 默认值   | 描述                                             |
|---------------------|-------|------------------------------------------------|
| channel             | -     | 与之相连的Channel。仅可在“properties.properties”文件中配置。  |
| type                | -     | 类型，需配置为“avro”。仅可在“properties.properties”文件中配置。 |
| hostname            | -     | 绑定关联的主机名或IP地址。                                 |
| port                | -     | 监测端口。                                          |
| batch-size          | 1000  | 批次发送的Event个数。                                  |
| ssl                 | false | 是否使用SSL加密。                                     |
| truststore-type     | JKS   | Java信任库类型。                                     |
| truststore          | -     | Java信任库文件。                                     |
| truststore-password | -     | Java信任库密码。                                     |
| keystore-type       | JKS   | 密钥存储类型。                                        |
| keystore            | -     | 密钥存储文件。                                        |
| keystore-password   | -     | 密钥存储密码                                         |

- **HBase Sink**

HBase Sink将数据写入到HBase中。常用配置如[表7-31](#)所示：



图 7-58 HBase Sink

### HBase Sink-HBase Sink

|                    |                                                                         |
|--------------------|-------------------------------------------------------------------------|
| * 名称               | <input type="text"/>                                                    |
| * table            | <input type="text"/>                                                    |
| * columnFamily     | <input type="text"/>                                                    |
| montime            | <input type="text"/>                                                    |
| batchSize          | <input type="text" value="1000"/>                                       |
| coalesceIncrements | <input type="checkbox"/> true <input checked="" type="checkbox"/> false |
| kerberosPrincipal  | <input type="text"/>                                                    |
| kerberosKeytab     | <input type="text"/>                                                    |
| additional-items   | <input type="text"/>                                                    |

表 7-31 HBase Sink 常用配置

| 参数                | 默认值    | 描述                                              |
|-------------------|--------|-------------------------------------------------|
| channel           | -      | 与之相连的Channel。仅可在“properties.properties”文件中配置。   |
| type              | -      | 类型，需配置为“hbase”。仅可在“properties.properties”文件中配置。 |
| table             | -      | HBase表名称。                                       |
| monTime           | 0（不开启） | 线程监控阈值，更新时间大于阈值时重新启动该Sink，单位：秒。                 |
| columnFamily      | -      | HBase列族名称。                                      |
| batchSize         | 1000   | 每次写入HBase的Event个数。                              |
| kerberosPrincipal | -      | 认证HBase的Kerberos用户名，未启用Kerberos认证集群不配置。         |

| 参数             | 默认值 | 描述                                            |
|----------------|-----|-----------------------------------------------|
| kerberosKeytab | -   | 认证HBase的Kerberos keytab路径，未启用Kerberos认证集群不配置。 |

- **Kafka Sink**

Kafka Sink将数据写入到Kafka中。常用配置如表7-32所示：

图 7-59 Kafka Sink

### Kafka Sink-Kafka Sink

\* 名称

kafka.topic

flumeBatchSize

kafka.bootstrap.servers

kafka.security.protocol

ignoreLongMessage  true  false

montime

additional-items

—

表 7-32 Kafka Sink 常用配置

| 参数      | 默认值 | 描述                                            |
|---------|-----|-----------------------------------------------|
| channel | -   | 与之相连的Channel。仅可在“properties.properties”文件中配置。 |

| 参数                              | 默认值                 | 描述                                                                                   |
|---------------------------------|---------------------|--------------------------------------------------------------------------------------|
| type                            | -                   | 类型，需配置为“org.apache.flume.sink.kafka.Kafka Sink”。<br>仅可在“properties.properties”文件中配置。 |
| kafka.bootstrap.servers         | -                   | Kafkabrokers列表，多个用英文逗号分隔。                                                            |
| monTime                         | 0（不开启）              | 线程监控阈值，更新时间大于阈值时重新启动该Sink，单位：秒。                                                      |
| kafka.topic                     | default-flume-topic | 数据写入的topic。                                                                          |
| flumeBatchSize                  | 1000                | 每次写入Kafka的Event个数。                                                                   |
| kafka.security.protocol         | SASL_PLAINTEXT      | Kafka安全协议，未启用Kerberos认证集群下须配置为“PLAINTEXT”。                                           |
| kafka.kerberos.domain.name      | -                   | Kafka Domain名称。安全集群必填。<br>仅可在“properties.properties”文件中配置。                           |
| Other Kafka Producer Properties | -                   | 其他Kafka配置，可以接受任意Kafka支持的生产参数配置，配置需要加前缀“.kafka”。<br>仅可在“properties.properties”文件中配置。  |

## 7.8.2 Flume 业务配置指南

本章节适用于MRS 3.x及之后版本。

该操作指导用户完成Flume常用业务的配置。其他一些不太常用的Source、Channel、Sink的配置请参考Flume社区提供的用户手册（<http://flume.apache.org/releases/1.9.0.html>）。

### 📖 说明

- 各个表格中所示参数，黑体加粗的参数为必选参数。
- Sink的BatchSize参数必须小于Channel的transactionCapacity。
- 集群Flume配置工具界面篇幅有限，Source、Channel、Sink只展示部分参数，详细请参考如下常用配置。
- 集群Flume配置工具界面上所展示Customer Source、Customer Channel及Customer Sink需要用户根据自己开发的代码来进行配置，下述常用配置不再展示。

### 常用 Source 配置

- **Avro Source**

Avro Source监测Avro端口，接收外部Avro客户端数据并放入配置的Channel中。  
常用配置如下表所示：

表 7-33 Avro Source 常用配置

| 参数                  | 默认值   | 描述                                                                                         |
|---------------------|-------|--------------------------------------------------------------------------------------------|
| channels            | -     | 与之相连的channel，可以配置多个。                                                                       |
| type                | avro  | avro source的类型，必须为avro。                                                                    |
| bind                | -     | 监测主机名/IP。                                                                                  |
| port                | -     | 绑定监测端口，该端口需未被占用。                                                                           |
| threads             | -     | source工作的最大线程数。                                                                            |
| compression-type    | none  | 消息压缩格式：“none”或“deflate”。“none”表示不压缩，“deflate”表示压缩。                                         |
| compression-level   | 6     | 数据压缩级别（1-9），数值越高，压缩率越高。                                                                    |
| ssl                 | false | 是否使用SSL加密。设置为true时还必须指定“密钥(keystore)”和“密钥存储密码(keystore-password)”。                         |
| truststore-type     | JKS   | Java信任库类型，“JKS”或“PKCS12”。<br><b>说明</b><br>JKS的密钥库和私钥采用不同的密码进行保护，而PKCS12的密钥库和私钥采用相同密码进行保护。  |
| truststore          | -     | Java信任库文件。                                                                                 |
| truststore-password | -     | Java信任库密码。                                                                                 |
| keystore-type       | JKS   | ssl启用后密钥存储类型，“JKS”或“PKCS12”。<br><b>说明</b><br>JKS的密钥库和私钥用不同的密码进行保护，而PKCS12的密钥库和私钥用相同密码进行保护。 |

| 参数                | 默认值   | 描述                                                                                                                                            |
|-------------------|-------|-----------------------------------------------------------------------------------------------------------------------------------------------|
| keystore          | -     | ssl启用后密钥存储文件路径，开启ssl后，该参数必填。                                                                                                                  |
| keystore-password | -     | ssl启用后密钥存储密码，开启ssl后，该参数必填。                                                                                                                    |
| trust-all-certs   | false | 是否关闭SSL server证书检查。设置为“true”时将不会检查远端source的SSL server证书，不建议在生产中使用。                                                                            |
| exclude-protocols | SSLv3 | 排除的协议列表，用空格分开。默认排除SSLv3协议。                                                                                                                    |
| ipFilter          | false | 是否开启ip过滤。                                                                                                                                     |
| ipFilter.rules    | -     | 定义N网络的ipFilters，多个主机或IP地址用逗号分隔。ipFilter设置为“true”时，配置规则有允许和禁止两种，配置格式如下：<br>ipFilterRules=allow:ip:127.*,<br>allow:name:localhost,<br>deny:ip:* |

- **SpoolDir Source**

Spool Dir Source监控并传输目录下新增的文件，可实现实时数据传输。常用配置如下表所示：

表 7-34 Spooling Directory Source 常用配置

| 参数         | 默认值        | 描述                                               |
|------------|------------|--------------------------------------------------|
| channels   | -          | 与之相连的channel，可以配置多个。                             |
| type       | spooldir   | spooling source的类型，必须设置为spooldir。                |
| spoolDir   | -          | Spooldir source的监控目录，flume运行用户需要对该目录具有可读可写可执行权限。 |
| monTime    | 0（不开启）     | 线程监控阈值，更新时间超过阈值后，重新启动该Source，单位：秒。               |
| fileSuffix | .COMPLETED | 文件传输完成后添加的后缀。                                    |

| 参数                         | 默认值         | 描述                                                                                                                                                                   |
|----------------------------|-------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| deletePolicy               | never       | 文件传输完成后源文件删除策略，never或immediate。“never”表示不删除已完成传输的源文件，“immediate”表示传输完成后立刻删除源文件。                                                                                      |
| ignorePattern              | ^\$         | 忽略文件的正则表达式表示。默认为“^\$”，表示忽略空格。                                                                                                                                        |
| includePattern             | ^.*\$       | 包含文件的正则表达式表示。可以与ignorePattern同时使用，如果一个文件既满足ignorePattern也满足includePattern，则该文件会被忽略。另外，以“.”开头的文件不会被过滤。                                                                |
| trackerDir                 | .flumespool | 传输过程中元数据存储路径。                                                                                                                                                        |
| batchSize                  | 1000        | 批次写入Channel的Event数量。                                                                                                                                                 |
| decodeErrorPolicy          | FAIL        | 编码错误策略。<br><b>说明</b><br>如果文件中有编码错误，请配置“decodeErrorPolicy”为“REPLACE”或“IGNORE”，Flume遇到编码错误将跳过编码错误，继续采集后续日志。                                                            |
| deserializer               | LINE        | 文件解析器，值为“LINE”或“BufferedLine”。<br><ul style="list-style-type: none"> <li>配置为“LINE”时，对从文件读取的字符逐个转码。</li> <li>配置为“BufferedLine”时，对文件读取的一行或多行的字符进行批量转码，性能更优。</li> </ul> |
| deserializer.maxLineLength | 2048        | 按行解析最大长度。                                                                                                                                                            |
| deserializer.maxBatchLine  | 1           | 按行解析最多行数，如果行数设置为多行，maxLineLength也应该设置为相应的倍数。<br><b>说明</b><br>用户设置Interceptor时，需要考虑多行合并后的场景，否则会造成数据丢失。如果Interceptor无法处理多行合并场景，请将该配置设置为1。                              |
| selector.type              | replicating | 选择器类型，“replicating”或“multiplexing”。“replicating”表示将数据复制多份，分别传递给每一个channel，每个channel接收到的数据都是相同的，而“multiplexing”表示根据event中header的value来选择特定的channel，每个channel中的数据是不同的。 |
| interceptors               | -           | 拦截器。多个拦截器用空格分开。                                                                                                                                                      |

| 参数                       | 默认值    | 描述                                                                                                                                                                                                                 |
|--------------------------|--------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| inputCharset             | UTF-8  | 读取文件的编码格式。须与读取数据源文件编码格式相同，否则字符解析可能会出错。                                                                                                                                                                             |
| fileHeader               | false  | 是否把文件名（包含路径）添加到event的header中。                                                                                                                                                                                      |
| fileHeaderKey            | -      | 设置header中数据存储结构为<key,value>模式，需要fileHeaderKey与fileHeader配合使用。若fileHeader设置为true，可参考如下示例。<br>示例：将fileHeaderKey定义为file，当读取到文件名为/root/a.txt的内容时，header中以file=/root/a.txt的形式存在。                                        |
| basenameHeader           | false  | 是否把文件名（不包含路径）添加到event的header中。                                                                                                                                                                                     |
| basenameHeaderKey        | -      | 设置header中数据存储结构为<key,value>模式，需要basenameHeaderKey与basenameHeader配合使用。若basenameHeader设置为true，可参考如下示例。<br>示例：将basenameHeaderKey定义为file，当读取到文件名为a.txt的内容时，header中以file=a.txt的形式存在。                                    |
| pollDelay                | 500    | 轮询监控目录下新文件时的时延。单位：毫秒。                                                                                                                                                                                              |
| recursiveDirectorySearch | false  | 是否监控配置的目录下子目录中的新文件。                                                                                                                                                                                                |
| consumeOrder             | oldest | 监控目录下文件的消耗次序。如果配置为oldest或者youngest，会根据监控目录下文件的最后修改时间来决定，当目录下有大量文件时，会消耗较长时间去寻找oldest或者youngest的文件。需要注意的是，如果配置为random，创建比较早的文件有可能长时间未被读取。如果配置为oldest或者youngest，那么进程会需要较多时间来查找最新的或最旧的文件。可选值：random, youngest, oldest。 |
| maxBackoff               | 4000   | 当Channel满了以后，尝试再次去写Channel所等待的最大时间。超过这个时间，则会发生异常。对应的Source会以一个较小的时间开始，然后每尝试一次，该时间数字指数增长直到达到当前指定的值，如果还不能成功写入，则认为失败。时间单位：秒。                                                                                          |

| 参数             | 默认值  | 描述                                                                                                                                                  |
|----------------|------|-----------------------------------------------------------------------------------------------------------------------------------------------------|
| emptyFileEvent | true | 是否采集空文件信息发送到Sink端，默认值为true，表示将空文件信息发送到Sink端。该参数只对HDFS Sink有效，其他Sink该参数无效。以HDFS Sink为例，当参数为true时，如果spoolDir路径下存在空文件，那么HDFS的hdfs.path路径下就会创建一个同名的空文件。 |

### 说明

SpoolDir Source在按行读取过程中会忽略掉每一个event的最后一个换行符，该换行符所占用的数据量指标不会被Flume统计。

- **Kafka Source**

Kafka Source从Kafka的topic中消费数据，可以设置多个Source消费同一个topic的数据，每个Source会消费topic的不同partitions。常用配置如下表所示：

表 7-35 Kafka Source 常用配置

| 参数                      | 默认值                                       | 描述                                                                                                                                                                  |
|-------------------------|-------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| channels                | -                                         | 与之相连的channel，可以配置多个。                                                                                                                                                |
| type                    | org.apache.flume.source.kafka.KafkaSource | kafka source的类型，必须设置为org.apache.flume.source.kafka.KafkaSource。                                                                                                     |
| kafka.bootstrap.servers | -                                         | Kafka的bootstrap地址端口列表。如果集群已安装Kafka并且配置已经同步，服务端可以不配置此项，默认值为Kafka集群中所有的broker列表。客户端必须配置该项，多个值用逗号分隔。端口和安全协议的匹配规则必须为：21007匹配安全模式（SASL_PLAINTEXT），9092匹配普通模式（PLAINTEXT）。 |
| kafka.topics            | -                                         | 订阅的Kafka topic列表，用逗号分隔。                                                                                                                                             |
| kafka.topics.regex      | -                                         | 符合正则表达式的topic会被订阅，优先级高于“kafka.topics”，如果存在将覆盖“kafka.topics”。                                                                                                        |
| monTime                 | 0（不开启）                                    | 线程监控阈值，更新时间超过阈值后，重新启动该Source，单位：秒。                                                                                                                                  |
| nodatotime              | 0（不开启）                                    | 告警阈值，从Kafka中订阅不到数据的时长超过阈值时发送告警，单位：秒。该参数可在配置文件properties.properties进行设置。                                                                                             |



| 参数                              | 默认值            | 描述                                                                                                                                                           |
|---------------------------------|----------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------|
| batchSize                       | 1000           | 批次写入Channel的Event数量。                                                                                                                                         |
| batchDuration<br>Millis         | 1000           | 批次消费topic数据的最大时长，单位：ms。                                                                                                                                      |
| keepTopicInHeader               | false          | 是否在Event Header中保存topic。设置为true，则Kafka Sink配置的topic将无效。                                                                                                      |
| setTopicHeader                  | true           | 当设置为true时，会将“topicHeader”中定义的topic名称存储到Header中。                                                                                                              |
| topicHeader                     | topic          | 当setTopicHeader属性设置为true，此参数用于定义存储接收的topic名称。如果与Kafka Sink的topicHeader属性结合使用，应该注意，避免将消息循环发送到同一主题。                                                            |
| useFlumeEventFormat             | false          | 默认情况下，event会以字节的形式从kafka topic传递到event的body体中。设置为true，则会以Flume的Avro二进制格式来读取Event。与KafkaSink或KafkaChannel 中同名的parseAsFlumeEvent参数一起使用时，会保留从数据源产生的任何设定的Header。 |
| keepPartitionInHeader           | false          | 是否在Event Header中保存partitionID。设置为true，则Kafka Sink将写入对应的Partition。                                                                                            |
| kafka.consumer.group.id         | flume          | Kafka消费组ID。多个源或代理中设置相同的ID表示它们是同一个consumer group。                                                                                                             |
| kafka.security.protocol         | SASL_PLAINTEXT | Kafka安全协议，普通模式集群下须配置为“PLAINTEXT”。端口和安全协议的匹配规则必须为：21007匹配安全模式（SASL_PLAINTEXT），9092匹配普通模式（PLAINTEXT）。                                                          |
| Other Kafka Consumer Properties | -              | 其他Kafka配置，可以接受任意Kafka支持的消费配置，配置需要加前缀“kafka.”。                                                                                                                |

- **Taildir Source**

Taildir Source监控目录下文件的变化并自动读取文件内容，可实现实时数据传输，常用配置如下表所示：

表 7-36 Taildir Source 常用配置

| 参数                                     | 默认值            | 描述                                                                                                     |
|----------------------------------------|----------------|--------------------------------------------------------------------------------------------------------|
| channels                               | -              | 与之相连的channel，可以配置多个。                                                                                   |
| type                                   | TAILDIR        | taildir source的类型，必须为TAILDIR。                                                                          |
| filegroups                             | -              | 设置采集文件目录分组名字，分组名字中间使用空格间隔。                                                                             |
| filegroups.<filegroupName>.parentDir   | -              | 父目录，需要配置为绝对路径。                                                                                         |
| filegroups.<filegroupName>.filePattern | -              | 相对父目录的文件路径，可以包含目录，支持正则表达式，须与父目录联合使用。                                                                   |
| positionFile                           | -              | 传输过程中元数据存储路径。                                                                                          |
| headers.<filegroupName>.<headerKey>    | -              | 设置某一个分组采集数据时event中的key-value值。                                                                         |
| byteOffsetHeader                       | false          | 是否在每一个event头中携带该event在源文件中的位置信息。设置为true，则该信息保存在byteoffset变量中。                                          |
| maxBatchCount                          | Long.MAX_VALUE | 控制从一个文件中连续读取的最大批次。如果监控目录会一直读取多个文件，且其中一个文件以非常快的速率在写入，那么其他文件可能会无法处理。因为高速写入的这个文件会陷入无限读取的循环中。这种情况下，应该降低此值。 |
| skipToEnd                              | false          | Flume在重启后是否直接定位到文件最新的位置处读取最新的数据。设置为true，则重启后直接定位到文件最新位置读取最新数据。                                         |
| idleTimeout                            | 120000         | 设置读取文件的空闲时间，单位：毫秒，如果在该时间内文件内容没有变更，关闭掉该文件，关闭后如果该文件有数据写入，重新打开并读取数据。                                      |
| writePosInterval                       | 3000           | 设置将元数据写入到文件的周期，单位：毫秒。                                                                                  |
| batchSize                              | 1000           | 批次写入Channel的Event数量。                                                                                   |
| monTime                                | 0（不开启）         | 线程监控阈值，更新时间超过阈值后，重新启动该Source，单位：秒。                                                                     |
| fileHeader                             | false          | 是否把文件名（包含路径）添加到event的header中。                                                                          |

| 参数            | 默认值  | 描述                                                                                                                                                                          |
|---------------|------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| fileHeaderKey | file | 设置header中数据存储结构为<key,value>模式，需要fileHeaderKey与fileHeader配合使用。若fileHeader设置为true，可参考如下示例。<br>示例：将fileHeaderKey定义为file，当读取到文件名为/root/a.txt的内容时，header中以file=/root/a.txt的形式存在。 |

- **Http Source**

Http Source接收外部HTTP客户端发送过来的数据，并放入配置的Channel中，常用配置如下表所示：

表 7-37 Http Source 常用配置

| 参数                    | 默认值                                      | 描述                                                                                                                               |
|-----------------------|------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------|
| channels              | -                                        | 与之相连的channel，可以配置多个。                                                                                                             |
| type                  | http                                     | http source的类型，必须为http。                                                                                                          |
| bind                  | -                                        | 监测主机名/IP。                                                                                                                        |
| port                  | -                                        | 绑定监测端口，该端口需未被占用。                                                                                                                 |
| handler               | org.apache.flume.source.http.JSONHandler | http请求的消息解析方式，支持Json格式解析（org.apache.flume.source.http.JSONHandler）和二进制Blob块解析（org.apache.flume.sink.solr.morphline.BlobHandler）。 |
| handler.*             | -                                        | 设置handler的参数。                                                                                                                    |
| exclude-protocols     | SSLv3                                    | 排除的协议列表，用空格分开。默认排除SSLv3协议。                                                                                                       |
| include-cipher-suites | -                                        | 包含的协议列表，用空格分开。如果设置为空，则默认支持所有协议。                                                                                                  |
| enableSSL             | false                                    | http协议是否启用SSL。设置为true时还必须指定“密钥(keystore)”和“密钥存储密码(keystore-password)”。                                                           |
| keystore-type         | JKS                                      | Keystore类型，可以为JKS或者PKCS12。                                                                                                       |
| keystore              | -                                        | http启用SSL后设置keystore的路径。                                                                                                         |

| 参数               | 默认值 | 描述                       |
|------------------|-----|--------------------------|
| keystorePassword | -   | http启用SSL后设置keystore的密码。 |

- **Thrift Source**

Thrift Source监测thrift端口，接收外部Thrift客户端数据并放入配置的Channel中。常用配置如下表所示：

| 参数                | 默认值    | 描述                                                                                                                                                                                                                                                                             |
|-------------------|--------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| channels          | -      | 与之相连的channel，可以配置多个。                                                                                                                                                                                                                                                           |
| type              | thrift | thrift source的类型，必须设置为thrift。                                                                                                                                                                                                                                                  |
| bind              | -      | 监测主机名/IP。                                                                                                                                                                                                                                                                      |
| port              | -      | 绑定监测端口，该端口需未被占用。                                                                                                                                                                                                                                                               |
| threads           | -      | 允许运行的最大的worker线程数目。                                                                                                                                                                                                                                                            |
| kerberos          | false  | 是否启用Kerberos认证。                                                                                                                                                                                                                                                                |
| agent-keytab      | -      | 服务端使用的keytab文件地址，必须使用机机账号。建议使用Flume服务安装目录下flume/conf/flume_server.keytab。                                                                                                                                                                                                      |
| agent-principal   | -      | 服务端使用的安全用户的Principal，必须使用机机账户。建议使用Flume服务默认用户flume_server/hadoop.<系统域名>@<系统域名><br><b>说明</b><br>“flume_server/hadoop.<系统域名>”为用户名，用户的用户名所包含的系统域名所有字母为小写。例如“本端域”参数为“9427068F-6EFA-4833-B43E-60CB641E5B6C.COM”，用户名为“flume_server/hadoop.9427068f-6efa-4833-b43e-60cb641e5b6c.com”。 |
| compression-type  | none   | 消息压缩格式：“none”或“deflate”。“none”表示不压缩，“deflate”表示压缩。                                                                                                                                                                                                                             |
| ssl               | false  | 是否使用SSL加密。设置为true时还必须指定“密钥(keystore)”和“密钥存储密码(keystore-password)”。                                                                                                                                                                                                             |
| keystore-type     | JKS    | SSL启用后密钥存储类型。                                                                                                                                                                                                                                                                  |
| keystore          | -      | SSL启用后密钥存储文件路径，开启SSL后，该参数必填。                                                                                                                                                                                                                                                   |
| keystore-password | -      | SSL启用后密钥存储密码，开启ssl后，该参数必填。                                                                                                                                                                                                                                                     |

## 常用 Channel 配置

- **Memory Channel**

Memory Channel使用内存作为缓存区，Events存放在内存队列中。常用配置如下表所示：

表 7-38 Memory Channel 常用配置

| 参数                           | 默认值         | 描述                                                                                                                                                  |
|------------------------------|-------------|-----------------------------------------------------------------------------------------------------------------------------------------------------|
| type                         | -           | memory channel的类型，必须设置为memory。                                                                                                                      |
| capacity                     | 10000       | 缓存在channel中的最大Event数。                                                                                                                               |
| transactionCapacity          | 1000        | 每次存取的最大Event数。<br><b>说明</b> <ul style="list-style-type: none"> <li>• 此参数值需要大于source和sink的batchSize。</li> <li>• 事务缓存容量必须小于或等于Channel缓存容量。</li> </ul> |
| channelFullcount             | 10          | channel full次数，达到该次数后发送告警。                                                                                                                          |
| keep-alive                   | 3           | 当事务缓存或Channel缓存满时，Put、Take线程等待时间。单位：秒。                                                                                                              |
| byteCapacity                 | JVM最大内存的80% | channel中最多能容纳所有event body的总字节数，默认是JVM最大可用内存（-Xmx）的80%，单位：bytes。                                                                                     |
| byteCapacityBufferPercentage | 20          | channel中字节容量百分比（%）。                                                                                                                                 |

- **File Channel**

File Channel使用本地磁盘作为缓存区，Events存放在设置的dataDirs配置项文件夹中。常用配置如下表所示：

表 7-39 File Channel 常用配置

| 参数   | 默认值 | 描述                         |
|------|-----|----------------------------|
| type | -   | file channel的类型，必须设置为file。 |

| 参数                   | 默认值                                                                                              | 描述                                                                                                                                              |
|----------------------|--------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------|
| checkpointDir        | \${BIGDATA_DATA_HOME}/<br>hadoop/data1~N/flume/<br>checkpoint<br><br><b>说明</b><br>此路径随自定义数据路径变更。 | 检查点存放路径。                                                                                                                                        |
| dataDirs             | \${BIGDATA_DATA_HOME}/<br>hadoop/data1~N/flume/data<br><br><b>说明</b><br>此路径随自定义数据路径变更。           | 数据缓存路径，设置多个路径可提升性能，中间用逗号分开。                                                                                                                     |
| maxFileSize          | 2146435071                                                                                       | 单个缓存文件的最大值，单位：bytes。                                                                                                                            |
| minimumRequiredSpace | 524288000                                                                                        | 缓冲区空闲空间最小值，单位：bytes。                                                                                                                            |
| capacity             | 1000000                                                                                          | 缓存在channel中的最大Event数。                                                                                                                           |
| transactionCapacity  | 10000                                                                                            | 每次存取的最大Event数。<br><b>说明</b> <ul style="list-style-type: none"> <li>此参数值需要大于source和sink的batchSize。</li> <li>事务缓存容量必须小于或等于Channel缓存容量。</li> </ul> |
| channelFullCount     | 10                                                                                               | channel full次数，达到该次数后发送告警。                                                                                                                      |
| useDualCheckpoints   | false                                                                                            | 是否备份检查点。设置为“true”时，必须设置backupCheckpointDir的参数值。                                                                                                 |
| backupCheckpointDir  | -                                                                                                | 备份检查点路径。                                                                                                                                        |
| checkpointInterval   | 30000                                                                                            | 检查点间隔时间，单位：秒。                                                                                                                                   |
| keep-alive           | 3                                                                                                | 当事务缓存或Channel缓存满时，Put、Take线程等待时间。单位：秒。                                                                                                          |
| use-log-replay-v1    | false                                                                                            | 是否启用旧的回复逻辑。                                                                                                                                     |
| use-fast-replay      | false                                                                                            | 是否使用队列回复。                                                                                                                                       |

| 参数                | 默认值  | 描述                 |
|-------------------|------|--------------------|
| checkpointOnClose | true | channel关闭时是否创建检查点。 |

- **Memory File Channel**

Memory File Channel同时使用内存和本地磁盘作为缓存区，消息可持久化，性能优于File Channel，接近Memory Channel的性能。此Channel目前处于试验阶段，可靠性不够高，不建议在生产环境使用。常用配置如下表所示：

表 7-40 Memory File Channel 常用配置

| 参数                   | 默认值                                        | 描述                                                                                                                                                                                                                                                                                                                                                                                                            |
|----------------------|--------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| type                 | org.apache.flume.channel.MemoryFileChannel | memory file channel的类型，必须设置为“org.apache.flume.channel.MemoryFileChannel”。                                                                                                                                                                                                                                                                                                                                     |
| capacity             | 50000                                      | Channel缓存容量：缓存在Channel中的最大Event数。                                                                                                                                                                                                                                                                                                                                                                             |
| transactionCapacity  | 5000                                       | 事务缓存容量：一次事务能处理的最大Event数。<br><b>说明</b> <ul style="list-style-type: none"> <li>• 此参数值需要大于source和sink的batchSize。</li> <li>• 事务缓存容量必须小于或等于Channel缓存容量。</li> </ul>                                                                                                                                                                                                                                                 |
| subqueueByteCapacity | 20971520                                   | 每个subqueue最多保存多少byte的Event，单位：byte。<br>Memory File Channel采用queue和subqueue两级缓存，event保存在subqueue，subqueue保存在queue。<br>subqueue能保存多少event，由“subqueueCapacity”和“subqueueInterval”两个参数决定，“subqueueCapacity”限制subqueue内的Event总容量，“subqueueInterval”限制subqueue保存Event的时长，只有subqueue达到“subqueueCapacity”或“subqueueInterval”上限时，subqueue内的Event才会发往目的地。<br><b>说明</b> “subqueueByteCapacity”必须大于一个batchsize内的Event总容量。 |
| subqueueInterval     | 2000                                       | 每个subqueue最多保存一段多长时间的Event，单位：毫秒。                                                                                                                                                                                                                                                                                                                                                                             |

| 参数               | 默认值         | 描述                                                 |
|------------------|-------------|----------------------------------------------------|
| keep-alive       | 3           | 当事务缓存或Channel缓存满时，Put、Take线程等待时间。<br>单位：秒。         |
| dataDir          | -           | 缓存本地文件存储目录。                                        |
| byteCapacity     | JVM最大内存的80% | Channel缓存容量。<br>单位：bytes。                          |
| compression-type | None        | 消息压缩格式：“none”或“deflate”。“none”表示不压缩，“deflate”表示压缩。 |
| channelFullCount | 10          | channel full次数，达到该次数后发送告警。                         |

Memory File Channel配置样例：

```
server.channels.c1.type = org.apache.flume.channel.MemoryFileChannel
server.channels.c1.dataDir = /opt/flume/mfdata
server.channels.c1.subqueueByteCapacity = 20971520
server.channels.c1.subqueueInterval=2000
server.channels.c1.capacity = 500000
server.channels.c1.transactionCapacity = 40000
```

- **Kafka Channel**

Kafka Channel使用Kafka集群缓存数据，Kafka提供高可用、多副本，以防Flume或Kafka Broker崩溃，Channel中的数据会立即被Sink消费。

表 7-41 Kafka channel 常用配置

| Parameter | Default Value | Description                                                          |
|-----------|---------------|----------------------------------------------------------------------|
| type      | -             | kafka channel的类型，必须设置为“org.apache.flume.channel.kafka.KafkaChannel”。 |



| Parameter                        | Default Value  | Description                                                                                                                                                                      |
|----------------------------------|----------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| kafka.bootstrap.servers          | -              | Kafka的bootstrap地址<br>端口列表。<br>如果集群已安装Kafka并且配置已经同步，则服务端可以不配置此项，默认值为Kafka集群中所有的broker列表。客户端必须配置该项，多个值用逗号分隔。端口和安全协议的匹配规则必须为：21007匹配安全模式（ SASL_PLAINTEXT ），9092匹配普通模式（ PLAINTEXT ）。 |
| kafka.topic                      | flume-channel  | channel用来缓存数据的topic。                                                                                                                                                             |
| kafka.consumer.group.id          | flume          | 从kafka中获取数据的组标识，此参数不能为空。                                                                                                                                                         |
| parseAsFlumeEvent                | true           | 是否解析为Flume event。                                                                                                                                                                |
| migrateZookeeperOffsets          | true           | 当Kafka没有存储offset时，是否从ZooKeeper中查找，并提交到Kafka。                                                                                                                                     |
| kafka.consumer.auto.offset.reset | latest         | 当没有offset记录时从什么位置消费，可选为“earliest”、“latest”或“none”。“earliest”表示将offset重置为初始点，“latest”表示将offset置为最新位置点，“none”表示若没有offset则发生异常。                                                     |
| kafka.producer.security.protocol | SASL_PLAINTEXT | Kafka生产安全协议。端口和安全协议的匹配规则必须为：21007匹配安全模式（ SASL_PLAINTEXT ），9092匹配普通模式（ PLAINTEXT ）。<br><b>说明</b><br>若该参数没有显示，请单击弹窗左下角的“+”显示全部参数。                                                  |

| Parameter                        | Default Value  | Description                                                                 |
|----------------------------------|----------------|-----------------------------------------------------------------------------|
| kafka.consumer.security.protocol | SASL_PLAINTEXT | 同上，但用于消费。端口和安全协议的匹配规则必须为：21007匹配安全模式（SASL_PLAINTEXT），9092匹配普通模式（PLAINTEXT）。 |
| pollTimeout                      | 500            | consumer调用poll()函数能接受的最大超时时间，单位：毫秒。                                         |
| ignoreLongMessage                | false          | 是否丢弃超大消息。                                                                   |
| messageMaxLength                 | 1000012        | Flume写入Kafka的消息的最大长度。                                                       |

## 常用 Sink 配置

- **HDFS Sink**

HDFS Sink将数据写入Hadoop分布式文件系统（HDFS）。常用配置如下表所示：

表 7-42 HDFS Sink 常用配置

| 参数                | 默认值    | 描述                                                         |
|-------------------|--------|------------------------------------------------------------|
| channel           | -      | 与之相连的channel。                                              |
| type              | hdfs   | hdfs sink的类型，必须设置为hdfs。                                    |
| hdfs.path         | -      | HDFS上数据存储路径，必须以“hdfs://hacluster/”开头。                      |
| monTime           | 0（不开启） | 线程监控阈值，更新时间超过阈值后，重新启动该Sink，单位：秒。                           |
| hdfs.inUseSuffix  | .tmp   | 正在写入的hdfs文件后缀。                                             |
| hdfs.rollInterval | 30     | 按时间滚动文件，单位：秒，同时需将“hdfs.fileCloseByEndEvent”设置为“false”。     |
| hdfs.rollSize     | 1024   | 按大小滚动文件，单位：bytes，同时需将“hdfs.fileCloseByEndEvent”设置为“false”。 |

| 参数                       | 默认值                          | 描述                                                                                                                                                                                                                                         |
|--------------------------|------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| hdfs.rollCount           | 10                           | 按Event个数滚动文件，同时需将“hdfs.fileCloseByEndEvent”设置为“false”。<br><b>说明</b><br>参数“rollInterval”、“rollSize”和“rollCount”可同时配置，三个参数采取优先原则，哪个参数值先满足，优先按照哪个参数进行压缩。                                                                                      |
| hdfs.idleTimeout         | 0                            | 自动关闭空闲文件超时时间，单位：秒。                                                                                                                                                                                                                         |
| hdfs.batchSize           | 1000                         | 批次写入HDFS的Event个数。                                                                                                                                                                                                                          |
| hdfs.kerberosPrincipal   | -                            | 认证HDFS的Kerberos principal，普通模式集群不配置，安全模式集群必须配置。                                                                                                                                                                                            |
| hdfs.kerberosKeytab      | -                            | 认证HDFS的Kerberos keytab，普通模式集群不配置，安全模式集群中，用户必须对jaas.conf文件中的keyTab路径有访问权限。                                                                                                                                                                  |
| hdfs.fileCloseByEndEvent | true                         | 收到源文件的最后一个Event时是否关闭hdfs文件。                                                                                                                                                                                                                |
| hdfs.batchCallTimeout    | -                            | 批次写入HDFS超时控制时间，单位：毫秒。<br>当不配置此参数时，对每个Event写入HDFS进行超时控制。当“hdfs.batchSize”大于0时，配置此参数可以提升写入HDFS性能。<br><b>说明</b><br>“hdfs.batchCallTimeout”设置多长时间需要考虑“hdfs.batchSize”的大小，“hdfs.batchSize”越大，“hdfs.batchCallTimeout”也要调整更长时间，设置过短时间容易导致写HDFS失败。 |
| serializer.appendNewline | true                         | 将一个Event写入HDFS后是否追加换行符（'\n'），如果追加该换行符，该换行符所占用的数据量指标不会被HDFS Sink统计。                                                                                                                                                                         |
| hdfs.filePrefix          | over_<br>%<br>{base<br>name} | 数据写入hdfs后文件名的前缀。                                                                                                                                                                                                                           |
| hdfs.fileSuffix          | -                            | 数据写入hdfs后文件名的后缀。                                                                                                                                                                                                                           |
| hdfs.inUsePrefix         | -                            | 正在写入的hdfs文件前缀。                                                                                                                                                                                                                             |

| 参数                     | 默认值        | 描述                                                                                                                                                                        |
|------------------------|------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| hdfs.fileType          | DataStream | hdfs文件格式，包括“SequenceFile”、“DataStream”以及“CompressedStream”。<br><b>说明</b><br>“SequenceFile”和“DataStream”不压缩输出文件，不能设置参数“codeC”，“CompressedStream”压缩输出文件，必须设置“codeC”参数值配合使用。 |
| hdfs.codeC             | -          | 文件压缩格式，包括gzip、bzip2、lzo、lzop、snappy。                                                                                                                                      |
| hdfs.maxOpenFiles      | 5000       | 最大允许打开的hdfs文件数，当打开的文件数达到该值时，最早打开的文件将会被关闭。                                                                                                                                 |
| hdfs.writeFormat       | Writable   | 文件写入格式，“Writable”或者“Text”。                                                                                                                                                |
| hdfs.callTimeout       | 10000      | 写入HDFS超时控制时间，单位：毫秒。                                                                                                                                                       |
| hdfs.threadsPoolSize   | -          | 每个HDFS sink用于HDFS io操作的线程数。                                                                                                                                               |
| hdfs.rollTimerPoolSize | -          | 每个HDFS sink用于调度定时文件滚动的线程数。                                                                                                                                                |
| hdfs.round             | false      | 时间戳是否四舍五入。若设置为true，则会影响所有基于时间的转义序列（%t除外）。                                                                                                                                 |
| hdfs.roundUnit         | second     | 时间戳四舍五入单位，可选为“second”、“minute”或“hour”，分别对应为秒、分钟和小时。                                                                                                                       |
| hdfs.useLocalTimeStamp | true       | 是否启用本地时间戳，建议设置为“true”。                                                                                                                                                    |
| hdfs.closeTries        | 0          | hdfs sink尝试关闭重命名文件的最大次数。默认为0表示sink会一直尝试重命名，直至重命名成功。                                                                                                                       |
| hdfs.retryInterval     | 180        | 尝试关闭hdfs文件的时间间隔，单位：秒。<br><b>说明</b><br>每个关闭请求都会有多个RPC往返Namenode，因此设置的太低可能导致Namenode超负荷。如果设置0，如果第一次尝试失败的话，该Sink将不会尝试关闭文件，并且把文件打开，或者用“.tmp”作为扩展名。                            |
| hdfs.failcount         | 10         | 数据写入hdfs失败的次数。该参数作为sink写入hdfs失败次数的阈值，当超过该阈值后上报数据传输异常告警。                                                                                                                   |

- **Avro Sink**

Avro Sink把events转化为Avro events并发送到配置的主机的监测端口。常用配置如下表所示：

**表 7-43 Avro Sink 常用配置**

| 参数          | 默认值     | 描述                                                                                                                                                                                                                                                                                                       |
|-------------|---------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| channel     | -       | 与之相连的channel。                                                                                                                                                                                                                                                                                            |
| type        | -       | avro sink的类型，必须设置为avro。                                                                                                                                                                                                                                                                                  |
| hostname    | -       | 绑定的主机名/IP。                                                                                                                                                                                                                                                                                               |
| port        | -       | 监测端口，该端口需未被占用。                                                                                                                                                                                                                                                                                           |
| batch-size  | 1000    | 批次发送的Event个数。                                                                                                                                                                                                                                                                                            |
| client.type | DEFAULT | <p>客户端实例类型，根据所配置的模型实际使用到的通信协议设置。该值可选值包括：</p> <ul style="list-style-type: none"> <li>• DEFAULT，返回AvroRPC类型的客户端实例。</li> <li>• OTHER，返回NULL。</li> <li>• THRIFT，返回Thrift RPC类型的客户端实例。</li> <li>• DEFAULT_LOADBALANCING，返回LoadBalancing RPC客户端实例。</li> <li>• DEFAULT_FAILOVER，返回Failover RPC客户端实例。</li> </ul> |
| ssl         | false   | 是否使用SSL加密。设置为true时还必须指定“密钥(keystore)”和“密钥存储密码(keystore-password)”。                                                                                                                                                                                                                                       |

| 参数                        | 默认值   | 描述                                                                                        |
|---------------------------|-------|-------------------------------------------------------------------------------------------|
| truststore-type           | JKS   | Java信任库类型，“JKS”或“PKCS12”。<br><b>说明</b><br>JKS的密钥库和私钥采用不同的密码进行保护，而PKCS12的密钥库和私钥采用相同密码进行保护。 |
| truststore                | -     | Java信任库文件。                                                                                |
| truststore-password       | -     | Java信任库密码。                                                                                |
| keystore-type             | JKS   | ssl启用后密钥存储类型。                                                                             |
| keystore                  | -     | ssl启用后密钥存储文件路径，开启ssl后，该参数必填。                                                              |
| keystore-password         | -     | ssl启用后密钥存储密码，开启ssl后，该参数必填。                                                                |
| connect-timeout           | 20000 | 第一次连接的超时时间，单位：毫秒。                                                                         |
| request-timeout           | 20000 | 第一次请求后一次请求的最大超时时间，单位：毫秒。                                                                  |
| reset-connection-interval | 0     | 一次断开连接后，等待多少时间后进行重新连接，单位：秒。默认为0表示不断尝试。                                                    |
| compression-type          | none  | 批数据压缩类型，“none”或“deflate”，“none”表示不压缩，“deflate”表示压缩。该值必须与AvroSource的compression-type匹配。    |
| compression-level         | 6     | 批数据压缩级别（1-9），数值越高，压缩率越高。                                                                  |
| exclude-protocols         | SSLv3 | 排除的协议列表，用空格分开。默认排除SSLv3协议。                                                                |

- **HBase Sink**

HBase Sink将数据写入到HBase中。常用配置如下表所示：

表 7-44 HBase Sink 常用配置

| 参数                 | 默认值    | 描述                                                                               |
|--------------------|--------|----------------------------------------------------------------------------------|
| channel            | -      | 与之相连的channel。                                                                    |
| type               | -      | hbase sink的类型，必须设置为hbase。                                                        |
| table              | -      | HBase表名称。                                                                        |
| columnFamily       | -      | HBase列族。                                                                         |
| monTime            | 0（不开启） | 线程监控阈值，更新时间超过阈值后，重新启动该Sink，单位：秒。                                                 |
| batchSize          | 1000   | 批次写入HBase的Event个数。                                                               |
| kerberosPrincipal  | -      | 认证HBase的Kerberos principal，普通模式集群不配置，安全模式集群必须配置。                                 |
| kerberosKeytab     | -      | 认证HBase的Kerberos keytab，普通模式集群不配置，安全模式集群中，flume运行用户必须对jaas.cof文件中的keyTab路径有访问权限。 |
| coalesceIncrements | true   | 是否在同一处理批次中，合并对同一个hbase cell多个操作。设置为true有利于提高性能。                                  |

- **Kafka Sink**

Kafka Sink将数据写入到Kafka中。常用配置如下表所示：

表 7-45 Kafka Sink 常用配置

| 参数                      | 默认值    | 描述                                                                                                                                                                   |
|-------------------------|--------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| channel                 | -      | 与之相连的channel。                                                                                                                                                        |
| type                    | -      | kafka sink的类型，必须设置为org.apache.flume.sink.kafka.KafkaSink。                                                                                                            |
| kafka.bootstrap.servers | -      | Kafka 的bootstrap 地址端口列表。如果集群安装有kafka并且配置已经同步，服务端可以不配置此项，默认值为Kafka集群中所有的broker列表，客户端必须配置该项，多个用逗号分隔。端口和安全协议的匹配规则必须为：21007匹配安全模式（SASL_PLAINTEXT），9092匹配普通模式（PLAINTEXT）。 |
| monTime                 | 0（不开启） | 线程监控阈值，更新时间超过阈值后，重新启动该Sink，单位：秒。                                                                                                                                     |

| 参数                              | 默认值            | 描述                                                                                                                                                                              |
|---------------------------------|----------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| kafka.producer.acks             | 1              | 必须收到多少个replicas的确认信息才认为写入成功。0表示不需要接收确认信息，1表示只等待leader的确认信息。-1表示等待所有的relicas的确认信息。设置为-1，在某些leader失败的场景中可以避免数据丢失。                                                                 |
| kafka.topic                     | -              | 数据写入的topic，必须填写。                                                                                                                                                                |
| flumeBatchSize                  | 1000           | 批次写入Kafka的Event个数。                                                                                                                                                              |
| kafka.security.protocol         | SASL_PLAINTEXT | Kafka安全协议，普通模式集群下须配置为“PLAINTEXT”。端口和安全协议的匹配规则必须为：21007匹配安全模式（SASL_PLAINTEXT），9092匹配普通模式（PLAINTEXT）。                                                                             |
| ignoreLongMessage               | false          | 是否丢弃超大消息的开关。                                                                                                                                                                    |
| messageMaxLength                | 1000012        | Flume写入Kafka的消息的最大长度。                                                                                                                                                           |
| defaultPartitionId              | -              | 用于指定channel中的events被传输到哪一个Kafka partition ID，此值会被partitionIdHeader覆盖。默认情况下，如果此参数不设置，会由Kafka Producer's partitioner 进行events分发(可以通过指定key或者kafka.partition.class自定义的partitioner)。 |
| partitionIdHeader               | -              | 设置时，对应的Sink 将从Event 的Header 中获取使用此属性的值命名的字段的值，并将消息发送到主题的指定分区。如果该值无对应的有效分区，则会发生EventDeliveryException。如果Header 值已经存在，则此设置将覆盖参数defaultPartitionId。                                |
| Other Kafka Producer Properties | -              | 其他Kafka配置，可以接受任意Kafka支持的生产配置，配置需要加前缀.kafka。                                                                                                                                     |

- **Thrift Sink**

Thrift Sink把events转化为Thrift events并发送到配置的主机的监测端口。常用配置如下表所示：

**表 7-46 Thrift Sink 常用配置**

| 参数      | 默认值 | 描述            |
|---------|-----|---------------|
| channel | -   | 与之相连的channel。 |



| 参数                        | 默认值    | 描述                                                         |
|---------------------------|--------|------------------------------------------------------------|
| type                      | thrift | thrift sink的类型，必须设置为thrift。                                |
| hostname                  | -      | 绑定的主机名/IP。                                                 |
| port                      | -      | 监测端口，该端口需未被占用。                                             |
| batch-size                | 1000   | 批次发送的Event个数。                                              |
| connect-timeout           | 20000  | 第一次连接的超时时间，单位：毫秒。                                          |
| request-timeout           | 20000  | 第一次请求后一次请求的最大超时时间，单位：毫秒。                                   |
| kerberos                  | false  | 是否启用Kerberos认证。                                            |
| client-keytab             | -      | 客户端使用的keytab文件地址，flume运行用户必须对认证文件具有访问权限。                   |
| client-principal          | -      | 客户端使用的安全用户的Principal。                                      |
| server-principal          | -      | 服务端使用的安全用户的Principal。                                      |
| compression-type          | none   | Flume发送数据的压缩类型，“none”或“deflate”，“none”表示不压缩，“deflate”表示压缩。 |
| maxConnections            | 5      | Flume发送数据时的最大连接池大小。                                        |
| ssl                       | false  | 是否使用SSL加密。                                                 |
| truststore-type           | JKS    | Java信任库类型。                                                 |
| truststore                | -      | Java信任库文件。                                                 |
| truststore-password       | -      | Java信任库密码。                                                 |
| reset-connection-interval | 0      | 一次断开连接后，等待多少时间后进行重新连接，单位：秒。默认为0表示不断尝试。                     |

## 注意事项

- Flume可靠性保障措施有哪些？
  - Source&Channel、Channel&Sink之间的事务机制。
  - Sink Processor支持配置failover、load\_balance机制，例如负载均衡示例如下，详细参考<http://flume.apache.org/releases/1.9.0.html>。  

```
server.sinkgroups=g1
server.sinkgroups.g1.sinks=k1 k2
server.sinkgroups.g1.processor.type=load_balance
server.sinkgroups.g1.processor.backoff=true
server.sinkgroups.g1.processor.selector=random
```
- Flume多agent聚合级联时的注意事项？
  - 级联时需要使用Avro或者Thrift协议进行级联。
  - 聚合端存在多个节点时，连接配置尽量配置均衡，不要聚合到单节点上。

## 7.8.3 Flume 日志介绍

### 日志描述

**日志路径：**Flume相关日志的默认存储路径为“/var/log/Bigdata/角色名”。

- FlumeServer：“/var/log/Bigdata/flume/flume”
- FlumeClient：“/var/log/Bigdata/flume-client-n/flume”
- MonitorServer：“/var/log/Bigdata/flume/monitor”

**日志归档规则：**Flume日志启动了自动压缩归档功能，缺省情况下，当日志大小超过50MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd\_hh-mm-ss>.[编号].log.zip”。最多保留最近的20个压缩文件，压缩文件保留个数可以在Manager界面中配置。

表 7-47 Flume 日志列表

| 日志类型 | 日志文件名                          | 描述                     |
|------|--------------------------------|------------------------|
| 运行日志 | /flume/flumeServer.log         | FlumeServer运行环境信息日志。   |
|      | /flume/install.log             | FlumeServer安装日志。       |
|      | /flume/flumeServer-gc.log.<编号> | FlumeServer进程的GC归档日志。  |
|      | /flume/prestartDvietail.log    | Flume启动前的工作日志。         |
|      | /flume/startDetail.log         | Flume进程启动工作日志。         |
|      | /flume/stopDetail.log          | Flume进程停止日志。           |
|      | /monitor/monitorServer.log     | MonitorServer运行环境信息日志。 |
|      | /monitor/startDetail.log       | MonitorServer进程启动工作日志。 |
|      | /monitor/stopDetail.log        | MonitorServer进程停止日志。   |

| 日志类型 | 日志文件名                          | 描述               |
|------|--------------------------------|------------------|
|      | function.log                   | 外部函数调用日志。        |
|      | /flume/flume-用户名-日期-pid-gc.log | Flume进程的GC日志。    |
|      | /flume/Flume-audit.log         | Flume客户端的审计日志。   |
|      | /flume/startAgent.out          | Flume启动前的进程参数日志。 |

## 日志级别

Flume提供了如表7-48所示的日志级别。

运行日志的级别优先级从高到低分别是FATAL、ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 7-48 日志级别

| 日志类型 | 级别    | 描述                      |
|------|-------|-------------------------|
| 运行日志 | FATAL | FATAL表示系统运行的致命错误信息。     |
|      | ERROR | ERROR表示系统运行的错误信息。       |
|      | WARN  | WARN表示当前事件处理存在异常信息。     |
|      | INFO  | INFO表示记录系统及各事件正常运行状态信息。 |
|      | DEBUG | DEBUG表示记录系统及系统的调试信息。    |

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 请参考[修改集群服务配置参数](#)，进入Flume的“全部配置”页面。
- 步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤3** 选择所需修改的日志级别。
- 步骤4** 保存配置，在弹出窗口中单击“确定”使配置生效。

----结束

### 说明

配置完成后即生效，不需要重启服务。

## 日志格式

Flume的日志格式如下所示：

表 7-49 日志格式

| 日志类型 | 格式                                                                                                    | 示例                                                                                                                                                            |
|------|-------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 运行日志 | <yyyy-MM-dd<br>HH:mm:ss,SSS> <Log Level><<br>产生该日志的线程名字> <log<br>中的message> <日志事件的发<br>生位置>           | 2014-12-12 11:54:57,316   INFO<br>  [main]   log4j dynamic load is<br>start.  <br>org.apache.flume.tools.LogDyn<br>amicLoad.start(LogDynamicLoa<br>d.java:59) |
|      | <yyyy-MM-dd<br>HH:mm:ss,SSS><User<br>Name><User<br>IP><Time><Operation><Reso<br>urce><Result><Detail> | 2014-12-12 23:04:16,572   INFO<br>  [SinkRunner-PollingRunner-<br>DefaultSinkProcessor]  <br>SRCIP=null OPERATION=close                                       |

### 7.8.4 查看 Flume 客户端日志

**步骤1** 安装Flume客户端。

**步骤2** 进入Flume客户端日志目录，默认为“/var/log/Bigdata”。

**步骤3** 执行如下命令查看日志文件列表。

**ls -lR flume-client-\***

日志文件示例如下：

```
flume-client-1/flume:
total 7672
-rw-----, 1 root root 0 Sep 8 19:43 Flume-audit.log
-rw-----, 1 root root 1562037 Sep 11 06:05 FlumeClient.2017-09-11_04-05-09.[1].log.zip
-rw-----, 1 root root 6127274 Sep 11 14:47 FlumeClient.log
-rw-----, 1 root root 2935 Sep 8 22:20 flume-root-20170908202009-pid72456-gc.log.0.current
-rw-----, 1 root root 2935 Sep 8 22:27 flume-root-20170908202634-pid78789-gc.log.0.current
-rw-----, 1 root root 4382 Sep 8 22:47 flume-root-20170908203137-pid84925-gc.log.0.current
-rw-----, 1 root root 4390 Sep 8 23:46 flume-root-20170908204918-pid103920-gc.log.0.current
-rw-----, 1 root root 3196 Sep 9 10:12 flume-root-20170908215351-pid44372-gc.log.0.current
-rw-----, 1 root root 2935 Sep 9 10:13 flume-root-20170909101233-pid55119-gc.log.0.current
-rw-----, 1 root root 6441 Sep 9 11:10 flume-root-20170909101631-pid59301-gc.log.0.current
-rw-----, 1 root root 0 Sep 9 11:10 flume-root-20170909111009-pid119477-gc.log.0.current
-rw-----, 1 root root 92896 Sep 11 13:24 flume-root-20170909111126-pid120689-gc.log.0.current
-rw-----, 1 root root 5588 Sep 11 14:46 flume-root-20170911132445-pid42259-gc.log.0.current
-rw-----, 1 root root 2576 Sep 11 13:24 prestartDetail.log
-rw-----, 1 root root 3303 Sep 11 13:24 startDetail.log
-rw-----, 1 root root 1253 Sep 11 13:24 stopDetail.log

flume-client-1/monitor:
total 8
-rw-----, 1 root root 141 Sep 8 19:43 flumeMonitorChecker.log
-rw-----, 1 root root 294 Sep 11 13:24 flumeMonitor.log
```

其中**FlumeClient.log**即为Flume客户端的运行日志。

----结束

## 7.8.5 查看 Flume 客户端监控信息

### 操作场景

集群外的Flume客户端也是端到端数据采集的一环，与集群内Flume服务端一起都需要监控，用户通过FusionInsight Manager可以对Flume客户端进行监控，可以查看客户端的Source、Sink、Channel的监控指标以及客户端的进程状态。

本章节适用于MRS 3.x及之后版本。

### 操作步骤

**步骤1** 登录FusionInsight Manager。

**步骤2** 选择“集群 > 待操作集群的名称 > 服务 > Flume > Flume管理”，即可查看当前Flume客户端列表及进程状态。

图 7-60 Flume 管理



**步骤3** 选择“实例ID”，进入客户端监控列表，在“实时”区域框中，可查看客户端的各监控指标。

**步骤4** 选择“历史”进入历史监控数据查询界面。筛选时间段，单击“查看”可显示该时间段内的监控数据。

----结束

## 7.8.6 停止或卸载 Flume 客户端

### 操作场景

指导运维工程师停止、启动Flume客户端，以及在不需要Flume数据采集通道时，卸载Flume客户端。

### 操作步骤

- 停止Flume角色的客户端。  
假设Flume客户端安装路径为“/opt/FlumeClient”，执行以下命令，停止Flume客户端：

```
cd /opt/FlumeClient/fusioninsight-flume-Flume组件版本号/bin
./flume-manage.sh stop
```

执行脚本后，显示如下信息，说明成功的停止了Flume客户端：

```
Stop Flume PID=120689 successful..
```

### 📖 说明

Flume客户端停止后会自动重启，如果不需自动重启，请执行以下命令：

```
./flume-manage.sh stop force
```

需要启动时，可执行以下命令：

```
./flume-manage.sh start force
```

- 卸载Flume角色的客户端。

假设Flume客户端安装路径为“/opt/FlumeClient”，执行以下命令，卸载Flume客户端：

```
cd /opt/FlumeClient/fusioninsight-flume-Flume组件版本号/inst
./uninstall.sh
```

## 7.9 Flume 常见问题

### 7.9.1 如何查看 Flume 日志

Flume日志保存在/var/log/Bigdata/flume/flume/flumeServer.log 里。绝大多数数据传输异常、数据传输不成功，在日志里都可以看到提示。可以直接输入以下命令查看：

```
tailf /var/log/Bigdata/flume/flume/flumeServer.log
```

- 问题：当配置文件上传后，发现异常，重新上传配置文件，发现仍然没有满足场景要求，但日志上没有任何异常。

解决方法：重启此flume进程，**kill -9 进程代码**，再看日志。

- 问题：连接HDFS出现java.lang.IllegalArgumentException: Keytab is not a readable file: /opt/test/conf/user.keytab。

解决方法：添加Flume运行用户读写权限。

- 问题：执行Flume客户端连接Kafka报如下错误：

```
Caused by: java.io.IOException: /opt/FlumeClient/fusioninsight-flume-1.9.0/cof//jaas.conf (No such file or directory)
```

解决方法：新增jaas.conf配置文件并保存到flume client的conf路径下。

#### vi jaas.conf

```
KafkaClient {
com.sun.security.auth.module.Krb5LoginModule required
useKeyTab=true
keyTab="/opt/test/conf/user.keytab"
principal="flume_hdfs@<系统域名>"
useTicketCache=false
storeKey=true
debug=true;
};
```

参数keyTab和principal根据实际情况修改。

- 问题：执行Flume客户端连接HBase报如下错误：

```
Caused by: java.io.IOException: /opt/FlumeClient/fusioninsight-flume-1.9.0/cof//jaas.conf (No such file or directory)
```

解决方法：新增jaas.conf配置文件并保存到flume client的conf路径下。

#### vi jaas.conf

```
Client {
 com.sun.security.auth.module.Krb5LoginModule required
 useKeyTab=true
 keyTab="/opt/test/conf/user.keytab"
 principal="flume_hbase@<系统域名>"
 useTicketCache=false
 storeKey=true
 debug=true;
};
```

参数keyTab和principal根据实际情况修改。

- 问题：一旦提交配置文件后，flume agent即在占用资源运行，如何恢复到没有上传配置文件的状况？

解决方法：提交一个内容为空的properties.properties文件。

## 7.9.2 如何在 Flume 配置文件中 使用环境变量

本章节描述如何在配置文件“properties.properties”中使用环境变量。

本章节适用于MRS 3.x及之后版本。

**步骤1** 安装Flume客户端。

**步骤2** 以root用户登录安装Flume客户端所在节点。

**步骤3** 切换到以下目录。

```
cd Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/conf
```

**步骤4** 在该目录下的“flume-env.sh”文件中添加环境变量。

- 格式：

```
export 变量名=变量值
```

- 示例：

```
JAVA_OPTS="-Xms2G -Xmx4G -XX:CMSFullGCsBeforeCompaction=1 -XX:+UseConcMarkSweepGC -
XX:+CMSParallelRemarkEnabled -XX:+UseCMSCompactAtFullCollection -
DpropertiesImplementation=org.apache.flume.node.EnvVarResolverProperties"
export TAILDIR_PATH=/tmp/flumetest/201907/20190703/1/*.log.*
```

**步骤5** 重启Flume实例进程。

1. 登录FusionInsight Manager。
2. 选择“集群 > 服务 > Flume > 实例”，勾选Flume实例，选择“更多 > 重启实例”输入密码，单击“确定”等待实例重启成功。

### 须知

服务端flume-env.sh生效后不能通过Manager界面重启整个Flume服务，否则用户自定义环境变量丢失，仅需在Manager界面重启对应实例即可。

**步骤6** 在“Flume客户端安装目录/fusioninsight-flume-Flume组件版本号/conf/properties.properties”配置文件中 使用“\${变量名}”格式引用变量，示例如下：

```
client.sources.s1.type = TAILDIR
client.sources.s1.filegroups = f1
client.sources.s1.filegroups.f1 = ${TAILDIR_PATH}
client.sources.s1.positionFile = /tmp/flumetest/201907/20190703/1/taildir_position.json
client.sources.s1.channels = c1
```

**须知**

- 必须保证“flume-env.sh”生效之后，再执行**步骤6**配置“properties.properties”文件。
- 若在本地配置该文件，配置完成后可参考如下步骤在Manager界面上上传配置文件。若操作顺序不规范，可能造成用户自定义环境变量丢失。
  1. 登录FusionInsight Manager。
  2. 选择“集群 > 服务 > Flume > 配置”，勾选Flume实例，在“flume.config.file”后单击“上传文件”，上传“properties.properties”文件。

----结束

### 7.9.3 如何开发 Flume 第三方插件

该操作指导用户进行第三方插件二次开发。

本章节适用于MRS 3.x及之后版本。

**步骤1** 将自主研发的代码打成jar包。

**步骤2** 安装Flume服务端或者客户端，如安装目录为“/opt/flumeclient”。

**步骤3** 建立插件目录布局。

1. 进入“Flume客户端安装目录/fusionInsight-flume-\*/plugins.d”路径下，使用以下命令建立目录，可根据实际业务进行命名，无固定名称：

```
cd /opt/flumeclient/fusioninsight-flume-1.9.0/plugins.d
```

```
mkdir thirdPlugin
```

```
cd thirdPlugin
```

```
mkdir lib libext native
```

显示结果如下：

```
[root@... plugins.d]#mkdir thirdPlugin
[root@... plugins.d]#ll
total 8
drwxr-x-- 3 root root 4096 ... native
drwxr-xr-x 2 root root 4096 ... thirdPlugin
[root@... plugins.d]#cd thirdPlugin/
[root@... thirdPlugin]#mkdir lib libext native
[root@... thirdPlugin]#ll
total 12
drwxr-xr-x 2 root root 4096 ... lib
drwxr-xr-x 2 root root 4096 ... libext
drwxr-xr-x 2 root root 4096 ... native
[root@... thirdPlugin]#
```

2. 将第三方jar包放入“Flume客户端安装目录/fusionInsight-flume-\*/plugins.d/thirdPlugin/lib”路径下，若该jar包依赖其他jar包，则将所依赖的jar包放入“Flume客户端安装目录/fusionInsight-flume-\*/plugins.d/thirdPlugin/libext”文件夹中，“Flume客户端安装目录/fusionInsight-flume-\*/plugins.d/thirdPlugin/native”放置本地库文件。



**步骤4** 配置“*Flume客户端安装目录*/fusionInsight-flume-\*/conf/properties.properties”文件。

具体properties.properties参数配置方法，参考[配置Flume非加密传输数据采集任务](#)和[配置Flume加密传输数据采集任务](#)对应典型场景中properties.properties文件参数列表的说明。

----结束

## 7.9.4 如何配置 Flume 定制脚本

Flume支持定制脚本，支持在传输前或者传输后执行指定的脚本，用于执行准备工作。

本章节适用于MRS 3.x及之后版本。

### 未安装 Flume 客户端

**步骤1** 获取软件包。

登录FusionInsight Manager，选择“*集群 > 待操作集群的名称 > 服务 > Flume*”进入Flume服务界面，在右上角选择“*更多 > 下载客户端*”，选择“*选择客户端类型*”为“*完整客户端*”，下载Flume服务客户端文件。

客户端文件名称为“*FusionInsight\_Cluster\_<集群ID>\_Flume\_Client.tar*”，本章节以“*FusionInsight\_Cluster\_1\_Flume\_Client.tar*”为例进行描述。

**步骤2** 上传软件包。以user用户将软件包上传到将要安装Flume服务客户端的节点目录上，例如“/opt/client”

#### 说明

user用户为安装和运行Flume客户端的用户。

**步骤3** 解压软件包。

以user用户登录将要安装Flume服务客户端的节点。进入安装包所在目录，例如“/opt/client”，执行如下命令解压安装包到当前目录。

```
cd /opt/client
```

```
tar -xvf FusionInsight_Cluster_1_Flume_Client.tar
```

**步骤4** 校验软件包。

执行sha256sum -c命令校验解压得到的文件，返回“OK”表示校验通过。例如：

```
sha256sum -c FusionInsight_Cluster_1_Flume_ClientConfig.tar.sha256
```

```
FusionInsight_Cluster_1_Flume_ClientConfig.tar: OK
```

**步骤5** 解压文件。

```
tar -xvf FusionInsight_Cluster_1_Flume_ClientConfig.tar
```

**步骤6** 在客户端/opt/client/FusionInsight\_Cluster\_1\_Flume\_ClientConfig/Flume/FlumeClient/flume/conf/flume-check.properties文件中配置client.per-check.shell，指向plugin.sh的绝对路径。

配置如下：

```
client.per-check.shell=/opt/client/
FusionInsight_Cluster_1_Flume_ClientConfig/Flume/FlumeClient/flume/
plugins.s/plugin.sh
```

```
plugins = com.huawei.flume.services.FlumePreTransmitService
```

```
flume.check.default.interval = 15
```

**步骤7** 配置/opt/client/FusionInsight\_Cluster\_1\_Flume\_ClientConfig/Flume/FlumeClient/flume/conf/plugin.conf文件，定义具体调用的脚本、相关参数。

配置如下：

```
RUN_PLUGIN="PLUGIN_LIST_1"
```

```
LOG_TO_HDFS_PATH="/yxs"
```

```
LOG_TO_HDFS_ENCODE_PATH="${LOG_TO_HDFS_PATH}/Flume_Encoded/"
```

```
PLUGIN_LINK_DIR="/tmp/yxs1"
```

```
PLUGIN_MV_TARGET_DIR="/tmp/yxs2"
```

```
PLUGIN_SUFFIX="COMPLETED"
```

```
PLUGIN_LIST_1="mv_complete.sh --linkdir ${PLUGIN_LINK_DIR} --mvtargetdir
${PLUGIN_MV_TARGET_DIR} --suffix ${PLUGIN_SUFFIX}"
```

**步骤8** 安装并启动Flume客户端。安装客户端详细操作请参考[安装Flume客户端](#)。

----结束

## 已安装 Flume 客户端

**步骤1** 在客户端flume-check.properties文件中配置client.per-check.shell，指向plugin.sh的绝对路径。

例如Flume客户端安装路径为“/opt/FlumeClient”，则flume-check.properties文件所在目录为/opt/FlumeClient/fusioninsight-flume-1.9.0/conf，

配置如下：

```
client.per-check.shell=/opt/FlumeClient/fusioninsight-flume-1.9.0/plugins.s/
plugin.sh
```

```
plugins = com.huawei.flume.services.FlumePreTransmitService
```

```
flume.check.default.interval = 15
```

**步骤2** 配置plugin.conf，定义具体调用的脚本、相关参数。

例如Flume客户端安装路径为“/opt/FlumeClient”，则plugin.conf配置文件所在目录为/opt/FlumeClient/fusioninsight-flume-1.9.0/conf，

配置如下：

```
RUN_PLUGIN="PLUGIN_LIST_1"
```

```
LOG_TO_HDFS_PATH="/yxs"
```

```
LOG_TO_HDFS_ENCODE_PATH="${LOG_TO_HDFS_PATH}/Flume_Encoded/"
```

```
PLUGIN_LINK_DIR="/tmp/yxs1"
```

```
PLUGIN_MV_TARGET_DIR="/tmp/yxs2"
```

```
PLUGIN_SUFFIX="COMPLETED"
```

```
PLUGIN_LIST_1="mv_complete.sh --linkdir ${PLUGIN_LINK_DIR} --mvtargetdir
${PLUGIN_MV_TARGET_DIR} --suffix ${PLUGIN_SUFFIX}"
```

**步骤3** 在客户端安装路径bin目录执行以下命令，重启Flume客户端，例如“/opt/FlumeClient/fusioninsight-flume-1.9.0/bin”。

```
./flume-manage.sh restart
```

----结束

# 8 使用 HBase

## 8.1 创建 HBase 权限角色

### 操作场景

该任务指导MRS集群管理员在Manager创建并设置HBase的角色。HBase角色可设置HBase管理员权限以及HBase表和列族的读（R）、写（W）、创建（C）、执行（X）或管理（A）权限。

用户需要在HBase中对指定的数据库或表设置权限，才能够创建表、查询数据、删除数据、插入数据、更新数据以及授权他人访问HBase表。

#### 说明

- 本章节适用于MRS 3.x及之后版本。
- 仅开启了Kerberos认证的集群（安全模式）支持创建HBase角色。
- 如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加HBase的Ranger访问权限策略](#)。

### 前提条件

- MRS集群管理员已明确业务需求。
- 已登录Manager。

### 创建 HBase 角色

**步骤1** 在Manager界面，选择“系统 > 权限 > 角色”。

## 权限



- 用户
- 用户组
- **角色**
- 安全策略
- 域和互信
- 第三方AD

**步骤2** 单击“添加角色”，然后在“角色名称”和“描述”输入角色名字与描述。

\* 角色名称:

配置资源权限:

| 所有资源 | 描述   |
|------|------|
| 所有资源 | 集群管理 |

描述:

**步骤3** 设置角色“配置资源权限”请参见[表8-1](#)。



HBase权限:

- HBase Scope: 对HBase表授权，最小支持设置列的读（R）和写（W）权限。
- HBase管理员权限: HBase管理员权限。

说明

用户对自己创建的表具有读（R）、写（W）、创建（C）、执行（X）或管理（A）权限。

表 8-1 设置 HBase 角色资源权限

| 任务场景            | 角色授权操作                                                                                                                                                                                                                    |
|-----------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 设置HBase管理员权限    | 在“配置资源权限”的表格中选择“待操作集群的名称 > HBase”，勾选“HBase管理员权限”。                                                                                                                                                                         |
| 设置用户创建表的权限      | <ol style="list-style-type: none"> <li>1. 在“配置资源权限”的表格中选择“待操作集群的名称 &gt; HBase &gt; HBase Scope”。</li> <li>2. 单击“global”。</li> <li>3. 在指定命名空间的“权限”列，勾选“创建”和“执行”。例如勾选默认命名空间“default”的“创建”和“执行”。</li> </ol>                  |
| 设置用户写入数据的权限     | <ol style="list-style-type: none"> <li>1. 在“配置资源权限”的表格中选择“待操作集群的名称 &gt; HBase &gt; HBase Scope &gt; global”。</li> <li>2. 在指定命名空间的“权限”列，勾选“写”。例如勾选默认命名空间“default”的“写”。HBase子对象默认可从父对象继承权限，此时已授予向命名空间中的表写入数据的权限。</li> </ol> |
| 设置用户读取数据的权限     | <ol style="list-style-type: none"> <li>1. 在“配置资源权限”的表格中选择“待操作集群的名称 &gt; HBase &gt; HBase Scope &gt; global”。</li> <li>2. 在指定命名空间的“权限”列，勾选“读”。例如勾选默认命名空间“default”的“读”。HBase子对象默认可从父对象继承权限，此时已授予从命名空间中的表读取数据的权限。</li> </ol> |
| 设置用户管理命名空间或表的权限 | <ol style="list-style-type: none"> <li>1. 在“配置资源权限”的表格中选择“待操作集群的名称 &gt; HBase &gt; HBase Scope &gt; global”。</li> <li>2. 在指定命名空间的“权限”列，勾选“管理”。例如勾选默认命名空间“default”的“管理”。</li> </ol>                                        |

| 任务场景        | 角色授权操作                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
|-------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 设置列的读取或写入权限 | <ol style="list-style-type: none"> <li>在“配置资源权限”的表格中选择“待操作集群的名称 &gt; HBase &gt; HBase Scope &gt; global”，单击指定命名空间显示命名空间的表。</li> <li>单击指定的表。</li> <li>单击指定的列族。</li> <li>确认是否是新建角色？ <ul style="list-style-type: none"> <li>是，在“资源名称”的输入框输入列名称，多个列用英文逗号分隔，勾选“读”或“写”。如果HBase表中不存在同名的列，则创建同名的列后角色将拥有该列的权限。列权限设置完成。</li> <li>否，修改已有HBase角色的列权限，表格将显示已单独设置权限的列，执行<a href="#">步骤3.5</a>。</li> </ul> </li> <li>角色新增列权限，在“资源名称”的输入框输入列名称并设置列的权限。角色修改列权限，可以在“资源名称”的输入框输入列名称并设置列权限，也可以在表格中直接修改列的权限。若在表格中修改了列权限，又同时增加了同名的列权限，则无法保存。角色修改列权限，建议直接修改列的权限。支持搜索功能。</li> </ol> |

**步骤4** 单击“确定”完成，返回“角色”。

---结束

## 8.2 HBase 客户端使用实践

### 操作场景

该任务指导用户在运维场景或业务场景中使用HBase客户端。

### 前提条件

- 已安装客户端。例如安装目录为“/opt/hadoopclient”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- 各组件业务用户由MRS集群管理员根据业务需要创建。  
“机机”用户需要下载keytab文件，“人机”用户第一次登录时需修改密码。
- 非root用户使用HBase客户端，请确保该HBase客户端目录的属主为该用户，否则请参考如下命令修改属主。

```
chown user:group -R 客户端安装目录/HBase
```

### 使用 HBase 客户端（MRS 3.x 之前版本）

**步骤1** 安装客户端，具体请参考[安装客户端](#)章节。

**步骤2** 以客户端安装用户，登录安装客户端的节点。

**步骤3** 执行以下命令切换到客户端目录。

```
cd /opt/hadoopclient
```

**步骤4** 执行以下命令配置环境变量。

```
source bigdata_env
```

**步骤5** 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建HBase表的权限，具体请参见[创建角色](#)配置拥有对应权限的角色，参考[创建用户](#)章节，为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit 组件业务用户
```

例如，`kinit hbaseuser`。

**步骤6** 直接执行HBase组件的客户端命令。

```
hbase shell
```

```
----结束
```

## 使用 HBase 客户端（MRS 3.x 及之后版本）

**步骤1** 安装客户端，具体请参考[安装客户端](#)章节。

**步骤2** 以客户端安装用户，登录安装客户端的节点。

**步骤3** 执行以下命令切换到客户端目录。

```
cd /opt/hadoopclient
```

**步骤4** 执行以下命令配置环境变量。

```
source bigdata_env
```

**步骤5** 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建HBase表的权限，具体请参见[角色管理](#)配置拥有对应权限的角色，参考[创建用户](#)章节，为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit 组件业务用户
```

例如，`kinit hbaseuser`。

**步骤6** 直接执行HBase组件的客户端命令。

```
hbase shell
```

```
----结束
```

## HBase 客户端常用命令

常用的HBase客户端命令如下表所示。更多命令可参考<http://hbase.apache.org/2.2/book.html>。



表 8-2 HBase 客户端命令

| 命令       | 说明                                                                                         |
|----------|--------------------------------------------------------------------------------------------|
| create   | 创建一张表，例如create 'test', 'f1', 'f2', 'f3'。                                                   |
| disable  | 停止指定的表，例如disable 'test'。                                                                   |
| enable   | 启动指定的表，例如enable 'test'。                                                                    |
| alter    | 更改表结构。可以通过alter命令增加、修改、删除列族信息以及表相关的参数值，例如alter 'test', {NAME => 'f3', METHOD => 'delete'}。 |
| describe | 获取表的描述信息，例如describe 'test'。                                                                |
| drop     | 删除指定表。删除前表必须已经是停止状态，例如drop 'test'。                                                         |
| put      | 写入指定cell的value。Cell的定位由表、rowk、列组合起来唯一决定，例如put 'test','r1','f1:c1','myvalue1'。              |
| get      | 获取行的值或者行的指定cell的值。例如get 'test','r1'。                                                       |
| scan     | 查询表数据，参数中需指定表名和scanner，例如scan 'test'。                                                      |

## 8.3 快速使用 HBase 进行离线数据分析

HBase是一个高可靠性、高性能、面向列、可伸缩的分布式存储系统。本章节提供从零开始使用HBase的操作指导，通过客户端实现创建表，往表中插入数据，修改表，读取表数据，删除表中数据以及删除表的功能。

### 背景信息

假定用户开发一个应用程序，用于管理企业中的使用A业务的用户信息，使用HBase客户端实现A业务操作流程如下：

- 创建用户信息表user\_info。
- 在用户信息中新增用户的学历、职称信息。
- 根据用户编号查询用户姓名和地址。
- 根据用户姓名进行查询。
- 用户销户，删除用户信息表中该用户的数据。
- A业务结束后，删除用户信息表。

表 8-3 用户信息

| 编号          | 姓名 | 性别 | 年龄 | 地址  |
|-------------|----|----|----|-----|
| 12005000201 | A  | 男  | 19 | A城市 |
| 12005000202 | B  | 女  | 23 | B城市 |
| 12005000203 | C  | 男  | 26 | C城市 |

| 编号          | 姓名 | 性别 | 年龄 | 地址  |
|-------------|----|----|----|-----|
| 12005000204 | D  | 男  | 18 | D城市 |
| 12005000205 | E  | 女  | 21 | E城市 |
| 12005000206 | F  | 男  | 32 | F城市 |
| 12005000207 | G  | 女  | 29 | G城市 |
| 12005000208 | H  | 女  | 30 | H城市 |
| 12005000209 | I  | 男  | 26 | I城市 |
| 12005000210 | J  | 男  | 25 | J城市 |

## 前提条件

已安装客户端，例如安装目录为“/opt/client”。以下操作的客户端目录只是举例，请根据实际安装目录修改。在使用客户端前，需要先下载并更新客户端配置文件，确认 Manager 的主管理节点后才能使用客户端。

## 操作步骤

**MRS 3.x以前版本集群执行以下操作：**

**步骤1** 下载客户端配置文件。

1. 登录MRS Manager页面，具体请参见[访问集群Manager](#)，然后选择“服务管理”。
2. 单击“下载客户端”。  
“客户端类型”选择“仅配置文件”，“下载路径”选择“服务器端”，单击“确定”开始生成客户端配置文件，文件生成后默认保存在主管理节点“/tmp/MRS-client”。文件保存路径支持自定义。

图 8-1 下载客户端的配置文件

### 下载客户端

**警告：**生成客户端会占用大量的磁盘IO，不建议在集群处于安装中、启动中、打补丁中等非稳态场景进行“下载客户端”操作。

\* 客户端类型  完整客户端  仅配置文件

\* 下载路径  服务器端  远端主机

仅保存到服务器如下路径，如果存在客户端文件，会覆盖路径下已有的客户端文件。

\* 客户端路径

确定

取消

**步骤2** 登录MRS Manager的主管理节点。

1. 在集群详情的“节点信息”页签中查看节点名称，名称中包含“master1”的节点为Master1节点，名称中包含“master2”的节点为Master2节点。

MRS Manager的主管理节点默认安装在集群Master节点上。在主备模式下，由于Master1和Master2之间会切换，Master1节点不一定是MRS Manager的主管理节点，需要在Master1节点中执行命令，确认MRS Manager的主管理节点。命令请参考[步骤2.4](#)。

2. 以root用户使用密码方式登录Master1节点。操作方法，请参见[登录集群节点](#)章节。
3. 切换至omm用户。

```
sudo su - root
```

```
su - omm
```

4. 执行以下命令确认MRS Manager的主管理节点。

```
sh ${BIGDATA_HOME}/om-0.0.1/sbin/status-oms.sh
```

回显信息中“HAActive”参数值为“active”的节点为主管理节点（如下例中“mgtomsdat-sh-3-01-1”为主管理节点），参数值为“standby”的节点为备管理节点（如下例中“mgtomsdat-sh-3-01-2”为备管理节点）。

```
Ha mode
double
NodeName HostName HAVersion StartTime HAActive
HAAllResOK HARunPhase
192-168-0-30 mgtomsdat-sh-3-01-1 V100R001C01 2019-11-18 23:43:02
active normal Activated
192-168-0-24 mgtomsdat-sh-3-01-2 V100R001C01 2019-11-21 07:14:02
standby normal Deactivated
```

5. 使用root用户登录MRS Manager的主管理节点，例如“192-168-0-30”节点，并执行以下命令切换到omm用户。

```
sudo su - omm
```

**步骤3** 执行以下命令切换到客户端安装目录，例如“/opt/client”。

```
cd /opt/client
```

**步骤4** 执行以下命令，更新主管理节点的客户端配置。

```
sh refreshConfig.sh /opt/client 客户端配置文件压缩包完整路径
```

例如，执行命令：

```
sh refreshConfig.sh /opt/client /tmp/MRS-client/MRS_Services_Client.tar
```

界面显示以下信息表示配置刷新更新成功：

```
ReFresh components client config is complete.
Succeed to refresh components client config.
```

**说明**

步骤[步骤1](#)~[步骤4](#)的操作也可以参考[更新客户端](#)页面的方法二操作。

**步骤5** 在Master节点使用客户端。

1. 在已更新客户端的主管理节点，例如“192-168-0-30”节点，执行以下命令切换到客户端目录。

```
cd /opt/client
```

2. 执行以下命令配置环境变量。  
**source bigdata\_env**
3. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建HBase表的权限，具体请参见[创建角色](#)配置拥有对应权限的角色，参考[创建用户](#)为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行此命令。  
**kinit MRS集群用户**  
例如，**kinit hbaseuser**。
4. 直接执行HBase组件的客户端命令。  
**hbase shell**

**步骤6** 运行HBase客户端命令，实现A业务。

1. 根据[表8-3](#)创建用户信息表user\_info并添加相关数据。  
**create 'user\_info',{NAME => 'i'}**  
以增加编号12005000201的用户信息为例，其他用户信息参照如下命令依次添加：  
**put 'user\_info','12005000201','i:name','A'**  
**put 'user\_info','12005000201','i:gender','Male'**  
**put 'user\_info','12005000201','i:age','19'**  
**put 'user\_info','12005000201','i:address','City A'**
2. 在用户信息表user\_info中新增用户的学历、职称信息。  
以增加编号为12005000201的用户的学历、职称信息为例，其他用户类似。  
**put 'user\_info','12005000201','i:degree','master'**  
**put 'user\_info','12005000201','i:pose','manager'**
3. 根据用户编号查询用户姓名和地址。  
以查询编号为12005000201的用户姓名和地址为例，其他用户类似。  
**scan'user\_info',**  
**{STARTROW=>'12005000201',STOPROW=>'12005000201',COLUMNS=>['i:name','i:address']}**
4. 根据用户姓名进行查询。  
以查询A用户信息为例，其他用户类似。  
**scan'user\_info',{FILTER=>"SingleColumnValueFilter('i','name',=,'binary:A')"**
5. 删除用户信息表中该用户的数据。  
所有用户的数据都需要删除，以删除编号为12005000201的用户数据为例，其他用户类似。  
**delete'user\_info','12005000201','i'**
6. 删除用户信息表。  
**disable'user\_info'**  
**drop 'user\_info'**

----结束

**MRS 3.x及之后版本集群执行以下操作：**

- 步骤1** 在主管理节点使用客户端。

1. 安装客户端，具体请参考[安装客户端](#)章节。
2. 以客户端安装用户登录客户端安装节点，执行以下命令切换到客户端目录。  
**cd /opt/client**
3. 执行以下命令配置环境变量。  
**source bigdata\_env**
4. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建HBase表的权限，具体请参见[创建HBase权限角色](#)配置拥有对应权限的角色，参考[创建用户](#)为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行此命令。  
**kinit MRS集群用户**  
例如，**kinit hbaseuser**。
5. 直接执行HBase组件的客户端命令。  
**hbase shell**

## 步骤2 运行HBase客户端命令，实现A业务。

1. 根据[表8-3](#)创建用户信息表user\_info并添加相关数据。  
**create 'user\_info',{NAME => 'i'}**  
以增加编号12005000201的用户信息为例，其他用户信息参照如下命令依次添加：  
**put 'user\_info','12005000201','i:name','A'**  
**put 'user\_info','12005000201','i:gender','Male'**  
**put 'user\_info','12005000201','i:age','19'**  
**put 'user\_info','12005000201','i:address','City A'**
2. 在用户信息表user\_info中新增用户的学历、职称信息。  
以增加编号为12005000201的用户的学历、职称信息为例，其他用户类似。  
**put 'user\_info','12005000201','i:degree','master'**  
**put 'user\_info','12005000201','i:pose','manager'**
3. 根据用户编号查询用户姓名和地址。  
以查询编号为12005000201的用户姓名和地址为例，其他用户类似。  
**scan 'user\_info',**  
**{STARTROW=>'12005000201',STOPROW=>'12005000201',COLUMNS=>['i:name','i:address']}**
4. 根据用户姓名进行查询。  
以查询A用户信息为例，其他用户类似。  
**scan 'user\_info',{FILTER=>"SingleColumnValueFilter('i','name',=,'binary:A')"**
5. 删除用户信息表中该用户的数据。  
所有用户的数据都需要删除，以删除编号为12005000201的用户数据为例，其他用户类似。
  - 依次删除编号为12005000201的用户的的所有数据字段，以删除“age”字段为例：  
**delete 'user\_info','12005000201','i:age'**
  - 删除编号为12005000201的用户的的所有数据：  
**deleteall 'user\_info','12005000201'**

6. 删除用户信息表。  
`disable 'user_info'`  
`drop 'user_info'`  
----结束

## 8.4 使用 BulkLoad 工具向 HBase 迁移数据

HBase的数据都是存储在HDFS中的，数据导入即是加载存放在HDFS中的数据到HBase表中。Apache HBase提供了“Import”和“ImportTsv”工具用于批量导入HBase数据。

- “Import”通过“org.apache.hadoop.hbase.mapreduce.Import”方法导入已导出至HDFS中的HBase数据。
- “ImportTsv”通过“org.apache.hadoop.hbase.mapreduce.ImportTsv”可将TSV格式的数据加载到HBase中。

更多详细信息请参见：<http://hbase.apache.org/2.2/book.html#tools>。

## 8.5 HBase 数据操作

### 8.5.1 创建 HBase 索引进行数据查询

#### 操作场景

HBase是一个Key-Value类型的分布式存储数据库，HIndex为HBase提供了按照某些列的值进行索引的能力，缩小搜索范围并缩短时延。

#### 使用约束

- 列族应以“;”分隔。
- 列和数据类型应包含在“[]”中。
- 列数据类型在列名称后使用“->”指定。
- 如果未指定列数据类型，则使用默认数据类型（字符串）。
- “#”用于在两个索引详细信息之间进行分隔。
- 以下是一个可选参数：  
-Dscan.caching：在扫描数据表时的缓存行数。  
如果不设置该参数，则默认值为1000。
- 为单个Region构建索引是为了修复损坏的索引。  
此功能不应用于生成新索引。

#### 创建 HBase HIndex

**步骤1** 安装HBase客户端，详情参见[HBase客户端使用实践](#)。

**步骤2** 进入客户端安装路径，例如“/opt/client”

```
cd /opt/client
```

**步骤3** 配置环境变量。

**source bigdata\_env**

**步骤4** 如果集群为安全模式，执行以下命令进行用户认证。普通模式集群无需执行用户认证。

**kinit 组件业务用户**

**步骤5** 执行以下命令访问Hindex。

**hbase org.apache.hadoop.hbase.hindex.mapreduce.TableIndexer**

表 8-4 HIndex 常用命令

| 功能            | 命令                                                                                                                                                                                          |
|---------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 增加索引          | TableIndexer-Dtablename.to.index=table1 -<br>Dindexspecs.to.add='IDX1=>cf1:[q1->datatype],[q2],[q3];cf2:<br>[q1->datatype],[q2->datatype]#IDX2=>cf1:[q5]'                                   |
| 构建索引          | TableIndexer -Dtablename.to.index=table1 -<br>Dindexnames.to.build='IDX1#IDX2'                                                                                                              |
| 删除索引          | TableIndexer -Dtablename.to.index=table1 -<br>Dindexnames.to.drop='IDX1#IDX2'                                                                                                               |
| 禁用索引          | TableIndexer -Dtablename.to.index=table1 -<br>Dindexnames.to.disable='IDX1#IDX2'                                                                                                            |
| 同时添加和构建索引     | TableIndexer -Dtablename.to.index=table1 -<br>Dindexspecs.to.add='IDX1=>cf1:[q1->datatype],[q2],[q3];cf2:<br>[q1->datatype],[q2->datatype]#IDX2=>cf1:[q5]' -<br>Dindexnames.to.build='IDX1' |
| 为单个Region构建索引 | TableIndexer -Dtablename.to.index=table1 -<br>Dregion.to.index=regionEncodedName -<br>Dindexnames.to.build='IDX1#IDX2'                                                                      |

**说明**

- **IDX1**: 索引名称。
- **cf1**: 列族名称。
- **q1**: 列名称。
- **datatype**: 数据类型，包括String、Integer、Double、Float、Long、Short、Byte、Char。

----结束

## 8.5.2 配置 HBase 数据压缩格式和编码

### 操作场景

HBase可以通过对HFile中的data block编码，减少Key-Value中Key的重复部分，从而减少空间的使用。目前对data block的编码方式有：NONE、PREFIX、DIFF、

FAST\_DIFF和ROW\_INDEX\_V1，其中NONE表示不使用编码。另外，HBase还支持使用压缩算法对HFile文件进行压缩，默认支持的压缩算法有：NONE、GZ、SNAPPY和ZSTD，其中NONE表示HFile不压缩。

这两种方式都是作用在HBase的列簇上，可以同时使用，也可以单独使用。

## 前提条件

- 已安装HBase客户端。例如，客户端安装目录为“/opt/client”。
- 如果集群开启了Kerberos认证，操作的用户还需要具备对应的操作权限。即创建表时需要具备对应的namespace或更高级别的创建（C）或者管理（A）权限，修改表时需要具备已创建的表或者更高级别的创建（C）或者管理（A）权限。具体的授权操作请参考[创建HBase权限角色](#)章节。

## 配置 HBase 数据压缩格式和编码

创建时设置data block encoding和压缩算法。

- **方法一：使用hbase shell。**
  - a. 以客户端安装用户，登录安装客户端的节点。
  - b. 执行以下命令切换到客户端目录。  
**cd /opt/client**
  - c. 执行以下命令配置环境变量。  
**source bigdata\_env**
  - d. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。  
**kinit 组件业务用户**  
例如，**kinit hbaseuser**。
  - e. 执行以下命令登录HBase客户端：  
**hbase shell**
  - f. 创建表。  
**create 't1', {NAME => 'f1', COMPRESSION => 'SNAPPY', DATA\_BLOCK\_ENCODING => 'FAST\_DIFF'}**

### 📖 说明

- t1：表名。
  - f1：列簇名。
  - SNAPPY：该列簇使用的压缩算法为SNAPPY。
  - FAST\_DIFF：使用的编码方式为FAST\_DIFF。
  - {}内的参数为指定列簇的参数，多个列簇可以用多个{}，然后用逗号隔开。关于建表语句的更多使用说明可以在**hbase shell**中执行**help 'create'**进行查看。
- **方法二：使用Java API。**

以下代码片段仅展示如何在建表时设置列簇的编码和压缩方式，完整的建表代码以及如何通过代码建表请参考HBase开发指南的[创建HBase表](#)章节。

```
TableDescriptorBuilder htd = TableDescriptorBuilder.newBuilder(TableName.valueOf("t1")); // 创建t1表的descriptor.
ColumnFamilyDescriptorBuilder hcd =
ColumnFamilyDescriptorBuilder.newBuilder(Bytes.toBytes("f1")); // 创建列簇f1的builder.
hcd.setDataBlockEncoding(DataBlockEncoding.FAST_DIFF); // 设置列簇f1的编码方式为FAST_DIFF.
```



```
hcd.setCompressionType(Compression.Algorithm.SNAPPY);// 设置列簇f1的压缩算法为SNAPPY
htd.setColumnFamily(hcd.build())// 将列簇f1添加到t1表的descriptor.
```

### 对已存在的表设置或修改data block encoding和压缩算法

- **方法一：使用hbase shell。**
  - a. 以客户端安装用户，登录安装客户端的节点。
  - b. 执行以下命令切换到客户端目录。  
**cd /opt/client**
  - c. 执行以下命令配置环境变量。  
**source bigdata\_env**
  - d. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。  
**kinit 组件业务用户**  
例如，**kinit hbaseuser**。
  - e. 执行以下命令登录HBase客户端：  
**hbase shell**
  - f. 执行以下命令修改HBase表：  
**alter 't1', {NAME => 'f1', COMPRESSION => 'SNAPPY',  
DATA\_BLOCK\_ENCODING => 'FAST\_DIFF'}**

- **方法二：使用Java API。**

以下代码片段仅展示如何修改指定表的已有列簇的编码和压缩方式，完整的修改表代码以及如何通过代码修改表请参考HBase应用开发指南的[修改HBase表](#)章节：

```
TableDescriptor htd = admin.getDescriptor(TableName.valueOf("t1"));// 获取表t1的descriptor
ColumnFamilyDescriptor originCF = htd.getColumnFamily(Bytes.toBytes("f1"));// 获取列簇f1的
descriptor
builder.ColumnFamilyDescriptorBuilder hcd = ColumnFamilyDescriptorBuilder.newBuilder(originCF);//
通过已有的列簇属性构造一个builder
hcd.setDataBlockEncoding(DataBlockEncoding.FAST_DIFF);// 重新设置列簇的编码方式为FAST_DIFF
hcd.setCompressionType(Compression.Algorithm.SNAPPY);// 重新设置列簇的压缩算法为SNAPPY
admin.modifyColumnFamily(TableName.valueOf("t1"), hcd.build());// 提交到服务端修改列簇f1的属性
```

修改后完成后，已有的HFile的编码和压缩方式需要在下次做完compaction后才会生效。

## 8.6 HBase 企业级能力增强

### 8.6.1 配置 HBase 本地二级索引提升查询效率

#### 8.6.1.1 HBase 本地二级索引介绍

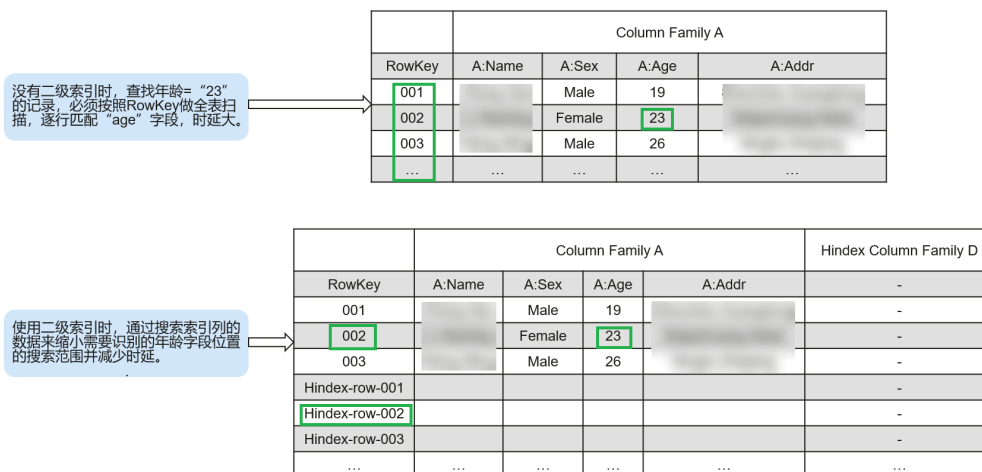
##### 场景介绍

HBase是基于Key-Value的分布式存储数据库，基于rowkeys对表中的数据按照字典进行排序。如果您根据指定的rowkey查询数据，或者扫描指定rowkey范围内的数据，HBase可以快速查找到需要读取的数据，从而提高效率。在大多数实际情况下，会需要查询列值为XXX的数据。HBase提供了Filter功能来查询具有特定列值的数据：所有数据按RowKey的顺序进行扫描，然后将数据与特定的列值进行匹配，直到找到所需的

数据。过滤器功能会scan一些不必要的的数据以获取所需的数据。因此，Filter功能不能满足高性能标准频繁查询的要求。

这就是HBase HIndex产生的背景。如图8-2所示，HBase HIndex为HBase提供了能够根据特定的列值进行索引的能力，使得查询会变得更快速。

图 8-2 HBase HIndex

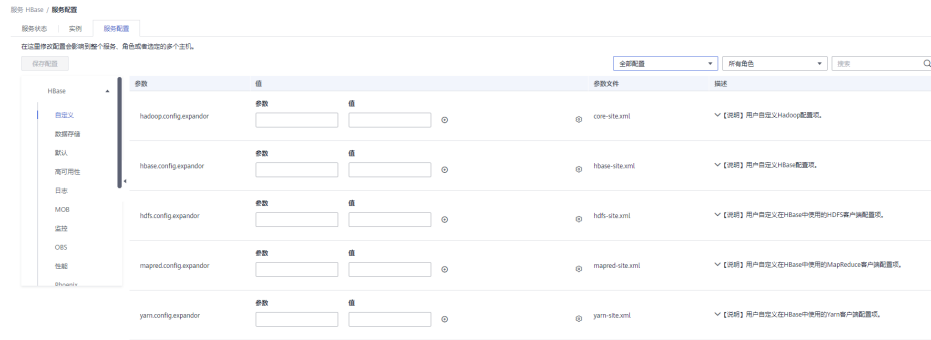


### 说明

- 索引数据不支持滚动升级。
- 复合索引：用户必须将所有参与复合索引的列全部放入/删除，否则会导致数据不一致。
- 用户不应将任何split policy显式地配置到已建立索引的数据表中。
- 不支持mutation操作，如increment、append。
- 不支持列索引的版本maxVersions > 1。
- 添加索引的列值不应超过32KB。
- 当用户数据由于列族级TTL失效而被删除时，相应的索引数据不会立即删除。索引数据将在major compaction期间被删除。
- 创建索引后，不应更改用户列族的TTL。
  - 如果在创建索引后将列族TTL更改为更高值，则应删除并重新创建索引，否则某些已生成的索引数据可能比用户数据先删除。
  - 如果在创建索引后将列族TTL更改为较低值，则索引可能会晚于用户数据被删除。
- HBase表启动容灾之后，主集群新建二级索引，索引表变更不会自动同步到备集群。要实现该容灾场景，必须执行以下操作：
  1. 在主表创建二级索引之后，需要在备集群使用相同方法创建结构、名称完全相同的二级索引。
  2. 在主集群手动将索引列族（默认是d）的REPLICATION\_SCOPE设置为1。

## 配置 HBase 本地二级索引

1. 登录MRS控制台，单击集群名称，选择“组件管理”。
2. 在组件列表中选择“HBase > 服务配置”，在下拉列表中将“基础配置”切换为“全部配置”，进入HBase服务参数“全部配置”界面。



3. 在HBase全部配置界面查看参数。

| 配置入口           | 配置项                              | 默认值                                                                                                                                                                                                                                                                                                                 | 描述                                                                |
|----------------|----------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------|
| “HMaster > 系统” | hbase.coprocessor.master.classes | org.apache.hadoop.hbase.hindex.server.master.HIndexMasterCoprocesor,com.xxx.hadoop.hbase.backup.services.RecoveryCoprocesor,org.apache.ranger.authorization.hbase.RangerAuthorizationCoprocesor,org.apache.hadoop.hbase.security.access.ReadOnlyClusterEnabler,org.apache.hadoop.hbase.rsgroup.RSGroupAdminEndpoint | 该协处理器用于在启用Hindex功能后处理Master级的操作，比如创建索引meta表，添加索引，删除索引，删除表删除索引元数据。 |

| 配置入口                                | 配置项                                    | 默认值                                                                                                                                                                                       | 描述                                                 |
|-------------------------------------|----------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------|
| “RegionServer<br>><br>RegionServer” | hbase.coprocessor.regionserver.classes | org.apache.hadoop.hbase.hindex.server.regionserver.HIndexRegionServerCoprocessor,org.apache.hadoop.hbase.JMXListener,org.apache.ranger.authorization.hbase.RangerAuthorizationCoprocessor | 该协处理器用于在启用Hindex功能后实际上处理Master下发到RegionServer上的操作。 |

| 配置入口 | 配置项                              | 默认值                                                                                                                                                                                                                                                                                                                                                                                                                                              | 描述                                  |
|------|----------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------|
|      | hbase.coprocessor.region.classes | org.apache.hadoop.hbase.hindex.server.regionserver.HIndexRegionCoprocessor,org.apache.hadoop.hbase.security.token.TokenProvider,com.xxx.hadoop.hbase.backup.services.RecoveryCoprocessor,org.apache.ranger.authorization.hbase.RangerAuthorizationCoprocessor,org.apache.hadoop.hbase.security.access.SecureBulkLoadEndpoint,org.apache.hadoop.hbase.security.access.ReadOnlyClusterEnabler,org.apache.hadoop.hbase.coprocessor.MetaTableMetrics | 该协处理器用于在启用Hindex功能后实际上操作Region上的数据。 |

| 配置入口 | 配置项                           | 默认值                                                                                                                                                                                       | 描述                                                                                    |
|------|-------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------|
|      | hbase.coprocessor.wal.classes | org.apache.hadoop.hbase.hindex.server.regionserver.HIndexRegionServerCoprocessor,org.apache.hadoop.hbase.JMXListener,org.apache.ranger.authorization.hbase.RangerAuthorizationCoprocessor | 该协处理器用于Replication，其会过滤掉索引数据以避免索引数据发送到对等集群中，对等集群中的数据索引数据将会自己生成。<br>该参数仅MRS 3.x之前版本支持。 |

#### 说明

- 上述默认值为启用HBase HIndex功能后需额外配置的值，当前支持HBase HIndex功能的MRS集群默认已配置。
- 必须确保Master参数配置在HMster上，region/regionserver参数配置在RegonServer上。

## 相关接口

使用HIndex的API都在类org.apache.hadoop.hbase.hindex.client.HIndexAdmin中，相关接口介绍如下：

| 操作   | 接口                   | 描述                                                                                                                                                     | 注意事项                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
|------|----------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 添加索引 | addIndices()         | 将索引添加到没有数据的表中。调用此接口会将用户指定的索引添加到表中，但会跳过生成索引数据。因此，在此操作之后，索引不能用于scan/filter操作。该接口的使用场景为用户想要在具有大量预先存在用户数据的表上批量添加索引，其具体操作为使用诸如TableIndexer工具之类的外部工具来构建索引数据。 | <ul style="list-style-type: none"> <li>索引一旦添加则不能修改。若要修改，则需先删除旧的索引然后重新创建。</li> <li>应注意不要在具有不同索引名称的相同列上创建两个索引，否则会导致存储和处理的资源浪费。</li> <li>索引不能添加到系统表中。</li> <li>向索引列put数据时不支持append和increment操作。</li> <li>如果客户端出现任何故障，除非发生DoNotRetryIOException，否则应该重试。</li> <li>索引列族按以下优先级从数据表中已存在的列族选取，优先级从高到低依次为：<br/>d、#、@、\$、%、#0、@0、\$0、%0、#1、@1 ...上至#255、@255、\$255和%255<br/>创建索引时，系统会在表中按以上优先级顺序检查是否存在以上列族，如果不存在，则将第一个不存在的列族设为索引列族。<br/>例如： <ul style="list-style-type: none"> <li>数据表中仅存在d列族，则索引列族默认为#。</li> <li>数据表中已存在d和#列族，则默认索引列族默认为@。</li> </ul> </li> </ul> |
|      | addIndicesWithData() | 将索引添加到有数据的表中。此方法将用户指定的索引添加到表中，并会对已经存在的用户数据创建对应的索引数据，也可先调用该方法生成索引再存入用户数据的同时生成索引数据。在此操作之后，这些索引立即可用于scan/filter操作。                                        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |

| 操作   | 接口                    | 描述                                                                                                                                                                                          | 注意事项                                                                                                                                                                                                                                                                                                                                                                   |
|------|-----------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|      |                       |                                                                                                                                                                                             | <ul style="list-style-type: none"> <li>- 数据表中已存在 d、#和\$列族，则索引列族默认为@。</li> <li>• 可以通过HIndex TableIndexer工具添加索引而无需建立索引数据。</li> </ul>                                                                                                                                                                                                                                     |
| 删除索引 | dropIndices()         | <p>仅删除索引。该API从表中删除用户指定的索引，但跳过相应的索引数据。在此操作之后，索引不能用于scan/filter操作。集群在major compaction期间会自动删除旧的索引数据。</p> <p>此API使用场景为表中包含大量索引数据且dropIndicesWithData()不可行。另外，也可以通过TableIndexer工具删除索引以及索引数据。</p> | <ul style="list-style-type: none"> <li>• 在索引的状态为ACTIVE, INACTIVE和DROPPING时，允许禁用索引的操作。</li> <li>• 对于使用dropIndices()删除索引的操作，用户必须确保在将索引添加到具有相同索引名的表之前，相应的索引数据已被删除（即major compaction已完成）。</li> <li>• 用户删除相应的索引会删除： <ul style="list-style-type: none"> <li>- 一个带有索引的列族。</li> <li>- 组合索引所有列族中的任一个列族。</li> </ul> </li> <li>• 索引可以通过HIndex TableIndexer工具与索引数据一起删除。</li> </ul> |
|      | dropIndicesWithData() | <p>删除索引数据。此API删除用户指定的索引，并删除用户表中与这些索引对应的所有索引数据。在此操作之后，删除的索引完全从表中删除，不再可用于scan/filter操作。</p>                                                                                                   |                                                                                                                                                                                                                                                                                                                                                                        |



| 操作       | 接口               | 描述                                    | 注意事项                                                                                                                                                                                                                                                                 |
|----------|------------------|---------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 启用/禁用索引  | disableIndices() | 该API禁用所有用户指定的索引，使其不再可用于scan/filter操作。 | <ul style="list-style-type: none"> <li>在索引的状态为ACTIVE, INACTIVE和BUILDING时允许启用索引的操作。</li> <li>在索引的状态为ACTIVE和INACTIVE时允许禁用索引操作。</li> <li>在禁用索引之前，用户必须确保索引数据与用户数据一致。如果在索引处于禁用状态期间没有在表中添加新的数据，索引数据与用户数据将保持一致。</li> <li>启用索引时，可以通过使用TableIndexer工具构建索引来保证数据一致性。</li> </ul> |
|          | enableIndices()  | 该API启用所有用户指定的索引，使其可用于scan/filter操作。   |                                                                                                                                                                                                                                                                      |
| 查看已创建的索引 | listIndices()    | 该API可用于列出给定表中的所有索引。                   | 无                                                                                                                                                                                                                                                                    |

## 基于 HBase 本地二级索引查询数据

在具有索引的用户表中，可以使用Filter来查询数据。对于创建单索引和组合索引的用户表，使用过滤器查询的结果与没有使用索引的表相同，但数据查询性能高于没有使用索引的表。

索引的使用规则如下：

- 对于一个或多个列创建单个索引的情况：
  - 当将此列用于AND或OR查询筛选时，使用索引可以提高查询性能。  
例如，Filter\_Condition ( IndexCol1 ) AND / OR Filter\_Condition ( IndexCol2 ) 。
  - 当在查询中使用“索引列和非索引列”进行过滤时，此索引可以提高查询性能。  
例如，Filter\_Condition ( IndexCol1 ) AND Filter\_Condition ( IndexCol2 ) AND Filter\_Condition ( NonIndexCol1 ) 。
  - 当在查询中使用“索引列或非索引列”进行筛选时，但不使用索引，查询性能不会提高。  
例如，Filter\_Condition ( IndexCol1 ) AND / OR Filter\_Condition ( IndexCol2 ) OR Filter\_Condition ( NonIndexCol1 ) 。

- 对于为多个列创建组合索引的情况：
  - 当用于查询的列是组合索引的全部或部分列并且与组合索引具有相同的顺序时，使用索引会提高查询性能。  
例如，为C1，C2和C3创建组合索引。
    - 该索引在以下情况下生效：  
Filter\_Condition ( IndexCol1 ) AND Filter\_Condition ( IndexCol2 )  
AND Filter\_Condition ( IndexCol3 )  
Filter\_Condition ( IndexCol1 ) AND Filter\_Condition ( IndexCol2 )  
FILTER\_CONDITION ( IndexCol1 )
    - 该索引在下列情况下不生效：  
Filter\_Condition ( IndexCol2 ) AND Filter\_Condition ( IndexCol3 )  
Filter\_Condition ( IndexCol1 ) AND Filter\_Condition ( IndexCol3 )  
FILTER\_CONDITION ( IndexCol2 )  
FILTER\_CONDITION ( IndexCol3 )
  - 当在查询中使用“索引列和非索引列”进行过滤时，使用索引可提高查询性能。  
例如：  
Filter\_Condition ( IndexCol1 ) AND Filter\_Condition ( NonIndexCol1 )  
Filter\_Condition ( IndexCol1 ) AND Filter\_Condition ( IndexCol2 ) AND  
Filter\_Condition ( NonIndexCol1 )
  - 当在查询中使用“索引列或非索引列”进行筛选时，但不使用索引，查询性能不会提高。  
例如：  
Filter\_Condition ( IndexCol1 ) OR Filter\_Condition ( NonIndexCol1 )  
( Filter\_Condition ( IndexCol1 ) AND Filter\_Condition ( IndexCol2 ) ) OR  
( Filter\_Condition ( NonIndexCol1 ) )
  - 当多个列用于查询时，只能为组合索引中的最后一列指定值范围，而其他列只能设置为指定值。  
例如，为C1，C2和C3创建组合索引。在范围查询中，只能为C3设置数值范围，过滤条件为“C1 = XXX，C2 = XXX，C3 = 数值范围”。

## HBase 本地二级索引查询策略选择

使用SingleColumnValueFilter或SingleColumnRangeFilter，它会在一个在过滤条件中提供确定值column\_family:qualifierpair（称该列为col1）。

若col1作为表上的第一个索引列，那么该表上的任何索引都可以成为查询期间使用的候选索引。例如：

如果有col1上的索引，可以将此索引作为候选索引，因为col1是此索引的第一列也是唯一的列；如果在col1和col2上有另一个索引，可以将此索引视为候选索引，因为col1是索引列表中的第一列。另一方面，如果在col2和col1上有一个索引，则不能将此索引作为候选索引，因为索引列表中的第一列不是col1。

现在最适合使用索引的方法是，当有多个候选索引时，需要从可能的候选索引中选择最适合scan数据的索引。

可借助以下方案来了解如何选择索引策略：

- 可以完全匹配。  
场景：有两个索引可用，一个用于col1 & col2，另一个单独用于col1。  
在上面的场景中，第二个索引会比第一个索引更好，因为它会使scan的较少索引数据。
- 如果有多个候选多列索引，则选择具有较少索引列的索引。  
场景：有两个索引可用，一个用于col1 & col2，另一个用于col1 & col2 & col3。  
在这种情况下，使用col1和col2上的索引，因为它会使scan的较少索引数据。

#### 📖 说明

- 基于索引查询时索引的状态必须为ACTIVE（可通过调用listIndices() API查看索引的状态）。
- 为了保证基于索引查询数据的正确性，用户应确保索引数据与用户数据的一致性。
- 使用以下命令可通过HBase shell客户端执行复杂查询（假定此时已为指定列建立索引）。  
**scan 'tablename', {FILTER => "SingleColumnValueFilter(family, qualifier, compareOp, comparator, filterIfMissing, latestVersionOnly)"}**  
例如：**scan 'test', {FILTER => "SingleColumnValueFilter('info', 'age', =, 'binary:26', true, true)"}**  
（在以上场景中，用户希望在结果中保存没有查询到的列所在行时，不应该在任何这样的列上创建任何索引，因为如果查询的列不存在于其中时，使用SCVF扫描索引列会过滤出一行。而使用filterIfMissingset为false（这是默认值）的SCVF扫描非索引列时，也将会在结果中返回没有查询列的行。因此，为避免查询结果不一致，建议在为索引列创建SCVF后将filterIfMissing设置为true。）
- 在hbase shell中可以通过以下命令查看为用户数据建立的索引数据。  
**scan 'tablename', {ATTRIBUTES => {'FETCH\_INDEX\_DATA' => 'true'}}**

### 8.6.1.2 批量加载 HBase 数据并生成本地二级索引

#### 场景介绍

HBase本身提供了ImportTsv&LoadIncremental工具来批量加载用户数据。当前提供了HIndexImportTsv来支持加载用户数据的同时可以完成对索引数据的批量加载。HIndexImportTsv继承了HBase批量加载数据工具ImportTsv的所有功能。此外，若在执行HIndexImportTsv工具之前未建表，直接运行该工具，将会在创建表时创建索引，并在生成用户数据的同时生成索引数据。

#### 前提条件

- 已安装客户端，具体请参考[安装客户端](#)章节。
- 已根据业务需要创建具有相应权限的组件业务用户。“机机”用户需要下载keytab文件，“人机”用户第一次登录时需修改密码。

#### 使用 HIndexImportTsv 批量生成 HBase 本地二级索引数据

1. 以客户端安装用户登录安装了客户端的节点。
2. 执行以下命令配置环境变量并认证用户：  
**cd 客户端安装目录**  
**source bigdata\_env**  
**kinit 组件业务用户**（未开启Kerberos认证的集群请跳过该操作）

3. 将数据导入到HDFS中。

```
hdfs dfs -mkdir <inputdir>
```

```
hdfs dfs -put <local_data_file> <inputdir>
```

例如定义数据文件“data.txt”，内容如下：

```
12005000201,Zhang San,Male,19,City a, Province a
12005000202,Li Wanting,Female,23,City b, Province b
12005000203,Wang Ming,Male,26,City c, Province c
12005000204,Li Gang,Male,18,City d, Province d
12005000205,Zhao Enru,Female,21,City e, Province e
12005000206,Chen Long,Male,32,City f, Province f
12005000207,Zhou Wei,Female,29,City g, Province g
12005000208,Yang Yiwen,Female,30,City h, Province h
12005000209,Xu Bing,Male,26,City i, Province i
12005000210,Xiao Kai,Male,25,City j, Province j
```

执行以下命令：

```
hdfs dfs -mkdir /datadiriImport
```

```
hdfs dfs -put data.txt /datadiriImport
```

4. 执行以下命令创建表bulkTable：

```
hbase shell
```

```
create 'bulkTable', {NAME => 'info',COMPRESSION => 'SNAPPY',
DATA_BLOCK_ENCODING => 'FAST_DIFF'},{NAME=>'address'}
```

命令执行完成后执行!**quit**退出**hbase shell**。

5. 执行如下命令，生成HFile文件（StoreFiles）：

```
hbase org.apache.hadoop.hbase.index.mapreduce.HIndexImportTsv -
Dimporttsv.separator=<separator>
```

```
-Dimporttsv.bulk.output=</path/for/output> -
```

```
Dindexspecs.to.add=<indexspecs> -Dimporttsv.columns=<columns>
tableName <inputdir>
```

- **-Dimport.separator**：分隔符，例如，**-Dimport.separator='|'**。
- **-Dimport.bulk.output=</path/for/output>**：表示执行结果输出路径，需指定一个不存在的路径。
- **<columns>**：表示导入数据在表中的对应关系，例如，**-Dimporttsv.columns=HBASE\_ROW\_KEY,info:name,info:gender,info:age,address:city,address:province**。
- **<tablename>**：表示要操作的表名。
- **<inputdir>**：表示要批量导入的数据目录。
- **-Dindexspecs.to.add=<indexspecs>**：表示索引名与列的映射，例如**-Dindexspecs.to.add='index\_bulk=>info:[age->String]'**。其构成如下所示：  
**indexNameN=>familyN :[columnQualifierN-> columnQualifierDataType],  
[columnQualifierM-> columnQualifierDataType];familyM:  
[columnQualifierO-> columnQualifierDataType]# indexNameN=>  
familyM: [columnQualifierO-> columnQualifierDataType]**

其中：

- 列限定符用逗号（,）分隔，例如：  
**index1 => f1:[c1-> String], [c2-> String]**
- 列族由分号（;）分隔，例如：  
**index1 => f1:[c1-> String], [c2-> String]; f2:[c3-> Long]**

- 多个索引由#号分隔，例如：  
index1 => f1:[c1-> String], [c2-> String]; f2:[c3-> Long] # index2 => f2:[c3-> Long]
- 列限定的数据类型  
可用的数据类型有：STRING，INTEGER，FLOAT，LONG，DOUBLE，SHORT，BYTE，CHAR。

#### 📖 说明

数据类型也可以用小写传递。

例如执行以下命令：

```
hbase org.apache.hadoop.hbase.hindex.mapreduce.HIndexImportTsv -
Dimporttsv.separator=', ' -Dimporttsv.bulk.output=/dataOutput -
Dindexspecs.to.add='index_bulk=>info:[age->String]' -
Dimporttsv.columns=HBASE_ROW_KEY,info:name,info:gender,info:age,address:
city,address:province bulkTable /datadirImport/data.txt
```

输出：

```
[root@shap000000406 opt]# hbase org.apache.hadoop.hbase.hindex.mapreduce.HIndexImportTsv -
Dimporttsv.separator=', ' -Dimporttsv.bulk.output=/dataOutput -Dindexspecs.to.add='index_bulk=>info:
[age->String]' -
Dimporttsv.columns=HBASE_ROW_KEY,info:name,info:gender,info:age,address:city,address:province
bulkTable /datadirImport/data.txt
2018-05-08 21:29:16,059 INFO [main] mapreduce.HFileOutputFormat2: Incremental table bulkTable
output configured.
2018-05-08 21:29:16,069 INFO [main] client.ConnectionManager$HConnectionImplementation:
Closing master protocol: MasterService
2018-05-08 21:29:16,069 INFO [main] client.ConnectionManager$HConnectionImplementation:
Closing zookeeper sessionId=0x80007c2cb4fd5b4d
2018-05-08 21:29:16,072 INFO [main] zookeeper.ZooKeeper: Session: 0x80007c2cb4fd5b4d closed
2018-05-08 21:29:16,072 INFO [main-EventThread] zookeeper.ClientCnxn: EventThread shut down
for session: 0x80007c2cb4fd5b4d
2018-05-08 21:29:16,379 INFO [main] client.ConfiguredRMFailoverProxyProvider: Failing over to 147
2018-05-08 21:29:17,328 INFO [main] input.FileInputFormat: Total input files to process : 1
2018-05-08 21:29:17,413 INFO [main] mapreduce.JobSubmitter: number of splits:1
2018-05-08 21:29:17,430 INFO [main] Configuration.deprecation: io.bytes.per.checksum is
deprecated. Instead, use dfs.bytes-per-checksum
2018-05-08 21:29:17,687 INFO [main] mapreduce.JobSubmitter: Submitting tokens for job:
job_1525338489458_0002
2018-05-08 21:29:18,100 INFO [main] impl.YarnClientImpl: Submitted application
application_1525338489458_0002
2018-05-08 21:29:18,136 INFO [main] mapreduce.Job: The url to track the job: http://
shap000000407:8088/proxy/application_1525338489458_0002/
2018-05-08 21:29:18,136 INFO [main] mapreduce.Job: Running job: job_1525338489458_0002
2018-05-08 21:29:28,248 INFO [main] mapreduce.Job: Job job_1525338489458_0002 running in uber
mode : false
2018-05-08 21:29:28,249 INFO [main] mapreduce.Job: map 0% reduce 0%
2018-05-08 21:29:38,344 INFO [main] mapreduce.Job: map 100% reduce 0%
2018-05-08 21:29:51,421 INFO [main] mapreduce.Job: map 100% reduce 100%
2018-05-08 21:29:51,428 INFO [main] mapreduce.Job: Job job_1525338489458_0002 completed
successfully
2018-05-08 21:29:51,523 INFO [main] mapreduce.Job: Counters: 50
```

6. 执行如下命令将生成的HFile导入HBase中：

```
hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles </
path/for/output> <tablename>
```

例如执行以下命令：

```
hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles /
dataOutput bulkTable
```

输出：

```
[root@shap000000406 opt]# hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles /
dataOutput bulkTable
2018-05-08 21:30:01,398 WARN [main] mapreduce.LoadIncrementalHFiles: Skipping non-directory
hdfs://hacluster/dataOutput/_SUCCESS
2018-05-08 21:30:02,006 INFO [LoadIncrementalHFiles-0] hfile.CacheConfig: Created cacheConfig:
CacheConfig:disabled
2018-05-08 21:30:02,006 INFO [LoadIncrementalHFiles-2] hfile.CacheConfig: Created cacheConfig:
CacheConfig:disabled
2018-05-08 21:30:02,006 INFO [LoadIncrementalHFiles-1] hfile.CacheConfig: Created cacheConfig:
CacheConfig:disabled
2018-05-08 21:30:02,085 INFO [LoadIncrementalHFiles-2] compress.CodecPool: Got brand-new
decompressor [.snappy]
2018-05-08 21:30:02,120 INFO [LoadIncrementalHFiles-0] mapreduce.LoadIncrementalHFiles: Trying
to load hfile=hdfs://hacluster/dataOutput/address/042426c252f74e859858c7877b95e510
first=12005000201 last=12005000210
2018-05-08 21:30:02,120 INFO [LoadIncrementalHFiles-2] mapreduce.LoadIncrementalHFiles: Trying
to load hfile=hdfs://hacluster/dataOutput/info/f3995920ae0247a88182f637aa031c49
first=12005000201 last=12005000210
2018-05-08 21:30:02,128 INFO [LoadIncrementalHFiles-1] mapreduce.LoadIncrementalHFiles: Trying
to load hfile=hdfs://hacluster/dataOutput/d/c53b252248af42779f29442ab84f86b8 first=\x00index_bulk
\x00\x00\x00\x00\x00\x00\x00\x0018\x00\x0012005000204 last=\x00index_bulk
\x00\x00\x00\x00\x00\x00\x00\x0032\x00\x0012005000206
2018-05-08 21:30:02,231 INFO [main] client.ConnectionManager$HConnectionImplementation:
Closing master protocol: MasterService
2018-05-08 21:30:02,231 INFO [main] client.ConnectionManager$HConnectionImplementation:
Closing zookeeper sessionId=0x81007c2cf0f55cc5
2018-05-08 21:30:02,235 INFO [main] zookeeper.ZooKeeper: Session: 0x81007c2cf0f55cc5 closed
2018-05-08 21:30:02,235 INFO [main-EventThread] zookeeper.ClientCnxn: EventThread shut down
for session: 0x81007c2cf0f55cc5
```

### 8.6.1.3 使用 TableIndexer 工具生成 HBase 本地二级索引

#### 场景介绍

为了快速对数据创建索引，HBase提供了可通过MapReduce功能创建索引的TableIndexer工具，该工具可实现添加、构建和删除索引。具体使用场景如下：

- 在表中预先存在大量数据的情况下，可能希望在某个列上添加索引。但是，使用addIndicesWithData() API添加索引会生成与相关数据对应的索引数据，这将花费大量时间。另一方面，使用addIndices()创建的索引不会构建与表数据对应的索引数据。因此，可以使用TableIndexer工具来完成索引的构建。
- 如果索引数据与表数据不一致，该工具可用于重新构建索引数据。  
如果暂时禁用索引并且在此期间，向禁用的索引列执行新的put操作，直接将索引从禁用状态启用可能会导致索引数据与用户数据不一致。因此，必须注意在再次使用之前重新构建所有索引数据。
- 对于大量现有的索引数据，可以使用TableIndexer工具将索引数据从表中完全删除。
- 对于未建立索引的表，该工具允许用户同时添加和构建索引。

#### TableIndexer 工具使用方法

- 添加新的索引到用户表

命令如下所示：

```
hbase org.apache.hadoop.hbase.hindex.mapreduce.TableIndexer -
Dtablename.to.index=tablename -Dindexspecs.to.add='idx_0=>cf_0:[q_0-
>string],[q_1];cf_1:[q_2],[q_3]#idx_1=>cf_1:[q_4]'
```

相关参数如下：

- **tablename.to.index**: 表示创建索引的表名称。
- **indexspecs.to.add**: 表示索引名对应表的列的映射关系。
- **scan.caching** (可选): 包含一个整数值, 表示在扫描数据表时将传递给扫描器的缓存行数。

上述命令中的参数描述如下:

- **idx\_1**: 表示索引名称
- **cf\_0**: 表示列族名称
- **q\_0**: 表示列名称
- **string**: 表示数据类型, 支持STRING, INTEGER, FLOAT, LONG, DOUBLE, SHORT, BYTE或CHAR。

#### 📖 说明

- '#'用于分隔索引, ';'用于分隔列族, ','用于分隔列限定符。
- 列名及其数据类型应包含在'[]'中。
- 列名及其数据类型通过' ->'分隔。
- 如果未指定具体列的数据类型, 则使用默认数据类型 (string)。
- 如果未设置可选参数scan.caching, 则将采用默认值1000。
- 用户表必须存在。
- 表中指定的索引不能存在。
- 如果用户表中已经存在名称为'd'的ColumnFamily, 则用户必须使用TableIndexer工具构建索引数据。

在执行以上的命令之后, 指定的索引将被添加到表中并且将处于INACTIVE状态。该行为与addIndices() API类似。

- **为用户表中的现有索引构建索引数据**

该命令如下:

```
hbase org.apache.hadoop.hbase.hindex.mapreduce.TableIndexer -
Dtablename.to.index=tablename -Dindexnames.to.build='idx_0 # idx_1'
```

相关参数如下:

- **tablename.to.index**: 表示创建索引的表的名称
- **indexspecs.to.build**: 表示与索引名称
- **scan.caching** (可选): 包含一个整数值, 表示在扫描数据表时将传递给扫描器的缓存行数

上述命令中的参数描述如下:

- **idx\_1**: 表示索引名称

#### 📖 说明

- '#'用于分隔索引名称。
- 如果未设置可选参数scan.caching, 则将采用默认值1000。
- 用户表必须存在。

在执行以上的命令之后, 指定的索引将被设置为ACTIVE状态。用户扫描数据时可以使用它们。

- **从用户表中删除现有索引及其数据**

该命令如下:

```
hbase org.apache.hadoop.hbase.hindex.mapreduce.TableIndexer -
Dtablename.to.index=tablename -Dindexnames.to.drop='idx_0 # idx_1'
```

相关参数如下：

- **tablename.to.index**：表示创建索引的表的名称
- **indexnames.to.drop**：表示应该和其数据一起删除的索引的名称（必须存在于表中）
- **scan.caching**（可选）：其中包含一个整数值，指示在扫描数据表时将传递给扫描器的缓存行数

上述命令中的参数描述如下：

- **idx\_1**：表示索引名称

#### 📖 说明

- '#'用于分隔索引名称。
- 如果未设置可选参数scan.caching，则将采用默认值1000。
- 用户表必须存在。

在执行前面的命令之后，指定的索引将从表中删除。

- **为用户表添加新的索引以及基于现有数据的数据构建**

该命令如下：

```
hbase org.apache.hadoop.hbase.hindex.mapreduce.TableIndexer -
Dtablename.to.index=tablename -Dindexspecs.to.add='idx_0 => cf_0:
[q_0-> string],[q_1];cf_1:[q_2], [q_3] # idx_1 => cf_1:[q_4]' -
Dindexnames.to.build='idx_0'
```

#### 📖 说明

- 用户表必须存在。
- **indexspecs.to.add**中指定的索引不能已存在于表中。
- **indexnames.to.build**中指定的索引名称必须已经存在于表中，或者应该是**indexspecs.to.add**的一部分。

在执行前面的命令之后，indexspecs.to.add中指定的所有索引都将添加到该表中，并且将为通过indexnames.to.build为指定的所有索引构建索引数据。

## 8.6.1.4 迁移 HBase 索引数据

### 操作场景

MRS 1.7及其以后版本中使用的索引与以前MRS版本中HBase使用的二级索引都不兼容。因此，为了将索引数据从以前的版本（MRS 1.5及其以前版本）迁移到MRS 1.7及其以后版本，需要遵循以下步骤。

### 前提条件

1. 迁移数据时旧版本集群应为MRS1.5及其以前的版本，新版本集群应为MRS1.7及其以后的版本。
2. 迁移数据前用户应该有旧的索引数据。
3. 安全集群需配置跨集群互信和启用集群间拷贝功能，普通集群仅需启用集群间拷贝功能。详情请参见[配置跨集群互信](#)。



## 操作步骤

把旧集群中的用户数据迁移至新集群中。迁移数据需单表手动同步新旧集群的数据，通过Export、distcp、Import来完成。

例如，当前旧集群有用户表（t1，索引名为idx\_t1）及其对应的索引表（t1\_idx）。迁移数据的操作步骤如下：

1. 从旧集群导出表中数据。

```
hbase org.apache.hadoop.hbase.mapreduce.Export -Dhbase.mapreduce.include.deleted.rows=true
<tableName> <path/for/data>
```

- <tableName>: 指的是表名。例如，t1。
- <path/for/data>: 指的是保存源数据的路径，例如“/user/hbase/t1”。

例如，**hbase org.apache.hadoop.hbase.mapreduce.Export -Dhbase.mapreduce.include.deleted.rows=true t1 /user/hbase/t1**

2. 把导出的数据按如下步骤复制到新集群中。

```
hadoop distcp <path/for/data> hdfs://ActiveNameNodeIP:9820/<path/for/newData>
```

- <path/for/data>: 指的是旧集群保存源数据的路径。例如，/user/hbase/t1。
- <path/for/newData>: 指的是新集群保存源数据的路径。例如，/user/hbase/t1。

其中，ActiveNameNodeIP是新集群中主NameNode节点的IP地址。

例如，**hadoop distcp /user/hbase/t1 hdfs://192.168.40.2:9820/user/hbase/t1**

### 说明

- 可手动把导出的数据复制到新集群HDFS中，如上路径：“/user/hbase/t1”。
3. 使用新集群HBase表用户，在新集群中生成HFiles。  

```
hbase org.apache.hadoop.hbase.mapreduce.Import -Dimport.bulk.output=<path/for/hfiles>
<tableName><path/for/newData>
```

    - <path/for/hfiles>: 指的是新集群生成HFiles的路径。例如，/user/hbase/output\_t1。
    - <tableName>: 指的是表名。例如，t1。
    - <path/for/newData>: 指的是新集群保存源数据的路径。例如，/user/hbase/t1。

例如，

**hbase org.apache.hadoop.hbase.mapreduce.Import -Dimport.bulk.output=/user/hbase/output\_t1 t1 /user/hbase/t1**

4. 把生成的HFiles导入新集群相应表中。

命令如下：

```
hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles <path/for/hfiles> <tableName>
```

- <path/for/hfiles>: 指的是新集群生成HFiles的路径。例如，/user/hbase/output\_t1。
- <tableName>: 指的是表名。例如，t1。

例如，

**hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles /user/hbase/output\_t1 t1**

### 📖 说明

1. 以上为迁移用户数据的过程，旧集群的索引数据迁移只需按照前三步操作，并更改相应表名为索引表名（如，t1\_idx）。
  2. 迁移索引数据时无需执行4。
5. 向新集群表中导入索引数据。
- a. 在新集群的用户表中添加与之前版本用户表相同的索引（名称为'd'的列族不应该已经存在于用户表中）。

命令如下所示：

```
hbase org.apache.hadoop.hbase.index.mapreduce.TableIndexer -
Dtablename.to.index=<tableName> -Dindexspecs.to.add=<indexspecs>
```

- -Dtablename.to.index=<tableName>：指的是表名。例如， -Dtablename.to.index=t1。
- -Dindexspecs.to.add=<indexspecs>：指的是索引名与列的映射，例如-Dindexspecs.to.add='idx\_t1=>info:[name->String]'。

例如，

```
hbase org.apache.hadoop.hbase.index.mapreduce.TableIndexer -
Dtablename.to.index=t1 -Dindexspecs.to.add='idx_t1=>info:[name->
>String]'
```

### 📖 说明

如果用户表中已经存在名称为'd'的ColumnFamily，则用户必须使用TableIndexer工具构建索引数据。

- b. 运行LoadIncrementalHFiles工具加载索引数据，将旧集群索引数据加载到新集群表中。

命令如下：

```
hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles </path/for/hfiles>
<tableName>
```

- </path/for/hfiles>：指的是索引数据在HDFS上的路径（其为-Dimport.bulk.output中指定的索引生成路径）。例如， /user/hbase/output\_t1\_idx。
- <tableName>：指的是新集群中表名，例如， t1。

例如，

```
hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles /
user/hbase/output_t1_idx t1
```

## 8.6.2 增强 HBase BulkLoad 工具数据迁移能力

### 8.6.2.1 使用 BulkLoad 工具批量导入 HBase 数据

#### 操作场景

您可以按照自定义的方式，通过命令批量导入数据到HBase中并创建索引。

您可以在“configuration.xml”文件中定义多个方式来批量导入数据。导入数据时可不创建索引。

## 📖 说明

- 列的名称不能包含特殊字符，只能由字母、数字和下划线组成。
- 大任务下MapReduce任务运行失败，请参考[MapReduce任务运行失败，ApplicationMaster出现物理内存溢出异常](#)进行处理。
- BulkLoad支持的数据源格式为带分隔符的文本文件。
- 已安装客户端。例如安装目录为“/opt/hadoopclient”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- 若批量导入数据时创建二级索引，还需注意：
  - 当将列的类型设置为string时，不能设置其长度。例如“<column index="1" type="string" length="1" >COLOUMN\_1</column>”，此类型不支持。
  - 当将列的类型设置为date时，不能设置其日期格式。例如“<column index="13" type="date" format="yyyy-MM-dd hh:mm:ss">COLOUMN\_13</column>”，此类型不支持。
  - 不能针对组合列建立二级索引。

## 使用 BulkLoad 工具批量导入 HBase 数据

**步骤1** 以客户端安装用户，登录安装客户端的节点。

**步骤2** 执行以下命令切换到客户端目录。

```
cd /opt/client
```

**步骤3** 执行以下命令配置环境变量。

```
source bigdata_env
```

**步骤4** 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建HBase表的权限和HDFS的操作权限：

```
kinit 组件业务用户
```

如果当前集群未启用Kerberos认证，则执行以下命令设置Hadoop用户名：

```
export HADOOP_USER_NAME=hbase
```

**步骤5** 将数据导入到HDFS中。

```
hdfs dfs -mkdir<inputdir>
```

```
hdfs dfs -put<local_data_file> <inputdir>
```

例如定义数据文件“data.txt”，内容如下：

```
001,Hadoop,citya
002,HBaseFS,cityb
003,HBase,cityc
004,Hive,cityd
005,Streaming,citye
006,MapReduce,cityf
007,Kerberos,cityg
008,LdapServer,cityh
```

执行以下命令：

```
hdfs dfs -mkdir /datadirImport
```

```
hdfs dfs -put data.txt /datadirImport
```

**步骤6** 进入 **hbase shell**，创建表 **ImportTable** 并创建 “configuration.xml” 文件（该文件可以参考模板文件进行编辑，模板文件获取路径为：“/opt/client/HBase/hbase/conf/import.xml.template”）。

例如执行以下命令建表：

```
create 'ImportTable', {NAME => 'f1', COMPRESSION => 'SNAPPY',
DATA_BLOCK_ENCODING => 'FAST_DIFF'}, {NAME => 'f2'}
```

例如自定义导入模板文件 “configuration.xml”，内容如下：

### 📖 说明

- column\_num 要和数据文件中的列的数量对应。
- family 的指定要和表的列族名称对应。
- 仅当批量导入数据时创建二级索引才需配置以下参数，且索引类型的首字母需要大写，例如 **type="String"**；以下片段中 **length="30"** 表示索引列 “H\_ID” 的列值不能超过 30 个字符：

```
<indices>
 <index name="IDX1">
 <index_column family="f1">
 <qualifier type="String" length="30">H_ID</qualifier>
 </index_column>
 </index>
</indices>
```

```
<?xml version="1.0" encoding="UTF-8"?>
<configuration>
 <import id="first" column_num="3">
 <columns>
 <column index="1" type="int">SMS_ID</column>
 <column index="2" type="string">SMS_NAME</column>
 <column index="3" type="string">SMS_ADDRESS</column>
 </columns>
 <rowkey>
 SMS_ID+_'+substring(SMS_NAME,1,4)+'_'+reverse(SMS_ADDRESS)
 </rowkey>
 <qualifiers>
 <normal family="f1">
 <qualifier column="SMS_ID">H_ID</qualifier>
 <qualifier column="SMS_NAME">H_NAME</qualifier>
 <qualifier column="SMS_ADDRESS">H_ADDRESS</qualifier>
 </normal>
 <!-- Define composite columns -->
 <composite family="f2">
 <qualifier class="com.huawei.H_COMBINE_1">H_COMBINE_1</qualifier>
 <columns>
 <column>SMS_ADDRESS</column>
 <column>SMS_NAME</column>
 </columns>
 </composite>
 </qualifiers>
 <indices>
 <index name="IDX1">
 <index_column family="f1">
 <qualifier type="String" length="30">H_ID</qualifier>
 </index_column>
 </index>
 </indices>
 <badlines>SMS_ID < 7000 && SMS_NAME == 'HBase'</badlines>
 </import>
</configuration>
```

**步骤7** 执行如下命令，生成HFile文件。

```
hbase com.huawei.hadoop.hbase.tools.bulkload.ImportData -
Dimport.skip.bad.lines=true-Dimport.separator=<separator>
Dimport.bad.lines.output=</path/badlines/output>-Dimport.hfile.output=</
path/for/output> <configuration xmlfile> <tablename> <inputdir>
```

- -Dimport.skip.bad.lines: 指定值为“false”，表示遇到不适用的行则停止执行。指定值为“true”，表示遇到不适用的数据行则跳过该行继续执行，如果没有在“configuration.xml”中定义不适用行，该参数不需要添加。
- -Dimport.separator: 分隔符，例如，-Dimport.separator=','。
- -Dimport.bad.lines.output=</path/badlines/output>: 指的是不适用的数据行输出路径，如果没有在configuration.xml中定义不适用行，该参数不需要添加。
- -Dimport.hfile.output=< /path/for/output>: 指的是执行结果输出路径。
- <configuration xmlfile>: 指向configuration配置文件。
- <tablename>: 表示要操作的表名。
- <inputdir>: 表示要批量上传的数据目录。

例如执行以下命令：

- **hbase com.huawei.hadoop.hbase.tools.bulkload.ImportData -**  
**Dimport.skip.bad.lines=true -Dimport.separator=',' -**  
**Dimport.bad.lines.output=/badline -Dimport.hfile.output=/hfile**  
**configuration.xml ImportTable /datadirImport**
- **hbase com.huawei.hadoop.hbase.tools.bulkload.IndexImportData -**  
**Dimport.skip.bad.lines=true -Dimport.separator=',' -**  
**Dimport.bad.lines.output=/badline -Dimport.hfile.output=/hfile**  
**configuration\_index.xml IndexImportTable /datadirIndexImport**

#### 须知

- 当HBase已经配置透明加密后，在执行bulkload命令生成HFile时，“-Dimport.hfile.output”指定的HFile路径必须为“/HBase根目录/extdata”的子目录，例如“/hbase/extdata/bulkloadTmp/hfile”。
- 当HBase已经配置透明加密后，执行bulkload命令的HBase用户需要添加到对应集群的hadoop用户组（非FusionInsight Manager下第一个安装的集群，用户组为“c<集群ID>\_hadoop”，例如“c2\_hadoop”），且具有HBase根目录的加密key的读权限。
- 检查目录/tmp/hbase的权限，需要手动添加当前用户对该目录的写权限。

**步骤8** 执行如下命令将HFile导入HBase。

- 批量导入数据：  
**hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles </**  
**path/for/output> <tablename>**
- 批量导入数据时创建二级索引：  
**hbase**  
**org.apache.hadoop.hbase.hindex.mapreduce.HIndexLoadIncrementalHFile**  
**s </path/for/output> <tablename>**

例如执行以下命令：

- `hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles /hfile ImportTable`
  - `hbase org.apache.hadoop.hbase.hindex.mapreduce.HIndexLoadIncrementalHFiles /hfile IndexImportTable`
- 结束

## 8.6.2.2 使用 BulkLoad 工具批量更新 HBase 数据

### 操作场景

HBase BulkLoad工具支持根据RowKey的命名规则、RowKey的范围、字段名以及字段值进行批量更新数据。

### 使用 BulkLoad 工具批量更新 HBase 数据

执行如下命令更新从“row\_start”到“row\_stop”的行，并且把输出结果定向到“/output/destdir/”。

```
hbase com.huawei.hadoop.hbase.tools.bulkload.UpdateData
-Dupdate.rowkey.start="row_start"
-Dupdate.rowkey.stop="row_stop"
-Dupdate.hfile.output=/user/output/
-Dupdate.qualifier=f1:c1,f2
-Dupdate.qualifier.new.value=0,a
'table1'
```

- `-Dupdate.rowkey.start="row_start"`：表示开始行号为“row\_start”。
- `-Dupdate.rowkey.stop="row_stop"`：表示结束行号为“row\_stop”。
- `-Dupdate.hfile.output=/user/output/`：表示执行结果输出路径为“/user/output/”。

#### 须知

当HBase已经配置透明加密后，“批量更新”操作注意事项请参考[步骤7](#)。

执行以下命令，加载HFiles：

```
hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles <path/for/output> <tablename>
```

### 注意事项

- 批量更新会把满足条件的行对应的字段值替换为要更新的值。
- 如果要更新的字段上建有索引，批量更新是不允许的。
- 如果不设置执行结果输出文件，默认是（/tmp/updatedata/表名）。

### 8.6.2.3 使用 BulkLoad 工具批量删除 HBase 数据

#### 操作场景

BulkLoad工具支持根据rowkey的取值模式、范围、字段名、字段值对HBase做批量删除。

#### 使用 BulkLoad 工具批量删除 HBase 数据

执行如下命令删除从“row\_start”到“row\_stop”的行，并且把输出结果定向到“/output/destdir/”。

```
hbase com.huawei.hadoop.hbase.tools.bulkload.DeleteData
-Ddelete.rowkey.start="row_start"
-Ddelete.rowkey.stop="row_stop"
-Ddelete.hfile.output="/output/destdir/"
-Ddelete.qualifier="cf1,cf0:vch,cf0:lng:1000"
'table1'
```

- -Ddelete.rowkey.start="row\_start": 表示开始行号为“row\_start”。
- -Ddelete.rowkey.stop="row\_stop": 表示结束行号为“row\_stop”。
- -Ddelete.hfile.output="/output/destdir/": 表示执行结果输出到“/output/destdir/”目录下。
- -Ddelete.qualifier="cf1,cf0:vch,cf0:lng:1000": 表示删除column family cf1中所有列，column family cf0中列为vch的列，column family cf0中列lng中值为1000的列。

#### 须知

当HBase已经配置透明加密后，“批量删除”操作注意事项请参考[步骤7](#)。

执行以下命令，加载HFiles。

```
hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles <path/for/output> <tablename>
```

#### 注意事项

- 如果column qualifier上建有索引，在该字段的批量删除是会失败的，即不允许在建有索引的字段上执行批量删除。
- 如果不设置执行结果输出数据文件（delete.hfile.output），默认是/tmp/deletedata/表名。

### 8.6.2.4 使用 BulkLoad 工具查询 HBase 表的行统计数

#### 操作场景

HBase BulkLoad工具支持根据rowkey的命名规则、rowkey的范围、字段名以及字段值统计符合条件的行数。

## 操作步骤

直接执行如下命令统计满足如下条件的行数。rowkey在从“row\_start”到“row\_stop”的范围，字段“f3:age”的值为“25”，rowkey的前两个字符为“mi”的行数。

```
hbase com.huawei.hadoop.hbase.tools.bulkload.RowCounter -Dcounter.rowkey.start="row_start" -Dcounter.rowkey.stop="row_stop" -Dcounter.qualifier="f3:age:25" -Dcounter.rowkey.value="substring(0,2) == 'mi'" table1
```

- -Dcounter.rowkey.start="row\_start": 表示开始的rowkey为"row\_start"。
- -Dcounter.rowkey.stop="row\_stop": 表示结束的rowkey为"row\_stop"。
- -Dcounter.qualifier="f3:age:25": 表示列族f3中列为age的列值为25。
- -Dcounter.rowkey.value="substring(0,2) == 'mi'": 表示rowkey的值中前两个为mi。

### 说明

如果指定了“row\_start”和“row\_stop”，则统计的为大于等于“row\_start”并且小于“row\_stop”的数据。

## 8.6.2.5 BulkLoad 工具配置文件说明

### 配置自定义的组合 rowkey

使用BulkLoad工具批量导入HBase数据时，支持用户自定义组合rowkey。BulkLoad组合rowkey即通过一些规则将多个列名经过一些自定义处理，组合生成新的rowkey。

### 说明

列的名称不能包含特殊字符，只能由字母、数字和下划线组成。

关于组合rowkey在“configuration.xml”文件中的配置如下所示，该样例中定义组合rowkey为列“SMS\_ID”、“SMS\_NAME”的取第二个字符开始的三个字符以及“SMS\_SERAIL”的反转（各部分用'\_'连接）。

```
<columns>
 <column index="1" type="int">SMS_ID</column>
 <column index="2" type="string">SMS_NAME</column>
 <column index="3" type="string">SMS_ADDRESS</column>
</columns>

<rowkey>
 SMS_ID+'_'+substring(SMS_NAME,1,4)+'_'+reverse(SMS_ADDRESS)
</rowkey>
```

表 8-5 rowkey 字段处理函数

函数原型	描述	示例
format(data,"DataType")	格式化字符串数据。	例如，format(data,"0.000")是指将数据按照"0.000"格式输出。



函数原型	描述	示例
<code>converse(data,"yyyy-MM-dd","yyyyMMdd")</code>	转化日期格式。	例如， <code>converse(data,"yyyy-MM-dd","yyyyMMdd")</code> 是指将日期格式从"yyyy-MM-dd"转化为"yyyyMMdd"。
<code>rand</code>	随机一个整数，只支持int类型。	无
<code>replace(data,"A","B")</code>	数据替换。	例如， <code>replace(data,"A","B")</code> 是指将A用B替换。
<code>reverse(data)</code>	将字符串反转。	例如， <code>reverse(ABC)</code> 将"ABC"反转成"CBA"。
<code>substring(data,Length1,Length2)</code> , or <code>substring(data,Length3)</code>	截取字符串。	例如， <code>substring(data,1,5)</code> , or <code>substring(data,3)</code> 是指将data字符串进行截取[1,5)或[3,data.length)。
<code>to_number("data")</code>	将字符串转化成数值型，支持返回Long类型。	例如， <code>to_number("123")</code> 是指将"123"转化为123，注意当前data必须为数值。

## 配置自定义 rowkey 实现

使用BulkLoad工具批量导入HBase数据时，支持用户自定义的组合rowkey实现。用户可编写rowkey实现代码，导入时根据该代码逻辑进行组合rowkey导入。

配置自定义rowkey实现步骤如下：

**步骤1** 用户编写自定义rowkey的实现类，需要继承接口，该接口所在的Jar包路径为“客户端安装目录/HBase/hbase/lib/hbase-it-bulk-load-\*.jar”：

`[com.huawei.hadoop.hbase.tools.bulkload.RowkeyHandlerInterface]`，

实现接口中方法：

`byte[] getRowkeyBytes(String[] colsValues, RegulationDomain regulation)`

其中：

- 传入参数“colsValues”为原始数据中的一行数据集，每个元素为一列。
- 传入参数“regulation”为配置导入文件信息（一般情况下并不需要使用）。

**步骤2** 将该实现类与其依赖包同时打包成Jar文件，保存到HBase客户端所在节点的任意位置并确保执行命令的用户具有读取和执行该Jar包的权限。

**步骤3** 在执行导入命令时，增加两个参数配置项：

`-Dimport.rowkey.jar= "第二步中Jar包的全路径"`

```
-Dimport.rowkey.class= "用户实现类的全类名"
```

----结束

## 配置自定义组合字段

BulkLoad支持自定义组合字段，把多个列通过追加的方式即多个列串到一块组合成一个列。

### 说明

列的名称不能包含特殊字符，只能由字母、数字和下划线组成。

关于组合字段H\_COMBINE\_1的定义如下所示，该样例中H\_COMBINE\_1由字段“SMS\_ADDRESS”、“SMS\_SNAME”构成。

```
<!-- Define composite columns -->
<composite family="f2">
 <!-- 定义拼接字段的类名，且该类必须在客户应用中不存在 -->
 <qualifier class="com.huawei.H_COMBINE_1">H_COMBINE_1</qualifier>
 <columns>
 <column>SMS_ADDRESS</column>
 <column>SMS_NAME</column>
 </columns>
</composite>
```

## 指定字段数据类型

HBase BulkLoad支持读取原生态数据文件，把数据文件的每个字段映射为HBase定义的字段，并对该字段的数据类型做定义。

您可以在“configuration.xml”文件中定义多个方式来批量导入数据。

### 说明

列的名称不能包含特殊字符，只能由字母、数字和下划线组成。

指定字段数据类型的配置如下所示，该样例中对“SMS\_ID”、“SMS\_NAME”、“SMS\_ADDRESS”列指定数据类型。

```
<columns>
 <column index="1" type="int">SMS_ID</column>
 <column index="2" type="string">SMS_NAME</column>
 <column index="3" type="string">SMS_ADDRESS</column>
</columns>
```

### 说明

支持的数据类型有：short、int、long、float、double、boolean和string。

## 定义不适用的数据行

BulkLoad支持定义不适用数据行的功能，不适用数据行不会存储到HBase中，这些数据会被保存到指定的文件中。

您可以在“configuration.xml”文件中定义多个方式来批量导入数据。

### 说明

列的名称不能包含特殊字符，只能由字母、数字和下划线组成。

定义不适用的行，配置样例如下所示，即SMS\_ID < 7000 && SMS\_NAME == 'HBase':

```
<!-- Define bad line filter rule -->
<badlines>SMS_ID < 7000 && SMS_NAME == 'HBase'</badlines>
```

针对“<badlines>”标签中的算符和对应的参数类型如表8-6所示。

表 8-6 算符和对应的参数类型

算符类型	参数类型
&&	对应的参数类型应为布尔型。
&	对应的参数类型应为整数。
	对应的参数类型应为整数。
^	对应的参数类型应为整数。
/	对应的参数类型应为数字。
==	对应的参数类型应为字符串。
>=	对应的参数类型应为数字。
>	对应的参数类型应为数字。
<<	对应的参数类型应为整数。
<=	对应的参数类型应为数字。
<	对应的参数类型应为数字。
%	对应的参数类型应为数字。
*	对应的参数类型应为数字。
!=	对应的参数类型应为字符串。
	对应的参数类型应为布尔型。
+	对应的参数类型应为数字和字符串。
>>	对应的参数类型应为整数。
-	对应的参数类型应为字符串。
>>>	对应的参数类型应为整数。

### 8.6.3 配置 RSGroup 管理 RegionServer 资源

#### 操作场景

HBase服务的数据节点较多，需要根据不同的业务规模将数据节点资源分配给特定的业务，从而达到资源独占使用的目的。当AZ容灾特性被开启时，为了保证AZ容灾生效，保障业务可靠性，在为RSGroup分配RegionServer时，需遵循分配结果能使该RSGroup在每个AZ下都存在RegionServer实例的规则。

### 📖 说明

本章节内容仅适用于MRS 3.1.2及之后版本。

## 前提条件

- 已登录Manager。
- 登录角色拥有Manager管理员权限。
- 将RSGroup最小节点数设置为下述三种情况的最大值。
  - 为了保证服务的可靠性，RSGroup内的RegionServer节点数量需要配置一定的冗余量，确保冗余节点数  $> (\text{RSGroup内业务表region总数}/2000) * 50\%$ 。
  - 如果系统表在单独的RSGroup，需要确保该RSGroup的节点数量  $> 2$ 。
  - 为了不影响滚动重启功能，如果RegionServer节点总数在300以内，那么单个RSGroup的节点数量不应小于3。如果RegionServer节点总数大于等于300，那么单个RSGroup的节点数量不应小于 $(\text{节点数} * 1\%) + 1$ 。

## 可能的影响

- 由于RSGroup约束了region转移可用的RegionServer节点，如果RSGroup内部分节点故障或者滚动重启，可能会触发region超过阈值的告警，也可能导致业务性能下降。
- 当提交修改RSGroup请求产生大量region转移任务时，如果进行相关RSGroup操作会面临失败。需先观察WebUI页面的region转移情况，等待转移任务结束后再进行后续操作。

## 配置 RSGroup

### 创建RSGroup

**步骤1** 在FusionInsight Manager界面，选择“集群 > 服务 > HBase > RSGroup管理”。

**步骤2** 单击“添加RSGroup”按钮，在弹出的添加RSGroup页面填写新增的RSGroup名称，RSGroup名称包括数字、字母或下划线（\_），长度为1-120个字符。然后单击“确定”。

### 查看RSGroup

**步骤3** 选择待操作的RSGroup，在操作列单击“查看”，即可在弹出框中查看该RSGroup的RegionServers详情和Tables详情。

### 📖 说明

default RSGroup是HBase的默认RSGroup，所有已启动并且未手动添加到其他RSGroup的RegionServer节点都会添加到default RSGroup。

### 修改RSGroup名称

**步骤4** 选择待操作的RSGroup，在操作列单击“修改名称”。在修改RSGroup名称弹出框中填写RSGroup新名称，新名称不能与已存在名称相同，单击“确定”。

### 修改RSGroup

**步骤5** 单击待操作的RSGroup名称，跳转到修改RSGroup页面。

**步骤6** 勾选欲分配的RegionServer实例，单击“下一步”。

### 📖 说明

- 一次分配操作仅允许勾选来自同一RSGroup的一个或多个RegionServer实例，且default组中的RegionServer的运行状态不为良好时不允许被勾选分配。若想要分配来自不同RSGroup的RegionServer实例，请分多次修改操作进行分配。
- 开启跨AZ特性时，分配操作需要保证分配结果能使每个AZ中均存在该RSGroup的RegionServer实例，而且无法对开启前已分配的RSGroup进行AZ约束校验。

**步骤7** 勾选欲分配的表，单击“下一步”。

### 📖 说明

- 一次分配操作仅允许勾选来自同一RSGroup的一个或多个表。若想要分配来自不同RSGroup的RegionServer实例，请分多次修改来进行分配。
- 当修改RSGroup操作中同时勾选了分配RegionServer和表时，RegionServer和表需来自同一RSGroup。
- 当修改RSGroup操作中只勾选了分配表，且分配前该RSGroup下不存在RegionServer，则将修改失败。

**步骤8** 单击“提交”。修改成功后，提示修改结果，页面将跳转至RSGroup列表展示界面。

当提示“任务入队”相关信息时，页面将跳转至RSGroup列表展示界面。此次提交的修改RSGroup请求，已进入任务队列中，请按照界面提示，观察原生界面region转移完成，确认入队任务执行成功，再进行后续操作。

### 删除RSGroup

**步骤9** 在RSGroup管理页面，勾选需要删除的RSGroup，然后选择“删除RSGroup > 确定”。

### 📖 说明

RSGroup删除失败可能原因及解决方法：

- “default”组不允许被删除。
- 该RSGroup中仍包含RegionServer或Table，请将该RSGroup中RegionServer或Table分配给别的RSGroup组后，再进行删除。

----结束

## 8.7 HBase 性能调优

### 8.7.1 提升 HBase BulkLoad 工具批量加载效率

#### 操作场景

批量加载功能采用了MapReduce jobs直接生成符合HBase内部数据格式的文件，然后把生成的StoreFiles文件加载到正在运行的集群。使用批量加载相比直接使用HBase的API会节约更多的CPU和网络资源。

ImportTSV是一个HBase的表数据加载工具。

### 📖 说明

本章节适用于MRS 3.x及之后版本。

## 前提条件

在执行批量加载时需要通过“Dimporttsv.bulk.output”参数指定文件的输出路径。

## 操作步骤

参数入口：执行批量加载任务时，在BulkLoad命令行中加入如下参数。

表 8-7 增强 BulkLoad 效率的配置项

参数	描述	配置的值
- Dimporttsv.map per.class	<p>用户自定义mapper通过把键值对的构造从mapper移动到reducer以帮助提高性能。mapper只需要把每一行的原始文本发送给reducer，reducer解析每一行的每一条记录并创建键值对。</p> <p><b>说明</b> 当该值配置为“org.apache.hadoop.hbase.mapreduce.TsvImporterByteMapper”时，只在执行没有HBASE_CELL_VISIBILITY OR HBASE_CELL_TTL选项的批量加载命令时使用。使用“org.apache.hadoop.hbase.mapreduce.TsvImporterByteMapper”时可以得到更好的性能。</p>	<p>org.apache.hadoop.hbase.mapreduce.TsvImporterByteMapper 和 org.apache.hadoop.hbase.mapreduce.TsvImporterTextMapper</p>

## 8.7.2 提升 HBase 连续 Put 数据场景性能

### 操作场景

对大批量、连续put的场景，配置下面的两个参数为“false”时能大量提升性能。

- “hbase.regionserver.wal.durable.sync”
- “hbase.regionserver.hfile.durable.sync”

当提升性能时，缺点是对DataNode（默认是3个）同时故障时，存在小概率数据丢失的现象。对数据可靠性要求高的场景请慎重配置。

#### 说明

本章节适用于MRS 3.x及之后版本。

### 操作步骤

参数入口：

在FusionInsight Manager系统中，选择“集群 > 待操作集群的名称 > 服务 > HBase > 配置”，单击“全部配置”。在搜索框中输入参数名称，并进行修改。

表 8-8 提升连续 put 场景性能的参数

参数	描述	配置值
hbase.wal.hsync	设置是否启用WAL文件持久性以将WAL数据持久化到磁盘。若将该参数设置为true，则性能将受到影响，原因是每个WAL的编辑都会被hadoop fsync同步到磁盘上。	false
hbase.hfile.hsync	设置是否启用Hfile持久性以将数据持久化到磁盘。若将该参数设置为true，则性能将受到影响，原因是每个Hfile写入时都会被hadoop fsync同步到磁盘上。	false

## 8.7.3 提升 HBase Put 和 Scan 性能综合调优

### 操作场景

HBase有很多与读写性能相关的配置参数。读写请求负载不同的情况下，配置参数需要进行相应的调整，本章节旨在指导用户通过修改RegionServer配置参数进行读写性能调优。

#### 📖 说明

本章节适用于MRS 3.x及之后版本。

### 操作步骤

登录FusionInsight Manager界面，选择“集群 > 服务 > HBase > 配置”，配置以下相关参数以提升HBase数据读写性能调优。

- JVM GC参数

RegionServer GC\_OPTS参数设置建议：

- -Xms与-Xmx设置相同的值，需要根据实际情况设置，增大内存可以提高读写性能，可以参考参数“hfile.block.cache.size”（见表8-10）和参数“hbase.regionserver.global.memstore.size”（见表8-9）的介绍进行设置。
- -XX:NewSize与-XX:MaxNewSize设置相同值，建议低负载场景下设置为“512M”，高负载场景下设置为“2048M”。
- -XX:CMSInitiatingOccupancyFraction建议设置为“100 \* (hfile.block.cache.size + hbase.regionserver.global.memstore.size + 0.05)”，最大值不超过90。
- -XX:MaxDirectMemorySize表示JVM使用的堆外内存，建议低负载情况下设置为“512M”，高负载情况下设置为“2048M”。

#### 📖 说明

GC\_OPTS参数中-XX:MaxDirectMemorySize默认没有配置，如需配置，用户可在GC\_OPTS参数中自定义添加。

- Put相关参数
  - RegionServer处理put请求的数据，会将数据写入memstore和hlog，
    - 当memstore大小达到设置的“hbase.hregion.memstore.flush.size”参数值大小时，memstore就会刷新到HDFS生成HFile。
    - 当当前region的列簇的HFile数量达到“hbase.hstore.compaction.min”参数值时会触发compaction。
    - 当当前region的列簇HFile数达到“hbase.hstore.blockingStoreFiles”参数值时会阻塞memstore刷新生成HFile的操作，导致put请求阻塞。

表 8-9 Put 相关参数

参数	描述	默认值
hbase.wal.hsync	每一条wal是否持久化到硬盘。 参考 <a href="#">提升HBase连续Put数据场景性能</a> 。	true
hbase.hfile.hsync	hfile写是否立即持久化到硬盘。 参考 <a href="#">提升HBase连续Put数据场景性能</a> 。	true
hbase.hregion.memstore.flush.size	若MemStore的大小（单位：Byte）超过指定值，MemStore将被冲洗至磁盘。该参数值将被运行每个hbase.server.thread.wakefrequency的线程所检验。建议设置为HDFS块大小的整数倍，在内存足够put负载大情况下可以调整增大。	134217728
hbase.regionserver.global.memstore.size	更新被锁定以及强制冲洗发生之前一个RegionServer上支持的所有MemStore的大小。建议设置为“hbase.hregion.memstore.flush.size * 写活跃region数 / RegionServer GC -Xmx”。默认值为“0.4”，表示使用RegionServer GC -Xmx的40%。	0.4
hbase.hstore.flusher.count	memstore的flush线程数，在put高负载场景下可以适当调大。	2
hbase.regionserver.thread.compaction.small	小压缩线程数，在put高负载情况下可以适当调大。	10



参数	描述	默认值
hbase.hstore.blockingStoreFiles	若一个Store内的HStoreFile文件数量超过指定值，则针对此HRegion的更新将被锁定直到一个压缩完成或者base.hstore.blockingWaitTime被超过。每冲洗一次MemStore一个StoreFile文件被写入。在put高负载场景下可以适当调大。	15

- Scan相关参数

表 8-10 Scan 相关参数

参数	描述	默认值
hbase.client.scanner.timeout.period	客户端和RegionServer端参数，表示客户端执行scan的租约超时时间。建议设置为60000ms的整数倍，在读高负载情况下可以适当调大。单位：毫秒。	60000
hfile.block.cache.size	数据缓存所占的RegionServer GC -Xmx百分比，在读高负载情况下可以适当调大以增大缓存命中率以提高性能。表示分配给HFile/StoreFile所使用的块缓存的最大heap（-Xmx setting）的百分比。	当offheap关闭时，默认值为0.25，当offheap开启时，默认值是0.1。

- Handler相关参数

表 8-11 Handler 相关参数

参数	描述	默认值
hbase.regionserver.handler.count	RegionServer上的RPC侦听器实例数，建议设置为200 ~ 400之间。	200
hbase.regionserver.metahandler.count	RegionServer中处理优先请求的程序实例的数量，建议设置为200 ~ 400之间。	200

## 8.7.4 提升 HBase 实时写数据效率

### 操作场景

需要把数据实时写入到HBase中或者对于大批量、连续put的场景。

#### 📖 说明

本章节适用于MRS 3.x及之后版本。

### 前提条件

调用HBase的put或delete接口，把数据保存到HBase中。

### 操作步骤

- **写数据服务端调优**  
参数入口：登录FusionInsight Manager，选择“集群 > 服务 > HBase > 配置 > 全部配置”，进入HBase服务参数“全部配置”界面，具体操作请参考[修改集群服务配置参数](#)章节。

表 8-12 影响实时写数据配置项

配置参数	描述	默认值
hbase.wal.hsync	控制HLog文件在写入到HDFS时的同步程度。如果为true，HDFS在把数据写入到硬盘后才返回；如果为false，HDFS在把数据写入OS的缓存后就返回。 把该值设置为false比true在写入性能上会更优。	true
hbase.hfile.hsync	控制HFile文件在写入到HDFS时的同步程度。如果为true，HDFS在把数据写入到硬盘后才返回；如果为false，HDFS在把数据写入OS的缓存后就返回。 把该值设置为false比true在写入性能上会更优。	true

配置参数	描述	默认值
GC_OPTS	<p>HBase利用内存完成读写操作。提高HBase内存可以有效提高HBase性能。GC_OPTS主要需要调整HeapSize的大小和NewSize的大小。调整HeapSize大小的时候，建议将Xms和Xmx设置成相同的值，这样可以避免JVM动态调整HeapSize大小的时候影响性能。调整NewSize大小的时候，建议把其设置为HeapSize大小的1/8。</p> <ul style="list-style-type: none"> <li>• HMaster：当HBase集群规模越大、Region数量越多时，可以适当调大HMaster的GC_OPTS参数。</li> <li>• RegionServer：RegionServer需要的内存一般比HMaster要大。在内存充足的情况下，HeapSize可以相对设置大一些。</li> </ul> <p><b>说明</b> 主HMaster的HeapSize为4G的时候，HBase集群可以支持100000 region数的规模。根据经验值，集群每增加35000个region，HeapSize增加2G，主HMaster的HeapSize不建议超过32GB。</p>	<ul style="list-style-type: none"> <li>• HMaster -server - Xms4G - Xmx4G - XX:NewSize= 512M - XX:MaxNewSi ze=512M - XX:Metaspac eSize=128M - XX:MaxMetas paceSize=512 M - XX:+UseConc MarkSweepG C - XX:+CMSPara llelRemarkEn abled - XX:CMSInitiat ingOccupancy Fraction=65 - XX:+PrintGCD etails - Dsun.rmi.dgc. client.gcInter val=0x7FFFFFF FFFFFFFFFE - Dsun.rmi.dgc. server.gcInter val=0x7FFFFFF FFFFFFFFFE - XX:- OmitStackTra ceInFastThro w - XX:+PrintGCT imeStamps - XX:+PrintGCD ateStamps - XX:+UseGCLo gFileRotation - XX:NumberO fGLogFiles= 10 - XX:GLogFile Size=1M</li> </ul>

配置参数	描述	默认值
		<ul style="list-style-type: none"> <li>• Region Server</li> <li>-server -</li> <li>Xms6G -</li> <li>Xmx6G -</li> <li>XX:NewSize=1024M -</li> <li>XX:MaxNewSize=1024M -</li> <li>XX:MetaspaceSize=128M -</li> <li>XX:MaxMetaspaceSize=512M -</li> <li>XX:+UseConcMarkSweepGC -</li> <li>XX:+CMSParallelRemarkEnabled -</li> <li>XX:CMSInitiatingOccupancyFraction=65 -</li> <li>XX:+PrintGCDetails -</li> <li>Dsun.rmi.dgc.client.gcInterval=0x7FFFFFFF -</li> <li>Dsun.rmi.dgc.server.gcInterval=0x7FFFFFFF -</li> <li>XX:-OmitStackTraceInFastThrow -</li> <li>XX:+PrintGCTimeStamps -</li> <li>XX:+PrintGCDateStamps -</li> <li>XX:+UseGCLogFileRotation -</li> <li>XX:NumberOfGCLogFiles=10 -</li> <li>XX:GCLogFileSize=1M</li> </ul>

配置参数	描述	默认值
hbase.regionserver.handler.count	<p>表示在RegionServer上启动的RPC侦听器实例数。如果设置过高会导致激烈线程竞争，如果设置过小，请求将会在RegionServer长时间等待，降低处理能力。根据资源情况，适当增加处理线程数。</p> <p>建议根据CPU的使用情况，可以选择设置为100至300之间的值。</p>	200
hbase.hregion.max.filesize	<p>HStoreFile的最大大小（单位：Byte）。若任何一个列族HStoreFile超过此参数值，则托管Hregion将会一分为二。</p>	10737418240
hbase.hregion.memstore.flush.size	<p>在RegionServer中，当写操作内存中存在超过memstore.flush.size大小的memstore，则MemStoreFlusher就启动flush操作将该memstore以hfile的形式写入对应的store中。</p> <p>如果RegionServer的内存充足，而且活跃Region数量也不是很多的时候，可以适当增大该值，可以减少compaction的次数，有助于提升系统性能。</p> <p>同时，这种flush产生的时候，并不是紧急的flush，flush操作可能会有一定延迟，在延迟期间，写操作还可以进行，Memstore还会继续增大，最大值为“memstore.flush.size” * “hbase.hregion.memstore.block.multiplier”。当超过最大值时，将会阻塞操作。适当增大“hbase.hregion.memstore.block.multiplier”可以减少阻塞，减少性能波动。单位：字节。</p>	134217728

配置参数	描述	默认值
hbase.regionserver.global.memstore.size	<p>更新被锁定以及强制冲洗发生之前一个RegionServer上支持的所有MemStore的大小。RegionServer中，负责flush操作的是MemStoreFlusher线程。该线程定期检查写操作内存，当写操作占用内存总量达到阈值，MemStoreFlusher将启动flush操作，按照从大到小的顺序，flush若干相对较大的memstore，直到所占用内存小于阈值。</p> <p>阈值 = “hbase.regionserver.global.memstore.size” * “hbase.regionserver.global.memstore.size.lower.limit” * “HBase_HEAPSIZE”</p> <p><b>说明</b> 该配置与“hfile.block.cache.size”的和不能超过0.8，也就是写和读操作的内存不能超过HeapSize的80%，这样可以保证除读和写外其它操作的正常运行。</p>	0.4
hbase.hstore.blockingStoreFiles	<p>在region flush前首先判断file文件个数，是否大于hbase.hstore.blockingStoreFiles。</p> <p>如果大于需要先compaction并且让flush延时90s（这个值可以通过hbase.hstore.blockingWaitTime进行配置），在延时过程中，将会继续写从而使得Memstore还会继续增大超过最大值“memstore.flush.size” * “hbase.hregion.memstore.block.multiplier”，导致写操作阻塞。当完成compaction后，可能就会产生大量写入。这样就导致性能激烈震荡。</p> <p>增加hbase.hstore.blockingStoreFiles，可以减低BLOCK机率。</p>	15
hbase.regionserver.thread.compaction.throttle	<p>大于此参数值的压缩将被大线程池执行，单位：Byte。控制一次Minor Compaction时，进行compaction的文件总大小的阈值。Compaction时的文件总大小会影响这一次compaction的执行时间，如果太大，可能会阻塞其它的compaction或flush操作。</p>	1610612736

配置参数	描述	默认值
hbase.hstore.compaction.min	每次执行minor compaction的HStoreFile的最小数量。当一个Store文件超过该值时，会进行compact，适当增大该值，可以减少文件被重复执行compaction。但是如果过大，会导致Store文件数过多而影响读取的性能。	6
hbase.hstore.compaction.max	每次执行minor compaction的HStoreFile的最大数量。与“hbase.hstore.compaction.max.size”的作用基本相同，主要是控制一次compaction操作的时间不要太长。	10
hbase.hstore.compaction.max.size	如果一个HFile文件的大小大于该值，那么在Minor Compaction操作中不会选择这个文件进行compaction操作，除非进行Major Compaction操作。 这个值可以防止较大的HFile参与compaction操作。在禁止Major Compaction后，一个Store中可能存在几个HFile，而不会合并成为一个HFile，这样不会对数据读取造成太大的性能影响。单位：字节。	9223372036854775807
hbase.hregion.majorcompaction	单个区域内所有HStoreFile文件主压缩的时间间隔，单位：毫秒。由于执行Major Compaction会占用较多的系统资源，如果正在处于系统繁忙时期，会影响系统的性能。 如果业务没有较多的更新、删除、回收过期数据空间时，可以把该值设置为0，以禁止Major Compaction。 如果必须要执行Major Compaction，以回收更多的空间，可以适当增加该值来调大Major Compaction的执行周期，减少对资源的频繁占用。单位：毫秒。	604800000

配置参数	描述	默认值
<ul style="list-style-type: none"> <li>hbase.regionserver.maxlogs</li> <li>hbase.regionserver.hlog.blocksize</li> </ul>	<ul style="list-style-type: none"> <li>表示一个RegionServer上未进行Flush的Hlog的文件数量的阈值，如果大于该值，RegionServer会强制进行flush操作。</li> <li>表示每个HLog文件的最大大小。如果HLog文件大小大于该值，就会滚动出一个新的HLog文件，旧的将被禁用并归档。</li> </ul> <p>这两个参数共同决定了RegionServer中可以存在的未进行Flush的hlog数量。当这个数据量小于MemStore的总大小的时候，会出现由于HLog文件过多而触发的强制flush操作。这个时候可以适当调整这两个参数的大小，以避免出现这种强制flush的情况。单位：字节。</p>	<ul style="list-style-type: none"> <li>32</li> <li>134217728</li> </ul>

• **写数据客户端调优**

写数据时，在场景允许的情况下，需要使用Put List的方式，可以极大的提升写性能。每一次Put的List的长度，需要结合单条Put的大小，以及实际环境的一些参数进行设定。建议在选定之前先做一些基础的测试。

• **写数据表设计调优**

表 8-13 影响实时写数据相关参数

配置参数	描述	默认值
COMPRESSION	<p>配置数据的压缩算法，这里的压缩是HFile中block级别的压缩。对于可以压缩的数据，配置压缩算法可以有效减少磁盘的IO，从而达到提高性能的目的。</p> <p><b>说明</b> 并非所有数据都可以进行有效压缩。例如一张图片的数据，因为图片一般已经是压缩后的数据，所以压缩效果有限。常用的压缩算法是SNAPPY，因为它有较好的Encoding/Decoding速度和可以接受的压缩率。</p>	NONE
BLOCKSIZE	<p>配置HFile中block块的大小，不同的block块大小，可以影响HBase读写数据的效率。越大的block块，配合压缩算法，压缩的效率就越好；但是由于HBase的读取数据是以block块为单位的，所以越大的block块，对于随机读的情况，性能可能会比较差。</p> <p>如果要提升写入的性能，一般扩大到128KB或者256KB，可以提升写数据的效率，也不会影响太大的随机读性能。单位：字节</p>	65536



配置参数	描述	默认值
IN_MEMORY	配置这个表的数据优先缓存在内存中，这样可以有效提升读取的性能。对于一些小表，而且需要频繁进行读取操作的，可以设置此配置项。	false

## 8.7.5 提升 HBase 实时读数据效率

### 操作场景

需要读取HBase数据场景。

### 前提条件

调用HBase的get或scan接口，从HBase中实时读取数据。

### 操作步骤

- **读数据服务端调优**  
参数入口：登录FusionInsight Manager，选择“集群 > 服务 > HBase > 配置 > 全部配置”，进入HBase服务参数“全部配置”界面，具体操作请参考[修改集群服务配置参数](#)章节。

表 8-14 影响实时读数据配置项

配置参数	描述	默认值
GC_OPTS	<p>HBase利用内存完成读写操作。提高HBase内存可以有效提高HBase性能。</p> <p>GC_OPTS主要需要调整HeapSize的大小和NewSize的大小。调整HeapSize大小的时候，建议将Xms和Xmx设置成相同的值，这样可以避免JVM动态调整HeapSize大小的时候影响性能。调整NewSize大小的时候，建议把其设置为HeapSize大小的1/8。</p> <ul style="list-style-type: none"> <li>• HMaster: 当HBase集群规模越大、Region数量越多时，可以适当调大HMaster的GC_OPTS参数。</li> <li>• RegionServer: RegionServer需要的内存一般比HMaster要大。在内存充足的情况下，HeapSize可以相对设置大一些。</li> </ul> <p><b>说明</b> 主HMaster的HeapSize为4G的时候，HBase集群可以支持100000 region数的规模。根据经验值，集群每增加35000个region，HeapSize增加2G，主HMaster的HeapSize不建议超过32GB。</p>	<p>MRS 3.x之前版本:</p> <ul style="list-style-type: none"> <li>• HMaster: <ul style="list-style-type: none"> <li>-server -</li> <li>Xms2G -</li> <li>Xmx2G -</li> <li>XX:NewSize=256M -</li> <li>XX:MaxNewSize=256M -</li> <li>XX:MetaspaceSize=128M -</li> <li>XX:MaxMetaspaceSize=512M -</li> <li>XX:MaxDirectMemorySize=512M -</li> <li>XX:+UseConcMarkSweepGC -</li> <li>XX:+CMSParallelRemarkEnabled -</li> <li>XX:CMSInitiatingOccupancyFraction=65 -</li> <li>XX:+PrintGCDetails -</li> <li>Dsun.rmi.dgc.client.gcInterval=0x7FFFFFFF -</li> <li>Dsun.rmi.dgc.server.gcInterval=0x7FFFFFFF -</li> <li>XX:-OmitStackTraceInFastThread -</li> <li>XX:+PrintGCTimeStamps -</li> </ul> </li> </ul>

配置参数	描述	默认值
		XX:+PrintGC DateStamps - XX:+UseGCL ogFileRotati on - XX:Number OfGCLogFil es=10 - XX:GCLogFil eSize=1M ● RegionServe r: -server - Xms4G - Xmx4G - XX:NewSize =512M - XX:MaxNew Size=512M - XX:Metaspa ceSize=128 M - XX:MaxMet aspaceSize= 512M - XX:MaxDire ctMemorySi ze=512M - XX:+UseCon cMarkSwee pGC - XX:+CMSPar allelRemark Enabled - XX:CMSIniti atingOccup ancyFractio n=65 - XX:+PrintGC Details - Dsun.rmi.dg c.client.gcln terval=0x7F FFFFFFFFFF FFE - Dsun.rmi.dg c.server.gcln terval=0x7F

配置参数	描述	默认值
		<p>FFFFFFFFF FFE -XX:- OmitStackTr aceInFastTh row - XX:+PrintGC TimeStamps - XX:+PrintGC DateStamps - XX:+UseGCL ogFileRotati on - XX:Number OfGCLogFil es=10 - XX:GCLogFil eSize=1M</p> <p>MRS 3.x及之后 版本：</p> <ul style="list-style-type: none"> <li>• HMaster -server - Xms4G - Xmx4G - XX:NewSize =512M - XX:MaxNew Size=512M - XX:Metaspa ceSize=128 M - XX:MaxMet aspaceSize= 512M - XX:+UseCon cMarkSwee pGC - XX:+CMSPar allelRemark Enabled - XX:CMSIniti atingOccup ancyFractio n=65 - XX:+PrintGC Details - Dsun.rmi.dg</li> </ul>

配置参数	描述	默认值
		<p>c.client.gclnterval=0x7FFFFFFF - Dsun.rmi.dgc.server.gclnterval=0x7FFFFFFF -XX:-OmitStackTracerInFastThrow -XX:+PrintGCTimeStamps -XX:+PrintGCDateStamps -XX:+UseGCLogFileRotation -XX:NumberOfGCLogFiles=10 -XX:GCLogFileSize=1M</p> <ul style="list-style-type: none"> <li>Region Server                     <ul style="list-style-type: none"> <li>-server -Xms6G -Xmx6G -XX:NewSize=1024M -XX:MaxNewSize=1024M -XX:MetaspaceSize=128M -XX:MaxMetaspaceSize=512M -XX:+UseConcMarkSweepGC -XX:+CMSParallelRemarkEnabled -XX:CMSInitiatingOccup</li> </ul> </li> </ul>

配置参数	描述	默认值
		ancyFraction=65 - XX:+PrintGC Details - Dsun.rmi.dgc.client.gcInterval=0x7FFFFFFF - Dsun.rmi.dgc.server.gcInterval=0x7FFFFFFF -XX:- OmitStackTraceInFastThrow - XX:+PrintGC TimeStamps - XX:+PrintGC DateStamps - XX:+UseGCLog Rotation - XX:Number OfGCLogFiles=10 - XX:GCLogFile Size=1M
hbase.regionserver.handler.count	表示RegionServer在同一时刻能够并发处理多少请求。如果设置过高会导致激烈线程竞争，如果设置过小，请求将会在RegionServer长时间等待，降低处理能力。根据资源情况，适当增加处理线程数。 建议根据CPU的使用情况，可以选择设置为100至300之间的值。	200
hfile.block.cache.size	HBase缓存区大小，主要影响查询性能。根据查询模式以及查询记录分布情况来决定缓存区的大小。如果采用随机查询使得缓存区的命中率较低，可以适当降低缓存区大小。	当offheap关闭时，默认值为0.25。当offheap开启时，默认值是0.1。

 说明

如果同时存在读和写的操作，这两种操作的性能会互相影响。如果写入导致的flush和Compaction操作频繁发生，会占用大量的磁盘IO操作，从而影响读取的性能。如果写入导致阻塞较多的Compaction操作，就会出现Region中存在多个HFile的情况，从而影响读取的性能。所以如果读取的性能不理想的时候，也要考虑写入的配置是否合理。

- **读数据客户端调优**

Scan数据时需要设置caching（一次从服务端读取的记录条数，默认是1），若使用默认值读性能会降到极低。

当不需要读一条数据所有的列时，需要指定读取的列，以减少网络IO。

只读取RowKey时，可以为Scan添加一个只读取RowKey的filter（FirstKeyOnlyFilter或KeyOnlyFilter）。

- **读数据表设计调优**

表 8-15 影响实时读数据相关参数

配置参数	描述	默认值
COMPRESSION	配置数据的压缩算法，这里的压缩是HFile中block级别的压缩。对于可以压缩的数据，配置压缩算法可以有效减少磁盘的IO，从而达到提高性能的目的。 <b>说明</b> 并非所有数据都可以进行有效压缩。例如一张图片的数据，因为图片一般已经是压缩后的数据，所以压缩效果有限。常用的压缩算法是SNAPPY，因为它有较好的Encoding/Decoding速度和可以接受的压缩率。	NONE
BLOCKSIZE	配置HFile中block块的大小，不同的block块大小，可以影响HBase读写数据的效率。越大的block块，配合压缩算法，压缩的效率就越好；但是由于HBase的读取数据是以block块为单位的，所以越大的block块，对于随机读的情况，性能可能会比较差。 如果要提升写入的性能，一般扩大到128KB或者256KB，可以提升写数据的效率，也不会影响太大的随机读性能。单位：字节。	65536
DATA_BLOCK_ENCODING	配置HFile中block块的编码方法。当一行数据中存在多列时，一般可以配置为“FAST_DIFF”，可以有效地节省数据存储的空间，从而提供性能。	NONE

## 8.7.6 HBase JVM 参数优化说明

### 操作场景

当集群数据量达到一定规模后，JVM的默认配置将无法满足集群的业务需求，轻则集群变慢，重则集群服务不可用。所以需要根据实际的业务情况进行合理的JVM参数配置，提高集群性能。

### 操作步骤

#### 参数入口：

HBase角色相关的JVM参数需要配置在安装有HBase服务的节点的“\${BIGDATA\_HOME}/FusionInsight\_HD\_\*/install/FusionInsight-HBase-2.2.3/hbase/conf/”目录下的“hbase-env.sh”文件中。

每个角色都有各自的JVM参数配置变量，如表8-16。

表 8-16 HBase 相关 JVM 参数配置变量

变量名	变量影响的角色
HBASE_OPTS	该变量中设置的参数，将影响HBase的所有角色。
SERVER_GC_OPTS	该变量中设置的参数，将影响HBase Server端的所有角色，例如：Master、RegionServer等。
CLIENT_GC_OPTS	该变量中设置的参数，将影响HBase的Client进程。
HBASE_MASTER_OPTS	该变量中设置的参数，将影响HBase的Master。
HBASE_REGIONSERVER_OPTS	该变量中设置的参数，将影响HBase的RegionServer。
HBASE_THRIFT_OPTS	该变量中设置的参数，将影响HBase的Thrift。

#### 配置方式举例：

```
export HADOOP_NAMENODE_OPTS="-Dhadoop.security.logger=${HADOOP_SECURITY_LOGGER:-INFO,RFAS} -Dhdfs.audit.logger=${HDFS_AUDIT_LOGGER:-INFO,NullAppender} $HADOOP_NAMENODE_OPTS"
```

## 8.8 HBase 运维管理

### 8.8.1 HBase 日志介绍

#### 日志描述

**日志存储路径：**HBase相关日志的默认存储路径为“/var/log/Bigdata/hbase/角色名”。

- HMaster：“/var/log/Bigdata/hbase/hm”（运行日志），“/var/log/Bigdata/audit/hbase/hm”（审计日志）。



- RegionServer: “/var/log/Bigdata/hbase/rs”（运行日志），“/var/log/Bigdata/audit/hbase/rs”（审计日志）。
- ThriftServer: “/var/log/Bigdata/hbase/ts2”（运行日志，ts2为具体实例名称），“/var/log/Bigdata/audit/hbase/ts2”（审计日志，ts2为具体实例名称）。

**日志归档规则：**HBase的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过30MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd\_hh-mm-ss>.[编号].log.zip”。最多保留最近的20个压缩文件，压缩文件保留个数可以在Manager界面中配置。

表 8-17 HBase 日志列表

日志类型	日志文件名	描述
运行日志	hbase-<SSH_USER>-<process_name>-<hostname>.log	HBase系统日志，主要包括启动时间，启动参数信息以及HBase系统运行时候所产生的大部分日志。
	hbase-<SSH_USER>-<process_name>-<hostname>.out	HBase运行环境信息日志。
	<process_name>-<SSH_USER>-<DATE>-<PID>-gc.log	HBase服务垃圾回收日志。
	checkServiceDetail.log	HBase服务启动是否成功的检查日志。
	hbase.log	HBase服务健康检查脚本以及部分告警检查脚本执行所产生的日志。
	sendAlarm.log	HBase告警检查脚本上报告警信息日志。
	hbase-haCheck.log	HMaster主备状态检测日志。
	stop.log	HBase服务进程启停操作日志。
审计日志	hbase-audit-<process_name>.log	HBase安全审计日志。

## 日志级别

HBase中提供了如表8-18所示的日志级别。日志级别优先级从高到低分别是FATAL、ERROR、WARN、INFO、DEBUG。程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 8-18 日志级别

级别	描述
FATAL	FATAL表示当前事件处理出现严重错误信息，可能导致系统崩溃。
ERROR	ERROR表示当前事件处理出现错误信息，系统运行出错。
WARN	WARN表示当前事件处理存在异常信息，但认为是正常范围，不会导致系统出错。
INFO	INFO表示记录系统及各事件正常运行状态信息
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

**步骤1** 进入HBase服务参数“全部配置”界面，具体操作请参考[修改集群服务配置参数](#)。

**步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。

 HBase（服务）

 RegionServer（角色）

Compaction

自定义

In-memory flush & compaction

日志

mapreduce

监控

**步骤3** 选择所需修改的日志级别。

**步骤4** 保存配置，在弹出窗口中单击“确定”使配置生效。

**说明**

配置完成后立即生效，不需要重启服务。

---结束

## 日志格式

HBase的日志格式如下所示：

**表 8-19** 日志格式

日志类型	组件	格式	示例
运行日志	HMaster	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置>	2020-01-19 16:04:53,558   INFO   main   env:HBASE_THRIFT_OPTS=   org.apache.hadoop.hbase.util.ServerCommandLine.logProcessInfo(ServerCommandLine.java:113)
	RegionServer	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置>	2020-01-19 16:05:18,589   INFO   regionserver16020-SendThread(linux-k6da:2181)   Client will use GSSAPI as SASL mechanism.   org.apache.zookeeper.client.ZooKeeperSaslClient\$1.run(ZooKeeperSaslClient.java:285)
	ThriftServer	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置>	2020-02-16 09:42:55,371   INFO   main   loaded properties from hadoop-metrics2.properties   org.apache.hadoop.metrics2.impl.MetricsConfig.loadFirst(MetricsConfig.java:111)
审计日志	HMaster	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置>	2020-02-16 09:42:40,934   INFO   master:linux-k6da:16000   Master: [master:linux-k6da:16000] start operation called.   org.apache.hadoop.hbase.master.HMaster.run(HMaster.java:581)

日志类型	组件	格式	示例
	RegionServer	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置>	2020-02-16 09:42:51,063   INFO   main   RegionServer: [regionserver16020] start operation called.   org.apache.hadoop.hbase.regionserver.HRegionServer.startRegionServer(HRegionServer.java:2396)
	ThriftServer	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置>	2020-02-16 09:42:55,512   INFO   main   thrift2 server start operation called.   org.apache.hadoop.hbase.thrift2.ThriftServer.main(ThriftServer.java:421)

## 8.8.2 HBase 常用参数配置

### 📖 说明

该章节操作仅适用于MRS 3.x之前版本集群。

当MRS服务中默认的参数配置不足以满足用户需要时，用户可以自定义修改参数配置来适应自身需求。

**步骤1** 登录集群详情页面，选择“组件管理”。

### 📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

**步骤2** 选择“HBase > 服务配置”，将“基础配置”切换为“全部配置”，进入HBase配置界面修改参数配置。

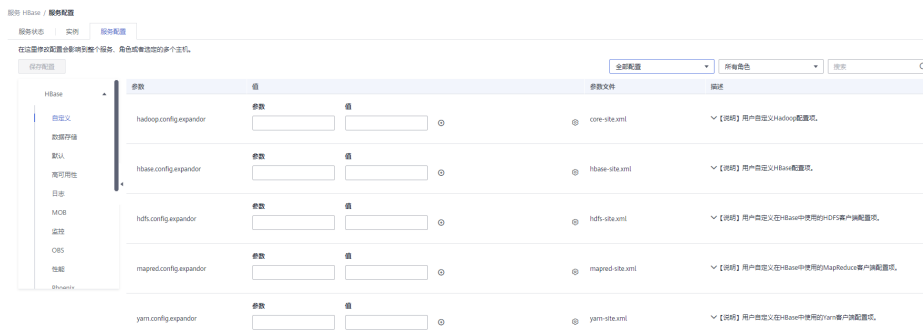


表 8-20 HBase 参数说明

参数	参数说明	参数值
hbase.regionserver.hfile.durable.sync	设置是否启用Hfile持久性以将数据持久化到磁盘。若将该参数设置为true，由于每个Hfile写入HBase时都会被hadoop fsync同步到磁盘上，则HBase性能将受到影响。 该参数仅在MRS 1.9.2及之前版本存在。	取值范围： <ul style="list-style-type: none"> <li>• true</li> <li>• false</li> </ul> 默认值为true
hbase.regionserver.wal.durable.sync	设置是否启用WAL文件持久性以将WAL数据持久化到磁盘。若将该参数设置为true，由于每个WAL的编辑都会被hadoop fsync同步到磁盘上，则HBase性能将受到影响。 该参数仅在MRS 1.9.2及之前版本存在。	取值范围： <ul style="list-style-type: none"> <li>• true</li> <li>• false</li> </ul> 默认值为true

----结束

### 8.8.3 配置 Region Transition 恢复线程

#### 操作场景

在故障环境中，由于诸如region服务器响应慢，网络不稳定，ZooKeeper节点版本不匹配等各种原因，有可能导致region长时间处于transition下。在region transition下，由于一些region不能对外提供服务，客户端操作可能无法正常执行。

#### 启用 Region Transition 恢复功能

在HMaster上设置chore服务，用于识别和恢复长期处于transition的region。

登录FusionInsight Manager界面，选择“集群 > 服务 > HBase > 配置”，下表是用于启用此功能的配置参数。

表 8-21 参数描述

参数	描述	默认值
hbase.region.assignment.auto.recovery.enabled	配置该参数以启用或禁用region分配恢复线程功能。	true

## 8.8.4 启用集群间拷贝功能备份集群数据

### 操作场景

当用户需要将保存在HDFS中的数据从当前集群备份到另外一个集群时，需要使用DistCp工具。DistCp工具依赖于集群间拷贝功能，该功能默认未启用。两个集群都需要配置。

该任务指导MRS集群管理员在MRS修改参数以启用集群间拷贝功能。

### 对系统的影响

启用集群间复制功能需要重启Yarn，服务重启期间无法访问。

### 前提条件

两个集群HDFS的参数“hadoop.rpc.protection”需使用相同的数据传输方式。设置为“privacy”表示加密，“authentication”表示不加密。

#### 说明

可登录FusionInsight Manager界面，选择“集群 > 服务 > HDFS > 配置”，搜索hadoop.rpc.protection查看。

针对MRS 3.x之前版本，在集群详情页选择“组件管理 > HDFS > 服务配置”，将“基础配置”切换为“全部配置”，搜索hadoop.rpc.protection查看。

### 操作步骤

**步骤1** 进入Yarn服务参数“全部配置”界面，具体操作请参考[修改集群服务配置参数](#)。

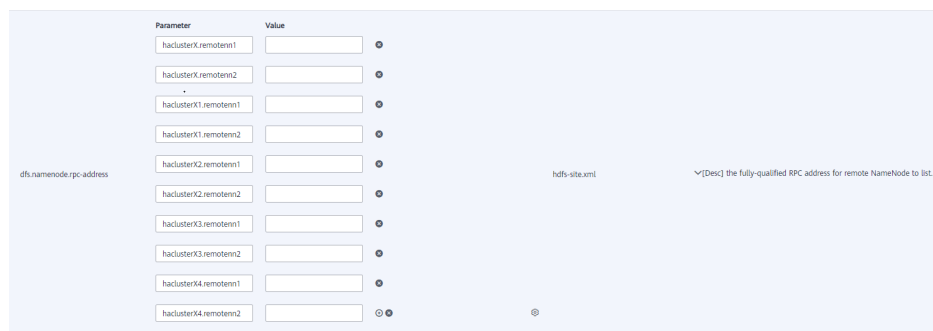
#### 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

**步骤2** 左边菜单栏中选择“Yarn > 集群间拷贝”。



**步骤3** 设置“dfs.namenode.rpc-address”参数的“haclusterX.remotenn1”值为对端集群其中一个NameNode实例的业务IP和RPC端口，设置“haclusterX.remotenn2”值为对端集群另外一个NameNode实例的业务IP和RPC端口。按照“IP:port”格式填写。



### 说明

针对MRS 3.x版本集群，登录FusionInsight Manager页面，选择“集群 > 待操作集群的名称 > 服务 > HDFS > 实例”，获取NameNode实例的业务IP。

针对MRS 3.x之前版本，在集群详情页选择“组件管理 > HDFS > 实例”，获取NameNode实例的业务IP。

“dfs.namenode.rpc-address.haclusterX.remotenn1”和“dfs.namenode.rpc-address.haclusterX.remotenn2”不区分主备NameNode。NameNode实例的业务IP可登录FusionInsight Manager页面，选择“集群 > 服务 > HDFS > 实例”获取；NameNode RPC端口可进入到HDFS服务配置页面搜索“dfs.namenode.rpc.port”参数获取，不支持通过Manager修改。

修改后参数值例如：“10.1.1.1:9820”和“10.1.1.2:9820”。

**步骤4** 保存配置并在概览页面选择“更多 > 重启服务”，重启Yarn服务。



界面提示“操作成功。”，单击“完成”，Yarn服务启动成功。

**步骤5** 登录另外一个集群，重复以上操作。

---结束

## 8.8.5 配置 HBase 主备集群数据自动备份

### 前提条件

1. 主备集群已经安装并且启动。
2. 主备集群上的时间必须一致，而且主备集群上的NTP服务必须使用同一个时间源。
3. 当主集群HBase服务关闭时，ZooKeeper和HDFS服务应该启动并运行。
4. 该工具应该由启动HBase进程的系统用户运行。
5. 如果处于安全模式，请确保备用集群的HBase系统用户具有主集群HDFS的读取权限。因为它将更新HBase系统ZooKeeper节点和HDFS文件。
6. 主集群HBase故障后，主集群的ZooKeeper，文件系统和网络依然可用。

### 场景介绍

Replication机制可以使用WAL将一个集群的状态与另一个集群的状态保持同步。启用HBase备份后，若主集群出现故障，ReplicationSyncUp工具会使用来自ZooKeeper的信息将主集群中的启用HBase备份功能的数据增量同步到备集群中。数据同步完成后，备集群可以作为主集群使用。



## 参数配置

参数	描述	默认值
hbase.replication.bulkload.enabled	是否开启批量加载数据复制功能。参数值类型为Boolean。开启批量加载数据复制功能后该参数须在主集群中设置为true。	false
hbase.replication.cluster.id	源HBase集群ID。开启批量加载数据复制功能时必须设置该参数，在源集群定义。参数值类型为String。	-

## 使用 ReplicationSyncUp 工具

在主集群hbase shell中输入如下命令使用：

```
hbase org.apache.hadoop.hbase.replication.regionserver.ReplicationSyncUp -Dreplication.sleep.before.failover=1
```

### 说明

replication.sleep.before.failover是指在RegionServer启动失败时备份其剩余数据前需要的休眠时间。由于30秒（默认值）的睡眠时间没有任何意义，因此将其设置为1（s），使备份过程更快触发。

## 注意事项

1. 当主集群关闭时，此工具将从ZooKeeper节点（RS znode）获得WAL的处理进度以及WAL的处理队列，并将未复制的队列复制到备集群中。
2. 每个主集群的RegionServer在备集群ZooKeeper上的replication节点下都有自己的znode。它包含每个对等集群的一个znode。
3. 当RegionServer故障时，主集群的每个RegionServer都会通过watcher收到通知，并尝试锁定故障RegionServer的znode，包含它的队列。成功创建的RegionServer会将所有队列转移到自己队列的znode下。队列传输后，它们将从旧位置删除。
4. 在主集群关闭期间，ReplicationSyncUp工具将使用来自ZooKeeper节点的信息同步主备集群的数据，并且RegionServer znode的wals将被移动到备集群下。

## 限制和约束

如果备集群处于关闭状态或关闭了对等关系，该工具正常运行，只有该对等关系复制不会发生。

## 8.8.6 HBase 集群容灾高可用

### 8.8.6.1 配置 HBase 主备集群容灾

#### 操作场景

HBase集群容灾作为提高HBase集群系统高可用性的一个关键特性，为HBase提供了实时的异地数据容灾功能。它对外提供了基础的运维工具，包含灾备关系维护，重建，

数据校验，数据同步进展查看等功能。为了实现数据的实时容灾，可以把本HBase集群中的数据备份到另一个集群。支持HBase表普通写数据与Bulkload批量写数据场景下的容灾。

#### 📖 说明

本章节适用于MRS 3.x及之后版本。

## 前提条件

- 主备集群都已经安装并启动成功，且获取集群的管理员权限。
- 必须保证主备集群间的网络畅通和端口的使用。
- 如果主集群部署为安全模式且不由一个FusionInsight Manager管理，主备集群必须已配置跨集群互信。如果主集群部署为普通模式，不需要配置跨集群互信。
- 主备集群必须已配置跨集群拷贝。
- 主备集群上的时间必须一致，而且主备集群上的NTP服务必须使用同一个时间源。
- 必须在主备集群的所有节点的hosts文件中，配置主备集群所有机器的机器名与业务IP地址的对应关系。

#### 📖 说明

若主集群的客户端安装在集群外的节点上，也需在该节点的hosts文件中配置主备集群所有机器的机器名与业务IP地址的对应关系。

- 主备集群间的网络带宽需要根据业务流量而定，不应少于最大的可能业务流量。
- 主备集群安装的MRS版本需要保持一致。
- 备集群规模不小于主集群规模。

## 使用约束

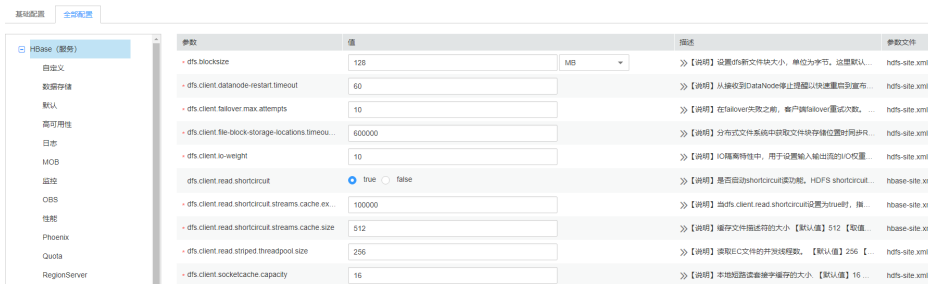
- 尽管容灾提供了实时的数据复制功能，但实际的数据同步进展，由多方面的因素决定的，例如，当前主集群业务的繁忙程度，备集群进程的健康状态等。因此，在正常情形下，备集群不应该接管业务。极端情形下是否可以接管业务，可由系统维护人员以及决策人员根据当前的数据同步指标来决定。
- 容灾功能当前仅支持一主一备。
- 通常情况下，不允许对备集群的灾备表进行表级别的操作，例如修改表属性、删除表等，一旦误操作备集群后会造成主集群数据同步失败、备集群对应表的数据丢失。
- 主集群的HBase表已启用容灾功能同步数据，用户每次修改表的结构时，需要手动修改备集群的灾备表结构，保持与主集群表结构一致。

## 操作步骤

### 配置主集群普通写数据容灾参数。

**步骤1** 登录主集群的Manager。

**步骤2** 选择“集群 > 待操作集群的名称 > 服务 > HBase > 配置”，单击“全部配置”，进入HBase配置界面。



**步骤3**（可选）如表8-22所示，为HBase容灾操作过程中的可选配置项，您可以根据描述来进行参数配置，或者使用缺省提供的值。

表 8-22 可选配置项

配置入口	配置项	缺省值	描述
“HMaster > 性能”	hbase.master.logcleaner.ttl	600000	指定HLog的保存期限。如果配置值为“604800000”（单位：毫秒），表示HLog的保存期限为7天。
	hbase.master.cleaner.interval	60000	HMaster清理过去HLog文件的周期，即超过设置的时间的HLog会被自动删除。建议尽可能配置大的值来保留更多的HLog。
“RegionServer > Replication”	replication.source.size.capacity	16777216	edits最大大小。单位为byte。如果edit大小超过这个值Hlog edits将会发送到备集群。
	replication.source.nb.capacity	25000	edits最大数目，这是另一个触发Hlog edits到备集群的条件。当主集群同步数据到备集群中时，主集群会从HLog中读取数据，此时会根据本参数配置的个数读取并发送。与“replication.source.size.capacity”一起配置使用。
	replication.source.maxretriesmultiplier	10	replication出现异常时的最大重试次数。
	replication.source.sleepforretries	1000	每次重试的sleep时间。（单位：毫秒）
	hbase.regionserver.replication.handler.count	6	RegionServer上的replication RPC服务器实例数。

**配置主集群Bulkload批量写数据容灾参数。**

**步骤4** 是否启用Bulkload批量写数据容灾功能？  
是，执行**步骤5**。

否，执行**步骤8**。

**步骤5** 选择“集群 > 待操作集群的名称 > 服务 > HBase > 配置”，单击“全部配置”，进入HBase配置界面。

**步骤6** 搜索并修改“hbase.replication.bulkload.enabled”参数，将配置项的值修改为“true”，启用Bulkload批量写数据容灾功能。

**步骤7** 搜索并修改“hbase.replication.cluster.id”参数，表示标识主集群HBase的ID，用于备集群连接主集群。参数值支持大小写字母、数字和下划线（\_），长度不超过30。

**重启HBase服务并安装客户端。**

**步骤8** 单击“保存”，保存配置。在弹出的窗口中单击“确定”。重启HBase服务。

**步骤9** 在主备集群，选择“集群 > 待操作集群的名称 > 服务 > HBase > 更多 > 下载客户端”，下载客户端并安装。

**添加主备集群容灾关系。**

**步骤10** 以“hbase”用户进入主集群的HBase shell界面。

hbase用户的初始密码为“Hbase@123”，详情请参考[用户账号一览表](#)。

**步骤11** 在HBase shell中执行如下命令，创建主集群HBase与备集群HBase之间的容灾同步关系。

```
add_peer '备集群ID', CLUSTER_KEY => "备集群ZooKeeper业务ip地址", CONFIG => {"hbase.regionserver.kerberos.principal" => "备集群RegionServer principal", "hbase.master.kerberos.principal" => "备集群HMaster principal"}
```

- 备集群ID表示主集群识别备集群使用的id，请重新指定id值。可以任意指定，建议使用数字。
- 备集群ZooKeeper地址信息包含ZooKeeper业务IP地址、侦听客户端连接的端口和备集群的HBase在ZooKeeper上的根目录。
- hbase.master.kerberos.principal、hbase.regionserver.kerberos.principal在备集群HBase hbase-site.xml配置文件中查找。

例如，添加主备集群容灾关系，执行：`add_peer '备集群ID', CLUSTER_KEY => "192.168.40.2,192.168.40.3,192.168.40.4:24002:/hbase", CONFIG => {"hbase.regionserver.kerberos.principal" => "hbase/hadoop.hadoop.com@HADOOP.COM", "hbase.master.kerberos.principal" => "hbase/hadoop.hadoop.com@HADOOP.COM"}`

**步骤12**（可选）如果启用Bulkload批量写数据容灾功能，主集群HBase客户端配置必须拷贝到备集群。

- 在备集群HDFS创建目录/hbase/replicationConf/*主集群*  
*hbase.replication.cluster.id*
- 主集群HBase客户端配置文件，拷贝到备集群HDFS目录/hbase/replicationConf/*主集群**hbase.replication.cluster.id*

例如：`hdfs dfs -put HBase/hbase/conf/core-site.xml HBase/hbase/conf/hdfs-site.xml HBase/hbase/conf/yarn-site.xml hdfs://NameNode IP.25000/hbase/replicationConf/source_cluster`

**启用HBase容灾功能同步数据。**

**步骤13** 检查备集群的HBase服务实例中，是否已存在一个命名空间，与待启用容灾功能的HBase表所属的命名空间名称相同？

- 是，存在同名的命名空间，执行**步骤14**。
- 否，不存在同名的命名空间，需先在备集群的HBase shell中，创建同名的命名空间，然后执行**步骤14**。

**步骤14** 在主集群的HBase shell中，以“hbase”用户执行以下命令，启用将主集群表的数据实时容灾功能，确保后续主集群中修改的数据能够实时同步到备集群中。

一次只能针对一个HTable进行数据同步。

**enable\_table\_replication '表名'**

#### 📖 说明

- 若备集群中不存在与要开启实时同步的表同名的表，则该表会自动创建。
- 若备集群中存在与要开启实时同步的表同名的表，则两个表的结构必须一致。
- 若'表名'设置了加密算法SMS4或AES，则不支持对此HBase表启用将数据从主集群实时同步到备集群的功能。
- 若备集群不在线，或备集群中已存在同名但结构不同的表，启用容灾功能将失败。
- 若主集群中部分Phoenix表启用容灾功能同步数据，则备集群中不能存在与主集群Phoenix表同名的普通HBase表，否则启用容灾功能失败或影响备集群的同名表正常使用。
- 若主集群中Phoenix表启用容灾功能同步数据，还需要对Phoenix表的元数据表启用容灾功能同步数据。需配置的元数据表包含SYSTEM.CATALOG、SYSTEM.FUNCTION、SYSTEM.SEQUENCE和SYSTEM.STATS。
- 若主集群的HBase表启用容灾功能同步数据，用户每次为HBase表增加新的索引，需要手动在备集群的灾备表增加二级索引，保持与主集群二级索引结构一致。

**步骤15**（可选）如果HBase没有使用Ranger，在主集群的HBase shell中，以“hbase”用户执行以下命令，启用主集群的HBase表权限控制信息数据实时容灾功能。

**enable\_table\_replication 'hbase:acl'**

#### 创建用户

**步骤16** 登录备集群的FusionInsight Manager，选择“系统 > 权限 > 角色 > 添加角色”创建一个角色，并根据主集群HBase源数据表的权限，为角色添加备数据表的相同权限。

**步骤17** 选择“系统 > 权限 > 用户 > 添加用户”创建一个用户，根据业务需要选择用户类型为“人机”或“机机”，并将用户加入创建的角色。使用新创建的用户，访问备集群的HBase容灾数据。

#### 📖 说明

- 主集群HBase源数据表修改权限时，如果备集群需要正常读取数据，请修改备集群角色的权限。
- 如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加HBase的Ranger访问权限策略](#)。

#### 同步主集群表数据。

**步骤18** 检查配置HBase容灾并启用数据同步后，主集群是否已存在表及数据，且历史数据需要同步到备集群？

- 是，存在表且需要同步数据，以HBase表用户登录安装主集群HBase客户端的节点，并执行kinit用户名认证身份。该用户需要拥有表的读写权限，以及“hbase:meta”表的执行权限。然后执行**步骤19**。

- 否，不需要同步数据，任务结束。

**步骤19** 配置HBase容灾时不支持自动同步表中的历史数据，需要对主集群的历史数据进行备份，然后再手动恢复历史数据到备集群中。

手动恢复即单表的恢复，单表手动恢复通过Export、distcp、Import来完成。

单表手动恢复操作步骤：

1. 从主集群导出表中数据。

**hbase org.apache.hadoop.hbase.mapreduce.Export -**  
**Dhbase.mapreduce.include.deleted.rows=true** 表名 保存源数据的目录

例如，**hbase org.apache.hadoop.hbase.mapreduce.Export -**  
**Dhbase.mapreduce.include.deleted.rows=true t1 /user/hbase/t1**

2. 把导出的数据复制到备集群。

**hadoop distcp** 主集群保存源数据的目录 *hdfs://ActiveNameNodeIP:8020/备集群保存源数据的目录*

其中，ActiveNameNodeIP是备集群中主NameNode节点的IP地址。

例如，**hadoop distcp /user/hbase/t1 hdfs://192.168.40.2:8020/user/hbase/t1**

3. 使用备集群HBase表用户，在备集群中导入数据。

在备集群HBase shell界面，使用“hbase”用户执行以下命令保持写数据状态：

**set\_clusterState\_active**

界面提示以下信息表示执行成功：

```
hbase(main):001:0> set_clusterState_active
=> true
```

**hbase org.apache.hadoop.hbase.mapreduce.Import -Dimport.bulk.output=备集群保存输出的目录 表名 备集群保存源数据的目录**

**hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles** 备集群保存输出的目录 表名

例如：

```
hbase(main):001:0> set_clusterState_active
=> true
```

**hbase org.apache.hadoop.hbase.mapreduce.Import -**  
**Dimport.bulk.output=/user/hbase/output\_t1 t1 /user/hbase/t1**

**hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles /user/hbase/output\_t1 t1**

**步骤20** 在HBase客户端执行以下命令，校验主备集群同步的数据。启用容灾功能同步功能后，也可以执行该命令检验新的同步数据是否一致。

**hbase org.apache.hadoop.hbase.mapreduce.replication.VerifyReplication --**  
**starttime=开始时间--endtime=结束时间列族名称 备集群ID 表名**

#### 说明

- 开始时间必须早于结束时间
- 开始时间和结束时间需要填写时间戳的格式，例如执行date -d "2015-09-30 00:00:00" +%s将普通时间转化为时间戳格式。

**指定主备集群写数据状态。**

**步骤21** 在主集群HBase shell界面，使用“hbase”用户执行以下命令保持写数据状态。

**set\_clusterState\_active**

界面提示以下信息表示执行成功：

```
hbase(main):001:0> set_clusterState_active
=> true
```

**步骤22** 在备集群HBase shell界面，使用“hbase”用户执行以下命令保持只读数据状态。

**set\_clusterState\_standby**

界面提示以下信息表示执行成功：

```
hbase(main):001:0> set_clusterState_standby
=> true
```

---结束

## 相关命令

表 8-23 HBase 容灾

操作	命令	描述
建立灾备关系	<pre>add_peer '备集群ID', CLUSTER_KEY =&gt; "备集群 ZooKeeper业务ip地址", CONFIG =&gt; {"hbase.regionserver.kerberos.principal" =&gt; "备集群RegionServer principal", "hbase.master.kerberos.principal" =&gt; "备集群HMaster principal"} add_peer '1','zk1,zk2,zk3:2181:/hbase1' 2181表示集群中ZooKeeper的端口号。</pre>	<p>建立主集群与备集群的关系，让其互相对应。</p> <p>如果启用Bulkload批量写数据容灾：</p> <ul style="list-style-type: none"> <li>在备集群HDFS创建目录/hbase/replicationConf/<i>主集群hbase.replication.cluster.id</i></li> <li>主集群HBase客户端配置文件，拷贝到备集群HDFS目录/hbase/replicationConf/<i>主集群hbase.replication.cluster.id</i></li> </ul>
移除灾备关系	<pre>remove_peer '备集群ID' 示例： remove_peer '1'</pre>	在主集群中移除备集群的信息。
查询灾备关系	<pre>list_peers</pre>	在主集群中查询已经设置的备集群的信息，主要为Zookeeper信息。
启用用户表实时同步	<pre>enable_table_replication '表名' 示例： enable_table_replication 't1'</pre>	在主集群中，设置已存在的表同步到备集群。
禁用用户表实时同步	<pre>disable_table_replication '表名' 示例： disable_table_replication 't1'</pre>	在主集群中，设置已存在的表不同步到备集群。

操作	命令	描述
主备集群数据校验	<code>bin/hbase org.apache.hadoop.hbase.mapreduce.replication.VerifyReplication --starttime=<i>开始时间</i>--endtime=<i>结束时间</i> <i>列族名称</i> <i>备集群ID</i> <i>表名</i></code>	<p>检查指定的表在主备集群间的数据是否一致。</p> <p>命令行中参数说明如下：</p> <ul style="list-style-type: none"> <li>开始时间：如果未设置，则取默认的开始时间为0。</li> <li>结束时间：如果未设置，则取默认的结束时间为当前操作提交的时间。</li> <li>表名：如果未输入表名，则默认校验所有的启用了实时同步的用户表。</li> </ul>
切换数据写入状态	<code>set_clusterState_active</code> <code>set_clusterState_standby</code>	设置集群HBase表是否可写入数据。
新增或更新已经在对端集群保存的主集群中HDFS配置	<code>hdfs dfs -put -f HBase/hbase/conf/core-site.xml HBase/hbase/conf/hdfs-site.xml HBase/hbase/conf/yarn-site.xml hdfs://<i>备集群NameNode IP:PORT</i>/hbase/replicationConf/<i>主集群hbase.replication.cluster.id</i></code>	<p>启用包含Bulkload数据的容灾，在主集群修改HDFS参数时，新的参数值默认不会从主集群自动同步到备集群，需要手动执行命令同步。受影响的参数如下：</p> <ul style="list-style-type: none"> <li>“fs.defaultFS”</li> <li>“dfs.client.failover.proxy.provider.hacluster”</li> <li>“dfs.client.failover.connection.retries.on.timeouts”</li> <li>“dfs.client.failover.connection.retries”</li> </ul> <p>例如，“fs.defaultFS”修改为“hdfs://hacluster_sale”，主集群HBase客户端配置文件，重新拷贝到备集群HDFS目录/hbase/replicationConf/<i>主集群hbase.replication.cluster.id</i></p>

### 8.8.6.2 HBase 容灾集群主备倒换

#### 操作场景

当前环境HBase已经是容灾集群，因为某些原因，需要将主备集群互换，即备集群变成主集群，主集群变成备集群。

#### 📖 说明

本章节适用于MRS 3.x及之后版本。



## 对系统的影响

主备集群互换后，原先主集群将不能再写入数据，原先备集群将变成主集群，接管上层业务。

## 操作步骤

### 确保上层业务已经停止

**步骤1** 确保上层业务已经停止，如果没有停止，先执行 参考[HBase容灾集群业务切换指导](#)。

### 关闭主集群写功能

**步骤2** 下载并安装HBase客户端。

具体请参考[安装客户端（3.x及之后版本）](#)章节。

**步骤3** 在备集群HBase客户端，以**hbase**用户执行以下命令指定备集群写数据状态关闭。

```
kinit hbase
```

```
hbase shell
```

```
set_clusterState_standby
```

界面提示以下信息表示执行成功：

```
hbase(main):001:0> set_clusterState_standby
=> true
```

### 检查当前主备同步是否完成

**步骤4** 执行以下命令，确保当前数据已经同步，要求SizeOfLogQueue=0，SizeOfLogToReplicate=0，如果不为零，等待，重复执行以下命令，直到等于0。

```
status 'replication'
```

### 关闭主备集群同步

**步骤5** 查询所有的同步集群，获取PEER\_ID。

```
list_peers
```

**步骤6** 删除所有同步集群。

```
remove_peer '备集群ID'
```

示例：

```
remove_peer '1'
```

**步骤7** 查询所有同步的table。

```
list_replicated_tables
```

**步骤8** 分别disable上面查询到的所有同步的table。

```
disable_table_replication '表名'
```

示例：

```
disable_table_replication 't1'
```

### 切换主备

**步骤9** 重新配置HBase容灾，参考[配置HBase主备集群容灾](#)。

----结束

## 8.8.6.3 HBase 容灾集群业务切换指导

### 操作场景

MRS集群管理员可配置HBase集群容灾功能，以提高系统可用性。容灾环境中的主集群完全故障影响HBase上层应用连接时，需要为HBase上层应用配置备集群信息，才可以使得该应用在备集群上运行。

#### 说明

本章节适用于MRS 3.x及之后版本。

### 对系统的影响

切换业务后，写入备集群的数据默认不会同步到主集群。主集群故障修复后，备集群新增的数据需要通过备份恢复的方式同步到主集群。如果需要自动同步数据，需要切换HBase容灾主备集群。

### 操作步骤

**步骤1** 登录备集群FusionInsight Manager。

**步骤2** 下载并安装HBase客户端。

**步骤3** 在备集群HBase客户端，以**hbase**用户执行以下命令指定备集群写数据状态启用。

```
kinit hbase
```

```
hbase shell
```

```
set_clusterState_active
```

界面提示以下信息表示执行成功：

```
hbase(main):001:0> set_clusterState_active
=> true
```

**步骤4** 确认HBase上层应用中原有的配置文件“hbase-site.xml”、“core-site.xml”和“hdfs-site.xml”是否为适配应用运行修改或新增过配置内容。

- 是，将相关内容同步更新到新的配置文件中，并替换旧的配置文件。
- 否，使用新的配置文件替换HBase上层应用中原有的配置文件。

**步骤5** 配置HBase上层应用所在主机与备集群的网络连接。

#### 说明

当客户端所在主机不是集群中的节点时，配置客户端网络连接，可避免执行客户端命令时出现错误。

1. 确保客户端所在主机能与客户端安装包文件解压目录下的“hosts”文件中所列出的集群各主机在网络上互通。

2. 当客户端所在主机不是集群中的节点时，需要在客户端所在节点的“/etc/hosts”文件中设置主机名和IP地址（业务平面）映射。主机名和IP地址请保持一一对应。

**步骤6** 配置HBase上层应用所在主机的时间与备集群的时间保持一致，时间差要小于5分钟。

**步骤7** 检查主集群的认证模式。

- 若为安全模式，执行**步骤8**。
- 若为普通模式，任务结束。

**步骤8** 获取HBase上层应用用户的keytab文件和krb5.conf配置文件。

1. 在备集群FusionInsight Manager界面，选择“系统 > 权限 > 用户”。
2. 在用户所在行的“操作”列单击“更多 > 下载认证凭据”，下载keytab文件到本地。
3. 解压得到“user.keytab”和“krb5.conf”。

**步骤9** 使用“user.keytab”和“krb5.conf”两个文件替换HBase上层应用中原有的文件。

**步骤10** 停止上层业务。

**步骤11** 是否需要切换HBase主备集群，即主变成备，备变成主。如果不切换，数据将不再同步。

- 是，先执行HBase容灾主备集群倒换，具体请参考[HBase容灾集群主备倒换](#)，然后再执行**步骤12**。
- 否，直接执行**步骤12**。

**步骤12** 启动上层业务。

----结束

## 8.9 HBase 常见问题

### 8.9.1 结束 BulkLoad 客户端程序导致作业执行失败

#### 问题

执行BulkLoad程序导入数据时，如果结束客户端程序，为什么有时会导致已提交的作业执行失败？

#### 回答

BulkLoad程序在客户端启动时会生成一个partitioner文件，用于划分Map任务数据输入的范围。

此文件在BulkLoad客户端退出时会被自动删除。

一般来说当所有Map任务都启动运行以后，退出BulkLoad客户端也不会导致已提交的作业失败。但由于Map任务存在重试机制和推测执行机制；Reduce任务下载一个已运行完成的Map任务的数据失败次数过多时，Map任务也会被重新执行。

如果此时BulkLoad客户端已经退出，则重试的Map任务会因为找不到partitioner文件而执行失败，导致作业执行失败。

因此，强烈建议BulkLoad程序在数据导入期间不要结束客户端程序。

## 8.9.2 如何修复长时间处于 RIT 状态的 Region

### 问题

在HBase WebUI界面看到有长时间处于RIT状态的Region，如何修复？

### 回答

登录HMaster WebUI，在导航栏选择“Procedure & Locks”，查看是否有处于Waiting状态的process id。如果有，需要执行以下命令将procedure lock释放：

```
hbase hbck -j 客户端安装目录/HBase/hbase/tools/hbase-hbck2-*.jar bypass -o pid
```

查看State是否处于Bypass状态，如果界面上的procedures一直处于RUNNABLE(Bypass)状态，需要进行主备切换。执行**assigns**命令使region重新上线。

```
hbase hbck -j 客户端安装目录/HBase/hbase/tools/hbase-hbck2-*.jar assigns -o regionName
```

## 8.9.3 HMaster 等待 NameSpace 表上线时超时退出

### 问题

为什么在等待namespace表上线时超时HMaster退出？

### 回答

在HMaster主备倒换或启动期间，HMaster为先前失败/停用的RegionServer执行WAL splitting及region恢复。

在后台运行有多个监控HMaster启动进程的线程：

- **TableNamespaceManager**  
这是一个帮助类，用于在HMaster主备倒换或启动期间，管理namespace表及监控表region的分配。如果namespace表在规定时间内（`hbase.master.namespace.init.timeout`，默认为3600000ms）内没有上线，那么它就会异常中断HMaster进程。
- **InitializationMonitor**  
这是一个主HMaster初始化线程监控类，用于监控主Master的初始化。如果在规定时间内（`hbase.master.initializationmonitor.timeout`，默认为3600000ms）内初始化线程失败，该线程会异常终止HMaster（如果该`hbase.master.initializationmonitor.haltontimeout`被启动，默认为false）。

在HMaster主备倒换或启动期间，如果WAL hlog文件存在，它会初始化WAL splitting任务。如果WAL hlog splitting任务完成，它将初始化表region分配任务。

HMaster通过ZooKeeper协调log splitting任务和有效的RegionServer，并追踪任务的发展。如果主HMaster在log splitting任务期间退出，新的主HMaster会尝试重发没有完成的任务，RegionServer从头启动log splitting任务。

HMaster初始化工作情况会由于很多原因被延迟：

- 间歇性的网络故障。

- 磁盘瓶颈。
- log split任务工作负荷较大，RegionServer运行缓慢。
- RegionServer（region opening）响应缓慢。

在以上场景中，为使HMaster更早完成恢复任务，建议增加以下配置参数，否则Master将退出导致整个恢复进程被更大程度地延迟。

- 增加namespace表在线等待超时周期，保证Master有足够的时间协调RegionServer workers split任务，避免一次次重复相同的任务。  
“hbase.master.namespace.init.timeout”（默认为3600000ms）
- 通过RegionServer worker增加并行split任务执行数，保证RegionServer worker能并行处理split work（RegionServer需要有更多的核心）。在“客户端安装路径/HBase/hbase/conf/hbase-site.xml”中添加参数：  
“hbase.regionserver.wal.max.splitters”（默认为2）
- 如果所有的恢复过程都需要时间，增加初始化监控线程超时时间。  
“hbase.master.initializationmonitor.timeout”（默认为3600000ms）

## 8.9.4 客户端查询 HBase 出现 SocketTimeoutException 异常

### 问题

使用HBase客户端操作表数据的时候客户端出现类似如下异常：

```
2015-12-15 02:41:14,054 | WARN | [task-result-getter-2] | Lost task 2.0 in stage 58.0 (TID 3288, linux-175):
org.apache.hadoop.hbase.client.RetriesExhaustedException: Failed after attempts=36, exceptions:
Tue Dec 15 02:41:14 CST 2015, null, java.net.SocketTimeoutException: callTimeout=60000,
callDuration=60303:
row 'xxxxxx' on table 'xxxxxx',\x05\x1E
\x80\x00\x00\x00\x80\x00\x00\x00\x00\x00\x00\x00\x80\x00\x00\x00\x00\x00\x00\x00\x80\x00\x00\x00
0\x80\x00\x00\x00\x80\x00\x00,
1449912620868.6a6b7d0c272803d8186930a3bfdb10a9., hostname=xxxxxx,16020,1449941841479,
seqNum=5
at
org.apache.hadoop.hbase.client.RpcRetryingCallerWithReadReplicas.throwEnrichedException(RpcRetryingCall
erWithReadReplicas.java:275)
at org.apache.hadoop.hbase.client.ScannerCallableWithReplicas.call(ScannerCallableWithReplicas.java:223)
at org.apache.hadoop.hbase.client.ScannerCallableWithReplicas.call(ScannerCallableWithReplicas.java:61)
at org.apache.hadoop.hbase.client.RpcRetryingCaller.callWithoutRetries(RpcRetryingCaller.java:200)
at org.apache.hadoop.hbase.client.ClientScanner.call(ClientScanner.java:323)
```

同时，在RegionServer上出现类似如下日志：

```
2015-12-15 02:45:44,551 | WARN | PriorityRpcServer.handler=7,queue=1,port=16020 | (responseTooSlow):
{"call": "Scan(org.apache.hadoop.hbase.protobuf.generated.ClientProtos$ScanRequest)
","starttimems":1450118730780,"responsesize":416,"method": "Scan","processingtimems":13770,"client": "10.9
1.8.175:41182","queuetimems":0,"class": "HRegionServer"} |
org.apache.hadoop.hbase.ipc.RpcServer.logResponse(RpcServer.java:2221)
2015-12-15 02:45:57,722 | WARN | PriorityRpcServer.handler=3,queue=1,port=16020 | (responseTooSlow):
{"call": "Scan(org.apache.hadoop.hbase.protobuf.generated.ClientProtos
$ScanRequest)","starttimems":1450118746297,"responsesize":416,
"method": "Scan","processingtimems":11425,"client": "10.91.8.175:41182","queuetimems":1746,"class": "HRegi
onServer"} | org.apache.hadoop.hbase.ipc.RpcServer.logResponse(RpcServer.java:2221)
2015-12-15 02:47:21,668 | INFO | LruBlockCacheStatsExecutor | totalSize=7.54 GB, freeSize=369.52 MB,
max=7.90 GB, blockCount=406107,
accesses=35400006, hits=16803205, hitRatio=47.47%, , cachingAccesses=31864266, cachingHits=14806045,
cachingHitsRatio=46.47%,
evictions=17654, evicted=16642283, evictedPerRun=942.69189453125 |
org.apache.hadoop.hbase.io.hfile.LruBlockCache.logStats(LruBlockCache.java:858)
2015-12-15 02:52:21,668 | INFO | LruBlockCacheStatsExecutor | totalSize=7.51 GB, freeSize=395.34 MB,
max=7.90 GB, blockCount=403080,
```

```
accesses=35685793, hits=16933684, hitRatio=47.45%, , cachingAccesses=32150053, cachingHits=14936524,
cachingHitsRatio=46.46%,
evictions=17684, evicted=16800617, evictedPerRun=950.046142578125 |
org.apache.hadoop.hbase.io.hfile.LruBlockCache.logStats(LruBlockCache.java:858)
```

## 回答

出现该问题的主要原因为RegionServer分配的内存过小、Region数量过大导致在运行过程中内存不足，服务端对客户端的响应过慢。在RegionServer的配置文件“hbase-site.xml”中需要调整如下对应的内存分配参数。

表 8-24 RegionServer 内存调整参数

参数	描述	默认值
GC_OPTS	在启动参数中给RegionServer分配的初始内存和最大内存。	-Xms8G -Xmx8G
hfile.block.cache.size	分配给HFile/StoreFile所使用的块缓存的最大 heap ( -Xmx setting ) 的百分比。	当offheap关闭时，默认值为0.25。当offheap开启时，默认值是0.1。

## 8.9.5 在启动 HBase shell 时，报错 “java.lang.UnsatisfiedLinkError: Permission denied”

### 问题

在启动HBase shell时，为什么会发生“java.lang.UnsatisfiedLinkError: Permission denied”异常？

### 回答

在执行HBase shell期间，JRuby会在“java.io.tmpdir”路径下创建一个临时文件，该路径的默认值为“/tmp”。如果为“/tmp”目录设置NOEXEC权限，然后HBase shell会启动失败并发生“java.lang.UnsatisfiedLinkError: Permission denied”异常。

因此，如果为“/tmp”目录设置了NOEXEC权限，那么“java.io.tmpdir”必须设置为HBASE\_OPTS/CLIENT\_GC\_OPTS中不同的路径。

## 8.9.6 停止运行的 RegionServer，在 HMaster WebUI 中显示的 “Dead Region Servers” 信息什么时候会被清除掉

### 问题

在HMaster Web UI中显示处于“Dead Region Servers”状态的RegionServer什么时候会被清除掉？

## 回答

当一个在线的RegionServer突然运行停止，会在HMaster Web UI中显示处于“Dead Region Servers”状态。当停止运行的RegionServer重启并且向HMaster上报成功信息，在HMaster Web UI中会清除掉“Dead Region Servers”信息。

当HMaster主备倒换操作成功执行时，在HMaster Web UI中也会清除掉“Dead Region Servers”信息。

以防掌控有一些region的主用HMaster突然停止响应，备用的HMaster将会成为新的主用HMaster，同时显示先前主用HMaster变成dead RegionServer。当HMaster主备倒换操作成功执行，在HMaster Web UI中也会清除掉“Dead Region Servers”。

## 8.9.7 访问 HBase Phoenix 提示权限不足如何处理

### 问题

使用租户访问Phoenix提示权限不足。

### 回答

创建租户的时候需要关联HBase服务和Yarn队列。

租户要操作Phoenix还需要额外操作的权限，即Phoenix系统表的RWX权限。

例如：

创建好的租户为**hbase**，使用**admin**用户登录hbase shell，执行**scan 'hbase:acl'**命令查询租户对应的角色为**hbase\_1450761169920**（格式为：租户名\_时间戳）。

执行以下命令进行授权（如果还没有生成Phoenix系统表，请用**admin**用户登录Phoenix客户端后再回到hbase shell里授权）：

```
grant '@hbase_1450761169920','RWX','SYSTEM.CATALOG'
```

```
grant '@hbase_1450761169920','RWX','SYSTEM.FUNCTION'
```

```
grant '@hbase_1450761169920','RWX','SYSTEM.SEQUENCE'
```

```
grant '@hbase_1450761169920','RWX','SYSTEM.STATS'
```

新建用户**phoenix**并绑定租户**hbase**，该用户**phoenix**就可以用来访问Phoenix客户端。

## 8.9.8 租户使用 HBase BulkLoad 功能提示权限不足如何处理

### 问题

租户使用HBase bulkload功能提示权限不足。

### 回答

创建租户的时候需要关联HBase服务和Yarn队列。

例如：

新建用户**user**并绑定租户同名的角色。

用户user需要使用bulkload功能还需要额外权限。

以下以用户user为例：

参见“批量导入数据”章节举例，以下是一些差异点。

1. 将数据文件目录建在“/tmp”目录下，执行以下命令：

```
hdfs dfs -mkdir /tmp/datadirImport
```

```
hdfs dfs -put data.txt /tmp/datadirImport
```

2. 生成HFile的时候使用HDFS的“/tmp”目录：

```
hbase com.huawei.hadoop.hbase.tools.bulkload.ImportData -
Dimport.skip.bad.lines=true -Dimport.separator=';' -
Dimport.bad.lines.output=/tmp/badline -Dimport.hfile.output=/tmp/hfile
configuration.xml ImportTable /tmp/datadirImport
```

3. 导入HFile的时候使用HDFS的“/tmp”目录：

```
hbase org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles /tmp/
hfile ImportTable
```

## 8.9.9 如何修复 Overlap 状态的 HBase Region

### 问题

MRS 3.x及之后版本，使用hbck工具检查Region状态，若日志中存在“ERROR: (regions region1 and region2) There is an overlap in the region chain.”或者“ERROR: (region region1) Multiple regions have the same startkey: xxx”信息，表示某些region存在overlap的问题，需要如何解决？

### 回答

修复步骤如下：

- 步骤1** 执行**hbase hbck -j \${CLIENT\_HOME}/HBase/hbase/tools/hbase-hbck2-1.1.0-h0.cbu.mrs.\*.jar fixInconsistencies *tableName***命令修复存在overlap的表。
  - 步骤2** 执行**hbase hbck -j \${CLIENT\_HOME}/HBase/hbase/tools/hbase-hbck2-1.1.0-h0.cbu.mrs.\*.jar listInconsistencies -run *tableName***命令检查修复的表是否还存在overlap。
    - 如果不存在overlap，执行**步骤3**。
    - 如果存在overlap，执行**步骤1**。
  - 步骤3** 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > HBase > 更多 > 执行HMaster倒换”，完成HMaster主备倒换。
  - 步骤4** 执行**hbase hbck -j \${CLIENT\_HOME}/HBase/hbase/tools/hbase-hbck2-1.1.0-h0.cbu.mrs.\*.jar listInconsistencies -run *tableName***命令检查修复的表是否还存在overlap。
    - 如果不存在overlap，修复完成。
    - 如果存在overlap，从**步骤1**开始重新执行修复步骤。
- 结束



## 8.9.10 Phoenix BulkLoad Tool 使用限制说明

### 问题

当更新索引字段数据时，若用户表已经存在一批数据，则BulkLoad工具不能更新全局和局部可变索引。

### 回答

#### 问题分析

1. 创建表。

```
CREATE TABLE TEST_TABLE(
DATE varchar not null,
NUM integer not null,
SEQ_NUM integer not null,
ACCOUNT1 varchar not null,
ACCOUNTDES varchar,
FLAG varchar,
SALL double,
CONSTRAINT PK PRIMARY KEY (DATE,NUM,SEQ_NUM,ACCOUNT1)
);
```

2. 创建全局索引

```
CREATE INDEX TEST_TABLE_INDEX ON
TEST_TABLE(ACCOUNT1,DATE,NUM,ACCOUNTDES,SEQ_NUM);
```

3. 插入数据

```
UPSERT INTO TEST_TABLE
(DATE,NUM,SEQ_NUM,ACCOUNT1,ACCOUNTDES,FLAG,SALL) values
('20201001',30201001,13,'367392332','sffa1','');
```

4. 执行BulkLoad任务更新数据

```
hbase org.apache.phoenix.mapreduce.CsvBulkLoadTool -t TEST_TABLE -
i /tmp/test.csv, test.csv内容如下:
```

```
20201001 30201001 13 367392332 sffa888 1231243 23
```

5. 问题现象：无法直接更新之前存在的索引数据，导致存在两条索引数据。

```
+-----+-----+-----+-----+-----+
|:ACCOUNT1 | :DATE | :NUM | 0:ACCOUNTDES |:SEQ_NUM |
+-----+-----+-----+-----+-----+
| 367392332 | 20201001 | 30201001 | sffa1 | 13 |
| 367392332 | 20201001 | 30201001 | sffa888 | 13 |
+-----+-----+-----+-----+-----+
```

#### 解决方法

步骤1 删除旧的索引表。

```
DROP INDEX TEST_TABLE_INDEX ON TEST_TABLE;
```

步骤2 异步方式创建新的索引表。

```
CREATE INDEX TEST_TABLE_INDEX ON
TEST_TABLE(ACCOUNT1,DATE,NUM,ACCOUNTDES,SEQ_NUM) ASYNC;
```

步骤3 索引重建。

```
hbase org.apache.phoenix.mapreduce.index.IndexTool --data-table
TEST_TABLE --index-table TEST_TABLE_INDEX --output-path /user/test_table
```

----结束

## 8.9.11 CTBase 对接 Ranger 权限插件提示权限不足

### 问题

在MRS集群中，CTBase访问启用Ranger插件的HBase服务时，如果创建聚簇表，提示权限不足。

报错信息如下：

```
ERROR: Create ClusterTable failed. Error: org.apache.hadoop.hbase.security.AccessDeniedException:
Insufficient permissions for user 'ctbase2@HADOOP.COM' (action=create)
at org.apache.ranger.authorization.hbase.AuthorizationSession.publishResults(AuthorizationSession.java:278)
at
org.apache.ranger.authorization.hbase.RangerAuthorizationCoproprocessor.authorizeAccess(RangerAuthorizatio
nCoproprocessor.java:654)
at
org.apache.ranger.authorization.hbase.RangerAuthorizationCoproprocessor.requirePermission(RangerAuthorizati
onCoproprocessor.java:772)
at
org.apache.ranger.authorization.hbase.RangerAuthorizationCoproprocessor.preCreateTable(RangerAuthorization
Coproprocessor.java:943)
at
org.apache.ranger.authorization.hbase.RangerAuthorizationCoproprocessor.preCreateTable(RangerAuthorization
Coproprocessor.java:428)
at org.apache.hadoop.hbase.master.MasterCoproprocessorHost$12.call(MasterCoproprocessorHost.java:351)
at org.apache.hadoop.hbase.master.MasterCoproprocessorHost$12.call(MasterCoproprocessorHost.java:348)
at org.apache.hadoop.hbase.coprocessor.CoproprocessorHost
$ObserverOperationWithoutResult.callObserver(CoproprocessorHost.java:581)
at org.apache.hadoop.hbase.coprocessor.CoproprocessorHost.execOperation(CoproprocessorHost.java:655)
at
org.apache.hadoop.hbase.master.MasterCoproprocessorHost.preCreateTable(MasterCoproprocessorHost.java:348)
at org.apache.hadoop.hbase.master.HMaster$5.run(HMaster.java:2192)
at
org.apache.hadoop.hbase.master.procedure.MasterProcedureUtil.submitProcedure(MasterProcedureUtil.java:1
34)
at org.apache.hadoop.hbase.master.HMaster.createTable(HMaster.java:2189)
at org.apache.hadoop.hbase.master.MasterRpcServices.createTable(MasterRpcServices.java:711)
at org.apache.hadoop.hbase.shaded.protobuf.generated.MasterProtos$MasterService
$2.callBlockingMethod(MasterProtos.java)
at org.apache.hadoop.hbase.ipc.RpcServer.call(RpcServer.java:458)
at org.apache.hadoop.hbase.ipc.CallRunner.run(CallRunner.java:133)
at org.apache.hadoop.hbase.ipc.RpcExecutor$Handler.run(RpcExecutor.java:338)
at org.apache.hadoop.hbase.ipc.RpcExecutor$Handler.run(RpcExecutor.java:318)
```

### 回答

确认当前使用的账号是否具有足够的权限。

需要CTBase用户在Ranger界面配置权限策略，赋予CTBase元数据表\_ctmeta\_、聚簇表和索引表RWCAE（READ，WRITE，EXEC，CREATE，ADMIN）权限。

Ranger界面配置权限操作请参考[使用Ranger（MRS 3.x）](#)。

## 8.10 HBase 故障排除

## 8.10.1 HBase 客户端连接服务端时，长时间无法连接成功

### 问题

在HBase服务端出现问题，无法提供服务，此时HBase客户端进行表操作，会出现该操作挂起，长时间无任何反应。

### 回答

#### 问题分析

当HBase服务端出现问题，HBase客户端进行表操作的时候，会进行重试，并等待超时。该超时默认值为Integer.MAX\_VALUE (2147483647 ms)，所以HBase客户端会在这么长的时间内一直重试，造成挂起表象。

#### 解决方法

HBase客户端提供两个配置项来控制客户端的重试超时方式，如表8-25。

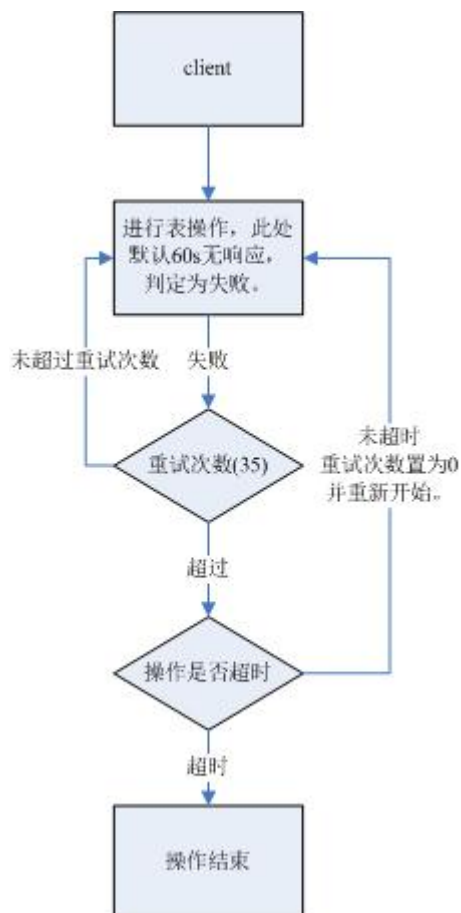
在“客户端安装路径/HBase/hbase/conf/hbase-site.xml”配置文件中配置如下参数。

表 8-25 HBase 客户端操作重试超时相关配置

配置参数	描述	默认值
hbase.client.operation.timeout	客户端操作超时时间。需在配置文件中手动添加。	2147483647 ms
hbase.client.retries.number	最大重试次数。用于表示所有可重试操作所支持的最大重试次数。	35

这两个参数的重试超时的配合方式如图8-3所示。

图 8-3 HBase 客户端操作重试超时流程



从该流程可以看出，如果未对这两个配置参数根据具体使用场景进行配置，会造成挂起迹象。建议根据使用场景，配置合适的超时时间，如果是长时间操作，则把超时时间设置长一点；如果是短时间操作，则把超时时间设置短一点。而重试次数可以设置为：“(hbase.client.retries.number)\*60\*1000(ms)”。刚好大于“hbase.client.operation.timeout”设置的超时时间。

## 8.10.2 在 HBase 连续对同一个表名做删除创建操作时出现创建表异常

### 问题

在HBase连续对同一个表名做删除创建操作时，可能出现创建表异常。

### 回答

执行过程：Disable Table > Drop Table > Create Table > Disable Table > Drop Table >...

1. 在Disable表时，HMaster会发送RPC请求到RegionServer，RegionServer会将相关Region下线。当RegionServer上的Region关闭所需的时间超过HBase的HMaster等待Region处于RIT状态的超时时间，HMaster会默认该Region下线，实际上该Region可能还处在flush memstore阶段。

2. 发送RPC请求关闭Region之后，HMaster会判断该表的所有Region是否下线，上述1的情况下关闭超时也会认为是下线，然后HMaster返回关闭成功。
3. 关闭成功之后，删除表，HBase表对应的数据目录被删掉。
4. 在删除表之后，该数据目录会被还处于flush memstore阶段的Region重新创建。
5. 再创建该表时，将temp目录拷贝到HBase数据目录时，由于HBase数据目录不为空，导致调用HDFS rename接口时，数据目录变为temp目录最后一层追加到HBase的数据目录下，如\$rootDir/data/\$nameSpace/\$tableName/\$tableName，那样创建表就会失败。

#### 解决办法：

出现该问题时，请检查该表对应的HBase数据目录是否存在，如果存在请将该目录重命名。

HBase数据目录由\$rootDir/data/\$nameSpace/\$tableName组成，例如“hdfs://hacluster/hbase/data/default/TestTable”，其中\$rootDir是HBase的根目录，该值通过在“hbase-site.xml”中配置hbase.rootdir.perms得到，data目录是HBase的固定目录，\$nameSpace是nameSpace名字，\$tableName是表名。

### 8.10.3 HBase 占用网络端口，连接数过大会导致其他服务不稳定

#### 问题

HBase占用网络端口，连接数过大会导致其他服务不稳定。

#### 回答

使用操作系统命令*lsof*或者*netstat*发现大量TCP连接处于CLOSE\_WAIT状态，且连接持有者为HBase RegionServer，可能导致网络端口耗尽或HDFS连接超限，那样可能会导致其他服务不稳定。HBase CLOSE\_WAIT现象为HBase机制。

HBase CLOSE\_WAIT产生原因：HBase数据以HFile形式存储在HDFS上，这里可以叫StoreFiles，HBase作为HDFS的客户端，HBase在创建StoreFile或启动加载StoreFile时创建了HDFS连接，当创建StoreFile或加载StoreFile完成时，HDFS方面认为任务已完成，将连接关闭权交给HBase，但HBase为了保证实时响应，有请求时就可以连接对应数据文件，需要保持连接，选择不关闭连接，所以连接状态为CLOSE\_WAIT（需客户端关闭）。

什么时候会创建StoreFile：当HBase执行Flush时。

什么时候执行Flush：HBase写入数据首先会存在内存memstore，只有内存使用达到阈值或手动执行*flush*命令时会触发flush操作，将数据写入HDFS。

#### 解决方法：

由于HBase连接机制，若想减小HBase端口占用，则需控制StoreFile数量，具体可以通过触发HBase的compaction动作完成，即触发HBase文件合并，方法如下：

方法1：使用HBase shell客户端，在客户端手动执行*major\_compact*操作。

方法2：编写HBase客户端代码，调用HBaseAdmin类中的compact方法触发HBase的compaction动作。

如果compact无法解决HBase端口占用现象，说明HBase使用情况已经达到瓶颈，需考虑如下几点：

- table的Region数初始设置是否合适。
- 是否存在无用数据。

若存在无用数据，可删除对应数据以减小HBase存储文件数量，若以上情况都不满足，则需考虑扩容。

## 8.10.4 有 210000 个 map 和 10000 个 reduce 的 HBase BulkLoad 任务运行失败

### 问题

MRS 3.x及之后版本HBase bulkLoad任务（单个表有26T数据）有210000个map和10000个reduce，任务失败。

### 回答

**ZooKeeper IO瓶颈观测手段：**

1. 通过Manager的监控页面查看单个节点上ZooKeeper请求监控，判断是否严重超出规格限制。
2. 通过观测ZooKeeper的日志以及HBase的日志，查看是否有大量的IO Exception Timeout或者SocketTimeout Exception异常。

**调优建议：**

1. 将ZooKeeper实例个数调整为5个及以上，可以通过设置peerType=observer来增加observer的数目。
2. 通过控制单个任务并发的map数或减少每个节点下运行task的内存，降低节点负载。
3. 升级ZooKeeper数据磁盘，如SSD等。

## 8.10.5 使用 scan 命令仍然可以查询到已修改和已删除的数据

### 问题

为什么使用如下scan命令仍然可以查询到已修改和已删除的数据？

```
scan '<table_name>',{FILTER=>"SingleColumnValueFilter('<column_family>','column',=,'binary:<value>')"} }
```

### 回答

由于HBase的可扩展性，在查询表的时候，默认情况下会匹配被查询列的所有版本的值，即使被删除或被修改的值也可以查询出来。对于命中列失败的行（即在某一行中不存在该列），HBase会将该行查询出来。

如果用户仅需查询该表的最新值和命中列成功的行，可使用如下查询语句：

```
scan '<table_name>',
{FILTER=>"SingleColumnValueFilter('<column_family>','column',=,'binary:<value>',true,true)"} }
```

使用该命令，不但可以过滤掉命中列失败的行，而且查询的是表的当前数据的最新版本的值，即不查询被修改之前的值和被删除的值。

## 📖 说明

过滤器SingleColumnValueFilter的相关参数说明如下：

SingleColumnValueFilter(final byte[] family, final byte[] qualifier, final CompareOp compareOp, ByteArrayComparable comparator, final boolean filterIfMissing, final boolean latestVersionOnly)

参数说明：

- family：需要查询的列所在的列族；
- qualifier：需要查询的列；
- compareOp：比较符，如“=”、“>”等；
- comparator：需要查找的目标值；
- filterIfMissing：如果某一行不存在该列，是否过滤，默认值为false；
- latestVersionOnly：是否仅查询最新版本的值，默认值为false。

## 8.10.6 如何处理由于 Region 处于 FAILED\_OPEN 状态而造成的建表失败异常

### 问题

如何处理由于Region处于FAILED\_OPEN状态而造成的建表失败异常。

### 回答

建表过程中如果发生网络故障、HDFS故障或者Active HMaster故障等情况时，可能会造成部分Region上线失败而处于FAILED\_OPEN状态，导致建表失败。

由于Region上线失败而处于FAILED\_OPEN状态造成的建表失败异常不能直接修复，需要删除该表后重新建表。

操作步骤如下：

1. 在集群客户端使用如下命令修复表的状态。  
**hbase hbck -j \${CLIENT\_HOME}/HBase/hbase/tools/hbase-hbck2-1.1.0-h0.cbu.mrs.\*.jar setTableState <table\_name> ENABLED**
2. 进入HBase shell并执行以下命令完成表的清理。  
**disable '<table\_name>'**  
**drop '<table\_name>'**
3. 使用建表命令重新创建该表。

## 8.10.7 如何清理由于建表失败残留在 ZooKeeper 中的 table-lock 节点下的表名

### 问题

安全模式下，由于建表失败，在ZooKeeper的table-lock节点（默认路径/hbase/table-lock）下残留有新建的表名，请问该如何清理？

### 回答

操作步骤如下：

1. 在安装好客户端的环境下，使用hbase用户进行kinit认证。
2. 执行 **hbase zkcli**命令进入ZooKeeper命令行。
3. 在ZooKeeper命令行中执行 **ls /hbase/table**，查看新建的表名是否存在。
  - 是，执行 **ls /hbase/table-lock**查看新建的表名是否存在，若存在新建的表时使用 **delete /hbase/table-lock/ <table>**命令删除该表，其中<table>为残留的表名。
  - 否，结束。

## 8.10.8 为什么给 HBase 使用的 HDFS 目录设置 quota 会造成 HBase 故障

### 问题

为什么给HDFS上的HBase使用的目录设置quota会造成HBase故障？

### 回答

表的flush操作是在HDFS中写memstore数据。

如果HDFS目录没有足够的磁盘空间quota，flush操作会失败，这样region server将会终止。

```
Caused by: org.apache.hadoop.hdfs.protocol.DSQuotaExceededException: The DiskSpace quota of /hbase/
data/<namespace>/<tableName> is exceeded: quota = 1024 B = 1 KB but disk space consumed = 402655638
B = 384.00 MB
?at
org.apache.hadoop.hdfs.server.namenode.DirectoryWithQuotaFeature.verifyStorageSpaceQuota(DirectoryWith
QuotaFeature.java:211)
?at
org.apache.hadoop.hdfs.server.namenode.DirectoryWithQuotaFeature.verifyQuota(DirectoryWithQuotaFeatu
re.java:239)
?at org.apache.hadoop.hdfs.server.namenode.FSDirectory.verifyQuota(FSDirectory.java:882)
?at org.apache.hadoop.hdfs.server.namenode.FSDirectory.updateCount(FSDirectory.java:711)
?at org.apache.hadoop.hdfs.server.namenode.FSDirectory.updateCount(FSDirectory.java:670)
?at org.apache.hadoop.hdfs.server.namenode.FSDirectory.addBlock(FSDirectory.java:495)
```

上述异常中，表“/hbase/data/<namespace>/<tableName>”的磁盘空间quota值为1KB，但是memstore数据为384.00MB，所以flush操作失败并且region server会终止。

在region server终止时，HMaster对终止的region server的WAL文件进行replay操作以恢复数据。由于限制了磁盘空间quota值，导致WAL文件的replay操作失败进而导致HMaster进程异常退出。

```
2016-07-28 19:11:40,352 | FATAL | MASTER_SERVER_OPERATIONS-10-91-9-131:16000-0 | Caught throwable
while processing event M_SERVER_SHUTDOWN |
org.apache.hadoop.hbase.master.HMaster.abort(HMaster.java:2474)
java.io.IOException: failed log splitting for 10-91-9-131,16020,1469689987884, will retry
?at
org.apache.hadoop.hbase.master.handler.ServerShutdownHandler.resubmit(ServerShutdownHandler.java:365
)
?at
org.apache.hadoop.hbase.master.handler.ServerShutdownHandler.process(ServerShutdownHandler.java:220)
?at org.apache.hadoop.hbase.executor.EventHandler.run(EventHandler.java:129)
?at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1142)
?at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:617)
?at java.lang.Thread.run(Thread.java:745)
Caused by: java.io.IOException: error or interrupted while splitting logs in [hdfs://hacluster/hbase/WALs/<RS-
HostName>,<RS-Port>,<startcode>-splitting] Task = installed = 6 done = 3 errors = 3
?at org.apache.hadoop.hbase.master.SplitLogManager.splitLogDistributed(SplitLogManager.java:290)
```



```
?at org.apache.hadoop.hbase.master.MasterFileSystem.splitLog(MasterFileSystem.java:402)
?at org.apache.hadoop.hbase.master.MasterFileSystem.splitLog(MasterFileSystem.java:375)
```

因此，不支持用户对HDFS上的HBase目录进行quota值设置。上述问题可通过下述步骤解决：

- 步骤1** 在客户端命令提示符下运行 `kinit 用户名` 命令，使HBase用户获得安全认证。
- 步骤2** 运行 `hdfs dfs -count -q/hbase/data/<namespace>/<tableName>` 命令检查分配的磁盘空间quota。
- 步骤3** 使用下列命令取消quota值限制，恢复HBase。

```
hdfs dfsadmin -clrSpaceQuota/hbase/data/<namespace>/<tableName>
```

----结束

## 8.10.9 使用 OfflineMetaRepair 工具重新构建元数据后 HMaster 启动失败

### 问题

为什么在使用OfflineMetaRepair工具重新构建元数据后，HMaster启动的时候会等待namespace表分配超时，最后启动失败？

且HMaster将输出下列FATAL消息表示中止：

```
2017-06-15 15:11:07,582 FATAL [Hostname:16000.activeMasterManager] master.HMaster: Unhandled
exception. Starting shutdown.
java.io.IOException: Timedout 120000ms waiting for namespace table to be assigned
 at org.apache.hadoop.hbase.master.TableNamespaceManager.start(TableNamespaceManager.java:98)
 at org.apache.hadoop.hbase.master.HMaster.initNamespace(HMaster.java:1054)
 at org.apache.hadoop.hbase.master.HMaster.finishActiveMasterInitialization(HMaster.java:848)
 at org.apache.hadoop.hbase.master.HMaster.access$600(HMaster.java:199)
 at org.apache.hadoop.hbase.master.HMaster$2.run(HMaster.java:1871)
 at java.lang.Thread.run(Thread.java:745)
```

### 回答

当通过OfflineMetaRepair工具重建元数据时，HMaster在启动期间等待所有region server的WAL分割，以避免数据不一致问题。一旦WAL分割完成，HMaster将进行用户region的分配。所以当在集群异常的场景下，WAL分割可能需要很长时间，这取决于多个因素，例如太多的WALs，较慢的I/O，region servers不稳定等。

为确保HMaster能够成功完成所有region server WAL分割，请执行以下步骤：

1. 确保集群稳定，不存在其他问题。如有任何问题，请先修复。
2. 为“`hbase.master.initializationmonitor.timeout`”参数配置一个较大的值，默认值为“3600000”毫秒。
3. 重启HBase服务。

## 8.10.10 HMaster 日志中频繁打印出 FileNotFoundException 信息

### 问题

当集群重启后会进行split WAL操作，在splitWAL期间，HMaster出现不能close log，日志中频繁打印出FileNotFoundException及no lease信息。



## 8.10.11 ImportTsv 工具执行失败报“Permission denied”异常

### 问题

当使用与Region Server相同的Linux用户（例如omm用户）但不同的kerberos用户（例如admin用户）时，为什么ImportTsv工具执行失败报“Permission denied”的异常？

```
Exception in thread "main" org.apache.hadoop.security.AccessControlException: Permission denied:
user=admin, access=WRITE, inode="/user/omm-bulkload/hbase-staging/
partitions_cab16de5-87c2-4153-9cca-a6f4ed4278a6":hbase:hadoop:drwx--x--x
 at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:342)
 at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:315)
 at
org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:23
1)
 at
com.xxx.hadoop.adapter.hdfs.plugin.HWAccessControlEnforce.checkPermission(HWAccessControlEnforce.java:
69)
 at
org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:19
0)
 at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1789)
 at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1773)
 at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkAncestorAccess(FSDirectory.java:1756)
 at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.startFileInternal(FSNamesystem.java:2490)
 at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.startFileInt(FSNamesystem.java:2425)
 at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.startFile(FSNamesystem.java:2308)
 at
org.apache.hadoop.hdfs.server.namenode.NameNodeRpcServer.create(NameNodeRpcServer.java:745)
 at
org.apache.hadoop.hdfs.protocolPB.ClientNamenodeProtocolServerSideTranslatorPB.create(ClientNamenodeP
rotocolServerSideTranslatorPB.java:434)
 at org.apache.hadoop.hdfs.protocol.proto.ClientNamenodeProtocolProtos$ClientNamenodeProtocol
$2.callBlockingMethod(ClientNamenodeProtocolProtos.java)
 at org.apache.hadoop.ipc.ProtobufRpcEngine$Server
$ProtoBufRpcInvoker.call(ProtobufRpcEngine.java:616)
 at org.apache.hadoop.ipc.RPC$Server.call(RPC.java:973)
 at org.apache.hadoop.ipc.Server$Handler$1.run(Server.java:2260)
 at org.apache.hadoop.ipc.Server$Handler$1.run(Server.java:2256)
 at java.security.AccessController.doPrivileged(Native Method)
 at javax.security.auth.Subject.doAs(Subject.java:422)
 at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1781)
 at org.apache.hadoop.ipc.Server$Handler.run(Server.java:2254)
```

### 回答

ImportTsv工具在“客户端安装路径/HBase/hbase/conf/hbase-site.xml”文件中“hbase.fs.tmp.dir”参数所配置的HBase临时目录中创建partition文件。因此客户端（kerberos用户）应该在指定的临时目录上具有rwx的权限来执行ImportTsv操作。“hbase.fs.tmp.dir”参数的默认值为“/user/\${user.name}/hbase-staging”（例如“/user/omm/hbase-staging”），此处“\${user.name}”是操作系统用户名（即omm用户），客户端（kerberos用户，例如admin用户）不具备该目录的rwx权限。

上述问题可通过执行以下步骤解决：

1. 在客户端将“hbase.fs.tmp.dir”参数设置为当前kerberos用户的目录（如“/user/admin/hbase-staging”），或者为客户端（kerberos用户）提供已配置的目录所必需的rwx权限。
2. 重试ImportTsv操作。

## 8.10.12 使用 HBase BulkLoad 导入数据成功，执行相同的查询时却可能返回不同的结果

### 问题

在使用HBase bulkload导入数据时，如果导入的数据存在相同的rowkey值，数据可以导入成功，但是执行相同的查询时可能返回不同的结果。

### 回答

正常情况下，相同rowkey值的数据加载到HBase是有先后顺序的，HBase以最近的时间戳的数据为最新数据，一般的默认查询中，没有指定时间戳的，就会对相同rowkey值的数据仅返回最新数据。

使用bulkload加载数据，由于数据在内存中处理生成HFile，速度是很快的，很可能出现相同rowkey值的数据具有相同时间戳，从而造成查询结果混乱的情况。

建议在建表和数据加载时，设计好rowkey值，尽量避免在同一个数据文件中存在相同rowkey值的情况。

## 8.10.13 HBase 数据恢复任务报错回滚失败

### 问题

HBase恢复任务执行失败后系统自动回滚数据，若页面详情中提示“Rollback recovery failed”信息，表示回滚失败。由于回滚失败后就不会处理数据，所以有可能产生垃圾数据，需要如何解决？

### 回答

在下次执行备份或恢复任务前，需要手动清除这些垃圾数据。

**步骤1** 安装集群客户端，例如安装目录为“/opt/client”。

**步骤2** 使用客户端安装用户，执行`source /opt/client/bigdata_env`命令配置环境变量。

**步骤3** 执行`kinit admin`命令。

**步骤4** 执行`zkCli.sh -server ZooKeeper节点业务IP地址:2181`连接ZooKeeper。

**步骤5** 执行`deleteall /recovering`删除垃圾数据。然后执行`quit`退出ZooKeeper连接。

#### 📖 说明

执行该命令会导致数据丢失，请谨慎操作。

**步骤6** 执行`hdfs dfs -rm -f -r /user/hbase/backup`删除临时数据。

**步骤7** 登录FusionInsight Manager界面，选择“运维 > 备份恢复 > 恢复管理”，在任务列表中对应任务的“操作”列，单击“查询历史”，在弹出的窗口中，在指定一次执行记录前单击▼，即可查看相关的快照名称信息：

```
Snapshot [snapshot name] is created successfully before recovery.
```

**步骤8** 切换到客户端，执行`hbase shell`，然后运行`delete_all_snapshot 'snapshot name.*'`删除临时快照。

----结束

## 8.10.14 HBase RegionServer GC 参数 Xms 和 Xmx 的配置为 31GB，导致 RegionServer 启动失败

### 问题

MRS 3.x及之后版本，查看RegionServer启动失败节点的hbase-omm-\*.out日志，发现日志中存在“An error report file with more information is saved as: /tmp/hs\_err\_pid\*.log”，查看/tmp/hs\_err\_pid\*.log发现日志存在“#Internal Error (vtableStubs\_aarch64.cpp:213), pid=9456, tid=0x0000ffff97fdd200”和“#guarantee(\_\_pc() <= s->code\_end()) failed: overflowed buffer”，表示此问题是由JDK导致，需要如何解决？

### 回答

修复步骤如下：

- 步骤1** 在RegionServer启动失败的某个节点执行 `su - omm`，切换到omm用户。
  - 步骤2** 在omm用户下执行 `java -XX:+PrintFlagsFinal -version |grep HeapBase`，出现如下类似结果。
- ```
uintx HeapBaseMinAddress = 2147483648 {pd product}
```
- 步骤3** 修改“GC_OPTS”中“-Xms”和“-Xmx”的值使其不在32G-HeapBaseMinAddress和32G的值之间，不包括32G和32G-HeapBaseMinAddress的值。
 - 步骤4** 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > HBase > 实例”，选择失败实例，选择“更多 > 重启实例”来重启失败实例。

----结束

8.10.15 在集群内节点使用 LoadIncrementalHFiles 批量导入数据，报错权限不足

问题

在普通集群中手动创建Linux用户，并使用集群内DataNode节点执行批量导入时，为什么LoadIncrementalHFiles工具执行失败报“Permission denied”的异常？

```
2020-09-20 14:53:53,808 WARN [main] shortcircuit.DomainSocketFactory: error creating DomainSocket
java.net.ConnectException: connect(2) error: Permission denied when trying to connect to '/var/run/
FusionInsight-HDFS/dn_socket'
    at org.apache.hadoop.net.unix.DomainSocket.connect0(Native Method)
    at org.apache.hadoop.net.unix.DomainSocket.connect(DomainSocket.java:256)
    at org.apache.hadoop.hdfs.shortcircuit.DomainSocketFactory.createSocket(DomainSocketFactory.java:168)
    at org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.nextDomainPeer(BlockReaderFactory.java:804)
    at
org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.createShortCircuitReplicaInfo(BlockReaderFactory.java
:526)
    at org.apache.hadoop.hdfs.shortcircuit.ShortCircuitCache.create(ShortCircuitCache.java:785)
    at org.apache.hadoop.hdfs.shortcircuit.ShortCircuitCache.fetchOrCreate(ShortCircuitCache.java:722)
    at
org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.getBlockReaderLocal(BlockReaderFactory.java:483)
    at org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.build(BlockReaderFactory.java:360)
    at org.apache.hadoop.hdfs.DFSInputStream.getBlockReader(DFSInputStream.java:663)
    at org.apache.hadoop.hdfs.DFSInputStream.blockSeekTo(DFSInputStream.java:594)
    at org.apache.hadoop.hdfs.DFSInputStream.readWithStrategy(DFSInputStream.java:776)
    at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:845)
    at java.io.DataInputStream.readFully(DataInputStream.java:195)
```

```
at org.apache.hadoop.hbase.io.hfile.FixedFileTrailer.readFromStream(FixedFileTrailer.java:401)
at org.apache.hadoop.hbase.io.hfile.HFile.isHFileFormat(HFile.java:651)
at org.apache.hadoop.hbase.io.hfile.HFile.isHFileFormat(HFile.java:634)
at org.apache.hadoop.hbase.tool.LoadIncrementalHFiles.visitBulkHFiles(LoadIncrementalHFiles.java:1090)
at
org.apache.hadoop.hbase.tool.LoadIncrementalHFiles.discoverLoadQueue(LoadIncrementalHFiles.java:1006)
at
org.apache.hadoop.hbase.tool.LoadIncrementalHFiles.prepareHFileQueue(LoadIncrementalHFiles.java:257)
at org.apache.hadoop.hbase.tool.LoadIncrementalHFiles.doBulkLoad(LoadIncrementalHFiles.java:364)
at org.apache.hadoop.hbase.tool.LoadIncrementalHFiles.run(LoadIncrementalHFiles.java:1263)
at org.apache.hadoop.hbase.tool.LoadIncrementalHFiles.run(LoadIncrementalHFiles.java:1276)
at org.apache.hadoop.hbase.tool.LoadIncrementalHFiles.run(LoadIncrementalHFiles.java:1311)
at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:76)
at org.apache.hadoop.hbase.tool.LoadIncrementalHFiles.main(LoadIncrementalHFiles.java:1333)
```

回答

如果LoadIncrementalHFiles工具依赖的Client在集群内安装，且和DataNode在相同的节点上，在工具执行过程中HDFS会创建短路读提高性能。短路读依赖“/var/run/FusionInsight-HDFS”目录（“dfs.domain.socket.path”），该目录默认权限是750。而当前Linux用户没有权限操作该目录。

上述问题可通过执行以下方法解决：

方法一：创建新用户(推荐使用)。

步骤1 通过Manager页面创建新的用户，该用户属组中默认包含ficommon组。

```
[root@xxx-xxx-xxx-xxx ~]# id test
uid=20038(test) gid=9998(ficommon) groups=9998(ficommon)
```

步骤2 重新执行ImportData。

----结束

方法二：修改当前用户的属组。

步骤1 将该用户添加到ficommon组中。

```
[root@xxx-xxx-xxx-xxx ~]# usermod -a -G ficommon test
[root@xxx-xxx-xxx-xxx ~]# id test
uid=2102(test) gid=2102(test) groups=2102(test),9998(ficommon)
```

步骤2 重新执行ImportData。

----结束

8.10.16 使用 Phoenix Ssqline 脚本报 import argparse 错误

问题

在客户端使用ssqline脚本时，报import argparse错误。

回答

步骤1 以root用户登录安装HBase客户端的节点，使用hbase用户进行安全认证。

步骤2 进入HBase客户端ssqline脚本所在目录执行python3 ssqline.py命令。

----结束

9 使用 HDFS

9.1 HDFS 文件系统目录简介

HDFS是Hadoop的分布式文件系统（Hadoop Distributed File System），实现大规模数据可靠的分布式读写。HDFS针对的使用场景是数据读写具有“一次写，多次读”的特征，而数据“写”操作是顺序写，也就是在文件创建时的写入或者在现有文件之后的添加操作。HDFS保证一个文件在一个时刻只被一个调用者执行写操作，而可以被多个调用者执行读操作。

HDFS文件系统中目录结构如下表所示。

表 9-1 HDFS 文件系统目录结构（适用于 MRS 3.x 之前版本）

| 路径 | 类型 | 简略功能 | 是否可以删除 | 删除的后果 |
|--------------------------------|------|---|--------|---------------------------|
| /tmp/spark/sparkhive-scratch | 固定目录 | 存放Spark JDBCServer中 metastore session临时文件 | 否 | 任务运行失败 |
| /tmp/sparkhive-scratch | 固定目录 | 存放Spark cli方式运行 metastore session临时文件 | 否 | 任务运行失败 |
| /tmp/carbon/ | 固定目录 | 数据导入过程中，如果存在异常CarbonData数据，则将异常数据放在此目录下 | 是 | 错误数据丢失 |
| /tmp/Loader-\${作业名}_\${MR作业id} | 临时目录 | 存放Loader Hbase bulkload 作业的region信息，作业完成后自动删除 | 否 | Loader Hbase Bulkload作业失败 |

| 路径 | 类型 | 简略功能 | 是否可以删除 | 删除的后果 |
|--|------|---|--------|--------------|
| /tmp/logs | 固定目录 | MR任务日志在HDFS上的聚合路径 | 是 | MR任务日志丢失 |
| /tmp/archived | 固定目录 | MR任务日志在HDFS上的归档路径 | 是 | MR任务日志丢失 |
| /tmp/hadoop-yarn/staging | 固定目录 | 保存AM运行作业运行日志、作业概要信息和作业配置属性 | 否 | 任务运行异常 |
| /tmp/hadoop-yarn/staging/history/done_intermediate | 固定目录 | 所有任务运行完成后，临时存放/tmp/hadoop-yarn/staging目录下文件 | 否 | MR任务日志丢失 |
| /tmp/hadoop-yarn/staging/history/done | 固定目录 | 周期性扫描线程定期将done_intermediate的日志文件转移到done目录 | 否 | MR任务日志丢失 |
| /tmp/mr-history | 固定目录 | 存储预加载历史记录文件的路径 | 否 | MR历史任务日志数据丢失 |
| /tmp/hive | 固定目录 | 存放Hive的临时文件 | 否 | 导致Hive任务失败 |
| /tmp/hive-scratch | 固定目录 | Hive运行时生成的临时数据，如会话信息等 | 否 | 当前执行的任任务会失败 |
| /user/{user}/.sparkStaging | 固定目录 | 存储SparkJDBCServer应用临时文件 | 否 | executor启动失败 |
| /user/spark/jars | 固定目录 | 存放Spark executor运行依赖包 | 否 | executor启动失败 |

| 路径 | 类型 | 简略功能 | 是否可以删除 | 删除的后果 |
|--|------|--|--------|---------------------------|
| /user/loader | 固定目录 | 存放loader的作业脏数据以及HBase作业数据的临时存储目录 | 否 | HBase作业失败或者脏数据丢失 |
| /user/loader/
etl_dirty_data_dir | | | | |
| /user/loader/
etl_hbase_putlist_t
mp | | | | |
| /user/loader/
etl_hbase_tmp | | | | |
| /user/mapred | 固定目录 | 存放Hadoop相关的文件 | 否 | 导致Yarn启动失败 |
| /user/hive | 固定目录 | Hive相关数据存储的默认路径，包含依赖的spark lib包和用户默认表数据存储位置等 | 否 | 用户数据丢失 |
| /user/omm-
bulkload | 临时目录 | HBase批量导入工具临时目录 | 否 | HBase批量导入任务失败 |
| /user/hbase | 临时目录 | HBase批量导入工具临时目录 | 否 | HBase批量导入任务失败 |
| /sparkJobHistory | 固定目录 | Spark eventlog数据存储目录 | 否 | HistoryServer服务不可用，任务运行失败 |
| /flume | 固定目录 | Flume采集到HDFS文件系统
中的数据存储目录 | 否 | Flume工作异常 |
| /mr-history/tmp | 固定目录 | MapReduce作业产生的日志
存放位置 | 是 | 日志信息丢失 |
| /mr-history/done | 固定目录 | MR JobHistory Server管理的
日志的存放位置 | 是 | 日志信息丢失 |

| 路径 | 类型 | 简略功能 | 是否可以删除 | 删除的后果 |
|------------------|---------|--|--------|------------------|
| /tenant | 添加租户时创建 | 配置租户在HDFS中的存储目录，系统默认将自动在“/tenant”目录中以租户名称创建文件夹。例如租户“ta1”，默认HDFS存储目录为“tenant/ta1”。第一次创建租户时，系统自动在HDFS根目录创建“/tenant”目录。支持自定义存储路径。 | 否 | 租户不可用 |
| /apps{1~5}/ | 固定目录 | WebHCat使用到Hive的包的路径 | 否 | 执行WebHCat任务会失败 |
| /hbase | 固定目录 | HBase数据存储目录 | 否 | HBase用户数据丢失 |
| /hbaseFileStream | 固定目录 | HFS文件存储目录 | 否 | HFS文件丢失，且无法恢复 |
| /ats/active | 固定目录 | HDFS路径，用于存储活动的应用程序的timeline数据 | 否 | 删除后会导导致tez任务运行失败 |
| /ats/done | 固定目录 | HDFS路径，用于存储完成的应用程序的timeline数据 | 否 | 删除后会自动创建 |
| /flink | 固定目录 | 存放checkpoint任务数据 | 否 | 删除会导致运行任务失败 |

表 9-2 HDFS 文件系统目录结构（适用于 MRS 3.x 及之后版本）

| 路径 | 类型 | 简略功能 | 是否可以删除 | 删除的后果 |
|--|------|---|--------|---------------------------|
| /tmp/spark2x/sparkhive-scratch | 固定目录 | 存放Spark2x JDBCServer中 metastore session临时文件 | 否 | 任务运行失败 |
| /tmp/sparkhive-scratch | 固定目录 | 存放Spark2x cli方式运行 metastore session临时文件 | 否 | 任务运行失败 |
| /tmp/logs/ | 固定目录 | 存放container日志文件 | 是 | container日志不可查看 |
| /tmp/carbon/ | 固定目录 | 数据导入过程中，如果存在异常CarbonData数据，则将异常数据放在此目录下 | 是 | 错误数据丢失 |
| /tmp/Loader-\${作业名}_\${MR作业id} | 临时目录 | 存放Loader Hbase bulkload 作业的region信息，作业完成后自动删除 | 否 | Loader Hbase Bulkload作业失败 |
| /tmp/hadoop-omm/yarn/system/rmstore | 固定目录 | ResourceManager运行状态信息 | 是 | ResourceMan ager重启后状态信息丢失 |
| /tmp/archived | 固定目录 | MR任务日志在HDFS上的归档路径 | 是 | MR任务日志丢失 |
| /tmp/hadoop-yarn/staging | 固定目录 | 保存AM运行作业运行日志、作业概要信息和作业配置属性 | 否 | 任务运行异常 |
| /tmp/hadoop-yarn/staging/history/done_intermediate | 固定目录 | 所有任务运行完成后，临时存放/tmp/hadoop-yarn/staging目录下文件 | 否 | MR任务日志丢失 |
| /tmp/hadoop-yarn/staging/history/done | 固定目录 | 周期性扫描线程定期将 done_intermediate的日志文件转移到done目录 | 否 | MR任务日志丢失 |

| 路径 | 类型 | 简略功能 | 是否可以删除 | 删除的后果 |
|--|------|--|--------|------------------|
| /tmp/mr-history | 固定目录 | 存储预加载历史记录文件的路径 | 否 | MR历史任务日志数据丢失 |
| /tmp/hive-scratch | 固定目录 | Hive运行时生成的临时数据，如会话信息等 | 否 | 当前执行的任务会失败 |
| /user/{user}/.sparkStaging | 固定目录 | 存储SparkJDBCServer应用临时文件 | 否 | executor启动失败 |
| /user/spark2x/jars | 固定目录 | 存放Spark2x executor运行依赖包 | 否 | executor启动失败 |
| /user/loader | 固定目录 | 存放loader的作业脏数据以及HBase作业数据的临时存储目录 | 否 | HBase作业失败或者脏数据丢失 |
| /user/loader/etl_dirty_data_dir | | | | |
| /user/loader/etl_hbase_putlist_tmp | | | | |
| /user/loader/etl_hbase_tmp | | | | |
| /user/oozie | 固定目录 | 存放oozie运行时需要的依赖库，需用户手动上传 | 否 | oozie调度失败 |
| /user/mapred/hadoop-mapreduce-3.1.1.tar.gz | 固定文件 | MR分布式缓存功能使用的各jar包 | 否 | MR分布式缓存功能无法使用 |
| /user/hive | 固定目录 | Hive相关数据存储的默认路径，包含依赖的spark lib包和用户默认表数据存储位置等 | 否 | 用户数据丢失 |
| /user/omm-bulkload | 临时目录 | HBase批量导入工具临时目录 | 否 | HBase批量导入任务失败 |

| 路径 | 类型 | 简略功能 | 是否可以删除 | 删除的后果 |
|----------------------|---------|--|--------|---------------------------|
| /user/hbase | 临时目录 | HBase批量导入工具临时目录 | 否 | HBase批量导入任务失败 |
| /spark2xJobHistory2x | 固定目录 | Spark2x eventlog数据存储目录 | 否 | HistoryServer服务不可用，任务运行失败 |
| /flume | 固定目录 | Flume采集到HDFS文件系统
中的数据存储空间 | 否 | Flume工作异常 |
| /mr-history/tmp | 固定目录 | MapReduce作业产生的日志
存放位置 | 是 | 日志信息丢失 |
| /mr-history/done | 固定目录 | MR JobHistory Server管理的
日志的存放位置 | 是 | 日志信息丢失 |
| /tenant | 添加租户时创建 | 配置租户在HDFS中的存储目录，系统默认将自动在“/tenant”目录中以租户名称创建文件夹。例如租户“ta1”，默认HDFS存储目录为“tenant/ta1”。第一次创建租户时，系统自动在HDFS根目录创建“/tenant”目录。支持自定义存储路径。 | 否 | 租户不可用 |
| /apps{1~5}/ | 固定目录 | WebHCat使用到Hive的包的路径 | 否 | 执行WebHCat任务会失败 |
| /hbase | 固定目录 | HBase数据存储目录 | 否 | HBase用户数据丢失 |
| /hbaseFileStream | 固定目录 | HFS文件存储目录 | 否 | HFS文件丢失，且无法恢复 |

9.2 HDFS 用户权限管理

9.2.1 创建 HDFS 权限角色

操作场景

该任务指导MRS集群管理员在FusionInsight Manager创建并设置HDFS的角色。HDFS角色可设置HDFS目录或文件的读、写和执行权限。

用户在HDFS中对自己创建的目录或文件拥有完整权限，可直接读取、写入以及授权他人访问此HDFS目录与文件。

说明

- 本章节适用于MRS 3.x及后续版本。
- 安全模式支持创建HDFS角色，普通模式不支持创建HDFS角色。
- 如果当前组件使用了Ranger进行权限控制，须基于Ranger配置HDFS相关策略进行权限管理，具体操作可参考[添加HDFS的Ranger访问权限策略](#)。

操作步骤

步骤1 登录FusionInsight Manager，选择“系统 > 权限 > 角色”。

步骤2 单击“添加角色”，然后在“角色名称”和“描述”中输入角色名字与描述。

步骤3 配置资源权限，请参见[表9-3](#)。

表 9-3 设置角色

| 任务场景 | 角色授权操作 |
|------------------------|---|
| 设置HDFS管理员权限 | 在“配置资源权限”的表格中选择“待操作集群的名称 > HDFS”，勾选“集群管理操作权限”。
说明
设置HDFS管理员权限需要重启HDFS服务才可生效。 |
| 设置用户执行HDFS检查和HDFS修复的权限 | 1. 在“配置资源权限”的表格中选择“待操作集群的名称 > HDFS > 文件系统”。
2. 定位到指定目录或文件在HDFS中保存的位置。
3. 在指定目录或文件的“权限”列，勾选“读”和“执行”。 |
| 设置用户读取其他用户的目录或文件的权限 | 1. 在“配置资源权限”的表格中选择“待操作集群的名称 > HDFS > 文件系统”。
2. 定位到指定目录或文件在HDFS中保存的位置。
3. 在指定目录或文件的“权限”列，勾选“读”和“执行”。 |

| 任务场景 | 角色授权操作 |
|-----------------------------|---|
| 设置用户在其他用户的文件写入数据的权限 | <ol style="list-style-type: none">1. 在“配置资源权限”的表格中选择“待操作集群的名称 > HDFS > 文件系统”。2. 定位到指定文件在HDFS中保存的位置。3. 在指定文件的“权限”列，勾选“写”和“执行”。 |
| 设置用户在其他用户的目录新建或删除子文件、子目录的权限 | <ol style="list-style-type: none">1. 在“配置资源权限”的表格中选择“待操作集群的名称 > HDFS > 文件系统”。2. 定位到指定目录在HDFS中保存的位置。3. 在指定目录的“权限”列，勾选“写”和“执行”。 |
| 设置用户在其他用户的目录或文件执行的权限 | <ol style="list-style-type: none">1. 在“配置资源权限”的表格中选择“待操作集群的名称 > HDFS > 文件系统”。2. 定位到指定目录或文件在HDFS中保存的位置。3. 在指定目录或文件的“权限”列，勾选“执行”。 |
| 设置子目录继承上级目录权限 | <ol style="list-style-type: none">1. 在“配置资源权限”的表格中选择“待操作集群的名称 > HDFS > 文件系统”。2. 定位到指定目录或文件在HDFS中保存的位置。3. 在指定目录或文件的“权限”列，勾选“递归”。 |

步骤4 单击“确定”完成，返回“角色”页面。

---结束

9.2.2 配置 HDFS 用户访问 HDFS 文件权限

配置 HDFS 目录权限

默认情况下，某些HDFS的文件目录权限为777或者750，存在安全风险。建议您在安装完成后修改该HDFS目录的权限，增加用户的安全性。

在HDFS客户端中，使用具有HDFS管理员权限的用户，执行如下命令，将“/user”的目录权限进行修改。

此处将权限修改为“1777”，即在权限处增加“1”，表示增加目录的粘性，即只有创建的用户才可以删除此目录。

```
hdfs dfs -chmod 1777 /user
```

为了系统文件的安全，建议用户将非临时目录进行安全加固，例如：

- /user:777
- /mr-history:777
- /mr-history/tmp:777

- /mr-history/done:777
- /user/mapred:755

配置 HDFS 文件和目录的权限

HDFS支持用户进行文件和目录默认权限的修改。HDFS默认用户创建文件和目录的权限的掩码为“022”，如果默认权限满足不了用户的需求，可以通过配置项进行默认权限的修改。

参数入口：

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-4 参数说明

| 参数 | 描述 | 默认值 |
|---------------------------|---|-----|
| fs.permissions.umask-mode | <p>当客户端在HDFS上创建文件和目录时使用此umask值（用户掩码）。类似于linux上的文件权限掩码。</p> <p>可以使用八进制数字也可以使用符号，例如：“022”（八进制，等同于以符号表示的u=rwx,g=r-x,o=r-x），或者“u=rwx,g=rwx,o=”（符号法，等同于八进制的“007”）。</p> <p>说明
8进制的掩码，和实际权限设置值正好相反，建议使用符号表示法，描述更清晰。</p> | 022 |

9.3 HDFS 客户端使用实践

操作场景

该任务指导用户在运维场景或业务场景中使用HDFS客户端。

前提条件

- 已安装客户端。
例如安装目录为“/opt/client”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- 各组件业务用户由MRS集群管理员根据业务需要创建。安全模式下，“机机”用户需要下载keytab文件。“人机”用户第一次登录时需修改密码。（普通模式不涉及）

使用 HDFS 客户端

步骤1 安装客户端，详细操作请参考[使用MRS客户端](#)。

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤4 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤5 如果集群为安全模式，执行以下命令进行用户认证。普通模式集群无需执行用户认证。

```
kinit 组件业务用户
```

步骤6 执行HDFS Shell命令。例如：

```
hdfs dfs -ls /
```

----结束

HDFS 客户端常用命令

常用的HDFS客户端命令如下表所示。

更多命令可参考https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/CommandsManual.html#User_Commands

表 9-5 HDFS 客户端常用命令

| 命令 | 说明 | 样例 |
|--|-----------------|--|
| <code>hdfs dfs -mkdir 文件夹名称</code> | 创建文件夹 | <code>hdfs dfs -mkdir /tmp/mydir</code> |
| <code>hdfs dfs -ls 文件夹名称</code> | 查看文件夹 | <code>hdfs dfs -ls /tmp</code> |
| <code>hdfs dfs -put 客户端节点上本地文件 HDFS指定路径</code> | 上传本地文件到HDFS指定路径 | <code>hdfs dfs -put /opt/test.txt /tmp</code>
上传客户端节点“/opt/test.txt”文件到HDFS的“/tmp”路径下 |
| <code>hdfs dfs -get HDFS指定文件 客户端节点上指定路径</code> | 下载HDFS文件到本地指定路径 | <code>hdfs dfs -get /tmp/test.txt /opt/</code>
下载HDFS的“/tmp/test.txt”文件到客户端节点的“/opt”路径下 |
| <code>hdfs dfs -rm -r -f HDFS指定文件夹</code> | 删除文件夹 | <code>hdfs dfs -rm -r -f /tmp/mydir</code> |
| <code>hdfs dfs -chmod 权限参数 文件目录</code> | 为用户设置HDFS目录权限 | <code>hdfs dfs -chmod 700 /tmp/test</code> |

客户端常见使用问题

- 问题一：当执行HDFS客户端命令时，客户端程序异常退出，报“java.lang.OutOfMemoryError”的错误。
该问题是由于HDFS客户端运行时所需的内存超过了HDFS客户端设置的内存上限（默认为128MB）。可以通过修改“<客户端安装路径>/HDFS/

component_env”中的“CLIENT_GC_OPTS”来修改HDFS客户端的内存上限。例如，需要设置该内存上限为1GB，则设置：

```
CLIENT_GC_OPTS="-Xmx1G"
```

在修改完后，使用如下命令刷新客户端配置，使之生效：

```
source <客户端安装路径>/bigdata_env
```

- 问题二：如何设置HDFS客户端运行时的日志级别？

HDFS客户端运行时的日志是默认输出到管理控制台的，其级别默认为INFO。如果需要开启DEBUG级别日志，可以通过导出一个环境变量来设置，命令如下：

```
export HADOOP_ROOT_LOGGER=DEBUG,console
```

在执行完上面命令后，再执行HDFS Shell命令时，即可打印出DEBUG级别日志。

如果想恢复INFO级别日志，可执行如下命令：

```
export HADOOP_ROOT_LOGGER=INFO,console
```

- 问题三：如何彻底删除HDFS文件？

由于HDFS的回收站机制，一般删除HDFS文件后，文件会移动到HDFS的回收站中。如果确认文件不再需要并且需要立马释放存储空间，可以继续清理对应的回收站目录（例如：`hdfs://hacluster/user/xxx/.Trash/Current/xxx`）。

9.4 快速使用 Hadoop

本章节提供从零开始使用Hadoop提交wordcount作业的操作指导，wordcount是最经典的Hadoop作业，它用来统计海量文本的单词数量。

操作步骤

步骤1 准备wordcount程序。

开源的Hadoop的样例程序包含多个例子，其中包含wordcount。可以从<https://dist.apache.org/repos/dist/release/hadoop/common/>中下载Hadoop的样例程序。

例如，选择hadoop-x.x.x版本，下载“hadoop-x.x.x.tar.gz”，解压后在“hadoop-x.x.x\share\hadoop\mapreduce”路径下获取“hadoop-mapreduce-examples-x.x.x.jar”，即为Hadoop的样例程序。“hadoop-mapreduce-examples-x.x.x.jar”样例程序包含了wordcount程序。

说明

hadoop-x.x.x表示Hadoop的版本号，具体以实际为准。

步骤2 准备数据文件。

数据文件无格式要求，准备一个或多个txt文件即可，如下内容为txt文件样例：

```
qwdsfhoedfrffrofhuncckgktpmhutopmma  
jjpsffjorgjgtyiuyjmhombmbogohoyhm  
jhheyeombdhuaqqiquyebchdhmamdhdemmj  
doeyhjwedcrftgbmojjyhqssdddddfkf  
kjhjhkehdeiyrudjfhfhffooqweopuyyyy
```

步骤3 上传数据至OBS。

1. 登录OBS控制台。

- 选择“并行文件系统 > 创建并行文件系统”，创建一个名称为wordcount01的文件系统。

wordcount01仅为示例，文件系统名称必须全局唯一，否则会创建并行文件系统失败。



- 在OBS文件系统列表中单击文件系统名称wordcount01，选择“文件 > 新建文件夹”，分别创建program、input文件夹，创建完成后如图9-1所示。

图 9-1 wordcount01 文件系统文件夹列表



- program：存放用户程序
- input：存放用户数据文件

- 进入program文件夹，选择“上传文件 > 添加文件”，从本地选择步骤1中下载的程序包，然后单击“上传”，上传完成后如图9-2所示。

图 9-2 程序列表



- 进入input文件夹，将步骤2中准备的数据文件上传到input文件夹，上传完成后如图9-3所示。

图 9-3 数据文件列表



| <input type="checkbox"/> | 名称 | 存储类别 | 大小 | 加密状态 |
|--------------------------|----------------|------|----------|------|
| ← 返回上一级 | | | | |
| <input type="checkbox"/> | wordcount2.txt | 标准存储 | 23 Bytes | 未加密 |
| <input type="checkbox"/> | wordcount1.txt | 标准存储 | 29 Bytes | 未加密 |

步骤4 登录MRS控制台，在左侧导航栏选择“现有集群”，单击集群名称，该集群需要包含Hadoop组件，且已为MRS集群绑定具有OBS文件系统操作权限的IAM权限委托。

查看或绑定委托的操作如下：

1. 登录MRS集群的“概览”页面，查看“委托”参数是否有值，且绑定的委托具有OBS文件系统操作权限。



- 是，集群已绑定委托。
 - 否，执行**步骤4.2**。
2. 单击“管理委托”，为集群绑定具有OBS文件系统操作权限的委托。
您可以直接选择系统默认的“MRS_ECS_DEFAULT_AGENCY”，也可以单击“新建委托”自行创建其他具有OBS文件系统操作权限的委托。

步骤5 提交wordcount作业。

在MRS控制台选择“作业管理”页签，单击“添加”，进入“添加作业”页面，具体请参见[运行MapReduce作业](#)。

图 9-4 wordcount 作业

添加作业

* 作业类型: MapReduce

* 作业名称: mr_01

* 执行程序路径: obs://[redacted]/program/hadoop-mapreduce-examples-2.7.5.jar [HDFS] [OBS]

执行程序参数: wordcount obs://wordcount01/input/ obs://wordcount01/output/ [HDFS] [OBS]

服务配置参数: 参数 值

命令参考: yarn jar obs://[redacted]/program/hadoop-mapreduce-examples-2.7.5.jar wordcount obs://wordcount01/input/ obs://wordcount01/output/

[确定] [取消]

- 作业类型选择“MapReduce”。
- 作业名称为“mr_01”。
- 执行程序路径配置为OBS上存放程序的地址。例如：obs://wordcount01/program/hadoop-mapreduce-examples-x.x.x.jar。
- 执行程序参数中填写的参数为：wordcount obs://wordcount01/input/ obs://wordcount01/output/。

说明

- 参数“obs://wordcount01/input/”中的OBS文件系统名需要替换为实际环境创建的文件系统名。
- 参数“obs://wordcount01/output/”中的OBS文件系统名需要替换为实际环境创建的文件系统名，目录output为一个不存在的目录，具体以实际为准。
- 服务配置参数无需填写。

只有集群处于“运行中”状态时才能提交作业。

作业提交成功后默认为“已接受”状态，不需要用户手动执行作业。

步骤6 查看作业执行结果。

1. 进入“作业管理”页面，查看作业是否执行完成。
作业运行需要时间，作业运行结束后，刷新作业列表，查看作业列表如图9-5所示。

图 9-5 作业列表

| 作业名称/ID | 用户名称 | 作业类型 | 状态 | 执行结果 | 作业提交时间 | 持续时间(分钟) | 操作 |
|--|------------|-----------|-----|------|-------------------------------|----------|--------------|
| mr_01
8befb25-d5b-46f4-8189-0510c754f3a | [redacted] | MapReduce | 已接受 | 成功 | 2020/04/26 17:49:34 GMT+08:00 | 3.0 | 查看日志 查看详情 更多 |

作业执行成功或失败后都不能再次执行，只能新增或者复制作业，配置作业参数后重新提交作业。

2. 登录OBS控制台，进入OBS路径，查看作业输出信息。

进入到**步骤5**中创建的output路径查看相关的output文件，需要下载到本地以文本方式打开进行查看，如**图9-6**所示。

图 9-6 输出文件列表



----结束

9.5 配置 HDFS 文件回收站机制

配置场景

在HDFS中，如果删除HDFS的文件，删除的文件将被移动到回收站（trash）中，不会被立即清除，以便在误操作的情况下恢复被删除的数据。被删除的文件在超过老化时间后将变为老化文件，会基于系统机制清除或用户手动清除。

您可以设置文件保留在回收站中的时间阈值，一旦文件保存时间超过此阈值，将从回收站中永久地删除。如果回收站被清空，回收站中的所有文件将被永久删除。

配置描述

参数入口：

请参考**修改集群服务配置参数**，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-6 参数说明

| 参数 | 描述 | 默认值 |
|-------------------|--|------|
| fs.trash.interval | 以分钟为单位的垃圾回收时间，垃圾站中数据超过此时间，会被删除。取值范围：1440 ~ 259200。 | 1440 |

| 参数 | 描述 | 默认值 |
|------------------------------|--|-----|
| fs.trash.checkpoint.interval | <p>垃圾检查点间的间隔。单位：分钟。应小于等于“fs.trash.interval”的值。检查点程序每次运行时都会创建一个新的检查点并会移除fs.trash.interval分钟前创建的检查点。例如，系统每10分钟检测是否存在老化文件，如果发现有老化文件，则删除。对于未老化文件，则会存储在checkpoint列表中，等待下一次检查。</p> <p>如果此参数的值设置为0，则表示系统不会检查老化文件，所有老化文件会被保存在系统中。</p> <p>取值范围：0 ~ fs.trash.interval。</p> <p>说明
不推荐将此参数值设置为0，这样系统的老化文件会一直存储下去，导致集群的磁盘空间不足。</p> | 60 |

9.6 配置 HDFS DataNode 数据均衡

操作场景

说明

本章节适用于MRS 3.x及后续版本。

HDFS集群可能出现DataNode节点间磁盘利用率不平衡的情况，比如集群中添加新数据节点的场景。如果HDFS出现数据不平衡的状况，可能导致多种问题，比如MapReduce应用程序无法很好地利用本地计算的优势、数据节点之间无法达到更好的网络带宽使用率或节点磁盘无法利用等等。所以MRS集群管理员需要定期检查并保持DataNode数据平衡。

HDFS提供了一个容量均衡程序Balancer。通过运行这个程序，可以使得HDFS集群达到一个平衡的状态，使各DataNode磁盘使用率与HDFS集群磁盘使用率的偏差不超过阈值。图9-7和图9-8分别是Balance前后DataNode的磁盘使用率变化。

图 9-7 执行均衡操作前 DataNode 的磁盘使用率

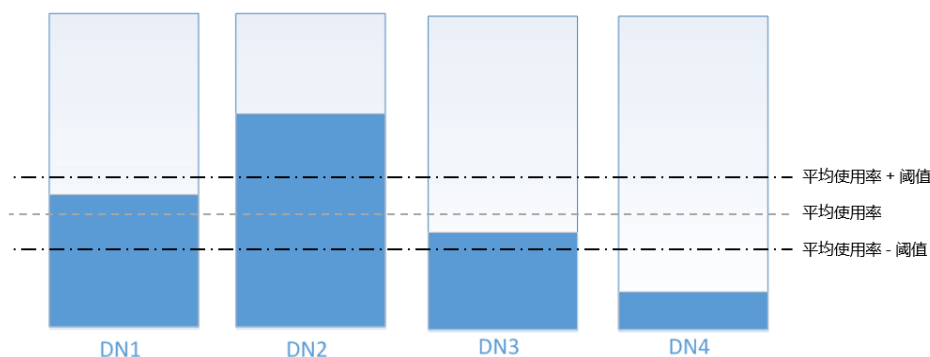
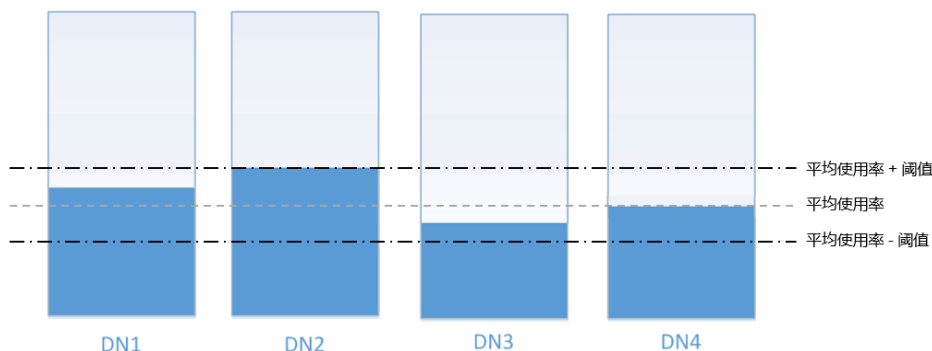


图 9-8 执行均衡操作后 DataNode 的磁盘使用率



均衡操作时间估算受两个因素影响：

1. 需要迁移的总数据量：

每个DataNode节点的数据量应大于（平均使用率-阈值）*平均数据量，小于（平均使用率+阈值）*平均数据量。若实际数据量小于最小值或大于最大值即存在不平衡，系统选择所有DataNode节点中偏差最多的数据量作为迁移的总数据量。

2. Balancer的迁移是按迭代（iteration）方式串行顺序处理的，每个iteration迁移数据量不超过10GB，每个iteration重新计算使用率的情况。

因此针对集群情况，可以大概估算每个iteration耗费的时间（可以通过执行Balancer的日志观察到每次iteration的时间），并用总数据量除以10GB估算任务执行时间。

由于按iteration处理，Balancer可以随时启动或者停止。

对系统的影响

- 执行Balance操作时会占用DataNode的网络带宽资源，请根据业务需求在维护期间执行任务。
- 默认使用带宽控制为20MB/s，如果重新设置带宽流量或加大数据量，Balance操作可能会对正在运行的业务产生影响。

前提条件

已安装HDFS客户端。

配置 Balance 任务

步骤1 使用客户端安装用户登录客户端所在节点。执行命令切换到客户端安装目录，例如“/opt/client”。

```
cd /opt/client
```

📖 说明

如果集群为普通模式，需先执行su - omm切换为omm用户。

步骤2 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤3 如果集群为安全模式，执行以下命令认证hdfs身份。

kinit hdfs

步骤4 是否调整带宽控制？

- 是，执行**步骤5**。
- 否，执行**步骤6**。

步骤5 执行以下命令，修改Balance的最大带宽，然后执行**步骤6**。

```
hdfs dfsadmin -setBalancerBandwidth <bandwidth in bytes per second>
```

<bandwidth in bytes per second>表示带宽控制的数值，单位为字节。例如要设置带宽控制为20MB/s，对应值为20971520，完整命令为：

```
hdfs dfsadmin -setBalancerBandwidth 20971520
```

📖 说明

- 默认为20MB/s，适用于当前集群使用万兆网络，且有业务正在执行的场景。若没有足够的业务空闲时间窗用于Balance维护，可适当增加该值以缩短Balance时间，如增大到209715200（即200MB/s）。
- 这个参数的调整要看组网情况，如果集群负载较高，可以改为209715200(200MB/s)；如果集群空闲，可以改为1073741824 (1GB/s)。
- 如果DataNode节点的带宽无法达到指定的最大带宽，可以在FusionInsight Manager修改HDFS的参数“dfs.datanode.balance.max.concurrent.moves”，将每个DataNode节点执行均衡的线程数修改为“32”，并重启HDFS服务。

步骤6 执行以下命令，启动Balance任务。

```
bash /opt/client/HDFS/hadoop/sbin/start-balancer.sh -threshold <threshold of balancer>
```

-threshold表示HDFS数据达到平衡状态时DataNode磁盘使用率偏差值，各个DataNode节点磁盘的使用率和整体HDFS集群的磁盘空间平均使用率偏差小于此阈值时，系统认为HDFS集群已经达到了平衡的状态并结束Balance任务。

例如，需要设置偏差率为5%，则执行：

```
bash /opt/client/HDFS/hadoop/sbin/start-balancer.sh -threshold 5
```

📖 说明

- “/opt/client”为客户端安装目录，具体请以实际路径替换。
- 上述命令会在后台执行该任务，相关日志可以通过客户端安装目录“/opt/client/HDFS/hadoop/logs”下的hadoop-root-balancer-*主机名*.out查看。
- 如果需要停止Balance任务，请执行以下命令：

```
bash /opt/client/HDFS/hadoop/sbin/stop-balancer.sh
```

- 如果只需要对部分节点进行数据均衡，可以在脚本上加上-include参数指定要移动的节点。具体参数使用方法，可通过命令行查看。

```
例如执行：bash /opt/client/HDFS/hadoop/sbin/start-balancer.sh -threshold 5 -include IP1,IP2,IP3
```

- 如果该命令执行失败，在日志中看到的错误信息为“Failed to APPEND_FILE /system/balancer.id”，则需要执行如下命令强行删除“/system/balancer.id”，再次执行start-balancer.sh脚本即可。

```
hdfs dfs -rm -f /system/balancer.id
```

步骤7 用户在执行了**步骤6**的脚本后，会在客户端安装目录“/opt/client/HDFS/hadoop/logs”目录下生成名为hadoop-root-balancer-主机名.out日志。打开该日志可以看到如下字段信息：

- Time Stamp: 时间戳
- Bytes Already Moved: 已经移动的字节数
- Bytes Left To Move: 待移动的字节数
- Bytes Being Moved: 正在移动的字节数

日志出现“Balancing took xxx seconds”信息表示均衡操作已完成。

----结束

设置自动执行 Balance 任务

步骤1 登录FusionInsight Manager。

步骤2 选择“集群 > 服务 > HDFS > 配置”，选择“全部配置”，搜索以下参数名并修改参数值。

- “dfs.balancer.auto.enable”表示是否启用自动执行Balance任务，默认值为“false”表示不启用，修改为“true”表示启用。
- “dfs.balancer.auto.cron.expression”表示任务执行的时间，默认值“0 1 * * 6”表示在每周六的1点执行任务。仅在启用自动执行Balance功能时有效。修改此参数时，表达式介绍如**表9-7**所示。支持“*”表示连续的时间段。

表 9-7 执行表达式参数解释

| 列 | 说明 |
|-----|--------------------|
| 第1列 | 分钟，参数值为0~59。 |
| 第2列 | 小时，参数值为0~23。 |
| 第3列 | 日期，参数值为1~31。 |
| 第4列 | 月份，参数值为1~12。 |
| 第5列 | 星期，参数值为0~6，0表示星期日。 |

- “dfs.balancer.auto.stop.cron.expression”表示任务自动停止的时间，默认值为空，表示不自动停止正在运行的Balancer任务。以“0 5 * * 6”为例，则表示在每周六的5点停止正在运行的Balancer任务。仅在启用自动执行Balance功能时有效。

修改此参数时，表达式介绍如**表9-7**所示。支持“*”表示连续的时间段。

步骤3 修改自动Balancer的运行参数，如**表9-8**所示：

表 9-8 自动 Balancer 运行参数

| 参数名 | 参数介绍 | 默认值 |
|-------------------------------------|---|------|
| dfs.balancer.aut.threshold | 表示磁盘容量百分比的均衡阈值。仅当dfs.balancer.auto.enable设置为true时才有效。 | 10 |
| dfs.balancer.auto.exclude.datanodes | 不需要执行磁盘自动均衡的DataNode列表，用逗号分隔。仅当dfs.balancer.auto.enable设置为true时才有效。 | 默认为空 |
| dfs.balancer.auto.bandwidthPerSec | 每个DataNode可用于负载均衡的最大带宽量（单位：MB/s）。 | 20 |
| dfs.balancer.auto.maxIdleIterations | Balancer的最大连续空闲迭代次数。一次空闲迭代为没有Block块被移动的迭代，当连续空闲迭代次数达到最大连续空闲迭代次数时，本次Balancer结束。当取值为-1时，代表无穷大。 | 5 |
| dfs.balancer.auto.maxDataNodesNum | 该参数用来控制进行自动Balancer的DataNode数量。假设该参数值为N，当N大于0，则选择剩余空间比例最高的N个DataNode和最低的N个DataNode之间进行数据均衡；当N等于0，则对集群中所有DataNode进行数据均衡。 | 5 |

步骤4 单击“保存”使配置生效。无需重启HDFS服务。

任务执行日志保存在主NameNode节点中，请查看“/var/log/Bigdata/hdfs/nn/hadoop-omm-balancer-主机名.log”。

----结束

9.7 配置 HDFS DiskBalancer 磁盘均衡

配置场景

DiskBalancer是一个在线磁盘均衡器，旨在根据各种指标重新平衡正在运行的DataNode上的磁盘数据。工作方式与HDFS的Balancer工具类似。不同的是，HDFS Balancer工具用于DataNode节点间的数据均衡，而HDFS DiskBalancer用于单个DataNode节点上各磁盘之间的数据均衡。

长时间运行的集群会因为曾经删除过大量的文件，或者集群中的节点做磁盘扩容等操作导致节点上出现磁盘间数据不均衡的现象。磁盘间数据不均衡会引起HDFS整体并发读写性能的下降或者因为不恰当的HDFS写策略导致业务故障。此时需要平衡节点磁盘间的数据密度，防止异构的小磁盘成为该节点的性能瓶颈。

 说明

本章节适用于MRS 3.x及后续版本。

配置描述

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-9 参数说明

| 参数 | 描述 | 默认值 |
|---|---|-----------|
| dfs.disk.balancer.auto.enabled | 是否开启自动执行HDFS DiskBalancer特性。默认值为“false”，表示关闭该特性。 | false |
| dfs.disk.balancer.auto.cron.expression | HDFS 磁盘均衡操作的CRON表达式，用于控制均衡操作的开始时间。仅当dfs.disk.balancer.auto.enabled设置为true时才有效。默认值“0 1 * * 6”表示在每周六的1点执行任务。表达式的具体含义可参见 表9-10 。 | 0 1 * * 6 |
| dfs.disk.balancer.max.disk.throughputInMBperSec | 执行磁盘数据均衡时可使用的最大磁盘带宽。单位为MB/s，默认值为10，可依据集群的实际磁盘条件设置。 | 10 |
| dfs.disk.balancer.max.disk.errors | 设置能够容忍的在指定的移动过程中出现的最大错误次数，超过此阈值则移动失败。 | 5 |
| dfs.disk.balancer.block.tolerance.percent | 设置磁盘之间进行数据均衡操作时，各个磁盘的数据存储量与理想状态之间的差异阈值。例如，各个磁盘的理想数据存储量为1TB，此参数设置为10。那么，当目标磁盘的数据存储量达到900GB时，就认为该磁盘的存储状态就已经足够好了。取值范围[1-100]。 | 10 |
| dfs.disk.balancer.plan.threshold.percent | 设置在磁盘数据均衡中可容忍的两磁盘之间的数据密度阈值差。如果任意两个磁盘数据密度差值的绝对值超过了此阈值，意味着对应的磁盘应该进行数据均衡。取值范围[1-100]。 | 10 |
| dfs.disk.balancer.top.nodes.number | 该参数用来指定集群中需要执行磁盘数据均衡的Top N 节点。 | 5 |

[表9-10](#)为HDFS磁盘均衡操作的CRON表达式。使用此功能时，需要先将参数dfs.disk.balancer.auto.enabled设置为true。其它参数依据集群状况设置。

表 9-10 CRON 表达式解释

| 列 | 说明 |
|-----|--------------------|
| 第1列 | 分钟，参数值为0~59。 |
| 第2列 | 小时，参数值为0~23。 |
| 第3列 | 日期，参数值为1~31。 |
| 第4列 | 月份，参数值为1~12。 |
| 第5列 | 星期，参数值为0~6，0表示星期日。 |

使用限制

1. 只支持同类型磁盘之间的数据移动，例如SSD->SSD，DISK->DISK等。
2. 执行该特性会占用涉及节点的磁盘IO资源、网络带宽资源，请尽量在业务不繁忙的时候使用。
3. 参数dfs.disk.balancer.top.nodes.number指定Top N节点返回的DataNode列表是不断重新计算的，因此不必设置的过大。
4. 如果要在HDFS客户端通过命令行使用DiskBalancer功能，其接口如下：

表 9-11 DiskBalancer 功能的接口说明

| 命令格式 | 说明 |
|---|--|
| hdfs diskbalancer -report -top <N> | N可以指定为大于0的整数，先利用此条命令查询集群中最需要执行磁盘数据均衡的Top N节点。 |
| hdfs diskbalancer -plan <Hostname IP Address> | 此条命令可以根据传入的DataNode生成一个Json文件，该文件包含了数据移动的源磁盘、目标磁盘、待移动的块等信息。同时，该命令还支持指定一些其他网络带宽参数等。 |
| hdfs diskbalancer -query <Hostname:\$dfs.datanode.ipc.port> | 集群默认的port值为9867。此条命令可以查询当前节点上运行的DiskBalancer任务的运行状态。 |
| hdfs diskbalancer -execute <planfile> | 此命令中的planfile指第二条命令中生成的Json文件，请使用绝对路径。 |
| hdfs diskbalancer -cancel <planfile> | 取消正在运行的planfile，同样需要使用绝对路径。 |

说明

- 在客户端执行此命令时，用户需要具备supergroup权限。可以使用HDFS服务的系统用户hdfs。或者在集群上创建一个具有supergroup权限的用户，再在客户端中执行此命令。
- [表9-11](#)只说明了命令接口的含义及使用方法，实际每个接口提供了更多的配置参数。具体信息可通过`hdfs diskbalancer -help <command>`命令查看。
- 在集群运维过程中，排查性能类问题时，可查看集群的事件信息中是否有HDFS磁盘均衡任务事件发生，如果有，可以排查集群中是否开启了DiskBalancer。
- 自动执行磁盘均衡的特性开启以后，会在此次数据均衡执行完成之后才会退出。无法在执行均衡中途取消本次执行任务。
- 如果想要灵活选择某些指定节点进行数据均衡，可以在客户端手动指定执行。

9.8 配置 HDFS Mover 命令迁移数据

配置场景

Mover是一个新的数据迁移工具，工作方式与HDFS的Balancer接口工作方式类似。Mover能够基于设置的数据存储策略，将集群中的数据重新分布。

通过运行Mover，周期性地检测HDFS文件系统中用户指定的HDFS文件或目录，判断该文件或目录是否满足设置的存储策略，如果不满足，则进行数据迁移，使目标目录或文件满足设定的存储策略。

说明

本章节适用于MRS 3.x及后续版本。

配置描述

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-12 参数说明

| 参数 | 描述 | 默认值 |
|----------------------------------|---|-----------|
| dfs.mover.auto.enable | 是否开启数据副本迁移功能，该功能支持多种。默认值为“false”，表示关闭该特性。 | false |
| dfs.mover.auto.cron.expression | HDFS执行自动数据迁移的CRON表达式，用于控制数据迁移操作的开始时间。仅当dfs.mover.auto.enable设置为true时才有效。默认值“0 * * * *”表示在每个整点执行任务。表达式的具体含义可参见 表9-13 。 | 0 * * * * |
| dfs.mover.auto.hdfsfiles_or_dirs | 指定集群执行自动副本迁移的HDFS文件或目录列表，以空格分隔。仅当dfs.mover.auto.enable设置为true时才有效。 | - |

表 9-13 Cron 表达式解释

| 列 | 说明 |
|-----|--------------------|
| 第1列 | 分钟，参数值为0~59。 |
| 第2列 | 小时，参数值为0~23。 |
| 第3列 | 日期，参数值为1~31。 |
| 第4列 | 月份，参数值为1~12。 |
| 第5列 | 星期，参数值为0~6，0表示星期日。 |

使用说明

若要在HDFS的客户端通过命令行执行mover功能，其命令格式如下：

```
hdfs mover -p <HDFS文件全路径或目录路径>
```

说明

在客户端执行此命令时，用户需要具备supergroup权限。可以使用HDFS服务的系统用户hdfs。或者在集群上创建一个具有supergroup权限的用户，再在客户端中执行此命令。

9.9 配置 HDFS 文件目录标签策略（NodeLabel）

配置场景

用户需要通过数据特征灵活配置HDFS文件数据块的存储节点。通过设置HDFS目录/文件对应一个标签表达式，同时设置每个DataNode对应一个或多个标签，从而给文件的数据块存储指定了特定范围的DataNode。

当使用基于标签的数据块摆放策略，为指定的文件选择DataNode节点进行存放时，会根据文件的标签表达式选择出DataNode节点范围，然后在这些DataNode节点范围内，选择出合适的存放节点。

说明

本章节适用于MRS 3.x及后续版本。

- 场景1 DataNodes分区场景。

场景说明：

用户需要让不同的应用数据运行在不同的节点，分开管理，就可以通过标签表达式，来实现不同业务的分离，指定业务存放对应的节点上。

通过配置NodeLabel特性使得：

- /HBase下的数据存储DN1、DN2、DN3、DN4节点上。
- /Spark下的数据存储DN5、DN6、DN7、DN8节点上。

图 9-9 DataNode 分区场景



说明

- 通过 `hdfs nodelabel -setLabelExpression -expression 'LabelA[fallback=NONE]' -path /Hbase` 命令，给 Hbase 目录设置表达式。从图 9-9 中可知，“/Hbase” 文件的数据块副本会被放置在有 LabelA 标签的节点上，即 DN1、DN2、DN3、DN4。同理，通过 `hdfs nodelabel -setLabelExpression -expression 'LabelB[fallback=NONE]' -path /Spark` 命令，给 Spark 目录设置表达式。在“/Spark” 目录下文件对应的数据块副本只能放置到 LabelB 标签上的节点，如 DN5、DN6、DN7、DN8。
- 设置数据节点的标签参考 [配置描述](#)。
- 如果同一个集群上存在多个机架，每个标签下可以有多个机架的 DataNodes，以确保数据块摆放的可靠性。

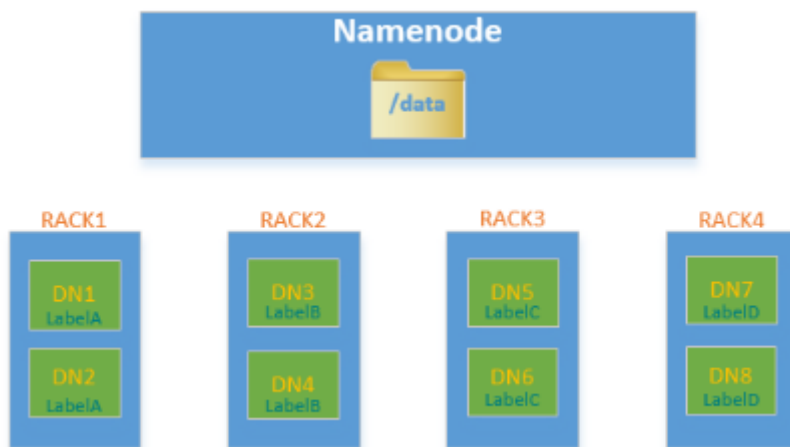
场景2 多机架下指定副本位置场景

场景说明：

在异构集群中，需要分配一些特定的具有高可靠性的节点用以存放重要的商业数据，可以通过标签表达式指定副本位置，指定文件数据块的其中一个副本存放到高可靠性的节点上。

“/data” 目录下的数据块，默认三副本情况下，其中至少有一个副本会被存放到 RACK1 或 RACK2 机架的节点上（RACK1 和 RACK2 机架的节点为高可靠性节点），另外两个副本会被分别存放到 RACK3 和 RACK4 机架的节点上。

图 9-10 场景样例



说明

通过 `hdfs nodelabel -setLabelExpression -expression 'LabelA|LabelB[fallback=NONE],LabelC,LabelD' -path /data` 命令给 “/data” 目录设置表达式。

当向 “/data” 目录下写数据时，至少有一个数据块副本存放在LabelA或者LabelB标签的节点中，剩余的两个数据块副本会被存放在有LabelC和LabelD标签的节点上。

配置描述

- DataNode节点标签配置

请参考 [修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-14 参数说明

| 参数 | 描述 | 默认值 |
|---|---|---|
| <code>dfs.block.replicator.classname</code> | 配置HDFS的DataNode原则策略。
如果需要开启NodeLabel功能，需要将该值设置为
<code>org.apache.hadoop.hdfs.server.blockmanagement.BlockPlacementPolicyWithNodeLabel</code> 。 | <code>org.apache.hadoop.hdfs.server.blockmanagement.AvailableSpaceBlockPlacementPolicy</code> |
| <code>host2tags</code> | 配置DataNode主机与标签的对应关系。
主机名称支持配置IP扩展表达式（如192.168.1.[1-128]或者192.168.[2-3].[1-128]，且IP必须为业务IP），或者为前后加上 / 的主机名的正则表达式（如/datanode-[123]/或者/datanode-\d{2}/）。
标签配置名称不允许包含 = / \ 字符。【注意】配置IP时必须是业务IP。 | - |

说明

- host2tags配置项内容详细说明：

例如某MRS集群有20个DataNode：dn-1到dn-20，对应的IP地址为10.1.120.1到10.1.120.20，host2tags配置文件内容可以使用如下的表示方式。

主机名正则表达式：

“/dn-\d/ = label-1”表示dn-1到dn-9对应的标签为label-1，即dn-1 = label-1，dn-2 = label-1，...dn-9 = label-1。

“/dn-((1[0-9]\$)|(20\$))/ = label-2”表示dn-10到dn-20对应的标签为label-2，即dn-10 = label-2，dn-11 = label-2，...dn-20 = label-2。

IP地址范围表示方式：

“10.1.120.[1-9] = label-1”表示10.1.120.1到10.1.120.9对应的标签为label-1，即10.1.120.1 = label-1，10.1.120.2 = label-1，...10.1.120.9 = label-1。

“10.1.120.[10-20] = label-2”表示10.1.120.10到10.1.120.20对应的标签为label-2，即10.1.120.10 = label-2，10.1.120.11 = label-2，...10.1.120.20 = label-2。

- 基于标签的数据块摆放策略支持扩容减容场景：

当集群中新增加DataNode节点时，如果该DataNode对应的IP匹配host2tags配置项中的IP地址范围，或者该DataNode的主机名匹配host2tags配置项中的主机名正则表达式，则该DataNode节点会被设置成对应的标签。

例如“host2tags”配置值为10.1.120.[1-9] = label-1，而当前集群只有10.1.120.1到10.1.120.3三个数据节点。进行扩容后，又添加了10.1.120.4这个数据节点，则该数据节点会被设置成label-1的标签；如果10.1.120.3这个数据节点被删除或者退出服务后，数据块不会再被分配到该节点上。

- 设置目录/文件的标签表达式
 - 在HDFS参数配置页面配置“path2expression”，配置HDFS目录与标签的对应关系。当配置的HDFS目录不存在时，也可以配置成功，新建不存在的同名目录，已设置的标签对应关系将在30分钟之内被继承。设置了标签的目录被删除后，新增一个同名目录，原有的对应关系也将在30分钟之内被继承。
 - 命令行设置方式请参考**hdfs nodelabel -setLabelExpression**命令。
 - Java API设置方式通过NodeLabelFileSystem实例化对象调用**setLabelExpression(String src, String labelExpression)**方法。*src*为HDFS上的目录或文件路径，“labelExpression”为标签表达式。
- 开启NodeLabel特性后，可以通过命令**hdfs nodelabel -listNodeLabels**查看每个DataNode的标签信息。

块副本位置选择

NodeLabel支持对各个副本的摆放采用不同的策略，如表达式

“label-1,label-2,label-3”，表示3个副本分别放到含有label-1、label-2、label-3的DataNode中，不同的副本策略用逗号分隔。

如果label-1，希望放2个副本，可以这样设置表达式：

“label-1[replica=2],label-2,label-3”。这种情况下，如果默认副本数是3，则会选择2个带有label-1和一个label-2的节点；如果默认副本数是4，会选择2个带有label-1、一个label-2以及一个label-3的节点。可以注意到，副本数是从左到右依次满足各个副本策略的，但也有副本数超过表达式表述的情况，当默认副本数为5时，多出来的一个副本会放到最后一个节点中，也就是label-3的节点里。

当启用ACLs功能并且用户无权访问表达式中使用的标签时，将不会为副本选择属于该标签的DataNode。

多余块副本删除选择

如果块副本数超过参数“dfs.replication”值（即用户指定的文件副本数，可以参考[修改集群服务配置参数](#)进入HDFS服务全部配置页面，搜索对应参数查看），HDFS会删除多余块副本来保证集群资源利用率。

删除规则如下：

- 优先删除不满足任何表达式的副本。

示例：文件默认副本数为3

/test标签表达式为“LA[replica=1],LB[replica=1],LC[replica=1]”；

/test文件副本分布的四个节点（D1~D4）以及对应标签（LA~LD）：

```
D1:LA
D2:LB
D3:LC
D4:LD
```

则选择删除D4节点上的副本块。

- 如果所有副本都满足表达式，删除多于表达式指定的数量的副本。

示例：文件默认副本数为3

/test标签表达式为“LA[replica=1],LB[replica=1],LC[replica=1]”；

/test文件副本分布的四个节点以及对应标签：

```
D1:LA
D2:LA
D3:LB
D4:LC
```

则选择删除D1或者D2上的副本块。

- 如果文件所有者或文件所有者的组不能访问某个标签，则优先删除映射到该标签的DataNode中的副本。

基于标签的数据块摆放策略样例

例如某MRS集群有六个DataNode：dn-1，dn-2，dn-3，dn-4，dn-5以及dn-6，对应的IP为10.1.120.[1-6]。有六个目录需要配置标签表达式，Block默认备份数为3。

- 下面给出3种DataNode标签信息在“host2labels”文件中的表示方式，其作用是一样的。

- 主机名正则表达式

```
/dn-[1456]/ = label-1,label-2
/dn-[26]/ = label-1,label-3
/dn-[3456]/ = label-1,label-4
/dn-5/ = label-5
```

- IP地址范围表示方式

```
10.1.120.[1-6] = label-1
10.1.120.1 = label-2
10.1.120.2 = label-3
10.1.120.[3-6] = label-4
10.1.120.[4-6] = label-2
10.1.120.5 = label-5
10.1.120.6 = label-3
```

- 普通的主机名表达式

```
/dn-1/ = label-1, label-2
/dn-2/ = label-1, label-3
/dn-3/ = label-1, label-4
/dn-4/ = label-1, label-2, label-4
/dn-5/ = label-1, label-2, label-4, label-5
/dn-6/ = label-1, label-2, label-3, label-4
```

- 目录的标签表达式设置结果如下：

```
/dir1 = label-1  
/dir2 = label-1 && label-3  
/dir3 = label-2 || label-4[replica=2]  
/dir4 = (label-2 || label-3) && label-4  
/dir5 = !label-1  
/sdir2.txt = label-1 && label-3[replica=3,fallback=NONE]  
/dir6 = label-4[replica=2],label-2
```

📖 说明

标签表达式设置方式请参考 `hdfs nodelabel -setLabelExpression` 命令。

文件的数据块存放结果如下：

- “/dir1” 目录下文件的数据块可存放在dn-1, dn-2, dn-3, dn-4, dn-5和dn-6六个节点中的任意一个。
- “/dir2” 目录下文件的数据块可存放在dn-2和dn-6节点上。Block默认备份数为3，表达式只匹配了两个DataNode节点，第三个副本会在集群上剩余的节点中选择一个DataNode节点存放。
- “/dir3” 目录下文件的数据块可存放在dn-1, dn-3, dn-4, dn-5和dn-6中的任意三个节点上。
- “/dir4” 目录下文件的数据块可存放在dn-4, dn-5和dn-6。
- “/dir5” 目录下文件的数据块没有匹配到任何一个DataNode，会从整个集群中任意选择三个节点存放（和默认选块策略行为一致）。
- “/sdir2.txt” 文件的数据块，两个副本存放在dn-2和dn-6节点上，虽然还缺失一个备份节点，但由于使用了 `fallback=NONE` 参数，所以只存放两个备份。
- “/dir6” 目录下文件的数据块在具备label-4的节点中选择2个节点(dn-3 -- dn-6)，然后在label-2中选择一个节点，如果用户指定“/dir6”下文件副本数大于3，则多出来的副本均在label-2。

使用限制

配置文件中，“key”、“value”是以“=”、“:”及空白字符作为分隔的。因此，“key”对应的主机名中间请勿包含以上字符，否则会被误认为分隔符。

9.10 配置 NameNode 内存参数

配置场景

在HDFS中，每个文件对象都需要在NameNode中注册相应的信息，并占用一定的存储空间。随着文件数的增加，当原有的内存空间无法存储相应的信息时，需要修改内存大小的设置。

配置描述

参数入口：

请参考[修改集群服务配置参数](#)，进入HDFS“全部配置”页面。

表 9-15 参数说明

| 配置参数 | 说明 | 默认值 |
|------------|---|--|
| GC_PROFILE | <p>NameNode所占内存主要由FsImage大小决定。
FsImage Size = 文件数 * 900 Bytes，根据计算结果可估算hdfs的NameNode应设内存大小。
该参数项的内存大小取值如下：</p> <ul style="list-style-type: none"> • high: 4G • medium: 2G • low: 256M • custom: 根据实际数据量大小在GC_OPTS中设置内存大小。 | custom |
| GC_OPTS | <p>JVM用于gc的参数。仅当GC_PROFILE设置为custom时该配置才会生效。需确保GC_OPT参数设置正确，否则进程启动会失败。</p> <p>须知
请谨慎修改该项。如果配置不当，将造成服务不可用。</p> | <p>-Xms2G -Xmx4G -
XX:NewSize=128M -
XX:MaxNewSize=256M -
XX:MetaspaceSize=128M -
XX:MaxMetaspaceSize=128M -
XX:+UseConcMarkSweepGC -
XX:+CMSParallelRemarkEnabled -
-
XX:CMSInitiatingOccupancyFraction=65 -XX:+PrintGCDetails -
Dsun.rmi.dgc.client.gcInterval=0x7FFFFFFF -
Dsun.rmi.dgc.server.gcInterval=0x7FFFFFFF -XX:-
OmitStackTraceInFastThrow -
XX:+PrintGCDateStamps -
XX:+UseGCLogFileRotation -
XX:NumberOfGCLogFiles=10 -
XX:GCLogFileSize=1M -
Djdk.tls.ephemeralDHKeySize=2048</p> |

9.11 设置 HBase 和 HDFS 的句柄数限制

操作场景

当打开一个HDFS文件时，句柄数限制导出，出现如下错误：

```
IOException (Too many open files)
```

此时可以参考该章节设置HBase和HDFS的句柄数。

设置 HBase 和 HDFS 的句柄数限制

联系集群管理员增加各用户的句柄数。该配置为操作系统的配置，并非HBase或者HDFS的配置。建议集群管理员根据HBase和HDFS的业务量及各操作系统用户的权限进行句柄数设置。如果某一个用户需对业务量很大的HDFS进行很频繁且很多的操作，则为此用户设置较大的句柄数，避免出现以上错误。

步骤1 使用root用户登录集群所有节点机器或者客户端机器的操作系统，并进入“/etc/security”目录。

步骤2 执行如下命令编辑“limits.conf”文件。

```
vi limits.conf
```

新增如下内容：

```
hdfs -    nofile 32768
hbase -   nofile 32768
```

其中“hdfs”和“hbase”表示业务中用到的操作系统用户名称。

说明

- 只有root用户有权限编辑“limits.conf”文件。
- 如果修改的配置不生效，请确认“/etc/security/limits.d”目录下是否有针对操作系统用户的其他nofile值。这样的值可能会覆盖“/etc/security/limits.conf”中配置的值。
- 如果用户需要对HBase进行操作，建议将该用户的句柄数设置为“10000”以上。如果用户需要对HDFS进行操作，建议根据业务量大小设置对应的句柄数，建议不要给太小的值。如果用户需要对HBase和HDFS操作，建议设置较大的值，例如“32768”。

步骤3 使用如下命令查看某一用户的句柄数限制。

```
su - user_name
```

```
ulimit -n
```

界面会返回此用户的句柄数限制值。如下所示：

```
8194
```

```
----结束
```

9.12 配置 HDFS 单目录文件数量

操作场景

通常一个集群上部署了多个服务，且大部分服务的存储都依赖于HDFS文件系统。当集群运行时，不同组件（例如Spark、Yarn）或客户端可能会向同一个HDFS目录不断写入文件。但HDFS系统支持的单目录文件数目是有上限的，因此用户需要提前做好规划，防止单个目录下的文件数目超过阈值，导致任务出错。

HDFS提供了“dfs.namenode.fs-limits.max-directory-items”参数设置单个目录下可以存储的文件数目。

操作步骤

步骤1 请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面。

步骤2 搜索配置项“dfs.namenode.fs-limits.max-directory-items”。

表 9-16 参数说明

| 参数名称 | 描述 | 默认值 |
|--|------------------------------------|---------|
| dfs.namenode.fs-limits.max-directory-items | 定义目录中包含的最大条目数。
取值范围：1 ~ 6400000 | 1048576 |

步骤3 设置单个HDFS目录下最大可容纳的文件数目。保存修改的配置。保存完成后请重新启动配置过期的服务或实例以使配置生效。

说明

用户尽量将数据做好存储规划，可以按时间、业务类型等分类，不要单个目录下直属的文件过多，建议使用默认值，单个目录下约100万条。

---结束

9.13 HDFS 企业级能力增强

9.13.1 配置 DataNode 节点容量不一致时的副本放置策略

操作场景

默认情况下，NameNode会随机选择DataNode节点写文件。当集群内某些数据节点的磁盘容量不一致（某些节点的磁盘总容量大，某些总容量小），会导致磁盘总容量小的节点先写满。通过修改集群默认的DataNode写数据时的磁盘选择策略为“节点磁盘可用空间块放置策略”，可提高将块数据写到磁盘可用空间较大节点的概率，解决因为数据节点磁盘容量不一致导致的节点使用率不均衡的情况。

对系统的影响

修改磁盘选择策略为“节点磁盘可用空间块放置策略（org.apache.hadoop.hdfs.server.blockmanagement.AvailableSpaceBlockPlacement Policy）”，经过测试验证，在该测试结果中，修改前后，HDFS写文件性能影响范围在3%以内。

说明

NameNode默认的副本存储策略为：

- 第一副本：存放客户端所在节点。
- 第二副本：远端机架的数据节点。
- 第三副本：存放客户端所在节点的不同机架的不同节点。

如还有更多副本，则随机选择其它DataNode。

“节点磁盘可用空间块放置策略”的副本选择机制为：

- 第一个副本：存放在客户端所在DataNode（和默认的存放策略一样）。
- 第二个副本：
 - 选择存储节点的时候，先挑选2个满足要求的数据节点。
 - 比较这2个节点磁盘空间使用比例，如果磁盘空间使用率的相差小于5%，随机存放到第一个节点。
 - 如果磁盘空间使用率相差超过5%，即有60%（由dfs.namenode.available-space-block-placement-policy.balanced-space-preference-fraction指定，默认值0.6）的概率写到磁盘空间使用率低的节点。
- 第三副本等其他后续副本的存储情况，也参考第二个副本的选择方式。

前提条件

集群里DataNode节点的磁盘总容量偏差不能超过100%。

操作步骤

步骤1 请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面。

步骤2 调整HDFS写数据时的依据的磁盘选择策略参数。搜索“dfs.block.replicator.classname”参数，并将参数的值改为“org.apache.hadoop.hdfs.server.blockmanagement.AvailableSpaceBlockPlacementPolicy”。

表 9-17 参数描述

| 参数 | 参数说明 |
|--------------------------------|--|
| dfs.block.replicator.classname | 选择副本放置的DataNode的策略。
默认值为
“org.apache.hadoop.hdfs.server.blockmanagement.AvailableSpaceBlockPlacementPolicy”。 |

步骤3 保存修改的配置。保存完成后请重新启动配置过期的服务或实例以使配置生效。

----结束

9.13.2 配置 DataNode 预留磁盘百分比

配置场景

当Yarn本地目录和DataNode目录配置在同一个磁盘时，具有较大容量的磁盘可以运行更多的任务，因此将有更多的中间数据存储在Yarn本地目录。

目前DataNode支持通过配置“dfs.datanode.du.reserved”来配置预留磁盘空间大小。配置较小的数值不能满足更大的磁盘要求。但对于更小的磁盘配置更大的数值将浪费大量的空间。

为了避免这种情况，添加一个新的参数“dfs.datanode.du.reserved.percentage”来配置预留磁盘空间占总磁盘空间大小的百分比，那样可以基于总的磁盘空间来预留磁盘百分比。

说明

- 如果用户同时配置“dfs.datanode.du.reserved.percentage”和“dfs.datanode.du.reserved”，则采用这两个参数较大的数值作为DataNode的预留空间大小。
- 建议基于磁盘空间设置“dfs.datanode.du.reserved”或者“dfs.datanode.du.reserved.percentage”。

配置描述

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-18 参数描述

| 参数 | 描述 | 默认值 |
|-------------------------------------|--|-----|
| dfs.datanode.du.reserved.percentage | DataNode预留空间占总磁盘空间大小的百分比。DataNode会永久预留由此百分比计算得出的磁盘空间大小。
整数值，取值范围是0~100。 | 10 |

9.13.3 配置 NameNode 黑名单功能

配置场景

说明

本章节适用于MRS 3.x及后续版本。

在现有的缺省DFSClient failover proxy provider中，一旦某进程中的一个NameNode发生故障，在同一进程中的所有HDFS client实例都会尝试再次连接NameNode，导致应用长时间等待超时。

当位于同一JVM进程中的客户端对无法访问的NameNode进行连接时，会对系统造成负担。为了避免这种负担，MRS集群搭载了NameNode blacklist功能。

在新的Blacklisting DFSClient failover provider中，故障的NameNode将被记录至一个列表中。DFSClient会利用这些信息，防止客户端再次连接这些NameNode。该功能被称为NameNode blacklisting。

例如，如下集群配置：

NameNode: nn1、nn2

dfs.client.failover.connection.retries: 20

单JVM中的进程：10个客户端

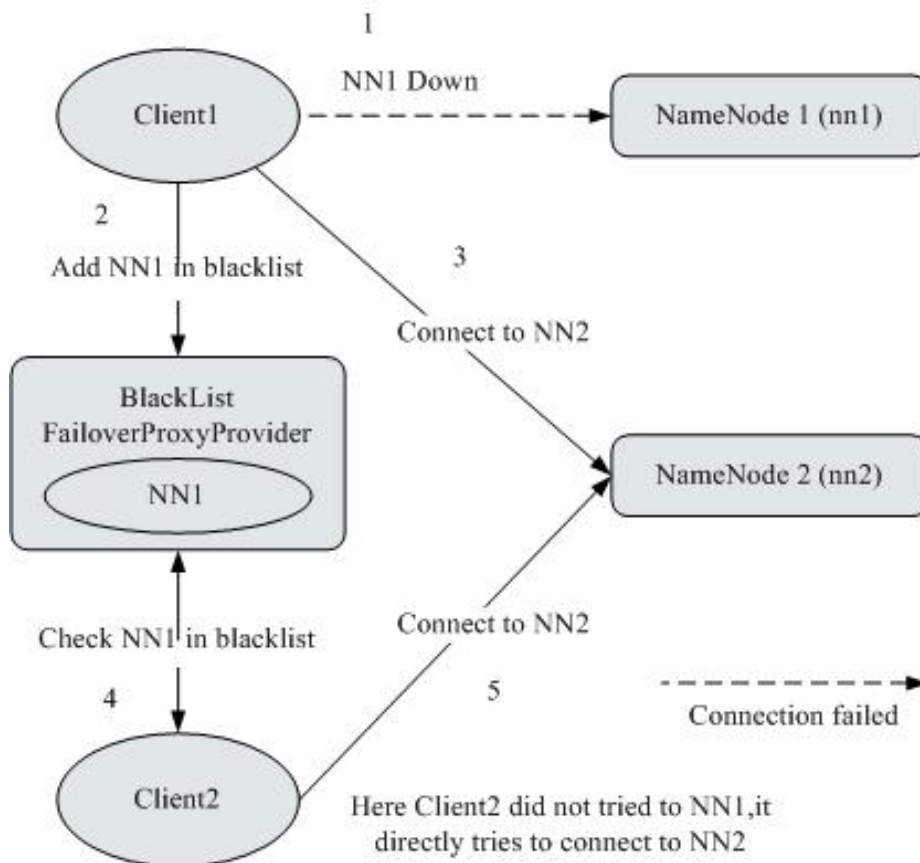
在上述集群中，如果当前处于active状态的nn1无法访问，client1将会对nn1进行20次重新连接，之后发生故障转移，client1将会连接至nn2。与此相同，client2至client10也会在对nn1进行20次重新连接后连接至nn2。这样会延长NameNode的整体故障恢复时间。

针对该情况，当client1试图连接当前处于active状态的nn1，但其已经发生故障时，nn1将会被添加至blacklist。这样其余client就不会连接已被添加至blacklist的nn1，而是会选择连接nn2。

说明

若在任一时刻，所有NameNode都被添加至blacklist，则其内容会被清空，client会按照初始的NameNode list重新尝试连接。若再次出现任何故障，NameNode仍会被添加至blacklist。

图 9-11 NameNode blacklisting 状态图



配置描述

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-19 NameNode blacklisting 的相关参数

| 参数 | 描述 | 默认值 |
|---|---|---|
| dfs.client.failover.proxy.provider.
[nameservice ID] | 利用已通过的协议创建namenode代理的Client Failover proxy provider类。
将参数值设置为
“org.apache.hadoop.hdfs.server.namenode.ha.BlackListingFailoverProxyProvider”，
可使用从NameNode支持读的特性。 | org.apache.hadoop.hdfs.server.namenode.ha.AdaptiveFailoverProxyProvider |

9.13.4 配置 Hadoop 数据传输加密

配置场景

安全加密通道是HDFS中RPC通信的一种加密协议，当用户调用RPC时，用户的login name会通过RPC头部传递给RPC，之后RPC使用Simple Authentication and Security Layer（SASL）确定一个权限协议（支持Kerberos和DIGEST-MD5两种），完成RPC授权。用户在部署安全集群时，需要使用安全加密通道，配置如下参数。安全Hadoop RPC相关信息请参考：

MRS 3.2.0之前版本：https://hadoop.apache.org/docs/r3.1.1/hadoop-project-dist/hadoop-common/SecureMode.html#Data_Encryption_on_RPC

MRS 3.2.0及之后版本：https://hadoop.apache.org/docs/r3.3.1/hadoop-project-dist/hadoop-common/SecureMode.html#Data_Encryption_on_RPC

配置描述

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-20 参数说明

| 参数 | 描述 | 默认值 |
|-----------------------|---|--|
| hadoop.rpc.protection | <p>须知</p> <ul style="list-style-type: none">• 设置后需要重启服务生效，且不支持滚动重启。• 设置后需要重新下载客户端配置，否则HDFS无法提供读写服务。 <p>设置Hadoop中各模块的RPC通道是否加密。通道包括：</p> <ul style="list-style-type: none">• 客户端访问HDFS的RPC通道。• HDFS中各模块间的RPC通道，如DataNode与NameNode间的RPC通道。• 客户端访问Yarn的RPC通道。• NodeManager和ResourceManager间的RPC通道。• Spark访问Yarn，Spark访问HDFS的RPC通道。• Mapreduce访问Yarn，Mapreduce访问HDFS的RPC通道。• HBase访问HDFS的RPC通道。 <p>说明</p> <p>用户可在HDFS组件的配置界面中设置该参数的值，设置后全局生效，即Hadoop中各模块的RPC通道的加密属性全部生效。</p> <p>对RPC的加密方式，有如下三种取值：</p> <ul style="list-style-type: none">• “authentication”：普通模式默认值，指数据在鉴权后直接传输，不加密。这种方式能保证性能但存在安全风险。• “integrity”：指数据直接传输，即不加密也不鉴权。为保证数据安全，请谨慎使用这种方式。• “privacy”：安全模式默认值，指数据在鉴权及加密后再传输。这种方式会降低性能。 | <ul style="list-style-type: none">• 安全模式：
privacy• 普通模式：
authentication |

9.14 HDFS 性能调优

9.14.1 提升 HDFS 写数据性能

操作场景

在HDFS中，通过调整属性的值，使得HDFS集群更适应自身的业务情况，从而提升HDFS的写性能。

说明

本章节适用于MRS 3.x及后续版本。

操作步骤

参数入口：

在FusionInsight Manager系统中，选择“集群 > 服务 > HDFS > 配置”，选择“全部配置”。在搜索框中输入参数名称。

表 9-21 HDFS 写性能优化配置

| 参数 | 描述 | 默认值 |
|--------------------------------------|--|--------|
| dfs.datanode.drop.cache.behind.reads | <p>表示是否让DataNode将在缓冲区中的数据传递给客户端后自动清除缓冲区中的所有数据。</p> <ul style="list-style-type: none"> • true：表示丢弃缓存的数据（需要在DataNode中配置）。当同一份数据，重复读取的次数较少时，建议设置为true，使得缓存能够被其他操作使用。 • false：重复读取的次数较多时，设置为false能够提升重复读取的速度。 <p>说明
在提升写性能操作中，该参数为可选参数，请根据实际需要进行修改。</p> | false |
| dfs.client-write-packet-size | <p>客户端写包的大小。当HDFS Client往DataNode写数据时，将数据生成一个包。然后将这个包在网络上传出。此参数指定传输数据包的大小，可以通过各Job来指定。单位：字节。</p> <p>在万兆网部署下，可适当增大该参数值，来提升传输的吞吐量。</p> | 262144 |

9.14.2 配置 HDFS 客户端元数据缓存提高读取性能

操作场景

通过使用客户端缓存元数据块的位置来提高HDFS读取性能。

📖 说明

此功能仅用于读取不经常修改的文件。因为在服务器端由某些其他客户端完成的数据修改，对于高速缓存的客户端将是不可见的，这可能导致从缓存中拿到的元数据是过期的。

本章节适用于MRS 3.x及后续版本。

操作步骤

设置参数的路径：

在FusionInsight Manager页面中，选择“集群 > 服务 > HDFS > 配置”，选择“全部配置”，并在搜索框中输入参数名称。

表 9-22 参数配置

| 参数 | 描述 | 默认值 |
|---------------------------------------|--|-------|
| dfs.client.metadata.cache.enabled | 启用/禁用块位置元数据的客户端缓存。将此参数设置为“true”，搭配“dfs.client.metadata.cache.pattern”参数以启用缓存。 | false |
| dfs.client.metadata.cache.pattern | 需要缓存的文件路径的正则表达式模式。只有这些文件的块位置元数据被缓存，直到这些元数据过期。此配置仅在参数“dfs.client.metadata.cache.enabled”设置为“true”时有效。
示例：“/test.*”表示读取其路径是以“/test”开头的所有文件。
说明 <ul style="list-style-type: none"> 为确保一致性，配置特定模式以仅缓存其他客户端不经常修改的文件。 正则表达式模式将仅验证URI的path部分，而不验证在Fully Qualified路径情况下的schema和authority。 | - |
| dfs.client.metadata.cache.expiry.sec | 缓存元数据的持续时间。缓存条目在该持续时间过期后失效。即使在缓存过程中经常使用的元数据也会发生失效。
配置值可采用时间后缀s/m/h表示，分别表示秒，分钟和小时。
说明
若将该参数配置为“0s”，将禁用缓存功能。 | 60s |
| dfs.client.metadata.cache.max.entries | 缓存一次最多可保存的非过期数据条目。 | 65536 |

说明

要在过期前完全清除客户端缓存，可调用`DFSClient#clearLocatedBlockCache()`。
用法如下所示。

```
FileSystem fs = FileSystem.get(conf);
DistributedFileSystem dfs = (DistributedFileSystem) fs;
DFSClient dfsClient = dfs.getClient();
dfsClient.clearLocatedBlockCache();
```

9.14.3 使用活动缓存提升 HDFS 客户端连接性能

操作场景

HDFS部署在具有多个NameNode实例的HA（High Availability）模式中，HDFS客户端需要依次连接到每个NameNode，以确定当前活动的NameNode是什么，并将其用于客户端操作。

一旦识别出来，当前活动的NameNode的详细信息就可以被缓存并共享给在客户端机器中运行的所有客户端。这样，每个新客户端可以首先尝试从缓存加载活动的Name

Node的详细信息，并将RPC调用保存到备用的NameNode。在异常情况下有很多优势，例如当备用的NameNode连接长时间不响应时。

当发生故障，将另一个NameNode切换为活动状态时，缓存的详细信息将被更新为当前活动的NameNode的信息。

📖 说明

本章节适用于MRS 3.x及后续版本。

操作步骤

设置参数的路径如下：

在FusionInsight Manager页面中，选择“集群 > 服务 > HDFS > 配置”，选择“全部配置”，并在搜索框中输入参数名称。

表 9-23 配置参数

| 参数 | 描述 | 默认值 |
|---|--|---|
| dfs.client.failover.proxy.provider.
[nameservice ID] | 用已通过的协议创建namenode代理的Client Failover proxy provider类。配置成org.apache.hadoop.hdfs.server.namenode.ha.BlackListingFailoverProxyProvider，可在HDFS客户端使用NameNode黑名单特性。配置成org.apache.hadoop.hdfs.server.namenode.ha.ObserverReadProxyProvider，可使用从NameNode支持读的特性。 | org.apache.hadoop.hdfs.server.namenode.ha.AdaptiveFailoverProxyProvider |
| dfs.client.failover.activeinfo.share.flag | 启用缓存并将当前活动的NameNode的详细信息共享给其他客户端。若要启用缓存，需将其设置为“true”。 | false |
| dfs.client.failover.activeinfo.share.path | 指定将在机器中的所有客户端创建的共享文件的本地目录。如果要为不同用户共享缓存，该文件夹应具有必需的权限（如在给定目录中创建，读写缓存文件）。 | /tmp |
| dfs.client.failover.activeinfo.share.io.timeout.sec | 控制超时的可选配置。用于在读取或写入缓存文件时获取锁定。如果在该时间内无法获取缓存文件上的锁定，则放弃尝试读取或更新缓存。单位为秒。 | 5 |

📖 说明

由HDFS客户端创建的缓存文件必须由其他客户端重新使用。因此，这些文件永远不会从本地系统中删除。若禁用该功能，可能需要进行手动清理。

9.14.4 HDFS 网络不稳定场景调优

配置场景

在网络不稳定的情况下，调整如下参数，降低客户端应用运行异常概率。

配置描述

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-24 参数说明

| 参数 | 描述 | 默认值 |
|--|--|--------|
| ha.health-monitor.rpc-timeout.ms | zkfc对NameNode健康状态检查的超时时间。增大该参数值，可以防止出现双Active NameNode，降低客户端应用运行异常的概率。
单位：毫秒。取值范围：30000~3600000 | 180000 |
| ipc.client.connect.max.retries.on.timeouts | 客户端与服务端建立Socket连接超时，客户端的重试次数。
取值范围：1~256 | 45 |
| ipc.client.connect.timeout | 客户端与服务端建立socket连接的超时时间。增大该参数值，可以增加建立连接的超时时间。
单位：毫秒。取值范围：1~3600000 | 20000 |

9.14.5 优化 HDFS NameNode RPC 的服务质量

配置场景

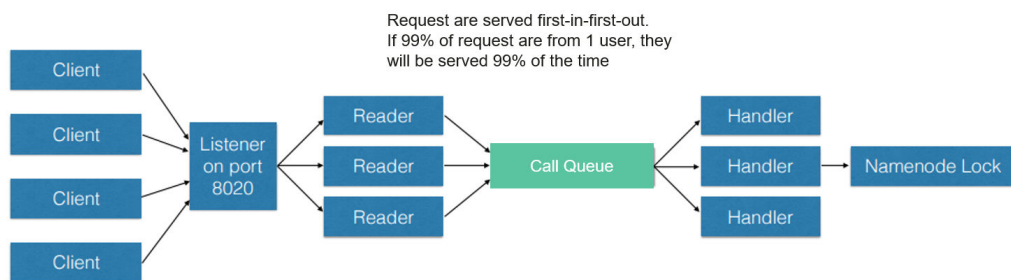
说明

本章节适用于MRS 3.x及后续版本。

数个成品Hadoop集群由于NameNode超负荷运行并失去响应而发生故障。

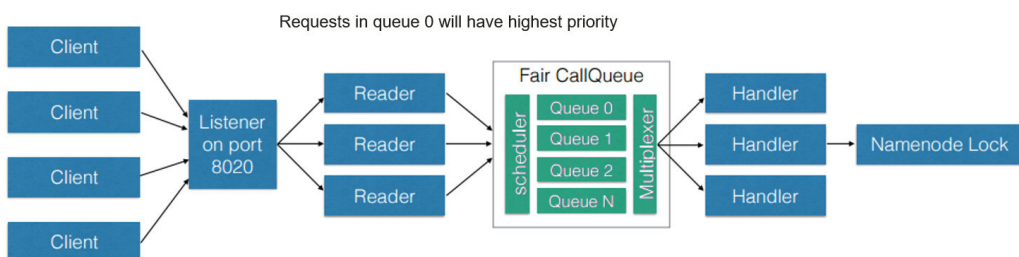
这种阻塞现象是由于Hadoop的初始设计造成的。在Hadoop中，NameNode作为单独的机器，在其namespace内协调HDFS的各种操作。这些操作包括获取数据块位置，列出目录及创建文件。NameNode接受HDFS的操作，将其视作RPC调用并置入FIFO调用队列，供读取线程处理。虽然FIFO在先到先服务的情况下足够公平，但如果用户执行的I/O操作较多，相比I/O操作较少的用户，将获得更多的服务。在这种情况下，FIFO有失公平并且会导致延迟增加。

图 9-12 基于 FIFO 调用队列的 NameNode 请求处理



如果将FIFO队列替换为一种被称作FairCallQueue的新型队列，这种情况就能够得到改善。按照这种方法，FAIR队列会根据调用者的调用规模将传入的RPC调用分配至多个队列中。调度模块会跟踪最新的调用，并为调用量较小的用户分配更高的优先级。

图 9-13 基于 FAIRCallQueue 的 NameNode 请求处理



配置描述

- FairCallQueue通过在内部调整RPC调用的顺序确保服务质量。该队列由以下三部分组成：
 - a. 调度模块（DecayRpcScheduler）用于提供从0至N的优先值数字（0的优先级最高）。
 - b. 多级队列（位于FairCallQueue内部）保持调用在内部按优先级排列。
 - c. 多路转换器（提供有WeightedRoundRobinMultiplexer）为队列选择提供逻辑控制。

在对FairCallQueue进行配置后，由控制模块决定将收到的调用分配至哪个子队列。当前调度模块为DecayRpcScheduler。该模块仅持续对各类调用的优先级数字进行追踪，并周期性地对这些数字进行减小处理。

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-25 Fair 调用队列参数

| 参数 | 描述 | 默认值 |
|---------------------------|--|--|
| ipc.<port>.callqueue.impl | 队列的实现类。用户需要通过“org.apache.hadoop.ipc.FairCallQueue”启用QoS特性。 | java.util.concurrent.LinkedBlockingQueue |

- **RPC BackOff**
Backoff是FairCallQueue的功能之一，要求客户端在一段时间后重试操作（如创建，删除，打开文件等）。当Backoff发生时，RCP服务器将发生RetriableException异常。FairCallQueue在以下两种情况时进行Backoff。
 - 当队列已满，即队列中有许多客户端调用时。
 - 当队列的响应时间大于配置的阈值（由参数“ipc.<port>.decay-scheduler.backoff.responsetime.thresholds”决定）时。

表 9-26 RPC BackOff 配置

| 参数 | 描述 | 默认值 |
|--|---|-------------------------|
| ipc.<port>.backoff.enable | 启用Backoff配置参数。当前，如果应用程序中包含较多的用户调用，假设没有达到操作系统的连接限制，则RPC请求将处于阻塞状态。或者，当RPC或NameNode在重负载时，可以基于某些策略将一些明确定义的异常抛回给客户端，客户端将理解这种异常并进行指数回退，以此作为类RetryInvocationHandler的另一个实现。 | false |
| ipc.<port>.decay-scheduler.backoff.responsetime.enable | 根据队列平均响应时间启用Backoff。 | false |
| ipc.<port>.decay-scheduler.backoff.responsetime.thresholds | 配置每个队列的响应时间阈值。ResponseTime阈值必须与优先级数目（ipc.<port>.faircallqueue.priority-levels）相匹配。单位：毫秒。 | 10000,20000,30000,40000 |

 说明

- <port>表示在NameNode上配置的RPC端口。
- 只有在“ipc.<port>.backoff.enable”为“true”时，响应时间backoff功能才会起作用。

9.14.6 优化 HDFS DataNode RPC 的服务质量

配置场景

当客户端写入HDFS的速度大于DataNode的硬盘带宽时，硬盘带宽会被占满，导致DataNode失去响应。客户端只能通过取消或恢复通道进行规避，这会导致写入失败及不必要的通道恢复操作。

说明

本章节适用于MRS 3.x及后续版本。

配置步骤

MRS引入配置参数“dfs.pipeline.ecn”。当该配置启用时，DataNode会在写入通道超出负荷时从其中发出信号。客户端可以基于该阻塞信号进行退避，从而防止系统超出负荷。引入该配置参数的目的是为了使通道更加稳定，并减少不必要的取消或恢复操作。收到信号后，客户端会退避一定的时间（5000ms），然后根据相关过滤器调整退避时间（单次退避最长时间为50000ms）。

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-27 NameNode ECN 配置

| 参数 | 描述 | 缺省值 |
|------------------|------------------------------|-------|
| dfs.pipeline.ecn | 进行该配置后，DataNode能够向客户端发送阻塞通知。 | false |

9.14.7 执行 HDFS 文件并发操作命令

操作场景

集群内并发修改文件和目录的权限及访问控制的工具。

说明

本章节适用于MRS 3.x及后续版本。

对系统的影响

因为集群内使用文件并发修改命令会对集群性能造成较大负担，所以在集群空闲时使用文件并发操作命令。

前提条件

- 已安装HDFS客户端或者包括HDFS的客户端。例如安装目录为“/opt/client”。
- 各组件业务用户由MRS集群管理员根据业务需要创建。安全模式下，“机机”用户需要下载keytab文件。“人机”用户第一次登录时需修改密码（普通模式不涉及）。

操作步骤

步骤1 以客户端安装用户，登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 如果集群为安全模式，执行的用户所属的用户组必须为**supergroup**组，且执行以下命令进行用户认证。普通模式集群无需执行用户认证。

kinit 组件业务用户

步骤5 增大客户端的JVM大小，防止OOM，方法如下。（1亿文件建议**32G**）

 **说明**

若执行HDFS客户端命令时，客户端程序异常退出，并且报“java.lang.OutOfMemoryError”错误。

这个问题是由于HDFS客户端运行时的所需的内存超过了HDFS客户端设置的内存上限（默认128M）。可通过修改“<客户端安装路径>/HDFS/component_env”中的“CLIENT_GC_OPTS”来修改HDFS客户端的内存上限。例如，需要设置内存上限为1GB，则设置：

```
CLIENT_GC_OPTS="-Xmx1G"
```

在修改完后，使用如下命令刷新客户端配置，使之生效：

```
source <客户端安装路径>/bigdata_env
```

步骤6 直接执行并发命令，命令详情如下表。

| 命令 | 参数及说明 | 命令作用 |
|---|---|-------------------|
| hdfs quickcmds [-t threadsNumber] [-p principal] [-k keytab] -setrep <rep> <path> ... | <ul style="list-style-type: none"> threadsNumber: 并发线程数，默认为本机CPU核数 principal: Kerberos用户 keytab: Keytab文件 rep: 副本数 path: HDFS目录 | 多并发设置目录中所有文件的副本数。 |
| hdfs quickcmds [-t threadsNumber] [-p principal] [-k keytab] -chown [owner][: [group]] <path> ... | <ul style="list-style-type: none"> threadsNumber: 并发线程数，默认为本机CPU核数 principal: Kerberos用户 keytab: Keytab文件 owner: 所属用户 group: 所属组 path: HDFS目录 | 多并发设置目录中所有文件的属组。 |
| hdfs quickcmds [-t threadsNumber] [-p principal] [-k keytab] -chmod <mode> <path> ... | <ul style="list-style-type: none"> threadsNumber: 并发线程数，默认为本机CPU核数 principal: Kerberos用户 keytab: Keytab文件 mode: 权限（如754） path: HDFS目录 | 多并发设置目录中所有文件的权限。 |

| 命令 | 参数及说明 | 命令作用 |
|--|---|---------------------|
| <code>hdfs quickcmds [-t threadsNumber] [-p principal] [-k keytab] -setfacl [{-b -k} {-m -x <acl_spec>} <path> ...] [--set <acl_spec> <path> ...]</code> | <ul style="list-style-type: none">• threadsNumber: 并发线程数, 默认为本机CPU核数• principal: Kerberos用户• keytab: Keytab文件• acl_spec: 逗号分隔的ACL列表• path: HDFS目录 | 多并发设置目录中所有文件的ACL信息。 |

----结束

9.14.8 使用 LZC 压缩算法存储 HDFS 文件

配置场景

文件压缩可以减少储存文件的空间，并且提高数据从磁盘读取和网络传输的速度。HDFS有Gzip和Snappy这两种默认压缩格式。本章节为HDFS新增加的压缩格式LZC（Lempel-Ziv Compression）提供配置方法。这种压缩格式增强了Hadoop压缩能力。有关Snappy的详细信息，请参阅<http://code.google.com/p/snappy/>。

说明

本章节适用于MRS 3.x及后续版本。

配置描述

为了使LZC压缩生效，需要在客户端的配置文件“core-site.xml”中（例如“客户端安装路径/HDFS/hadoop/etc/hadoop/”）配置如下参数。

表 9-28 参数描述

| 参数 | 描述 | 默认值 |
|--------------------------------|---|---|
| io.compression.codecs | 为了使LZC压缩格式生效，在现有的压缩格式列表中增加如下值：
“com.huawei.hadoop.datasight.io.compress.lzc.ZCodec”
说明
若配置了多于一种的压缩格式需要使用英文逗号分隔。 | org.apache.hadoop.io.compress.BZip2Codec,org.apache.hadoop.io.compress.DefaultCodec,org.apache.hadoop.io.compress.DeflateCodec,org.apache.hadoop.io.compress.Lz4Codec,org.apache.hadoop.io.compress.SnappyCodec,org.apache.hadoop.io.compress.GzipCodec,org.apache.hadoop.io.compress.ZStandardCodec,com.huawei.hadoop.datasight.io.compress.lzc.ZCodec |
| io.compression.codec.lzc.class | 为了使LZC压缩格式生效，使用该参数默认值，配置参数值为“com.huawei.hadoop.datasight.io.compress.lzc.ZCodec”。 | com.huawei.hadoop.datasight.io.compress.lzc.ZCodec |

📖 说明

- LZC压缩格式不支持FSImage和SequenceFile压缩。
- 当前HDFS提供了多种压缩算法，包括Gzip、LZ4、Snappy、Bzip2等。这几种压缩算法的压缩比和解压速度可参考如下：
压缩比排序：Bzip2>Gzip>LZ4>Snappy
解压速度排序：LZ4>Snappy>Gzip>Bzip2
- 使用场景建议：
 - 追求速度的场景（如Mapreduce任务中间数据的存储等）——建议使用LZ4和Snappy（高可靠场景，建议使用Snappy）。
 - 追求压缩比，而对压缩速度要求不高的场景（如冷数据的保存）——建议使用Bzip2或Gzip。
- 上述压缩算法除LZC外，皆支持Native（基于C语言实现）实现，压缩和解压缩效率较高。建议根据业务场景优先选用具备Native实现的压缩算法。

9.15 HDFS 运维管理

9.15.1 HDFS 常用配置参数

参数入口

请参考[修改集群服务配置参数](#)进入HDFS服务配置页面。

参数说明

表 9-29 HDFS 参数说明

| 参数 | 参数说明 | 默认值 |
|--------------------------|---|--|
| fs.obs.security.provider | <p>指定获取访问OBS文件系统的实现方式。</p> <p>参数取值：</p> <ul style="list-style-type: none"> com.huawei.mrs.MrsObsCredentialsProvider：通过MRS云服务委托获取凭证。 com.obs.services.EcsObsCredentialsProvider：通过ECS云服务获取AK/SK信息。 com.obs.services.BasicObsCredentialsProvider：使用用户传入OBS的AK/SK信息。 com.obs.services.EnvironmentVariableObsCredentialsProvider：从环境变量中读取AK/SK信息。 | com.huawei.mrs.MrsObsCredentialsProvider |

9.15.2 HDFS 日志介绍

日志描述

日志存储路径： HDFS相关日志的默认存储路径为“/var/log/Bigdata/hdfs/角色名”。

- NameNode： “/var/log/Bigdata/hdfs/nn”（运行日志）， “/var/log/Bigdata/audit/hdfs/nn”（审计日志）。
- DataNode： “/var/log/Bigdata/hdfs/dn”（运行日志）， “/var/log/Bigdata/audit/hdfs/dn”（审计日志）。
- ZKFC： “/var/log/Bigdata/hdfs/zkfc”（运行日志）， “/var/log/Bigdata/audit/hdfs/zkfc”（审计日志）。
- JournalNode： “/var/log/Bigdata/hdfs/jn”（运行日志）， “/var/log/Bigdata/audit/hdfs/jn”（审计日志）。
- Router： “/var/log/Bigdata/hdfs/router”（运行日志）， “/var/log/Bigdata/audit/hdfs/router”（审计日志）。
- HttpFS： “/var/log/Bigdata/hdfs/httpfs”（运行日志）， “/var/log/Bigdata/audit/hdfs/httpfs”（审计日志）。

日志归档规则：HDFS的日志启动了自动压缩归档功能，默认情况下，当日志大小超过100MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd_hh-mm-ss>.[编号].log.zip”。最多保留最近的100个压缩文件，压缩文件保留个数可以在Manager界面中配置。

表 9-30 HDFS 日志列表

| 日志类型 | 日志文件名 | 描述 |
|-----------------|---|---|
| 运行日志 | hadoop-<SSH_USER>-<process_name>-<hostname>.log | HDFS系统日志，记录HDFS系统运行时候所产生的大部分日志。 |
| | hadoop-<SSH_USER>-<process_name>-<hostname>.out | HDFS运行环境信息日志。 |
| | hadoop.log | Hadoop客户端操作日志。 |
| | hdfs-period-check.log | 周期运行的脚本的日志记录。包括：自动均衡、数据迁移、JournalNode数据同步检测等。 |
| | <process_name>-<SSH_USER>-<DATE>-<PID>-gc.log | 垃圾回收日志。 |
| | postinstallDetail.log | HDFS服务安装后启动前工作日志。 |
| | hdfs-service-check.log | HDFS服务启动是否成功的检查日志。 |
| | hdfs-set-storage-policy.log | HDFS数据存储策略日志。 |
| | cleanupDetail.log | HDFS服务卸载时候的清理日志。 |
| | prestartDetail.log | HDFS服务启动前集群操作的记录日志。 |
| | hdfs-recover-fsimage.log | NameNode元数据恢复日志。 |
| | datanode-disk-check.log | 集群安装过程和使用过程中磁盘状态检测的记录日志。 |
| | hdfs-availability-check.log | HDFS服务是否可用日志。 |
| | hdfs-backup-fsimage.log | NameNode元数据备份日志。 |
| startDetail.log | HDFS服务启动的详细日志。 | |

| 日志类型 | 日志文件名 | 描述 |
|----------|--|------------------------|
| | hdfs-blockplacement.log | HDFS块放置策略记录日志。 |
| | upgradeDetail.log | 升级日志。 |
| | hdfs-clean-acls-java.log | HDFS清除已删除角色的ACL信息的日志。 |
| | hdfs-haCheck.log | NameNode主备状态获取脚本运行日志。 |
| | <process_name>-jvmpause.log | 进程运行中，记录JVM停顿的日志。 |
| | hadoop-<SSH_USER>-balancer-<hostname>.log | HDFS自动均衡的运行日志。 |
| | hadoop-<SSH_USER>-balancer-<hostname>.out | HDFS运行自动均衡的环境信息日志。 |
| | hdfs-switch-namenode.log | HDFS主备倒换运行日志。 |
| | hdfs-router-admin.log | 管理挂载表操作的运行日志。 |
| | threadDump-<DATE>.log | 实例进程堆栈日志。 |
| Tomcat日志 | hadoop-omm-host1.out, httpfs-catalina.<DATE>.log, httpfs-host-manager.<DATE>.log, httpfs-localhost.<DATE>.log, httpfs-manager.<DATE>.log, localhost_access_web_log.log | Tomcat运行日志。 |
| 审计日志 | hdfs-audit-<process_name>.log
ranger-plugin-audit.log | HDFS操作审计日志（例如：文件增删改查）。 |
| | SecurityAuth.audit | HDFS安全审计日志。 |

日志级别

HDFS中提供了如表9-31所示的日志级别，日志级别优先级从高到低分别是FATAL、ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 9-31 日志级别

| 级别 | 描述 |
|-------|---------------------|
| FATAL | FATAL表示系统运行的致命错误信息。 |

| 级别 | 描述 |
|-------|-----------------------|
| ERROR | ERROR表示系统运行的错误信息。 |
| WARN | WARN表示当前事件处理存在异常信息。 |
| INFO | INFO表示系统及各事件正常运行状态信息。 |
| DEBUG | DEBUG表示系统及系统调试信息。 |

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面。
- 步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤3** 选择所需修改的日志级别。
- 步骤4** 保存配置，在弹出窗口中单击“确定”使配置生效。

说明

配置完成后立即生效，不需要重启服务。

----结束

日志格式

HDFS的日志格式如下所示：

表 9-32 日志格式

| 日志类型 | 格式 | 示例 |
|------|---|--|
| 运行日志 | <yyyy-MM-dd
HH:mm:ss,SSS> <Log
Level> <产生该日志的线
程名字> <log中的
message> <日志事件的发
生位置> | 20xx-01-26 18:43:42,840
 INFO IPC Server
handler 40 on 8020
Rolling edit logs
org.apache.hadoop.hdfs.s
erver.namenode.FSEditLo
g.rollEditLog(FSEditLog.j
ava:1096) |

| 日志类型 | 格式 | 示例 |
|------|---|--|
| 审计日志 | <yyyy-MM-dd
HH:mm:ss,SSS> <Log
Level> <产生该日志的线
程名字> <log中的
message> <日志事件的发
生位置> | 20xx-01-26 18:44:42,607
 INFO IPC Server
handler 32 on 8020
allowed=true ugi=hbase
(auth:SIMPLE) ip=/
10.177.112.145
cmd=getfileinfo src=/
hbase/WALs/
hghoulaslx410,16020,142
1743096083/
hghoulaslx410%2C16020
%2C1421743096083.142
2268722795 dst=null
perm=null
org.apache.hadoop.hdfs.s
erver.namenode.FSName
system
\$DefaultAuditLogger.log
AuditMessage(FSNamesy
stem.java:7950) |

9.15.3 查看 HDFS 容量状态

HDFS DataNode以Block的形式，保存用户的文件和目录，同时在NameNode中生成一个文件对象，对应DataNode中每个文件、目录和Block。

NameNode文件对象需要占用一定的内存，消耗内存大小随文件对象的生成而线性递增。DataNode实际保存的文件和目录越多，NameNode文件对象总量增加，需要消耗更多的内存，使集群现有硬件可能会难以满足业务需求，且导致集群难以扩展。

规划存储大量文件的HDFS系统容量，就是规划NameNode的容量规格和DataNode的容量规格，并根据容量设置参数。

容量规格

以下相关参数可以参考[修改集群服务配置参数](#)进入HDFS服务全部配置页面，搜索对应参数查看。

- NameNode容量规格

在NameNode中，每个文件对象对应DataNode中的一个文件、目录或Block。

一个文件至少占用一个Block，默认每个Block大小为“134217728”即128MB，对应参数为“dfs.blocksize”。默认情况下一个文件小于128MB时，只占用一个Block；文件大于128MB时，占用Block数为：文件大小/128MB。目录不占用Block。

根据“dfs.blocksize”，NameNode的文件对象数计算方法如下：

表 9-33 NameNode 文件对象数计算

| 单个文件大小 | 文件对象数 |
|-----------------|---|
| 小于128MB | 1（对应文件）+1（对应Block）=2 |
| 大于128MB（例如128G） | 1（对应文件）+1,024（对应128GB/128MB=1024 Block）=1,025 |

主备NameNode支持最大文件对象的数量为300,000,000（最多对应150,000,000个小文件）。“dfs.namenode.max.objects”规定当前系统可生成的文件对象数，默认值为“0”表示不限制。

- DataNode容量规格
在HDFS中，Block以副本的形式存储在DataNode中，默认副本数为“3”，对应参数为“dfs.replication”。
集群中所有DataNode角色实例保存的Block总数为：HDFS Block * 3。集群中每个DataNode实例平均保存的Blocks= HDFS Block * 3/DataNode节点数。

表 9-34 DataNode 支持规格

| 项目 | 规格 |
|--------------------------------|-----------|
| 单个DataNode实例支持最大Block数 | 5,000,000 |
| 单个DataNode实例上单个磁盘支持最大Block数 | 500,000 |
| 单个DataNode实例支持最大Block数需要的最小磁盘数 | 10 |

表 9-35 DataNode 节点数规划

| HDFS Block数 | 最少DataNode角色实例数 |
|-------------|------------------------------------|
| 10,000,000 | $10,000,000 * 3 / 5,000,000 = 6$ |
| 50,000,000 | $50,000,000 * 3 / 5,000,000 = 30$ |
| 100,000,000 | $100,000,000 * 3 / 5,000,000 = 60$ |

内存参数设置

以下相关参数可以参考[修改集群服务配置参数](#)进入HDFS服务全部配置页面，搜索对应参数查看。

- NameNode JVM参数配置规则
NameNode JVM参数“GC_OPTS”默认值为：
-Xms2G -Xmx4G -XX:NewSize=128M -XX:MaxNewSize=256M -
XX:MetaspaceSize=128M -XX:MaxMetaspaceSize=128M -

```
XX:+UseConcMarkSweepGC -XX:+CMSParallelRemarkEnabled -
XX:CMSInitiatingOccupancyFraction=65 -XX:+PrintGCDetails -
Dsun.rmi.dgc.client.gcInterval=0x7FFFFFFFFFFFFFFE -
Dsun.rmi.dgc.server.gcInterval=0x7FFFFFFFFFFFFFFE -XX:-
OmitStackTraceInFastThrow -XX:+PrintGCDateStamps -
XX:+UseGCLogFileRotation -XX:NumberOfGCLogFiles=10 -
XX:GCLogFileSize=1M -Djdk.tls.ephemeralDHKeySize=3072 -
Djdk.tls.rejectClientInitiatedRenegotiation=true -Djava.io.tmpdir=$
{Bigdata_tmp_dir}
```

NameNode文件数量和NameNode使用的内存大小成比例关系，文件对象变化时请修改默认值中的“-Xms2G -Xmx4G -XX:NewSize=128M -XX:MaxNewSize=256M”。参考值如下表所示。

表 9-36 NameNode JVM 配置

| 文件对象数量 | 参考值 |
|-------------|--|
| 10,000,000 | “-Xms6G -Xmx6G -XX:NewSize=512M -XX:MaxNewSize=512M” |
| 20,000,000 | “-Xms12G -Xmx12G -XX:NewSize=1G -XX:MaxNewSize=1G” |
| 50,000,000 | “-Xms32G -Xmx32G -XX:NewSize=3G -XX:MaxNewSize=3G” |
| 100,000,000 | “-Xms64G -Xmx64G -XX:NewSize=6G -XX:MaxNewSize=6G” |
| 200,000,000 | “-Xms96G -Xmx96G -XX:NewSize=9G -XX:MaxNewSize=9G” |
| 300,000,000 | “-Xms164G -Xmx164G -XX:NewSize=12G -XX:MaxNewSize=12G” |

- DataNode JVM参数配置规则

DataNode JVM参数“GC_OPTS”默认值为：

```
-Xms2G -Xmx4G -XX:NewSize=128M -XX:MaxNewSize=256M -
XX:MetaspaceSize=128M -XX:MaxMetaspaceSize=128M -
XX:+UseConcMarkSweepGC -XX:+CMSParallelRemarkEnabled -
XX:CMSInitiatingOccupancyFraction=65 -XX:+PrintGCDetails -
Dsun.rmi.dgc.client.gcInterval=0x7FFFFFFFFFFFFFFE -
Dsun.rmi.dgc.server.gcInterval=0x7FFFFFFFFFFFFFFE -XX:-
OmitStackTraceInFastThrow -XX:+PrintGCDateStamps -
XX:+UseGCLogFileRotation -XX:NumberOfGCLogFiles=10 -
XX:GCLogFileSize=1M -Djdk.tls.ephemeralDHKeySize=3072 -
Djdk.tls.rejectClientInitiatedRenegotiation=true -Djava.io.tmpdir=$
{Bigdata_tmp_dir}
```

集群中每个DataNode实例平均保存的Blocks= HDFS Block * 3/DataNode节点数，单个DataNode实例平均Block数量变化时请修改默认值中的“-Xms2G -Xmx4G -XX:NewSize=128M -XX:MaxNewSize=256M”。参考值如下表所示。

表 9-37 DataNode JVM 配置

| 单个DataNode实例平均Block数量 | 参考值 |
|-----------------------|--|
| 2,000,000 | "-Xms6G -Xmx6G -XX:NewSize=512M -XX:MaxNewSize=512M" |
| 5,000,000 | "-Xms12G -Xmx12G -XX:NewSize=1G -XX:MaxNewSize=1G" |

Xmx内存值对应DataNode节点块数阈值，每GB对应500000块数，用户可根据需要调整内存值。

查看 HDFS 容量状态

- NameNode信息

MRS 3.x之前版本：登录MRS控制台，选择“组件管理 > HDFS > NameNode(主)”，单击“Overview”，查看“Summary”显示的当前HDFS文件对象、文件数量、目录数量和Block数量信息。

Summary

| | |
|--|--------------------------------|
| Security is off. | |
| Safemode is off. | |
| 488 files and directories, 337 blocks = 825 total filesystem object(s). | |
| Heap Memory used 243.83 MB of 1.95 GB Heap Memory. Max Heap Memory is 3.95 GB. | |
| Non Heap Memory used 87.28 MB of 89 MB Committed Non Heap Memory. Max Non Heap Memory is 488 MB. | |
| Configured Capacity: | 250.68 GB |
| DFS Used: | 1.24 GB (0.5%) |
| Non DFS Used: | 0 B |
| DFS Remaining: | 233.14 GB (93%) |
| Block Pool Used: | 1.24 GB (0.5%) |
| DataNodes usages% (Min/Median/Max/stdDev): | 0.25% / 0.57% / 0.66% / 0.18% |
| Live Nodes | 3 (Decommissioned: 0) |
| Dead Nodes | 0 (Decommissioned: 0) |
| Decommissioning Nodes | 0 |
| Total Datanode Volume Failures | 0 (0 B) |
| Number of Under-Replicated Blocks | 0 |
| Number of Blocks Pending Deletion | 0 |
| Block Deletion Start Time | Tue Sep 28 16:19:45 +0800 2021 |
| Last Checkpoint Time | Wed Sep 29 09:19:55 +0800 2021 |

MRS 3.x及后续版本：登录FusionInsight Manager，选择“集群 > 服务 > HDFS > NameNode(主)”，单击“Overview”，查看“Summary”显示的当前HDFS文件对象、文件数量、目录数量和Block数量信息。

- DataNode信息

MRS 3.x之前版本：登录MRS控制台，选择“组件管理 > HDFS > NameNode(主)”，单击“Datanodes”，查看所有告警DataNode节点的Block数量信息。

| Node | Http Address | Last contact | Capacity | Blocks | Block pool used | Version |
|--|--------------|--------------|----------|--------|-------------------|-----------------|
| node-ana-corewAlX.34bF0d99-9d47-4e53-a911-d730af226491.com:9866 (39-9866) | | 0s | 83.56 GB | 229 | 568.74 MB (0.66%) | 2.8.3-mrs-1.9.0 |
| node-ana-corewAlG.34bF0d99-9d47-4e53-a911-d730af226491.com:9866 (64-9866) | | 0s | 83.56 GB | 216 | 217.68 MB (0.25%) | 2.8.3-mrs-1.9.0 |
| node-ana-corewAlm.34bF0d99-9d47-4e53-a911-d730af226491.com:9866 (128-9866) | | 1s | 83.56 GB | 219 | 488.4 MB (0.57%) | 2.8.3-mrs-1.9.0 |

Showing 1 to 3 of 3 entries

Previous 1 Next

MRS 3.x及后续版本：登录FusionInsight Manager，选择“集群 > 服务 > HDFS > NameNode(主)”，单击“DataNodes”，查看所有告警DataNode节点的Block数量信息。

- 告警信息
监控ID为14007、14008、14009的告警是否产生，根据业务需要修改告警阈值。

9.15.4 更改 DataNode 的存储目录

操作场景

说明

本章节适用于MRS 3.x及后续版本。

HDFS DataNode定义的存储目录不正确或HDFS的存储规划变化时，MRS集群管理员需要在FusionInsight Manager中修改DataNode的存储目录，以保证HDFS正常工作。适用于以下场景：

- 更改DataNode角色的存储目录，所有DataNode实例的存储目录将同步修改。
- 更改DataNode单个实例的存储目录，只对单个实例生效，其他节点DataNode实例存储目录不变。

对系统的影响

- 更改DataNode角色的存储目录需要停止并重新启动HDFS服务，集群未完全启动前无法提供服务。
- 更改DataNode单个实例的存储目录需要停止并重新启动实例，该节点DataNode实例未启动前无法提供服务。
- 服务参数配置如果使用旧的存储目录，需要更新为新目录。

前提条件

- 在各个数据节点准备并安装好新磁盘，并格式化磁盘。
- 规划好新的目录路径，用于保存旧目录中的数据。
- 已安装好HDFS客户端。
- 准备好业务用户hdfs。
- 更改DataNode单个实例的存储目录时，保持活动的DataNode实例数必须大于“dfs.replication”的值。

操作步骤

检查环境

步骤1 以root用户登录安装HDFS客户端的服务器，执行以下命令配置环境变量。

```
source HDFS客户端安装目录/bigdata_env
```

步骤2 如果集群为安全模式，执行以下命令认证用户身份。

```
kinit hdfs
```

步骤3 在HDFS客户端执行以下命令，检查HDFS根目录下全部目录和文件是否状态正常。

```
hdfs fsck /
```

检查fsck显示结果：

- 显示如下信息，表示无文件丢失或损坏，执行**步骤4**。
The filesystem under path '/' is HEALTHY
- 显示其他信息，表示有文件丢失或损坏，执行**步骤5**。

步骤4 登录FusionInsight Manager，选择“集群 > 服务”查看HDFS的状态“运行状态”是否为“良好”。



- 是，执行**步骤6**。
- 否，HDFS状态不健康，执行**步骤5**。

步骤5 修复HDFS异常的具体操作，任务结束。

步骤6 确定修改DataNode的存储目录场景。

- 更改DataNode角色的存储目录，执行**步骤7**。
- 更改DataNode单个实例的存储目录，执行**步骤12**。

更改DataNode角色的存储目录

步骤7 选择“集群 > 服务 > HDFS > 停止服务”，停止HDFS服务。

步骤8 以root用户登录到安装HDFS服务的各个数据节点中，执行如下操作：

1. 创建目标目录（data1,data2为集群原有目录）。
例如目标目录为“\${BIGDATA_DATA_HOME}/hadoop/data3/dn”，执行以下命令：

```
mkdir -p ${BIGDATA_DATA_HOME}/hadoop/data3/dn
```
2. 挂载目标目录到新磁盘。例如挂载“\${BIGDATA_DATA_HOME}/hadoop/data3”到新磁盘。

- 修改新目录的权限。
例如新目录路径为“`${BIGDATA_DATA_HOME}/hadoop/data3/dn`”，执行以下命令：

```
chmod 700 ${BIGDATA_DATA_HOME}/hadoop/data3/dn -R
chown omm:wheel ${BIGDATA_DATA_HOME}/hadoop/data3/dn -R
```

- 将数据复制到目标目录。
例如旧目录为“`${BIGDATA_DATA_HOME}/hadoop/data1/dn`”，目标目录为“`${BIGDATA_DATA_HOME}/hadoop/data3/dn`”，执行以下命令：

```
cp -af ${BIGDATA_DATA_HOME}/hadoop/data1/dn/* $
${BIGDATA_DATA_HOME}/hadoop/data3/dn
```

步骤9 在FusionInsight Manager管理界面，选择“集群 > 服务 > HDFS > 配置 > 全部配置”，打开HDFS服务配置页面。

将配置项“`dfs.datanode.data.dir`”从默认值“`%{@auto.detect.datapart.dn}`”修改为新的目标目录，例如“`${BIGDATA_DATA_HOME}/hadoop/data3/dn`”。

例如：原有的数据存储目录为“`/srv/BigData/hadoop/data1`”，“`/srv/BigData/hadoop/data2`”，如需将data1目录的数据迁移至新建的“`/srv/BigData/hadoop/data3`”目录，则将服务级别的此参数替换为现有的数据存储目录，如果有多个存储目录，用“`,`”隔开。则本示例中，为“`/srv/BigData/hadoop/data2,/srv/BigData/hadoop/data3`”。

步骤10 单击“保存”。然后在“集群 > 服务”界面启动集群中各个停止的服务。

步骤11 启动HDFS成功以后，在HDFS客户端执行以下命令，检查HDFS根目录下全部目录和文件是否复制正确。

```
hdfs fsck /
```

检查fsck显示结果：

- 显示如下信息，表示无文件丢失或损坏，数据复制成功，操作结束。
The filesystem under path '/' is HEALTHY
- 显示其他信息，表示有文件丢失或损坏，则检查8.4是否正确，并执行以下命令：
`hdfs fsck 损坏的文件名称 -delete`

更改DataNode单个实例的存储目录

步骤12 选择“集群 > 服务 > HDFS > 实例”，勾选需要修改存储目录的DataNode单个实例，选择“更多 > 停止实例”。



步骤13 以root用户登录到这个DataNode节点，执行如下操作。

1. 创建目标目录。

例如目标目录为“`${BIGDATA_DATA_HOME}/hadoop/data3/dn`”，执行以下命令：

```
mkdir -p ${BIGDATA_DATA_HOME}/hadoop/data3/dn
```

2. 挂载目标目录到新磁盘。

例如挂载“`${BIGDATA_DATA_HOME}/hadoop/data3`”到新磁盘。

3. 修改新目录的权限。

例如新目录路径为“`${BIGDATA_DATA_HOME}/hadoop/data3/dn`”，执行以下命令：

```
chmod 700 ${BIGDATA_DATA_HOME}/hadoop/data3/dn -R
```

```
chown omm:wheel ${BIGDATA_DATA_HOME}/hadoop/data3/dn -R
```

4. 将数据复制到目标目录。

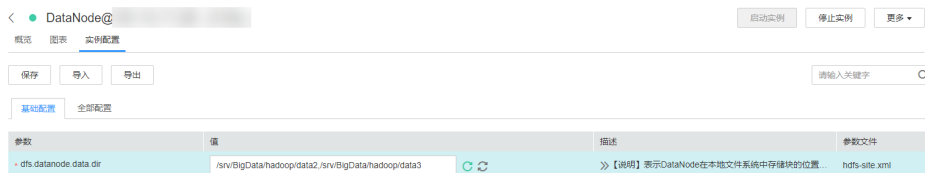
例如旧目录为“`${BIGDATA_DATA_HOME}/hadoop/data1/dn`”，目标目录为“`${BIGDATA_DATA_HOME}/hadoop/data3/dn`”，执行以下命令：

```
cp -af ${BIGDATA_DATA_HOME}/hadoop/data1/dn/* $  
{BIGDATA_DATA_HOME}/hadoop/data3/dn
```

步骤14 在FusionInsight Manager管理界面，选择“集群 > 服务 > HDFS > 实例”，单击指定的DataNode实例并切换到“实例配置”页签。

将配置项“`dfs.datanode.data.dir`”从默认值“`%{@auto.detect.datapart.dn}`”修改为新的目标目录，例如“`${BIGDATA_DATA_HOME}/hadoop/data3/dn`”。

示例：原有的数据存储目录为“`/srv/BigData/hadoop/data1,/srv/BigData/hadoop/data2`”，此处如需将data1目录的数据迁移至新建的`/srv/BigData/hadoop/data3`目录，则将该参数修改为“`/srv/BigData/hadoop/data2,/srv/BigData/hadoop/data3`”。



步骤15 单击“保存”，单击“确定”。

界面提示“操作成功。”，单击“完成”。

步骤16 选择“更多 > 重启实例”，重启DataNode实例。

---结束

9.15.5 调整 DataNode 磁盘坏卷信息

配置场景

在开源版本中，如果为DataNode配置多个数据存放卷，默认情况下其中一个卷损坏，则DataNode将不再提供服务。用户可以通过修改配置项

“`dfs.datanode.failed.volumes.tolerated`”的值，指定失败的个数，小于该个数，DataNode可以继续提供服务。

配置描述

参数入口:

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-38 参数说明

| 参数 | 描述 | 默认值 |
|---------------------------------------|--|----------------------------------|
| dfs.datanode.failed.volumes.tolerated | DataNode停止提供服务前允许失败的卷数。默认情况下，必须至少有一个有效卷。值-1表示有效卷的最小值是1。大于等于0的值表示允许失败的卷数。 | MRS 3.x之前版本：0
MRS 3.x及之后版本：-1 |

9.15.6 配置 HDFS token 的最大存活时间

配置场景

安全模式下，HDFS中用户可以对Token的最大存活时间和Token renew的时间间隔进行灵活地设置，根据集群的具体需求合理地配置。

配置描述

参数入口:

请参考[修改集群服务配置参数](#)，进入HDFS的“全部配置”页面，在搜索框中输入参数名称。

表 9-39 参数说明

| 参数 | 描述 | 默认值 |
|--|---|---------------|
| dfs.namenode.delegation.token.max-lifetime | 该参数为服务器端参数，设置Token的最大存活时间，单位为毫秒。取值范围：10000~100000000000000。 | 60480000
0 |
| dfs.namenode.delegation.token.renew-interval | 该参数为服务器端参数，设置Token renew的时间间隔，单位为毫秒。取值范围：10000~100000000000000。 | 86400000 |

9.15.7 使用 distcp 命令跨集群复制 HDFS 数据

操作场景

distcp是一种在集群间或集群内部拷贝大量数据的工具。它利用MapReduce任务实现大量数据的分布式拷贝。

前提条件

- 已安装Yarn客户端或者包括Yarn的客户端。例如安装目录为“/opt/client”。
- 各组件业务用户由MRS集群管理员根据业务需要创建。安全模式下，“机机”用户需要下载keytab文件。“人机”用户第一次登录时需修改密码。（普通模式不涉及）
- 如需在集群间拷贝数据，拷贝数据的集群双方都需要启用集群间拷贝数据功能。

操作步骤

步骤1 登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 如果集群为安全模式，执行distcp命令的用户所属的用户组必须为supergroup组，且执行以下命令进行用户认证。普通模式集群无需执行用户认证。

```
kinit 组件业务用户
```

步骤5 直接执行distcp命令。例如：

```
hadoop distcp hdfs://hacluster/source hdfs://hacluster/target
```

----结束

distcp 常见用法

1. 最常见的distcp用法，示例如下：

```
hadoop distcp -numListstatusThreads 40 -update -delete -prbugpaxtq hdfs://cluster1/source hdfs://cluster2/target
```

📖 说明

在上述命令中：

- -numListstatusThreads指定了40个构建被拷贝文件的列表的线程数；
- -update -delete表示将源位置和目标位置的文件同步，删除掉目标位置多余的文件，注意如果需要增量拷贝文件，请将-delete删掉；
- -prbugpaxtq与-update配合，表示被拷贝文件的状态信息也会被更新；
- hdfs://cluster1/source、hdfs://cluster2/target分别表示源位置和目标位置。

2. 集群间的数据拷贝，示例如下：

```
hadoop distcp hdfs://cluster1/foo/bar hdfs://cluster2/bar/foo
```

📖 说明

集群cluster1和集群cluster2之间的网络必须保持互通，且两个集群需要使用相同或兼容的HDFS版本。

3. 多个源目录的数据拷贝，示例如下：

```
hadoop distcp hdfs://cluster1/foo/a \  
hdfs://cluster1/foo/b \  
hdfs://cluster2/bar/foo
```

上面的命令的效果是将集群cluster1的文件夹a、b拷贝到集群cluster2的“/bar/foo”目录下，它的效果等效于下面的命令：

```
hadoop distcp -f hdfs://cluster1/srclist \  
hdfs://cluster2/bar/foo
```

其中srclist里面的内容如下。注意运行distcp命令前，需要将srclist文件上传到HDFS上。

```
hdfs://cluster1/foo/a  
hdfs://cluster1/foo/b
```

4. update和overwrite选项的用法。

- -update用于被拷贝的文件在目标位置中不存在，或者更新目标位置中被拷贝文件的内容；
- -overwrite用于覆盖在目标位置中已经存在的文件。

不加选项和加两个选项中任一个选项的区别，示例如下：

假设，源位置的文件结构如下：

```
hdfs://cluster1/source/first/1  
hdfs://cluster1/source/first/2  
hdfs://cluster1/source/second/10  
hdfs://cluster1/source/second/20
```

不加选项的命令：

```
hadoop distcp hdfs://cluster1/source/first hdfs://cluster1/source/second  
hdfs://cluster2/target
```

上述命令默认会在目标位置创建文件夹first、second，所以拷贝结果如下：

```
hdfs://cluster2/target/first/1  
hdfs://cluster2/target/first/2  
hdfs://cluster2/target/second/10  
hdfs://cluster2/target/second/20
```

加两个选项中任一个选项的命令，例如加update选项：

```
hadoop distcp -update hdfs://cluster1/source/first hdfs://cluster1/source/  
second hdfs://cluster2/target
```

上述命令只会将源位置的内容拷贝到目标位置，所以拷贝结果如下：

```
hdfs://cluster2/target/1  
hdfs://cluster2/target/2  
hdfs://cluster2/target/10  
hdfs://cluster2/target/20
```

说明

- 如果多个源位置有相同名称的文件，则distcp命令会失败。
- 在不使用update和overwrite选项的情况下，如果被拷贝文件在目标位置中已经存在，则该文件会跳过。
- 在使用update选项的情况下，如果被拷贝文件在目标位置中已经存在，但文件内容不同，则目标位置的文件内容会被更新。
- 在使用overwrite选项的情况下，如果被拷贝文件在目标位置中已经存在，目标位置的文件依然会被覆盖。

5. 其它命令选项：

表 9-40 其他命令选项

| 选项 | 描述 |
|-----------------------------------|---|
| -p[rbugpcaxtq] | 当同时使用-update选项时，即使被拷贝文件的内容没有被更新，它的状态信息也会被更新。
r: 副本数, b: 块大小, u: 所属用户, g: 所属用户组, p: 许可, c: 校验和类型, a: 访问控制, t: 时间戳, q: Quota信息 |
| -i | 拷贝过程中忽略失败。 |
| -log <logdir> | 指定日志路径。 |
| -v | 指定日志中的额外信息。 |
| -m <num_maps> | 最大的同时运行的执行拷贝的任务数。 |
| -numListstatusThreads | 构建被拷贝文件的文件列表时所用的线程数，该选项会提高distcp的运行速度。 |
| -overwrite | 覆盖目标位置的文件。 |
| -update | 如果源位置和目标位置的文件的的大小，校验和不同，则更新目标位置的文件。 |
| -append | 当同时使用-update选项时，追加源位置的文件内容到目标位置的文件。 |
| -f <urilist_uri> | 将<urilist_uri>文件的内容作为需要拷贝的文件列表。 |
| -filters | 指定一个本地文件，其文件内容是多条正则表达式。当被拷贝的文件与某条正则表达式匹配时，则该文件不会被拷贝。 |
| -async | 异步运行distcp命令。 |
| -atomic {-tmp <tmp_dir>} | 指定一次原子性的拷贝，可以添加一个临时目录的选项，作为拷贝过程中的暂存目录。 |
| -bandwidth | 指定每个拷贝任务的传输带宽，单位MB/s。 |
| -delete | 删除掉目标位置中存在，但源位置不存在的文件。该选项通常会和-update配合使用，表示将源位置和目标位置的文件同步，删除掉目标位置多余的文件。 |
| -diff <oldSnapshot> <newSnapshot> | 将新旧版本之间的差异内容，拷贝到目标位置的旧版本文件中。 |
| -skipcrccheck | 是否跳过源文件和目标文件之间的CRC校验。 |
| -strategy {dynamic uniformsize} | 指定拷贝任务的拷贝策略，默认策略是uniformsize，即每个拷贝任务复制相同的字节数。 |

distcp 常见使用问题

- 问题一：当使用distcp命令时，如果某些被拷贝的文件内容较大时，建议修改执行拷贝任务的mapreduce的超时时间。可以通过在distcp命令中指定 **mapreduce.task.timeout**选项实现。例如，修改超时时间为30分钟，则命令如下：

```
hadoop distcp -Dmapreduce.task.timeout=1800000 hdfs://cluster1/source hdfs://cluster2/target
```

您也可以使用选项filters，不对这种大文件进行拷贝，命令示例如下：

```
hadoop distcp -filters /opt/client/filterfile hdfs://cluster1/source hdfs://cluster2/target
```

其中filterfile是本地文件，它的内容是多条用于匹配不拷贝文件路径的正则表达式，它的内容示例如下：

```
.*excludeFile1.*  
.*excludeFile2.*
```

- 问题二：当使用distcp命令时，命令异常退出，报“java.lang.OutOfMemoryError”的错误。

这个问题的原因是拷贝任务运行时所需的内存超过了客户端设置的内存上限（默认为128MB）。可以通过修改“<客户端安装路径>/HDFS/component_env”中的“CLIENT_GC_OPTS”来修改客户端的内存上限。例如，需要设置该内存上限为1GB，则设置：

```
CLIENT_GC_OPTS="-Xmx1G"
```

在修改完后，使用如下命令刷新客户端配置，使之生效：

```
source <客户端安装路径>/bigdata_env
```

- 问题三：使用dynamic策略执行distcp命令时，命令异常退出，报“Too many chunks created with splitRatio”的错误。

这个问题的原因是“distcp.dynamic.max.chunks.tolerable”的值（默认值为20000）小于“distcp.dynamic.split.ratio”的值（默认为2）乘以Map数。即一般出现在Map数超过10000的情况。可以通过-m参数降低Map数小于10000：

```
hadoop distcp -strategy dynamic -m 9500 hdfs://cluster1/source hdfs://cluster2/target
```

或通过-D参数指定更大的“distcp.dynamic.max.chunks.tolerable”的值：

```
hadoop distcp -Ddistcp.dynamic.max.chunks.tolerable=30000 -strategy dynamic hdfs://cluster1/source hdfs://cluster2/target
```

9.15.8 配置 NFS 服务器存储 NameNode 元数据

操作场景

📖 说明

本章节适用于MRS 3.x及后续版本。

用户在部署集群前，可根据需要规划Network File System（简称NFS）服务器，用于存储NameNode元数据，以提高数据可靠性。

如果您已经部署NFS服务器，并已配置NFS服务，本操作提供集群侧的配置指导，为可选任务。

操作步骤

- 步骤1** 在NFS服务器上检查NFS的共享目录权限，确认服务器可以访问MRS集群的NameNode。

步骤2 以root用户登录NameNode主节点。

步骤3 执行如下命令，创建目录并赋予目录写权限。

```
mkdir ${BIGDATA_DATA_HOME}/namenode-nfs
chown omm:wheel ${BIGDATA_DATA_HOME}/namenode-nfs
chmod 750 ${BIGDATA_DATA_HOME}/namenode-nfs
```

步骤4 执行如下命令，挂载NFS到NameNode主节点。

```
mount -t nfs -o rsize=8192,wsiz=8192,soft,nolock,timeo=3,intr NFS服务器IP地址:共享目录 ${BIGDATA_DATA_HOME}/namenode-nfs
```

例如，NFS服务器的IP为“192.168.0.11”，共享目录为“/opt/Hadoop/NameNode”，则执行命令：

```
mount -t nfs -o rsize=8192,wsiz=8192,soft,nolock,timeo=3,intr
192.168.0.11:/opt/Hadoop/NameNode ${BIGDATA_DATA_HOME}/namenode-
nfs
```

步骤5 在NameNode备节点上执行**步骤2**~**步骤4**。

📖 说明

主备NameNode节点在NFS服务器上创建的共享目录名称（如“/opt/Hadoop/NameNode”）不能相同。

步骤6 登录FusionInsight Manager系统，选择“集群 > 服务 > HDFS > 配置 > 全部配置”。

步骤7 在界面右侧的“搜索”框中输入“dfs.namenode.name.dir”搜索，在其值中增加“\${BIGDATA_DATA_HOME}/namenode-nfs”路径，多个路径间使用“,”隔开，然后单击“保存”。

步骤8 单击“确定”。在概览页面选择“更多 > 重启服务”，重启服务。

----结束

9.16 HDFS 常见问题

9.16.1 执行 distcp 命令报错如何处理

问题

为何distcp命令在安全集群上执行失败并发生异常？

客户端出现异常：

```
Invalid arguments:Unexpected end of file from server
```

服务器端出现异常：

```
javax.net.ssl.SSLException:Unrecognized SSL message, plaintext connection?
```

回答

当用户在distcp命令中使用webhdfs://时，会发生上述异常，是由于集群所使用的HTTP政策为HTTPS，即配置在“hdfs-site.xml”（文件路径为“客户端安装目录/

HDFS/hadoop/etc/hadoop”) 的 “dfs.http.policy” 值为 “HTTPS_ONLY” 。所以要避免出现此异常，应使用 `swebhdfs://` 替代 `webhdfs://`。

例如：

```
./hadoop distcpswebhdfs://IP:PORT/testfile hdfs://IP:PORT/testfile1
```

9.16.2 HDFS 执行 Balance 时被异常停止如何处理

问题

在HDFS客户端启动一个Balance进程，该进程被异常停止后，再次执行Balance操作，操作会失败。

回答

通常，HDFS执行Balance操作结束后，会自动释放 “/system/balancer.id” 文件，可再次正常执行Balance。

但在上述场景中，由于第一次的Balance操作是被异常停止的，所以第二次进行Balance操作时，“/system/balancer.id” 文件仍然存在，则会触发 `append /system/balancer.id` 操作，进而导致Balance操作失败。

- 如果 “/system/balancer.id” 文件的释放时间超过了软租期60s，则第二次执行Balance操作的客户端的append操作会抢占租约，此时最后一个block处于under construction或者under recovery状态，会触发block的恢复操作，那么 “/system/balancer.id” 文件必须等待block恢复完成才能关闭，所以此次append操作失败。

`append /system/balancer.id`操作失败后，客户端发生 `RecoveryInProgressException` 异常：

```
org.apache.hadoop.ipc.RemoteException(org.apache.hadoop.protocol.RecoveryInProgressException):  
Failed to APPEND_FILE /system/balancer.id for DFSClient because lease recovery is in progress. Try  
again later.
```

- 如果该文件的释放时间没有超过默认设置60s，原有客户端会继续持有该租约，则会发生 `AlreadyBeingCreatedException` 异常，实际上向客户端返回的是null，导致客户端出现如下异常：
`java.io.IOException: Cannot create any NameNode Connectors.. Exiting...`

可通过以下方法避免上述问题：

- 方案1：等待硬租期超过1小时后，原有客户端释放租约，再执行第二次Balance操作。
- 方案2：执行第二次Balance操作之前删除 “/system/balancer.id” 文件。

9.16.3 访问 HDFS WebUI 时，界面提示无法显示此页

问题

通过IE 9、IE 10和IE 11等浏览器访问HDFS的原生UI界面，偶尔出现访问失败情况。

现象

访问页面失败，浏览器无法显示此页，如下图所示：

无法显示此页

在高级设置中启用 SSL 3.0、TLS 1.0、TLS 1.1 和 TLS 1.2，然后尝试再次连接

原因

IE 9、IE 10、IE 11浏览器的某些版本在处理SSL握手有问题导致访问失败。

解决方法

重新刷新页面即可。

9.16.4 HDFS WebUI 无法正常刷新损坏数据的信息

问题

1. 当DataNode的“dfs.datanode.data.dir”所配置的目录因权限或者磁盘损坏发生错误时，HDFS Web UI没有显示损坏数据的信息。
2. 当此错误被修复后，HDFS Web UI没有及时移除损坏数据的相关信息。

回答

1. DataNode只有在执行文件操作发生错误时，才会去检查磁盘是否正常，若发现数据损坏，则将此错误上报至NameNode，此时NameNode才会在HDFS Web UI显示数据损坏信息。
2. 当错误修复后，需要重启DataNode。当重启DataNode时，会检查所有数据状态并上传损坏数据信息至NameNode。所以当此错误被修复后，只有重启DataNode后，才会不显示损坏数据信息。

9.16.5 NameNode 节点长时间满负载导致客户端无响应

问题

当NameNode节点处于满负载、NameNode所在节点的CPU 100%耗尽时，导致NameNode无法响应，对于新连接到该NameNode的HDFS客户端，能够主备切换连接到另一个NameNode，进行正常的操作，而对于已经连接到该NameNode节点的HDFS客户端可能会长时间无响应，无法进行下一步操作。

回答

目前出现上述问题时使用的是默认配置，如表9-41所示，HDFS客户端到NameNode的RPC连接存在keep alive机制，保持连接不会超时，尽力等待服务器的响应，因此导致已经连接的HDFS客户端的操作会长时间无响应。

对于已经长时间无响应的HDFS客户端，可以进行如下操作：

- 等待NameNode响应，一旦NameNode所在节点的CPU利用率回落，NameNode可以重新获得CPU资源时，HDFS客户端即可得到响应。
- 如果无法等待更长时间，需要重启HDFS客户端所在的应用程序进程，使得HDFS客户端重新连接空闲的NameNode。

解决措施：

为了避免该问题出现，可以在“客户端安装路径/HDFS/hadoop/etc/hadoop/core-site.xml”中做如下配置。

表 9-41 参数说明

| 参数 | 描述 | 默认值 |
|-------------------|---|-------|
| ipc.client.ping | 当配置为true时，客户端会尽力等待服务端响应，定期发送ping消息，使得连接不会因为tcp timeout而断开。
当配置为false时，客户端会使用配置项“ipc.ping.interval”对应的值，作为timeout时间，在该时间内没有得到响应，即会超时。
在上述问题场景下，建议配置为false。 | true |
| ipc.ping.interval | 当“ipc.client.ping”配置为true时，表示发送ping消息的周期。
当“ipc.client.ping”设置为false时，表示连接的超时时间。
在上述问题场景下，建议配置一个较大的超时时间，避免服务繁忙时的超时，建议配置为900000，单位为ms。 | 60000 |

9.16.6 为什么主 NameNode 重启后系统出现双备现象

问题

为什么主NameNode重启后系统出现双备现象？

出现该问题时，查看ZooKeeper和ZKFC的日志，发现ZooKeeper服务端与客户端（ZKFC）通信时所使用的session不一致，ZooKeeper服务端的sessionId为0x164cb2b3e4b36ae4，ZKFC的sessionId为0x144cb2b3e4b36ae4。这意味着ZooKeeper服务端与客户端（ZKFC）之间数据交互失败。

ZooKeeper日志，如下所示：

```
2015-04-15 21:24:54,257 | INFO | CommitProcessor:22 | Established session 0x164cb2b3e4b36ae4 with negotiated timeout 45000 for client /192.168.0.117:44586 | org.apache.zookeeper.server.ZooKeeperServer.finishSessionInit(ZooKeeperServer.java:623)
2015-04-15 21:24:54,261 | INFO | NIOServerCxn.Factory:192-168-0-114/192.168.0.114:2181 | Successfully authenticated client: authenticationID=hdfs/hadoop@<系统域名>; authorizationID=hdfs/hadoop@<系统域名>. | org.apache.zookeeper.server.auth.SaslServerCallbackHandler.handleAuthorizeCallback(SaslServerCallbackHandler.java:118)
2015-04-15 21:24:54,261 | INFO | NIOServerCxn.Factory:192-168-0-114/192.168.0.114:2181 | Setting authorizedID: hdfs/hadoop@<系统域名> | org.apache.zookeeper.server.auth.SaslServerCallbackHandler.handleAuthorizeCallback(SaslServerCallbackHandler.java:134)
2015-04-15 21:24:54,261 | INFO | NIOServerCxn.Factory:192-168-0-114/192.168.0.114:2181 | adding SASL authorization for authorizationID: hdfs/hadoop@<系统域名> | org.apache.zookeeper.server.ZooKeeperServer.processSasl(ZooKeeperServer.java:1009)
2015-04-15 21:24:54,262 | INFO | ProcessThread(sid:22 cport:-1): | Got user-level KeeperException when processing sessionId:0x164cb2b3e4b36ae4 type:create cxid:0x3 zxid:0x20009fafc txntype:-1 reqpath:n/a Error Path:/hadoop-ha/hacluster/ActiveStandbyElectorLock Error:KeeperErrorCode = NodeExists for /hadoop-
```

```
ha/hacluster/ActiveStandbyElectorLock |  
org.apache.zookeeper.server.PrepareRequestProcessor.pRequest(PrepareRequestProcessor.java:648)
```

ZKFC日志，如下所示：

```
2015-04-15 21:24:54,237 | INFO | main-SendThread(192-168-0-114:2181) | Socket connection established to  
192-168-0-114/192.168.0.114:2181, initiating session | org.apache.zookeeper.ClientCnxn  
$SendThread.primeConnection(ClientCnxn.java:854)  
2015-04-15 21:24:54,257 | INFO | main-SendThread(192-168-0-114:2181) | Session establishment complete  
on server 192-168-0-114/192.168.0.114:2181, sessionid = 0x144cb2b3e4b36ae4, negotiated timeout =  
45000 | org.apache.zookeeper.ClientCnxn$SendThread.onConnected(ClientCnxn.java:1259)  
2015-04-15 21:24:54,260 | INFO | main-EventThread | EventThread shut down |  
org.apache.zookeeper.ClientCnxn$EventThread.run(ClientCnxn.java:512)  
2015-04-15 21:24:54,262 | INFO | main-EventThread | Session connected. |  
org.apache.hadoop.ha.ActiveStandbyElector.processWatchEvent(ActiveStandbyElector.java:547)  
2015-04-15 21:24:54,264 | INFO | main-EventThread | Successfully authenticated to ZooKeeper using SASL. |  
org.apache.hadoop.ha.ActiveStandbyElector.processWatchEvent(ActiveStandbyElector.java:573)
```

回答

- 原因分析

NameNode的主节点重启后，之前在ZooKeeper上建立的临时节点（/hadoop-ha/hacluster/ActiveStandbyElectorLock）就会被清理。同时，NameNode备节点发现该信息后进行抢占希望升主，所以它重新在ZooKeeper上建立了active的节点/hadoop-ha/hacluster/ActiveStandbyElectorLock。但是NameNode备节点通过客户端（ZKFC）与ZooKeeper建立连接时，由于网络问题、CPU使用率高、集群压力大等原因，出现了客户端（ZKFC）的session（0x144cb2b3e4b36ae4）与ZooKeeper服务端的session（0x164cb2b3e4b36ae4）不一致的问题，导致NameNode备节点的watcher没有感知到自己已经成功建立临时节点，依然认为自己还是备。而NameNode主节点启动后，发现/hadoop-ha/hacluster目录下已经有active的节点，所以也无法升主，导致两个节点都为备。

- 解决方法

建议通过在FusionInsight Manager界面上重启HDFS的两个ZKFC加以解决。

9.16.7 为什么 DataNode 无法正常上报数据块

问题

DataNode正常，但无法正常上报数据块，导致存在的数据块无法使用。

回答

当某个数据目录中的数据块数量超过4倍的数据块限定值“1M”时，可能会出现该错误。DataNode会产生相应的错误日志记录，如下所示：

```
2015-11-05 10:26:32,936 | ERROR | DataNode: [[[DISK]file:/srv/BigData/hadoop/data1/dn/]] heartbeating to  
vm-210/10.91.8.210:8020 | Exception in BPOfferService for Block pool  
BP-805114975-10.91.8.210-1446519981645  
(Datanode Uuid bcada350-0231-413b-bac0-8c65e906c1bb) service to vm-210/10.91.8.210:8020 |  
BPServiceActor.java:824  
java.lang.IllegalStateException: com.google.protobuf.InvalidProtocolBufferException: Protocol message was  
too large. May  
be malicious. Use CodedInputStream.setSizeLimit() to increase the size limit. at  
org.apache.hadoop.hdfs.protocol.BlockListAsLongs$BufferDecoder$1.next(BlockListAsLongs.java:369)  
at org.apache.hadoop.hdfs.protocol.BlockListAsLongs$BufferDecoder$1.next(BlockListAsLongs.java:347) at  
org.apache.hadoop.hdfs.  
protocol.BlockListAsLongs$BufferDecoder.getBlockListAsLongs(BlockListAsLongs.java:325) at  
org.apache.hadoop.hdfs.protocolPB.DatanodeProtocolClientSideTranslatorPB.  
blockReport(DatanodeProtocolClientSideTranslatorPB.java:190) at  
org.apache.hadoop.hdfs.server.datanode.BPServiceActor.blockReport(BPServiceActor.java:473)
```

```
at org.apache.hadoop.hdfs.server.datanode.BPServiceActor.offerService(BPServiceActor.java:685) at
org.apache.hadoop.hdfs.server.datanode.BPServiceActor.run(BPServiceActor.java:822)
at java.lang.Thread.run(Thread.java:745) Caused
by:com.google.protobuf.InvalidProtocolBufferException:Protocol message was too large.May be
malicious.Use CodedInputStream.setSizeLimit()
to increase the size limit. at
com.google.protobuf.InvalidProtocolBufferException.sizeLimitExceeded(InvalidProtocolBufferException.java:1
10) at com.google.protobuf.CodedInputStream.refillBuffer(CodedInputStream.java:755)
at com.google.protobuf.CodedInputStream.readRawByte(CodedInputStream.java:769) at
com.google.protobuf.CodedInputStream.readRawVarint64(CodedInputStream.java:462) at
com.google.protobuf.
CodedInputStream.readSint64(CodedInputStream.java:363) at
org.apache.hadoop.hdfs.protocol.BlockListAsLongs$BufferDecoder$1.next(BlockListAsLongs.java:363)
```

数据目录中数据块的数量会显示为Metric。用户可以通过以下URL对该值进行监视
`http://<datanode-ip>:<http-port>/jmx`，如果该值超过4倍的限定值（4*1M），建议
用户配置多个驱动器并重新启动HDFS。

恢复步骤：

1. 在DataNode上配置多个数据目录。

示例：在原先只配置了/data1/datadir的位置

```
<property> <name>dfs.datanode.data.dir</name> <value>/data1/datadir</value> </property>
```

按照如下内容进行配置。

```
<property> <name>dfs.datanode.data.dir</name> <value>/data1/datadir,/data2/datadir,/data3/
datadir</value> </property>
```

📖 说明

建议多个数据目录应该配置到多个磁盘中，否则所有的数据都将写入同一个磁盘，对性能有很大的影响。

2. 重新启动HDFS。
3. 按照如下方法将数据移动至新的数据目录。

```
mv /data1/datadir/current/finalized/subdir1 /data2/datadir/current/finalized/
subdir1
```

4. 重新启动HDFS。

9.16.8 是否可以手动调整 DataNode 数据存储目录

问题

- 数据块在DataNode上的存储目录由“dfs.datanode.data.dir”配置项指定，是否可以修改该配置项来修改数据存储目录？
- 是否可以手动拷贝数据存储目录下的文件？

回答

“dfs.datanode.data.dir”配置项用于指定数据块在DataNode上的存储目录，在系统安装时需要指定根目录，并且可以指定多个根目录。

- 请谨慎修改该配置项，可以添加新的数据根目录。
- 禁止删除原有存储目录，否则会造成数据块丢失，导致文件无法正常读写。
- 禁止手动删除或修改存储目录下的数据块，否则可能会造成数据块丢失。

说明

NameNode和JournalNode存在类似的配置项，也同样禁止删除原有存储目录，禁止手动删除或修改存储目录下的数据块。

- dfs.namenode.edits.dir
- dfs.namenode.name.dir
- dfs.journalnode.edits.dir

9.16.9 DataNode 的容量计算出错如何处理

问题

当多个data.dir被配置在一个磁盘分区内，DataNode的容量计算将会出错。

回答

目前容量计算是基于磁盘的，类似于Linux里面的`df`命令。理想状态下，用户不会在同一个磁盘内配置多个data.dir，否则所有的数据都将写入一个磁盘，在性能上会有很大的影响。

因此配置如下：

例如，如果机器有如下磁盘：

```
host-4:~ # df -h
Filesystem      Size  Used Avail Use% Mounted on
/dev/sda1       352G  11G  324G   4%  /
udev            190G  252K  190G   1%  /dev
tmpfs           190G  72K   190G   1%  /dev/shm
/dev/sdb1       2.7T  74G   2.5T   3%  /data1
/dev/sdc1       2.7T  75G   2.5T   3%  /data2
/dev/sdd1       2.7T  73G   2.5T   3%  /da
```

建议的配置方式：

```
<property>
<name>dfs.datanode.data.dir</name>
<value>/data1/datadir1,/data2/datadir1,/data3/datadir1</value>
</property>
```

不建议的配置方式：

```
<property>
<name>dfs.datanode.data.dir</name>
<value>/data1/datadir1,/data2/datadir1,/data3/datadir1,/data1/datadir2,data1/datadir3,/data2/datadir2,/
data2/datadir3,/data3/datadir2,/data3/datadir3</value>
</property>
```

9.16.10 为什么存储小文件过程中，缓存中的数据会丢失

问题

在存储小文件过程中，系统断电，缓存中的数据丢失。

回答

由于断电，当写操作完成之后，缓存中的block不会立即被写入磁盘，如果要同步地将缓存的block写入磁盘，用户需要将“客户端安装路径/HDFS/hadoop/etc/hadoop/hdfs-site.xml”中的“dfs.datanode.synconclose”设置为“true”。

默认情况下，“dfs.datanode.synconclose”为“false”，虽然性能很高，但是断电之后，存储在缓存中的数据会丢失。将“dfs.datanode.synconclose”设置为“true”，可以解决此问题，但对性能有很大影响。请根据具体的应用场景决定是否开启该参数。

9.16.11 当分级存储策略为 LAZY_PERSIST 时为什么文件的副本的存储类型为 DISK

问题

当文件的存储策略为LAZY_PERSIST时，文件的第一副本的存储类型应为RAM_DISK，其余副本为DISK。

为什么文件的所有副本的存储类型都是DISK？

回答

当用户写入存储策略为LAZY_PERSIST的文件时，文件的三个副本会逐一写入。第一副本会优先选择客户端所在的DataNode节点，在以下情况下，当文件的存储策略为LAZY_PERSIST时，文件的所有副本的存储类型都是DISK：

- 当客户端所在的DataNode节点没有RAM_DISK时，则会写入客户端所在的DataNode节点的DISK磁盘，其余副本会写入其他节点的DISK磁盘。
- 当客户端所在的DataNode节点有RAM_DISK，但“dfs.datanode.max.locked.memory”参数值未设置或设置过小（小于“dfs.blocksize”参数值）时（对应参数值可登录Manager，选择“集群 > 服务 > HDFS > 配置 > 全部配置”搜索该参数获取），则会写入客户端所在的DataNode节点的DISK磁盘，其余副本会写入其他节点的DISK磁盘。

9.16.12 为什么 NameNode UI 上显示有一些块缺失

问题

回滚成功后，为什么NameNode UI上显示有一些块缺失？

回答

原因：具有新id/genstamps的块可能存在于DataNode上。DataNode中的块文件可能具有与NameNode的回滚image中不同的生成标记和长度，所以NameNode会拒绝DataNode中的这些块，并将文件标记为已损坏。

场景如下：

1. 升级前：
客户端A ->将一些数据写入文件X（假设已写入“A”字节）
2. 升级过程中：
客户端A ->仍然将数据写入文件X（现在文件中的数据是“A + B”字节）
3. 升级完成：
客户端A ->完成写入文件。最终数据为“A + B”字节。
4. 回滚开始：

将回滚到步骤1（升级前）的状态。因此，NameNode中的文件X将具有“A”字节，但DataNode中的块文件将具有“A + B”字节。

恢复步骤：

1. 从NameNode Web UI中获取已损坏的文件列表，或者通过下面的命令获取。

```
hdfs fsck <filepath> -list-corruptfileblocks
```

2. 对于不需要的文件，请使用以下命令删除文件。

```
hdfs fsck <corrupt file path> - delete
```

📖 说明

删除文件为高危操作，在执行操作前请务必确认对应文件是否不再需要。

3. 对于所需的文件，执行fsck命令来获取块列表和块的顺序。
 - 在fsck中给出的块序列列表中，使用块ID搜索DataNode中的数据目录，并从DataNode下载相应的块。
 - 按照序列以追加的方式写入所有这样的块文件，并构造成原始文件。例如：
File 1--> blk_1, blk_2, blk_3
通过组合来自同一序列的所有三个块文件的内容来创建文件。
 - 从HDFS中删除旧文件并重写新构建的文件。

9.17 HDFS 故障排除

9.17.1 往 HDFS 写数据时报错 “java.net.SocketException”

问题

为什么在往HDFS写数据时报“java.net.SocketException: No buffer space available”异常？

这个问题发生在往HDFS写文件时。查看客户端和DataNode的错误日志。

客户端日志如下：

图 9-14 客户端日志

```
2017-07-05 21:58:06.459 INFO [htable-pool3-t1] ipc.AbstractPqccClient: RPC Server Kerberos principal name for service=ClientService is hbase/hadoop.hadoop123.com@H4000P12
2017-07-05 21:58:06.893 WARN [main] mapreduce.LoadIncrementalHFiles: Skipping non-directory hdfs://hacluster/HBaseTest/bulkload_output/_SUCCESS
2017-07-05 21:59:13.211 WARN [main] hdfs.BlockReaderFactory: I/O error constructing remote block reader.
java.net.SocketException: No buffer space available
    at sun.nio.ch.Net.connect(Native Method)
    at sun.nio.ch.Net.connect(Net.java:454)
    at sun.nio.ch.SocketChannelImpl.connect(SocketChannelImpl.java:648)
    at org.apache.hadoop.net.SocketIOWithTimeout.connect(SocketIOWithTimeout.java:192)
    at org.apache.hadoop.net.NetUtils.connect(NetUtils.java:531)
    at org.apache.hadoop.hdfs.DFSClient.newConnectedPeer(DFSClient.java:3345)
    at org.apache.hadoop.hdfs.BlockReaderFactory.nextTcpPeer(BlockReaderFactory.java:789)
    at org.apache.hadoop.hdfs.BlockReaderFactory.getRemoteBlockReaderFromTcp(BlockReaderFactory.java:706)
    at org.apache.hadoop.hdfs.BlockReaderFactory.build(BlockReaderFactory.java:369)
    at org.apache.hadoop.hdfs.DFSInputStream.getBlockReader(DFSInputStream.java:713)
    at org.apache.hadoop.hdfs.DFSInputStream.blockSeekTo(DFSInputStream.java:663)
    at org.apache.hadoop.hdfs.DFSInputStream.readWithStrategy(DFSInputStream.java:919)
    at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:973)
    at java.io.DataInputStream.readFully(DataInputStream.java:195)
    at org.apache.hadoop.hbase.io.hfile.FixedFileTrailer.readFromStream(FixedFileTrailer.java:391)
    at org.apache.hadoop.hbase.io.hfile.HFile.isHFileFormat(HFile.java:578)
    at org.apache.hadoop.hbase.io.hfile.HFile.isHFileFormat(HFile.java:560)
    at org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles.visitBulkHFiles(LoadIncrementalHFiles.java:229)
    at org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles.discoverLoadQueue(LoadIncrementalHFiles.java:281)
    at org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles.prepareHFileQueue(LoadIncrementalHFiles.java:452)
    at org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles.doBulkLoad(LoadIncrementalHFiles.java:365)
    at org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles.doBulkLoad(LoadIncrementalHFiles.java:331)
    at org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles.run(LoadIncrementalHFiles.java:1187)
    at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:70)
    at org.apache.hadoop.hbase.mapreduce.LoadIncrementalHFiles.main(LoadIncrementalHFiles.java:1114)
2017-07-05 21:59:13.215 WARN [main] hdfs.DFSClient: Failed to connect to /192.168.152.128:25009 for block BP-1989348819-192.168.199.5-1497961637591:blk_1107301222_335745
ffer space available
java.net.SocketException: No buffer space available
    at sun.nio.ch.Net.connect(Native Method)
    at sun.nio.ch.Net.connect(Net.java:454)
    at sun.nio.ch.SocketChannelImpl.connect(SocketChannelImpl.java:648)
    at org.apache.hadoop.net.SocketIOWithTimeout.connect(SocketIOWithTimeout.java:192)
    at org.apache.hadoop.net.NetUtils.connect(NetUtils.java:531)
    at org.apache.hadoop.hdfs.DFSClient.newConnectedPeer(DFSClient.java:3345)
```


DataNode日志如下：

```
2017-07-24 20:43:39,269 | ERROR | DataXceiver for client DFSCient_NONMAPREDUCE_996005058_86
at /192.168.164.155:40214 [Receiving block
BP-1287143557-192.168.199.6-1500707719940:blk_1074269754_528941 with io weight 10] |
DataNode{data=FSDataset{dirpath='[/srv/BigData/hadoop/data1/dn/current, /srv/BigData/hadoop/
data2/dn/current, /srv/BigData/hadoop/data3/dn/current, /srv/BigData/hadoop/data4/dn/current, /srv/
BigData/hadoop/data5/dn/current, /srv/BigData/hadoop/data6/dn/current, /srv/BigData/hadoop/data7/dn/
current]'}, localName='192-168-164-155:9866', datanodeUuid='a013e29c-4e72-400c-bc7b-bbbf0799604c',
xmitsInProgress=0}:Exception transferring block
BP-1287143557-192.168.199.6-1500707719940:blk_1074269754_528941 to mirror 192.168.202.99:9866:
java.net.SocketException: No buffer space available | DataXceiver.java:870
2017-07-24 20:43:39,269 | INFO | DataXceiver for client DFSCient_NONMAPREDUCE_996005058_86
at /192.168.164.155:40214 [Receiving block
BP-1287143557-192.168.199.6-1500707719940:blk_1074269754_528941 with io weight 10] | opWriteBlock
BP-1287143557-192.168.199.6-1500707719940:blk_1074269754_528941 received exception
java.net.SocketException: No buffer space available | DataXceiver.java:933
2017-07-24 20:43:39,270 | ERROR | DataXceiver for client DFSCient_NONMAPREDUCE_996005058_86
at /192.168.164.155:40214 [Receiving block
BP-1287143557-192.168.199.6-1500707719940:blk_1074269754_528941 with io weight 10] |
192-168-164-155:9866:DataXceiver error processing WRITE_BLOCK operation src: /192.168.164.155:40214
dst: /192.168.164.155:9866 | DataXceiver.java:304 java.net.SocketException: No buffer space available
at sun.nio.ch.Net.connect0(Native Method)
at sun.nio.ch.Net.connect(Net.java:454)
at sun.nio.ch.Net.connect(Net.java:446)
at sun.nio.ch.SocketChannelImpl.connect(SocketChannelImpl.java:648)
at org.apache.hadoop.net.SocketIOWithTimeout.connect(SocketIOWithTimeout.java:192)
at org.apache.hadoop.net.NetUtils.connect(NetUtils.java:531)
at org.apache.hadoop.net.NetUtils.connect(NetUtils.java:495)
at org.apache.hadoop.hdfs.server.datanode.DataXceiver.writeBlock(DataXceiver.java:800)
at org.apache.hadoop.hdfs.protocol.datatransfer.Receiver.opWriteBlock(Receiver.java:138)
at org.apache.hadoop.hdfs.protocol.datatransfer.Receiver.processOp(Receiver.java:74)
at org.apache.hadoop.hdfs.server.datanode.DataXceiver.run(DataXceiver.java:265)
at java.lang.Thread.run(Thread.java:748)
```

回答

上述问题可能是因为网络内存枯竭而导致的。

问题的解决方案是根据实际场景适当增大网络设备的阈值级别。

例如：

```
[root@xxxx ~]# cat /proc/sys/net/ipv4/neigh/default/gc_thresh*
128
512
1024
[root@xxxx ~]# echo 512 > /proc/sys/net/ipv4/neigh/default/gc_thresh1
[root@xxxx ~]# echo 2048 > /proc/sys/net/ipv4/neigh/default/gc_thresh2
[root@xxxx ~]# echo 4096 > /proc/sys/net/ipv4/neigh/default/gc_thresh3
[root@xxxx ~]# cat /proc/sys/net/ipv4/neigh/default/gc_thresh*
512
2048
4096
```

还可以将以下参数添加到“/etc/sysctl.conf”中，即使主机重启，配置依然能生效。

```
net.ipv4.neigh.default.gc_thresh1 = 512
net.ipv4.neigh.default.gc_thresh2 = 2048
net.ipv4.neigh.default.gc_thresh3 = 4096
```

9.17.2 删除大量文件后重启 NameNode 耗时长

问题

删除大量文件之后立刻重启NameNode（例如删除100万个文件），NameNode启动慢。

回答

由于在删除了大量文件之后，DataNode需要时间去删除对应的Block。当立刻重启NameNode时，NameNode会去检查所有DataNode上报的Block信息，发现已删除的Block时，会输出对应的INFO日志信息，如下所示：

```
2015-06-10 19:25:50,215 | INFO | IPC Server handler 36 on 25000 | BLOCK* processReport:
blk_1075861877_2121067 on node 10.91.8.218:9866 size 10249 does not belong to any file |
org.apache.hadoop.hdfs.server.blockmanagement.BlockManager.processReport(BlockManager.java:1854)
```

每一个被删除的Block会产生一条日志信息，一个文件可能会存在一个或多个Block。当删除的文件数过多时，NameNode会花大量的时间打印日志，然后导致NameNode启动慢。

当出现这种现象时，您可以通过如下方式提升NameNode的启动速度。

1. 删除大量文件时，不要立刻重启NameNode，待DataNode删除了对应的Block后重启NameNode，即不会存在这种情况。
您可以通过 `hdfs dfsadmin -report` 命令来查看磁盘空间，检查文件是否删除完毕。
2. 如已大量出现以上日志，您可以将NameNode的日志级别修改为ERROR，NameNode不会再打印此日志信息。
等待NameNode启动完毕后，再将此日志级别修改为INFO。修改日志级别后无需重启服务。

9.17.3 EditLog 不连续导致 NameNode 启动失败

问题

在JournalNode节点有断电，数据目录磁盘占满，网络异常时，会导致JournalNode上的EditLog不连续。此时如果重启NameNode，很可能会失败。

现象

重启NameNode会失败。在NameNode运行日志中会报如下的错误：

```
2019-11-08 16:30:28,399 | ERROR | main | Failed to start namenode. | NameNode.java:1732
java.io.IOException: There appears to be a gap in the edit log. We expected txid 13698019, but got txid 13698088.
at org.apache.hadoop.hdfs.server.namenode.MetaRecoveryContext.editLogLoaderPrompt(MetaRecoveryContext.java:94)
at org.apache.hadoop.hdfs.server.namenode.FSEditLogLoader.loadEditRecords(FSEditLogLoader.java:278)
at org.apache.hadoop.hdfs.server.namenode.FSEditLogLoader.loadFSEdits(FSEditLogLoader.java:188)
at org.apache.hadoop.hdfs.server.namenode.FSImage.loadEdits(FSImage.java:924)
at org.apache.hadoop.hdfs.server.namenode.FSImage.loadFSImage(FSImage.java:771)
at org.apache.hadoop.hdfs.server.namenode.FSImage.recoverTransitionRead(FSImage.java:331)
at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.loadFSImage(FSNamesystem.java:1108)
at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.loadFromDisk(FSNamesystem.java:727)
at org.apache.hadoop.hdfs.server.namenode.NameNode.loadNamesystem(NameNode.java:638)
at org.apache.hadoop.hdfs.server.namenode.NameNode.initialize(NameNode.java:700)
at org.apache.hadoop.hdfs.server.namenode.NameNode.<init>(NameNode.java:943)
at org.apache.hadoop.hdfs.server.namenode.NameNode.<init>(NameNode.java:916)
at org.apache.hadoop.hdfs.server.namenode.NameNode.createNameNode(NameNode.java:1655)
at org.apache.hadoop.hdfs.server.namenode.NameNode.main(NameNode.java:1725)
```

解决方法

1. 找到重启前的主NameNode，进入其数据目录（查看配置项“dfs.namenode.name.dir”可获取，例如/srv/BigData/namenode/current），得到最新的FSImage文件的序号。一般如下：

```

-rw-----, 1 omm wheel      574 Oct  2 01:12 edits_000000000013259401-0000000000:
-rw-----, 1 omm wheel      575 Oct  2 01:13 edits_000000000013259409-0000000000:
-rw-----, 1 omm wheel       42 Oct  2 01:13 edits_000000000013259417-0000000000:
-rw-----, 1 omm wheel 1048576 Nov  8 16:01 edits_inprogress_000000000013698088
-rw-----, 1 omm wheel  314803 Nov  8 15:53 fsimage_000000000013698018
-rw-----, 1 omm wheel       62 Nov  8 15:53 fsimage_000000000013698018.md5
-rw-----, 1 omm wheel  314803 Nov  8 15:56 fsimage_000000000013698050
-rw-----, 1 omm wheel       62 Nov  8 15:56 fsimage_000000000013698050.md5
-rw-----, 1 omm wheel  314803 Nov  8 15:59 fsimage_000000000013698066
-rw-----, 1 omm wheel       62 Nov  8 15:59 fsimage_000000000013698066.md5
-rw-----, 1 omm wheel       9 Oct  2 01:13 seen_txid
-rw-----, 1 omm wheel      187 Nov  8 15:59 VERSION

```

- 查看各JournalNode的数据目录（查看配置项“dfs.journalnode.edits.dir”可获取，例如/srv/BigData/journalnode/hacluster/current），查看序号从第一步获取到的序号开始的edits文件，看是否有不连续的情况（即前一个edits文件的最后一个序号和后一个edits文件的第一个序号不是连续的，如下图中的edits_000000000013259231-000000000013259237和后一个edits_000000000013259239-000000000013259246就是不连续的）。

```

-rw-----, 1 omm wheel      576 Oct  2 00:41 edits_000000000013259151-000000000013259158
-rw-----, 1 omm wheel      575 Oct  2 00:43 edits_000000000013259159-000000000013259166
-rw-----, 1 omm wheel      576 Oct  2 00:43 edits_000000000013259167-000000000013259174
-rw-----, 1 omm wheel      575 Oct  2 00:45 edits_000000000013259175-000000000013259182
-rw-----, 1 omm wheel      575 Oct  2 00:45 edits_000000000013259183-000000000013259190
-rw-----, 1 omm wheel      576 Oct  2 00:47 edits_000000000013259191-000000000013259198
-rw-----, 1 omm wheel      575 Oct  2 00:48 edits_000000000013259199-000000000013259206
-rw-----, 1 omm wheel      575 Oct  2 00:49 edits_000000000013259207-000000000013259214
-rw-----, 1 omm wheel      575 Oct  2 00:50 edits_000000000013259215-000000000013259222
-rw-----, 1 omm wheel      573 Oct  2 00:51 edits_000000000013259223-000000000013259230
-rw-----, 1 omm wheel      571 Oct  2 00:52 edits_000000000013259231-000000000013259237
-rw-----, 1 omm wheel      576 Oct  2 00:53 edits_000000000013259239-000000000013259246
-rw-----, 1 omm wheel      575 Oct  2 00:54 edits_000000000013259247-000000000013259254
-rw-----, 1 omm wheel      576 Oct  2 00:55 edits_000000000013259255-000000000013259262
-rw-----, 1 omm wheel       42 Oct  2 00:56 edits_000000000013259263-000000000013259264
-rw-----, 1 omm wheel  1107 Oct  2 00:57 edits_000000000013259265-000000000013259278
-rw-----, 1 omm wheel       42 Oct  2 00:58 edits_000000000013259279-000000000013259280
-rw-----, 1 omm wheel  1109 Oct  2 00:59 edits_000000000013259281-000000000013259294
-rw-----, 1 omm wheel       42 Oct  2 01:00 edits_000000000013259295-000000000013259296
-rw-----, 1 omm wheel  1299 Oct  2 01:01 edits_000000000013259297-000000000013259312
-rw-----, 1 omm wheel      260 Oct  2 01:02 edits_000000000013259313-000000000013259316
-rw-----, 1 omm wheel      984 Oct  2 01:03 edits_000000000013259317-000000000013259328
-rw-----, 1 omm wheel      572 Oct  2 01:04 edits_000000000013259329-000000000013259336
-rw-----, 1 omm wheel      575 Oct  2 01:05 edits_000000000013259337-000000000013259344
-rw-----, 1 omm wheel      983 Oct  2 01:06 edits_000000000013259345-000000000013259356

```

- 如果有这种不连续的edits文件，则需要查看其它的JournalNode的数据目录或NameNode数据目录中，有没有连续的该序号相关的连续的edits文件。如果可以找到，复制一个连续的片段到该JournalNode。
- 如此把所有的不连续的edits文件全部都修复。
- 重启NameNode，观察是否成功。如还是失败，请联系技术支持。

9.17.4 当备 NameNode 存储元数据时，断电后备 NameNode 启动失败

问题

当Standby NameNode存储元数据（命名空间）时，出现断电的情况，Standby NameNode启动失败并发生如下错误信息。

```
2015-12-04 11:49:12,121 | ERROR | main | Failed to load image from FS
ImageFile(file=/srv/BigData/namenode/current/fsimage_0000000000000096
080,
cpktTxId=0000000000000096080) | FSImage.java:685
java.io.IOException: Invalid MD5 file /srv/BigData/namenode/current/f
simage_0000000000000096080.md5:
the content "棍斤拷棍斤拷棍斤拷棍斤拷棍[1m^A!棍 does not match the expecte
d pattern.
at org.apache.hadoop.hdfs.util.MD5FileUtils.readStoredMd5(MD5FileUtil
s.java:92)
at org.apache.hadoop.hdfs.util.MD5FileUtils.readStoredMd5ForFile(MD5F
ileUtils.java:109)
at org.apache.hadoop.hdfs.server.namenode.FSImage.loadFSImage(FSImage
.java:975)
at org.apache.hadoop.hdfs.server.namenode.FSImage.loadFSImageFile(FSI
mage.java:744)
at org.apache.hadoop.hdfs.server.namenode.FSImage.loadFSImage(FSImage
.java:682)
at org.apache.hadoop.hdfs.server.namenode.FSImage.recoverTransitionRe
ad(FSImage.java:300)
at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.loadFSImage(FS
Namesystem.java:968)
at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.loadFromDisk(F
SNamesystem.java:675)
at org.apache.hadoop.hdfs.server.namenode.NameNode.loadNamesystem(Nam
eNode.java:625)
at org.apache.hadoop.hdfs.server.namenode.NameNode.initialize(NameNod
e.java:685)
at org.apache.hadoop.hdfs.server.namenode.NameNode.<init>(NameNode.ja
va:889)
at org.apache.hadoop.hdfs.server.namenode.NameNode.<init>(NameNode.ja
va:872)
at org.apache.hadoop.hdfs.server.namenode.NameNode.createNameNode(Nam
eNode.java:1580)
at org.apache.hadoop.hdfs.server.namenode.NameNode.main(NameNode.java
:1654)
```

回答

当Standby NameNode存储元数据（命名空间）时，出现断电的情况，Standby NameNode启动失败，MD5文件会损坏。通过移除损坏的fsimage，然后启动Standby NameNode，可以修复此问题。Standby NameNode会加载先前的fsimage并重现所有的edits。

修复步骤：

1. 移除损坏的fsimage。

```
rm -rf ${BIGDATA_DATA_HOME}/namenode/current/
fsimage_0000000000000096
```
2. 启动Standby NameNode。

9.17.5 dfs.datanode.data.dir 中定义的磁盘数量等于 dfs.datanode.failed.volumes.tolerated 的值时，DataNode 启动失败

问题

当“dfs.datanode.data.dir”中定义的磁盘数量等于“dfs.datanode.failed.volumes.tolerated”的值时，DataNode启动失败。

回答

默认情况下，单个磁盘的故障将会引起HDFS DataNode进程关闭，导致NameNode为每一个存在DataNode上的block调度额外的副本，在没有故障的磁盘中引起不必要的块复制。

为了防止此情况，用户可以通过配置DataNodes来承受dfs.data.dir目录的故障。登录Manager，选择“集群 > 服务 > HDFS > 配置 > 全部配置”搜索参数“dfs.datanode.failed.volumes.tolerated”。例如：如果该参数值为3，DataNode只有在4个或者更多个目录故障之后才会出现故障。该值会影响到DataNode的启动。

如果想要DataNode不出现故障，配置的“dfs.datanode.failed.volumes.tolerated”一定要小于所配置的卷数，也可以将“dfs.datanode.failed.volumes.tolerated”设置成-1，相当于设置该值为n-1（n为卷数），那样DataNode就不会出现启动失败。

9.17.6 HDFS 调用 FileInputFormat 的 getsplit 的时候出现数组越界

问题

HDFS调用FileInputFormat的getSplit方法的时候，出现ArrayIndexOutOfBoundsException: 0，日志如下：

```
java.lang.ArrayIndexOutOfBoundsException: 0
at org.apache.hadoop.mapred.FileInputFormat.identifyHosts(FileInputFormat.java:708)
at org.apache.hadoop.mapred.FileInputFormat.getSplitHostsAndCachedHosts(FileInputFormat.java:675)
at org.apache.hadoop.mapred.FileInputFormat.getSplits(FileInputFormat.java:359)
at org.apache.spark.rdd.HadoopRDD.getPartitions(HadoopRDD.scala:210)
at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:239)
at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:237)
at scala.Option.getOrElse(Option.scala:120)
at org.apache.spark.rdd.RDD.partitions(RDD.scala:237)
at org.apache.spark.rdd.MapPartitionsRDD.getPartitions(MapPartitionsRDD.scala:35)
```

回答

每个block对应的机架信息组成为：/default/rack0/;/default/rack0/datanodeip:port。

该问题是由于某个block块损坏或者丢失，导致该block对应的机器ip和port为空引起的，出现该问题的时候使用**hdfs fsck**检查对应文件块的健康状态，删除损坏或者恢复丢失的块，重新进行任务计算即可。

10 使用 Hive

10.1 Hive 用户权限管理

10.1.1 Hive 用户权限说明

Hive是建立在Hadoop上的数据仓库框架，提供类似SQL的HQL操作结构化数据。

MRS提供用户、用户组和角色，集群中的各类权限需要先授予角色，然后将用户或者用户组与角色绑定。用户只有绑定角色或者加入绑定角色的用户组，才能获得权限。Hive授权相关信息请参考：<https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Authorization>。

说明

- Hive在安全模式下需要进行权限管理，在普通模式下无需进行权限管理。
- MRS 3.x及后续版本支持Ranger，如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加Hive的Ranger访问权限策略](#)。

Hive 权限模型

使用Hive组件，必须对Hive数据库和表（含外表和视图）拥有相应的权限。在MRS中，完整的Hive权限模型由Hive元数据权限与HDFS文件权限组成。使用数据库或表时所需要的各种权限都是Hive权限模型中的一种。

- Hive元数据权限。
与传统关系型数据库类似，MRS的Hive数据库包含“建表”和“查询”权限，Hive表和列包含“查询”、“插入”和“删除”权限。Hive中还包含拥有者权限“OWNERSHIP”和“Hive管理员权限”。
- Hive数据文件权限，即HDFS文件权限。
Hive的数据库、表对应的文件保存在HDFS中。默认创建的数据库或表保存在HDFS目录“/user/hive/warehouse”。系统自动以数据库名称和数据库中表的名称创建子目录。访问数据库或者表，需要在HDFS中拥有对应文件的权限，包含“读”、“写”和“执行”权限。

用户对Hive数据库或表执行不同操作时，需要关联不同的元数据权限与HDFS文件权限。例如，对Hive数据表执行查询操作，需要关联元数据权限“查询”，以及HDFS文件权限“读”和“写”。

使用Manager界面图形化的角色管理功能来管理Hive数据库和表的权限，只需要设置元数据权限，系统会自动关联HDFS文件权限，减少界面操作，提高效率。

Hive 用户对象

MRS提供了用户和角色来使用Hive，比如创建表、在表中插入数据或者查询表。Hive中定义了“USER”类，对应用户实例；定义了“GROUP”类，对应角色实例。

使用Manager设置Hive用户对象的权限，只支持在角色中设置，用户或用户组需要绑定角色才能获得权限。支持授予Hive管理员权限、访问数据库、表和列的权限。

Hive 使用场景及对应权限

用户使用Hive并创建数据库需要加入hive组，不需要角色授权。用户在Hive和HDFS中对自己创建的数据库或表拥有完整权限，可直接创建表、查询数据、删除数据、插入数据、更新数据以及授权他人访问表与对应HDFS目录与文件。

如果用户访问别人创建的表或数据库，需要授予权限。所以根据Hive使用场景的不同，用户需要的权限可能也不相同。

表 10-1 Hive 使用场景

主要场景	用户需要的权限
使用Hive表、列或数据库	使用其他用户创建的Hive表、列或数据库，不同的场景需要不同的Hive权限，例如： <ul style="list-style-type: none">• 创建表，需要“建表”。• 查询数据，需要“查询”。• 插入数据，需要“插入”。• 删除数据，需要“删除”。
关联使用其他组件	部分场景除了Hive权限，还可能需要组件的权限，例如： <ul style="list-style-type: none">• 执行部分HQL命令，例如insert，count，distinct，group by，order by，sort by或join等语句时，需要设置YARN权限。建议为每个Hive用户的角色添加此权限。• 使用Hive over HBase，例如在Hive中查询HBase表数据，需要设置HBase权限。

在一些特殊Hive使用场景下，需要单独设置其他权限。

表 10-2 Hive 授权注意事项

可能场景	用户需要的权限
<p>创建Hive数据库、表、外表，或者为已经创建的Hive表或外表添加分区，且Hive用户指定数据文件保存在“/user/hive/warehouse”以外的HDFS目录。</p>	<p>需要此目录已经存在，Hive用户是目录的属主，且用户对目录拥有“读”、“写”和“执行”权限。同时用户对此目录上层的每一级目录都拥有“读”和“写”权限。然后管理员通过角色管理功能授予角色使用Hive的权限，会自动关联HDFS权限。</p>
<p>Hive用户使用load将指定目录下所有文件或者指定文件，导入数据到Hive表。</p>	<ul style="list-style-type: none"> ● 数据源为Linux本地磁盘，指定目录时需要此目录已经存在，系统用户“omm”对此目录以及此目录上层的每一级目录拥有“r”和“x”的权限。指定文件时需要此文件已经存在，“omm”对此文件拥有“r”的权限，同时对此文件上层的每一级目录拥有“r”和“x”的权限。 ● 数据源为HDFS，指定目录时需要此目录已经存在，Hive用户是目录属主，且用户对此目录及其子目录拥有“读”、“写”和“执行”权限，并且其上层的每一级目录拥有“读”和“写”权限。指定文件时需要此文件已经存在，Hive用户是文件属主，且用户对文件拥有“读”、“写”和“执行”权限，同时对此文件上层的每一级目录拥有“读”和“执行”权限。 <p>说明 使用load从Linux本地磁盘导入数据时，文件需上传到执行命令的HiveServer并修改权限。建议使用客户端执行命令，可查看客户端连接的HiveServer。例如，Hive客户端显示“0: jdbc:hive2://10.172.0.43:21066/>”，表示当前连接的HiveServer节点IP地址为“10.172.0.43”。</p>
<p>创建函数、删除函数或者修改任意数据库。</p>	<p>需要授予“Hive管理员权限”。</p>
<p>操作Hive中所有的数据库和表。</p>	<p>需加入到supergroup用户组，并且授予“Hive管理员权限”。</p>

可能场景	用户需要的权限
集群未启用Kerberos认证（普通模式） 开启Ranger鉴权	集群未启用Kerberos认证（普通模式） 默认关闭Ranger鉴权，如果启用Ranger鉴权，新增以下限制： <ul style="list-style-type: none"> 白名单限制：未配置的参数将不允许在客户端设置。白名单开关由Hive配置页面的“hive.security.whitelist.switch”参数控制，设置为“OFF”即可支持在客户端设置未配置的参数，存在安全风险，请谨慎操作。 不允执行reflect、reflect2、java_method和in_file函数。如果业务需要，可在HiveServer自定义参数中新增“hive.server2.builtin.udf.blacklist”参数项，值为“mpty_blacklist”来允许Hive执行这些函数，存在安全风险，请谨慎操作。

10.1.2 创建 Hive 角色

操作场景

该任务指导MRS集群管理员在Manager创建并设置Hive的角色。Hive角色可设置Hive管理员权限以及Hive数据表的数据操作权限。

用户使用Hive并创建数据库需要加入hive组，不需要角色授权。用户在Hive和HDFS中对自己创建的数据库或表拥有完整权限，可直接创建表、查询数据、删除数据、插入数据、更新数据以及授权他人访问表与对应HDFS目录与文件。默认创建的数据库或表保存在HDFS目录“/user/hive/warehouse”。

📖 说明

- 安全模式支持创建Hive角色，普通模式不支持创建Hive角色。
- MRS 3.x及后续版本支持Ranger，如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加Hive的Ranger访问权限策略](#)。

前提条件

- MRS集群管理员已明确业务需求。
- 已登录Manager。
- 已安装好Hive客户端。

操作步骤

MRS 3.x之前版本，创建Hive角色的操作如下：

- 步骤1** 登录MRS Manager。
- 步骤2** 选择“系统设置 > 权限配置 > 角色管理”。
- 步骤3** 单击“添加角色”，输入“角色名称”和“描述”。
- 步骤4** 设置角色“权限”请参见表10-3。



- “Hive Admin Privilege”：Hive管理员权限。如需使用该权限，在执行SQL语句时需要先执行set role admin来设置权限。
- “Hive Read Write Privileges”：Hive数据表管理权限，可设置与管理已创建的表的数据操作权限。根据需要勾选相应database的权限，如果要精确到表，可以单击database名称，勾选相应表的权限。

说明

- Hive角色管理支持授予Hive管理员权限、访问表和视图的权限，不支持数据库的授权。
- Hive管理员权限不支持管理HDFS的权限。
- 如果数据库中的表或者表中的文件数量比较多，在授权时可能需要等待一段时间。例如表的文件数量为1万时，可能需要等待2分钟。

表 10-3 设置角色

任务场景	角色授权操作
设置Hive管理员权限	<p>在“权限”的表格中单击“Hive”，勾选“Hive Admin Privilege”。</p> <p>说明 用户绑定Hive管理员角色后，在每个维护操作会话中，还需要执行以下操作：</p> <ol style="list-style-type: none"> 1. 请根据客户端所在位置，参考安装客户端章节，登录安装客户端的节点。 2. 执行以下命令配置环境变量。 例如，Hive客户端安装目录为“/opt/hiveclient”，执行source /opt/hiveclient/bigdata_env 3. 执行以下命令认证用户。 kinit Hive业务用户 4. 执行以下命令登录客户端工具。 beeline 5. 执行以下命令更新Hive用户的管理员权限。 set role admin;

任务场景	角色授权操作
设置在默认数据库中，查询其他用户表的权限	<ol style="list-style-type: none"> 1. 在“权限”的表格中选择“Hive > Hive Read Write Privileges”。 2. 在指定表的“权限”列，勾选“Select”。
设置在默认数据库中，插入其他用户表的权限	<ol style="list-style-type: none"> 1. 在“权限”的表格中选择“Hive > Hive Read Write Privileges”。 2. 在指定表的“权限”列，勾选“Insert”。
设置在默认数据库中，导入数据到其他用户表的权限	<ol style="list-style-type: none"> 1. 在“权限”的表格中选择“Hive > Hive Read Write Privileges”。 2. 在指定表的“权限”列，勾选“Delete”和“Insert”。
设置提交Hql命令到Yarn执行的权限	<p>部分业务需求使用的Hql命令将转化为MapReduce任务并提交到Yarn中执行，需要设置Yarn权限。例如运行的HQL使用了insert, count, distinct, group by, order by, sort by或join等语句的相关场景。</p> <ol style="list-style-type: none"> 1. 在“权限”的表格中选择“Yarn > Scheduler Queue > root”。 2. 在“default”队列的“权限”列，勾选“Submit”。

步骤5 单击“确定”，返回“角色”。

步骤6 选择“系统设置 > 用户管理 > 添加用户”。

步骤7 输入用户名，在“用户类型”选择“人机”类型，设置用户密码，在用户组添加一个绑定了Hive管理员角色的用户组，并绑定新创建的Hive角色，单击“确定”完成Hive用户创建。

步骤8 待用户生成后，即可使用该用户执行相应SQL语句。

----结束

MRS 3.x及后续版本，创建Hive角色的操作如下：

步骤1 登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。

步骤2 选择“系统 > 权限 > 角色”。

步骤3 单击“添加角色”，输入“角色名称”和“描述”。

步骤4 设置角色“配置资源权限”请参见[表10-4](#)。

- 设置HDFS目录的读和执行权限。
 - 选择“待操作集群的名称 > HDFS > 文件系统 > hdfs://hacluster/ > user”，在“hive”的“权限”列，勾选“读”和“执行”。

- 选择“待操作集群的名称 > HDFS > 文件系统 > hdfs://hacluster/ > user > hive”，在“warehouse”的“权限”列，勾选“读”和“执行”。
- 选择“待操作集群的名称 > HDFS > 文件系统 > hdfs://hacluster/ > tmp”，在“hive-scratch”的“权限”列，勾选“读”和“执行”。
- “Hive管理员权限”：Hive管理员权限。
- “Hive读写权限”：Hive数据表管理权限，可设置与管理已创建的表的数据操作权限。

说明

- MRS 3.1.0版本，Hive角色管理支持授予Hive管理员权限、访问表和视图的权限，不支持数据库的授权。
- MRS 3.1.2及之后版本，Hive角色管理支持授予管理员权限、访问库、表和视图的权限。
- Hive管理员权限不支持管理HDFS的权限。
- 如果数据库中的表或者表中的文件数量比较多，在授权时可能需要等待一段时间。例如表的文件数量为1万时，可能需要等待2分钟。

表 10-4 设置角色

任务场景	角色授权操作
设置Hive管理员权限	<p>在“配置资源权限”的表格中选择“待操作集群的名称 > Hive”，勾选“Hive管理员权限”。</p> <p>说明 用户绑定Hive管理员角色后，在每个维护操作会话中，还需要执行以下操作：</p> <ol style="list-style-type: none"> 1. 以客户端安装用户，登录安装Hive客户端的节点。 2. 执行以下命令配置环境变量。 例如，Hive客户端安装目录为“/opt/hiveclient”，执行source /opt/hiveclient/bigdata_env 3. 执行以下命令认证用户。 kinit Hive业务用户 4. 执行以下命令登录客户端工具。 beeline 5. 执行以下命令更新Hive用户的管理员权限。 set role admin;
设置在默认数据库中，查询其他用户表的权限	<ol style="list-style-type: none"> 1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive读写权限”。 2. 在数据库列表中单击指定的数据库名称，显示数据库中的表。 3. 在指定表的“权限”列，勾选“查询”。

任务场景	角色授权操作
设置在默认数据库中，插入其他用户表的权限	<ol style="list-style-type: none"> 1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive读写权限”。 2. 在数据库列表中单击指定的数据库名称，显示数据库中的表。 3. 在指定表的“权限”列，勾选“插入”。
设置在默认数据库中，导入数据到其他用户表的权限	<ol style="list-style-type: none"> 1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive读写权限”。 2. 在数据库列表中单击指定的数据库名称，显示数据库中的表。 3. 在指定表的“权限”列，勾选“删除”和“插入”。
设置提交Hql命令到Yarn执行的权限	<p>部分业务需求使用的Hql命令将转化为MapReduce任务并提交到Yarn中执行，需要设置Yarn权限。例如运行的HQL使用了insert, count, distinct, group by, order by, sort by或join等语句的相关场景。</p> <ol style="list-style-type: none"> 1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Yarn > 调度队列 > root”。 2. 在“default”队列的“权限”列，勾选“提交”。

步骤5 单击“确定”，返回“角色”。

步骤6 选择“权限 > 用户”，单击“添加用户”。

步骤7 输入用户名，选择“用户类型”选择“人机”类型，设置用户密码，在用户组添加Hive相应权限的用户组并选择主组，绑定新创建的角色，单击“确定”完成Hive用户创建。

步骤8 待用户生成后，即可使用该用户执行相应SQL语句。

----结束

10.1.3 配置 Hive 表、列或数据库的用户权限

操作场景

使用Hive表或者数据库时，如果用户访问别人创建的表或数据库，需要授予对应的权限。为了实现更严格权限控制，Hive也支持列级别的权限控制。如果要访问别人创建的表上某些列，需要授予列权限。以下介绍使用Manager角色管理功能在表授权、列授权和数据库授权三个场景下的操作。

说明

- 安全模式支持配置Hive表、列或数据库的权限，普通模式不支持配置Hive表、列或数据库的权限。
- MRS 3.x及后续版本支持Ranger，如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加Hive的Ranger访问权限策略](#)。

前提条件

- 获取一个拥有管理员权限的用户，例如“admin”。
- 请参考[创建Hive角色](#)，在Manager界面创建一个角色，例如“hrole”，不需要设置Hive权限，设置提交Hql命令到Yarn执行的权限。
- 在Manager界面创建两个使用Hive的“人机”用户并加入“hive”组，例如“huser1”和“huser2”。“huser2”需绑定“hrole”。使用“huser1”创建一个数据库“hdb”，并在此数据库中创建表“htable”。

操作步骤

- 表授权
用户在Hive和HDFS中对自己创建的表拥有完整权限，用户访问别人创建的表，需要授予权限。授予权限时只需要授予Hive元数据权限，HDFS文件权限将自动关联。以授予用户对应角色在表“htable”中查询、插入和删除数据的权限为例，操作步骤如下：

MRS 3.x之前版本，表授权的操作如下：

- a. 在MRS Manager界面，选择“系统设置 > 权限配置 > 角色管理”。
- b. 在角色“hrole”所在行，单击“修改”。

角色名	描述	创建时间	操作
hrole	The role of Hive.	2021/09/29 18:59:49 GMT+08:...	修改 删除

- c. 选择“Hive > Hive Read Write Privileges”。



- d. 在数据库列表中单击指定的数据库名称“hdb”，显示数据库中的表“htable”。
- e. 在表“htable”的“权限”列，勾选“Select”、“Insert”和“Delete”。
- f. 单击“确定”完成。

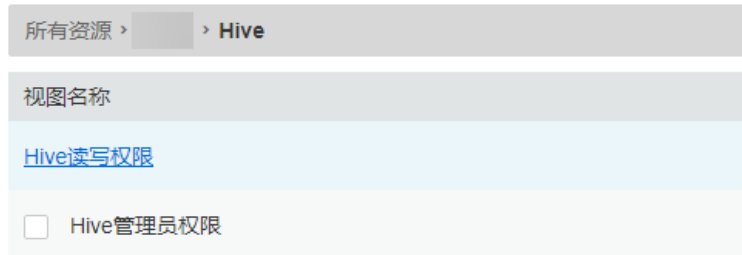
MRS3.x及后续版本，表授权的操作如下：

- 在FusionInsight Manager界面，选择“系统 > 权限 > 角色”。
- 在角色“hrole”所在行，单击“修改”。



- 选择“待操作的集群 > Hive > Hive读写权限”。

配置资源权限：



- 在数据库列表中单击指定的数据库名称“hdb”，显示数据库中的表“htable”。
- 在表“htable”的“权限”列，勾选“查询”、“插入”和“删除”。
- 单击“确定”完成。

📖 说明

在角色管理中，授予角色在Hive外表中查询、插入和删除数据的操作与Hive表相同，授予元数据权限将自动关联HDFS文件权限。

● 列授权

用户在Hive和HDFS中对自己创建的表拥有完整权限，用户没有权限访问别人创建的表。如果要访问别人创建的表上某些列，需要授予列权限。授予权限时只需要授予Hive元数据权限，HDFS文件权限将自动关联。以授予用户对应角色在表“htable”的列“hcol”中查询、插入数据的权限为例，操作步骤如下：

MRS 3.x之前版本，列授权的操作如下：

- 在MRS Manager界面，选择“系统设置 > 权限配置 > 角色管理”。
- 在角色“hrole”所在行，单击“修改”。
- 选择“Hive > Hive Read Write Privileges”。
- 在数据库列表中单击指定的数据库名称“hdb”，显示数据库中的表“htable”，单击表“htable”，显示表下的列“hcol”。
- 在列“hcol”的“权限”列，勾选“Select”和“Insert”。
- 单击“确定”完成。

MRS3.x及后续版本，列授权的操作如下：

- 在FusionInsight Manager界面，选择“系统 > 权限 > 角色”。
- 在角色“hrole”所在行，单击“修改”。
- 选择“待操作的集群 > Hive > Hive读写权限”。
- 在数据库列表中单击指定的数据库名称“hdb”，显示数据库中的表“htable”，单击表“htable”，显示表下的列“hcol”。
- 在列“hcol”的“权限”列，勾选“查询”和“插入”。
- 单击“确定”完成。

📖 说明

在权限管理中，授予元数据权限将自动关联HDFS文件权限，所以列授权后会增加表对应所有文件的HDFS ACL权限。

- 数据库授权

用户在Hive和HDFS中对自己创建的数据库拥有完整权限，用户访问别人创建的数据库，需要授予权限。授予权限时只需要授予Hive元数据权限，HDFS文件权限将自动关联。以授予用户对应角色在数据库“hdb”中查询和创建表的权限为例，操作步骤如下，不支持对角色授予数据库其他的操作权限：

MRS 3.x之前版本，数据库授权的操作如下：

- 在MRS Manager界面，选择“系统设置 > 权限配置 > 角色管理”。
- 在角色“hrole”所在行，单击“修改”。
- 选择“Hive > Hive Read Write Privileges”。
- 在数据库“hdb”的“权限”列，勾选“Select”和“Create”。
- 单击“确定”完成。

MRS3.x及后续版本，数据库授权的操作如下：

- 在FusionInsight Manager界面，选择“系统 > 权限 > 角色”。
- 在角色“hrole”所在行，单击“修改”。
- 选择“待操作的集群 > Hive > Hive读写权限”。
- 在数据库“hdb”的“权限”列，勾选“查询”和“建表”。
- 单击“确定”完成。

📖 说明

- 在权限管理中，为了方便用户使用，授予数据库下表的任意权限将自动关联该数据库目录的HDFS权限。为了避免产生性能问题，取消表的任意权限，系统不会自动取消数据库目录的HDFS权限，但对应的用户只能登录数据库和查看表名。
- 若为角色添加或删除数据库的查询权限，数据库中的表也将自动添加或删除查询权限。

相关概念

表 10-5 使用 Hive 表、列或数据库场景权限一览

操作场景	用户需要的权限
DESCRIBE TABLE	查询（Select）
SHOW PARTITIONS	查询（Select）
ANALYZE TABLE	查询（Select）、插入（Insert）
SHOW COLUMNS	查询（Select）
SHOW TABLE STATUS	查询（Select）
SHOW TABLE PROPERTIES	查询（Select）
SELECT	查询（Select）
EXPLAIN	查询（Select）

操作场景	用户需要的权限
CREATE VIEW	查询（ Select ）、Select授权（ Grant Of Select ）、建表（ Create ）
SHOW CREATE TABLE	查询（ Select ）、Select授权（ Grant Of Select ）
CREATE TABLE	建表（ Create ）
ALTER TABLE ADD PARTITION	插入（ Insert ）
INSERT	插入（ Insert ）
INSERT OVERWRITE	插入（ Insert ）、删除（ Delete ）
LOAD	插入（ Insert ）、删除（ Delete ）
ALTER TABLE DROP PARTITION	删除（ Delete ）
CREATE FUNCTION	Hive管理员权限（ Hive Admin Privilege ）
DROP FUNCTION	Hive管理员权限（ Hive Admin Privilege ）
ALTER DATABASE	Hive管理员权限（ Hive Admin Privilege ）

10.1.4 配置 Hive 业务使用其他组件的用户权限

操作场景

Hive业务还可能需关联使用其他组件，例如HQL语句触发MapReduce任务需要设置Yarn权限，或者Hive over HBase的场景需要HBase权限。以下介绍Hive关联Yarn和Hive over HBase两个场景下的操作。

说明

- 安全模式下Yarn和HBase的权限管理默认是开启的，因此在安全模式下默认需要配置Yarn和HBase权限。
- 在普通模式下，Yarn和HBase的权限管理默认是关闭的，即任何用户都有权限，因此普通模式下默认不需要配置Yarn和HBase权限。如果用户修改了YARN或者HBase的配置来开启权限管理，则修改后也需要配置Yarn和HBase权限。
- MRS 3.x及后续版本支持Ranger，如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加Hive的Ranger访问权限策略](#)。

前提条件

- 完成Hive客户端的安装。例如安装目录为“/opt/client”。
- 获取一个拥有管理员权限的用户，例如“admin”。

操作步骤

MRS 3.x之前版本，Hive关联Yarn

用户如果执行**insert**、**count**、**distinct**、**group by**、**order by**、**sort by**或**join**等语句时，将触发MapReduce任务，需要设置Yarn权限。以授予角色在表“thc”执行**count**语句的权限为例，操作步骤如下：

- 步骤1** 在MRS Manager角色界面创建一个角色。
- 步骤2** 在“权限”的表格中选择“Yarn > Scheduler Queue > root”。
- 步骤3** 在“default”队列的“权限”列，勾选“Submit”，单击“确定”保存。
- 步骤4** 在“权限”的表格中选择“Hive > Hive Read Write Privileges > default”，勾选表“thc”的“Select”，单击“确定”保存。

----结束

MRS 3.x及后续版本，Hive关联Yarn

用户如果执行**insert**、**count**、**distinct**、**group by**、**order by**、**sort by**或**join**等语句时，将触发MapReduce任务，需要设置Yarn权限。以授予角色在表“thc”执行**count**语句的权限为例，操作步骤如下：

- 步骤1** 在FusionInsight Manager角色界面创建一个角色。
- 步骤2** 在“配置资源权限”的表格中选择“待操作集群的名称 > Yarn > 调度队列 > root”。
- 步骤3** 在“default”队列的“权限”列，勾选“提交”，单击“确定”保存。
- 步骤4** 在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive读写权限 > default”，勾选表“thc”的“查询”，单击“确定”保存。

----结束

MRS 3.x之前版本，Hive over HBase授权

用户如果需要使用类似SQL语句的方式来操作HBase表，授予权限后可以在Hive中使用HQL命令访问HBase表。以授予用户在Hive中查询HBase表的权限为例，操作步骤如下

- 步骤1** 在MRS Manager角色管理界面创建一个HBase角色，例如“hive_hbase_create”，并授予创建HBase表的权限。
在“权限”的表格中选择“HBase > HBase Scope > global”，勾选命名空间“default”的“Create”，单击“确定”保存。
- 步骤2** 在MRS Manager用户管理界面创建一个“人机”用户，例如“hbase_creates_user”，加入“hive”组，绑定角色“hive_hbase_create”，用于创建Hive表和HBase表。
- 步骤3** 请根据客户端所在位置，参考[安装客户端](#)章节，登录安装客户端的节点。
- 步骤4** 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

- 步骤5** 执行以下命令，认证用户。

```
kinit hbase_creates_user
```

步骤6 执行以下命令，进入Hive客户端shell环境：

```
beeline
```

步骤7 执行以下命令，同时在Hive和HBase中创建表。例如创建表“thh”。

```
CREATE TABLE thh(id int, name string, country string) STORED BY  
'org.apache.hadoop.hive.hbase.HBaseStorageHandler' WITH  
SERDEPROPERTIES("hbase.columns.mapping" = "cf1:id,cf1:name,:key")  
TBLPROPERTIES ("hbase.table.name" = "thh");
```

创建好的Hive表和HBase表分别保存在Hive的数据库“default”和HBase的命名空间“default”。

步骤8 在MRS Manager角色管理界面创建一个角色，例如“hive_hbase_select”，并授予查询Hive表“thh”和HBase表“thh”的权限。

1. 在“权限”的表格中选择“HBase > HBase Scope > global > default”，勾选表“thh”的“read”，单击“确定”保存，授予HBase角色查询表的权限。
2. 编辑角色，在“权限”的表格中选择“HBase > HBase Scope > global > hbase”，勾选表“hbase:meta”的“Execute”，单击“确定”保存。
3. 编辑角色，在“权限”的表格中选择“Hive > Hive Read Write Privileges > default”，勾选表“thh”的“Select”，单击“确定”保存。

步骤9 在MRS Manager用户管理界面创建一个“人机”用户，例如“hbase_select_user”，加入“hive”组，绑定角色“hive_hbase_select”，用于查询Hive表和HBase表。

步骤10 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

步骤11 执行以下命令，认证用户。

```
kinit hbase_select_user
```

步骤12 执行以下命令，进入Hive客户端shell环境。

```
beeline
```

步骤13 执行以下命令，使用Hive的HQL语句查询HBase表的数据。

```
select * from thh;
```

----结束

MRS3.x及后续版本，Hive over HBase授权

用户如果需要使用类似SQL语句的方式来操作HBase表，授予权限后可以在Hive中使用HQL命令访问HBase表。以授予用户在Hive中查询HBase表的权限为例，操作步骤如下

步骤1 在FusionInsight Manager角色管理界面创建一个HBase角色，例如“hive_hbase_create”，并授予创建HBase表的权限。

在“配置资源权限”的表格中选择“待操作集群的名称 > HBase > HBase Scope > global”，勾选命名空间“default”的“创建”，单击“确定”保存。

步骤2 在FusionInsight Manager用户管理界面创建一个“人机”用户，例如“hbase_creates_user”，加入“hive”组，绑定角色“hive_hbase_create”，用于创建Hive表和HBase表。

步骤3 如果当前组件使用了Ranger进行权限控制，需给“hive_hbase_create”或“hbase_creates_user”配置“Create”权限，具体操作可参考[添加Hive的Ranger访问权限策略](#)。

步骤4 以客户端安装用户，登录安装客户端的节点。

步骤5 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

步骤6 执行以下命令，认证用户。

```
kinit hbase_creates_user
```

步骤7 执行以下命令，进入Hive客户端shell环境：

```
beeline
```

步骤8 执行以下命令，同时在Hive和HBase中创建表。例如创建表“thh”。

```
CREATE TABLE thh(id int, name string, country string) STORED BY  
'org.apache.hadoop.hive.hbase.HBaseStorageHandler' WITH  
SERDEPROPERTIES("hbase.columns.mapping" = "cf1:id,cf1:name,:key")  
TBLPROPERTIES ("hbase.table.name" = "thh");
```

创建好的Hive表和HBase表分别保存在Hive的数据库“default”和HBase的命名空间“default”。

步骤9 在FusionInsight Manager角色管理界面创建一个角色，例如“hive_hbase_select”，并授予查询Hive表“thh”和HBase表“thh”的权限。

1. 在“配置资源权限”的表格中选择“待操作集群的名称 > HBase > HBase Scope > global > default”，勾选表“thh”的“读”，单击“确定”保存，授予HBase角色查询表的权限。
2. 编辑角色，在“配置资源权限”的表格中选择“待操作集群的名称 > HBase > HBase Scope > global > hbase”，勾选表“hbase:meta”的“执行”，单击“确定”保存。
3. 编辑角色，在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive 读写权限 > default”，勾选表“thh”的“查询”，单击“确定”保存。

步骤10 在FusionInsight Manager用户管理界面创建一个“人机”用户，例如“hbase_select_user”，加入“hive”组，绑定角色“hive_hbase_select”，用于查询Hive表和HBase表。

步骤11 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

步骤12 执行以下命令，认证用户。

```
kinit hbase_select_user
```

步骤13 执行以下命令，进入Hive客户端shell环境。

```
beeline
```

步骤14 执行以下命令，使用Hive的HQL语句查询HBase表的数据。

```
select * from thh;
```

```
----结束
```

10.2 Hive 客户端使用实践

操作场景

该任务指导用户在运维场景或业务场景中使用Hive客户端。

前提条件

- 已安装客户端，例如安装目录为“/opt/hadoopclient”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- 各组件业务用户由MRS集群管理员根据业务需要创建。安全模式下，“机机”用户需要下载keytab文件。“人机”用户第一次登录时需修改密码。

使用 Hive 客户端（MRS 3.x 之前版本）

步骤1 安装客户端，具体请参考[安装客户端](#)章节。

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，切换到客户端安装目录。

```
cd /opt/hadoopclient
```

步骤4 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤5 根据集群认证模式，完成Hive客户端登录。

- 安全模式，则执行以下命令，完成用户认证并登录Hive客户端。

```
kinit 组件业务用户
```

```
beeline
```

- 普通模式，则执行以下命令，登录Hive客户端，如果不指定组件业务用户，则会以当前操作系统用户登录。

```
beeline -n 组件业务用户
```

说明

进行beeline连接后，可以编写并提交HQL语句执行相关任务。如需执行Catalog客户端命令，需要先执行!q命令退出beeline环境。

步骤6 使用以下命令，执行HCatalog的客户端命令。

```
hcat -e "cmd"
```

其中“cmd”必须为Hive DDL语句，如hcat -e "show tables"。

📖 说明

- 若要使用HCatalog客户端，必须从“组件管理”页面单击“下载客户端”，下载全部服务的客户端。Beeline客户端不受此限制。
- 由于权限模型不兼容，使用HCatalog客户端创建的表，在HiveServer客户端中不能访问，但可以使用WebHCat客户端访问。
- 在普通模式下使用HCatalog客户端，系统将以当前登录操作系统用户来执行DDL命令。
- 退出beeline客户端时请使用!**q**命令，不要使用“Ctrl + c”。否则会导致连接生成的临时文件无法删除，长期会累积产生大量的垃圾文件。
- 在使用beeline客户端时，如果需要在一行中输入多条语句，语句之间以“;”分隔，需要将“entireLineAsCommand”的值设置为“false”。

设置方法：如果未启动beeline，则执行**beeline --entireLineAsCommand=false**命令；如果已启动beeline，则在beeline中执行!**set entireLineAsCommand false**命令。

设置完成后，如果语句中含有不是表示语句结束的“;”，需要进行转义，例如**select concat_ws('\;', collect_set(col1)) from tbl**。

----结束

使用 Hive 客户端（MRS 3.x 及之后版本）

步骤1 安装客户端，具体请参考[安装客户端](#)章节。

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，切换到客户端安装目录。

```
cd /opt/hadoopclient
```

步骤4 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤5 根据集群认证模式，完成Hive客户端登录。

- 安全模式，则执行以下命令，完成用户认证并登录Hive客户端。

```
kinit 组件业务用户
```

```
beeline
```

- 普通模式，则执行以下命令，登录Hive客户端，如果不指定组件业务用户，则会以当前操作系统用户登录。

```
beeline -n组件业务用户
```

步骤6 使用以下命令，执行HCatalog的客户端命令。

```
hcat -e "cmd"
```

其中“cmd”必须为Hive DDL语句，如**hcat -e "show tables"**。

📖 说明

- 若要使用HCatalog客户端，必须从服务页面选择“更多 > 下载客户端”，下载全部服务的客户端。Beeline客户端不受此限制。
- 由于权限模型不兼容，使用HCatalog客户端创建的表，在HiveServer客户端中不能访问，但可以使用WebHCat客户端访问。
- 在普通模式下使用HCatalog客户端，系统将以当前登录操作系统用户来执行DDL命令。
- 退出beeline客户端时请使用!**q**命令，不要使用“Ctrl + C”。否则会导致连接生成的临时文件无法删除，长期会累积产生大量的垃圾文件。
- 在使用beeline客户端时，如果需要在一行中输入多条语句，语句之间以“;”分隔，需要将“entireLineAsCommand”的值设置为“false”。

设置方法：如果未启动beeline，则执行**beeline --entireLineAsCommand=false**命令；如果已启动beeline，则在beeline中执行!**set entireLineAsCommand false**命令。

设置完成后，如果语句中含有不是表示语句结束的“;”，需要进行转义，例如**select concat_ws('\;', collect_set(col1)) from tbl**。

----结束

Hive 客户端常用命令

常用的Hive Beeline客户端命令如下表所示。

更多命令可参考<https://cwiki.apache.org/confluence/display/Hive/HiveServer2+Clients#HiveServer2Clients-BeelineCommands>。

表 10-6 Hive Beeline 客户端常用命令

命令	说明
set <key>=<value>	设置特定配置变量（键）的值。 说明 若变量名拼错，Beeline不会显示错误。
set	打印由用户或Hive覆盖的配置变量列表。
set -v	打印Hadoop和Hive的所有配置变量。
add FILE[S] <filepath> <filepath>* add JAR[S] <filepath> <filepath>* add ARCHIVE[S] <filepath> <filepath>*	将一个或多个文件、JAR文件或ARCHIVE文件添加至分布式缓存的资源列表中。
add FILE[S] <ivyurl> <ivyurl>* add JAR[S] <ivyurl> <ivyurl>* add ARCHIVE[S] <ivyurl> <ivyurl>*	使用“ivy://goup:module:version?query_string”格式的Ivy URL，将一个或多个文件、JAR文件或ARCHIVE文件添加至分布式缓存的资源列表中。

命令	说明
list FILE[S] list JAR[S] list ARCHIVE[S]	列出已添加至分布式缓存中的资源。
list FILE[S] <filepath>* list JAR[S] <filepath>* list ARCHIVE[S] <filepath>*	检查给定的资源是否已添加至分布式缓存中。
delete FILE[S] <filepath>* delete JAR[S] <filepath>* delete ARCHIVE[S] <filepath>*	从分布式缓存中删除资源。
delete FILE[S] <ivyurl> <ivyurl>* delete JAR[S] <ivyurl> <ivyurl>* delete ARCHIVE[S] <ivyurl> <ivyurl>*	从分布式缓存中删除使用<ivyurl>添加的资源。
reload	使HiveServer2发现配置参数指定路径下JAR文件的变更“hive.reloadable.aux.jars.path”（无需重启HiveServer2）。更改操作包括添加、删除或更新JAR文件。
dfs <dfs command>	执行dfs命令。
<query string>	执行Hive查询，并将结果打印到标准输出。

10.3 快速使用 Hive 进行数据分析

Hive是基于Hadoop的一个数据仓库工具，可将结构化的数据文件映射成一张数据库表，并提供类SQL的功能对数据进行分析处理，通过类SQL语句快速实现简单的MapReduce统计，不必开发专门的MapReduce应用，十分适合数据仓库的统计分析。

背景信息

假定用户开发一个应用程序，用于管理企业中的使用A业务的用户信息，使用Hive客户端实现A业务操作流程如下：

普通表的操作：

- 创建用户信息表**user_info**。
- 在用户信息中新增用户的学历、职称信息。
- 根据用户编号查询用户姓名和地址。
- A业务结束后，删除用户信息表。

表 10-7 用户信息

编号	姓名	性别	年龄	地址
12005000201	A	男	19	A城市
12005000202	B	女	23	B城市
12005000203	C	男	26	C城市
12005000204	D	男	18	D城市
12005000205	E	女	21	E城市
12005000206	F	男	32	F城市
12005000207	G	女	29	G城市
12005000208	H	女	30	H城市
12005000209	I	男	26	I城市
12005000210	J	女	25	J城市

操作步骤

步骤1 下载客户端配置文件。

MRS 3.x之前版本，操作如下：

1. 登录MRS Manager页面，具体请参见[访问集群Manager](#)，然后选择“服务管理”。
2. 单击“下载客户端”。
“客户端类型”选择“仅配置文件”，“下载路径”选择“服务器端”，单击“确定”开始生成客户端配置文件，文件生成后默认保存在主管理节点“/tmp/MRS-client”。

图 10-1 仅下载客户端的配置文件

下载客户端

警告：生成客户端会占用大量的磁盘IO，不建议在集群处于安装中、启动中、打补丁中等非稳态场景进行“下载客户端”操作。

* 客户端类型 完整客户端 仅配置文件

* 下载路径 服务器端 远端主机

仅保存到服务器如下路径，如果存在客户端文件，会覆盖路径下已有的客户端文件。

* 客户端路径

确定

取消

MRS3.x及后续版本，操作如下：

1. 登录FusionInsight Manager页面，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群的名称 > 概览 > 更多 > 下载客户端”。
3. 下载集群客户端。

下载集群客户端

下载 的客户端，集群的客户端包括了所有服务

选择客户端类型：
 完整客户端 仅配置文件

选择平台类型：
 x86_64 aarch64

仅保存到如下路径：
 ?

“选择客户端类型”选择“仅配置文件”，选择平台类型，单击“确定”开始生成客户端配置文件，文件生成后默认保存在主管理节点“/tmp/FusionInsight-Client/”。

步骤2 登录Manager的主管理节点。

MRS 3.x之前版本，操作如下：

1. 在MRS控制台，选择“现有集群”，单击集群名称，在“节点管理”页签中查看节点名称，名称中包含“master1”的节点为Master1节点，名称中包含“master2”的节点为Master2节点。

MRS Manager的主备管理节点默认安装在集群Master节点上。在主备模式下，由于Master1和Master2之间会切换，Master1节点不一定是MRS Manager的主管理节点，需要在Master1节点中执行命令，确认MRS Manager的主管理节点。命令请参考[步骤2.4](#)。

2. 以root用户使用密码方式登录Master1节点。操作方法请参见[登录集群节点](#)章节。
3. 切换至omm用户。

```
sudo su - root
su - omm
```

4. 执行以下命令确认MRS Manager的主管理节点。

```
sh ${BIGDATA_HOME}/om-0.0.1/sbin/status-oms.sh
```

回显信息中“HAActive”参数值为“active”的节点为主管理节点（如下例中“mgtomsdat-sh-3-01-1”为主管理节点），参数值为“standby”的节点为备管理节点（如下例中“mgtomsdat-sh-3-01-2”为备管理节点）。

```
Ha mode
double
NodeName      HostName      HAVersion     StartTime     HAActive
HAAllResOK    HARunPhase
192-168-0-30  mgtomsdat-sh-3-01-1  V100R001C01  2014-11-18 23:43:02
```

```
active      normal      Activated
192-168-0-24 mgtomsdat-sh-3-01-2 V100R001C01 2014-11-21 07:14:02
standby     normal      Deactivated
```

5. 使用root用户登录Manager的主管理节点，例如“192-168-0-30”节点。

MRS3.x及后续版本，操作如下：

1. 以root用户登录任意部署Manager的节点。
2. 执行以下命令确认主备管理节点。

```
sh ${BIGDATA_HOME}/om-server/om/sbin/status-oms.sh
```

界面打印信息中“HAActive”参数值为“active”的节点为主管理节点（如下例中“node-master1”为主管理节点），参数值为“standby”的节点为备管理节点（如下例中“node-master2”为备管理节点）。

```
HAMode
double
NodeName      HostName      HAVersion      StartTime      HAActive
HAAllResOK    HARunPhase
192-168-0-30  node-master1 V100R001C01    2020-05-01 23:43:02 active
normal        Activated
192-168-0-24  node-master2 V100R001C01    2020-05-01 07:14:02 standby
normal        Deactivated
```

3. 以root用户登录主管理节点，并执行以下命令切换到omm用户。

```
sudo su - omm
```

步骤3 执行以下命令切换到客户端安装目录。

提前已安装集群客户端，以下客户端安装目录为举例，请根据实际情况修改。

```
cd /opt/client
```

步骤4 执行以下命令，更新主管理节点的客户端配置。

```
sh refreshConfig.sh /opt/client 客户端配置文件压缩包完整路径
```

例如，执行命令：

```
sh refreshConfig.sh /opt/client /tmp/FusionInsight-Client/
FusionInsight_Cluster_1_Services_Client.tar
```

界面显示以下信息表示配置刷新更新成功：

```
ReFresh components client config is complete.
Succeed to refresh components client config.
```

📖 说明

[步骤1~步骤4](#)的操作也可以参考[更新客户端](#)页面的方法二操作。

步骤5 在Master节点使用客户端。

1. 在已更新客户端的主管理节点，例如“192-168-0-30”节点，执行以下命令切换到客户端目录，客户端安装目录如：/opt/client。

```
cd /opt/client
```

2. 执行以下命令配置环境变量。

```
source bigdata_env
```

3. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，当前用户需要具有创建Hive表的权限，具体请参见[创建角色](#)配置拥有对应权限的角色，参考[创建用户](#)为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行此命令。

kinit *MRS集群用户*

例如，**kinit hiveuser**。

4. 直接执行Hive组件的客户端命令。

beeline

步骤6 运行Hive客户端命令，实现A业务。

内部表的操作：

1. 根据**表10-7**创建用户信息表user_info并添加相关数据，例如：

```
create table user_info(id string,name string,gender string,age int,addr string);
```

MRS 1.x和MRS3.x及后续版本，操作如下：

```
insert into table user_info(id,name,gender,age,addr) values("12005000201","A","男",19,"A城市");
```

MRS 2.x版本，操作如下：

```
insert into table user_info values("12005000201","A","男",19,"A城市");
```

2. 在用户信息表user_info中新增用户的学历、职称信息。
以增加编号为12005000201的用户的学历、职称信息为例，其他用户类似。

```
alter table user_info add columns(education string,technical string);
```

3. 根据用户编号查询用户姓名和地址。
以查询编号为12005000201的用户姓名和地址为例，其他用户类似。

```
select name,addr from user_info where id='12005000201';
```

4. 删除用户信息表。

```
drop table user_info;
```

外部分区表的操作：

创建外部分区表并导入数据：

1. 创建外部表数据存储路径：

```
hdfs dfs -mkdir /hive/
```

```
hdfs dfs -mkdir /hive/user_info
```

2. 建表：

```
create external table user_info(id string,name string,gender string,age int,addr string) partitioned by(year string) row format delimited fields terminated by ' ' lines terminated by '\n' stored as textfile location '/hive/user_info';
```

说明

fields terminated指明分隔的字符，如按空格分隔，' '。

lines terminated 指明分行的字符，如按换行分隔，'\n'。

/hive/user_info为数据文件的路径。

3. 导入数据。

- a. 使用insert语句插入数据。

```
insert into user_info partition(year="2018") values ("12005000201","A","男",19,"A城市");
```

- b. 使用load data命令导入文件数据。
 - i. 根据表10-7数据创建文件。如，文件名为txt.log，以空格拆分字段，以换行符作为行分隔符。
 - ii. 上传文件至hdfs。
hdfs dfs -put txt.log /tmp
 - iii. 加载数据到表中。
load data inpath '/tmp/txt.log' into table user_info partition (year='2011');
4. 查询导入数据。
select * from user_info;
5. 删除用户信息表。
drop table user_info;
6. 执行以下命令退出客户端。
!q
----结束

10.4 Hive 数据存储及加密配置

10.4.1 使用 HDFS Colocation 存储 Hive 表

操作场景

HDFS Colocation（同分布）是HDFS提供的数据分布控制功能，利用HDFS Colocation接口，可以将存在关联关系或者可能进行关联操作的数据存放在相同的存储节点上。Hive支持HDFS的Colocation功能，即在创建Hive表时，设置表文件分布的locator信息，当使用insert语句向该表中插入数据时会将该表的数据文件存放在相同的存储节点上（不支持其他数据导入方式），从而使后续的多表关联的数据计算更加方便和高效。表格式只支持TextFile和RCFile。

说明

本章节适用于MRS 3.x及后续版本。

操作步骤

- 步骤1** 使用客户端安装用户登录客户端所在节点。
- 步骤2** 执行以下命令，切换到客户端安装目录，如：/opt/client。
cd /opt/client
- 步骤3** 执行以下命令配置环境变量。
source bigdata_env
- 步骤4** 若集群为安全模式，执行以下命令认证用户。
kinit MRS用户名

步骤5 通过HDFS接口创建<groupid>

```
hdfs colocationadmin -createGroup -groupid <groupid> -locatorIds  
<locatorid1>,<locatorid2>,<locatorid3>
```

📖 说明

其中<groupid>为创建的group名称，该示例语句创建的group包含三个locator，用户可以根据需要定义locator的数量。

关于hdfs创建groupid，以及HDFS Colocation的详细介绍请参考hdfs的相关说明，这里不做赘述。

步骤6 执行以下命令进入Hive客户端：

```
beeline
```

步骤7 Hive使用colocation。

假设table_name1和table_name2是相关联的两张表，创建两表的语句如下：

```
CREATE TABLE <[db_name.]table_name1>[(col_name data_type , ...)] [ROW  
FORMAT <row_format>] [STORED AS <file_format>]  
TBLPROPERTIES("groupid"=" <group> ","locatorid"=" <locator1>");
```

```
CREATE TABLE <[db_name.]table_name2> [(col_name data_type , ...)] [ROW  
FORMAT <row_format>] [STORED AS <file_format>]  
TBLPROPERTIES("groupid"=" <group> ","locatorid"=" <locator1>");
```

当使用insert语句分别向table_name1和table_name2插入数据后，table_name1和table_name2的数据文件就会分布在hdfs的相同存储位置上，从而方便两表进行关联操作。

----结束

10.4.2 配置 Hive 分区元数据冷热存储

分区元数据冷热存储介绍

- 为了减轻元数据库压力，将长时间未使用过的指定范围的分区相关元数据移动到备份表，这一过程称为分区数据冻结，移动的分区数据称为冷分区，未冻结的分区称为热分区，存在冷分区的表称为冻结表。将被冻结的数据重新移回原元数据表，这一过程称为分区数据解冻。
- 一个分区从热分区变成冷分区，仅仅是在元数据中进行标识，其HDFS业务侧分区路径、数据文件内容并未发生变化。

📖 说明

本特性仅适用于MRS 3.1.2及之后版本。

冻结分区

支持创建表的用户按照条件过滤的方式对一个或多个分区进行冻结，格式为：freeze partitions 数据库名称.表名称 where 分区过滤条件

例如：

```
freeze partitions testdb.test where year <= 2021;
```

```
freeze partitions testdb.test where year<=2021 and month <= 5;
```

```
freeze partitions testdb.test where year<=2021 and month <= 5 and day <= 27;
```

解冻分区

支持创建表的用户按照条件过滤的方式对一个或多个分区进行解冻，格式为unfreeze partitions 数据库名称.表名称 where 分区过滤条件，如：

```
unfreeze partitions testdb.test where year <= 2021;
```

```
unfreeze partitions testdb.test where year<=2021 and month <= 5;
```

```
unfreeze partitions testdb.test where year<=2021 and month <= 5 and day <= 27;
```

查询含有冻结数据的表

- 查询当前数据库下的所有冻结表：
`show frozen tables;`
- 查询指定数据库下的所有冻结表：
`show frozen tables in 数据库名称;`

查询冻结表的冻结分区

查询冷分区：

```
show frozen partitions 表名;
```

📖 说明

- 默认元数据库冻结分区类型只支持int、string、varchar、date、timestamp类型。
- 外置元数据库只支持Postgres数据库，且冻结分区类型只支持int、string、varchar、timestamp类型。
- 对冻结后的表进行Msck元数据修复时，需要先解冻数据。如果对冻结表进行过备份后恢复操作，则可以直接执行Msck元数据修复操作，且解冻只能通过**msck repair**命令进行操作。
- 对冻结后的分区进行rename时，需要先解冻数据，否则会提示分区不存在。
- 删除存在冻结数据的表时，被冻结的数据会同步删除。
- 删除存在冻结数据的分区时，被冻结的分区信息不会被删除，HDFS业务数据也不会被删除。
- select查询数据时，会自动添加排查冷分区数据的过滤条件，查询结果将不包含冷分区的数据。
- show partitions table查询表下的分区数据时，查询结果将不包含冷分区，可通过show frozen partitions table进行冷分区查询。

10.4.3 Hive 支持 ZSTD 压缩格式

ZSTD（全称为Zstandard）是一种开源的无损数据压缩算法，其压缩性能和压缩比均优于当前Hadoop支持的其他压缩格式，本特性使得Hive支持ZSTD压缩格式的表。Hive支持基于ZSTD压缩的存储格式有常见的ORC，RCFile，TextFile，JsonFile，Parquet，Sequence，CSV。

📖 说明

本特性仅适用于MRS 3.1.2及之后版本。

ZSTD压缩格式的建表方式如下：

- ORC存储格式建表时可指定TBLPROPERTIES("orc.compress"="zstd")：
**create table tab_1(...) stored as orc
TBLPROPERTIES("orc.compress"="zstd");**
- Parquet存储格式建表可指定TBLPROPERTIES("parquet.compression"="zstd")：
**create table tab_2(...) stored as parquet
TBLPROPERTIES("parquet.compression"="zstd");**
- 其他格式或通用格式建表可执行设置参数指定compress,codec为“org.apache.hadoop.io.compress.ZStandardCode”：
**set hive.exec.compress.output=true;
set mapreduce.map.output.compress=true;
set
mapreduce.map.output.compress.codec=org.apache.hadoop.io.compress.Z
StandardCodec;
set mapreduce.output.fileoutputformat.compress=true;
set
mapreduce.output.fileoutputformat.compress.codec=org.apache.hadoop.i
o.compress.ZStandardCodec;
set hive.exec.compress.intermediate=true;
create table tab_3(...) stored as textfile;**

说明

ZSTD压缩格式的表和其他普通压缩表的SQL操作没有区别，可支持正常的增删查及聚合类SQL操作。

10.4.4 配置 Hive 列加密功能

操作场景

Hive支持对表的某一列或者多列进行加密；在创建Hive表时，可以指定要加密的列和加密算法。当使用insert语句向表中插入数据时，即可实现将对应列加密。列加密只支持存储在HDFS上的TextFile和SequenceFile文件格式的表。Hive列加密不支持视图以及Hive over HBase场景。

Hive列加密机制目前支持的加密算法有两种，在建表时指定：

- AES(对应加密类名称为：org.apache.hadoop.hive.serde2.AESRewriter)
- SMS4(对应加密类名称为：org.apache.hadoop.hive.serde2.SMS4Rewriter)

说明

将原始数据从普通Hive表导入到Hive列加密表后，在不影响其他业务情况下，建议删除普通Hive表上原始数据，因为保留一张未加密的表存在安全风险。

操作步骤

步骤1 在创建表时指定相应的加密列和加密算法：

```
create table <[db_name.]table_name> (<col_name1>  
<data_type>, <col_name2> <data_type>, <col_name3>
```



```
<data_type>,<col_name4> <data_type>) ROW FORMAT SERDE  
'org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe' WITH  
SERDEPROPERTIES ('column.encode.columns'='<col_name2>,<col_name3>',  
'column.encode.classname'='org.apache.hadoop.hive.serde2.AESRewriter')STO  
RED AS TEXTFILE;
```

或者使用如下语句：

```
create table <[db_name.]table_name> (<col_name1>  
<data_type>,<col_name2> <data_type>,<col_name3>  
<data_type>,<col_name4> <data_type>) ROW FORMAT SERDE  
'org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe' WITH  
SERDEPROPERTIES ('column.encode.indices'='1,2',  
'column.encode.classname'='org.apache.hadoop.hive.serde2.SMS4Rewriter')  
STORED AS TEXTFILE;
```

说明

- 使用序号指定加密列时，序号从0开始。0代表第1列，1代表第2列，依次类推。
- 创建列加密表时，表所在的目录必须是空目录。

步骤2 使用insert语法向设置列加密的表中导入数据。

假设test表已存在且有数据：

```
insert into table <table_name> select <col_list> from test;
```

----结束

10.5 Hive on HBase

10.5.1 配置跨集群互信下 Hive on HBase

两个开启Kerberos认证的互信集群中，使用Hive集群操作HBase集群，将目的端HBase集群的HBase关键配置项配置到源端Hive集群的HiveServer中。

前提条件

两个开启Kerberos认证的安全集群已完成跨集群互信配置。

跨集群配置 Hive on HBase

步骤1 下载HBase配置文件到本地，并解压。

1. 登录目的端HBase集群的FusionInsight Manager，选择“集群 > 服务 > HBase”。
2. 选择“更多 > 下载客户端”。

图 10-2 下载 HBase 客户端



3. 下载HBase配置文件，客户端类型选择仅配置文件。

图 10-3 下载 HBase 配置文件



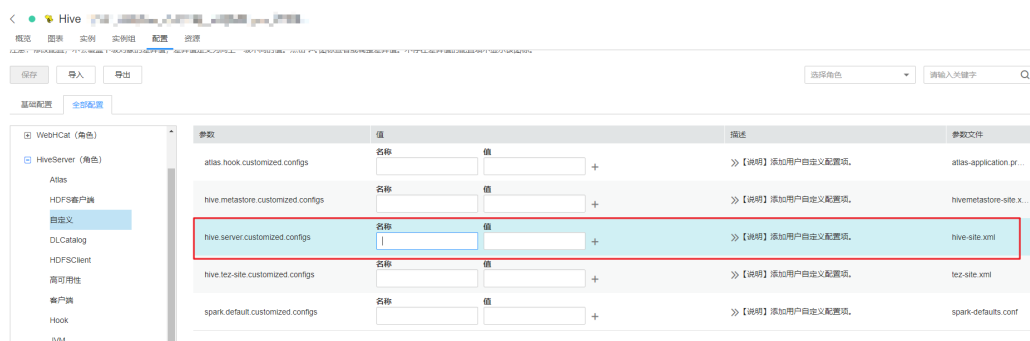
步骤2 登录源端Hive集群的FusionInsight Manager。

步骤3 选择“集群 > 服务 > Hive > 配置 > 全部配置”进入Hive服务配置页面，修改HiveServer角色的hive-site.xml自定义配置文件，增加HBase配置文件的如下配置项。

从已下载的HBase客户端配置文件的hbase-site.xml中，搜索并添加如下配置项及其取值到HiveServer中。

- hbase.security.authentication
- hbase.security.authorization
- hbase.zookeeper.property.clientPort
- hbase.zookeeper.quorum（域名需要转换为IP）
- hbase.regionserver.kerberos.principal
- hbase.master.kerberos.principal

图 10-4 HiveServer 角色的自定义配置



步骤4 保存配置并重启Hive服务。

----结束

10.5.2 删除 Hive on HBase 表中的单行记录

操作场景

由于底层存储系统的原因，Hive并不能支持对单条表数据进行删除操作，但在Hive on HBase功能中，MRS Hive提供了对HBase表的单条数据的删除功能，通过特定的语法，Hive可以将自己的HBase表中符合条件的一条或者多条数据清除。

表 10-8 删除 Hive on HBase 表中的单行记录所需权限

集群认证模式	用户所需权限
安全模式	“SELECT”、“INSERT”和“DELETE”
普通模式	无

操作步骤

步骤1 如果要删除某张HBase表中的某些数据，可以执行HQL语句：

remove table <table_name> where <expression>;

其中<expression>规定要删除数据的筛选条件；<table_name>为要删除数据的Hive on HBase表。

----结束

10.6 配置 Hive 读取关系型数据库数据

操作场景

Hive支持创建与其他关系型数据库关联的外表。该外表可以从关联到的关系型数据库中读取数据，并与Hive的其他表进行Join操作。

目前支持使用Hive读取数据的关系型数据库如下：

- DB2
- Oracle

📖 说明

本章节适用于MRS 3.x及后续版本。

前提条件

已安装Hive客户端。

操作步骤

步骤1 以Hive客户端安装用户登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录。

```
cd 客户端安装目录
```

例如安装目录为“/opt/client”，则执行以下命令：

```
cd /opt/client
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 集群认证模式是否为安全模式。

- 是，执行以下命令进行用户认证：

```
kinit Hive业务用户
```
- 否，执行**步骤5**。

步骤5 执行以下命令，将需要关联的关系型数据库驱动Jar包上传到HDFS目录下。

```
hdfs dfs -put Jar包所在目录 保存Jar包的HDFS目录
```

例如将“/opt”目录下ORACLE驱动Jar包上传到HDFS的“/tmp”目录下，则执行如下命令。

```
hdfs dfs -put /opt/ojdbc6.jar /tmp
```

步骤6 按照如下示例，在Hive客户端创建关联关系型数据库的外表。

📖 说明

如果是安全模式，建表的用户需要“ADMIN”权限，**ADD JAR**的路径请以实际路径为准。

```
-- 关联oracle linux6版本示例
-- 如果是安全模式，设置admin权限
set role admin;
-- 添加连接关系型数据库的驱动jar包,不同数据库有不同的驱动JAR
ADD JAR hdfs:///tmp/ojdbc6.jar;

CREATE EXTERNAL TABLE ora_test
-- hive表的列需比数据库返回结果多一列用于分页查询
(id STRING,rownum string)
STORED BY 'com.qubitproducts.hive.storage.jdbc.JdbcStorageHandler'
TBLPROPERTIES (
-- 关系型数据库类型
"qubit.sql.database.type" = "ORACLE",
-- 通过JDBC连接关系型数据库的url（不同数据库有不同的url格式）
"qubit.sql.jdbc.url" = "jdbc:oracle:thin:@//10.163.0.1:1521/mydb",
-- 关系型数据库驱动类名
"qubit.sql.jdbc.driver" = "oracle.jdbc.OracleDriver",
-- 在关系型数据库查询的sql语句,结果将返回hive表
"qubit.sql.query" = "select name from aaa",
-- hive表的列与关系型数据库表的列进行匹配（可忽略）
"qubit.sql.column.mapping" = "id=name",
-- 关系型数据库用户
"qubit.sql.dbcp.username" = "test",
-- 关系型数据库密码，命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的
history命令记录功能，避免信息泄露。
"qubit.sql.dbcp.password" = "xx");
```

----结束

10.7 Hive 企业级能力增强

10.7.1 配置 Hive 目录旧数据自动移除至回收站

操作场景

此功能适用于Hive组件。

开启此功能后，执行写目录：**insert overwrite directory "/path/" ...**，写成功之后，会将旧数据移除到回收站，并且同时限制该目录不能为Hive元数据库中已经存在的数据库路径。

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

📖 说明

- 若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。
- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 选择“HiveServer（角色）>自定义”，对参数文件“hive-site.xml”添加自定义参数，设置“名称”为“hive.override.directory.move.trash”，“值”为“true”，修改后重启所有Hive实例。

----结束

10.7.2 配置 Hive 插入数据到不存在的目录中

操作场景

此功能适用于Hive组件。

开启此功能后，在执行写目录：**insert overwrite directory** “/path1/path2/path3” ...时，其中“/path1/path2”目录权限为700且属主为当前用户，“path3”目录不存在，会自动创建“path3”目录，并写数据成功。

上述功能，在Hive参数“hive.server2.enable.doAs”为“true”时已经支持，本次增加当“hive.server2.enable.doAs”为“false”时的功能支持。

说明

本功能参数调整与[配置Hive目录旧数据自动移除至回收站](#)添加的自定义参数相同。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 选择“HiveServer（角色）>自定义”，对参数文件“hive-site.xml”添加自定义参数，设置“名称”为“hive.override.directory.move.trash”，“值”为“true”，修改后重启所有Hive实例。

----结束

10.7.3 配置创建 Hive 内部表时不能指定 Location

操作场景

此功能在MRS 3.x之前版本适用于Hive，Spark。在MRS3.x及后续版本适用于Hive，Spark2x。

开启此功能后，在创建Hive内部表时，不能指定location。即表创建成功之后，表的location路径会被创建在当前默认warehouse目录下，不能被指定到其他目录。如果创建内部表时指定location，则创建失败。

说明

开启本功能之后，创建Hive内部表不能执行location。因为对建表语句做了限制，如果数据库中已存在建表时指向非当前默认warehouse目录的表，在执行建库、表脚本迁移、重建元数据操作时需要特别注意，防止错误。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 选择“HiveServer（角色）> 自定义”，对参数文件“hive-site.xml”添加自定义参数，设置“名称”为“hive.internaltable.notallowlocation”，“值”为“true”，修改后重启所有Hive实例。

基础配置

全部配置

- Hive（服务）
- MetaStore（角色）
- WebHCat（角色）
- HiveServer（角色）

Atlas

HDFS客户端

自定义

DLCatalog

步骤3 是否需要在Spark/Spark2x客户端中启用此功能？

- 是，重新下载并安装Spark/Spark2x客户端。
- 否，操作结束。

----结束

10.7.4 配置用户在具有读和执行权限的目录中创建外表

操作场景

此功能在MRS 3.x之前版本适用于Hive，Spark。在MRS3.x及后续版本适用于Hive，Spark2x。

开启此功能后，允许有目录读权限和执行权限的用户和用户组创建外部表，而不必检查用户是否为该目录的属主，并且禁止外表的location目录在当前默认warehouse目录下。同时在外表授权时，禁止更改其location目录对应的权限。

📖 说明

开启本功能之后，外表功能变化大。请充分考虑实际应用场景，再决定是否做出调整。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 选择“HiveServer（角色）> 自定义”，对参数文件“hive-site.xml”添加自定义参数，设置“名称”为“hive.restrict.create.grant.external.table”，“值”为“true”。

步骤3 选择“MetaStore（角色）> 自定义”，对参数文件“hivemetastore-site.xml”添加自定义参数，设置“名称”为“hive.restrict.create.grant.external.table”，“值”为“true”，修改后重启所有Hive实例。

步骤4 是否需要在Spark/Spark2x客户端中启用此功能？

- 是，重新下载并安装Spark/Spark2x客户端。
- 否，操作结束。

----结束

10.7.5 配置基于 HTTPS/HTTP 协议的 REST 接口

操作场景

WebHCat为Hive提供了对外可用的REST接口，开源社区版本默认使用HTTP协议。

MRS Hive支持使用更安全的HTTPS协议，并且可以在两种协议间自由切换。

📖 说明

安全模式支持HTTPS和HTTP协议，普通模式只支持HTTP协议。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 修改Hive配置：

- MRS 3.x之前版本：在搜索框中输入参数名称，搜索“templeton.protocol.type”，修改参数值为HTTPS或者HTTP，修改后重启Hive服务即可使用对应的协议。
- MRS 3.x及后续版本：选择“WebHCat > 安全”，在该界面选择HTTPS或者HTTP，修改后重启Hive服务即可使用对应的协议。

基础配置

全部配置

Hive（服务）

MetaStore（角色）

WebHCat（角色）

基础配置

自定义

DLCatalog

客户端

JVM

日志

MetaDB

性能

安全

服务初始化

----结束

10.7.6 配置 Hive Transform 功能开关

操作场景

Hive开源社区版本禁用Transform功能。

MRS Hive提供配置开关，默认为禁用Transform功能，与开源社区版本保持一致。

用户可修改配置开关，开启Transform功能，当开启Transform功能时，存在一定的安全风险。

📖 说明

只有安全模式支持禁用Transform功能，普通模式不支持该功能。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。



步骤2 在搜索框中输入参数名称，搜索“hive.security.transform.disallow”，修改参数值为“true”或“false”，修改后重启所有HiveServer实例。

📖 说明

- 选择“true”时，禁用Transform功能，与开源社区版本保持一致。
- 选择“false”时，开启Transform功能，存在一定的安全风险。

----结束

10.7.7 切换 Hive 执行引擎为 Tez

操作场景

Hive支持使用Tez引擎处理数据计算任务，用户在执行任务前可手动切换执行引擎为Tez。

前提条件

集群已安装Yarn服务的TimelineServer角色，且角色运行正常。

客户端切换执行引擎为 Tez

步骤1 安装并登录Hive客户端，具体操作请参考[Hive客户端使用实践](#)。

步骤2 执行以下命令切换引擎并开启“yarn.timeline-service.enabled”参数：

```
set hive.execution.engine=tez;
```

```
set yarn.timeline-service.enabled=true;
```

📖 说明

- “yarn.timeline-service.enabled”参数开启后可以在Tez服务中通过TezUI查看Tez引擎执行任务的详细情况。开启后任务信息将上报TimelineServer，如果TimelineServer实例故障，会导致任务失败。
- 由于Tez使用ApplicationMaster缓冲池，“yarn.timeline-service.enabled”必须在提交Tez任务前开启，否则会导致此参数无法生效，需要重新登录客户端进行配置。
- 当执行引擎需要切换为其它引擎时，需要通过客户端执行set yarn.timeline-service.enabled=false命令关闭“yarn.timeline-service.enabled”参数。
- 如果需要指定Yarn运行队列，可以在客户端执行set tez.queue.name=default命令指定运行队列。

步骤3 提交并执行Tez任务。

步骤4 登录FusionInsight Manager界面，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)，选择“集群 > 待操作的集群 > 服务 > Tez > TezUI（主机名称）”，在TezUI界面查看任务执行情况。

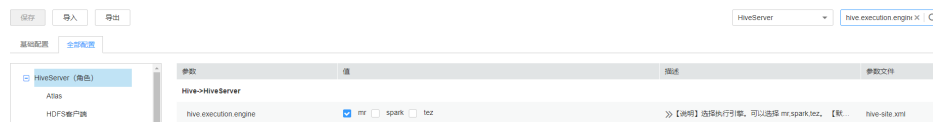


针对MRS 3.x之前版本，请登录MRS Manager界面，选择“服务管理 > Tez > Tez WebUI”，在TezUI界面查看任务执行情况。

----结束

切换 Hive 服务默认执行引擎为 Tez

步骤1 登录FusionInsight Manager界面，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)，选择“集群 > 待操作的集群 > 服务 > Hive > 配置 > 全部配置 > HiveServer（角色）”，搜索“hive.execution.engine”参数。



针对MRS 3.x之前版本，请登录MRS Manager界面，选择“服务管理 > Hive > 服务配置 > 全部配置 > HiveServer”，搜索“hive.execution.engine”参数。

步骤2 将“hive.execution.engine”参数设置为“tez”。

步骤3 选择“Hive（服务）> 自定义”，搜索“yarn.site.customized.configs”。

步骤4 在“yarn.site.customized.configs”参数后添加自定义参数，名称为“yarn.timeline-service.enabled”，值为“true”。

说明

- “yarn.timeline-service.enabled”开启后可以在Tez服务中通过TezUI查看Tez引擎执行任务详细情况。开启后任务信息将上报TimelineServer，如果TimelineServer实例故障，会导致任务失败。
- 由于Tez使用ApplicationMaster缓冲池，“yarn.timeline-service.enabled”必须在提交Tez任务前开启，否则会导致此参数无法生效，需要重新登录客户端配置。
- 当执行引擎需要切换为其它引擎时，需要将自定义参数“yarn.timeline-service.enabled”的值设置为“false”。

步骤5 单击“保存”在弹出窗口单击“确定”。

针对MRS 3.x之前版本，请单击“保存配置”在弹出窗口单击“是”。

步骤6 选择“概览 > 更多 > 重启服务”，重启Hive服务，输入密码开始重启服务。



针对MRS 3.x之前版本，请在“服务状态”页签选择“更多 > 重启服务”，重启Hive服务。

步骤7 安装并登录Hive客户端，具体操作请参考[Hive客户端使用实践](#)。

步骤8 提交并执行Tez任务。

步骤9 登录FusionInsight Manager界面，选择“集群 > 待操作的集群 > 服务 > Tez > TezUI（主机名称）”，跳转TezUI界面查看任务执行情况。

针对MRS 3.x之前版本，请登录MRS Manager界面，选择“服务管理 > Tez > Tez WebUI”，在TezUI界面查看任务执行情况。

----结束

10.7.8 Hive 负载均衡

10.7.8.1 配置 Hive 任务的最大 map 数

操作场景

- 此功能适用于Hive。
- 此功能用于从服务端限定Hive任务的最大map数，避免HiveServer服务过载而引发的性能问题。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 选择“MetaStore（角色）> 自定义”，对参数文件“hivemetastore-site.xml”添加自定义参数，设置“名称”为“hive.mapreduce.per.task.max.splits”，“值”为具体设定值，一般尽量设置大，修改后重启所有Hive实例。



----结束

10.7.8.2 配置用户租约隔离访问指定节点的 HiveServer

操作场景

- 此功能适用于Hive。
- 开启此功能可以限定指定用户访问指定节点上的HiveServer服务，实现对用户访问HiveServer服务的资源隔离。

📖 说明

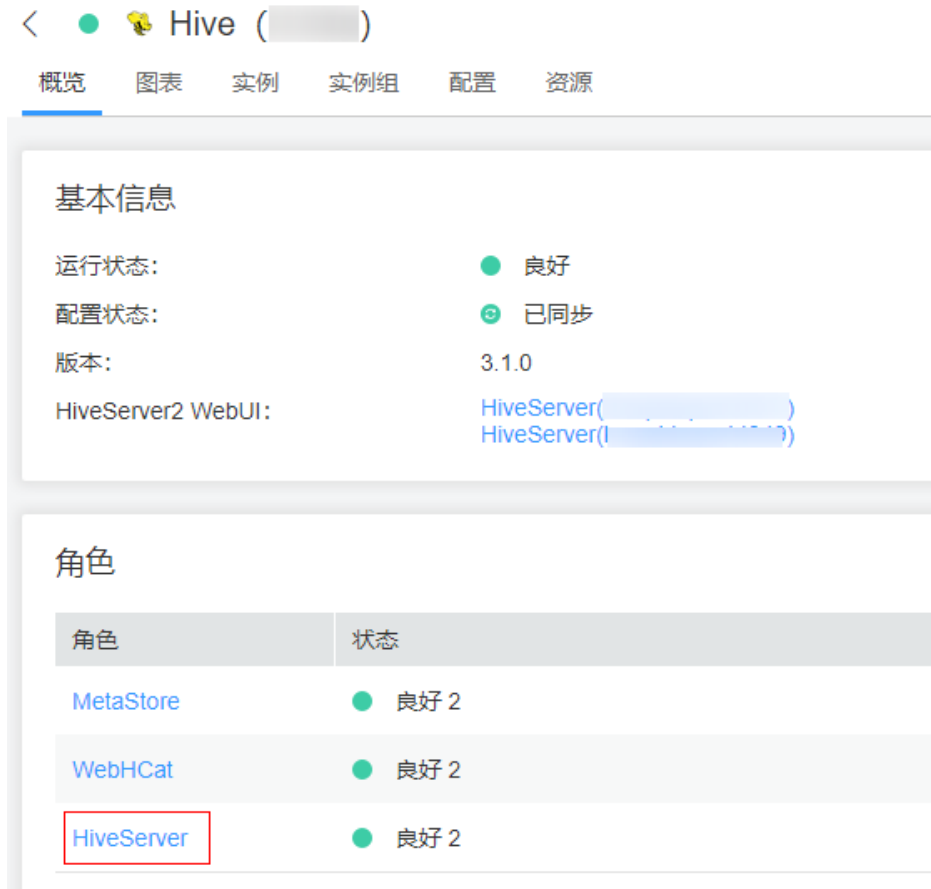
本章节适用于MRS 3.x及后续版本。

操作步骤

以对用户hiveuser设置租约隔离为例，选取Hive当前已有的或者新添加一个或者多个实例，此处选择已有的HiveServer实例：

步骤1 登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。

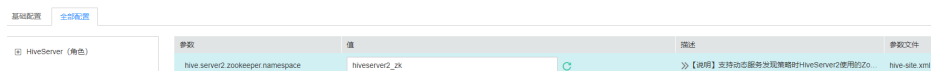
步骤2 选择“集群 > 待操作集群的名称 > 服务 > Hive > HiveServer”。



步骤3 在HiveServer列表里选择设置租约隔离的HiveServer，选择“HiveServer > 实例配置 > 全部配置”。



步骤4 在“全部配置”界面的右上角搜索“hive.server2.zookeeper.namespace”，“值”为具体设定值，比如为hiveserver2_zk。



步骤5 单击“保存”，在弹出对话框单击“确定”。

步骤6 选择“集群 > 待操作集群的名称 > 服务 > Hive”，选择“更多 > 重启服务”，输入密码开始重启服务。

步骤7 使用beeline -u 的方式登录客户端，执行以下命令：

```
beeline -u
"jdbc:hive2://10.5.159.13:2181/;serviceDiscoveryMode=zooKeeper;zooKeeperNameSpace=hiveserver2_zk;sasl.qop=auth-conf;auth=KERBEROS;principal=hive/hadoop.<系统域名>@<系统域名>"
```

执行命令时将“10.5.159.13”替换为任意一个ZooKeeper实例的IP地址，查找方式为“集群 > 待操作集群的名称 > 服务 > ZooKeeper > 实例”。

“zooKeeperNameSpace=”后面的“hiveserver2_zk”为**步骤4**中参数“hive.server2.zookeeper.namespace”设置的具体设定值。

结果将只会登录到被设置租约隔离的HiveServer。

📖 说明

- 开启本功能后，必须在登录时使用以上命令才可以访问这个被设置租约隔离的HiveServer。如果直接使用beeline命令登录客户端，将只会访问其他没有被设置租约隔离的HiveServer。
- 用户可登录FusionInsight Manager，选择“系统 > 权限 > 域和互信”，查看“本端域”参数，即为当前系统域名。“hive/hadoop.<系统域名>”为用户名，用户名所包含的系统域名所有字母为小写。

---结束

10.7.9 配置 Hive 单表动态视图的访问控制权限

操作场景

MRS中安全模式下Hive可以创建一个视图并控制用户访问权限，支持授权给不同的用户访问，又可以限定不同用户只能访问的不同数据。

在视图中，Hive可以通过获取当前客户端提交任务的用户的内置函数“current_user()”来进行过滤，这样被授权的用户，在访问视图时，即可被限定访问对应的数据。

📖 说明

- 在普通模式下“current_user()”函数无法区别客户端提交任务的用户，因此，当前访问控制仅对安全模式下的Hive有效。
- 如果已经在实际业务逻辑中使用了“current_user()”函数，那么，在安全模式与普通模式互转时，需要充分评估可能的风险。

操作示例

- 不采用“current_user”函数，要实现不同的用户，访问不同数据，需要创建不同的视图：
 - 将视图v1授权给用户hiveuser1，hiveuser1用户可以访问表table1中“type='hiveuser1'”的数据：

```
create view v1 as select * from table1 where type='hiveuser1';
```
 - 将视图v2授权给用户hiveuser2，hiveuser2用户可以访问表table1中“type='hiveuser2'”的数据：

```
create view v2 as select * from table1 where type='hiveuser2';
```

- 采用“current_user”函数，则只需要创建一个视图：
将视图v分别赋给用户hiveuser1、hiveuser2，当hiveuser1查询视图v时，“current_user()”被自动转化为hiveuser1，当hiveuser2查询视图v时，“current_user()”被自动转化为hiveuser2：

```
create view v as select * from table1 where type=current_user();
```

10.7.10 配置创建临时函数的用户不需要具有 ADMIN 权限

操作场景

Hive开源社区版本创建临时函数需要用户具备ADMIN权限。

MRS Hive提供配置开关，默认为创建临时函数需要ADMIN权限，与开源社区版本保持一致。

用户可修改配置开关，实现创建临时函数不需要ADMIN权限。当该选项配置成false时，存在一定的安全风险。

说明

安全模式支持配置创建临时函数是否需要ADMIN权限功能，而普通模式不支持该功能。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 在搜索框中输入参数名称，搜索“hive.security.temporary.function.need.admin”，修改参数值为“true”或“false”，修改后重启所有HiveServer实例。

说明

- 选择“true”时，创建临时函数需要ADMIN权限，与开源社区版本保持一致。
- 选择“false”时，创建临时函数不需要ADMIN权限。

----结束

10.7.11 配置具备表 select 权限的用户可查看表结构

操作场景

此功能在MRS3.x及后续版本适用于Hive，Spark2x。

开启此功能后，使用Hive建表时，其他用户被授予select权限后，可通过show create table查看表结构。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 选择“HiveServer（角色）> 自定义”，对参数文件“hive-site.xml”添加自定义参数，设置“名称”为“hive.allow.show.create.table.in.select.nogrant”，“值”为“true”，修改后重启所有Hive实例。

步骤3 是否需要在Spark/Spark2x客户端中启用此功能？

- 是，重新下载并安装Spark/Spark2x客户端。
- 否，操作结束。

----结束

10.7.12 配置仅 Hive 管理员用户能创建库和在 default 库建表

操作场景

此功能在MRS 3.x之前版本适用于Hive，Spark。在MRS3.x及后续版本适用于Hive，Spark2x。

开启此功能后，仅有Hive管理员可以创建库和在default库中建表，其他用户需通过Hive管理员授权才可使用库。

📖 说明

- 开启本功能之后，会限制普通用户新建库和在default库新建表。请充分考虑实际应用场景，再决定是否做出调整。
- 因为对执行用户做了限制，使用非管理员用户执行建库、表脚本迁移、重建元数据操作时需要特别注意，防止错误。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 选择“HiveServer（角色）>自定义”，对参数文件“hive-site.xml”添加自定义参数，设置“名称”为“hive.allow.only.admin.create”，“值”为“true”，修改后重启所有Hive实例。

步骤3 是否需要在Spark/Spark2x客户端中启用此功能？

- 是，执行**步骤4**。
- 否，操作结束。

步骤4 选择“SparkResource2x > 自定义”和“JDBCServer2x > 自定义”，对参数文件“hive-site.xml”添加自定义参数，设置“名称”为“hive.allow.only.admin.create”，“值”为“true”，修改后重启所有Spark2x实例。

步骤5 重新下载并安装Spark/Spark2x客户端。

----结束

10.7.13 配置 Hive 支持创建超过 32 个角色

操作场景

此功能适用于Hive。

因为操作系统用户组个数限制，导致Hive不能创建超过32个角色，开启此功能后，Hive将支持创建超过32个角色。

说明

- 开启本功能并对表库等授权后，对表库目录具有相同权限的角色将会用“|”合并。查询acl权限时，将显示合并后的结果，与开启该功能前的显示会有区别。此操作不可逆，请充分考虑实际应用场景，再决定是否做出调整。
- MRS3.x及后续版本支持Ranger，如果当前组件使用了Ranger进行权限控制，需基于Ranger配置相关策略进行权限管理，具体操作可参考[添加Hive的Ranger访问权限策略](#)。
- 开启此功能后，包括owner在内默认最大可支持512个角色，由MetaStore自定义参数“hive.supports.roles.max”控制，可考虑实际应用场景进行修改。

操作步骤

步骤1 进入Hive服务配置页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Hive > 配置 > 全部配置”。

步骤2 选择“MetaStore（角色）>自定义”，对参数文件“hivemetastore-site.xml”添加自定义参数，设置“名称”为“hive.supports.over.32.roles”，“值”为“true”。

步骤3 选择“HiveServer（角色）>自定义”，对参数文件“hive-site.xml”添加自定义参数，设置“名称”为“hive.supports.over.32.roles”，“值”为“true”，修改后重启所有Hive实例。

----结束

10.7.14 创建 Hive 用户自定义函数

当Hive的内置函数不能满足需要时，可以通过编写用户自定义函数UDF（User-Defined Functions）插入自己的处理代码并在查询中使用它们。

按实现方式，UDF分如下分类：

- 普通的UDF，用于操作单个数据行，且产生一个数据行作为输出。
- 用户定义聚集函数UDAF（User-Defined Aggregating Functions），用于接受多个输入数据行，并产生一个输出数据行。
- 用户定义表生成函数UDTF（User-Defined Table-Generating Functions），用于操作单个输入行，产生多个输出行。

按使用方法，UDF有如下分类：

- 临时函数，只能在当前会话使用，重启会话后需要重新创建。
- 永久函数，可以在多个会话中使用，不需要每次创建。

📖 说明

- 用户自定义函数需要用户控制函数中变量的内存、线程等资源的占用，如果控制不当可能会导致内存溢出、CPU使用高等问题。
- 若集群开启了Ranger鉴权，需要关闭Ranger鉴权后才能使用Python的UDF函数。

下面以编写一个AddDoublesUDF为例，说明UDF的编写和使用方法。

功能介绍

AddDoublesUDF主要用来对两个及多个浮点数进行相加，在该样例中可以掌握如何编写和使用UDF。

📖 说明

- 一个普通UDF必须继承自“org.apache.hadoop.hive.ql.exec.UDF”。
- 一个普通UDF必须至少实现一个evaluate()方法，evaluate函数支持重载。
- 开发自定义函数需要在工程中添加“hive-exec-*.jar”依赖包，可从Hive服务的安装目录下获取，例如在“\${BIGDATA_HOME}/components/FusionInsight_HD_*/Hive/disaster/plugin/lib/”目录下获取。

样例代码

以下为UDF示例代码：

其中，xxx通常为程序开发的组织名称。

```
package com.xxx.bigdata.hive.example.udf;
import org.apache.hadoop.hive.ql.exec.UDF;

public class AddDoublesUDF extends UDF {
    public Double evaluate(Double... a) {
        Double total = 0.0;
    }
}
```

```
// 处理逻辑部分.  
for (int i = 0; i < a.length; i++)  
    if (a[i] != null)  
        total += a[i];  
return total;  
}  
}
```

创建 Hive 用户自定义函数

步骤1 准备执行函数的用户。

1. 使用admin用户登录Manager界面，选择“集群 > 集群属性”，查看集群的“认证模式”并记录。
2. 选择“集群 > 服务 > Hive”，单击页面右上角的“更多”查看Hive是否启用Ranger鉴权。
3. 选择“系统 > 权限 > 用户”，单击“添加用户”，配置以下参数并单击“确定”，创建执行自定义函数的用户：
 - 用户名：填写用户名称，例如：test。
 - 用户类型：选择“人机”用户。
 - “密码”和“确认新密码”输入该用户对应的密码。
 - 用户组：单击“添加”，选择“hive”和“hadoop”用户组并单击“确定”。
4. 根据集群的认证模式及是否启用Ranger鉴权为新创建的用户赋权：
 - 集群的“认证模式”为“安全模式”：
 - “启用Ranger鉴权”按钮置灰（Hive已启用Ranger鉴权），执行[步骤 1.5](#)。
 - “停用Ranger鉴权”按钮置灰（Hive未启用Ranger鉴权），执行[步骤 1.6](#)。
 - 集群的“认证模式”为“普通模式”：
 - “启用Ranger鉴权”按钮置灰（Hive已启用Ranger鉴权），执行[步骤 1.5](#)。
 - “停用Ranger鉴权”按钮置灰（Hive未启用Ranger鉴权），执行[步骤 2](#)。
5. Hive使用Ranger鉴权，需使用Ranger管理员用户（安全模式集群为rangeradmin，普通模式集群为admin）登录Ranger管理界面为用户添加Hive权限控制策略。
 - a. 选择“集群 > 服务 > Ranger”，单击“Ranger Web UI”右侧的超链接登录Ranger WebUI页面。
 - b. 安全模式集群需单击页面右上角的用户名，在下拉框中单击“Log Out”退出当前用户，使用rangeradmin用户重新登录Ranger管理界面。
 - c. 在首页中单击“HADOOP SQL”区域的组件插件名称如“Hive”。
 - d. 在“Access”页签单击“Add New Policy”，配置以下参数并单击“Add”：
 - Policy Name：设置策略名称，例如：test_hive。

- database:
 - 永久函数：配置要添加函数的数据库名称，例如：default。
 - 临时函数：将“database”切换为“global”，并配置具体的函数名或设置为*。
 - table：切换为“udf”，并配置具体的函数名或设置为*。临时函数无需配置该参数。
 - 在“Allow Conditions”区域的“Select User”列选择新增的用户，在“Permissions”新增以下权限：
 - 永久函数：根据实际业务需求进行授权，例如，可新增“create”、“select”和“drop”权限。
 - 临时函数：添加“Temporary UDF Admin”权限。
6. Hive使用Manager角色鉴权，需创建具有Hive管理员权限的用户才能执行永久函数和临时函数。
- a. 在Manager页面首页，选择“系统 > 权限 > 角色”，单击“添加角色”，配置以下参数并单击“确定”：
 - 角色名称：填写角色名称，例如：test_role。
 - 配置资源权限：在“配置资源权限”列表中单击“待操作的集群名称 > Hive”，勾选“Hive管理员权限”。
 - b. 单击“用户”，单击[步骤1.3](#)新创建的用户所在行的“修改”。
 - c. 在修改用户页面，单击“角色”右侧的添加，添加新创建的具有Hive管理员权限的角色，单击“确定”。

步骤2 把以上程序打包成AddDoublesUDF.jar，并上传至客户端安装节点，例如“opt”目录下，再上传到HDFS指定目录下（例如“/user/hive_examples_jars”）。创建函数的用户与使用函数的用户都需要具有该文件的可读权限。

1. 切换至客户端安装目录并配置环境变量：

```
cd 客户端安装目录
source bigdata_env
```
2. 认证用户。
 - 集群已开启Kerberos认证（安全模式）：

```
kinit 业务用户
```
 - 集群未开启Kerberos认证（普通模式）：

```
export HADOOP_USER_NAME=业务用户
```
3. 上传UDF Jar包至HDFS目录中：

```
hdfs dfs -put /opt /user/hive_examples_jars
hdfs dfs -chmod 777 /user/hive_examples_jars
```

步骤3 登录Hive客户端。

- 集群已开启Kerberos认证（安全模式），执行如下命令：

```
beeline
```


📖 说明

如果用户绑定了Hive管理员角色，在每个beeline的维护操作会话中，都需要执行以下命令切换成admin角色再执行后续操作：

```
set role admin;
```

- 集群未开启Kerberos认证（普通模式），执行如下命令：

```
beeline -n Hive业务用户
```

步骤4 在Hive Server中执行以下命令定义该函数：

- 创建永久函数：

```
CREATE FUNCTION addDoubles AS  
'com.xxx.bigdata.hive.example.udf.AddDoublesUDF' using jar 'hdfs://  
hacluster/user/hive_examples_jars/AddDoublesUDF.jar';
```

其中*addDoubles*是该函数的别名，用于SELECT查询中使用；xxx通常为程序开发的组织名称。

- 创建临时函数：

```
CREATE TEMPORARY FUNCTION addDoubles AS  
'com.xxx.bigdata.hive.example.udf.AddDoublesUDF' using jar 'hdfs://  
hacluster/user/hive_examples_jars/AddDoublesUDF.jar';
```

- *addDoubles*是该函数的别名，用于SELECT查询中使用。

- 关键字TEMPORARY说明该函数只在当前这个Hive Server的会话过程中定义使用。

步骤5 在Hive Server中执行以下命令使用该函数：

```
SELECT addDoubles(1,2,3);
```

📖 说明

若重新连接客户端再使用函数出现[Error 10011]的错误，可执行**reload function;**命令后再使用该函数。

步骤6 在Hive Server中执行以下命令删除该函数：

```
DROP FUNCTION addDoubles;
```

----结束

扩展应用

无

10.7.15 配置 Hive Beeline 高可靠性

操作场景

- 在批处理任务运行过程中，beeline客户端由于网络异常等问题断线时，Hive能支持beeline在断线前已经提交的任务继续运行。当再次运行该批处理任务时，已经提交过的任务不再重新执行，直接从下一个任务开始执行。
- 在批处理任务运行过程中，HiveServer服务由于某些原因导致宕机时，Hive能支持当再次运行该批处理任务时，已经成功执行完成的任务不再重新执行，直接从HiveServer2宕机时正在运行的任务开始运行。

📖 说明

本章节适用于MRS 3.x及后续版本。

操作示例

1. beeline启动断线重连功能。

示例：

```
beeline -e "${SQL}" --hivevar batchid=xxxxx
```

2. beeline kill正在运行的任务。

示例：

```
beeline -e "" --hivevar batchid=xxxxx --hivevar kill=true
```

3. 登录beeline客户端，启动断线重连机制。

登录beeline客户端后，执行“set hivevar:batchid=xxxx”

📖 说明

使用说明：

- 其中“xxxx”表示每一次通过beeline提交任务的批次号，通过该批次号，可以识别出先提交的任务。如果提交任务时不带批次号，该特性功能不会启用。“xxxx”的值是执行任务时指定的，如下所示，“xxxx”值为“012345678901”：

```
beeline -f hdfs://hacluster/user/hive/table.sql --hivevar batchid=012345678901
```

- 如果运行的SQL脚本依赖数据的失效性，建议不启用断点重连机制，或者每次运行时使用新的batchid。因为重复执行时，可能由于某些SQL语句已经执行过了不再重新执行，导致获取到过期的数据。
- 如果SQL脚本中使用了一些内置时间函数，建议不启用断点重连机制，或者每次运行时使用新的batchid，理由同上。
- 一个SQL脚本里面会包含一个或多个子任务。如果SQL脚本中存在先创建再删除临时表的逻辑，建议将删除临时表的逻辑放到脚本的最后。假定删除临时表子任务的后续子任务执行失败，并且删除临时表的子任务之前的子任务用到了该临时表；当下一次以相同batchid执行该SQL脚本时，因为临时表在上一次执行时已被删除，则会导致删除临时表的子任务之前用到该临时表的子任务（不包括创建该临时表的子任务，因为上一次已经执行成功，本次不会再执行，仅可编译）编译失败。这种情况下，建议使用新的batchid执行脚本。

参数说明：

- zk.cleanup.finished.job.interval：执行清理任务的间隔时间，默认隔60s执行一次。
- zk.cleanup.finished.job.outdated.threshold：节点的过期时间，每个批次的任务都会生成对应节点，从当前批次任务的结束时间开始算，如果超过60分钟，则表示已经过期了，那么就清除节点。
- batch.job.max.retry.count：单批次任务的最大重试次数，当单批次的任务失败重试次数超过这个值，就会删除该任务记录，下次运行时将从头开始运行，默认是10次。
- beeline.reconnect.zk.path：存储任务执行进度的根节点，Hive服务默认是/beeline。

10.8 Hive 性能调优

10.8.1 建立 Hive 表分区提升查询效率

操作场景

Hive在做Select查询时，一般会扫描整个表内容，会消耗较多时间去扫描不关注的数
据。此时，可根据业务需求及其查询维度，建立合理的表分区，从而提高查询效率。

操作步骤

步骤1 MRS 3.x之前版本：

登录MRS控制台，在左侧导航栏选择“现有集群”，单击集群名称。选择“节点管理
> 节点名称”，进入弹性云服务器界面。单击“远程登录”按钮，完成Hive节点的登
录。

MRS3.x及后续版本：

以root用户登录已安装Hive客户端的节点。

步骤2 执行以下命令，进入客户端安装目录，例如“/opt/client”。

```
cd /opt/client
```

步骤3 执行source bigdata_env命令，配置客户端环境变量。

步骤4 在客户端中执行如下命令，执行登录操作。

```
kinit 用户名
```

步骤5 执行以下命令登录客户端工具。

```
beeline
```

步骤6 指定静态分区或者动态分区。

- 静态分区：

静态分区是手动输入分区名称，在创建表时使用关键字**PARTITIONED BY**指定分
区列名及数据类型。应用开发时，使用**ALTER TABLE ADD PARTITION**语句增加
分区，以及使用**LOAD DATA INTO PARTITION**语句将数据加载到分区时，只能静
态分区。

- 动态分区：通过查询命令，将结果插入到某个表的分区时，可以使用动态分区。
动态分区通过在客户端工具执行如下命令来开启：

```
set hive.exec.dynamic.partition=true;
```

动态分区默认模式是strict，也就是必须至少指定一列为静态分区，在静态分区下
建立动态子分区，可以通过如下设置来开启完全的动态分区：

```
set hive.exec.dynamic.partition.mode=nonstrict;
```

📖 说明

- 动态分区可能导致一个DML语句创建大量的分区，对应创建大量新文件夹，对系统性能可能带来影响。
- 在文件数量大的情况下，执行一个SQL语句启动时间较长，可以在执行SQL语句之前执行“set mapreduce.input.fileinputformat.list-status.num-threads = 100;”命令来缩短启动时间。“mapreduce.input.fileinputformat.list-status.num-threads”参数需要先添加到Hive的白名单才可设置。

----结束

10.8.2 Hive Join 数据优化

操作场景

使用Join语句时，如果数据量大，可能造成命令执行速度和查询速度慢，此时可进行Join优化。

Join优化可分为以下方式：

- Map Join
- Sort Merge Bucket Map Join
- Join顺序优化

Map Join

Hive的Map Join适用于能够在内存中存放下的小表（指表大小小于25MB），通过“hive.mapjoin.smalltable.filesize”定义小表的大小，默认为25MB。

Map Join的方法有两种：

- 使用/*+ MAPJOIN(join_table) */。
- 执行语句前设置如下参数，当前版本中该值默认为true。
set hive.auto.convert.join=true;

使用Map Join时没有Reduce任务，而是在Map任务前起了一个MapReduce Local Task，这个Task通过TableScan读取小表内容到本机，在本机以HashTable的形式保存并写入硬盘上传到DFS，并在distributed cache中保存，在Map Task中从本地磁盘或者distributed cache中读取小表内容直接与大表join得到结果并输出。

使用Map Join时需要注意小表不能过大，如果小表将内存基本用尽，会使整个系统性能下降甚至出现内存溢出的异常。

Sort Merge Bucket Map Join

使用Sort Merge Bucket Map Join必须满足以下2个条件：

- join的两张表都很大，内存中无法存放。
- 两张表都按照join key进行分桶（clustered by (column)）和排序（sorted by(column)），且两张表的分桶数正好是倍数关系。

通过如下设置，启用Sort Merge Bucket Map Join：

```
set hive.optimize.bucketmapjoin=true;  
set hive.optimize.bucketmapjoin.sortedmerge=true;
```

这种Map Join也没有Reduce任务，是在Map任务前启动MapReduce Local Task，将小表内容按桶读取到本地，在本机保存多个桶的HashTable备份并写入HDFS，并保存在Distributed Cache中，在Map Task中从本地磁盘或者Distributed Cache中按桶一个读取小表内容，然后与大表做匹配直接得到结果并输出。

Join 顺序优化

当有3张及以上的表进行Join时，选择不同的Join顺序，执行时间存在较大差异。使用恰当的Join顺序可以有效缩短任务执行时间。

Join顺序原则：

- Join出来结果较小的组合，例如表数据量小或两张表Join后产生结果较少，优先执行。
- Join出来结果大的组合，例如表数据量大或两张表Join后产生结果较多，在后面执行。

例如，customer表的数据量最多，orders表和lineitem表优先Join可获得较少的中间结果。

原有的Join语句如下：

```
select
  l_orderkey,
  sum(l_extendedprice * (1 - l_discount)) as revenue,
  o_orderdate,
  o_shippriority
from
  customer,
  orders,
  lineitem
where
  c_mktsegment = 'BUILDING'
  and c_custkey = o_custkey
  and l_orderkey = o_orderkey
  and o_orderdate < '1995-03-22'
  and l_shipdate > '1995-03-22'
limit 10;
```

Join顺序优化后如下：

```
select
  l_orderkey,
  sum(l_extendedprice * (1 - l_discount)) as revenue,
  o_orderdate,
  o_shippriority
from
  orders,
  lineitem,
  customer
where
  c_mktsegment = 'BUILDING'
  and c_custkey = o_custkey
  and l_orderkey = o_orderkey
  and o_orderdate < '1995-03-22'
  and l_shipdate > '1995-03-22'
limit 10;
```

注意事项

Join数据倾斜问题

执行任务的时候，任务进度长时间维持在99%，这种现象叫数据倾斜。

数据倾斜是经常存在的，因为有少量的Reduce任务分配到的数据量和其他Reduce差异过大，导致大部分Reduce都已完成任务，但少量Reduce任务还没完成的情况。

解决数据倾斜的问题，可通过设置“set hive.optimize.skewjoin=true”并调整hive.skewjoin.key的大小。hive.skewjoin.key是指Reduce端接收到多少个key即认为数据是倾斜的，并自动分发到多个Reduce。

10.8.3 Hive Group By 语句优化

操作场景

优化Group by语句，可提升命令执行速度和查询速度。

Group by的时候，Map端会先进行分组，分组完后分发到Reduce端，Reduce端再进行分组。可采用Map端聚合的方式来进行Group by优化，开启Map端初步聚合，减少Map的输出数据量。

操作步骤

在Hive客户端进行如下设置：

```
set hive.map.aggr=true;
```

注意事项

Group By数据倾斜

Group By也同样存在数据倾斜的问题，设置hive.groupby.skewindata为true，生成的查询计划会有两个MapReduce Job，第一个Job的Map输出结果会随机的分布到Reduce中，每个Reduce做聚合操作，并输出结果，这样的处理会使相同的Group By Key可能被分发到不同的Reduce中，从而达到负载均衡，第二个Job再根据预处理的结果按照Group By Key分发到Reduce中完成最终的聚合操作。

Count Distinct聚合问题

当使用聚合函数count distinct完成去重计数时，处理值为空的情况会使Reduce产生很严重的数据倾斜，可以将空值单独处理，如果是计算count distinct，可以通过where子句将该值排除掉，并在最后的count distinct结果中加1。如果还有其他计算，可以先将值为空的记录单独处理，再和其他计算结果合并。

10.8.4 Hive ORC 数据存储优化

操作场景

“ORC”是一种高效的列存储格式，在压缩比和读取效率上优于其他文件格式。

建议使用“ORC”作为Hive表默认的存储格式。

前提条件

已登录Hive客户端，具体操作请参见[Hive客户端使用实践](#)。

操作步骤

- 推荐：使用“SNAPPY”压缩，适用于压缩比和读取效率要求均衡场景。
Create table xx (col_name data_type) stored as orc tblproperties ("orc.compress"="SNAPPY");
- 可用：使用“ZLIB”压缩，适用于压缩比要求较高场景。
Create table xx (col_name data_type) stored as orc tblproperties ("orc.compress"="ZLIB");

 说明

xx为具体使用的Hive表名。

10.8.5 Hive SQL 逻辑优化

操作场景

在Hive上执行SQL语句查询时，如果语句中存在“(a&b) or (a&c)”逻辑时，建议将逻辑改为“a & (b or c)”。

样例

假设条件a为“p_partkey = l_partkey”，优化前样例如下所示：

```
select
    sum(l_extendedprice* (1 - l_discount)) as revenue
from
    lineitem,
    part
where
    (
        p_partkey = l_partkey
        and p_brand = 'Brand#32'
        and p_container in ('SM CASE', 'SM BOX', 'SM PACK', 'SM PKG')
        and l_quantity >= 7 and l_quantity <= 7 + 10
        and p_size between 1 and 5
        and l_shipmode in ('AIR', 'AIR REG')
        and l_shipinstruct = 'DELIVER IN PERSON'
    )
    or
    (
        p_partkey = l_partkey
        and p_brand = 'Brand#35'
        and p_container in ('MED BAG', 'MED BOX', 'MED PKG', 'MED PACK')
        and l_quantity >= 15 and l_quantity <= 15 + 10
        and p_size between 1 and 10
        and l_shipmode in ('AIR', 'AIR REG')
        and l_shipinstruct = 'DELIVER IN PERSON'
    )
    or
    (
        p_partkey = l_partkey
        and p_brand = 'Brand#24'
        and p_container in ('LG CASE', 'LG BOX', 'LG PACK', 'LG PKG')
        and l_quantity >= 26 and l_quantity <= 26 + 10
        and p_size between 1 and 15
        and l_shipmode in ('AIR', 'AIR REG')
        and l_shipinstruct = 'DELIVER IN PERSON'
    )
)
```

优化后样例如下所示：

```
select
    sum(l_extendedprice* (1 - l_discount)) as revenue
from
    lineitem,
    part
where p_partkey = l_partkey and
    ((
        p_brand = 'Brand#32'
        and p_container in ('SM CASE', 'SM BOX', 'SM PACK', 'SM PKG')
        and l_quantity >= 7 and l_quantity <= 7 + 10
        and p_size between 1 and 5
        and l_shipmode in ('AIR', 'AIR REG')
        and l_shipinstruct = 'DELIVER IN PERSON'
    ))
```

```
or
(
  p_brand = 'Brand#35'
  and p_container in ('MED BAG', 'MED BOX', 'MED PKG', 'MED PACK')
  and l_quantity >= 15 and l_quantity <= 15 + 10
  and p_size between 1 and 10
  and l_shipmode in ('AIR', 'AIR REG')
  and l_shipinstruct = 'DELIVER IN PERSON'
)
or
(
  p_brand = 'Brand#24'
  and p_container in ('LG CASE', 'LG BOX', 'LG PACK', 'LG PKG')
  and l_quantity >= 26 and l_quantity <= 26 + 10
  and p_size between 1 and 15
  and l_shipmode in ('AIR', 'AIR REG')
  and l_shipinstruct = 'DELIVER IN PERSON'
))
```

10.8.6 使用 Hive CBO 功能优化查询效率

操作场景

在Hive中执行多表Join时，Hive支持开启CBO（Cost Based Optimization），系统会自动根据表的统计信息，例如数据量、文件数等，选出合适计划提高多表Join的效率。Hive需要先收集表的统计信息后才能使CBO正确的优化。

说明

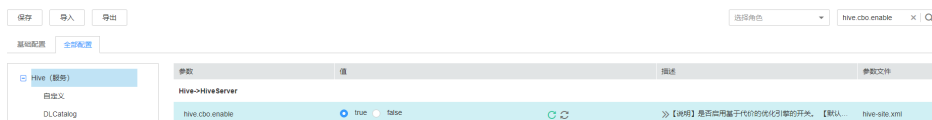
- CBO优化器会基于统计信息和查询条件，尽可能地使join顺序达到合适。但是也可能存在特殊情况导致join顺序调整不准确。例如数据存在倾斜，以及查询条件值在表中不存在等场景，可能调整出非优化的join顺序。
- 开启列统计信息自动收集时，需要在reduce侧做聚合统计。对于没有reduce阶段的insert任务，将会多出reduce阶段，用于收集统计信息。
- 本章节适用于MRS 3.x及后续版本。

前提条件

已登录Hive客户端，具体操作请参见[Hive客户端使用实践](#)。

操作步骤

步骤1 在Manager界面Hive组件的配置中搜索“hive.cbo.enable”参数，选中“true”永久开启功能。



步骤2 手动收集Hive表已有数据的统计信息。

执行以下命令，可以手动收集统计信息。仅支持统计一张表，如果需要统计不同的表需重复执行。

```
ANALYZE TABLE [db_name.]tablename [PARTITION(partcol1[=val1],  
partcol2[=val2], ...)]
```

```
COMPUTE STATISTICS
```


[FOR COLUMNS]

[NOSCAN];

📖 说明

- 指定FOR COLUMNS时，收集列级别的统计信息。
- 指定NOSCAN时，将只统计文件大小和个数，不扫描具体文件。

例如：

```
analyze table table_name compute statistics;
```

```
analyze table table_name compute statistics for columns;
```

步骤3 配置Hive自动收集统计信息。开启配置后，执行insert overwrite/into命令插入数据时才自动统计新数据的信息。

- 在Hive客户端执行以下命令临时开启收集：
set hive.stats.autogather = true;开启表/分区级别的统计信息自动收集。
set hive.stats.column.autogather = true;开启列级别的统计信息自动收集。

📖 说明

- 列级别统计信息的收集不支持复杂的数据类型，例如Map，Struct等。
- 表级别统计信息的自动收集不支持Hive on HBase表。
- 在Manager界面Hive的服务配置中，搜索参数“hive.stats.autogather”和“hive.stats.column.autogather”，选中“true”永久开启收集功能。

步骤4 执行以下命令可以查看统计信息。

```
DESCRIBE FORMATTED table_name[.column_name] PARTITION  
partition_spec;
```

例如：

```
desc formatted table_name;
```

```
desc formatted table_name id;
```

```
desc formatted table_name partition(time='2016-05-27');
```

📖 说明

分区表仅支持分区级别的统计信息收集，因此分区表需要指定分区来查询统计信息。

----结束

10.9 Hive 运维管理

10.9.1 Hive 常用配置参数

参数入口

- 对于MRS 3.x之前版本，登录MRS控制台，在左侧导航栏选择“现有集群”，单击集群名称，选择“组件管理 > Hive > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

- 对于MRS 3.x之后版本，登录FusionInsight Manager，选择“集群 > 服务 > Hive > 配置 > 全部配置”。

参数说明

表 10-9 Hive 参数说明

参数	参数说明	默认值
hive.auto.convert.join	Hive基于输入文件大小将普通join转为mapjoin的开关。 说明 在使用Hive进行联表查询，且关联的表无大小表的分别（小表数据<24M）时，建议将此参数值改为false，如果此时将此参数设置为true，执行联表查询时无法生成新的mapjoin。	取值范围： <ul style="list-style-type: none"> true false 默认值为true
hive.default.fileformat	Hive使用的默认文件格式。	MRS 3.x之前版本：TextFile MRS3.x及后续版本：RCFile
hive.exec.reducers.max	Hive提交的MR任务中reducer的最大个数。	999
hive.server2.thrift.max.worker.threads	HiveServer内部线程池，最大能启动的线程数量。	1000
hive.server2.thrift.min.worker.threads	HiveServer内部线程池，初始化时启动的线程数量。	5
hive.hbase.delete.mode.enabled	从Hive删除HBase记录的功能开关。如果启用，用户可以使用“remove table xx where xxx”从Hive中删除HBase记录。 说明 本参数适用于MRS 3.x及后续版本。	true
hive.metastore.server.min.threads	MetaStore启动的用于处理连接的线程数，如果超过设置的值之后，MetaStore就会一直维护不低于设定值的线程数，即常驻MetaStore线程池的线程会维护在指定值之上。	200
hive.server2.enable.doAs	HiveServer2在与其他服务（如YARN、HDFS等）会话时是否模拟客户端用户。如果将此配置项从false改成true，会导致只有列权限的用户访问相应表权限缺失。 说明 本参数适用于MRS 3.x及后续版本。	true

10.9.2 Hive 日志介绍

日志描述

日志路径：Hive相关日志的默认存储路径为“/var/log/Bigdata/hive/角色名”，Hive1相关日志的默认存储路径为“/var/log/Bigdata/hive1/角色名”，以此类推。

- HiveServer：“/var/log/Bigdata/hive/hiveserver”（运行日志），“/var/log/Bigdata/audit/hive/hiveserver”（审计日志）。
- MetaStore：“/var/log/Bigdata/hive/metastore”（运行日志），“/var/log/Bigdata/audit/hive/metastore”（审计日志）。
- WebHCat：“/var/log/Bigdata/hive/webhcat”（运行日志），“/var/log/Bigdata/audit/hive/webhcat”（审计日志）。

日志归档规则：Hive的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过20MB的时候（此日志文件大小可进行配置），会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd_hh-mm-ss>.[编号].log.zip”。最多保留最近的20个压缩文件，压缩文件保留个数和压缩文件阈值可以配置。

表 10-10 Hive 日志列表

日志类型	日志文件名	描述
运行日志	/hiveserver/hiveserver.out	HiveServer运行环境信息日志
	/hiveserver/hive.log	HiveServer进程的运行日志
	/hiveserver/hive-omm-<日期>-<PID>-gc.log.<编号>	HiveServer进程的GC日志
	/hiveserver/prestartDetail.log	HiveServer启动前的工作日志
	/hiveserver/check-serviceDetail.log	Hive服务启动是否成功的检查日志
	/hiveserver/cleanupDetail.log	HiveServer卸载的清理日志
	/hiveserver/startDetail.log	HiveServer进程启动日志
	/hiveserver/stopDetail.log	HiveServer进程停止日志
	/hiveserver/localtasklog/omm_<日期>_<任务ID>.log	Hive本地任务的运行日志
	/hiveserver/localtasklog/omm_<日期>_<任务ID>-gc.log.<编号>	Hive本地任务的GC日志
	/metastore/metastore.log	MetaStore进程的运行日志
	/metastore/hive-omm-<日期>-<PID>-gc.log.<编号>	MetaStore进程的GC日志

日志类型	日志文件名	描述
	/metastore/postinstallDetail.log	MetaStore安装后的工作日志
	/metastore/prestartDetail.log	MetaStore启动前的工作日志
	/metastore/cleanupDetail.log	MetaStore卸载的清理日志
	/metastore/startDetail.log	MetaStore进程启动日志
	/metastore/stopDetail.log	MetaStore进程停止日志
	/metastore/metastore.out	MetaStore运行环境信息日志
	/webhcat/webhcat-console.out	Webhcat进程启停正常日志
	/webhcat/webhcat-console-error.out	Webhcat进程启停异常日志
	/webhcat/prestartDetail.log	WebHCat启动前的工作日志
	/webhcat/cleanupDetail.log	Webhcat卸载时或安装前的清理日志
	/webhcat/hive-omm-<日期>-<PID>-gc.log.<编号>	WebHCat进程的GC日志
	/webhcat/webhcat.log	WebHCat进程的运行日志
审计日志	hive-audit.log	HiveServer审计日志
	hive-rangeraudit.log	
	metastore-audit.log	MetaStore审计日志
	webhcat-audit.log	WebHCat审计日志
	jetty-<日期>.request.log	Jetty服务的请求日志

日志级别

Hive提供了如表10-11所示的日志级别。

运行日志的级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 10-11 日志级别

级别	描述
ERROR	ERROR表示系统运行的错误信息。
WARN	WARN表示当前事件处理存在异常信息。
INFO	INFO表示记录系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 参考[修改集群服务配置参数](#)，进入Hive服务“全部配置”页面。
- 步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤3** 选择所需修改的日志级别并保存。

 说明

配置Hive日志级别后可立即生效，无需重启服务。

----结束

日志格式

Hive的日志格式如下所示：

表 10-12 日志格式

日志类型	格式	示例
运行日志	<yyyy-MM-dd HH:mm:ss,SSS> <LogLevel> <产生该日志的 线程名字> <log中的 message> <日志事件的发生 位置>	2014-11-05 09:45:01,242 INFO main Starting hive metastore on port 21088 org.apache.hadoop.hive.metas- tore.HiveMetaStore.main(Hive MetaStore.java:5198)
审计日志	<yyyy-MM-dd HH:mm:ss,SSS> <LogLevel> <产生该日志的 线程名字> <User Name><User IP><Time><Operation><Re- source><Result><Detail > < 日志事件的发生位置>	2018-12-24 12:16:25,319 INFO HiveServer2-Handler- Pool: Thread-185 UserName=hive UserIP=10.153.2.204 Time=2018/12/24 12:16:25 Operation=CloseSession Result=SUCCESS Detail= org.apache.hive.service.cli.thrif- t.ThriftCLIService.logAuditEven- t(ThriftCLIService.java:434)

10.10 Hive 常见 SQL 语法说明

10.10.1 Hive SQL 扩展语法说明

Hive SQL支持Hive-3.1.0版本中的所有特性，详情请参见<https://cwiki.apache.org/confluence/display/hive/languagemanual>。

系统提供的扩展Hive语句如表10-13所示。

表 10-13 扩展 Hive 语句

扩展语法	语法说明	语法示例	示例说明
<pre>CREATE [TEMPORARY] [EXTERNAL] TABLE [IF NOT EXISTS] [db_name.]table_ name (col_name data_type [COMMENT col_comment], ...) [ROW FORMAT row_format] [STORED AS file_format] STORED BY 'storage.handler.cl ass.name' [WITH SERDEPROPERTIE S (...)] [TBLPROPERTIES ("groupId"=" group1 ","locatorId"="loc ator1")] ...;</pre>	<p>创建一个hive表，并指定表数据文件分布的locator信息。详细说明请参见使用HDFS Colocation存储Hive表。</p>	<pre>CREATE TABLE tab1 (id INT, name STRING) row format delimited fields terminated by '\t' stored as RCFILE TBLPROPERTIES(" groupId"=" group1 ","locatorId"="loc ator1");</pre>	<p>创建表tab1，并指定tab1的表数据分布在locator1节点上。</p>

扩展语法	语法说明	语法示例	示例说明
<pre>CREATE [TEMPORARY] [EXTERNAL] TABLE [IF NOT EXISTS] [db_name.]table_ name (col_name data_type [COMMENT col_comment], ...) [ROW FORMAT row_format] [STORED AS file_format] STORED BY 'storage.handler.cl ass.name' [WITH SERDEPROPERTIE S (...)] ... [TBLPROPERTIES ('column.encode. columns'='col_na me1,col_name2'] 'column.encode.i ndices'='col_id1,c ol_id2','column.e ncode.classname' ='encode_classna me')]...;</pre>	<p>创建一个hive表，并指定表的加密列和加密算法。详细说明请参见使用Hive列加密功能。</p>	<pre>create table encode_test(id INT, name STRING, phone STRING, address STRING) ROW FORMAT SERDE 'org.apache.hadoop p.hive.serde2.lazy. LazySimpleSerDe' WITH SERDEPROPERTIE S ('column.encode.i ndices'='2,3', 'column.encode.cl assname'='org.apa che.hadoop.hive.s erde2.SMS4Rewrit er') STORED AS TEXTFILE;</pre>	<p>创建表 encode_test，并指定插入数据时对第2、3列加密，加密算法类为 org.apache.hadoop.p.hive.serde2.SMS4Rewriter。</p>
<pre>REMOVE TABLE hbase_tablename [WHERE where_condition];</pre>	<p>删除hive on hbase表中符合条件的数据。详细说明请参见删除Hive on HBase表中的单行记录。</p>	<pre>remove table hbase_table1 where id = 1;</pre>	<p>删除表中符合条件“id =1”的数据。</p>

扩展语法	语法说明	语法示例	示例说明
<pre>CREATE [TEMPORARY] [EXTERNAL] TABLE [IF NOT EXISTS] [db_name.]table_ name (col_name data_type [COMMENT col_comment], ...) [ROW FORMAT row_format] STORED AS inputformat 'org.apache.hado op.hive.contrib.fil eformat.Specifie dDelimiterInputF ormat' outputformat 'org.apache.hadoo p.hive.ql.io.HiveIg noreKeyTextOutpu tFormat';</pre>	<p>创建hive表，并设定表可以指定自定义行分隔符。详细说明请参见自定义行分隔符。</p>	<pre>create table blu(time string, num string, msg string) row format delimited fields terminated by ',' stored as inputformat 'org.apache.hado op.hive.contrib.fil eformat.Specifie dDelimiterInputF ormat' outputformat 'org.apache.hadoo p.hive.ql.io.HiveIg noreKeyTextOutpu tFormat';</pre>	<p>创建表blu，指定inputformat为SpecifiedDelimiterInputFormat，以便查询时可以指定表的查询行分隔符。</p>

10.10.2 自定义 Hive 表行分隔符

操作场景

通常情况下，Hive以文本文件存储的表会以回车作为其行分隔符，即在查询过程中，以回车符作为一行表数据的结束符。但某些数据文件并不是以回车分隔的规则文本格式，而是以某些特殊符号分隔其规则文本。

MRS Hive支持指定不同的字符或字符组作为Hive文本数据的行分隔符，即在创建表的时候，指定inputformat为SpecifiedDelimiterInputFormat，然后在每次查询前，都设置如下参数来指定分隔符，就可以以指定的分隔符查询表数据。

```
set hive.textinput.record.delimiter="";
```

📖 说明

- 当前版本的Hue组件，不支持导入文件到Hive表时设置多个分隔符。
- 本章节适用于MRS 3.x及后续版本。

操作步骤

步骤1 创建表时指定inputFormat和outputFormat：

```
CREATE [TEMPORARY] [EXTERNAL] TABLE [IF NOT EXISTS]
[db_name.]table_name [(col_name data_type [COMMENT col_comment], ...)]
```



```
[ROW FORMAT row_format] STORED AS inputformat
'org.apache.hadoop.hive.contrib.fileformat.SpecifiedDelimiterInputFormat'
outputformat 'org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat';
```

步骤2 查询之前指定配置项：

```
set hive.textinput.record.delimiter='!@!';
```

Hive会以‘!@!’为行分隔符查询数据。

----结束

10.10.3 Hive 支持的传统关系型数据库语法说明

概述

Hive支持如下传统关系型数据库语法：

- Grouping
- EXCEPT、INTERSECT

Grouping

语法简介：

- 当Group by语句带with rollup/cube选项时，Grouping才有意义。
- CUBE生成的结果集显示了所选列中值的所有组合的聚合。
- ROLLUP生成的结果集显示了所选列中值的某一层级结构的聚合。
- Grouping：当用CUBE或ROLLUP运算符添加行时，附加的列输出值为1；当所添加的行不是由CUBE或ROLLUP产生时，附加列值为0。

例如，Hive中有一张表“table_test”，表结构如下所示：

```
+-----+-----+--+
| table_test.id | table_test.value |
+-----+-----+--+
| 1             | 10                |
| 1             | 15                |
| 2             | 20                |
| 2             | 5                 |
| 2             | 13                |
+-----+-----+--+
```

执行如下语句：

```
select id,grouping(id),sum(value) from table_test group by id with rollup;
```

得到如下结果：

```
+-----+-----+-----+
| id | groupingresult | sum |
+-----+-----+-----+
| 1  | 0              | 25  |
| NULL | 1              | 63  |
| 2  | 0              | 38  |
+-----+-----+-----+
```

EXCEPT、INTERSECT

语法简介

- EXCEPT返回两个结果集的差（即从左查询中返回右查询没有找到的所有非重复值）。
- INTERSECT返回两个结果集的交集（即两个查询都返回的所有非重复值）。

例如，Hive中有两张表“test_table1”、“test_table2”。

“test_table1”表结构如下所示：

```
+-----+--+
| test_table1.id |
+-----+--+
| 1             |
| 2             |
| 3             |
| 4             |
+-----+--+
```

“test_table2”表结构如下所示：

```
+-----+--+
| test_table2.id |
+-----+--+
| 2             |
| 3             |
| 4             |
| 5             |
+-----+--+
```

- 执行如下的EXCEPT语句：
select id from test_table1 except select id from test_table2;

显示如下结果：

```
+-----+--+
| _alias_0.id |
+-----+--+
| 1           |
+-----+--+
```

- 执行INTERSECT语句：
select id from test_table1 intersect select id from test_table2;

显示如下结果：

```
+-----+--+
| _alias_0.id |
+-----+--+
| 2           |
| 3           |
| 4           |
+-----+--+
```

10.11 Hive 常见问题

10.11.1 如何删除所有 HiveServer 中的永久函数

问题

如果需要删除永久函数（Permanent UDF），如何在多个HiveServer之间同步删除？

回答

因为多个HiveServer之间共用一个MetaStore存储数据库，所以MetaStore存储数据库和HiveServer的内存之间数据同步有延迟。如果在单个HiveServer上删除永久函数，操作结果将无法同步到其他HiveServer上。

遇到如上情况，需要登录Hive客户端，连接到每个HiveServer，并分别删除永久函数。具体操作如下：

步骤1 以Hive客户端安装用户登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录。

cd 客户端安装目录

例如安装目录为“/opt/client”，则执行以下命令：

cd /opt/client

步骤3 执行以下命令配置环境变量。

source bigdata_env

步骤4 执行以下命令进行用户认证。

kinit Hive业务用户

📖 说明

登录的用户需具备Hive admin权限。

步骤5 执行如下命令，连接指定的HiveServer。

beeline -u "jdbc:hive2://10.39.151.74:21066/default;sas.qop=auth-conf;auth=KERBEROS;principal=hive/hadoop.<系统域名>@<系统域名>"

📖 说明

- 10.39.151.74为HiveServer所在节点的IP地址。
- 21066为HiveServer端口。HiveServer端口默认范围为21066~21070，用户需根据实际配置端口进行修改。
- hive为用户名。例如，使用Hive1实例时，则使用hive1。
- 用户可登录FusionInsight Manager，选择“系统 > 权限 > 域和互信”，查看“本端域”参数，即为当前系统域名。
- “hive/hadoop.<系统域名>”为用户名，用户的用户名所包含的系统域名所有字母为小写。

步骤6 执行如下命令，启用Hive admin权限。

set role admin;

步骤7 执行如下命令，删除永久函数。

drop function function_name;

📖 说明

- function_name为永久函数的函数名。
- 如果永久函数是在Spark中创建的，在Spark中删除该函数后需要在HiveServer中删除，即执行上述删除命令。

步骤8 确定是否已连接所有HiveServer并删除永久函数。

- 是，操作完毕。
- 否，执行**步骤5**。

----结束

10.11.2 为什么已备份的 Hive 表无法执行 drop 操作

问题

为什么已备份的Hive表执行drop操作会失败？

回答

由于已备份Hive表对应的HDFS目录创建了快照，导致HDFS目录无法删除，造成Hive表删除失败。

Hive表在执行备份操作时，会创建表对应的HDFS数据目录快照。而HDFS的快照机制有一个约束：如果一个HDFS目录已创建快照，则在快照完全删除之前，该目录无法删除或修改名称。Hive表（除EXTERNAL表外）执行drop操作时，会尝试删除该表对应的HDFS数据目录，如果目录删除失败，系统会提示表删除失败。

如果确实需要删除该表，可手动删除涉及到该表的所有备份任务。

10.11.3 如何在 Hive 自定义函数中操作本地文件

问题

在Hive自定义函数中需要操作本地文件，例如读取文件的内容，需要如何操作？

回答

默认情况下，可以在UDF中用文件的相对路径来操作文件，如下示例代码：

```
public String evaluate(String text) {  
    // some logic  
    File file = new File("foo.txt");  
    // some logic  
    // do return here  
}
```

在Hive中使用时，将UDF中用到的文件“foo.txt”上传到HDFS上，如上传到“hdfs://hacluster/tmp/foo.txt”，使用以下语句创建UDF，在UDF中就可以直接操作“foo.txt”文件了：

```
create function testFunc as 'some.class' using jar 'hdfs://hacluster/  
somejar.jar', file 'hdfs://hacluster/tmp/foo.txt';
```

例外情况下，如果“hive.fetch.task.conversion”参数的值为“more”，在UDF中不能再使用相对路径来操作文件，而要使用绝对路径，并且保证所有的HiveServer节点和NodeManager节点上该文件是存在的且omm用户对该文件有相应的权限，才能正常在UDF中操作本地文件。

10.11.4 如何强制停止 Hive 执行的 MapReduce 任务

问题

在Hive执行MapReduce任务长时间卡住的情况下想手动停止任务，需要如何操作？

回答

步骤1 登录FusionInsight Manager。

步骤2 选择“集群 > 待操作的集群名称 > 服务 > Yarn”。

步骤3 单击左侧页面的“ResourceManager(主机名称, 主)”按钮，登录Yarn界面。

步骤4 单击对应任务ID的按钮进入任务页面，单击界面左上角的“Kill Application”按钮，在弹框中单击“确认”停止任务。

----结束

10.11.5 Hive 不支持复杂类型字段名称中包含哪些特殊字符

问题

Hive复杂类型字段名称中包含特殊字符，导致建表失败。

回答

Hive不支持复杂类型字段名称中包含特殊字符。

特殊字符是指英文大小写字母、阿拉伯数字、中文字符、葡萄牙文字符以外的其他字符。

用户在创建相关字段时，应避免使用相关特殊字符。

10.11.6 如何对 Hive 表大小数据进行监控

问题

如何对Hive中的表大小数据进行监控？

回答

当用户要对Hive表大小数据进行监控时，可以通过HDFS的精细化监控对指定表目录进行监控，从而到达监控指定表大小数据的目的。

前提条件


- Hive、HDFS组件功能正常
- HDFS精细化监控功能正常

操作步骤

步骤1 登录FusionInsight Manager。

步骤2 通过“集群 > 待操作集群的名称 > 服务 > HDFS > 资源”，进入HDFS精细化页面。

步骤3 找到“资源使用（按目录）”监控项，单击该监控项左上角第一个图标。

资源使用（按目录）

步骤4 进入配置空间监控子页面，单击“添加”。

步骤5 在名称空格中填写监控的表名称（或其他用户自定义的别名），在路径中填写需要监控表的路径。单击“确定”。该监控的横坐标为时间，纵坐标为监控目录的大小。

----结束

10.11.7 如何防止 insert overwrite 语句误操作导致数据丢失

问题

如何对重点目录进行保护，防止“insert overwrite”语句误操作导致数据丢失？

回答

当用户要对Hive重点数据库、表或目录进行监控，防止“insert overwrite”语句误操作导致数据丢失时，可以利用Hive配置中的“hive.local.dir.confblacklist”进行目录保护。

该配置项已对“/opt/”，“/user/hive/warehouse”等目录进行了默认配置。

前提条件

Hive、HDFS组件功能正常。

操作步骤

步骤1 登录FusionInsight Manager。

步骤2 选择“集群 > 待操作集群的名称 > 服务 > Hive > 配置 > 全部配置”，搜索“hive.local.dir.confblacklist”配置项。

步骤3 在该配置项中添加用户要重点保护的数据库、表或目录路径。

步骤4 输入完成后，单击“保存”，保存配置项。

----结束

10.11.8 未安装 HBase 时 Hive on Spark 任务卡顿如何处理

操作场景

此功能适用于Hive组件。

按如下操作步骤设置参数后，在未安装HBase的环境执行Hive on Spark任务时，可避免任务卡顿。

📖 说明

Hive on Spark任务的Spark内核版本已经升级到Spark2x，可以支持在不安装Spark2x的情况下，执行Hive on Spark任务。如果没有安装HBase，默认在执行Spark任务时，会尝试去连接Zookeeper访问HBase，直到超时，这样会造成任务卡顿。

在未安装HBase的环境，要执行Hive on Spark任务，可以按如下操作处理。如果是从已有HBase低版本环境升级上来的，升级完成之后可不进行设置。

操作步骤

步骤1 登录FusionInsight Manager 。

步骤2 选择“集群 > 待操作集群的名称 > 服务 > Hive > 配置 > 全部配置”。

步骤3 选择“HiveServer（角色）> 自定义”，对参数文件“spark-defaults.conf”添加自定义参数，设置“名称”为“spark.security.credentials.hbase.enabled”，“值”为“false”。

步骤4 单击“保存”，在弹出对话框单击“确定”。

步骤5 选择“集群 > 待操作集群的名称 > 服务 > Hive > 实例”，勾选所有Hive实例，选择“更多 > 重启实例”，输入密码，单击“确定”。

----结束

10.11.9 Hive 使用 WHERE 条件查询超过 3.2 万分区的表报错

问题：

Hive创建超过3.2万分区的表，执行带有WHERE分区条件查询时出现异常，且“metastore.log”中打印的异常信息包含以下信息：

```
Caused by: java.io.IOException: Tried to send an out-of-range integer as a 2-byte value: 32970
    at org.postgresql.core.PGStream.SendInteger2(PGStream.java:199)
    at org.postgresql.core.v3.QueryExecutorImpl.sendParse(QueryExecutorImpl.java:1330)
    at org.postgresql.core.v3.QueryExecutorImpl.sendOneQuery(QueryExecutorImpl.java:1601)
    at org.postgresql.core.v3.QueryExecutorImpl.sendParse(QueryExecutorImpl.java:1191)
    at org.postgresql.core.v3.QueryExecutorImpl.execute(QueryExecutorImpl.java:346)
```

回答：

带有分区条件的查询，Hiveserver会对分区进行优化，避免全表扫描，需要查询元数据符合条件的所有分区。

而gaussDB中提供的接口sendOneQuery，调用的sendParse方法中对参数的限制为32767，如果分区条件数超过32767就会产生异常。

10.11.10 使用 IBM 的 JDK 访问 Beeline 客户端出现连接 HiveServer 失败

操作场景

查看客户端使用的jdk版本，如果是IBM JDK，则需要对Beeline客户端进行改造，否则会造成连接hiveserver失败。

操作步骤

步骤1 登录FusionInsight Manager 页面，选择“系统 > 权限 > 用户”，在待操作用户的“操作”栏下选择“更多 > 下载认证凭据”，选择集群信息后单击“确定”，下载keytab文件。

步骤2 解压keytab文件，使用WinSCP工具将解压得到的“user.keytab”文件上传到待操作节点的Hive客户端安装目录下，例如：“/opt/client”。

步骤3 使用以下命令打开hive客户端目录下面的配置文件Hive/component_env:

```
vi Hive客户端安装目录/Hive/component_env
```

在变量“export CLIENT_HIVE_URI”所在行后面添加如下内容:

```
\;user.principal=用户名@HADOOP.COM\;user.keytab=user.keytab文件所在路径/user.keytab
```

----结束

10.11.11 Hive 表的 Location 支持跨 OBS 和 HDFS 路径吗

问题

Hive表的location支持跨OBS和HDFS路径吗？

回答

1. Hive存储在OBS上的普通表，支持表location配置为hdfs路径。
2. 同一个Hive服务中可以分别创建存储在OBS上的表和存储在HDFS上的表。
3. Hive存储在OBS上的分区表，不支持将分区location配置为hdfs路径（存储在HDFS上的分区表也不支持修改分区location为OBS）。

10.11.12 MapReduce 引擎无法查询 Tez 引擎执行 union 语句写入的数据

问题

Hive通过Tez引擎执行union相关语句写入的数据，切换到Mapreduce引擎后进行查询，发现数据没有查询出来。

回答

由于Hive使用Tez引擎在执行union语句时，生成的输出文件会存在HIVE_UNION_SUBDIR目录，切回Mapreduce引擎后默认不读取目录下的文件，所以没有读取到HIVE_UNION_SUBDIR目录下的数据。

此时可以设置参数set **mapreduce.input.fileinputformat.input.dir.recursive=true**，开启union优化，决定是否读取目录下的数据。

10.11.13 Hive 是否支持对同一张表或分区进行并发写数据

问题

为什么通过接口并发对Hive表进行写数据会导致数据不一致？

回答

Hive不支持对同一张表或同一个分区进行并发数据插入，这样会导致多个任务操作同一个数据临时目录，一个任务将另一个任务的数据移走，导致任务数据异常。

解决方法是修改业务逻辑，单线程插入数据到同一张表或同一个分区。

说明

MRS 3.1.3及之后版本Hive支持对同一张表或分区进行并发写数据。

10.11.14 Hive 是否支持向量化查询

问题

当设置向量化参数hive.vectorized.execution.enabled=true时，为什么执行hive on Tez/Mapreduce/Spark时会偶现一些空指针或类型转化异常？

回答

当前MRS Hive不支持向量化执行。

向量化执行有很多社区问题引入目前没有稳定修复，默认hive.vectorized.execution.enabled=false，不建议将此参数打开。

10.11.15 Hive 表的 HDFS 数据目录被误删，但是元数据仍然存在，导致执行任务报错

问题

Hive表HDFS数据目录被误删，但是元数据仍然存在，导致执行任务报错。

回答

这是一种误操作的异常情况，需要手动删除对应表的元数据后重试。

例如：

执行以下命令进入控制台：

```
source ${BIGDATA_HOME}/FusionInsight_BASE_8.1.0.1/install/FusionInsight-  
dbservice-2.7.0/.dbservice_profile
```

```
gsql -p 20051 -U hive -d hivemeta -W HiveUser@
```

```
执行 delete from tbls where tbl_id='xxx';
```

10.11.16 如何关闭 Hive 客户端日志

问题

如何关闭Hive客户端的运行日志？

回答

步骤1 使用root用户登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录，例如“/opt/client”。

```
cd /opt/client
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 根据集群认证模式，完成Hive客户端登录。

- 安全模式，则执行以下命令，完成用户认证并登录Hive客户端。

```
kinit 组件业务用户
```

```
beeline
```

- 普通模式，则执行以下命令，登录Hive客户端。

- 使用指定组件业务用户登录Hive客户端。

```
beeline -n 组件业务用户
```

- 不指定组件业务用户登录Hive客户端，则会以当前操作系统用户登录。

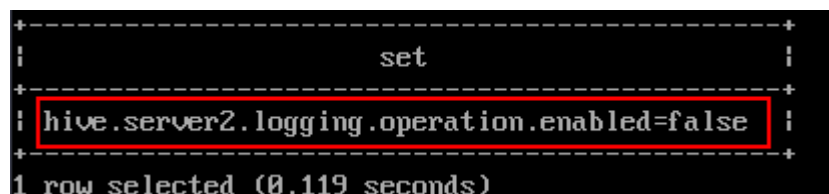
```
beeline
```

步骤5 执行以下命令关闭客户端日志：

```
set hive.server2.logging.operation.enabled=false;
```

步骤6 执行以下命令查看客户端日志是否已关闭，如下图所示即为关闭成功。

```
set hive.server2.logging.operation.enabled;
```



```
+-----+
|               set               |
+-----+
| hive.server2.logging.operation.enabled=false |
+-----+
1 row selected (0.119 seconds)
```

----结束

10.11.17 为什么在 Hive 自定义配置中添加 OBS 快删目录后不生效

问题

在配置MRS多用户访问OBS细粒度权限的场景中，在Hive自定义配置中添加OBS快删目录的配置后，删除Hive表。

执行结果为成功，但是OBS目录没有删掉。

回答

由于没有给用户配置快删目录的权限，导致数据不能被删除。

需要修改用户对应的委托的IAM自定义策略，在策略内容上，配置Hive快删目录的权限。

10.11.18 Hive 配置类问题

- Hive SQL执行报错：java.lang.OutOfMemoryError: Java heap space.
解决方案：
 - 对于MapReduce任务，增大下列参数：
set mapreduce.map.memory.mb=8192;
set mapreduce.map.java.opts=-Xmx6554M;
set mapreduce.reduce.memory.mb=8192;
set mapreduce.reduce.java.opts=-Xmx6554M;
 - 对于Tez任务，增大下列参数：
set hive.tez.container.size=8192;
- Hive SQL对列名as为新列名后，使用原列名编译报错：Invalid table alias or column reference 'xxx'.
解决方案：**set hive.cbo.enable=true;**
- Hive SQL子查询编译报错：Unsupported SubQuery Expression 'xxx': Only SubQuery expressions that are top level conjuncts are allowed.
解决方案：**set hive.cbo.enable=true;**
- Hive SQL子查询编译报错：CalciteSubquerySemanticException [Error 10249]: Unsupported SubQuery Expression Currently SubQuery expressions are only allowed as Where and Having Clause predicates.
解决方案：**set hive.cbo.enable=true;**
- Hive SQL编译报错：Error running query: java.lang.AssertionError: Cannot add expression of different type to set.
解决方案：**set hive.cbo.enable=false;**
- Hive SQL执行报错：java.lang.NullPointerException at org.apache.hadoop.hive.ql.udf.generic.GenericUDAFComputeStats \$GenericUDAFNumericStatsEvaluator.init.
解决方案：**set hive.map.aggr=false;**
- Hive SQL设置hive.auto.convert.join = true（默认开启）和hive.optimize.skewjoin=true执行报错：ClassCastException org.apache.hadoop.hive.ql.plan.ConditionalWork cannot be cast to org.apache.hadoop.hive.ql.plan.MapredWork.
解决方案：**set hive.optimize.skewjoin=false;**
- Hive SQL设置hive.auto.convert.join=true（默认开启）、hive.optimize.skewjoin=true和hive.exec.parallel=true执行报错：java.io.FileNotFoundException: File does not exist:xxx/reduce.xml.
解决方案：
 - 方法一：切换执行引擎为Tez，详情请参考[切换Hive执行引擎为Tez](#)。

- 方法二：**set hive.exec.parallel=false;**
- 方法三：**set hive.auto.convert.join=false;**
- Hive on Tez执行Bucket表Join报错：NullPointerException at org.apache.hadoop.hive.ql.exec.CommonMergeJoinOperator.mergeJoinComputeKeys
解决方案：**set tez.am.container.reuse.enabled=false;**

10.12 Hive 故障排除

10.12.1 如何对 insert overwrite 自读自写场景进行优化

场景说明

对于需要使用动态分区插入（使用历史分区更新）数据到目的表中，且和数据源表是同一张表时，由于直接在原表上执行insert overwrite可能会导致数据丢失或数据不一致的风险，建议先使用一个临时表来处理数据，再执行insert overwrite操作。

操作步骤

假设存在如下一张表：

```
user_data(user_group int, user_name string, update_time timestamp);
```

其中**user_group**是分区列，需要根据已有数据，按更新时间进行排序，刷新用户组信息。操作步骤如下：

步骤1 在Hive Beeline命令行执行以下命令开启Hive动态分区：

```
set hive.exec.dynamic.partition=true;  
set hive.exec.dynamic.partition.mode=nonstrict;
```

步骤2 执行以下命令创建一个临时表，用于存储去重后的数据：

```
CREATE TABLE temp_user_data AS  
SELECT * FROM (  
SELECT *,  
ROW_NUMBER() OVER(PARTITION BY user_group ORDER BY update_time  
DESC) as rank  
FROM user_data  
) tmp  
WHERE rank = 1;
```

步骤3 执行以下命令使用临时数据作为数据源，并插入到目的表中：

```
INSERT OVERWRITE TABLE user_data  
SELECT user_group, user_name, update_time  
FROM temp_user_data;
```

步骤4 执行以下命令清理临时表：

```
DROP TABLE IF EXISTS temp_user_data;
```

----结束

10.12.2 Hive SQL 运行变慢阶段如何排查

场景说明

对于一个Hive SQL任务，如果运行时间突然变长，可能有多种原因造成，如HiveServer编译变慢、访问HDFS变慢、访问Yarn变慢或访问元数据变慢。


操作步骤

以下为Hive SQL在哪个阶段执行变慢的排查方法：

- HiveServer编译变慢

查看HiveServer审计日志“/var/log/Bigdata/audit/hive/hiveserver/hive-audit.log”，搜索运行的SQL，如图10-5的SQL为“show databases”，过滤线程名包含“HiveServer2-Handler-Pool”的两条记录，第一条记录为开始编译时间，第二条记录为结束编译时间。根据SQL变慢前后的审计记录可判断是否是HiveServer编译变慢。

图 10-5 HiveServer 正在运行的 SQL



```
2024-04-22 17:05:42.113 [INFO] | hiveserver2-handler-pool: Thread-318 | user=hiveuser | userIP=XXXXXXXXXXXXXXXXXXXX | time=2024/04/22 17:05:42 | operation=ExecuteStatement | stmt=show databases | result=Details: conf | org.apache.hive.service.cli.thrift.ThriftTFCService | org.apache.hive.service.cli.thrift.ThriftTFCService | java:552  
2024-04-22 17:05:42.113 [INFO] | hiveserver2-handler-pool: Thread-318 | operation=ExecuteStatement | stmt=show databases | result=SUCCESS | org.apache.hive.service.cli.thrift.ThriftTFCService | org.apache.hive.service.cli.thrift.ThriftTFCService | java:552
```

- 访问HDFS变慢

- 方法一：

查看HiveServer运行日志“/var/log/Bigdata/hive/hiveserver/hive.log”，搜索运行的SQL对应的线程日志，再搜索number of splits日志，如果该日志前后间隔时间较长，则表示访问HDFS变慢。

- 方法二：

通过打印HiveServer进程jstack，查看相关线程是否卡顿在访问HDFS部分，是的话则表示访问HDFS变慢。

- 方法三：

查看HDFS RPC监控，看是否在SQL运行变慢期间HDFS RPC异常升高，是的话则大概率是访问HDFS变慢。

- 访问Yarn变慢

查看HiveServer运行日志“/var/log/Bigdata/hive/hiveserver/hive.log”，搜索运行的SQL对应的线程日志，再搜索Kill Command日志，如果该日志后面间隔较长时间才打出下一行日志，则表示访问Yarn变慢。

- 访问元数据变慢

在FusionInsight Manager界面，选择“集群 > 服务 > Hive > 实例 > 任一MetaStore实例 > 图表”，在“图表分类”选择“操作统计”，查看“create_table API元数据操作耗时情况”和“add_partitions_req api执行情况”等监控，看SQL运行慢是否是因为MetaStore访问变慢。

11 使用 Hudi

11.1 Hudi 表概述

Hudi 表类型

- Copy On Write
写时复制表也简称cow表，使用parquet文件存储数据，内部的更新操作需要通过重写原始parquet文件完成。
 - 优点：读取时，只读取对应分区的一个数据文件即可，较为高效。
 - 缺点：数据写入的时候，需要复制一个先前的副本再在其基础上生成新的数据文件，这个过程比较耗时。且由于耗时，读请求读取到的数据相对就会滞后。
- Merge On Read
读时合并表也简称mor表，使用列格式parquet和行格式Avro两种方式混合存储数据。其中parquet格式文件用于存储基础数据，Avro格式文件（也可叫做log文件）用于存储增量数据。
 - 优点：由于写入数据先写delta log，且delta log较小，所以写入成本较低。
 - 缺点：需要定期合并整理compact，否则碎片文件较多。读取性能较差，因为需要将delta log和老数据文件合并。

Hudi 表存储

Hudi在写入数据时会根据设置的存储路径、表名、分区结构等属性生成Hudi表。

Hudi表的数据文件，可以使用操作系统的文件系统存储，也可以使用HDFS这种分布式的文件系统存储。为了后续分析性能和数据的可靠性，一般使用HDFS进行存储。以HDFS存储来看，一个Hudi表的存储文件分为两类。

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<input type="checkbox"/>	drwxr-xr-x	testcz	hadoop	0 B	Apr 25 15:32	0	0 B	.hoodie
<input type="checkbox"/>	drwxr-xr-x	testcz	hadoop	0 B	Apr 25 15:30	0	0 B	americas
<input type="checkbox"/>	drwxr-xr-x	testcz	hadoop	0 B	Apr 25 15:30	0	0 B	asia

- “.hoodie”文件夹中存放了对应的文件合并操作相关的日志文件。

drwxr-xr-x	admintest	hadoop	0 B	Mar 30 09:44	0	0 B	.aux
drwxr-xr-x	admintest	hadoop	0 B	Mar 30 11:45	0	0 B	.temp
-rw-r--r--	admintest	hadoop	4.58 KB	Mar 30 09:44	3	128 MB	20210330094435.deltacommit
-rw-r--r--	admintest	hadoop	0 B	Mar 30 09:44	3	128 MB	20210330094435.deltacommit.inflight
-rw-r--r--	admintest	hadoop	0 B	Mar 30 09:44	3	128 MB	20210330094435.deltacommit.requested

- 包含 `_partition_key` 相关的路径是实际的数据文件和 metadata，按分区存储。Hudi 的数据文件使用 Parquet 文件格式的 base file 和 Avro 格式的 log file 存储。

-rw-r--r--	admintest	hadoop	93 B	Mar 30 09:44	3	128 MB	.hoodie_partition_metadata
-rw-r--r--	admintest	hadoop	441.77 KB	Mar 30 09:46	3	128 MB	2b4d098e-4dc8-4633-a22a-dc22f87c57d9-1_0-13-22_20210330094613.parquet
-rw-r--r--	admintest	hadoop	445.28 KB	Mar 30 09:44	3	128 MB	4010e8a8-1b20-4be7-8442-4e30af401e84-0_1-4-8_20210330094435.parquet

📖 说明

查看 Hudi 表：登录 FusionInsight Manager 页面，选择“集群 > 服务 > HDFS”，在“概览”页面单击 NameNode WebUI 后的链接，进入到 HDFS 的 WebUI 界面，选择“Utilities > Browse the file system”，即可查看 Hudi 表。

11.2 使用 Spark Shell 创建 Hudi 表

操作场景

本指南通过使用 spark-shell 简要介绍了 Hudi 功能。使用 Spark 数据源，将通过代码段展示如何插入和更新 Hudi 的默认存储类型数据集：COW 表。每次写操作之后，还将展示如何读取快照和增量数据。

前提条件

- 在 Manager 界面创建用户并添加 hadoop 和 hive 用户组，主组加入 hadoop。

操作步骤

步骤1 下载并安装 Hudi 客户端，具体请参考 [安装客户端（3.x 及之后版本）](#) 章节。

📖 说明

目前 Hudi 集成在 Spark2x 中，用户从 Manager 页面下载 Spark2x 客户端即可，例如客户端安装目录为：“/opt/client”。

步骤2 使用 root 登录客户端安装节点，执行如下命令：

```
cd /opt/client
```

步骤3 执行 source 命令加载环境变量：

```
source bigdata_env
```

```
source Hudi/component_env
```

```
kinit 创建的用户
```

📖 说明

- 新创建的用户需要修改密码，更改密码后重新kinit登录。
- 普通模式（未开启kerberos认证）无需执行kinit命令。
- 多服务场景下，在source bigdata_env之后，请先source Spark服务的component_env，再去source Hudi的component_env。

步骤4 使用spark-shell --master yarn-client，引入Hudi包生成测试数据：

- 引入需要的包

```
import org.apache.hudi.QuickstartUtils._
import scala.collection.JavaConversions._
import org.apache.spark.sql.SaveMode._
import org.apache.hudi.DataSourceReadOptions._
import org.apache.hudi.DataSourceWriteOptions._
import org.apache.hudi.config.HoodieWriteConfig._
```
- 定义表名，存储路径，生成测试数据

```
val tableName = "hudi_cow_table"
val basePath = "hdfs://hacluster/tmp/hudi_cow_table"
val dataGen = new DataGenerator
val inserts = convertToStringList(dataGen.generateInserts(10))
val df = spark.read.json(spark.sparkContext.parallelize(inserts, 2))
```

步骤5 写入Hudi表，模式为OVERWRITE。

```
df.write.format("org.apache.hudi").
options(getQuickstartWriteConfigs).
option(PRECOMBINE_FIELD_OPT_KEY, "ts").
option(RECORDKEY_FIELD_OPT_KEY, "uuid").
option(PARTITIONPATH_FIELD_OPT_KEY, "partitionpath").
option(TABLE_NAME, tableName).
mode(Overwrite).
save(basePath)
```

步骤6 查询Hudi表。

注册临时表并查询：

```
val roViewDF = spark.read.format("org.apache.hudi").load(basePath +
"/**/**/*")
roViewDF.createOrReplaceTempView("hudi_ro_table")
spark.sql("select fare, begin_lon, begin_lat, ts from hudi_ro_table where fare
> 20.0").show()
```

步骤7 生成更新数据并更新Hudi表，模式为APPEND。

```
val updates = convertToStringList(dataGen.generateUpdates(10))
```



```
val df = spark.read.json(spark.sparkContext.parallelize(updates, 1))
df.write.format("org.apache.hudi").
options(getQuickstartWriteConfigs).
option(PRECOMBINE_FIELD_OPT_KEY, "ts").
option(RECORDKEY_FIELD_OPT_KEY, "uuid").
option(PARTITIONPATH_FIELD_OPT_KEY, "partitionpath").
option(TABLE_NAME, tableName).
mode(Append).
save(basePath)
```

步骤8 查询Hudi表增量数据。

- 重新加载：

```
spark.read.format("org.apache.hudi").load(basePath + "/*/*/*/*").createOrReplaceTempView("hudi_ro_table")
```
- 进行增量查询：

```
val commits = spark.sql("select distinct(_hoodie_commit_time) as
commitTime from hudi_ro_table order by commitTime").map(k =>
k.getString(0)).take(50)
val beginTime = commits(commits.length - 2)
val incViewDF = spark.
read.
format("org.apache.hudi").
option(VIEW_TYPE_OPT_KEY, VIEW_TYPE_INCREMENTAL_OPT_VAL).
option(BEGIN_INSTANTTIME_OPT_KEY, beginTime).
load(basePath);
incViewDF.registerTempTable("hudi_incr_table")
spark.sql("select `_hoodie_commit_time`, fare, begin_lon, begin_lat, ts
from hudi_incr_table where fare > 20.0").show()
```

步骤9 进行指定时间点提交的查询。

```
val beginTime = "000"
val endTime = commits(commits.length - 2)
val incViewDF = spark.read.format("org.apache.hudi").
option(VIEW_TYPE_OPT_KEY, VIEW_TYPE_INCREMENTAL_OPT_VAL).
option(BEGIN_INSTANTTIME_OPT_KEY, beginTime).
option(END_INSTANTTIME_OPT_KEY, endTime).
load(basePath);
incViewDF.registerTempTable("hudi_incr_table")
spark.sql("select `_hoodie_commit_time`, fare, begin_lon, begin_lat, ts from
hudi_incr_table where fare > 20.0").show()
```

步骤10 删除数据。

- 准备删除的数据

```
val df = spark.sql("select uuid, partitionpath from hudi_ro_table limit 2")
val deletes = dataGen.generateDeletes(df.collectAsList())
```
- 执行删除操作

```
val df = spark.read.json(spark.sparkContext.parallelize(deletes, 2));
df.write.format("org.apache.hudi").
options(getQuickstartWriteConfigs).
option(OPERATION_OPT_KEY,"delete").
option(PRECOMBINE_FIELD_OPT_KEY, "ts").
option(RECORDKEY_FIELD_OPT_KEY, "uuid").
option(PARTITIONPATH_FIELD_OPT_KEY, "partitionpath").
option(TABLE_NAME, tableName).
mode(Append).
save(basePath);
```
- 重新查询数据

```
val roViewDFAfterDelete = spark.
read.
format("org.apache.hudi").
load(basePath + "/*/*/*")
roViewDFAfterDelete.createOrReplaceTempView("hudi_ro_table")
spark.sql("select uuid, partitionPath from hudi_ro_table").show()
```

----结束

11.3 使用 Hudi-Cli.sh 操作 Hudi 表

前提条件

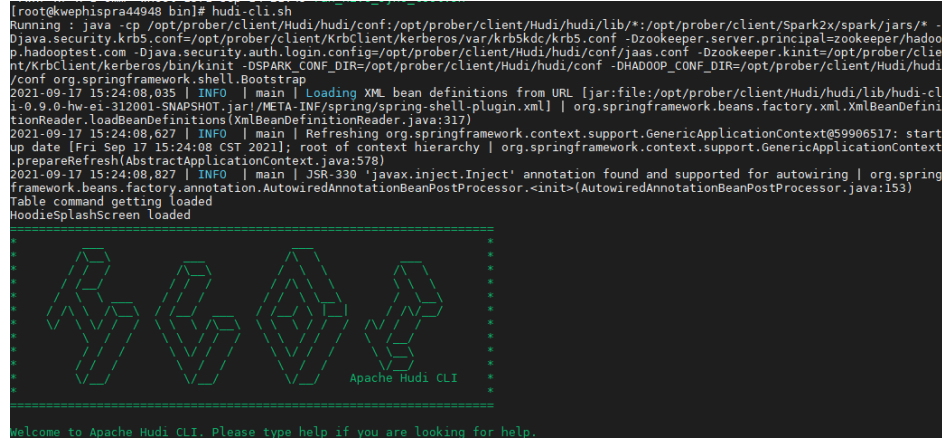
- 对于开启了Kerberos认证的安全模式集群，已在集群FusionInsight Manager界面创建一个用户并关联“hadoop”和“hive”用户组。
- 已下载并安装Hudi集群客户端。

基础操作

1. 使用root用户登录集群客户端节点，执行如下命令：

```
cd {客户端安装目录}
source bigdata_env
source Hudi/component_env
kinit 创建的用户
```
2. 执行hudi-cli.sh进入Hudi客户端，

```
cd {客户端安装目录}/Hudi/hudi/bin/
./hudi-cli.sh
```



3. 即可执行各种Hudi命令，执行示例（仅部分命令，全部命令请参考Hudi官网：<https://hudi.apache.org/docs/quick-start-guide/>）：

- 查看帮助：
help //查看hudi-cli的所有命令
help 'command' //查看某一个命令的帮助及参数列表。
- 连接表：
connect --path '/tmp/huditest/test_table'
- 查看表信息：
desc
- 查看compaction计划：
compactions show all
- 查看clean计划：
cleans show
- 执行clean：
cleans run
- 查看commit信息：
commits show
- 查看commit写入的分区：
commit showpartitions --commit 20210127153356

📖 说明

20210127153356表示commit的时间戳，下同。

- 查看指定commit写入的文件：
commit showfiles --commit 20210127153356
- 比较两个表的commit信息差异：
commits compare --path /tmp/hudimor/mytest100
- rollback指定提交（rollback每次只允许rollback最后一次commit）：
commit rollback --commit 20210127164905
- compaction调度：
**compaction schedule --hoodieConfigs
'hoodie.compaction.strategy=org.apache.hudi.table.action.compact.strateg**

- ```
y.BoundedIOCompactionStrategy,hoodie.compaction.target.io=1,hoodie.compact.inline.max.delta.commits=1'
```
- 执行compaction  
**compaction run --parallelism 100 --sparkMemory 1g --retry 1 --compactionInstant 20210602101315 --hoodieConfigs 'hoodie.compaction.strategy=org.apache.hudi.table.action.compact.strategy.BoundedIOCompactionStrategy,hoodie.compaction.target.io=1,hoodie.compact.inline.max.delta.commits=1' --propsFilePath hdfs://hacluster/tmp/default/tb\_test\_mor/.hoodie/hoodie.properties --schemaFilePath /tmp/default/tb\_test\_mor/.hoodie/compact\_tb\_base.json**
  - 创建savepoint  
**savepoint create --commit 20210318155750**
  - 回滚指定的savepoint  
**savepoint rollback --savepoint 20210318155750**

#### 注意

1. 若commit写入导致元数据冲突异常，执行commit rollback、savepoint rollback能回退数据，但不能回退Hive元数据，只能删除Hive表然后手动进行同步刷新。
2. commit rollback只能回退当前最新的一个commit，savepoint rollback只能回退到最新的一个savepoint。二者均不能随意指定进行回退。

## 11.4 Hudi 写操作

### 11.4.1 批量写入 Hudi 表

#### 操作场景

Hudi提供多种写入方式，具体见hoodie.datasource.write.operation配置项，这里主要介绍UPSERT、INSERT和BULK\_INSERT。

- INSERT（插入）：该操作流程和UPSERT基本一致，但是不需要通过索引去查询具体更新的文件分区，因此它的速度比UPSERT快。当数据源不包含更新数据时建议使用该操作，若数据源中存在更新数据，则在数据湖中会出现重复数据。
- BULK\_INSERT（批量插入）：用于初始数据集加载，该操作会对主键进行排序后直接以写普通parquet表的方式插入Hudi表，该操作性能是最高的，但是无法控制小文件，而UPSERT和INSERT操作使用启发式方法可以很好的控制小文件。
- UPSERT（插入更新）：默认操作类型。Hudi会根据主键进行判断，如果历史数据存在则update如果不存在则insert。因此在对于CDC之类几乎肯定包括更新的数据源，建议使用该操作。

#### 说明

- 由于INSERT时不会对主键进行排序，所以初始化数据集不建议使用INSERT。
- 在确定数据都为新数据时建议使用INSERT，当存在更新数据时建议使用UPSERT，当初始化数据集时建议使用BULK\_INSERT。

## 批量写入 Hudi 表

1. 引入Hudi包生成测试数据，参考[使用Spark Shell创建Hudi表](#)章节的[步骤2](#)到[步骤4](#)。

2. 写入Hudi表，写入命令中加入参数：

option("hoodie.datasource.write.operation", "bulk\_insert"), 指定写入方式为 bulk\_insert，如下所示：

```
df.write.format("org.apache.hudi").
options(getQuickstartWriteConfigs).
option("hoodie.datasource.write.precombine.field", "ts").
option("hoodie.datasource.write.recordkey.field", "uuid").
option("hoodie.datasource.write.partitionpath.field", "").
option("hoodie.datasource.write.operation", "bulk_insert").
option("hoodie.table.name", tableName).
option("hoodie.datasource.write.keygenerator.class",
"org.apache.hudi.keygen.NonpartitionedKeyGenerator").
option("hoodie.datasource.hive_sync.enable", "true").
option("hoodie.datasource.hive_sync.partition_fields", "").
option("hoodie.datasource.hive_sync.partition_extractor_class",
"org.apache.hudi.hive.NonPartitionedExtractor").
option("hoodie.datasource.hive_sync.table", tableName).
option("hoodie.datasource.hive_sync.use_jdbc", "false").
option("hoodie.bulkinsert.shuffle.parallelism", 4).
mode(Overwrite).
save(basePath)
```

### 📖 说明

- 示例中各参数介绍请参考[表11-4](#)。
- 使用spark datasource接口更新Mor表，Upsert写入小数据量时可能触发更新数据的小文件合并，使在Mor表的读优化视图中能查到部分更新数据。
- 当update的数据对应的base文件是小文件时，insert中的数据和update中的数据会被合在一起和base文件直接做合并产生新的base文件，而不是写log。

## 分区设置操作

Hudi支持多种分区方式，如多级分区、无分区、单分区、时间日期分区。用户可以根据实际需求选择合适的分区方式，接下来将详细介绍Hudi如何配置各种分区类型。

- 多级分区

多级分区即指定多个字段为分区键，需要注意的配置项：

| 配置项                                                   | 说明                                                      |
|-------------------------------------------------------|---------------------------------------------------------|
| hoodie.datasource.write.partitionpath.field           | 配置为多个分区字段，例如：p1，p2，p3。                                  |
| hoodie.datasource.hive_sync.partition_fields          | 和 hoodie.datasource.write.partitionpath.field的分区字段保持一致。 |
| hoodie.datasource.write.keygenerator.class            | 配置为 org.apache.hudi.keygen.ComplexKeyGenerator。         |
| hoodie.datasource.hive_sync.partition_extractor_class | 配置为 org.apache.hudi.hive.MultiPartKeyValueExtractor。    |

```
df.write.format("org.apache.hudi").
options(getQuickstartWriteConfigs).
option("hoodie.datasource.write.precombine.field", "ts").
option("hoodie.datasource.write.recordkey.field", "uuid").
option("hoodie.datasource.write.partitionpath.field", "p1,p2,p3").
option("hoodie.datasource.write.operation", "bulk_insert").
option("hoodie.table.name", tableName).
option("hoodie.datasource.write.keygenerator.class",
"org.apache.hudi.keygen.ComplexKeyGenerator").
option("hoodie.datasource.hive_sync.enable", "true").
option("hoodie.datasource.hive_sync.partition_fields", "p1,p2,p3").
option("hoodie.datasource.hive_sync.partition_extractor_class",
"org.apache.hudi.hive.MultiPartKeyValueExtractor").
option("hoodie.datasource.hive_sync.table", tableName).
option("hoodie.datasource.hive_sync.use_jdbc", "false").
option("hoodie.bulkinsert.shuffle.parallelism", 4).
mode(Overwrite).
save(basePath)
```

- 无分区

hudi支持无分区表，需要注意的配置项：

| 配置项                                                   | 说明                                                     |
|-------------------------------------------------------|--------------------------------------------------------|
| hoodie.datasource.write.partitionpath.field           | 配置为空。                                                  |
| hoodie.datasource.hive_sync.partition_fields          | 配置为空。                                                  |
| hoodie.datasource.write.keygenerator.class            | 配置为 org.apache.hudi.keygen.NonpartitionedKeyGenerator。 |
| hoodie.datasource.hive_sync.partition_extractor_class | 配置为 org.apache.hudi.hive.NonPartitionedExtractor。      |

```
df.write.format("org.apache.hudi").
options(getQuickstartWriteConfigs).
option("hoodie.datasource.write.precombine.field", "ts").
option("hoodie.datasource.write.recordkey.field", "uuid").
option("hoodie.datasource.write.partitionpath.field", "").
option("hoodie.datasource.write.operation", "bulk_insert").
option("hoodie.table.name", tableName).
option("hoodie.datasource.write.keygenerator.class",
"org.apache.hudi.keygen.NonpartitionedKeyGenerator").
option("hoodie.datasource.hive_sync.enable", "true").
option("hoodie.datasource.hive_sync.partition_fields", "").
option("hoodie.datasource.hive_sync.partition_extractor_class",
"org.apache.hudi.hive.NonPartitionedExtractor").
option("hoodie.datasource.hive_sync.table", tableName).
option("hoodie.datasource.hive_sync.use_jdbc", "false").
option("hoodie.bulkinsert.shuffle.parallelism", 4).
mode(Overwrite).
save(basePath)
```

- 单分区

和多级分区类似，需要配置项：

| 配置项                                                   | 说明                                                                                                     |
|-------------------------------------------------------|--------------------------------------------------------------------------------------------------------|
| hoodie.datasource.write.partitionpath.field           | 配置为一个字段，例如：p                                                                                           |
| hoodie.datasource.hive_sync.partition_fields          | 和 hoodie.datasource.write.partitionpath.field 分区字段保持一致。                                                |
| hoodie.datasource.write.keygenerator.class            | 默认可以配置为 org.apache.hudi.keygen.SimpleKeyGenerator 和 org.apache.hudi.keygen.ComplexKeyGenerator，也可以不配置。 |
| hoodie.datasource.hive_sync.partition_extractor_class | 配置为 org.apache.hudi.hive.MultiPartKeysValueExtractor。                                                  |

```

df.write.format("org.apache.hudi").
options(getQuickstartWriteConfigs).
option("hoodie.datasource.write.precombine.field", "ts").
option("hoodie.datasource.write.recordkey.field", "uuid").
option("hoodie.datasource.write.partitionpath.field", "p").
option("hoodie.datasource.write.operation", "bulk_insert").
option("hoodie.table.name", tableName).
option("hoodie.datasource.write.keygenerator.class",
"org.apache.hudi.keygen.ComplexKeyGenerator").
option("hoodie.datasource.hive_sync.enable", "true").
option("hoodie.datasource.hive_sync.partition_fields", "p").
option("hoodie.datasource.hive_sync.partition_extractor_class",
"org.apache.hudi.hive.MultiPartKeysValueExtractor").
option("hoodie.datasource.hive_sync.table", tableName).
option("hoodie.datasource.hive_sync.use_jdbc", "false").
option("hoodie.bulkinsert.shuffle.parallelism", 4).
mode(Overwrite).
save(basePath)

```

- 时间日期分区

即指定date类型字段作为分区字段，需要注意的配置项：

| 配置项                                          | 说明                                                                                                    |
|----------------------------------------------|-------------------------------------------------------------------------------------------------------|
| hoodie.datasource.write.partitionpath.field  | 配置为date类型字段。                                                                                          |
| hoodie.datasource.hive_sync.partition_fields | 配置为operationTime，和 hoodie.datasource.write.partitionpath.field 分区字段保持一致。                              |
| hoodie.datasource.write.keygenerator.class   | 默认配置为 org.apache.hudi.keygen.SimpleKeyGenerator，也可以不配置配置为 org.apache.hudi.keygen.ComplexKeyGenerator。 |

| 配置项                                                   | 说明                                                              |
|-------------------------------------------------------|-----------------------------------------------------------------|
| hoodie.datasource.hive_sync.partition_extractor_class | 配置 org.apache.hudi.hive.SlashEncodedDayPartitionValueExtractor。 |

#### 📖 说明

SlashEncodedDayPartitionValueExtractor存在以下约束：要求写入的日期格式为 yyyy/mm/dd。

- 分区排序：

| 配置项                                                  | 说明                          |
|------------------------------------------------------|-----------------------------|
| hoodie.bulkinsert.user.defined.partition_order.class | 指定分区排序类，可自行定义排序方法，具体参考样例代码。 |

#### 📖 说明

bulk\_insert默认字符排序，仅适用于StringType的主键。

## 11.4.2 流式写入 Hudi 表

### HoodieDeltaStreamer 流式写入

Hudi自带HoodieDeltaStreamer工具支持流式写入，也可以使用SparkStreaming以微批的方式写入。HoodieDeltaStreamer提供以下功能：

- 支持Kafka，DFS多种数据源接入。
- 支持管理检查点、回滚和恢复，保证exactly once语义。
- 支持自定义转换操作。

示例：

准备配置文件kafka-source.properties

```
#hudi配置
hoodie.datasource.write.recordkey.field=id
hoodie.datasource.write.partitionpath.field=age
hoodie.upsert.shuffle.parallelism=100
#hive config
hoodie.datasource.hive_sync.table=hudimor_deltastreamer_partition
hoodie.datasource.hive_sync.partition_fields=age
hoodie.datasource.hive_sync.partition_extractor_class=org.apache.hudi.hive.MultiPartKeyValueExtractor
hoodie.datasource.hive_sync.use_jdbc=false
hoodie.datasource.hive_sync.support_timestamp=true
Kafka Source topic
hoodie.deltastreamer.source.kafka.topic=hudimor_deltastreamer_partition
#checkpoint
hoodie.deltastreamer.checkpoint.provider.path=hdfs://hacluster/tmp/huditest/hudimor_deltastreamer_partition
Kafka props
The kafka cluster we want to ingest from
bootstrap.servers= xx.xx.xx.xx:xx
```



```
auto.offset.reset=earliest
#auto.offset.reset=latest
group.id=hoodie-delta-streamer
offset.rang.limit=10000
```

指定HoodieDeltaStreamer执行参数（具体参数配置，请查看官网<https://hudi.apache.org/>）执行如下命令：

```
spark-submit --master yarn
```

```
--jars /opt/hudi-java-examples-1.0.jar // 指定spark运行时需要的hudi jars路径
```

```
--driver-memory 1g
```

```
--executor-memory 1g --executor-cores 1 --num-executors 2 --conf spark.kryoserializer.buffer.max=128m
```

```
--driver-class-path /opt/client/Hudi/hudi/conf:/opt/client/Hudi/hudi/lib/*:/opt/client/Spark2x/spark/jars/*:/opt/hudi-examples-0.6.1-SNAPSHOT.jar:/opt/hudi-examples-0.6.1-SNAPSHOT-tests.jar // 指定spark driver需要的hudi jars路径
```

```
--class org.apache.hudi.utilities.deltastreamer.HoodieDeltaStreamer spark-internal
```

```
--props file:///opt/kafka-source.properties // 指定配置文件，注意：使用yarn-cluster模式提交任务时，请指定配置文件路径为HDFS路径。
```

```
--target-base-path /tmp/huditest/hudimor1_deltastreamer_partition // 指定hudi表路径
```

```
--table-type MERGE_ON_READ // 指定要写入的hudi表类型
```

```
--target-table hudimor_deltastreamer_partition // 指定hudi表名
```

```
--source-ordering-field name // 指定hudi表预合并列
```

```
--source-class org.apache.hudi.utilities.sources.JsonKafkaSource // 指定消费的数据源为JsonKafkaSource，该参数根据不同数据源指定不同的source类
```

```
--schemaprovider-class
```

```
com.huawei.bigdata.hudi.examples.DataSchemaProviderExample // 指定hudi表所需要的schema
```

```
--transformer-class com.huawei.bigdata.hudi.examples.TransformerExample // 指定如何处理数据源拉取来的数据，可根据自身业务需求做定制
```

```
--enable-hive-sync // 开启hive同步，同步hudi表到hive
```

```
--continuous // 指定流处理模式为连续模式
```

### 11.4.3 将 Hudi 表数据同步到 Hive

通过执行run\_hive\_sync\_tool.sh可以将Hudi表数据同步到Hive中。

例如：需要将HDFS上目录为hdfs://hacluster/tmp/huditest/hudimor1\_deltastreamer\_partition的Hudi表同步为Hive表，表名为table\_hive\_sync\_test3，使用unite、country和state为分区键，命令示例如下：

```
run_hive_sync_tool.sh --partitioned-by unite,country,state --base-path hdfs://hacluster/tmp/huditest/hudimor1_deltastreamer_partition --table
```

**hive\_sync\_test3 --partition-value-extractor  
org.apache.hudi.hive.MultiPartKeyValueExtractor --support-timestamp**

表 11-1 参数说明

| 命令                               | 描述                                                                                               | 必填 | 默认值                                    |
|----------------------------------|--------------------------------------------------------------------------------------------------|----|----------------------------------------|
| --database                       | Hive database名称                                                                                  | N  | default                                |
| --table                          | Hive表名                                                                                           | Y  | -                                      |
| --base-file-format               | 文件格式 (PARQUET或HFILE)                                                                             | N  | PARQUET                                |
| --user                           | Hive用户名                                                                                          | N  | -                                      |
| --pass                           | Hive密码                                                                                           | N  | -                                      |
| --jdbc-url                       | Hive jdbc connect url                                                                            | N  | -                                      |
| --base-path                      | 待同步的Hudi表存储路径                                                                                    | Y  | -                                      |
| --partitioned-by                 | 分区键-                                                                                             | N  | -                                      |
| --partition-value-extractor      | 分区类, 需实现PartitionValueExtractor, 可以从HDFS路径中提取分区值                                                 | N  | SlashEncodedDayPartitionValueExtractor |
| --assume-date-partitioning       | 以 yyyy/mm/dd进行分区从而支持向后兼容。                                                                        | N  | false                                  |
| --use-pre-apache-input-format    | 使用com.uber.hoodie包下的InputFormat替换org.apache.hudi包下的。除了从com.uber.hoodie迁移项目至org.apache.hudi外请勿使用。 | N  | false                                  |
| --use-jdbc                       | 使用Hive jdbc连接                                                                                    | N  | true                                   |
| --auto-create-database           | 自动创建Hive database                                                                                | N  | true                                   |
| --skip-ro-suffix                 | 注册时跳过读取_ro后缀的读优化视图                                                                               | N  | false                                  |
| --use-file-listing-from-metadata | 从Hudi的元数据中获取文件列表                                                                                 | N  | false                                  |
| --verify-metadata-file-listing   | 根据文件系统验证Hudi元数据中的文件列表                                                                            | N  | false                                  |
| --help、-h                        | 查看帮助                                                                                             | N  | false                                  |

| 命令                  | 描述                                              | 必填 | 默认值   |
|---------------------|-------------------------------------------------|----|-------|
| --support-timestamp | 将原始类型中'INT64'的TIMESTAMP_MICROS转换为Hive的timestamp | N  | false |
| --decode-partition  | 如果分区在写入过程中已编码，则解码分区值                            | N  | false |
| --batch-sync-num    | 指定每批次同步hive的分区数                                 | N  | 1000  |

### 📖 说明

Hive Sync时会判断表不存在时建外表并添加分区，表存在时对比表的schema是否存在差异，存在则替换，对比分区是否有新增，有则添加分区。

因此使用hive sync时有以下约束：

- 写入数据Schema只允许增加字段，不允许修改、删除字段。
- 分区目录只能新增，不会删除。
- Overwrite覆写Hudi表不支持同步覆盖Hive表。
- Hudi同步Hive表时，不支持使用timestamp类型作为分区列。

## 11.5 Hudi 读操作

### 11.5.1 读取 Hudi 数据概述

Hudi的读操作，作用于Hudi的三种视图之上，可以根据需求差异选择合适的视图进行查询。

Hudi 支持多种查询引擎Spark和Hive，具体支持矩阵见[表11-2](#)和[表11-3](#)。

表 11-2 cow 表

| 查询引擎                          | 实时视图/读优化视图 | 增量视图 |
|-------------------------------|------------|------|
| Hive                          | Y          | Y    |
| Spark ( SparkSQL )            | Y          | Y    |
| Spark ( SparkDataSource API ) | Y          | Y    |

表 11-3 mor 表

| 查询引擎 | 实时视图 | 增量视图 | 读优化视图 |
|------|------|------|-------|
| Hive | Y    | Y    | Y     |

| 查询引擎                                 | 实时视图 | 增量视图 | 读优化视图 |
|--------------------------------------|------|------|-------|
| Spark<br>( SparkSQL )                | Y    | Y    | Y     |
| Spark<br>( SparkDataSourc<br>e API ) | Y    | Y    | Y     |

**⚠ 注意**

- 当前Hudi使用Spark datasource接口读取时，不支持分区推断能力。比如bootstrap表使用datasource接口查询时，可能出现分区字段不显示，或者显示为null的情况。
- 增量视图，需设置set hoodie.hudicow.consume.mode = INCREMENTAL;，但该参数仅限于增量视图查询，不能用于Hudi表的其他类型查询，和其他表的查询。恢复配置可设置set hoodie.hudicow.consume.mode = SNAPSHOT;或任意值。

## 11.5.2 读取 Hudi cow 表视图

- 实时视图读取（Hive，SparkSQL为例）：直接读取Hive里面存储的Hudi表即可，**`\${table\_name}`**表示表名称。  

```
select count(*) from `${table_name}`;
```
- 实时视图读取（Spark dataSource API为例）：和读普通的数据Source表类似。必须指定查询类型QUERY\_TYPE\_OPT\_KEY 为 QUERY\_TYPE\_SNAPSHOT\_OPT\_VAL，**`\${table\_name}`**表示表名称。  

```
spark.read.format("hudi")
.option(QUERY_TYPE_OPT_KEY, QUERY_TYPE_SNAPSHOT_OPT_VAL) // 指定查询类型为实时视图模式
.load("/tmp/default/cow_bugx/") // 指定读取的hudi表路径
.createTempView("mycall")
spark.sql("select * from mycall").show(100)
```
- 增量视图读取（Hive为例，**`\${table\_name}`**表示表名称）：

```
set hoodie.`${table_name}`.consume.mode=INCREMENTAL; //设置增量读取模式
set hoodie.`${table_name}`.consume.max.commits=3; // 指定最大消费的commits数量
set hoodie.`${table_name}`.consume.start.timestamp=20201227153030; // 指定初始增量拉取commit
select count(*) from default.`${table_name}` where `hoodie_commit_time`>'20201227153030'; // 这个过滤条件必须加且值为初始增量拉取的commit。
```
- 增量视图读取（SparkSQL为例，**`\${table\_name}`**表示表名称）：

```
set hoodie.`${table_name}`.consume.mode=INCREMENTAL; //设置增量读取模式
set hoodie.`${table_name}`.consume.start.timestamp=20201227153030; // 指定初始增量拉取commit
set hoodie.`${table_name}`.consume.end.timestamp=20210308212318; // 指定增量拉取结束commit，如果不指定的话采用最新的commit
select count(*) from default.`${table_name}` where `hoodie_commit_time`>'20201227153030'; // 这个过滤条件必须加且值为初始增量拉取的commit。
```
- 增量视图读取（Spark dataSource API为例）：

必须指定查询类型QUERY\_TYPE\_OPT\_KEY 为增量模式  
QUERY\_TYPE\_INCREMENTAL\_OPT\_VAL

```
spark.read.format("hudi")
.option(QUERY_TYPE_OPT_KEY, QUERY_TYPE_INCREMENTAL_OPT_VAL) // 指定查询类型为增量模式
.option(BEGIN_INSTANTTIME_OPT_KEY, "20210308212004") // 指定初始增量拉取commit
.option(END_INSTANTTIME_OPT_KEY, "20210308212318") // 指定增量拉取结束commit
.load("/tmp/default/cow_bugx/") // 指定读取的hudi表路径
```

```
.createTempView("mycall") // 注册为spark临时表
spark.sql("select * from mycall where `_hoodie_commit_time`>'20210308211131'")// 开始查询, 和hive
增量查询语句一样
.show(100, false)
```

- 读优化视图：cow表读优化视图等同于实时视图。

### 11.5.3 读取 Hudi mor 表视图

mor表同步给Hive后, 会在Hive表中同步出：“表名+后缀\_rt”和“表名+后缀\_ro”两张表。其中后缀为rt表代表实时视图, 后缀为ro的表代表读优化视图。例如：同步给Hive的hudi表名为``${table_name}``, 同步Hive后hive表中多出两张表分别为``${table_name}_rt``和``${table_name}_ro``。

- 实时视图读取（Hive, SparkSQL为例）：直接读取Hive里面存储的后缀为\_rt的hudi表即可。

```
select count(*) from `${table_name}_rt`;
```
- 实时视图读取（Spark dataSource API为例）：和cow表一样, 请参考cow表相关操作。
- 增量视图读取（hive为例）：

```
set hive.input.format=org.apache.hudi.hadoop.hive.HoodieCombineHiveInputFormat; // sparksql 不需要指定
set hoodie.`${table_name}`.consume.mode=INCREMENTAL;
set hoodie.`${table_name}`.consume.max.commits=3;
set hoodie.`${table_name}`.consume.start.timestamp=20201227153030;
select count(*) from default.`${table_name}_rt` where `_hoodie_commit_time`>'20201227153030';
```
- 增量视图读取（SparkSQL为例）：

```
set hoodie.`${table_name}`.consume.mode=INCREMENTAL;
set hoodie.`${table_name}`.consume.start.timestamp=20201227153030; // 指定初始增量拉取commit
set hoodie.`${table_name}`.consume.end.timestamp=20210308212318; // 指定增量拉取结束commit, 如果不指定的话采用最新的commit
select count(*) from default.`${table_name}_rt` where `_hoodie_commit_time`>'20201227153030';
```
- 增量视图（Spark dataSource API为例）：和cow表一样, 请参考cow表相关操作。
- 读优化视图读取（Hive, SparkSQL为例）：直接读取Hive里面存储的后缀为\_ro的hudi表即可。

```
select count(*) from `${table_name}_ro`;
```
- 读优化视图读取（Spark dataSource API为例）：和读普通的数据Source表类似。  
必须指定查询类型`QUERY_TYPE_OPT_KEY`为`QUERY_TYPE_READ_OPTIMIZED_OPT_VAL`

```
spark.read.format("hudi")
.option(QUERY_TYPE_OPT_KEY, QUERY_TYPE_READ_OPTIMIZED_OPT_VAL) // 指定查询类型为读优化视图
.load("/tmp/default/mor_bugx/") // 指定读取的hudi表路径
.createTempView("mycall")
spark.sql("select * from mycall").show(100)
```

## 11.6 数据管理维护

### 11.6.1 Hudi Clustering 操作说明

#### 什么是 Clustering

即数据布局, 该服务可重新组织数据以提高查询性能, 也不会影响摄取速度。

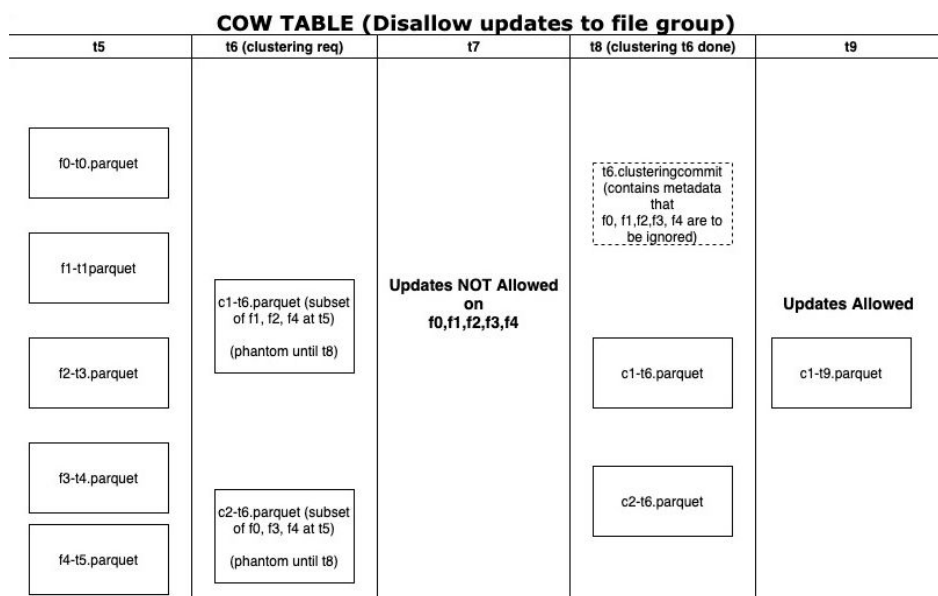
## Clustering 架构

Hudi通过其写入客户端API提供了不同的操作，如insert/upsert/bulk\_insert来将数据写入Hudi表。为了能够在文件大小和入湖速度之间进行权衡，Hudi提供了一个hoodie.parquet.small.file.limit配置来设置最小文件大小。用户可以将该配置设置为“0”，以强制新数据写入新的文件组，或设置为更高的值以确保新数据被“填充”到现有小的文件组中，直到达到指定大小为止，但其会增加摄取延迟。

为能够支持快速摄取的同时不影响查询性能，引入了Clustering服务来重写数据以优化Hudi数据湖文件的布局。

Clustering服务可以异步或同步运行，Clustering会添加了一种新的REPLACE操作类型，该操作类型将在Hudi元数据时间轴中标记Clustering操作。

Clustering服务基于Hudi的MVCC设计，允许继续插入新数据，而Clustering操作在后台运行以重新格式化数据布局，从而确保并发读写者之间的快照隔离。



总体而言Clustering分为两个部分：

- 调度Clustering：使用可插拔的Clustering策略创建Clustering计划。
  - a. 识别符合Clustering条件的文件：根据所选的Clustering策略，调度逻辑将识别符合Clustering条件的文件。
  - b. 根据特定条件对符合Clustering条件的文件进行分组。每个组的数据大小应为targetFileSize的倍数。分组是计划中定义的"策略"的一部分。此外还有一个选项可以限制组大小，以改善并行性并避免混排大量数据。
  - c. 将Clustering计划以avro元数据格式保存到时间线。
- 执行Clustering：使用执行策略处理计划以创建新文件并替换旧文件。
  - a. 读取Clustering计划，并获得ClusteringGroups，其标记了需要进行Clustering的文件组。
  - b. 对于每个组使用strategyParams实例化适当的策略类（例如：sortColumns），然后应用该策略重写数据。
  - c. 创建一个REPLACE提交，并更新HoodieReplaceCommitMetadata中的元数据。

## 如何执行 Clustering

1. 同步执行Clustering配置。

在写入时加上配置参数：

```
option("hoodie.clustering.inline", "true").
```

```
option("hoodie.clustering.inline.max.commits", "4").
```

```
option("hoodie.clustering.plan.strategy.target.file.max.bytes",
"1073741824").
```

```
option("hoodie.clustering.plan.strategy.small.file.limit", "629145600").
```

```
option("hoodie.clustering.plan.strategy.sort.columns",
"column1,column2").
```

2. 异步执行Clustering：

```
spark-submit --master yarn --class
```

```
org.apache.hudi.utilities.HoodieClusteringJob /opt/client/Hudi/hudi/lib/
hudi-utilities*.jar --schedule --base-path <table_path> --table-name
<table_name> --props /tmp/clusteringjob.properties --spark-memory 1g
```

```
spark-submit --master yarn --driver-memory 16G --executor-memory 12G
--executor-cores 4 --num-executors 4 --class
```

```
org.apache.hudi.utilities.HoodieClusteringJob /opt/client/Hudi/hudi/lib/
hudi-utilities*.jar --base-path <table_path> --instant-time
20210605112954 --table-name <table_name> --props /tmp/
clusteringjob.properties --spark-memory 12g
```

3. 指定clustering的排序方式和排序列：

当前clustering支持linear、z-order、hilbert 三种排序方式，可以通过option方式或者set方式来设置。

- linear：普通排序，默认排序，适合排序一个字段，或者多个低级字段。
- z-order和hilbert：多维排序，需要指定“hoodie.layout.optimize.strategy”为z-order或者hilbert。

适合排序多个字段，例如查询条件中涉及到多个字段。推荐排序字段的个数2到4个。

hilbert多维排序效果比z-order好，但是排序效率没z-order高。

详细配置请参考[Hudi常见配置参数](#)。

---

### 注意

1. Clustering的排序列不允许值存在null，是spark rdd的限制。
  2. 当target.file.max.bytes的值较大时，启动Clustering执行需要提高--spark-memory，否则会导致executor内存溢出。
  3. 当前clean不支持清理Clustering失败后的垃圾文件。
  4. Clustering后可能出现新文件大小不等引起数据倾斜的情况。
  5. cluster不支持和upsert并发。
  6. 如果clustering处于inflight状态，该FileGroup下的文件不支持Update操作。
-

## 11.6.2 Hudi Cleaning 操作说明

Cleaning用于清理不再需要的版本数据。

Hudi使用Cleaner后台作业，不断清除不需要的旧版本数据。通过配置hoodie.cleaner.policy和hoodie.cleaner.commits.retained可以使用不同的清理策略和保存的commit数量。

执行cleaning有两种方式：

- 同步clean由参数hoodie.clean.automatic控制，默认自动开启。  
关闭同步clean：  
datasource写入时可以通过`.option("hoodie.clean.automatic", "false")`来关闭自动clean。  
spark-sql写入时可以通过`set hoodie.clean.automatic=false;`来关闭自动clean。
- 异步clean可以使用spark-sql来执行。

更多clean相关参数请参考[compaction&cleaning配置](#)章节。

## 11.6.3 Hudi Compaction 操作说明

Compaction用于合并mor表Base和Log文件。

对于Merge-On-Read表，数据使用列式Parquet文件和行式Avro文件存储，更新被记录到增量文件，然后进行同步/异步compaction生成新版本的列式文件。Merge-On-Read表可减少数据摄入延迟，因而进行不阻塞摄入的异步Compaction很有意义。

- 异步Compaction会进行如下两个步骤：
  - a. 调度Compaction：由入湖作业完成，在这一步，Hudi扫描分区并选出待进行compaction的FileSlice，最后CompactionPlan会写入Hudi的Timeline。
  - b. 执行Compaction：一个单独的进程/线程将读取CompactionPlan并对FileSlice执行Compaction操作。
- 使用Compaction的方式分为同步和异步两种：
  - 同步方式由参数hoodie.compact.inline控制，默认为true，自动生成compaction调度计划并执行compaction：
    - 关闭同步compaction  
datasource写入时可以通过`.option("hoodie.compact.inline", "false")`来关闭自动compaction。  
spark-sql写入时可以通过`set hoodie.compact.inline=false;`来关闭自动compaction。
    - 仅同步生成compaction调度而不执行compaction
      - ·datasource写入时可以通过以下option参数来实现：  
`option("hoodie.compact.inline", "true").`  
`option("hoodie.schedule.compact.only.inline", "true").`  
`option("hoodie.run.compact.only.inline", "false").`
      - ·spark-sql写入时可以通过set 以下参数来实现：  
`set hoodie.compact.inline=true;`  
`set hoodie.schedule.compact.only.inline=true;`



```
set hoodie.run.compact.only.inline=false;
```

- 异步方式由spark-sql来实现。

如果需要在异步compaction时只执行已经产生的compaction调度计划而不创建新的调度计划，则需要通过set命令设置以下参数：

```
set hoodie.compact.inline=true;
```

```
set hoodie.schedule.compact.only.inline=false;
```

```
set hoodie.run.compact.only.inline=true;
```

更多compaction参数请参考[compaction&cleaning配置](#)章节。

#### 📖 说明

为了保证入湖的最高效率，推荐使用同步产生compaction调度计划，异步执行compaction调度计划的方式。

## 11.6.4 Hudi Savepoint 操作说明

Savepoint用于保存并还原自定义的版本数据。

Hudi提供的savepoint就可以将不同的commit保存起来以便清理程序不会将其删除，后续可以使用Rollback进行恢复。

使用spark-sql管理savepoint。

示例如下：

- 创建savepoint  

```
call create_savepoint('hudi_test1', '20220908155421949');
```
- 查看所有存在的savepoint  

```
call show_savepoints(table =>'hudi_test1');
```
- 回滚savepoint  

```
call rollback_to_savepoint('hudi_test1', '20220908155421949');
```

#### 📖 说明

MoR表暂时不支持savepoint。

## 11.7 Hudi 常见配置参数

### 11.7.1 写入操作配置

本章节介绍Hudi重要配置的详细信息，更多配置请参考hudi官网：<http://hudi.apache.org/cn/docs/configurations.html>。

表 11-4 写入操作重要配置项

| 参数                                 | 描述           | 默认值 |
|------------------------------------|--------------|-----|
| hoodie.datasource.write.table.name | 指定写入的hudi表名。 | 无   |

| 参数                                              | 描述                                                                                                                                                                                                                                                                                                                                                                                                                                                       | 默认值                                                     |
|-------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------|
| hoodie.datasource.write.operation               | <p>写hudi表指定的操作类型，当前支持upsert、delete、insert、bulk_insert等方式。</p> <ul style="list-style-type: none"> <li>• upsert: 更新插入混合操作</li> <li>• delete: 删除操作</li> <li>• insert: 插入操作</li> <li>• bulk_insert: 用于初始建表导入数据，注意初始建表禁止使用upsert、insert方式</li> <li>• insert_overwrite: 对静态分区执行insert overwrite</li> <li>• insert_overwrite_table: 动态分区执行insert overwrite，该操作并不会立刻删除全表做overwrite，会逻辑上重写hudi表的元数据，无用数据后续由hudi的clean机制清理。效率比bulk_insert + overwrite 高</li> </ul> | upsert                                                  |
| hoodie.datasource.write.table.type              | 指定hudi表类型，一旦这个表类型被指定，后续禁止修改该参数，可选值MERGE_ON_READ。                                                                                                                                                                                                                                                                                                                                                                                                         | COPY_ON_WRITE                                           |
| hoodie.datasource.write.precombine.field        | 该值用于在写之前对具有相同的key的行进行合并去重。                                                                                                                                                                                                                                                                                                                                                                                                                               | 指定为具体的表字段                                               |
| hoodie.datasource.write.payload.class           | 在更新过程中，该类用于提供方法将要更新的记录和更新的记录做合并，该实现可插拔，如要实现自己的合并逻辑，可自行编写。                                                                                                                                                                                                                                                                                                                                                                                                | org.apache.hudi.common.model.DefaultHoodieRecordPayload |
| hoodie.datasource.write.recordkey.field         | 用于指定hudi的主键，hudi表要求有唯一主键。                                                                                                                                                                                                                                                                                                                                                                                                                                | 指定为具体的表字段                                               |
| hoodie.datasource.write.partitionpath.field     | 用于指定分区键，该值配合hoodie.datasource.write.keygenerator.class使用可以满足不同的分区场景。                                                                                                                                                                                                                                                                                                                                                                                     | 无                                                       |
| hoodie.datasource.write.hive_style_partitioning | 用于指定分区方式是否和hive保持一致，建议该值设置为true。                                                                                                                                                                                                                                                                                                                                                                                                                         | true                                                    |

| 参数                                         | 描述                                                                                                                                                       | 默认值                                        |
|--------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------|
| hoodie.datasource.write.keygenerator.class | 配合 hoodie.datasource.write.partitionpath.field, hoodie.datasource.write.recordkey.field产生主键和分区方式。<br><b>说明</b><br>写入设置KeyGenerator与表保存的参数值不一致时将提示需要保持一致。 | org.apache.hudi.keygen.ComplexKeyGenerator |

## 11.7.2 同步 Hive 表配置

| 参数                                                    | 描述                                                                       | 默认值                                                         |
|-------------------------------------------------------|--------------------------------------------------------------------------|-------------------------------------------------------------|
| hoodie.datasource.hive_sync.enable                    | 是否同步hudi表信息到hive metastore。<br><b>注意</b><br>建议该值设置为true，统一使用hive管理hudi表。 | false                                                       |
| hoodie.datasource.hive_sync.database                  | 要同步给hive的数据库名。                                                           | default                                                     |
| hoodie.datasource.hive_sync.table                     | 要同步给hive的表名，建议这个值和 hoodie.datasource.write.table.name保证一致。               | unknown                                                     |
| hoodie.datasource.hive_sync.username                  | 同步hive时，指定的用户名。                                                          | hive                                                        |
| hoodie.datasource.hive_sync.password                  | 同步hive时，指定的密码。                                                           | hive                                                        |
| hoodie.datasource.hive_sync.jdbcurl                   | 连接hive jdbc指定的连接。                                                        | ""                                                          |
| hoodie.datasource.hive_sync.use_jdbc                  | 是否使用hive jdbc方式连接hive同步hudi表信息。建议该值设置为false，设置为false后 jdbc连接相关配置无效。      | true                                                        |
| hoodie.datasource.hive_sync.partition_fields          | 用于决定hive分区列。                                                             | ""                                                          |
| hoodie.datasource.hive_sync.partition_extractor_class | 用于提取hudi分区列值，将其转换成hive分区列。                                               | org.apache.hudi.hive.SlashEncodedDayPartitionValueExtractor |

| 参数                                            | 描述                                                                                                                                        | 默认值  |
|-----------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------|------|
| hoodie.datasource.hive_sync.support_timestamp | 当hudi表存在timestamp类型字段时，需指定此参数为true，以实现同步timestamp类型到hive元数据中。该值默认为false，默认将timestamp类型同步为bigInt，默认情况可能导致使用sql查询包含timestamp类型字段的hudi表出现错误。 | true |

### 11.7.3 index 相关配置

| 参数                             | 描述                                                                                                                                                                                                           | 默认值        |
|--------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------|
| hoodie.index.class             | 用户自定义索引的全路径名，索引类必须为HoodieIndex的子类，当指定该配置时，其会优先于hoodie.index.type配置。                                                                                                                                          | ""         |
| hoodie.index.type              | 使用的索引类型，默认为布隆过滤器。可能的选项是[BLOOM   HBASE   GLOBAL_BLOOM   SIMPLE   GLOBAL_SIMPLE]。布隆过滤器消除了对外部系统的依赖，并存储在Parquet数据文件的页脚中。                                                                                         | BLOOM      |
| hoodie.index.bloom.num_entries | 存储在布隆过滤器中的条目数。假设maxParquetFileSize为128MB，averageRecordSize为1024B，因此，一个文件中的记录总数约为130K。默认值（60000）大约是此近似值的一半。<br><b>注意</b><br>将此值设置的太低，将产生很多误报，并且索引查找将必须扫描比其所需的更多的文件；如果将其设置的非常高，将线性增加每个数据文件的大小（每50000个条目大约4KB）。 | 60000      |
| hoodie.index.bloom.fpp         | 根据条目数允许的误差率。用于计算应为布隆过滤器分配多少位以及哈希函数的数量。通常将此值设置得很低（默认值：0.00000001），在磁盘空间上进行权衡以降低误报率。                                                                                                                           | 0.00000001 |

| 参数                                       | 描述                                                                                                                                              | 默认值      |
|------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------|----------|
| hoodie.bloom.index.parallelism           | 索引查找的并行度，其中涉及 Spark Shuffle。默认情况下，根据输入的工作负载特征自动计算的。                                                                                             | 0        |
| hoodie.bloom.index.prune.by.ranges       | 为true时，从文件框定信息，可以加快索引查找的速度。如果键具有单调递增的前缀，例如时间戳，则特别有用。                                                                                            | true     |
| hoodie.bloom.index.use.caching           | 为true时，将通过减少用于计算并行度或受影响分区的IO来缓存输入的RDD以加快索引查找。                                                                                                   | true     |
| hoodie.bloom.index.use.trebased.filter   | 为true时，启用基于间隔树的文件过滤优化。与暴力模式相比，此模式可根据键范围加快文件过滤速度。                                                                                                | true     |
| hoodie.bloom.index.bucketized.checking   | 为true时，启用了桶式布隆过滤。这减少了在基于排序的布隆索引查找中看到的偏差。                                                                                                        | true     |
| hoodie.bloom.index.keys.per.bucket       | 仅在启用 bloomIndexBucketizedChecking并且索引类型为bloom的情况下适用。<br>此配置控制“存储桶”的大小，该大小可跟踪对单个文件进行的记录键检查的次数，并且是分配给执行布隆过滤器查找的每个分区的工作单位。较高的值将分摊将布隆过滤器读取到内存的固定成本。 | 10000000 |
| hoodie.bloom.index.update.partition.path | 仅在索引类型为 GLOBAL_BLOOM时适用。<br>为true时，当对一个已有记录执行包含分区路径的更新操作时，将会导致把新记录插入到新分区，而把原有记录从旧分区里删除。为false时，只对旧分区的原有记录进行更新。                                  | true     |
| hoodie.index.hbase.zk.quorum             | 仅在索引类型为HBASE时适用，必填选项。要连接的HBase ZK Quorum URL。                                                                                                   | 无        |
| hoodie.index.hbase.zk.port               | 仅在索引类型为HBASE时适用，必填选项。要连接的HBase ZK Quorum端口。                                                                                                     | 无        |

| 参数                                 | 描述                                                                                            | 默认值 |
|------------------------------------|-----------------------------------------------------------------------------------------------|-----|
| hoodie.index.hbase.zk<br>node.path | 仅在索引类型为HBASE时适用，必填选项。这是根znode，它将包含HBase创建及使用的所有znode。                                         | 无   |
| hoodie.index.hbase.ta<br>ble       | 仅在索引类型为HBASE时适用，必填选项。HBase表名称，用作索引。Hudi将row_key和 [partition_path, fileId, commitTime]映射存储在表中。 | 无   |

## 11.7.4 存储配置

| 参数                                     | 描述                                                                                       | 默认值                       |
|----------------------------------------|------------------------------------------------------------------------------------------|---------------------------|
| hoodie.parquet.max.fil<br>e.size       | Hudi写阶段生成的parquet文件的目标大小。对于DFS，这需要与基础文件系统块大小保持一致，以实现最佳性能。                                | 120 * 1024 * 1024<br>byte |
| hoodie.parquet.block.s<br>ize          | parquet页面大小，页面是parquet文件中的读取单位，在一个块内，页面被分别压缩。                                            | 120 * 1024 * 1024<br>byte |
| hoodie.parquet.compr<br>ession.ratio   | 当Hudi尝试调整新parquet文件的大小时，预期对parquet数据进行压缩的比例。如果bulk_insert生成的文件小于预期大小，请增加此值。              | 0.1                       |
| hoodie.parquet.compr<br>ession.codec   | parquet压缩编解码方式名称，默认值为gzip。可能的选项是 [gzip   snappy   uncompressed   lzo]                    | snappy                    |
| hoodie.logfile.max.size                | LogFile的最大值。这是在将日志文件移到下一个版本之前允许的最大值。                                                     | 1GB                       |
| hoodie.logfile.data.blo<br>ck.max.size | LogFile数据块的最大值。这是允许将单个数据块附加到日志文件的最大值。这有助于确保附加到日志文件的数据被分解为可调整大小的块，以防止发生OOM错误。此大小应大于JVM内存。 | 256MB                     |

| 参数                                          | 描述                                                                                  | 默认值  |
|---------------------------------------------|-------------------------------------------------------------------------------------|------|
| hoodie.logfile.to.parquet.compression.ratio | 随着记录从日志文件移动到parquet，预期会进行额外压缩的比例。用于merge_on_read存储，以将插入内容发送到日志文件中并控制压缩parquet文件的大小。 | 0.35 |

### 11.7.5 compaction&cleaning 配置

| 参数                                  | 描述                                                                                                                                | 默认值                     |
|-------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------|-------------------------|
| hoodie.clean.automat<br>ic          | 是否执行自动clean。                                                                                                                      | true                    |
| hoodie.cleaner.policy               | 要使用的清理策略。Hudi将删除旧版本的parquet文件以回收空间。任何引用此版本文件的查询和计算都将失败。需要确保数据保留的时间超过最大查询执行时间。                                                     | KEEP_LATEST_COMMIT<br>S |
| hoodie.cleaner.commi<br>ts.retained | 保留的提交数。因此，数据将保留为num_of_commits * time_between_commits（计划的），这也直接转化为逐步提取此数据集的数量。                                                    | 10                      |
| hoodie.keep.max.com<br>mits         | 触发归档操作的commit数阈值                                                                                                                  | 30                      |
| hoodie.keep.min.com<br>mits         | 归档操作保留的commit数。                                                                                                                   | 20                      |
| hoodie.commits.archiv<br>al.batch   | 这控制着批量读取并一起归档的提交即时的数量。                                                                                                            | 10                      |
| hoodie.parquet.small.f<br>ile.limit | 该值应小于maxFileSize，如果将其设置为0，会关闭此功能。由于批处理中分区中插入记录的数量众多，总会出现小文件。Hudi提供了一个选项，可以通过将该分区中的插入作为对现有小文件的更新来解决小文件的问题。此处的大小是被视为“小文件大小”的最小文件大小。 | 104857600 byte          |

| 参数                                      | 描述                                                                                                                                                      | 默认值    |
|-----------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------|--------|
| hoodie.copyonwrite.insert.split.size    | 插入写入并行度。为单个分区的总共插入次数。写出100MB的文件，至少1KB大小的记录，意味着每个文件有100K记录。默认值是超额配置为500K。为了改善插入延迟，请对其进行调整以匹配单个文件中的记录数。将此值设置为较小的值将导致文件变小（尤其是当compactionSmallFileSize为0时）。 | 500000 |
| hoodie.copyonwrite.insert.auto.split    | Hudi是否应该基于最后24个提交的元数据动态计算insertSplitSize，默认关闭。                                                                                                          | true   |
| hoodie.copyonwrite.record.size.estimate | 平均记录大小。如果指定，Hudi将使用它，并且不会基于最后24个提交的元数据动态地计算。没有默认值设置。这对于计算插入并行度以及将插入打包到小文件中至关重要。                                                                         | 1024   |
| hoodie.compact.inline                   | 当设置为true时，紧接在插入或插入更新或批量插入的提交或增量提交操作之后由摄取本身触发压缩。                                                                                                         | true   |
| hoodie.compact.inline.max.delta.commits | 触发内联压缩之前要保留的最大增量提交数。                                                                                                                                    | 5      |
| hoodie.compaction.lazy.block.read       | 当CompactedLogScanner合并所有日志文件时，此配置有助于选择是否应延迟读取日志块。选择true以使用I/O密集型延迟块读取（低内存使用），或者为false来使用内存密集型立即块读取（高内存使用）。                                              | true   |
| hoodie.compaction.reverse.log.read      | HoodieLogFormatReader会从pos=0到pos=file_length向前读取日志文件。如果此配置设置为true，则Reader会从pos=file_length到pos=0反向读取日志文件。                                               | false  |
| hoodie.cleaner.parallelism              | 如果清理变慢，请增加此值。                                                                                                                                           | 200    |



| 参数                                           | 描述                                                                                     | 默认值                                                                              |
|----------------------------------------------|----------------------------------------------------------------------------------------|----------------------------------------------------------------------------------|
| hoodie.compaction.strategy                   | 用来决定在每次压缩运行期间选择要压缩的文件组的压缩策略。默认情况下，Hudi选择具有累积最多未合并数据的日志文件。                              | org.apache.hudi.table.action.compact.strategy.<br>LogFileBasedCompactionStrategy |
| hoodie.compaction.target.io                  | LogFileBasedCompactionStrategy的压缩运行期间要花费的MB量。当压缩以内联模式运行时，此值有助于限制摄取延迟。                  | 500 * 1024 MB                                                                    |
| hoodie.compaction.daybased.target.partitions | 由 org.apache.hudi.io.compact.strategy.DayBasedCompactionStrategy使用，表示在压缩运行期间要压缩的最新分区数。 | 10                                                                               |
| hoodie.compaction.payload.class              | 这需要与插入/插入更新过程中使用的类相同。就像写入一样，压缩也使用记录有效负载类将日志中的记录彼此合并，再次与基本文件合并，并生成压缩后要写入的最终记录。          | org.apache.hudi.common.model.Defaulthoodierecordpayload                          |
| hoodie.schedule.compact.only.inline          | 在写入操作时，是否只生成压缩计划。在 hoodie.compact.inline=true时有效。                                      | false                                                                            |
| hoodie.run.compact.only.inline               | 通过Sql执行run compact命令时，是否只执行压缩操作，压缩计划不存在时直接退出。                                          | false                                                                            |

## 11.7.6 单表并发控制配置

| 参数                                        | 描述                                                                              | 默认值                                                                |
|-------------------------------------------|---------------------------------------------------------------------------------|--------------------------------------------------------------------|
| hoodie.write.lock.provider                | 指定lock provider，不建议使用默认值，使用 org.apache.hudi.hive.HiveMetastoreBasedLockProvider | org.apache.hudi.client.transaction.lock.ZookeeperBasedLockProvider |
| hoodie.write.lock.hive.metastore.database | Hive的database                                                                   | 无                                                                  |
| hoodie.write.lock.hive.metastore.table    | Hive的table name                                                                 | 无                                                                  |
| hoodie.write.lock.client.num_retries      | 重试次数                                                                            | 10                                                                 |

| 参数                                                  | 描述                                              | 默认值                                                                                     |
|-----------------------------------------------------|-------------------------------------------------|-----------------------------------------------------------------------------------------|
| hoodie.write.lock.client.wait_time_ms_between_retry | 重试间隔                                            | 10000                                                                                   |
| hoodie.write.lock.conflict.resolution.strategy      | lock provider类，必须是ConflictResolutionStrategy的子类 | org.apache.hudi.client.transaction.SimpleConcurrentFileWritesConflictResolutionStrategy |
| hoodie.write.lock.zookeeper.base_path               | 存放ZNodes的路径，同一张表的并发写入需配置一致                      | 无                                                                                       |
| hoodie.write.lock.zookeeper.lock_key                | ZNode的名称，建议与Hudi表名相同                            | 无                                                                                       |
| hoodie.write.lock.zookeeper.connection_timeout_ms   | zk连接超时时间                                        | 15000                                                                                   |
| hoodie.write.lock.zookeeper.port                    | zk端口号                                           | 无                                                                                       |
| hoodie.write.lock.zookeeper.url                     | zk的url                                          | 无                                                                                       |
| hoodie.write.lock.zookeeper.session_timeout_ms      | zk的session过期时间                                  | 60000                                                                                   |

## 11.8 Hudi 性能调优

### 性能调优方式

当前版本Hudi写入操作主推Spark，因此Hudi的调优和Spark比较类似，可参考[Spark Core性能调优](#)。

### 推荐资源配置

- mor表：  
由于其本质上是写增量文件，调优可以直接根据hudi的数据大小（dataSize）进行调整。  
dataSize如果只有几个G，推荐跑单节点运行spark，或者yarn模式但是只分配一个container。  
入湖程序的并行度p设置：建议  $p = (\text{dataSize}) / 128\text{M}$ ，程序分配core的数量保持和p一致即可。内存设置建议内存大小和core的比例大于1.5:1 即一个core配1.5G内存，堆外内存设置建议内存大小和core的比例大于0.5:1。
- cow表：

cow表的原理是重写原始数据，因此这种表的调优，要兼顾dataSize和最后重写的文件数量。总体来说core数量越大越好（和最后重写多少个文件数直接相关），并行度p和内存大小和mori设置类似。

## 11.9 Hudi 常见问题

### 11.9.1 数据写入

#### 11.9.1.1 写入更新数据时报错 Parquet/Avro schema

##### 问题

数据写入时报错：

```
org.apache.parquet.io.InvalidRecordException: Parquet/Avro schema mismatch: Avro field 'col1' not found
```

##### 回答

建议在使用Hudi时，schema应该以向后兼容的方式演进。此错误通常发生在使用向后不兼容的演进方式删除某些列如“col1”后，更新parquet文件中以旧的schema写入的列“col1”，在这种情况下，parquet尝试在传入记录中查找所有当前字段，当发现“col1”不存在时，发生上述异常。

解决这个问题的办法是使用所有schema演进版本来创建uber schema，并使用该schema作为target schema。用户可以从hive metastore中获取schema并将其与当前schema合并。

#### 11.9.1.2 写入更新数据时报错 UnsupportedOperationException

##### 问题

数据写入时报错：

```
java.lang.UnsupportedOperationException: org.apache.parquet.avro.AvroConverters$FieldIntegerConverter
```

##### 回答

因为schema演进以非向后兼容的方式进行，此错误将再次发生。基本上，如果已经写入Hudi数据集parquet文件的记录R有一些更新U。R包含字段F，该字段包含某类数据类型，也就是LONG。U具有相同的字段F，该字段的数据类型是INT。Parquet FS不支持这种不兼容的数据类型转换。

对于此类错误，请从源头数据采集的位置进行有效的数据类型转换。

#### 11.9.1.3 写入更新数据时报错 SchemaCompatabilityException

##### 问题

数据写入时报错：

```
org.apache.hudi.exception.SchemaCompatibilityException: Unable to validate the rewritten record <record>
against schema <schema>at
org.apache.hudi.common.util.HoodieAvroUtils.rewrite(HoodieAvroUtils.java:215)
```

## 回答

如果schema包含non-nullable字段但是值是不存在或者null，则可能会发生这种情况。

建议以使用向后兼容的演进schema。本质上，这意味着要么将每个新添加的字段设置为空值，要么为每个新字段设置为默认值。从Hudi版本0.5.1起，如果依赖字段的默认值，则该故障处理对此无效。

### 11.9.1.4 Hudi 在 upsert 时占用了临时文件夹中大量空间

#### 问题

Hudi在upsert时占用了临时文件夹中大量空间。

#### 回答

当UPSERT大量输入数据时，如果数据量达到合并的最大内存时，Hudi将溢出部分输入数据到磁盘。

如果有足够的内存，请增加spark executor的内存和添加“hoodie.memory.merge.fraction”选项，如：  
option("hoodie.memory.merge.fraction", "0.8")

### 11.9.1.5 Hudi 写入小精度 Decimal 数据失败

#### 问题

Hudi表初始入库采用BULK\_INSERT方式入库含有Decimal类型的数据，之后执行upsert，数据写入时报错：

```
java.lang.UnsupportedOperationException: org.apache.parquet.avro.AvroConverters$FieldFixedConverter
```

#### 回答

##### 原因：

Hudi表数据含有Decimal类型数据。

初始入库BULK\_INSERT方式会使用Spark内部parquet文件的写入类进行写入，Spark对不同精度的Decimal类型处理是不同的。

UPSERT操作时，Hudi使用Avro兼容的parquet文件写入类进行写入，这个和Spark的写入方式是不兼容的。

##### 解决方案：

执行BULK\_INSERT时指定设置“hoodie.datasource.write.row.writer.enable = false”，使hoodie采用Avro兼容的parquet文件写入类进行写入。

## 11.9.2 数据采集

### 11.9.2.1 使用 kafka 采集数据时报错 IllegalArgumentException

#### 问题

线程“main”报错 org.apache.kafka.common.KafkaException，构造kafka消费者失败，报错：

```
java.lang.IllegalArgumentException: Could not find a 'KafkaClient' entry in the JAAS configuration. System property 'java.security.auth.login.config' is not set
```

#### 回答

当试图从启用SSL的kafka数据源采集数据时，而安装程序无法读取jars.conf文件及其属性时，可能会发生这种情况。

要解决此问题，需要将所需的属性作为通过Spark提交的命令的一部分传递。如：--files jaas.conf,failed\_tables.json --conf 'spark.driver.extraJavaOptions=-Djava.security.auth.login.config=jaas.conf' --conf 'spark.executor.extraJavaOptions=-Djava.security.auth.login.config=jaas.conf'

### 11.9.2.2 采集数据时报错 HoodieException

#### 问题

数据采集时报错：

```
com.uber.hoodie.exception.HoodieException: created_at(Part -created_at) field not found in record. Acceptable fields were :[col1, col2, col3, id, name, dob, created_at, updated_at]
```

#### 回答

这种情况通常当标记为recordKey或partitionKey的字段在某些传入记录中不存在时发生。

请交叉验证传入记录。

### 11.9.2.3 采集数据时报错 HoodieKeyException

#### 问题

创建Hudi表时，是否可以使用包含空记录的可空字段作为主键？

#### 回答

不可以。

使用包含空记录的可空字段作为主键时会返回HoodieKeyException异常：

```
Caused by: org.apache.hudi.exception.HoodieKeyException: recordKey value: "null" for field: "name" cannot be null or empty.
at org.apache.hudi.keygen.SimpleKeyGenerator.getKey(SimpleKeyGenerator.java:58)
at org.apache.hudi.HoodieSparkSqlWriter$$anonfun$1.apply(HoodieSparkSqlWriter.scala:104)
at org.apache.hudi.HoodieSparkSqlWriter$$anonfun$1.apply(HoodieSparkSqlWriter.scala:100)
```

## 11.9.3 Hive 同步

### 11.9.3.1 Hive 同步数据报错 SQLException

#### 问题

Hive同步数据时报错：

```
Caused by: java.sql.SQLException: Error while processing statement: FAILED: Execution Error, return code 1 from org.apache.hadoop.hive.ql.exec.DDLTask. Unable to alter table. The following columns have types incompatible with the existing columns in their respective positions :
__col1,__col2
```

#### 回答

这种情况通常会发生当您试图使用HiveSyncTool.java类向现有hive表添加新列时。数据库通常不允许将列数据类型按照从高到低的顺序修改，或者数据类型可能与表中已存储/将要存储的数据冲突。若要修复相同的问题，请尝试设置以下属性：

设置hive.metastore.disallow.in compatible.col.type.changes为false。

### 11.9.3.2 Hive 同步数据报错 HoodieHiveSyncException

#### 问题

Hive同步数据时报错：

```
com.uber.hoodie.hive.HoodieHiveSyncException: Could not convert field Type from <type1> to <type2> for field col1
```

#### 回答

出现这种情况是因为HiveSyncTool目前只支持很少的兼容数据类型转换。进行任何其他不兼容的更改都会引发此异常。

请检查相关字段的数据类型演进，并验证它是否确实可以被视为根据Hudi代码库的有效数据类型转换。

### 11.9.3.3 Hive 同步数据报错 SemanticException

#### 问题

Hive同步数据时报错：

```
org.apache.hadoop.hive.ql.parse.SemanticException: Database does not exist: test_db
```

#### 回答

这种情况通常在试图对Hudi数据集执行Hive同步，但配置的hive\_sync数据库不存在时发生。

请在您的Hive集群上创建对应的数据库后重试。

# 12 使用 Hue（MRS 3.x 之前版本）

## 12.1 访问 Hue WebUI 界面

### 操作场景

MRS集群安装Hue组件后，用户可以通过Hue的WebUI，在图形化界面使用Hadoop与Hive。

该任务指导用户在MRS集群中打开Hue的WebUI。

#### 说明

Internet Explorer浏览器可能存在兼容性问题，建议更换兼容的浏览器访问Hue WebUI，例如Google Chrome浏览器50版本。

### 对系统的影响

第一次访问Manager和Hue WebUI，需要在浏览器中添加站点信任以继续打开Hue WebUI。

### 前提条件

启用Kerberos认证时，MRS集群管理员已分配用户使用Hive的权限。例如创建一个“人机”用户“hueuser”，并加入“hive”、“hadoop”、“supergroup”组和“manager\_view”角色，主组为“hive”。

该用户用于登录Hue WebUI。

### 操作步骤



**步骤1** 登录服务页面：单击集群名称，登录集群详情页面，选择“组件管理”。

#### 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

**步骤2** 选择“Hue”，在“Hue WebUI”右侧，单击链接，打开Hue的WebUI，以创建的“hueuser”用户登录Hue WebUI。

Hue的WebUI支持以下功能：

- 使用“Query Editors”执行Hive的查询语句。需要MRS集群已安装Hive。
- 使用“Data Browsers”管理Hive中的表。需要MRS集群已安装Hive。
- 使用查看HDFS中的目录和文件。需要MRS集群已安装HDFS。
- 使用查看MRS集群中所有作业。需要MRS集群已安装YARN。

#### 说明

- 使用创建的用户第一次登录Hue WebUI，需修改密码。
- 用户获取Hue WebUI的访问地址后，可以给其他无法访问Manager的用户用于访问Hue WebUI。
- 在Hue的WebUI操作但不操作Manager页面，重新访问Manager时需要输入已登录的账号密码。

---结束

## 12.2 使用 Hue WebUI 操作 Hive 表

Hue提供了文件浏览器功能，使用户可以通过界面图形化的方式查看Hive上文件及目录功能。

### 前提条件

已安装Hive以及Hue组件，且状态为运行中的Kerberos认证的集群。

### 操作步骤

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 打开Hue WebUI，然后选择“Query Editors > Hive”。


**步骤3** 在“Databases”选择一个Hive中的数据库，默认数据库为“default”。

系统将自动显示数据库中的所有表。可以输入表名关键字，系统会自动搜索包含此关键字的全部表。


**步骤4** 单击指定的表名，可以显示表中所有的列。

**步骤5** 在HiveQL语句编辑区输入HiveQL语句。


```
create table hue_table(id int,name string,company string) row format delimited fields terminated by ',' stored as textfile;
```

单击 并选择“Explain”，编辑器将分析输入的HiveQL语句是否有语法错误以及执行计划，如果存在语法错误则显示“Error while compiling statement”。

**步骤6** 单击, 选择HiveQL语句执行的引擎。

**步骤7** 单击 开始执行HiveQL语句。



**步骤8** 在命令输入框内输入 `show tables;`，单击  按钮，查看结果中有 **步骤5** 创建的表 `hue_table`。

----结束

## 12.3 在 Hue WebUI 使用 HiveQL 编辑器

### 操作场景


用户需要使用图形化界面在集群中执行 HiveQL 语句时，可以通过 Hue 完成任务。

### 访问“Query Editors”

**步骤1** 访问 Hue WebUI，请参考 [访问 Hue WebUI 界面](#)。

**步骤2** 选择“Query Editors > Hive”，进入“Hive”。


“Hive”支持以下功能：

- 执行和管理 HiveQL 语句。
- 在“Saved Queries”中查看当前访问用户已保存的 HiveQL 语句。
- 在“Query History”中查看当前访问用户执行过的 HiveQL 语句。
- 单击 ，在“Databases”下可以显示 Hive 中所有的数据库。

----结束


### 执行 HiveQL 语句

**步骤1** 选择“Query Editors > Hive”，进入“Hive”。


**步骤2** 单击 ，在“Databases”下选择一个数据库，默认数据库为“default”。

系统将自动显示数据库中的所有表。可以输入表名关键字，系统会自动搜索包含此关键字的全部表。

**步骤3** 单击指定的表名，可以显示表中所有的列。

光标移动到表所在的行，单击  可以查看列的详细信息。

**步骤4** 在 HiveQL 语句编辑区输入查询语句。

单击  并选择“Explain”，编辑器将分析输入的查询语句是否有语法错误以及执行计划，如果存在语法错误则显示“Error while compiling statement”。

**步骤5** 单击 ，选择 HiveQL 语句执行的引擎。







- “mr”表示语句使用 MapReduce 计算框架执行语句。
- “spark”表示语句使用 Spark 计算框架执行语句。
- “tez”表示语句使用 Tez 计算框架执行语句。

#### 说明

tez 适用于 MRS 1.9.x 及以后版本。

**步骤6** 单击  开始执行HiveQL语句。

#### 说明

- 如果希望下次继续使用已输入的HiveQL语句，请单击  保存。
- 格式化HiveQL语句，请单击  选择“Format”。
- 删除已输入的HiveQL语句，请单击  选择“Clear”。
- 清空已输入的语句并执行一个新的语句，请单击  选择“New query”。
- 查看历史：  
单击“Query History”，可查看HiveQL运行情况，支持显示所有语句或只显示保存的语句的运行情况。历史记录存在多个结果时，可以在输入框使用关键字进行搜索。
- 高级查询配置：  
单击右上角的 ，对文件、函数、设置等信息进行配置。
- 查看快捷键：  
单击右上角的 ，可查看所有快捷键信息。

----结束

## 查看执行结果

**步骤1** 在“Hive”的执行区，默认显示“Query History”。

**步骤2** 单击“Results”查看已执行语句的执行结果。

----结束

## 管理查询语句


**步骤1** 选择“Query Editors > Hive”，进入“Hive”。

**步骤2** 单击“Saved Queries”。

单击一条已保存的语句，系统会自动将其填充至编辑区中。

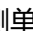
----结束


## 修改在 Hue 使用“Query Editors”的会话配置

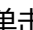
**步骤1** 在“Hive”页签，单击 。


**步骤2** 在“Files”的右侧单击 ，然后单击  指定该文件的存储目录。

可以单击  新增加一个文件资源。

**步骤3** 在“Functions”的右侧单击 ，输入用户自定义的名称和函数的类名称。

可以单击  新增加一个自定义函数。

**步骤4** 在“Settings”的右侧单击 ，在“Key”输入Hive的参数名，在“Value”输入对应的参数值，则当前Hive会话会以用户定义的配置连接Hive。

可以单击  新增加一个参数。

----结束

## 12.4 在 Hue WebUI 使用元数据浏览器

### 操作场景


用户需要使用图形化界面在集群中管理Hive的元数据，可以通过Hue完成任务。

### Metastore 管理器使用介绍




访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

选择“Data Browsers > Metastore Tables”，进入“Metastore Manager”。

- 查看Hive表的元数据

在左侧导航栏中，将鼠标放在某一表上，单击显示在其右侧的图标 ，界面将显示Hive表的元数据信息。



- 管理Hive表的元数据

在Hive表的元数据信息界面，单击右上角的  可导入数据，单击  可浏览数据，单击  可查看表文件的位置信息。

#### 注意

Hue界面主要用于文件、表等数据的查看与分析，禁止通过Hue界面对操作对象进行删除等高危管理操作。如需操作，建议在确认对业务没有影响后通过各组件的相应操作方法进行处理，例如使用HDFS客户端对HDFS文件进行操作，使用Hive客户端对Hive表进行操作。

- 管理Hive元数据表

选择右上角的  可在数据库中根据上传的文件创建一个新表，选择右上角的  可手动创建一个新表。

### 访问“Metastore Manager”

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 选择“Data Browsers > Metastore Tables”，进入“Metastore Manager”。

“Metastore Manager”支持以下功能：


- 使用文件创建一个Hive表
- 手动创建一个Hive表
- 查看Hive表元数据

----结束

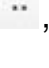
## 使用文件创建一个 Hive 表

**步骤1** 访问“Metastore Manager”，在“Databases”选择一个数据库。

默认数据库为“default”。

**步骤2** 单击 ，进入“Create a new table from a file”页面。

**步骤3** 选择文件。

1. 在“Table Name”填写Hive表的名称。  
支持字母、数字、下划线，首位必须为字母或数字，且长度不能超过128位。
2. 根据需要，在“Description”填写Hive表的描述信息。
3. 在“Input File or Location”单击 ，在HDFS中选择一个用于创建Hive表文件。  
此文件将存储Hive表的新数据。  
如果文件未在HDFS中保存，可以单击“Upload a file”从本地选择文件并上传。  
支持同时上传多个文件，文件不可为空。
4. 如果需要将文件中的数据导入Hive表，选择“Import data”作为“Load method”。默认选择“Import data”。  
选择“Create External Table”时，创建的是Hive外部表。

### 说明


当选择“Create External Table”时，参数“Input File or Location”需要选择为路径。  
选择“Leave Empty”则创建空的Hive表。

5. 单击“Next”。

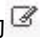
**步骤4** 设置分隔符。

1. 在“Delimiter”选择一个分隔符。  
如果分隔符不在列表中，选择“Other..”，然后输入新定义的分隔符。
2. 单击“Preview”查看数据处理预览。
3. 单击“Next”。

**步骤5** 定义字段列。

1. 单击“Use first row as column names”右侧的 ，则使用文件中第一行数据作为列名称。取消则不使用数据作为列名称。
2. 在“Column name”编辑每个列的名称。  
支持字母、数字、下划线，首位必须为字母或数字，且长度不能超过128位。

### 说明

单击“Bulk edit column names”右侧的 ，可批量对列重新命名。输入所有列的名称并使用逗号分隔。

3. 在“Column Type”选择每个列的类型。


**步骤6** 单击“Create Table”创建表，等待Hue显示Hive表的信息。

----结束

## 手工创建一个 Hive 表

**步骤1** 访问“Metastore Manager”，在“Databases”选择一个数据库。

默认数据库为“default”。

**步骤2** 单击，进入“Create a new table manually”页面。

**步骤3** 设置表名称。

1. 在“Table Name”填写Hive表的名称。  
支持字母、数字、下划线，首位必须为字母或数字，且长度不能超过128位。
2. 根据需要，在“Description”填写Hive表的描述信息。
3. 单击“Next”。

**步骤4** 选择一个存储数据的格式。

- 需要使用分隔符分隔数据时，选择“Delimited”，然后执行**步骤5**。
- 需要使用序列化格式保存数据时，选择“SerDe”，执行**步骤6**。

**步骤5** 配置分隔符。

1. 在“Field terminator”设置一个列分隔符。  
如果分隔符不在列表中，选择“Other..”，然后输入新定义的分隔符。
2. 在“Collection terminator”设置一个分隔符，用于分隔Hive中类型为“array”的列的数据集合。例如一个列为array类型，其中一个值需要保存“employee”和“manager”，用户指定分隔符为“:”，则最终的值为“employee:manager”。
3. 在“Map key terminator”设置一个分隔符，用于分隔Hive中类型为“map”的列的数据。例如某个列为map类型，其中一个值需要保存描述为“aaa”的“home”，和描述为“bbb”的“company”，用户指定分隔符为“|”，则最终的值为“home|aaa:company|bbb”。
4. 单击“Next”，执行**步骤7**。


**步骤6** 设置序列化属性。

1. 在“SerDe Name”输入序列化格式的类名称  
“org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe”。  
用户可扩展Hive支持更多自定义的序列化类。
2. 在“Serde properties”输入序列化的样式的值：“field.delim”=“,”  
“collection.delim”=“:” “mapkey.delim”=“|”。
3. 单击“Next”，执行**步骤7**。

**步骤7** 选择一个数据表的格式，并单击“Next”。

- “TextFile”表示使用文本类型文件存储数据。
- “SequenceFile”表示使用二进制类型文件存储数据。
- “InputFormat”表示使用自定义的输入输出格式来使用文件中的数据。  
用户可扩展Hive支持更多的自定义格式类。
  - a. 在“InputFormat Class”填写输入数据使用的类  
“org.apache.hadoop.hive.ql.io.RCFileInputFormat”。
  - b. 在“OutputFormat Class”填写输出数据使用的类  
“org.apache.hadoop.hive.ql.io.RCFileOutputFormat”。

**步骤8** 选择一个文件保存位置，并单击“Next”。

默认勾选“Use default location”。如果需要自定义存储位置，请取消选中状态并在“External location”单击  指定一个文件存储位置。

**步骤9** 设置Hive表的字段。

1. 在“Column name”设置列的名称。  
支持字母、数字、下划线，首位必须为字母或数字，且长度不能超过128位。
2. 在“Column type”选择一个数据类型。  
单击“Add a column”可增加新的列。
3. 单击“Add a partition”为Hive表增加分区，可提高查询效率。

**步骤10** 单击“Create Table”创建表，等待Hue显示Hive表的信息。

----结束

## 管理 Hive 表

**步骤1** 访问“Metastore Manager”，在“Databases”选择一个数据库，页面显示数据库中所有的表。

默认数据库为“default”。

**步骤2** 单击数据库中的表名称，打开表的详细信息。

支持导入数据、浏览数据或查看文件存储位置。查看数据库所有的表时，可以直接勾选表然后执行查看、浏览数据操作。

---

### 注意

Hue界面主要用于文件、表等数据的查看与分析，禁止通过Hue界面对操作对象进行删除等高危管理操作。如需操作，建议在确认对业务没有影响后通过各组件的相应操作方法进行处理，例如使用HDFS客户端对HDFS文件进行操作，使用Hive客户端对Hive表进行操作。

---

----结束

## 12.5 在 Hue WebUI 使用文件浏览器

### 操作场景

用户需要使用图形化界面管理HDFS文件时，可以通过Hue完成任务。

---

### 注意

Hue界面主要用于文件、表等数据的查看与分析，禁止通过Hue界面对操作对象进行删除等高危管理操作。如需操作，建议在确认对业务没有影响后通过各组件的相应操作方法进行处理，例如使用HDFS客户端对HDFS文件进行操作，使用Hive客户端对Hive表进行操作。

---

## 访问文件浏览器（File Browser）

**步骤1** 访问Hue WebUI。

**步骤2** 单击，进入“File Browser”。

默认进入当前登录用户的主目录。

文件浏览器将显示目录中的子目录或文件以下信息：

表 12-1 HDFS 文件属性介绍


| 属性名           | 描述            |
|---------------|---------------|
| “Name”        | 表示目录或文件的名称。   |
| “Size”        | 表示文件的大小。      |
| “User”        | 表示目录或文件的属主。   |
| “Group”       | 表示目录或文件的属组。   |
| “Permissions” | 表示目录或文件的权限设置。 |
| “Date”        | 表示目录或文件创建时间。  |

**步骤3** 在搜索框输入关键字，系统会在当前目录自动搜索目录或文件。

**步骤4** 清空搜索框的内容，系统会重新显示所有目录和文件。

----结束

## 执行动作

**步骤1** 单击，选择一个或多个目录或文件。

**步骤2** 单击“Actions”，在弹出菜单选择一个操作。

- “Rename”：表示重新命名一个目录或文件。
- “Move”：表示移动文件，在“移至”选择新的目录并单击“移动”完成移动。
- “Copy”：表示复制选中的文件或目录。
- “Change permissions”：表示修改选中目录或文件的访问权限。
  - 可以为属主、属组和其他用户设置“Read”、“Write”和“Excute”权限。
  - “Sticky”表示禁止HDFS的管理员、目录属主或文件属主以外的用户在目录中移动文件。
  - “Recursive”表示递归设置权限到子目录。
- “Storage policies”：表示设置目录或文件在HDFS中的存储策略。
- “Summary”：表示查看选中的文件或目录的HDFS存储信息。

----结束

## 访问其他目录

**步骤1** 单击目录名并输入需要访问的目录完整路径，例如“/mr-history/tmp”并按回车键进入目录。

需要当前登录Hue WebUI的用户拥有其他目录的访问权限。

**步骤2** 单击“Home”可进入用户的主目录。

**步骤3** 单击“History”可以显示最近访问目录的历史记录，并重新访问。

**步骤4** 单击“Trash”可以访问当前目录的回收站空间。

单击“Empty Trash”可清空回收站。

----结束

## 上传用户文件

**步骤1** 单击, 单击Upload。

**步骤2** 选择一个操作。

- “Files”：表示上传用户文件到当前用户。
- “Zip/Tgz/Bz2 file”：表示上传了一个压缩文件，在弹出框单击“Select ZIP, TGZ or BZ2 files”选择需要上传的压缩文件。系统会自动在HDFS中对文件解压。支持“ZIP”、“TGZ”和“BZ2”格式的压缩文件。

----结束

## 创建新文件或者目录

**步骤1** 单击, 单击“New”。

**步骤2** 选择一个操作。

- “File”：表示创建一个文件，输入文件名后单击“Create”完成。
- “Directory”：表示创建一个目录，输入目录名后单击“Create”完成

----结束

## 存储策略定义使用介绍

### 说明

若Hue的服务配置参数“fs\_defaultFS”配置为“viewfs://ClusterX”时，不能启用存储策略定义功能。

**步骤1** 登录MRS Manager。

**步骤2** 在MRS Manager界面，选择“系统设置 > 权限配置 > 角色管理 > 添加角色”：

1. 设置“角色名称”。
2. 选择“权限 > Hue”，勾选“Storage Policy Admin”，单击“确定”，为该角色赋予存储策略管理员的权限。


**步骤3** 选择“系统设置 > 权限配置 > 用户组管理 > 添加用户组”，设置“组名”，单击“角色”后的“选择添加角色”，在弹出的界面选择刚创建的角色，单击“确定”将该角色添加到组中。



**步骤4** 选择“系统设置 > 权限配置 > 用户管理 > 添加用户”：

1. 设置可以登录Hue的WebUI界面且有存储策略管理员权限的用户的“用户名”。
2. “用户类型”选择“人机”。
3. 设置登录Hue的WebUI界面的“密码”、“确认密码”。
4. 单击“用户组”后的“选择添加的用户组”，在弹出的界面选择创建的用户组、supergroup、hadoop和hive用户组，单击“确定”。
5. “主组”选择“hive”。
6. 单击“分配角色权限”右侧的“选择并绑定角色”，在弹出的界面选择刚刚创建的角色和System\_administrator角色，单击“确定”。
7. 再单击“确定”成功添加该用户。

**步骤5** 访问Hue WebUI。

**步骤6** 单击右上角的。

**步骤7** 勾选目录的复选框，单击页面上方的“Action”，选择“Storage policies”。

**步骤8** 在弹出的对话框中设置新的存储策略，单击“OK”。

----结束

## 12.6 在 Hue WebUI 使用作业浏览器

### 操作场景

用户需要使用图形化界面查看集群中所有作业时，可以通过Hue完成任务。

### 访问“Job Browser”

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 单击“Job Browser”。


默认显示当前集群的所有作业。

#### 说明

“Job Browser”显示的数字表示集群中所有作业的总数。

“Job Browser”将显示作业以下信息：

表 12-2 MRS 作业属性介绍

| 属性名                | 描述                                                                                                              |
|--------------------|-----------------------------------------------------------------------------------------------------------------|
| “Logs”             | 表示作业的日志信息。如果作业有输出日志，则显示  。 |
| “ID”               | 表示作业的编号，由系统自动生成。                                                                                                |
| “Name”             | 表示作业的名称。                                                                                                        |
| “Application Type” | 表示作业的类型。                                                                                                        |

| 属性名         | 描述                                                 |
|-------------|----------------------------------------------------|
| “Status”    | 表示作业的状态，包含“RUNNING”、“SUCCEEDED”、“FAILED”和“KILLED”。 |
| “User”      | 表示启动该作业的用户。                                        |
| “Maps”      | 表示作业执行Map过程的进度。                                    |
| “Reduces”   | 表示作业执行Reduce过程的进度。                                 |
| “Queue”     | 表示作业运行时使用的YARN队列。                                  |
| “Priority”  | 表示作业运行时的优先级。                                       |
| “Duration”  | 表示作业运行使用的时间。                                       |
| “Submitted” | 表示作业提交到MRS集群的时间。                                   |

#### 说明

如果MRS集群安装了Spark组件，则默认会启动一个作业“Spark-JDBCServer”，用于执行任务。

----结束

## 搜索作业

**步骤1** 在“Job Browser”的“Username”或“Text”，输入指定的字符，系统会自动搜索包含此关键字的全部作业。

**步骤2** 清空搜索框的内容，系统会重新显示所有作业。


----结束

## 查看作业详细信息

**步骤1** 在“Job Browser”的作业列表，单击作业所在的行，可以打开作业详情。

**步骤2** 在“Metadata”页签，可查看作业的元数据。

#### 说明

单击可打开作业运行时的日志。

----结束

# 12.7 Hue 常用配置参数

## 参数入口

参数入口，请参考[修改集群服务配置参数](#)。

## 参数说明

表 12-3 Hue 常用参数

| 配置参数                           | 说明             | 缺省值   | 范围                                                                                                         |
|--------------------------------|----------------|-------|------------------------------------------------------------------------------------------------------------|
| HANDLER_ACCESSLOG_LEVEL        | 表示Hue的访问日志级别。  | DEBUG | <ul style="list-style-type: none"> <li>• ERROR</li> <li>• WARN</li> <li>• INFO</li> <li>• DEBUG</li> </ul> |
| HANDLER_AUDITLOG_LEVEL         | 表示Hue的审计日志级别。  | DEBUG | <ul style="list-style-type: none"> <li>• ERROR</li> <li>• WARN</li> <li>• INFO</li> <li>• DEBUG</li> </ul> |
| HANDLER_ERRORLOG_LEVEL         | 表示Hue的错误日志级别。  | ERROR | <ul style="list-style-type: none"> <li>• ERROR</li> <li>• WARN</li> <li>• INFO</li> <li>• DEBUG</li> </ul> |
| HANDLER_LOGFILE_LEVEL          | 表示Hue的运行日志级别。  | INFO  | <ul style="list-style-type: none"> <li>• ERROR</li> <li>• WARN</li> <li>• INFO</li> <li>• DEBUG</li> </ul> |
| HANDLER_LOGFILE_MAXBACKUPINDEX | 表示Hue日志文件最大个数。 | 20    | 1 ~ 999                                                                                                    |
| HANDLER_LOGFILE_SIZE           | 表示Hue日志文件最大大小。 | 5MB   | -                                                                                                          |

# 13 使用 Hue（MRS 3.x 及之后版本）

## 13.1 访问 Hue WebUI 界面

### 操作场景

MRS集群安装Hue组件后，用户可以通过Hue的WebUI，在图形化界面使用Hadoop生态相关组件。

该任务指导用户在MRS集群中打开Hue的WebUI。

#### 说明

Internet Explorer浏览器可能存在兼容性问题，建议更换兼容的浏览器访问Hue WebUI，例如Google Chrome浏览器。

### 对系统的影响

第一次访问Manager和Hue WebUI，需要在浏览器中添加站点信任以继续打开Hue WebUI。

### 前提条件

启用Kerberos认证时，MRS集群管理员已分配用户使用Hive的权限，具体操作请参见[创建MRS集群用户](#)。例如创建一个“人机”用户“hueuser”，并加入“hive”、“hadoop”、“supergroup”组和“manager\_view”角色，主组为“hive”。

该用户用于登录Manager。

### 操作步骤









#### 步骤1 登录服务页面：

MRS 3.x之前版本，在MRS控制台单击集群名称，选择“组件管理 > Hue”。

MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)，选择“集群 > 服务 > Hue”。

#### 步骤2 在“Hue WebUI”右侧，单击链接，打开Hue的WebUI。

Hue的WebUI支持以下功能：

- 使用编辑器执行Hive、SparkSql的查询语句以及Notebook代码段。需要MRS集群已安装Hive、Spark2x。
- 使用计划程序提交Workflow任务、计划任务、Bundle任务。
- 使用文档查看、导入、导出在Hue页面上操作的任务，例如保存的Workflow任务、定时任务、Bundle任务等。
- 使用表管理Hive、SparkSql中的元数据。需要MRS集群已安装Hive、Spark2x。
- 使用文件查看HDFS中的目录和文件。需要MRS集群已安装HDFS。
- 使用作业查看MRS集群中所有作业。需要MRS集群已安装Yarn。
- 使用HBase创建/查询HBase表。需要MRS集群已安装HBase组件并添加Thrift1Server实例。
- 使用导入器通过“.csv”，“.txt”等格式的文件导入数据。

#### 说明

- 使用创建的用户第一次登录Hue WebUI，需修改密码。
- 用户获取Hue WebUI的访问地址后，可以给其他无法访问Manager的用户用于访问Hue WebUI。
- 在Hue的WebUI操作但不操作Manager页面，重新访问Manager时需要输入已登录的账号密码。

----结束

## 13.2 使用 Hue WebUI 操作 Hive 表


Hue汇聚了与大多数Apache Hadoop组件交互的接口，致力让用户通过界面图形化的方式轻松使用Hadoop组件。目前Hue支持HDFS、Hive、HBase、Yarn、MapReduce、Oozie和SparkSQL等组件的可视化操作。

### 前提条件

已安装Hue组件。

### 操作步骤

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 在左侧导航栏单击编辑器，然后选择“Hive”。

**步骤3** 在“Database”右侧下拉列表选择一个Hive中的数据库，默认数据库为“default”。

系统将自动显示数据库中的所有表。可以输入表名关键字，系统会自动搜索包含此关键字的全部表。

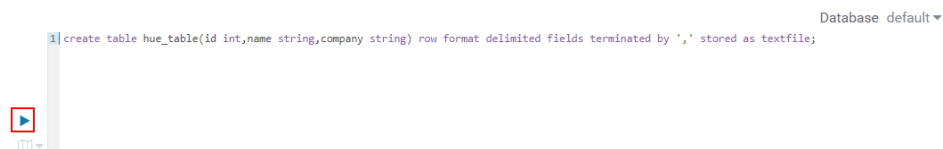
**步骤4** 单击指定的表名，可以显示表中所有的列。


**步骤5** 在HiveQL语句编辑区输入HiveQL语句。

```
create table hue_table(id int,name string,company string) row format delimited fields terminated by ',' stored as textfile;
```

**步骤6** 单击  开始执行HiveQL语句。

图 13-1 执行语句



**步骤7** 在命令输入框内输入show tables;，单击  按钮，查看“结果”中有**步骤5**创建的表 hue\_table。

----结束

## 13.3 创建 Hue 操作任务


### 13.3.1 在 Hue WebUI 使用 HiveQL 编辑器

#### 操作场景

用户需要使用图形化界面在集群中执行HiveQL语句时，可以通过Hue完成任务。

#### 访问编辑器

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 在左侧导航栏单击 ，然后选择“Hive”，进入“Hive”。

“Hive”支持以下功能：

- 执行和管理HiveQL语句。
- 在“保存的查询”中查看当前访问用户已保存的HiveQL语句。
- 在“查询历史记录”中查看当前访问用户执行过的HiveQL语句。


----结束

#### 执行 HiveQL 语句

**步骤1** 在“Database”右侧下拉列表选择一个Hive中的数据库，默认数据库为“default”。

系统将自动显示数据库中的所有表。可以输入表名关键字，系统会自动搜索包含此关键字的全部表。

**步骤2** 单击指定的表名，可以显示表中所有的列。

光标移动到表或列所在的行，单击  可以查看详细信息。





**步骤3** 在HiveQL语句编辑区输入查询语句。

**步骤4** 单击  开始执行HiveQL语句。

图 13-2 执行语句



### 说明

- 如果希望下次继续使用已输入的HiveQL语句，请单击  保存。
- 高级查询配置：  
单击右上角的 ，对文件、功能、设置等信息进行配置。
- 查看快捷键：  
单击右上角的 ，可查看语法和键盘快捷方式信息。
- 删除已输入的HiveQL语句，请单击  后的三角选择“清除”。
- 查看历史：  
单击“查询历史记录”，可查看HiveQL运行情况，支持显示所有语句或只显示保存的语句的运行情况。历史记录存在多个结果时，可以在输入框使用关键字进行搜索。

----结束

## 查看执行结果

**步骤1** 在“Hive”的执行区，默认显示“查询历史记录”。

**步骤2** 单击结果查看已执行语句的执行结果。

### 说明

Hue暂不支持大数据量展示，当SQL查询结果加载过量时可能出现页面卡顿，部分数据不显示等情况。目前建议查询结果加载不超过5000行。

----结束


## 管理查询语句



**步骤1** 单击“保存的查询”。

**步骤2** 单击一条已保存的语句，系统会自动将其填充至编辑区中。


----结束

## 修改在 Hue 使用编辑器的会话配置

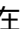
**步骤1** 在编辑器页面，单击 。


**步骤2** 在“文件”的右侧单击 ，然后单击  选择文件。

可以单击“文件”后的  新增加一个文件资源。

**步骤3** 在“功能” ，输入用户自定义的名称和函数的类名称。

可以单击“功能”后的  新增加一个自定义函数。

**步骤4** 在“设置” ，在“设置”的“键”输入Hive的参数名，在“值”输入对应的参数值，则当前Hive会话会以用户定义的配置连接Hive。

可以单击  新增加一个参数。

----结束

## 13.3.2 在 Hue WebUI 使用 SparkSql 编辑器

### 操作场景

用户需要使用图形化界面在集群中执行SparkSql语句时，可以通过Hue完成任务。

### 配置 Spark2x

使用SparkSql编辑器之前需要先修改Spark2x配置。

**步骤1** 进入Spark2x的全部配置页面，具体操作请参考[修改集群服务配置参数](#)。

**步骤2** 设置Spark2x多实例模式，搜索并修改Spark2x服务的以下参数：

| 参数名称                             | 值                             |
|----------------------------------|-------------------------------|
| spark.thriftserver.proxy.enabled | false                         |
| spark.scheduler.allocation.file  | #{conf_dir}/fairscheduler.xml |

**步骤3** 进入JDBCServer2x自定义界面，在“spark.core-site.customized.configs”参数内，添加如下两个自定义项：

表 13-1 自定义参数

| 名称                          | 值 |
|-----------------------------|---|
| hadoop.proxyuser.hue.groups | * |
| hadoop.proxyuser.hue.hosts  | * |

**步骤4** 保存配置，重启Spark2x服务。

----结束



## 访问编辑器

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 在左侧导航栏单击 ，然后选择“SparkSql”，进入“SparkSql”。

“SparkSql”支持以下功能：

- 执行和管理SparkSql语句。
- 在“保存的查询”中查看当前访问用户已保存的SparkSql语句。
- 在“查询历史记录”中查看当前访问用户执行过的SparkSql语句。

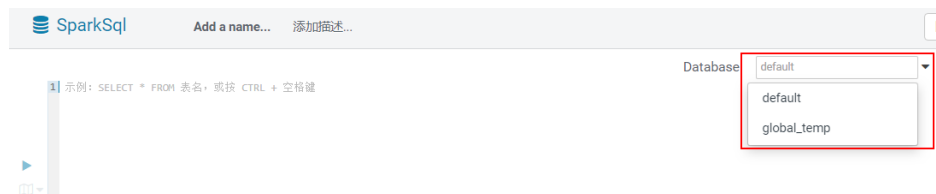
----结束

## 执行 SparkSql 语句


**步骤1** 在“Database”右侧下拉列表选择一个SparkSql中的数据库，默认数据库为“default”。

系统将自动显示数据库中的所有表。可以输入表名关键字，系统会自动搜索包含此关键字的全部表。


图 13-3 选择数据库



**步骤2** 单击指定的表名，可以显示表中所有的列。

光标移动到表所在的行，单击  可以查看列的详细信息。

**步骤3** 在SparkSql语句编辑区输入查询语句。

单击  后的三角并选择“解释”，编辑器将分析输入的查询语句是否有语法错误以及执行计划，如果存在语法错误则显示“Error while compiling statement”。







**步骤4** 单击  开始执行SparkSql语句。

图 13-4 执行语句



## 📖 说明

- 如果希望下次继续使用已输入的SparkSql语句，请单击  保存。
- 高级查询配置：  
单击右上角的 ，对文件、功能、设置等信息进行配置。
- 查看快捷键：  
单击右上角的 ，可查看语法和键盘快捷方式信息。
- 格式化SparkSql语句，请单击  后的三角选择“格式”
- 删除已输入的SparkSql语句，请单击  后的三角选择“清除”
- 查看历史：  
单击“查询历史记录”，可查看SparkSql运行情况，支持显示所有语句或只显示保存的语句的运行情况。历史记录存在多个结果时，可以在输入框使用关键字进行搜索。

----结束

## 查看执行结果

**步骤1** 在“SparkSql”的执行区，默认显示“查询历史记录”。

**步骤2** 单击结果查看已执行语句的执行结果。

----结束

## 管理查询语句

**步骤1** 单击“保存的查询”。

**步骤2** 单击一条已保存的语句，系统会自动将其填充至编辑区中。

----结束

## 13.3.3 在 Hue WebUI 使用元数据浏览器


### 操作场景

用户需要使用图形化界面在集群中管理Hive的元数据，可以通过Hue完成任务。

### 元数据管理器使用介绍

访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

- 查看Hive表的元数据


在左侧导航栏单击表 ，单击某一表名称，界面将显示Hive表的元数据信息。

- 管理Hive表的元数据

在Hive表的元数据信息界面：

- 单击右上角的“导入”可导入数据。
- 单击“概述”，在“属性”域可查看表文件的位置信息。

可查看Hive表各列字段的信息，并手动添加描述信息，注意此处添加的描述信息并不是Hive表中的字段注释信息（comment）。

- 单击“样本”可浏览数据。
- 管理Hive元数据表  
单击左侧列表中的  可在数据库中根据上传的文件创建一个新表，也可手动创建一个新表。

 **注意**

Hue界面主要用于文件、表等数据的查看与分析，禁止通过Hue界面对操作对象进行删除等高危管理操作。如需操作，建议在确认对业务没有影响后通过各组件的相应操作方法进行处理，例如使用HDFS客户端对HDFS文件进行操作，使用Hive客户端对Hive表进行操作。

## 13.3.4 在 Hue WebUI 使用文件浏览器

### 操作场景

用户需要使用图形化界面管理HDFS文件时，可以通过Hue完成任务。

 **注意**

Hue界面主要用于文件、表等数据的查看与分析，禁止通过Hue界面对操作对象进行删除等高危管理操作。如需操作，建议在确认对业务没有影响后通过各组件的相应操作方法进行处理，例如使用HDFS客户端对HDFS文件进行操作，使用Hive客户端对Hive表进行操作。

### 访问文件浏览器

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 在左侧导航栏单击文件 。进入“文件浏览器”页面。

“文件浏览器”的“主页”默认进入当前登录用户的主目录。界面将显示目录中的子目录或文件的以下信息：

表 13-2 HDFS 文件属性介绍

| 属性名 | 描述          |
|-----|-------------|
| 名称  | 表示目录或文件的名称。 |
| 大小  | 表示文件的大小。    |
| 用户  | 表示目录或文件的属主。 |
| 组   | 表示目录或文件的属组。 |

| 属性名 | 描述            |
|-----|---------------|
| 权限  | 表示目录或文件的权限设置。 |
| 日期  | 表示目录或文件创建时间。  |

**步骤3** 在搜索框输入关键字，系统会在当前目录自动搜索目录或文件。

**步骤4** 清空搜索框的内容，系统会重新显示所有目录和文件。

----结束

## 执行动作

**步骤1** 在“文件浏览器”界面，勾选一个或多个目录或文件。

**步骤2** 单击“操作”，在弹出菜单选择一个操作。

- 重命名：表示重新命名一个目录或文件。
- 移动：表示移动文件，在“移至”页面选择新的目录并单击“移动”完成移动。
- 复制：表示复制选中的文件或目录。
- 更改权限：表示修改选中目录或文件的访问权限。
  - 可以为属主、属组和其他用户设置“读取”、“写”和“执行”权限。
  - “易贴”表示禁止HDFS的管理员、目录属主或文件属主以外的用户在目录中移动文件。
  - “递归”表示递归设置权限到子目录。
- 存储策略：表示设置目录或文件在HDFS中的存储策略。
- 摘要：表示查看选中的文件或目录的HDFS存储信息。

----结束

## 上传用户文件

**步骤1** 在“文件浏览器”界面，单击“上传”。

**步骤2** 在弹出的上传文件窗口中单击“选择文件”或将文件拖至窗口中，完成文件上传。

----结束

## 创建新文件或者目录

**步骤1** 在“文件浏览器”界面，单击“新建”。

**步骤2** 选择一个操作。

- 文件：表示创建一个文件，输入文件名后单击“创建”完成。
- 目录：表示创建一个目录，输入目录名后单击“创建”完成。

----结束

## 存储策略定义使用介绍

### 说明

若Hue的服务配置参数“fs\_defaultFS”配置为“viewfs://ClusterX”时，不能启用存储策略定义功能。

**步骤1** 登录FusionInsight Manager。

**步骤2** 在FusionInsight Manager界面，选择“系统 > 权限 > 角色 > 添加角色”：

1. 设置“角色名称”。
2. 在“配置资源权限”下选择“待操作集群名称>Hue”，勾选“存储策略管理员”，单击“确定”，为该角色赋予存储策略管理员的权限。

**步骤3** 选择“系统 > 权限 > 用户组 > 添加用户组”，设置“组名”，单击“角色”后的“添加”，在弹出的界面选择**步骤2**创建的角色，单击“确定”将该角色添加到组中，单击“确定”完成用户组的创建。

**步骤4** 选择“系统 > 权限 > 用户 > 添加用户”：

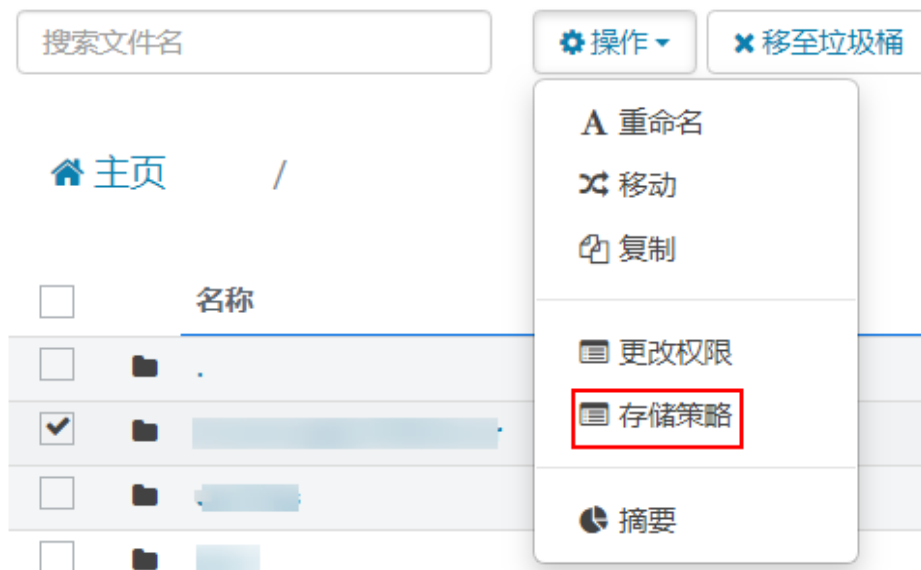
1. “用户名”填写待添加的用户名。
2. “用户类型”设置为“人机”。
3. 设置登录Hue的WebUI界面的“密码”、“确认密码”。
4. 单击“用户组”后的“添加”，在弹出的界面选择**步骤3**创建的用户组、supergroup、hadoop和hive用户组，单击“确定”。
5. “主组”选择“hive”。
6. 单击“角色”后的“添加”，在弹出的界面选择**步骤2**创建的角色和System\_administrator角色，单击“确定”。
7. 再单击“确定”，成功添加该用户。

**步骤5** 使用创建的用户访问Hue WebUI，具体操作请参考[访问Hue WebUI界面](#)。

**步骤6** 左侧导航栏单击文件。进入“文件浏览器”页面。

**步骤7** 勾选目录的复选框，单击页面上方的“操作”，单击“存储策略”。

图 13-5 存储策略



**步骤8** 在弹出的对话框中设置新的存储策略，单击“保存”。

----结束


## 13.3.5 在 Hue WebUI 使用作业浏览器

### 操作场景

用户需要使用图形化界面查看集群中所有作业时，可以通过Hue完成任务。

### 访问作业浏览器

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 单击作业。

默认显示当前集群的所有作业。

#### 说明

作业浏览器显示的数字表示集群中所有作业的总数。

“作业浏览器”将显示作业以下信息：

**表 13-3** MRS 作业属性介绍

| 属性名  | 描述                          |
|------|-----------------------------|
| 名称   | 表示作业的名称。                    |
| 用户   | 表示启动该作业的用户。                 |
| 类型   | 表示作业的类型。                    |
| 状态   | 表示作业的状态，包含“成功”、“正在运行”、“失败”。 |
| 进度   | 表示作业运行进度。                   |
| 组    | 表示作业所属组。                    |
| 开始   | 表示作业开始时间。                   |
| 持续时间 | 表示作业运行使用的时间。                |
| Id   | 表示作业的编号，由系统自动生成。            |

#### 说明

如果MRS集群安装了Spark组件，则默认会启动一个作业“Spark-JDBCServer”，用于执行任务。

----结束

## 搜索作业

**步骤1** 在“作业浏览器”的搜索栏，输入指定的字符，系统会按照ID、名称、用户自动搜索包含此关键字的全部作业。

**步骤2** 清空搜索框的内容，系统会重新显示所有作业。

----结束

## 查看作业详细信息

**步骤1** 在“作业浏览器”的作业列表，单击作业所在的行，可以打开作业详情。

**步骤2** 在“元数据”页签，可查看作业的元数据。

### 说明

单击“日志”可打开作业运行时的日志。

----结束

## 13.3.6 在 Hue WebUI 使用 HBase

### 操作场景

用户需要使用图形化界面在集群中创建或查询HBase表时，可以通过Hue完成任务。

### 说明

如需在Hue WebUI中操作HBase，当前MRS集群中必须部署HBase的Thrift1Server实例。

Thrift1Server实例默认不会安装，用户可在创建自定义类型的MRS集群时，选择HBase组件并通过调整集群自定义拓扑，添加Thrift1Server实例，详情请参考[购买自定义拓扑集群](#)。

如果当前集群支持手动添加服务，也可以在首次添加HBase服务时，选择部署Thrift1Server实例，服务添加成功后，需重启Hue服务，详情请参考[添加服务](#)。

### 访问作业浏览器

**步骤1** 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

**步骤2** 单击HBase ，进入“HBase Browser”页面。

----结束

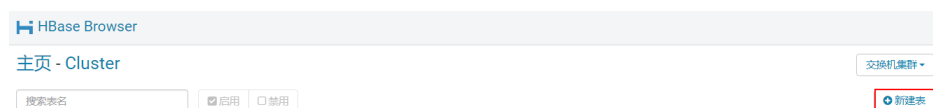
### 新建 HBase 表

**步骤1** 访问Hue WebUI。

**步骤2** 单击HBase ，进入“HBase Browser”页面。

**步骤3** 单击右侧“新建表”按钮，输入表名和列族参数，单击“提交”，完成HBase表创建。

图 13-6 新建表



----结束

## 查询 HBase 表数据

步骤1 访问Hue WebUI。

步骤2 单击HBase , 进入“HBase Browser”页面。

步骤3 单击需要查询的HBase表。可在上方的搜索栏后单击键值，对HBase表进行查询。

图 13-7 根据键值搜索



----结束

## 13.4 使用 Hue WebUI 典型场景

### 13.4.1 HDFS on Hue


Hue提供了文件浏览器功能，使用户可以通过界面图形化的方式使用HDFS。

#### 注意

Hue界面主要用于文件、表等数据的查看与分析，禁止通过Hue界面对操作对象进行删除等高危管理操作。如需操作，建议在确认对业务没有影响后通过各组件的相应操作方法进行处理，例如使用HDFS客户端对HDFS文件进行操作，使用Hive客户端对Hive表进行操作。

### 文件浏览器使用介绍

访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

然后单击 , 进入“文件浏览器”页面。您可以进行以下操作。

- 查看文件和目录

默认显示登录用户的目录及目录中的文件，可查看目录或文件的“名称”、“大小”、“用户”、“组”、“权限”和“日期”信息。

单击文件名，可查看文本文件的文本信息或二进制数据。支持编辑文件内容。

如果文件和目录数量比较多，可以在搜索框输入关键字，搜索特定的文件或目录。



- 创建文件或目录  
单击右上角的“新建”，选择“文件”创建文件，选择“目录”创建目录。
- 管理文件或目录  
勾选文件或目录的复选框，单击“操作”，选择“重命名”、“移动”、“复制”和“更改权限”等，实现文件或目录的重命名、移动、复制、更改权限等功能。
- 上传文件  
单击右上角的“上传”，单击“选择文件”或将文件拖至窗口中可进行文件上传。

## 存储策略定义使用介绍

### 说明

若Hue的服务配置参数“fs\_defaultFS”配置为“viewfs://ClusterX”时，不能启用存储策略定义功能。

存储策略定义在Hue的WebUI界面上分为两大类：

- 静态存储策略  
当前存储策略  
根据HDFS的文档访问频率、重要性，为HDFS目录指定存储策略，例如ONE\_SSD、ALL\_SSD等，此目录下的文件可被迁移到相应存储介质上保存。
- 动态存储策略  
为HDFS目录设置规则，系统可以根据文件的最近访问时间、最近修改时间自动修改存储策略、修改文件副本数、移动文件目录，详细的介绍请参见[配置HDFS冷热数据迁移](#)。

在Hue的WebUI界面设置动态存储策略之前，需先在Manager界面设置冷热数据迁移的CRON表达式，并启动自动冷热数据迁移特性。

操作方法为：

修改HDFS服务的NameNode的如下参数值。参数修改方法请参考[修改集群服务配置参数](#)。

| 参数                                  | 描述                                                                                                          | 取值示例      |
|-------------------------------------|-------------------------------------------------------------------------------------------------------------|-----------|
| dfs.auto.data.mover.enable          | 表示是否启用自动冷热数据迁移特性。默认值是“false”。                                                                               | true      |
| dfs.auto.data.mover.cron.expression | HDFS执行冷热数据迁移的CRON表达式，用于控制数据迁移操作的开始时间。仅当“dfs.auto.data.mover.enable”设置为“true”时才有效。默认值“0 * * * *”表示在每个整点执行任务。 | 0 * * * * |

修改参数“dfs.auto.data.mover.cron.expression”时，表达式介绍如[表13-4](#)所示。支持“\*”表示连续的时间段。

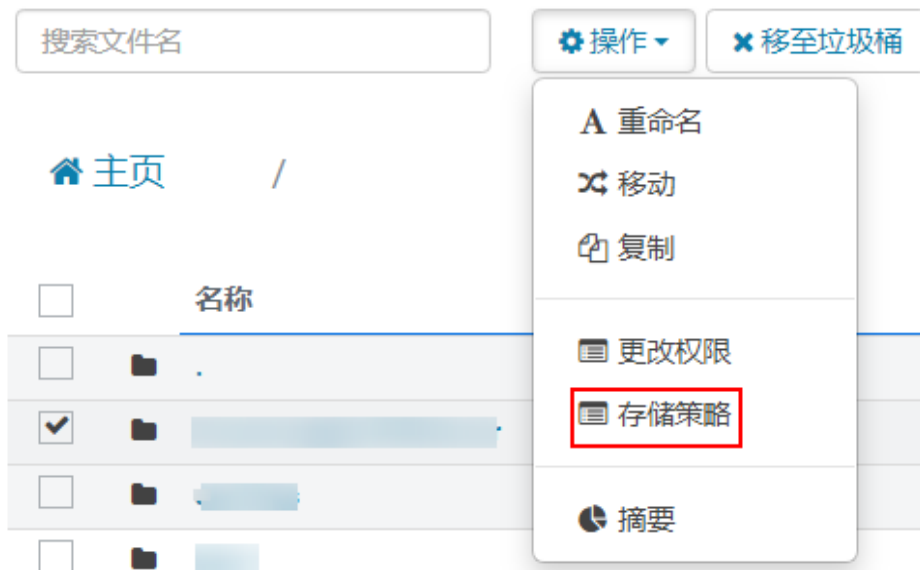
表 13-4 执行表达式参数解释

| 列   | 说明                 |
|-----|--------------------|
| 第1列 | 分钟，参数值为0~59。       |
| 第2列 | 小时，参数值为0~23。       |
| 第3列 | 日期，参数值为1~31。       |
| 第4列 | 月份，参数值为1~12。       |
| 第5列 | 星期，参数值为0~6，0表示星期日。 |

存储策略定义在WebUI界面上的操作如下：

- 步骤1** 登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
- 步骤2** 在FusionInsight Manager界面，选择“系统 > 权限 > 角色 > 添加角色”：
1. 设置“角色名称”。
  2. 在“配置资源权限”下选择“待操作集群名称>Hue”，勾选“存储策略管理员”，单击“确定”，为该角色赋予存储策略管理员的权限。
- 步骤3** 选择“系统 > 权限 > 用户组 > 添加用户组”，设置“组名”，单击“角色”后的“添加”，在弹出的界面选择**步骤2**创建的角色，单击“确定”将该角色添加到组中，单击“确定”完成用户组的创建。
- 步骤4** 选择“系统 > 权限 > 用户 > 添加用户”：
1. “用户名”填写待添加的用户名。
  2. “用户类型”设置为“人机”。
  3. 设置登录Hue的WebUI界面的“密码”、“确认密码”。
  4. 单击“用户组”后的“添加”，在弹出的界面选择**步骤3**创建的用户组、supergroup、hadoop和hive用户组，单击“确定”。
  5. “主组”选择“hive”。
  6. 单击“角色”后的“添加”，在弹出的界面选择**步骤2**创建的角色和System\_administrator角色，单击“确定”。
  7. 再单击“确定”，成功添加该用户。
- 步骤5** 使用创建的用户访问Hue WebUI。
- 步骤6** 左侧导航栏单击文件。进入“文件浏览器”页面。
- 步骤7** 勾选目录的复选框，单击页面上方的“操作”，单击“存储策略”。

图 13-8 存储策略



**步骤8** 在弹出的对话框中设置新的存储策略，单击“确定”。

- 在“静态存储策略”页签设置静态存储策略，单击“保存”。
- 在“动态存储策略”页签可创建、删除、修改动态存储策略，详细的参数介绍如表13-5所示。

表 13-5 动态存储策略参数介绍

| 分类 | 参数       | 说明                                      |
|----|----------|-----------------------------------------|
| 规则 | 文件最近访问时间 | 按照该文件最近一次访问时间。                          |
|    | 文件最近修改时间 | 按照该文件最近一次修改时间。                          |
| 操作 | 修改副本数    | 设置文件副本数。                                |
|    | 修改存储策略   | 修改存储策略，包括HOT、WARM、COLD、ONE_SSD、ALL_SSD。 |
|    | 移动到目录    | 移动该文件到其他目录。                             |

### 📖 说明

- 设置规则需要用户充分考虑合理性，例如多条规则之间是否有冲突，是否会对系统造成破坏等。
- 一个目录设置多个规则和动作时，规则被先触发的放在规则/动作列表的下面，规则被后触发的放在规则/动作列表的上面，避免动作反复执行。
- 系统每个小时整点扫描动态存储策略指定的目录下的文件是否符合规则，如果满足，则触发执行动作。执行日志记录在主NameNode的“/var/log/Bigdata/hdfs/nn/hadoop.log”目录下。

----结束

## 典型场景

通过Hue界面对HDFS以文本或二进制查看和编辑文件的操作如下：

### 查看文件

步骤1 访问Hue WebUI。

步骤2 左侧导航栏单击文件 。进入“文件浏览器”页面。

步骤3 单击需要查看的文件名。

步骤4 单击“以二进制格式查看”，可以切换视图从文本到二进制；单击“以文本格式查看”，可以切换视图从二进制到文本。

### 编辑文件

步骤5 单击“编辑文件”，显示文件内容可编辑。

步骤6 单击“保存”或“另存为”保存文件。

----结束

## 13.4.2 配置 HDFS 冷热数据迁移

### 配置场景

冷热数据迁移工具根据配置的策略移动HDFS文件。配置策略是条件或非条件规则的集合。如果规则匹配文件集，则该工具将对该文件执行一组行为操作。

冷热数据迁移工具支持以下规则和行。

- 迁移规则：
  - 根据文件的最后访问时间迁移数据
  - 根据年龄时间迁移数据（修改时间）
  - 无条件迁移数据

表 13-6 规则条件标签

| 条件标签                  | 描述            |
|-----------------------|---------------|
| <age operator="lt">   | 定义年龄/修改时间的条件。 |
| <atime operator="gt"> | 定义访问时间的条件。    |

### 说明

对于手动迁移规则，不需要条件。

- 行为列表：
  - 将存储策略设置为给定的数据层名称
  - 迁移到其他文件夹

- 为文件设置新的副本数
- 删除文件
- 设置节点标签（NodeLabel）

表 13-7 行为类型

| 行为类型               | 描述                                            | 所需参数                                                                                                                                                                                                          |
|--------------------|-----------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| MARK               | 为确定数据的冷热度并设置相应的数据存储策略。                        | <param><br><name>targettier</name><br><value>STORAGE_POLICY</value><br></param>                                                                                                                               |
| MOVE               | 为设置数据存储策略或 NodeLabel 并调用 HDFS Mover 工具。       | <param><br><name>targettier</name><br><value>STORAGE_POLICY</value><br></param><br><param><br><name>targetnodelabels</name><br><value>SOME_EXPRESSION</value><br></param><br><b>说明</b><br>用户可以配置其中任一参数或两者都配置。 |
| SET_REPL           | 为文件设置新的副本数。                                   | <param><br><name>replcount</name><br><value>INTEGER</value><br></param>                                                                                                                                       |
| MOVE_TARGET_FOLDER | 将文件移动到目标文件夹。如果“overwrite”参数为“true”，则目标路径将被覆盖。 | <param><br><name>target</name><br><value>PATH</value><br></param><br><param><br><name>overwrite</name><br><value>true/false</value><br></param><br><b>说明</b><br>“overwrite”是可选参数，如果未配置，则默认值为“false”。          |
| DELETE             | 删除文件。                                         | NA                                                                                                                                                                                                            |

## 配置描述

必须定期调用迁移工具，并需要在客户端的“hdfs-site.xml”文件中进行以下配置。

表 13-8 参数描述

| 参数                                    | 描述                                                              | 默认值                                                                  |
|---------------------------------------|-----------------------------------------------------------------|----------------------------------------------------------------------|
| dfs.auto-data-movement.policy.classes | 用于指定默认的数据迁移策略。<br><b>说明</b><br>当前只支持 DefaultDataMovementPolicy。 | com.huawei.hadoop.hdfs.datamovement.policy.DefaultDataMovementPolicy |
| dfs.auto.data.mover.id                | 冷热数据迁移输出（行为状态）文件的名称。                                            | 当前系统时间（毫秒）                                                           |
| dfs.auto.data.mover.output.dir        | 冷热数据迁移输出在HDFS中的目录名称。迁移工具将在此处写入行为状态文件。                           | /system/datamovement                                                 |

DefaultDataMovementPolicy拥有配置文件“default-datamovement-policy.xml”。用户需要定义所有基于age/accessTime的规则和在此文件中采取的行为操作，此文件必须存储在客户端的classpath中。

如下为“default-datamovement-policy.xml”配置文件的示例：

```
<policies>
<policy>
<fileset>
<file>
<name>/opt/data/1.txt</name>
</file>
<file>
<name>/opt/data/*/subpath/</name>
<excludes>
<name>/opt/data/some/subpath/sub1</name>
</excludes>
</file>
</fileset>
<rules>
<rule>
<age>2w</age>
<action>
<type>MOVE</type>
<params>
<param>
<name>targettier</name>
<value>HOT</value>
</param>
</params>
</action>
</rule>
</rules>
</policy>
</policies>
```

**说明**

在策略，规则和行为操作中使用的标签中，可以添加其他属性，例如“name”可用于管理用户界面（例如：Hue UI）和工具输入xml之间的映射。

示例：<policy name="Manage\_File1">

标签（Tag）说明如下：

表 13-9 配置标签（Tag）描述

标签（Tag）名称	描述	是否可重复使用
<policy>	<p>定义单一策略。</p> <ul style="list-style-type: none"> <li>idempotent属性：指定当策略中有多个规则时，如果满足当前规则，是否检查下一个规则。 示例：&lt;policy name="policy2" idempotent="true"&gt;。 其默认值为“true”，表示其中的规则和行为操作是幂等的，可以继续检查下一个规则。如果值为“false”，则将在当前规则处停止评估。</li> <li>hours_allowed属性：配置是否根据系统时间执行策略评估。hours_allowed的值是以逗号分隔的数字，范围从0到23，表示系统时间。 示例：&lt;policy name="policy1" hours_allowed="2-6,13-14"&gt; 如果当前系统时间在配置的范围之内，则继续评估。否则，将跳过评估。</li> </ul> <p><b>说明</b> 在输入XML中，每个文件仅支持一个策略。因此，文件中的所有规则必须由一个策略标签覆盖。</p>	Yes
<fileset>	为每个策略定义一组文件/文件夹。	No（在policy标签内）
<file>	定义文件和/或文件夹在<file>标签内被配置一个或者多个<name>标签。文件/文件夹名支持POSIX globs配置。	Yes（在fileset标签内）
<excludes>	在<file>标签内定义该标签，该标签下可以包含多个<name>标签，在<file>标签中配置的文件或文件夹范围下，<name>标签所包含的文件或文件夹将会被排除。文件或文件夹名支持POSIX globs配置。	No（在fileset标签内）
<rules>	针对策略定义多个规则。	No（在policy标签内）
<rule>	定义单一规则。	Yes（在rules标签内）

标签 (Tag) 名称	描述	是否可重复使用
<b>&lt;age&gt;or&lt;atime&gt;</b>	<p>定义在&lt;fileset&gt;中定义的文件age/accesstime。策略将匹配该age。age可以用 [num]y[num]m[num]w[num]d[num]h的格式表示。其中num表示数字。</p> <p>其中字母的意思如下：</p> <ul style="list-style-type: none"> <li>* y--年（一年是365天）。</li> <li>* m--月（一个月是30天）。</li> <li>* w--周（一周是7天）。</li> <li>* d--天。</li> <li>* h--小时。</li> </ul> <p>可以单独使用年，月，周，天或小时，也可以将时间组合。比如，1y2d表示1年零2天或者367天。</p> <p>如果没有单位（即数字后面没有任何上述字母），默认单位为天。</p> <p><b>说明</b> 用户可以在&lt;age&gt;和&lt;atime&gt;标签中配置“gt”（greater）和“lt”（less），默认运算符为“gt”。</p> <p>示例：&lt;age operator="lt"&gt;</p>	No（在rule标签内）
<b>&lt;action&gt;</b>	如果规则匹配，这个标签定义了要执行的action。	No（在rule标签内）
<b>&lt;type&gt;</b>	定义了action类型。当前支持的action类型是MOVE和MARK。	No（在action标签内）
<b>&lt;params&gt;</b>	定义与每个action相关的参数。	No（在action标签内）
<b>&lt;param&gt;</b>	<p>定义单个使用&lt;name&gt;和&lt;value&gt;标签的name-value格式参数。</p> <p>对于MARK和MOVE，只支持参数名“targettier”。该参数表示如果满足age规则，则指定数据存储策略。</p> <p>如果多个param中具有相同name的参数，则采用第一个参数值。</p> <p>对于MARK，支持的“targettier”参数值为“ALL_SSD”，“ONE_SSD”，“HOT”，“WARM”，“COLD”。</p> <p>对于MOVE，支持的“targettier”参数值为“ALL_SSD”，“ONE_SSD”，“HOT”，“WARM”和“COLD”。</p>	Yes（在params标签内）

对于在<file>标签下的文件/文件夹使用FileSystem#globStatus API，对于其他的使用GlobPattern类（被GlobFilter使用）。参照支持的API的细节。例如，对于



globStatus，“/opt/hadoop/\*”将匹配“/opt/hadoop”文件夹下的一切。“/opt/\*/hadoop”将匹配“opt”目录的子目录下的所有hadoop文件夹。

对于globStatus，分别匹配每个路径组件的glob模式，而对于其他的，直接匹配glob模式。

[https://hadoop.apache.org/docs/r3.1.1/api/org/apache/hadoop/fs/FileSystem.html#globStatus\(org.apache.hadoop.fs.Path\)](https://hadoop.apache.org/docs/r3.1.1/api/org/apache/hadoop/fs/FileSystem.html#globStatus(org.apache.hadoop.fs.Path))

Glob	Name	Matches
*	<i>asterisk</i>	Matches zero or more characters
?	<i>question mark</i>	Matches a single character
[ab]	<i>character class</i>	Matches a single character in the set {a, b}
[^ab]	<i>negated character class</i>	Matches a single character that is not in the set {a, b}
[a-b]	<i>character range</i>	Matches a single character in the (closed) range [a, b], where a is lexicographically less than or equal to b
[^a-b]	<i>negated character range</i>	Matches a single character that is not in the (closed) range [a, b], where a is lexicographically less than or equal to b
{a,b}	<i>alternation</i>	Matches either expression a or b
\c	<i>escaped character</i>	Matches character c when it is a metacharacter

## 行为操作示例

- MARK

```
<action>
<type>MARK</type>
<params>
<param>
<name>targettier</name>
<value>HOT</value>
</param>
</params>
</action>
```

- MOVE

```
<action>
<type>MOVE</type>
<params>
<param>
<name>targettier</name>
<value>HOT</value>
</param>
<param>
<name>targetnodeLabels</name>
<value>SOME_EXPRESSION</value>
</param>
</params>
</action>
```

- SET\_REPL

```
<action>
<type>SET_REPL</type>
<params>
<param>
<name>replcount</name>
<value>5</value>
</param>
</params>
</action>
```

- MOVE\_TO\_FOLDER

```
<action>
<type>MOVE_TO_FOLDER</type>
<params>
<param>
<name>target</name>
<value>path</value>
</param>
<param>
<name>overwrite</name>
<value>true</value>
</param>
</params>
</action>
```

#### 📖 说明

MOVE\_TO\_FOLDER操作只是将文件路径更改为目标文件夹，不会更改块位置。如果想要移动块，则需要配置一个独立的move策略。

- DELETE

```
<action>
<type>DELETE</type>
</action>
```

#### 📖 说明

- 在编写xml文件时，用户应该注意行为操作的配置和顺序。冷热数据迁移工具按照输入xml中给定的顺序执行规则。
- 如果只希望运行基于atime/age的一个规则，则按照时间逆序排列，且将idempotent属性设置为false。
- 如果为文件集配置删除操作，则在删除操作后不能再配置其他规则。
- 支持使用“-fs”选项，用于指定客户端默认的文件系统地址。

## 审计日志

冷热数据迁移工具支持以下操作的审计日志。

- 工具启动状态
- 行为类型及参数详细信息和状态
- 工具完成状态

对于启用审计日志工具，在“<HADOOP\_CONF\_DIR>/log4j.property”文件中添加以下属性。

```
autodatatool.logger=INFO, ADMTRFA
autodatatool.log.file=HDFSAutoDataMovementTool.audit
log4j.logger.com.huawei.hadoop.hdfs.datamovement.HDFSAutoDataMovementTool.audit=${autodatatool.logger}
log4j.additivity.com.huawei.hadoop.hdfs.datamovement.HDFSAutoDataMovementTool-audit=false
log4j.appender.ADMTRFA=org.apache.log4j.RollingFileAppender
log4j.appender.ADMTRFA.File=${hadoop.log.dir}/${autodatatool.log.file}
log4j.appender.ADMTRFA.layout=org.apache.log4j.PatternLayout
log4j.appender.ADMTRFA.layout.ConversionPattern=%d{ISO8601} %p %c: %m%n
log4j.appender.ADMTRFA.MaxBackupIndex=10
log4j.appender.ADMTRFA.MaxFileSize=64MB
```

#### 📖 说明

具体请参考“<HADOOP\_CONF\_DIR>/log4j\_autodata\_movment\_template.properties”文件。

## 13.4.3 Hive on Hue

Hue提供了Hive图形化管理功能，使用户可以通过界面的方式查询Hive的不同数据。


### 查询编辑器使用介绍

访问Hue WebUI，请参考[访问Hue WebUI界面](#)。


在左侧导航栏单击编辑器，然后选择“Hive”，进入“Hive”。

- 执行Hive HQL语句


在左侧选中目标数据库，也可通过单击右上角的 `default` ▾，输入目标数据库的名称以搜索目标数据库。

在文本编辑框输入Hive HQL语句，单击  或者按“Ctrl+Enter”，运行HQL语句，执行结果将在“结果”页签显示。

- 分析HQL语句

在左侧选中目标数据库，在文本编辑框输入Hive HQL语句，单击  编译HQL语句并显示语句是否正确，执行结果将在文本编辑框下方显示。


- 保存HQL语句

在文本编辑框输入Hive HQL语句，单击右上角的 ，并输入名称和描述。已保存的语句可以在“保存的查询”页签查看。


- 查看历史

单击“查询历史记录”，可查看HQL运行情况，支持显示所有语句或只显示保存的语句的运行情况。历史记录存在多个结果时，可以在输入框使用关键字进行搜索。

- 高级查询配置

单击右上角的 ，对文件、函数、设置等信息进行配置。


- 查看快捷键

单击右上角的 ，可查看所有快捷键信息。

### 元数据浏览器使用介绍

访问Hue WebUI。

- 查看Hive表的元数据

在左侧导航栏单击表 ，单击某一表名称，界面将显示Hive表的元数据信息。

- 管理Hive表的元数据


在Hive表的元数据信息界面：

- 单击右上角的“导入”可导入数据。
- 单击“概述”，在“属性”域可查看表文件的位置信息。

可查看Hive表各列字段的信息，并手动添加描述信息，注意此处添加的描述信息并不是Hive表中的字段注释信息（comment）。

- 单击“样本”可浏览数据。

- 管理Hive元数据表

单击左侧列表中的  可在数据库中根据上传的文件创建一个新表，也可手动创建一个新表。

 **注意**

Hue界面主要用于文件、表等数据的查看与分析，禁止通过Hue界面对操作对象进行删除等高危管理操作。如需操作，建议在确认对业务没有影响后通过各组件的相应操作方法进行处理，例如使用HDFS客户端对HDFS文件进行操作，使用Hive客户端对Hive表进行操作。

## 典型场景

通过Hue界面对Hive进行创建表的操作如下：

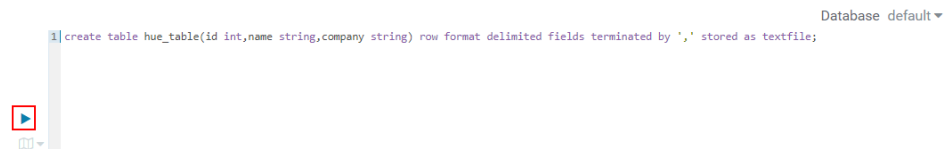
**步骤1** 单击Hue的WebUI界面左上角的 ，选择要操作的Hive实例，进入Hive命令的执行页面。

**步骤2** 在命令输入框内输入一条HQL语句，例如：

```
create table hue_table(id int,name string,company string) row format delimited fields terminated by ',' stored as textfile;
```

单击  执行HQL。

图 13-9 执行语句



**步骤3** 在命令输入框内输入：

```
show tables;
```


单击 ，查看“结果”中有创建的表hue\_table。

图 13-10 查看结果



----结束

## 13.4.4 Oozie on Hue

Hue提供了Oozie作业管理器功能，使用户可以通过界面图形化的方式使用Oozie。

### ⚠ 注意

Hue界面主要用于文件、表等数据的查看与分析，禁止通过Hue界面对操作对象进行删除等高危管理操作。如需操作，建议在确认对业务没有影响后通过各组件的相应操作方法进行处理，例如使用HDFS客户端对HDFS文件进行操作，使用Hive客户端对Hive表进行操作。

## Oozie 作业设计器使用介绍

访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

在左侧导航栏单击，选择“Workflow”。

在作业设计器，支持用户创建MapReduce、Java、Streaming、Fs、Ssh、Shell和DistCp作业。

## 仪表盘使用介绍

访问Hue WebUI。

选择右上角“作业”，进入“作业浏览器”。

支持查看Workflow、Coordinator和Bundles作业的运行情况。



## 编辑器使用介绍

访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

在左侧导航栏单击，然后选择“Workflow”。

支持创建Workflow、计划和Bundles的操作。支持提交运行、共享、复制和导出已创建的应用。

- 每个Workflow可以包含一个或多个作业，形成完整的工作流，用于实现指定的业务。  
创建Workflow时，可直接在Hue的编辑器设计作业，并添加到Workflow中。
- 每个计划可定义一个时间触发器，用于定时触发执行一个指定的Workflow。不支持多个Workflow。
- 每个Bundles可定义一个集合，用于触发执行多个计划，使不同Workflow批量执行。

## 13.5 Hue 常用配置参数

### 参数入口

参数入口，请参考[修改集群服务配置参数](#)进入Hue服务“全部配置”页面。

### 参数说明

Hue常用参数请参见[表13-10](#)。

表 13-10 Hue 常用参数

配置参数	说明	缺省值	范围
HANDLER_ACCESSLOG_LEVEL	Hue的访问日志级别。	DEBUG	<ul style="list-style-type: none"> <li>• ERROR</li> <li>• WARN</li> <li>• INFO</li> <li>• DEBUG</li> </ul>
HANDLER_AUDITLOG_LEVEL	Hue的审计日志级别。	DEBUG	<ul style="list-style-type: none"> <li>• ERROR</li> <li>• WARN</li> <li>• INFO</li> <li>• DEBUG</li> </ul>
HANDLER_ERRORLOG_LEVEL	Hue的错误日志级别。	ERROR	<ul style="list-style-type: none"> <li>• ERROR</li> <li>• WARN</li> <li>• INFO</li> <li>• DEBUG</li> </ul>

配置参数	说明	缺省值	范围
HANDLER_LOGFILE_LEVEL	Hue的运行日志级别。	INFO	<ul style="list-style-type: none"><li>• ERROR</li><li>• WARN</li><li>• INFO</li><li>• DEBUG</li></ul>
HANDLER_LOGFILE_MAXBACKUPINDEX	Hue日志文件最大个数。	20	1 ~ 999
HANDLER_LOGFILE_SIZE	Hue日志文件最大大小。	5MB	-

Hue自定义参数请参见表13-11。以下自定义参数仅MRS 3.1.2及之后版本适用。

表 13-11 Hue 自定义参数

配置参数	参数描述
dfs.customized.configs	添加全局hdfs-site.xml中用户自定义配置项
hbase.customized.configs	添加全局hbase-site.xml中用户自定义配置项
hive.customized.configs	添加全局hive-site.xml中用户自定义配置项

## 13.6 Hue 日志介绍

### 日志描述

**日志路径：**Hue相关日志的默认存储路径为“/var/log/Bigdata/hue”（运行日志），“/var/log/Bigdata/audit/hue”（审计日志）。

**日志归档规则：**Hue的日志启动了自动压缩归档功能，默认情况下，当“access.log”、“error.log”、“runcpserver.log”和“hue-audits.log”大小超过5MB的时候，会自动压缩。最多保留最近的20个压缩文件，压缩文件保留个数和压缩文件阈值可以配置。

表 13-12 Hue 日志列表

日志类型	日志文件名	描述
运行日志	access.log	访问日志。
	error.log	错误日志。
	gsdb_check.log	gaussDB检查日志。

日志类型	日志文件名	描述
	kt_renewer.log	Kerberos认证日志。
	kt_renewer.out.log	Kerberos认证日志的异常输出日志。
	runcpserver.log	操作记录日志。
	runcpserver.out.log	进程运行异常日志。
	supervisor.log	进程启动日志。
	supervisor.out.log	进程启动异常日志。
	dbDetail.log	数据库初始化日志
	initSecurityDetail.log	keytab文件下载初始化日志。
	postinstallDetail.log	Hue服务安装后工作日志。
	prestartDetail.log	Prestart日志。
	statusDetail.log	Hue服务健康状态日志。
	startDetail.log	启动日志。
	get-hue-ha.log	Hue HA状态日志。
	hue-ha-status.log	Hue HA状态监控日志。
	get-hue-health.log	Hue健康状态日志。
	hue-health-check.log	Hue健康检查日志。
	hue-refresh-config.log	Hue配置刷新日志。
	hue-script-log.log	Manager界面的Hue操作日志。
	hue-service-check.log	Hue服务状态监控日志。
	db_pwd.log	Hue连接DBService数据库密码修改日志
	modifyDBPwd_日期.log	-
	watch_config_update.log	参数更新日志。
审计日志	hue-audits.log	审计日志。

## 日志级别

Hue提供了如表13-13所示的日志级别。

日志的级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。



表 13-13 日志级别

级别	描述
ERROR	ERROR表示系统运行的错误信息。
WARN	WARN表示当前事件处理存在异常信息。
INFO	INFO表示记录系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1 参考[修改集群服务配置参数](#)进入Hue服务“全部配置”页面。
- 步骤2 在左侧导航栏选择需修改的角色所对应的“日志”菜单。
- 步骤3 在右侧选择所需修改的日志级别。
- 步骤4 保存配置，在弹出窗口中单击“确定”使配置生效。
- 步骤5 重新启动配置过期的服务或实例以使配置生效。

----结束

## 日志格式

Hue的日志格式如下所示：

表 13-14 日志格式

日志类型	格式	示例
运行日志	<dd-MM-yy HH:mm:ss,SSS><日志事件 的发生位置><log level><log中的message>	[03/Nov/2014 11:57:19 ] middleware   INFO   Unloading MimeTypeJSFileFixStrea mingMiddleware.
	<Log Level><时间格式 ><yyyy-MM-dd HH:mm:ss,SSS><日志事件 的发生位置><log中的 message>	INFO : CST 2014-11-06 11:22:52 hue-ha- status.sh : update 4 <= 15:myHostName=10.0.0. 250 ACTIVE=10.0.0.250

日志类型	格式	示例
审计日志	<UserName><yyyy-MM-dd HH:mm:ss,SSS><审计操作描述><资源参数><url><是否允许><审计操作><ip地址>	{ "username": "admin", "eventTime": "2014-11-06 10:28:34", "operationText": "Successful login for user: admin", "service": "accounts", "url": "/ accounts/login/", "allowed": true, "operation": "USER_LOGIN", "ipAddress": "10.0.0.250"} }

## 13.7 Hue 常见问题

### 13.7.1 使用 Hive 输入 use database 语句失效

#### 问题

使用Hive的时候，在输入框中输入了**use database**的语句切换数据库，重新在输入框内输入其他语句，为什么数据库没有切换过去？

#### 回答

在Hue上使用Hive有区别于用Hive客户端使用Hive，Hue界面上有选择数据库的按钮，当前SQL执行的数据库以界面上显示的数据库为准。

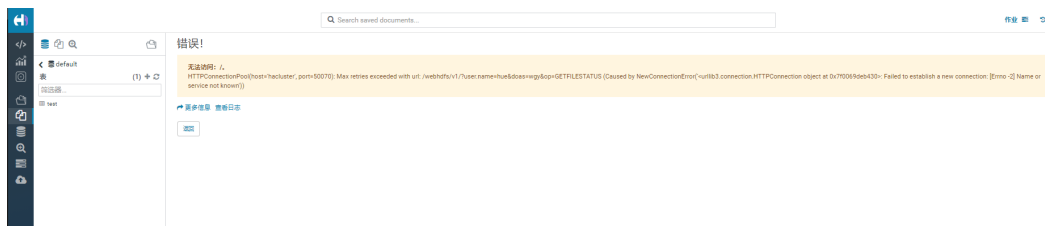
与此相关的还有设置参数等session级别的一次性操作，都应该使用界面功能进行设置，不建议使用输入语句进行操作。

若是必须使用输入语句进行操作，需保证所有语句在同一个输入框内。

### 13.7.2 使用 Hue WebUI 访问 HDFS 文件失败

#### 问题

在使用Hue WebUI访问HDFS文件时，报如下图所示无法访问的错误提示，该如何处理？



## 回答

1. 查看登录Hue WebUI的用户是否具有“hadoop”用户组权限。
2. 查看HDFS服务是否安装了HttpFS实例且运行正常。如果未安装HttpFS实例，需手动安装并重启Hue服务。

### 13.7.3 在 Hue 页面上传大文件失败

#### 问题

通过Hue页面上传大文件时，上传失败。

#### 回答

1. 不建议使用Hue文件浏览器上传大文件，大文件建议使用客户端通过命令上传。
2. 如果必须使用Hue上传，参考以下步骤修改Httpd的参数：

- a. 以omm用户登录主管理节点。
- b. 执行以下命令编辑“httpd.conf”配置文件。

```
vi $BIGDATA_HOME/om-server/Apache-httpd-*/conf/httpd.conf
```

- c. 搜索21201，在</VirtualHost>配置中加上“RequestReadTimeout handshake=0 header=0 body=0”，如下所示。

```
...
<VirtualHost *:21201>
 ServerName https://10.112.16.93:21201
 AllowEncodedSlashes On
 SSLProxyEngine On
 ProxyRequests Off
 TraceEnable off
 ProxyTimeout 1200
 RewriteEngine on
 RewriteMap proxylist dbm:${BIGDATA_ROOT_HOME}/om-server_*/Apache-httpd-*/conf/
 proxylist.dbm

 RewriteRule ^(\/*)$ ${proxylist:/Hue/Hue/21201}$1 [E=TARGET_PATH:$1,L,P]

 Header edit Location ^(!https://10.112.16.93:20009|https://
 10.112.16.93:21201)http[s]?://[^\/*]*$ https://10.112.16.93:21201$1

 ProxyPassReverseCookiePath / / interpolate

 SSLEngine On
 SSLProxyProtocol All +TLSv1.2 -SSLv2 -SSLv3 -TLSv1 -TLSv1.1
 SSLProtocol ALL +TLSv1.2 -SSLv2 -SSLv3 -TLSv1 -TLSv1.1
 SSLCipherSuite ECDHE-RSA-AES256-GCM-SHA384:ECDHE-ECDSA-AES256-GCM-
 SHA384:ECDHE-RSA-AES128-GCM-SHA256:ECDHE-ECDSA-AES128-GCM-SHA256:DHE-DSS-
 AES256-GCM-SHA384:DHE-RSA-AES256-GCM-SHA384:DHE-DSS-AES128-GCM-SHA256:DHE-
 RSA-AES128-GCM-SHA256
 SSLProxyCheckPeerName off
 SSLProxyCheckPeerCN off
 SSLCertificateFile "${BIGDATA_ROOT_HOME}/om-server_*/Apache-httpd-*/conf/security/
 proxy_ssl.cert"
 SSLCertificateKeyFile "${BIGDATA_ROOT_HOME}/om-server_*/Apache-httpd-*/conf/security/
 server.key"
 SSLProxyCACertificateFile ${BIGDATA_ROOT_HOME}/om-server_*/apache-tomcat-*/conf/
 security/tomcat.crt
 SSLCertificateChainFile "${BIGDATA_ROOT_HOME}/om-server_*/Apache-httpd-2.4.39/conf/
 security/proxy_chain.cert"
 RequestReadTimeout handshake=0 header=0 body=0
</VirtualHost>
...
```

- d. 执行 `pkill -9 httpd` 命令结束 `httpd` 进程，并等待自动重启 `httpd`。

## 13.7.4 Hue WebUI 中 Oozie 编辑器的时区设置问题

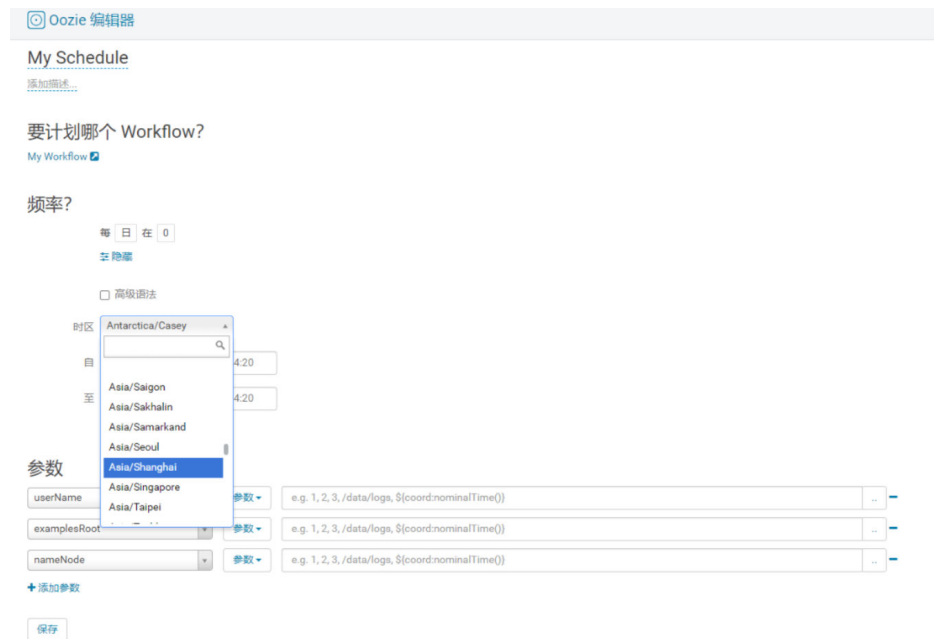
### 问题

在 Hue 设置 Oozie 工作流调度器的时区时，部分时区设置会导致任务提交失败。

### 回答

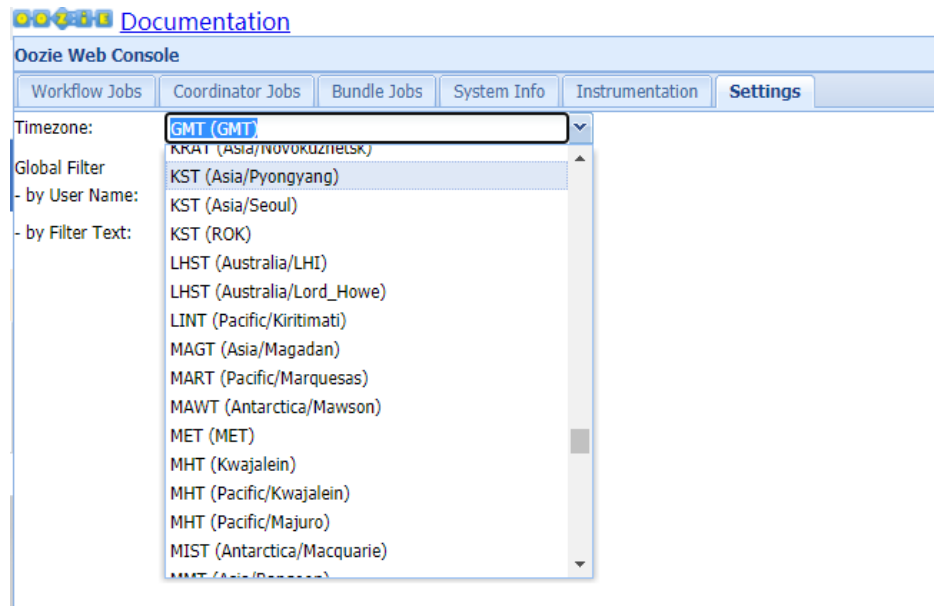
部分时区存在适配问题，建议时区选择“Asia/Shanghai”，如图 13-11 所示。

图 13-11 时区选择



支持的时区可以参考 Oozie WebUI 页面“Settings”页签的“Timezone”，如图 13-12。

图 13-12 时区参考



## 13.7.5 访问 Hue 原生页面时间长，文件浏览器报错 Read timed out

### 问题

访问 Hue 原生页面时页面加载时间较长，访问 Hue 的 HDFS 文件浏览器报错 Read timed out，如何解决。

### 回答

检查 HDFS 服务中是否安装 Httpfs 实例。

- 否，请联系运维人员处理。
- 是，重启 HttpFS 实例解决。

# 14 使用 Impala

## 14.1 Impala 客户端使用实践

Impala是用于处理存储在Hadoop集群中的大量数据的MPP（大规模并行处理）SQL查询引擎。它是一个用C++和Java编写的开源软件。与其他Hadoop的SQL引擎相比，它拥有高性能和低延迟的特点。

### 背景信息

假定用户开发一个应用程序，用于管理企业中的使用A业务的用户信息，使用Impala客户端实现A业务操作流程如下：

#### 普通表的操作：

- 创建用户信息表user\_info。
- 在用户信息中新增用户的学历、职称信息。
- 根据用户编号查询用户姓名和地址。
- A业务结束后，删除用户信息表。

表 14-1 用户信息

编号	姓名	性别	年龄	地址
12005000201	A	男	19	A城市
12005000202	B	女	23	B城市
12005000203	C	男	26	C城市
12005000204	D	男	18	D城市
12005000205	E	女	21	E城市
12005000206	F	男	32	F城市
12005000207	G	女	29	G城市
12005000208	H	女	30	H城市

编号	姓名	性别	年龄	地址
12005000209	I	男	26	I城市
12005000210	J	女	25	J城市

## 前提条件

- 已安装客户端，例如安装目录为“/opt/hadoopclient”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- MRS 3.x及之后版本的Impala客户端节点（Euler2.9及以上操作系统）需要安装Python2版本，具体请参考[Impala客户端安装Python2](#)。

## 操作步骤

**步骤1** 以客户端安装用户，登录安装客户端的节点。

**步骤2** 执行以下命令，切换到客户端安装目录。

```
cd /opt/hadoopclient
```

**步骤3** 执行以下命令配置环境变量。

```
source bigdata_env
```

**步骤4** 运行Impala客户端命令，实现A业务。

- **内部表的操作：**

直接执行Impala组件的客户端命令**impala-shell**

### 说明

默认情况下，**impala-shell**尝试连接到localhost的21000端口上的Impala守护程序。如需连接到其他主机，请使用**-i <host:port>**选项，例如：`impala-shell -i xxx.xxx.xxx.xxx:21000`。要自动连接到特定的Impala数据库，请使用**-d <database>**选项。例如，如果您的所有Kudu表都位于数据库“impala\_kudu”中，则**-d impala\_kudu**可以使用此数据库。要退出Impala Shell，请使用**quit**命令。

- 根据[表14-1](#)创建用户信息表user\_info并添加相关数据。

```
create table user_info(id string,name string,gender string,age int,addr string);
insert into table user_info(id,name,gender,age,addr) values("12005000201","A","男",19,"A城市");
.....（其他语句相同）
```
  - 在用户信息表user\_info中新增用户的学历、职称信息。  
以增加编号为12005000201的用户的学历、职称信息为例，其他用户类似。

```
alter table user_info add columns(education string,technical string);
```
  - 根据用户编号查询用户姓名和地址。  
以查询编号为12005000201的用户姓名和地址为例，其他用户类似。

```
select name,addr from user_info where id='12005000201';
```
  - 删除用户信息表。

```
drop table user_info;
```
- **外部分区表的操作：**  
创建外部分区表并导入数据

- a. 创建外部表数据存储路径。  
**kinit hive**（安全模式需要执行kinit，普通模式无需执行）

#### 📖 说明

用户hive需要具有Hive管理员权限。

```
hdfs dfs -mkdir /hive
hdfs dfs -mkdir /hive/user_info
```

- b. 建表。

```
impala-shell
```

#### 📖 说明

默认情况下，**impala-shell**尝试连接到localhost的21000端口上的Impala守护程序。如需连接到其他主机，请使用**-i <host:port>**选项，例如：`impala-shell -i xxx.xxx.xxx.xxx:21000`。要自动连接到特定的Impala数据库，请使用**-d <database>**选项。例如，如果您的所有Kudu表都位于数据库“impala\_kudu”中，则**-d impala\_kudu**可以使用此数据库。要退出Impala Shell，请使用**quit**命令。

```
create external table user_info(id string,name string,gender string,age int,addr string)
partitioned by(year string) row format delimited fields terminated by ' ' lines terminated by '\n'
stored as textfile location '/hive/user_info';
```

#### 📖 说明

- `fields terminated`指明分隔的字符,如按空格分隔, ' '。
- `lines terminated` 指明分行的字符, 如按换行分隔, '\n'。
- `/hive/user_info`为数据文件的路径。

- c. 导入数据。

- i. 使用insert语句插入数据。

```
insert into user_info partition(year="2018") values ("12005000201","A","男",19,"A城市");
```

- ii. 使用load data命令导入文件数据。

1) 根据表14-1数据创建文件。如，文件名为txt.log，以空格拆分字段，以换行符作为行分隔符。

2) 上传文件至hdfs。

```
hdfs dfs -put txt.log /tmp
```

3) 加载数据到表中。

```
load data inpath '/tmp/txt.log' into table user_info partition
(year='2018');
```

- d. 查询导入数据。

```
select * from user_info;
```

- e. 删除用户信息表。

```
drop table user_info;
```

----结束

## 14.2 访问 Impala WebUI 界面

用户可以通过Impala的WebUI，在图形化界面查看Impala作业的相关信息。Impala的WebUI根据实例不同分为如下三种：



- StateStore WebUI：用于管理节点。
- Catalog WebUI：用于查看元数据。

## 前提条件

已安装Impala服务的集群。

## 访问 StateStore WebUI

- 步骤1** 登录Manager页面，请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
- 步骤2** 选择“服务 > Impala”。
- 步骤3** 在“概览”页面，单击“基本信息”区域中“StateStore WebUI”后的“StateStore(Statestore)”，打开StateStore的WebUI页面。

图 14-1 StateStore WebUI



----结束

## 访问 Catalog WebUI

- 步骤1** 登录Manager页面，请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
- 步骤2** 选择“服务 > Impala”。
- 步骤3** 在“概览”页面，单击“基本信息”区域中“Catalog WebUI”后的“Catalog(Catalog)”，打开Catalog的WebUI页面。

----结束

## 访问 Impalad WebUI

**Impala 3.4.0版本:**

- 步骤1** 登录Manager页面，请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
- 步骤2** 选择“服务 > Impala > 实例”。

**步骤3** 移动鼠标至“角色”列Impalad实例上，在浏览器底部会显示Impalad实例的ID，并记录该ID。

图 14-2 Impalad 实例 ID



角色	运行状态	配置状态	主机名称	管理IP	业务IP
<input type="checkbox"/> Catalog	● 未启动	● 已同步	node-group-1YrJc0...	1	
<input type="checkbox"/> Catalog	● 未启动	● 已同步	node-group-1YrJc0...	1	
<input type="checkbox"/> Impalad	● 未启动	● 已同步	ClickHousekoHs00...	1	3
<input type="checkbox"/> Impalad	● 未启动	● 已同步	ClickHousekoHs00...	1	1
<input type="checkbox"/> StateStore	● 未启动	● 已同步	node-group-1YrJc0...	1	
<input type="checkbox"/> StateStore	● 未启动	● 已同步	node-group-1YrJc0...	1	

https://10.94.9.243:9022/mrsmanager/#/app/cluster/service/Impala/Impala/instance/Impalad/Impalad-DEFAULT/107/10.0.109/ClickHousekoHs0002.mrs-uot0.com/status/detail

**步骤4** 访问StateStore WebUI，具体参考[访问StateStore WebUI](#)。

**步骤5** 修改StateStore WebUI的URL地址中的“StateStore/xx”为“Impalad/xx”并访问修改后的URL，其中xx为**步骤3**中获取的ID，如下所示：

https://10.94.9.243:9022/component/Impala/StateStore/108/

修改为：

https://10.94.9.243:9022/component/Impala/Impalad/107/

----结束

**Impala 4.3.0版本：**

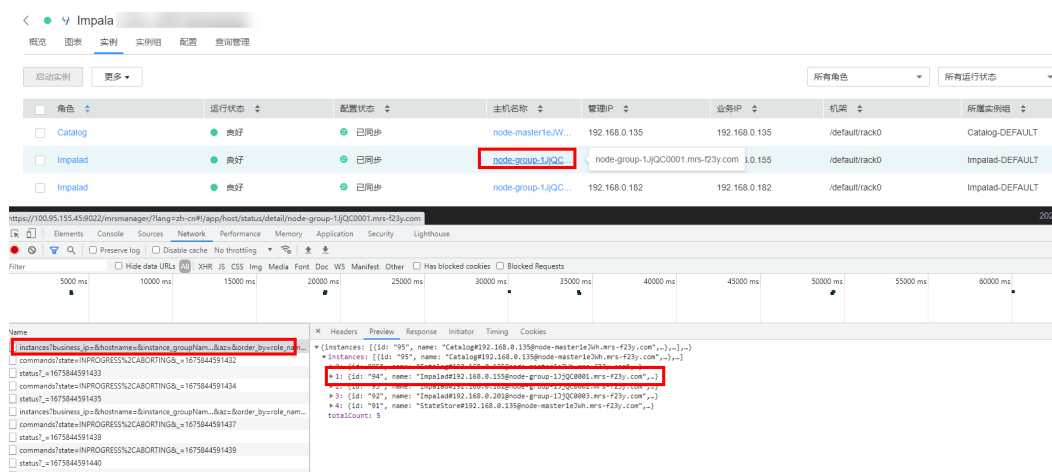
**步骤1** 登录Manager页面，请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。

**步骤2** 选择“服务 > Impala > 实例”。

**步骤3** 移动鼠标至“主机名称”列的主机名称上，键盘敲击“F12”，按照截图选择后页面显示如下内容，获取id后的数值，例如本例中的94。

其中94为样例值，实际值请以实际环境为准。

图 14-3 Impalad 实例



步骤4 参考访问StateStore WebUI。

步骤5 修改StateStore WebUI的URL地址中的“StateStore/xx”为“Impalad/xx”并访问修改后的URL，其中xx为步骤3中获取的数值。

---结束

## 14.3 使用 Impala 操作 Kudu 表

您可以使用Impala的SQL语法插入、查询、更新和删除Kudu中的数据，作为使用Kudu API构建自定义Kudu应用程序的替代方案。

### 前提条件

已安装集群完整客户端。例如安装目录为“/opt/Bigdata/client”，以下操作的客户端目录只是举例，请根据实际安装目录修改。

### Impala on Kudu

步骤1 登录安装客户端的节点。

步骤2 执行如下命令初始化环境变量。

```
source /opt/Bigdata/client/bigdata_env
```

步骤3 若集群开启Kerberos认证，请执行如下步骤认证用户。若集群未开启Kerberos认证请跳过该步骤。

```
kinit 业务用户
```

步骤4 执行如下命令登录impala客户端。

```
impala-shell
```

#### 说明

默认情况下，impala-shell尝试连接到localhost的21000端口上的Impala守护程序。如需连接到其他主机，请使用-i <host:port>选项。要自动连接到特定的Impala数据库，请使用-d <database>选项。例如，如果您的所有Kudu表都位于数据库“impala\_kudu”中，则-d impala\_kudu可以使用此数据库。要退出Impala Shell，请使用以下命令quit。

**步骤5** 执行如下命令创建Impala表并导入已准备好的数据，例如/tmp/data10。

```
create table dataorigin (name string,age string,pt string, date_p date) row
format delimited fields terminated by ',' stored as textfile;
```

```
load data inpath '/tmp/data10' overwrite into table dataorigin;
```

**步骤6** 执行如下命令创建Kudu表，其中kudu.master\_addresses地址为KuduMaster实例的IP，请根据实际集群地址填写。

```
create table dataorigin2 (name string,age string,pt string, date_p date,
primary key(name)) stored as kudu
TBLPROPERTIES('kudu.master_addresses'='192.168.190.164:7051,192.168.204.1
78:7051,192.168.244.63:7051');
```

若impala集群开启了Ranger鉴权，上述命令会报错，需要新增Impalad角色自定义配置--

```
kudu_master_hosts=192.168.190.164:7051,192.168.204.178:7051,192.168.244.63:70
51，然后重启Impala集群，使用如下命令创建kudu表：
```

```
create table dataorigin2 (name string,age string,pt string, date_p date, primary
key(name)) stored as kudu
```

**步骤7** 执行如下命令操作Kudu表。

1. 插入数据

```
insert into dataorigin2 select * from dataorigin;
```

2. 更新数据

```
UPDATE dataorigin2 SET date_p="2021-03-31" where age="73";
```

3. 更新或插入行

```
UPSERT INTO dataorigin2 VALUES ("spjted","75","28","2021-03-32");
UPSERT INTO dataorigin2 VALUES ("kwhakb","92","29","2021-03-33");
UPSERT INTO dataorigin2 VALUES ("oftrkf","13","30","2021-03-34");
UPSERT INTO dataorigin2 VALUES ("kiewti","36","31","2021-03-35");
UPSERT INTO dataorigin2 VALUES ("rknmql","98","32","2021-03-36");
UPSERT INTO dataorigin2 VALUES ("fwcoij","52","33","2021-03-37");
UPSERT INTO dataorigin2 VALUES ("pgvpdo","37","34","2021-03-35");
```

4. 删除行

```
DELETE FROM dataorigin2 WHERE date_p="2021-03-31";
```

----结束

## 14.4 Impala 对接外部 LDAP

本操作适用于MRS 3.1.0及之后版本。

**步骤1** 登录Manager。

**步骤2** 在Manager界面，选择“集群 > 待操作集群的名称 > 服务 > Impala > 配置 > 全部配置 > Impalad（角色） > LDAP”。

**步骤3** 配置如下参数的值。

表 14-2 参数配置

参数名称	参数描述	备注
--enable_ldap_auth	是否开启LDAP认证	【取值范围】 true或false
--ldap_bind_pattern	LDAP userDNPattern	例如： cn=#UID,ou=People,dc=huawei,dc=com或cn=%s,ou=People,dc=huawei,dc=com
--ldap_passwords_in_clear_ok	LDAP 密码是否以明文发送	如果设置为true，将允许LDAP密码在网络上明文发送 【取值范围】 true或false <b>说明</b> 当“--enable_ldap_auth”设置为“true”时，认证时默认没有开启Ldap TLS协议，所以需要“--ldap_passwords_in_clear_ok”参数设置为“true”，否则会导致Impalad角色启动失败。 如需开启Ldap TLS协议则需要Impalad角色的自定义配置中添加配置项“--ldap_tls”为“true”，配置之后密码将支持用密文传输。
--ldap_uri-ip	LDAP IP	-
--ldap_uri-port	LDAP 端口	【默认值】 389

**步骤4** 修改完成后，单击左上方“保存”，在弹出的对话框中单击“确定”保存配置。

**步骤5** 选择“集群 > 待操作集群的名称 > 服务 > Impala > 实例”，勾选配置状态为“配置过期”的实例，选择“更多 > 重启实例”重启受影响的Impala实例。

----结束

## 14.5 Impala 启用并配置动态资源池

本文介绍如何使用动态资源池控制impala并发。

### 问题背景

客户需要使用动态资源池控制impala并发。

### Pool Config

Property	Value
Max memory (cluster wide)	1048576
Max concurrent queries	-1
Max queue size	200
Queue Timeout (ms)	60000
Min Query MEM_LIMIT range	0
Max Query MEM_LIMIT range	0
Clamp MEM_LIMIT query option	true

1. 登录到集群的master1节点上，然后切换到omm用户下，在/home/omm目录下创建fair-scheduler.xml、llama-site.xml文件。

```
[omm@node-master1IoKo impala]$ ll
total 16
-rw-----. 1 omm wheel 708 May 11 23:40 fair-scheduler.xml
-rw-----. 1 omm wheel 1062 May 11 23:53 llama-site.xml
-rw-----. 1 omm wheel 1118 May 11 23:12 llama-site.xml.bak
-rw-----. 1 omm wheel 572 May 11 23:32 update_config.sh
[omm@node-master1IoKo impala]$
```

2. 打开fair-scheduler.xml文件，添加如下配置。

```
<allocations>
 <queue name="root">
 <aclSubmitApps> </aclSubmitApps>
 <queue name="default">
 <maxResources>4096 mb, 0 vcores</maxResources><!--参数仅供参考-->
 <aclSubmitApps>*</aclSubmitApps>
 </queue>
 <queue name="development">
 <maxResources>2048 mb, 0 vcores</maxResources><!--参数仅供参考-->
 <aclSubmitApps>admin</aclSubmitApps>
 </queue>
 <queue name="production">
 <maxResources>7168 mb, 0 vcores</maxResources><!--参数仅供参考-->
 <aclSubmitApps>omm</aclSubmitApps>
 </queue>
 </queue>
 <queuePlacementPolicy>
 <rule name="specified" create="false"/>
 <rule name="default" />
 </queuePlacementPolicy>
</allocations>
```

3. 打开llama-site.xml文件，添加如下配置：

```
<?xml version="1.0" encoding="UTF-8"?>
<configuration>
 <property>
 <name>llama.am.throttling.maximum.placed.reservations.root.default</name>
 <value>1</value>
 </property>
 <property>
 <name>llama.am.throttling.maximum.queued.reservations.root.default</name>
 <value>2</value><!--参数仅供参考-->
 </property>
 <property>
 <name>impala.admission-control.pool-default-query-options.root.default</name>
 <value>mem_limit=128m,query_timeout_s=20,max_io_buffers=10</value>
 </property>
 <property>
 <name>impala.admission-control.pool-queue-timeout-ms.root.default</name>
 <value>30000</value><!--参数仅供参考-->
 </property>
</configuration>
```

```
<property>
 <name>impala.admission-control.max-query-mem-limit.root.default</name>
 <value>307200000</value><!--3GB--><!--参数仅供参考-->
</property>
<property>
 <name>impala.admission-control.min-query-mem-limit.root.default</name>
 <value>204800000</value><!--2GB-->
</property>
</property>
 <name>impala.admission-control.clamp-mem-limit-query-option.root.default.regularPool</name>
 <value>true</value>
</property>
</configuration>
```

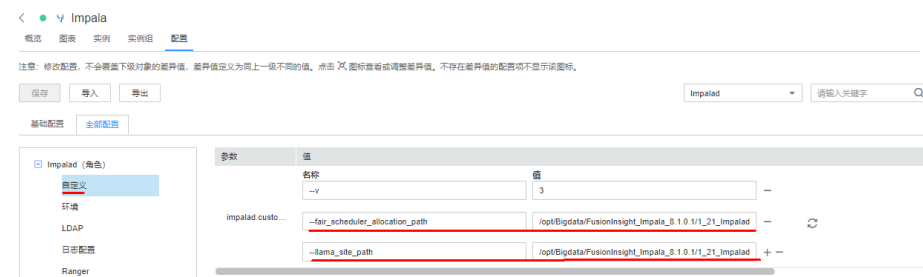
4. 执行如下命令分别将fair-scheduler.xml、llama-site.xml同步到所有的impalad节点的安装目录的etc文件夹下。

```
scp fair-scheduler.xml {impalad实例ip}:/opt/Bigdata/
FusionInsight_Impala_***/***/Impalad/etc/
scp llama-site.xml {impalad实例ip}:/opt/Bigdata/
FusionInsight_Impala_***/***/Impalad/etc/
```

```
+ scp fair-scheduler.xml 19:/opt/Bigdata/FusionInsight_Impala_8.1.0.1/1_21_Impalad/etc/
Warning: Permanently added '19' (ED25519) to the list of known hosts.
```

5. 登录到manager管理页面上，找到impala组件，然后在impalad实例中添加如下自定义配置项及值：

```
--fair_scheduler_allocation_path 值如: /opt/Bigdata/
FusionInsight_Impala_***/***/Impalad/etc/fair-scheduler.xml
--llama_site_path 值如: /opt/Bigdata/FusionInsight_Impala_***/
***/Impalad/etc/llama-site.xml
```



6. 重启impalad实例。



7. 登录到impala客户端所在的节点上，source环境变量，然后执行如下命令。  
**impala-shell -i {impalad实例ip:port} -Q request\_pool=root.default** (fair-scheduler.xml与llama-site.xml文件中配置的资源池)

```
[root@node-master1ioKo ~]# impala-shell -i 192.168.1.19:21000 -Q request_pool=root.default
Starting Impala Shell without Kerberos authentication
Opened TCP connection to 192.168.1.19:21000
Connected to 192.168.1.19:21000
Server version: impalad version 3.4.0-RELEASE RELEASE (build ac0f95df4baa94cfdc36ef370f6a432d582ac1f)

Welcome to the Impala shell.
(Impala Shell v3.4.0-RELEASE (f68c12e) built on Sat Jun 26 17:16:01 CST 2021)

The '-B' command line flag turns off pretty-printing for query results. Use this
flag to remove formatting from results you want to save for later, or to benchmark
Impala.

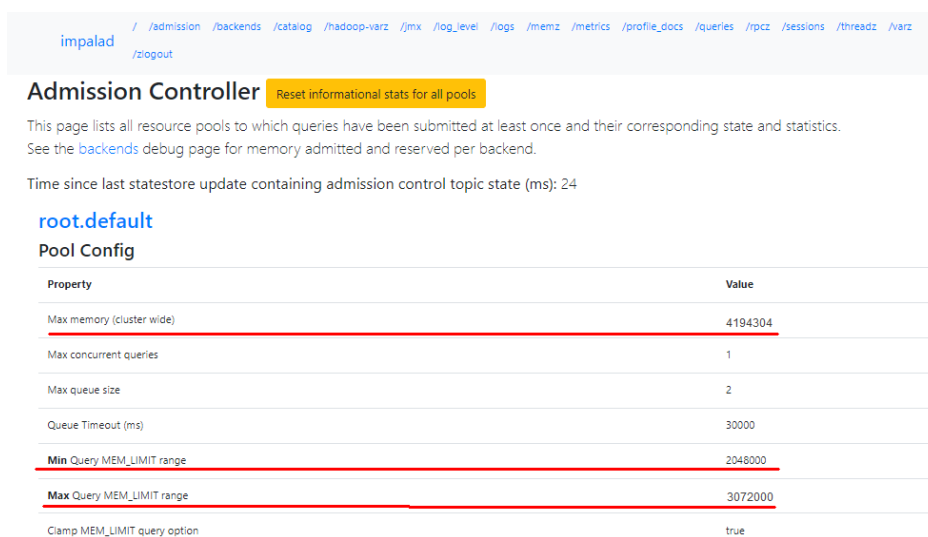
[192.168.1.19:21000] default>
```

执行SQL查询。

```
[192.168.1.19:21000] default> select * from test1;
Query: select * from test1
Query submitted at: 2022-05-12 10:01:01 (Coordinator: http://192.168.1.19:25000)
Query progress can be monitored at: http://192.168.1.19:25000/query_plan?query_id=97440454dbab28ea:35bd90d600000000
```

8. 登录到Impalad WebUI上查看资源池使用情况，确认配置已生效。

<https://{{集群控制台地址}}:9022/component/Impala/Impalad/95/>



## 14.6 使用 Impala 查询管理界面

### 操作场景

用户可以根据业务需要，在FusionInsight Manager中通过交互式查询，查看Impala的相关任务。

#### 说明

本章节内容仅适用于MRS 3.1.5及之后版本。

### 前提条件

已获取“admin”账号密码。“admin”密码在创建MRS集群时由用户指定。

### 操作步骤

- 步骤1** 进入Impala服务页面：登录FusionInsight Manager，然后选择“集群 > 待操作的集群名称 > 服务 > Impala”。
- 步骤2** 单击“查询管理”，列表默认显示所有正在进行的查询。



单击“已经结束的查询”可以查看已经完成查询的相关信息。



### 说明

用户可以根据实际情况按照慢查询运行时长、查询id、用户、所属数据库进行查询。用户可以通过“停止”操作手动停止正在进行查询的任务。

---结束

## 14.7 Impala 常见配置参数

本章节适用于MRS 3.x及后续版本。

### 参数入口

在Manager系统中，选择“集群 > 服务 > Impala > 配置”，选择“全部配置”。在搜索框中输入参数名称。

### 参数说明

#### 说明

下表仅列举了部分常用参数，实际参数以Manager页面为准，参数详情请参见官网[https://docs.cloudera.com/documentation/enterprise/6/properties/6.3/topics/cm\\_props\\_cdh630\\_impala.html](https://docs.cloudera.com/documentation/enterprise/6/properties/6.3/topics/cm_props_cdh630_impala.html)。

表 14-3 Impala 常用参数

配置参数	说明	默认值	范围
impalad.customized.configs	impalad进程的自定义配置项。	-	-
--enable_ldap_auth	是否开启ldap认证。	false	true或false
--ldap_bind_pattern	ldap userDNPattern 例如: cn=%s,ou=People,dc=huawei,dc=com	-	-
--ldap_passwords_in_clear_ok	如果设置为true，将允许ldap密码在网络上明文发送(不含TLS/SSL)。	false	true或false
--ldap_uri-ip	ldap ip	-	-
--ldap_uri-port	ldap port	389	-

配置参数	说明	默认值	范围
--max_log_files	进程日志的最大文件个数。	10	-
--max_log_size	进程的日志文件大小最大值，单位 MB。	200	-
statestored.customized.configs	Statestored进程的自定义配置项。	-	-
catalogd.customized.configs	Catalogd进程的自定义配置项。	-	-

## 14.8 Impala 常见问题

### 14.8.1 Impala 服务是否支持磁盘热插拔

#### 问题

MRS集群中Impala服务是否支持磁盘热插拔？

#### 回答

Impala服务的数据一般是存储在HDFS或者OBS（对象存储服务）中，无需直接使用本地节点的磁盘。

仅Impalad实例在业务查询执行过程中由于内存空间不足，才需要溢写到磁盘（由--scratch\_dirs指定）。

由于是非多副本存储的临时数据，不提供磁盘热插拔能力。

# 15 使用 Kafka

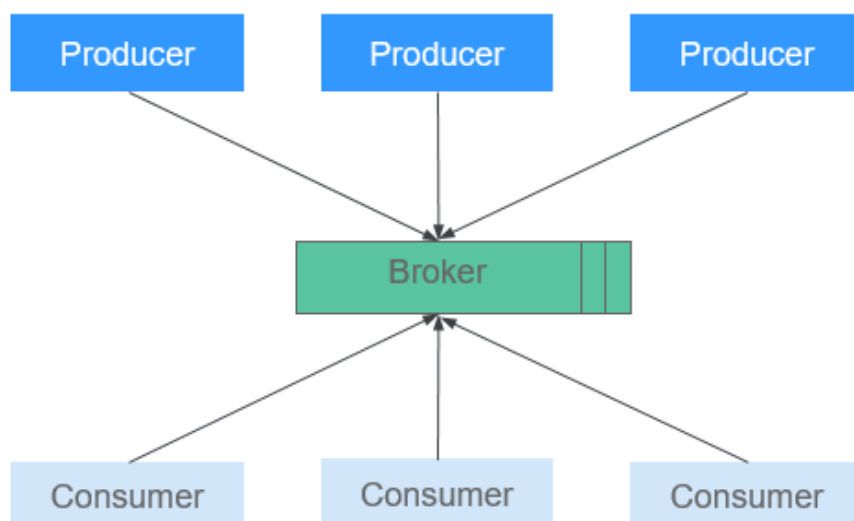
## 15.1 Kafka 数据消费概述

Kafka是一个分布式的、分区的、多副本的消息发布-订阅系统，它提供了类似于JMS的特性，但在设计上完全不同，它具有消息持久化、高吞吐、分布式、多客户端支持、实时等特性，适用于离线和在线的消息消费，如常规的消息收集、网站活性跟踪、聚合统计系统运营数据（监控数据）、日志收集等大量数据的互联网服务的数据收集场景。

### Kafka 结构

生产者（Producer）将消息发布到Kafka主题（Topic）上，消费者（Consumer）订阅这些主题并消费这些消息。在Kafka集群上一个服务器称为一个Broker。对于每一个主题，Kafka集群保留一个用于缩放、并行化和容错性的分区（Partition）。每个分区是一个有序、不可变的消息序列，并不断追加到提交日志文件。分区的消息每个也被赋值一个称为偏移顺序（Offset）的序列化编号。

图 15-1 Kafka 结构



## Kafka UI

Kafka UI提供Kafka Web服务，通过界面展示Kafka集群中Broker、Topic、Partition、Consumer等功能模块的基本信息，同时提供Kafka服务常用命令的界面操作入口。该功能作为Kafka Manager替代，提供符合安全规范的Kafka Web服务。

通过Kafka UI可以进行以下操作：

- 支持界面检查集群状态（主题，消费者，偏移量，分区，副本，节点）
- 支持界面执行集群内分区重新分配
- 支持界面选择配置创建主题
- 支持界面删除主题（Kafka服务设置了参数“delete.topic.enable = true”）
- 支持为已有主题增加分区
- 支持更新现有主题的配置
- 可以为分区级别和主题级别度量标准启用JMX查询

## 15.2 Kafka 用户权限管理

### 15.2.1 Kafka 用户权限说明

#### 操作场景

在启用Kerberos认证的集群中，用户使用Kafka前需要拥有对应的权限。MRS集群支持将Kafka的使用权限，授予不同用户。

Kafka默认用户组如表15-1所示。

#### 说明

在MRS 3.x及之后版本中，Kafka支持两种鉴权插件：“Kafka开源自带鉴权插件”和“Ranger鉴权插件”。

本章节描述的是基于“Kafka开源自带鉴权插件”的用户权限管理。若想使用“Ranger鉴权插件”，请参考[添加Kafka的Ranger访问权限策略](#)。

表 15-1 Kafka 默认用户组

用户组名称	描述
kafkaadmin	Kafka管理员用户组。添加入本组的用户，拥有所有主题的创建，删除，授权及读写权限。
kafkasuperuser	Kafka高级用户组。添加入本组的用户，拥有所有主题的读写权限。
kafka	Kafka普通用户组。添加入本组的用户，需要被kafkaadmin组用户授予特定主题的读写权限，才能访问对应主题。

## 前提条件

- 已安装客户端。
- 用户已明确业务需求，并准备一个属于kafkaadmin组的用户，作为Kafka管理员用户。例如“admin”。

## 操作步骤

**步骤1** 进入ZooKeeper实例页面：

- MRS3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > ZooKeeper > 实例”。

### 📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > ZooKeeper > 实例”。

**步骤2** 查看ZooKeeper角色实例的IP地址。

记录ZooKeeper角色实例其中任意一个的IP地址即可。

**步骤3** 根据业务情况，准备好客户端，登录安装客户端的节点。

请根据客户端所在位置，参考[使用MRS客户端](#)章节，登录安装客户端的节点。

**步骤4** 执行以下命令，切换到客户端目录，例如“/opt/client/Kafka/kafka/bin”。

```
cd /opt/client/Kafka/kafka/bin
```

**步骤5** 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

**步骤6** 执行以下命令，进行用户认证。

```
kinit 组件业务用户
```

**步骤7** MRS 3.x之前版本：选择业务需要对应的场景，管理Kafka用户权限。

- 查看某个主题的权限控制列表  

```
sh kafka-acls.sh --authorizer-properties zookeeper.connect=ZooKeeper角色实例所在节点IP地址:2181/kafka --list --topic 主题名称
```
- 为某个用户添加生产者的权限  

```
sh kafka-acls.sh --authorizer-properties zookeeper.connect=ZooKeeper角色实例所在节点IP地址:2181/kafka --add --allow-principal User:用户名 --producer --topic 主题名称
```
- 删除某个用户的生产者权限  

```
sh kafka-acls.sh --authorizer-properties zookeeper.connect=ZooKeeper角色实例所在节点IP地址:2181/kafka --remove --allow-principal User:用户名 --producer --topic 主题名称
```
- 为某个用户添加消费者的权限

```
sh kafka-acls.sh --authorizer-properties zookeeper.connect=ZooKeeper角色实例所在节点IP地址:2181/kafka --add --allow-principal User:用户名 --consumer --topic 主题名称 --group 消费者组名称
```

- 删除某个用户的消费者权限

```
sh kafka-acls.sh --authorizer-properties zookeeper.connect=ZooKeeper角色实例所在节点IP地址:2181/kafka --remove --allow-principal User:用户名 --consumer --topic 主题名称 --group 消费者组名称
```

#### 📖 说明

删除权限时需要输入两次“y”确认删除权限。

**步骤8** MRS 3.x及后续版本：使用“kafka-acl.sh”进行用户授权常用命令如下。

- 查看某Topic权限控制列表：

```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --list --topic <Topic名称>
```

```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --list --topic <Topic名称>
```

- 添加给某用户Producer权限：

```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --add --allow-principal User:<用户名> --producer --topic <Topic名称>
```

```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --add --allow-principal User:<用户名> --producer --topic <Topic名称>
```

- 给某用户批量添加Producer权限

```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --add --allow-principal User:<用户名> --producer --topic <Topic名称> --resource-pattern-type prefixed
```

```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --add --allow-principal User:<用户名> --producer --topic <Topic名称> --resource-pattern-type prefixed
```

- 删除某用户Producer权限：

```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --remove --allow-principal User:<用户名> --producer --topic <Topic名称>
```

```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --remove --allow-principal User:<用户名> --producer --topic <Topic名称>
```

- 批量删除某用户Producer权限：

```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --remove --allow-principal User:<用户名> --producer --topic <Topic名称> --resource-pattern-type prefixed
```

```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --remove --allow-principal User:<用户名> --producer --topic <Topic名称> --resource-pattern-type prefixed
```

- 添加给某用户Consumer权限：

- ```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --add --allow-principal User:<用户名> --consumer --topic <Topic名称> --group <消费者组名称>
```
- ```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --add --allow-principal User:<用户名> --consumer --topic <Topic名称> --group <消费者组名称>
```
- 给某用户批量添加Consumer权限

```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --add --allow-principal User:<用户名> --consumer --topic <Topic名称> --group <消费者组名称> --resource-pattern-type prefixed
```

```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --add --allow-principal User:<用户名> --consumer --topic <Topic名称> --group <消费者组名称> --resource-pattern-type prefixed
```
  - 删除某用户Consumer权限：

```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --remove --allow-principal User:<用户名> --consumer --topic <Topic名称> --group <消费者组名称>
```

```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --remove --allow-principal User:<用户名> --consumer --topic <Topic名称> --group <消费者组名称>
```
  - 批量删除某用户Consumer权限：

```
./kafka-acls.sh --authorizer-properties zookeeper.connect=<ZooKeeper的任意一个节点的业务IP:2181/kafka > --remove --allow-principal User:<用户名> --consumer --topic <Topic名称> --group <消费者组名称> --resource-pattern-type prefixed
```

```
./kafka-acls.sh --bootstrap-server <Kafka集群IP:21007> --command-config ../config/client.properties --remove --allow-principal User:<用户名> --consumer --topic <Topic名称> --group <消费者组名称> --resource-pattern-type prefixed
```
- 结束

## 15.2.2 创建 Kafka 权限角色

### 操作场景

该任务指导MRS集群管理员创建并设置Kafka的角色。

本章节内容适用于MRS 3.x及后续版本。

#### 说明

安全模式集群支持创建Kafka角色，普通模式集群不支持创建Kafka角色。

如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加Kafka的Ranger访问权限策略](#)。

## 操作步骤

- 步骤1** 登录FusionInsight Manager，选择“系统 > 权限 > 角色”。
- 步骤2** 单击“添加角色”，然后在“角色名称”和“描述”输入角色名字与描述。
- 步骤3** 在“配置资源权限”中，选择“待操作集群的名称 > Kafka”。
- 步骤4** 根据业务需求选择权限，具体配置项，请参见表15-2

表 15-2 配置项说明

任务场景	角色授权操作
设置Kafka管理员权限	在“配置资源权限”的表格中选择“待操作集群的名称 > Kafka > Kafka Manager权限”。 <b>说明</b> 设置此权限，拥有Topic的创建、删除等权限，但是不具备任何Topic的生产和消费权限。
设置用户对Topic的生产权限	1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Kafka > Kafka Topic生产和消费权限”。 2. 在指定Topic的“权限”列，勾选“Kafka生产者权限”。
设置用户对Topic的消费权限	1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Kafka > Kafka Topic生产和消费权限”。 2. 在指定Topic的“权限”列，勾选“Kafka消费者权限”。

- 步骤5** 单击“确定”完成，返回“角色”。

----结束

### 15.2.3 配置 Kafka 用户 Token 认证信息

#### 操作场景

使用Token认证机制时对Token的操作。

本章节内容适用于MRS 3.x及后续版本的启用Kerberos认证的集群。

#### 前提条件

- MRS集群管理员已明确业务需求，并准备一个系统用户。
- 已开启Token认证机制，详细操作请参考[Kafka服务端配置](#)。
- 已安装Kafka客户端。

#### 操作步骤

- 步骤1** 以客户端安装用户，登录安装Kafka客户端的节点。



**步骤2** 切换到Kafka客户端安装目录，例如“/opt/client”。

```
cd /opt/client
```

**步骤3** 执行以下命令，配置环境变量。

```
source bigdata_env
```

**步骤4** 执行以下命令，进行用户认证。

```
kinit 组件业务用户
```

**步骤5** 执行以下命令，切换到Kafka客户端安装目录。

```
cd Kafka/kafka/bin
```

**步骤6** 使用kafka-delegation-tokens.sh对Token进行操作

- 为用户生成Token:

```
./kafka-delegation-tokens.sh --create --bootstrap-server <IP1:PORT,
IP2:PORT,...> --max-life-time-period <Long: max life period in milliseconds>
--command-config <config file> --renewer-principal User:<user name>
```

例如: `./kafka-delegation-tokens.sh --create --bootstrap-server  
192.168.1.1:21007,192.168.1.2:21007,192.168.1.3:21007 --command-  
config ../config/producer.properties --max-life-time-period -1 --renewer-  
principal User:username`

- 列出归属在特定用户下的所有Token信息:

```
./kafka-delegation-tokens.sh --describe --bootstrap-server <IP1:PORT,
IP2:PORT,...> --command-config <config file> --owner-principal User:<user
name>
```

例如: `./kafka-delegation-tokens.sh --describe --bootstrap-server  
192.168.1.1:21007,192.168.1.2:21007,192.168.1.3:21007 --command-  
config ../config/producer.properties --owner-principal User:username`

- Token有效期刷新:

```
./kafka-delegation-tokens.sh --renew --bootstrap-server <IP1:PORT,
IP2:PORT,...> --renew-time-period <Long: renew time period in milliseconds>
--command-config <config file> --hmac <String: HMAC of the delegation
token>
```

例如: `./kafka-delegation-tokens.sh --renew --bootstrap-server  
192.168.1.1:21007,192.168.1.2:21007,192.168.1.3:21007 --renew-time-  
period -1 --command-config ../config/producer.properties --hmac  
ABCDEFGH`

- 销毁Token:

```
./kafka-delegation-tokens.sh --expire --bootstrap-server <IP1:PORT,
IP2:PORT,...> --expiry-time-period <Long: expiry time period in milliseconds>
--command-config <config file> --hmac <String: HMAC of the delegation
token>
```

例如: `./kafka-delegation-tokens.sh --expire --bootstrap-server  
192.168.1.1:21007,192.168.1.2:21007,192.168.1.3:21007 --expiry-time-  
period -1 --command-config ../config/producer.properties --hmac  
ABCDEFGH`

----结束

## 15.3 Kafka 客户端使用实践

### 操作场景

该任务指导用户在运维场景或业务场景中使用Kafka客户端。

本章节适用于MRS 3.x及后续版本。

### 前提条件

- 已安装集群客户端，例如安装目录为“/opt/client”。
- 各组件业务用户由MRS集群管理员根据业务需要创建。“机机”用户需要下载keytab文件，“人机”用户第一次登录时需修改密码。（普通模式不涉及）
- 在修改集群域名后，需要重新下载客户端，以保证客户端配置文件中kerberos.domain.name配置为正确的服务端域名。

### 使用 Kafka 客户端

**步骤1** 以客户端安装用户，登录安装客户端的节点。

**步骤2** 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

**步骤3** 执行以下命令配置环境变量。

```
source bigdata_env
```

**步骤4** 执行以下命令，进行用户认证。（普通模式跳过此步骤）

```
kinit 组件业务用户
```

**步骤5** 执行以下命令切换到Kafka客户端安装目录。

```
cd Kafka/kafka/bin
```

**步骤6** 执行以下命令使用客户端工具查看帮助并使用。

- `./kafka-console-consumer.sh`: Kafka消息读取工具
- `./kafka-console-producer.sh`: Kafka消息发布工具
- `./kafka-topics.sh`: Kafka Topic管理工具

**步骤7** MRS 3.x之前版本：执行以下命令，管理Kafka主题。

- 创建主题

```
sh kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份个数 --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

- 删除主题

```
sh kafka-topics.sh --delete --topic 主题名称 --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

## 📖 说明

- 主题分区数和主题备份个数不能大于Kafka角色实例数量。
- 默认情况下，ZooKeeper的“clientPort”为“2181”。
- ZooKeeper角色实例所在节点IP地址，填写三个角色实例其中任意一个的IP地址即可。

### 步骤8 MRS 3.x及后续版本：使用kafka-topics.sh管理Kafka主题。

- 创建主题：

Topic的Partition自动划分时，默认根据节点及磁盘上已有的Partition数进行均衡划分，如果期望根据磁盘容量进行Partition划分，那么需要修改Kafka服务配置“log.partition.strategy”为“capacity”。

Kafka创建Topic时，支持基于“机架感知”和“跨AZ特性”两种选项组合生成分区及副本的分配方案且支持“--zookeeper”和“--bootstrap-server”两种方式

- 禁用机架策略 & 禁用跨AZ特性（默认策略）。

基于此策略新建的Topic的副本会完全随机分配到集群中任意节点上。

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --zookeeper ZooKeeper的任意一个节点的业务IP:clientPort/kafka
```

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties
```

其中，使用“--bootstrap-server”方式创建Topic时，需配置“rack.aware.enable=false”和“az.aware.enable=false”。

- 启用机架策略 & 禁用跨AZ特性。

基于此策略新建的Topic的各个Partition的Leader会在集群节点上随机分配，但会确保同一Partition的不同Replica会分配在不同的机架上，所以当使用此策略时，需保证各个机架内的节点个数一致，否则会导致节点少的机架上的机器负载远高于集群平均水平。

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --zookeeper ZooKeeper的任意一个节点的业务IP:clientPort/kafka --enable-rack-aware
```

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties
```

其中，使用“--bootstrap-server”方式创建Topic时，需配置“rack.aware.enable=true”和“az.aware.enable=false”。

- 禁用机架策略 & 启用跨AZ特性。

基于此策略新建的Topic的各个Partition的Leader会在集群节点上随机分配，但会确保同一Partition的不同Replica会分配在不同的AZ上，所以当使用此策略时，需保证各个AZ内的节点个数一致，否则会导致节点少的AZ上的机器负载远高于集群平均水平。

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --zookeeper ZooKeeper的任意一个节点的业务IP:clientPort/kafka --enable-az-aware
```

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties
```

其中，使用 “--bootstrap-server” 方式创建Topic时，需配置 “rack.aware.enable=false” 和 “az.aware.enable=true”。

- 启用机架策略 & 启用跨AZ特性。

基于此策略新建的Topic的各个Partition的Leader会在集群节点上随机分配，但会确保同一Partition的不同Replica会分配到不同AZ内的不同RACK上，使用此策略需保证每个AZ内的每个RACK上的节点个数一致，否则会导致集群内负载不均衡。

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --zookeeper ZooKeeper的任意一个节点的IP:clientPort/kafka --enable-rack-aware --enable-az-aware
```

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties
```

使用 “--bootstrap-server” 方式创建Topic时，需配置 “rack.aware.enable=true” 和 “az.aware.enable=true”。

## 📖 说明

- Kafka创建Topic支持 “--zookeeper” 和 “--bootstrap-server” 两种方式，区别如下：
  - “--zookeeper” 方式由客户端生成副本分配方案，社区从一开始就支持这种方式，为了降低对Zookeeper组件的依赖，社区将在后续版本中删除对这种方式的支持。基于这种方式创建Topic时，可以通过 “--enable-rack-aware” 和 “--enable-az-aware” 这两个选项自由组合来选用副本分配策略。注意：使用 “--enable-az-aware” 选项的前提是服务端开启了跨AZ特性，即服务端启动参数 “az.aware.enable” 为 “true”，否则会执行失败。
  - “--bootstrap-server” 方式由服务端生成副本分配方案，后续版本，社区将只支持这种方式来进行Topic管理。基于这种方式创建Topic时，不支持 “--enable-rack-aware” 和 “--enable-az-aware” 选项来控制副本分配策略，支持 “rack.aware.enable” 和 “az.aware.enable” 这两个服务启动参数组合来控制副本分配策略，需注意的是 “az.aware.enable” 参数不可修改，在创建集群时，如果开启跨AZ特性，会自动配置为 “true”；“rack.aware.enable” 参数支持用户自定义修改。
- 罗列主题：
  - `./kafka-topics.sh --list --zookeeper ZooKeeper的任意一个节点的IP:clientPort/kafka`
  - `./kafka-topics.sh --list --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties`
- 查看主题：
  - `./kafka-topics.sh --describe --zookeeper ZooKeeper的任意一个节点的IP:clientPort/kafka --topic 主题名称`
  - `./kafka-topics.sh --describe --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties --topic 主题名称`
- 修改主题：
  - `./kafka-topics.sh --alter --topic 主题名称 --config 配置项=配置值 --zookeeper ZooKeeper的任意一个节点的IP:clientPort/kafka`
- 扩展分区：
  - `./kafka-topics.sh --alter --topic 主题名称 --zookeeper ZooKeeper的任意一个节点的IP:clientPort/kafka --command-config Kafka/kafka/config/client.properties --partitions 扩展后分区个数`

- `./kafka-topics.sh --alter --topic 主题名称 --bootstrap-server Kafka集群IP:21007 --command-config Kafka/kafka/config/client.properties --partitions 扩展后分区个数`
  - 删除主题：
    - `./kafka-topics.sh --delete --topic 主题名称 --zookeeper ZooKeeper的任意一个节点的IP:clientPort/kafka`
    - `./kafka-topics.sh --delete --topic 主题名称 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties`
- 结束

## 15.4 快速使用 Kafka 生产消费数据

### 操作场景

用户可以在集群客户端完成Topic的创建、查询、删除等基本操作。可参考[Kafka用户权限说明](#)设置用户权限，然后参考[使用Kafka客户端生产消费数据（MRS 3.x之前版本）](#)进行操作。

MRS 3.1.2及之后版本集群也可以通过登录KafkaUI查看当前集群的消费信息。详细操作请参考[使用KafkaUI查看消费信息（MRS 3.1.2及之后版本）](#)。

### 前提条件

- 使用Kafka客户端时：已安装客户端，例如安装目录为“/opt/client”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- 使用KafkaUI时：已创建具有KafkaUI页面访问权限的用户，如需在页面上进行相关操作，例如创建Topic，需同时授予用户相关权限，请参考[Kafka用户权限说明](#)。  
第一次访问Manager和KafkaUI，需要在浏览器中添加站点信任以继续访问KafkaUI。

### 使用 Kafka 客户端生产消费数据（MRS 3.x 之前版本）

**步骤1** 安装客户端，具体请参考[安装客户端](#)章节。

**步骤2** 进入ZooKeeper实例页面：

单击集群名称，登录集群详情页面，选择“组件管理 > ZooKeeper > 实例”。

#### 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

**步骤3** 查看ZooKeeper角色实例的IP地址。

记录ZooKeeper角色实例其中任意一个的IP地址即可。

**步骤4** 登录安装客户端的节点。

**步骤5** 执行以下命令，切换到客户端目录，例如“/opt/client/Kafka/kafka/bin”。

```
cd /opt/client/Kafka/kafka/bin
```

**步骤6** 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

**步骤7** 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinitKafka用户
```

**步骤8** 创建一个Topic：

```
sh kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份个数 --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

例如：

```
sh kafka-topics.sh --create --topic TopicTest --partitions 3 --replication-factor 3 --zookeeper 10.10.10.100:2181/kafka
```

**步骤9** 执行以下命令，查询集群中的Topic信息：

```
sh kafka-topics.sh --list --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

例如：

```
sh kafka-topics.sh --list --zookeeper 10.10.10.100:2181/kafka
```

**步骤10** 删除**步骤8**中创建的Topic：

```
sh kafka-topics.sh --delete --topic 主题名称 --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

例如：

```
sh kafka-topics.sh --delete --topic TopicTest --zookeeper 10.10.10.100:2181/kafka
```

输入 "y"，回车。

----结束

## 使用 Kafka 客户端生产消费数据（MRS 3.x 及之后版本）

**步骤1** 安装客户端，具体请参考[安装客户端](#)章节。

**步骤2** 进入ZooKeeper实例页面：

登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 服务 > ZooKeeper > 实例”。

**步骤3** 查看ZooKeeper角色实例的IP地址。

记录ZooKeeper角色实例其中任意一个的IP地址即可。

**步骤4** 登录安装客户端的节点。

**步骤5** 执行以下命令，切换到客户端目录，例如“/opt/client/Kafka/kafka/bin”。

```
cd /opt/client/Kafka/kafka/bin
```

**步骤6** 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

**步骤7** 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

### **kinitKafka用户**

**步骤8** 登录FusionInsight Manager，选择“集群 > 待操作的集群名称 > 服务 > ZooKeeper > 配置 > 全部配置”，搜索参数“clientPort”，记录“clientPort”的参数值。

**步骤9** 创建一个Topic：

```
sh kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份个数 --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

例如：`sh kafka-topics.sh --create --topic TopicTest --partitions 3 --replication-factor 3 --zookeeper 10.10.10.100:2181/kafka`

**步骤10** 执行以下命令，查询集群中的Topic信息：

```
sh kafka-topics.sh --list --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

例如：`sh kafka-topics.sh --list --zookeeper 10.10.10.100:2181/kafka`

**步骤11** 删除**步骤9**中创建的Topic：

```
sh kafka-topics.sh --delete --topic 主题名称 --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

例如：`sh kafka-topics.sh --delete --topic TopicTest --zookeeper 10.10.10.100:2181/kafka`

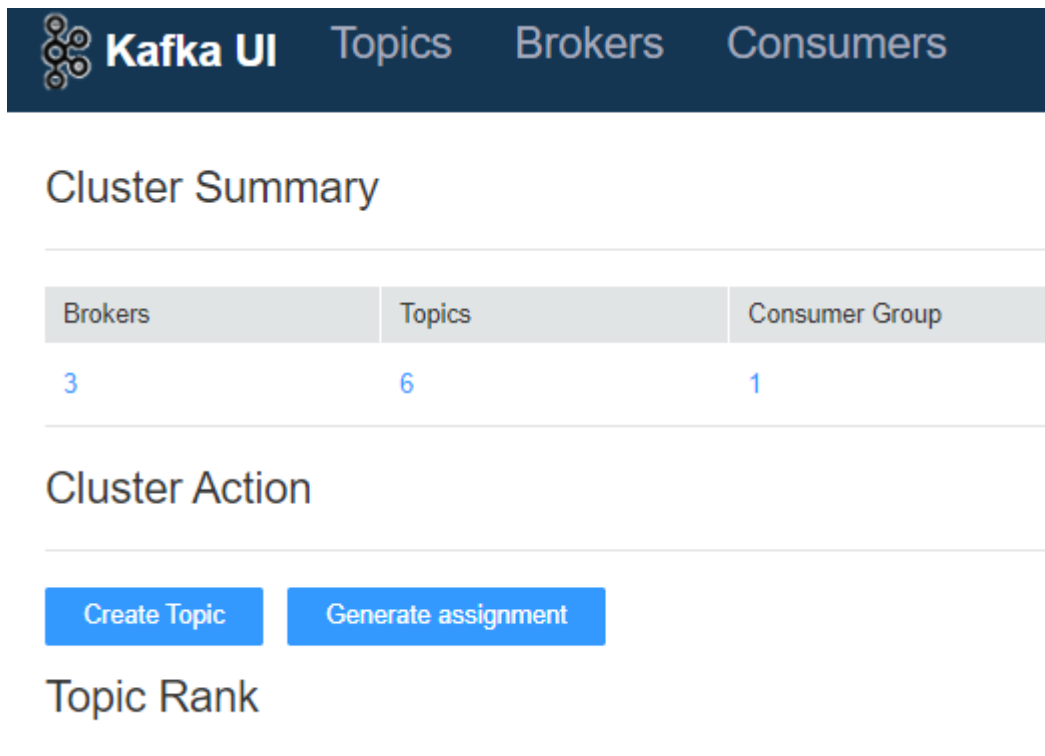
----结束

## 使用 KafkaUI 查看消费信息（MRS 3.1.2 及之后版本）

**步骤1** 进入KafkaUI界面。

1. 使用具有KafkaUI页面访问权限的用户登录FusionInsight Manager，选择“集群 > 服务 > Kafka”。  
如需在页面上进行相关操作，例如创建Topic，需同时授予用户相关权限，请参考[Kafka用户权限说明](#)。
2. 在“KafkaManager WebUI”右侧，单击URL链接，访问KafkaUI的页面。

**步骤2** 在“Cluster Summary”栏，可查看当前集群已有的Topic、Broker和Consumer Group数量。



**步骤3** 单击“Brokers”、“Topics”、“Consumer Group”下方的数字，可自动跳转至对应页面，查看并操作对应信息。

**步骤4** 在“Cluster Action”栏，可创建Topic与分区迁移，具体操作请参考[增加Kafka Topic分区](#)。

**步骤5** 在“Topic Rank”栏，可查看当前集群Topic日志条数、数据体积大小、数据流入量、数据流出量前十名的Topic。

Topic Rank

Topic Logsize Top 10			
RankID	TopicName	Logsize	Default Topic
1	test1	142171956	false
2	__consumer_offsets	15174	true
3	__default_metrics	14148	true
4	__KafkaMetric-Report	3477	true
5	cdi-connect-configs	20	false
6	test2	5	false
7	test	3	false
8	cdi-connect-offsets	0	false
9	cdi-connect-status	0	false
10			

Topic Capacity Top 10			
RankID	TopicName	Capacity	Default Topic
1	test1	15.9GB	false
2	__default_metrics	12.0MB	true
3	__consumer_offsets	2.9MB	true
4	__KafkaMetric-Report	679.5KB	true
5	cdi-connect-configs	3.8KB	false
6	test2	225.0B	false
7	test	147.0B	false
8	cdi-connect-offsets	0.0B	false
9	cdi-connect-status	0.0B	false
10			

**步骤6** 单击“TopicName”可进入到该Topic的详情页面中，在该页面的具体操作请参考[查看Kafka数据生产消费详情](#)。

----结束



## 15.5 创建 Kafka Topic

### 操作场景

用户可以根据业务需要，使用集群客户端创建Kafka的主题。启用Kerberos认证的集群，需要拥有管理Kafka主题的权限。

### 前提条件

已安装客户端。

### 使用 Kafka 客户端创建 Kafka Topic

**步骤1** 进入ZooKeeper实例页面：

- MRS3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > ZooKeeper > 实例”。

#### 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > ZooKeeper > 实例”。

**步骤2** 查看ZooKeeper角色实例的IP地址。

记录ZooKeeper角色实例其中任意一个的IP地址即可。

**步骤3** 根据业务情况，准备好客户端，登录安装客户端的节点。

请根据客户端所在位置，参考[使用MRS客户端](#)章节，登录安装客户端的节点。

**步骤4** 执行以下命令，切换到客户端目录，例如“/opt/client/Kafka/kafka/bin”。

```
cd /opt/client/Kafka/kafka/bin
```

**步骤5** 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

**步骤6** 执行以下命令，进行用户认证。（普通模式跳过此步骤）

```
kinit 组件业务用户
```

**步骤7** MRS 3.x之前版本：执行以下命令，创建Kafka主题。

- 创建主题

```
sh kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份个数 --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

- 删除主题

```
sh kafka-topics.sh --delete --topic 主题名称 --zookeeper ZooKeeper角色实例所在节点IP地址:clientPort/kafka
```

## 📖 说明

- 主题分区数和主题备份个数不能大于Kafka角色实例数量。
- 默认情况下，ZooKeeper的“clientPort”为“2181”。
- ZooKeeper角色实例所在节点IP地址，填写三个角色实例其中任意一个的IP地址即可。
- 使用Kafka主题管理消息，请参见[管理Kafka Topic中的消息](#)。

### 步骤8 MRS 3.x及后续版本：使用kafka-topics.sh创建Kafka主题。

- 创建主题：

Topic的Partition自动划分时，默认根据节点及磁盘上已有的Partition数进行均衡划分，如果期望根据磁盘容量进行Partition划分，那么需要修改Kafka服务配置“log.partition.strategy”为“capacity”。

Kafka创建Topic时，支持基于“机架感知”和“跨AZ特性”两种选项组合生成分区及副本的分配方案且支持“--zookeeper”和“--bootstrap-server”两种方式

- 禁用机架策略 & 禁用跨AZ特性（默认策略）。

基于此策略新建的Topic的副本会完全随机分配到集群中任意节点上。

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --zookeeper ZooKeeper的任意一个节点的业务IP:clientPort/kafka
```

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties
```

其中，使用“--bootstrap-server”方式创建Topic时，需配置“rack.aware.enable=false”和“az.aware.enable=false”。

- 启用机架策略 & 禁用跨AZ特性。

基于此策略新建的Topic的各个Partition的Leader会在集群节点上随机分配，但会确保同一Partition的不同Replica会分配在不同的机架上，所以当使用此策略时，需保证各个机架内的节点个数一致，否则会导致节点少的机架上的机器负载远高于集群平均水平。

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --zookeeper ZooKeeper的任意一个节点的业务IP:clientPort/kafka --enable-rack-aware
```

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties
```

其中，使用“--bootstrap-server”方式创建Topic时，需配置“rack.aware.enable=true”和“az.aware.enable=false”。

- 禁用机架策略 & 启用跨AZ特性。

基于此策略新建的Topic的各个Partition的Leader会在集群节点上随机分配，但会确保同一Partition的不同Replica会分配在不同的AZ上，所以当使用此策略时，需保证各个AZ内的节点个数一致，否则会导致节点少的AZ上的机器负载远高于集群平均水平。

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --zookeeper ZooKeeper的任意一个节点的业务IP:clientPort/kafka --enable-az-aware
```

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties
```

其中，使用“--bootstrap-server”方式创建Topic时，需配置“rack.aware.enable=false”和“az.aware.enable=true”。

- 启用机架策略 & 启用跨AZ特性。

基于此策略新建的Topic的各个Partition的Leader会在集群节点上随机分配，但会确保同一Partition的不同Replica会分配到不同AZ内的不同RACK上，使用此策略需保证每个AZ内的每个RACK上的节点个数一致，否则会导致集群内负载不均衡。

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --zookeeper ZooKeeper的任意一个节点的IP:clientPort/kafka --enable-rack-aware --enable-az-aware
```

```
./kafka-topics.sh --create --topic 主题名称 --partitions 主题占用的分区数 --replication-factor 主题的备份数 --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties
```

使用“--bootstrap-server”方式创建Topic时，需配置“rack.aware.enable=true”和“az.aware.enable=true”。

#### 📖 说明

- Kafka创建Topic支持“--zookeeper”和“--bootstrap-server”两种方式，区别如下：
  - “--zookeeper”方式由客户端生成副本分配方案，社区从一开始就支持这种方式，为了降低对Zookeeper组件的依赖，社区将在后续版本中删除对这种方式的支持。基于这种方式创建Topic时，可以通过“--enable-rack-aware”和“--enable-az-aware”这两个选项自由组合来选用副本分配策略。注意：使用“--enable-az-aware”选项的前提是服务端开启了跨AZ特性，即服务端启动参数“az.aware.enable”为“true”，否则会执行失败。
  - “--bootstrap-server”方式由服务端生成副本分配方案，后续版本，社区将只支持这种方式来进行Topic管理。基于这种方式创建Topic时，不支持“--enable-rack-aware”和“--enable-az-aware”选项来控制副本分配策略，支持“rack.aware.enable”和“az.aware.enable”这两个服务启动参数组合来控制副本分配策略，需注意的是“az.aware.enable”参数不可修改，在创建集群时，如果开启跨AZ特性，会自动配置为“true”；“rack.aware.enable”参数支持用户自定义修改。
- 查看主题：
  - `./kafka-topics.sh --describe --zookeeper ZooKeeper的任意一个节点的IP:clientPort/kafka --topic 主题名称`
  - `./kafka-topics.sh --describe --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties --topic 主题名称`

----结束

## 15.6 在 Kafka Topic 中接入消息

### 操作场景

用户可以根据业务需求，通过Kafka客户端或KafkaUI查看当前消费情况。

本章节内容适用于MRS 3.x及后续版本。

## 前提条件

如果当前使用Kafka客户端，需要满足以下条件：

- MRS集群管理员已明确业务需求，并准备一个系统用户。
- 已安装Kafka客户端。

## 使用 Kafka 客户端查看当前消费情况

**步骤1** 以客户端安装用户，登录安装Kafka客户端的节点。

**步骤2** 切换到Kafka客户端安装目录，例如“/opt/client”。

```
cd /opt/client
```

**步骤3** 执行以下命令，配置环境变量。

```
source bigdata_env
```

**步骤4** 执行以下命令，进行用户认证。（普通模式跳过此步骤）

```
kinit 组件业务用户
```

**步骤5** 执行以下命令，切换到Kafka客户端安装目录。

```
cd Kafka/kafka/bin
```

**步骤6** 使用kafka-consumer-groups.sh查看当前消费情况。

- 查看Offset保存在Kafka上的Consumer Group列表：

```
./kafka-consumer-groups.sh --list --bootstrap-server <Broker的任意一个节点的
业务IP:Kafka集群IP端口号> --command-config ../config/
consumer.properties
```

```
例如：./kafka-consumer-groups.sh --bootstrap-server 192.168.1.1:21007 --
list --command-config ../config/consumer.properties
```

- 查看Offset保存在Kafka上的Consumer Group消费情况：

```
./kafka-consumer-groups.sh --describe --bootstrap-server <Broker的任意一
个节点的
业务IP:Kafka集群IP端口号> --group 消费组名称 --command-config ../
config/consumer.properties
```

```
例如：./kafka-consumer-groups.sh --describe --bootstrap-server
192.168.1.1:21007 --group example-group --command-config ../
consumer.properties
```

---

### 须知

1. 确保当前consumer在线消费。
2. 确保配置文件consumer.properties中的group.id与命令中--group的参数均配置为待查询的group。
3. Kafka集群IP端口号安全模式下是21007，普通模式下是9092。

---

----结束

## 使用 KafkaUI 查看当前消费情况（MRS 3.1.2 及之后版本）

### 步骤1 进入KafkaUI界面。

1. 使用具有KafkaUI页面访问权限的用户登录FusionInsight Manager，选择“集群 > 服务 > Kafka”。

如需在页面上进行相关操作，例如创建Topic，需同时授予用户相关权限，请参考[Kafka用户权限说明](#)。

2. 在“KafkaManager WebUI”右侧，单击URL链接，访问KafkaUI的页面。

### 步骤2 单击“Consumers”，进入消费组详情页面，可以查看当前集群内的所有ConsumerGroups，并可以查看各个ConsumerGroups Coordinator所在节点IP，在页面右上角，用户可以输入ConsumerGroup来搜索指定的ConsumerGroup信息。

Consumer Summary			
Group	Topics	Coordinator	Active Topics
<a href="#">example-group11</a>	2	10.244.228.252	0
<a href="#">example-group4</a>	1	10.244.229.85	0
<a href="#">example-group5</a>	1	10.244.229.170	0
<a href="#">example-group6</a>	1	10.244.229.85	0
<a href="#">example-group7</a>	1	10.244.228.252	0
<a href="#">example-group8</a>	1	10.244.229.170	0
<a href="#">__KafkaMetricReportGroup</a>	1	10.244.228.252	0
<a href="#">example-group9</a>	1	10.244.229.85	0
<a href="#">example-group10</a>	1	10.244.228.89	0
<a href="#">example-group1</a>	1	10.244.229.85	0

### 步骤3 在Consumer Summary一栏，可查看当前集群已存在的消费组，单击消费组名称，可查看该消费组所消费过的Topic，消费过的Topic有两种状态：“pending”和“running”，分别表示“曾经消费过但现在未消费”和“现在正在消费”，在弹框右上角，可以输入Topic名来进行过滤。

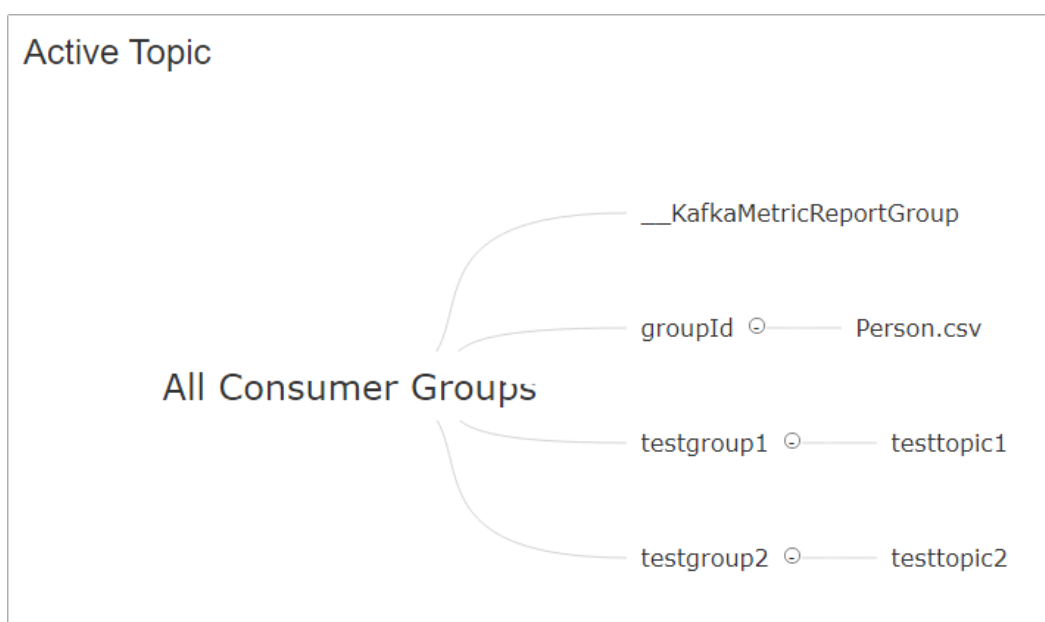
Consumer Topics	
Topic	Consumer Status
<a href="#">123456789012345678901234567890123456789...</a>	pending
<a href="#">test0</a>	pending

**步骤4** 单击Topic名称，进入Consumer Offsets页面，可查看Topic消费详情。

Consumer Offsets					
example-group11 : aaa					
Partition	Log End Offset	Current Offset	Lag	ConsumerID	Host
0	21683	18206	3477	consumer-example-group11-1-7c65fa74-01...	10.244.228.252
1	21498	18155	3343	consumer-example-group11-1-7c65fa74-01...	10.244.228.252

**步骤5** 查看消费关系图。

单击“Consumers”，进入消费组详情页面。在Active Topic 处可以查看当前集群所有的消费组，以及各个Consumer Group正在消费的Topic。



**说明**

MRS集群当前不支持单击消费组名称进行跳转。

----结束

## 15.7 管理 Kafka Topic

### 15.7.1 查看 Kafka Topic 信息

#### 操作场景

用户可以在Manager或KafkaUI上查看Kafka已创建的主题信息。

#### 在 Manager 查看 Kafka Topic 信息

**步骤1** 进入Kafka服务页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Kafka”。

#### 📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务 > Kafka”。

**步骤2** 单击“KafkaTopic监控”。

主题列表默认显示所有主题。可以查看主题的分区数和备份数。

**步骤3** 在主题列表单击指定主题的名称，可查看详细信息。

#### 📖 说明

如果执行过以下几种操作：

- Kafka或者Zookeeper进行过扩容或缩容操作。
- Kafka或者Zookeeper增加或者删除过实例。
- 重装Zookeeper服务。
- Kafka切换到了其他的Zookeeper服务。

可能导致Kafka Topic监控不显示，请按以下步骤恢复：

1. 登录到集群的主OMS节点，执行以下切换到omm用户：

```
su - omm
```

2. 重启cep服务：

```
restart_app cep
```

重启后等待3分钟，再次查看kafka Topic监控。

----结束

## 在 KafkaUI 查看 Kafka Topic 信息（MRS 3.1.2 及之后版本）

**步骤1** 进入KafkaUI界面。

1. 使用具有KafkaUI页面访问权限的用户登录FusionInsight Manager，选择“集群 > 服务 > Kafka”。

如需在页面上进行相关操作，例如创建Topic，需同时授予用户相关权限，请参考[Kafka用户权限说明](#)。

2. 在“KafkaManager WebUI”右侧，单击URL链接，访问KafkaUI的页面。

**步骤2** 单击“Brokers”，进入Broker详情页面。

**步骤3** 在“Broker Summary”一栏可查看Broker的“Broker ID”、“Host”、“Rack”、“Disk(Used|Total)”和“Memory(Used|Total)”。

Broker Summary				
Broker ID	Host	Rack	Disk(Used Total)	Memory(Used Total)
1	10.112.17.150	/default/rack0	40.2MB   9.1GB	4.4G   6G
2	10.112.17.189	/default/rack0	40.2MB   9.1GB	4.4G   6G
3	10.112.17.228	/default/rack0	41.3MB   9.1GB	4.4G   6G

**步骤4** 在“Brokers Metrics”处可查看Broker节点数据流量的jmx指标，包括在不同时段的时间窗口内，Broker节点平均每秒流入消息条数，每秒流入消息字节数，每秒流出消息字节数，每秒失败的请求数，每秒总的请求数和每秒生产的请求数。

Brokers Metrics ©

Window	Message in /sec	Bytes in /sec	Bytes out /sec	Failed fetch request /sec	Total fetch request /sec	Total produce request /sec
1 min	6067	6639249	10	0	106415	1339
5 min	16769	1855373	10	0	30536	372
15 min	5937	658534	136	0	11611	132
All time	1850	224273	170077	0	17220	122

**步骤5** 在页面右上角，用户可以输入主机IP地址或者机架配置信息搜索查看该Broker信息。

----结束

## 15.7.2 修改 Kafka Topic 配置

### 操作场景

用户可以根据业务需要，使用集群客户端创建Kafka Topic。启用Kerberos认证的集群，需要拥有管理Kafka主题的权限。也可以通过KafkaUI修改Topic Configs。

#### 📖 说明

- 安全模式下，KafkaUI对修改Topic Configs场景，需保证KafkaUI登录用户属于“kafkaadmin”用户组或者单独给用户授予对应操作权限，否则将会鉴权失败。
- 非安全模式下，KafkaUI对所有操作不作鉴权处理。
- 该章节仅适用MRS 3.x及之后版本。

### 使用 Kafka 客户端修改 Kafka Topic

**步骤1** 进入ZooKeeper实例页面：

登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 服务 > ZooKeeper > 实例”。

**步骤2** 查看ZooKeeper角色实例的IP地址。

记录ZooKeeper角色实例其中任意一个的IP地址即可。

**步骤3** 根据业务情况，准备好客户端，登录安装客户端的节点。

请根据客户端所在位置，参考[使用MRS客户端](#)章节，登录安装客户端的节点。

**步骤4** 执行以下命令，切换到客户端目录，例如“/opt/client/Kafka/kafka/bin”。

```
cd /opt/client/Kafka/kafka/bin
```

**步骤5** 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

**步骤6** 执行以下命令，进行用户认证。（普通模式跳过此步骤）

```
kinit 组件业务用户
```

**步骤7** 使用kafka-topics.sh修改主题：

```
./kafka-topics.sh --alter --topic 主题名称 --config 配置项=配置值 --zookeeper ZooKeeper的任意一个节点的业务IP:clientPort/kafka
```



**步骤8** 使用kafka-topics.sh查看修改后主题：

- `./kafka-topics.sh --describe --zookeeper ZooKeeper的任意一个节点的业务IP:clientPort/kafka --topic 主题名称`
- `./kafka-topics.sh --describe --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties --topic 主题名称`

----结束

## 使用 KafkaUI 修改 Kafka Topic（MRS 3.1.2 及之后版本）

**步骤1** 进入KafkaUI界面。

1. 使用具有KafkaUI页面访问权限的用户登录FusionInsight Manager，选择“集群 > 服务 > Kafka”。  
如需在页面上进行相关操作，例如创建Topic，需同时授予用户相关权限，请参考[Kafka用户权限说明](#)。
2. 在“KafkaManager WebUI”右侧，单击URL链接，访问KafkaUI的页面。

**步骤2** 单击“Topics”，进入Topic管理页面。

**步骤3** 在待修改项的“Operation”列单击“Action > Config”，弹出的页面中可修改Topic的“Key”和“Value”值，如需要添加多条，可单击+添加。

**步骤4** 单击“OK”完成修改。

----结束

## 15.7.3 增加 Kafka Topic 分区

### 操作场景

用户可以通过KafkaUI增加Kafka Topic分区。

#### 说明

- 安全模式集群下，执行分区迁移操作的用户需属于“kafkaadmin”用户组，否则将会由于鉴权失败导致操作失败。
- 非安全模式下，KafkaUI对任意操作不作鉴权处理。
- 本章节内容仅适用于MRS 3.1.2及之后版本。

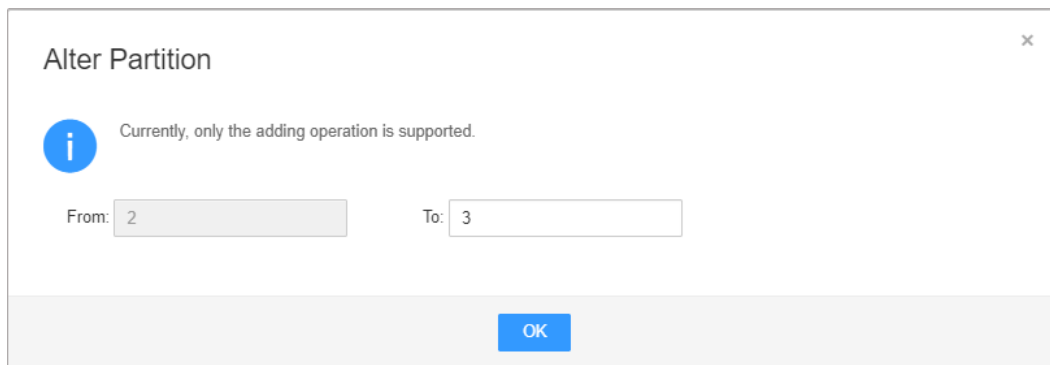
### 增加分区

**步骤1** 进入KafkaUI界面。

1. 使用具有KafkaUI页面访问权限的用户登录FusionInsight Manager，选择“集群 > 服务 > Kafka”。  
如需在页面上进行相关操作，例如创建Topic，需同时授予用户相关权限，请参考[Kafka用户权限说明](#)。
2. 在“KafkaManager WebUI”右侧，单击URL链接，访问KafkaUI的页面。

**步骤2** 单击“Topics”，进入Topic管理页面。

**步骤3** 在待修改项的“Operation”列单击“Action > Alter”，弹出的页面中修改Topic分区。



#### 📖 说明

目前集群只支持增加分区操作，即修改的分区个数要大于原设置的分区个数。

**步骤4** 单击“OK”完成修改。

----结束

## 15.7.4 管理 Kafka Topic 中的消息

### 操作场景

用户可以根据业务需要，使用MRS集群客户端，在Kafka主题中产生消息，或消费消息。

### 前提条件

- 已安装集群客户端。
- 启用Kerberos认证的集群，需要提前在Manager中创建业务用户，用户拥有在Kafka主题中执行相应操作的权限。

### 操作步骤

**步骤1** 进入Kafka服务页面：

- MRS3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理 > Kafka”。

#### 📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，然后选择“集群 > 待操作的集群名称 > 服务 > Kafka”。

**步骤2** 单击“实例”，查看Kafka Broker角色实例的IP地址。

记录Kafka角色实例其中任意一个的IP地址即可。

**步骤3** 根据业务情况，准备好客户端，登录安装客户端的节点。

请根据客户端所在位置，参考[使用MRS客户端](#)章节，登录安装客户端的节点。

**步骤4** 执行以下命令，切换到客户端目录，例如“/opt/client/Kafka/kafka/bin”。

```
cd /opt/client/Kafka/kafka/bin
```

**步骤5** 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
```

**步骤6** 启用Kerberos认证的集群，执行以下命令认证用户身份。未启用Kerberos认证的集群无需执行本步骤。

```
kinit Kafka用户
```

**步骤7** 根据业务需要，管理Kafka主题中的消息。

- 在主题中产生消息

```
sh kafka-console-producer.sh --broker-list Broker角色实例所在节点的IP地址:9092 --topic Topic名称 --producer.config /opt/client/Kafka/kafka/config/producer.properties
```

Topic需提前创建，用户可以输入指定的内容作为生产者产生的消息，输入完成后按回车发送消息。如果需要结束产生消息，使用“Ctrl + C”退出任务。

- 消费主题中的消息

重新打开一个客户端连接，执行以下命令消费主题中的消息。

```
cd /opt/client/Kafka/kafka/bin
```

```
source /opt/client/bigdata_env
```

```
sh kafka-console-consumer.sh --topic Topic名称 --bootstrap-server Broker角色实例所在节点的IP地址:9092 --consumer.config /opt/client/Kafka/kafka/config/consumer.properties
```

配置文件中“group.id”指定的消费者组默认为“example-group1”。用户可根据业务需要，自定义其他消费者组。每次消费时生效。

执行命令时默认会读取当前消费者组中未被处理的消息。如果在配置文件指定了新的消费者组且命令中增加参数“--from-beginning”，则会读取所有Kafka中未被自动删除的消息。

#### 说明

- Kafka角色实例所在节点IP地址，填写Broker角色实例其中任意一个的IP地址即可。
- 如果集群启用Kerberos认证，则端口需要修改为“21007”。
- 默认情况下，ZooKeeper的“clientPort”为“2181”。

----结束

## 15.7.5 查看 Kafka 数据生产消费详情

### 操作场景

用户可以通过KafkaUI查看Topic详情、修改Topic Configs、增加Topic分区个数、删除Topic，并可实时查看不同时段的生产数据条数。

## 说明

- 安全模式下，KafkaUI对查看Topic详情操作不作鉴权处理，即任何用户都可以查询Topic信息；对于修改Topic Configs、增加Topic分区个数、删除Topic场景，需保证KafkaUI登录用户属于“kafkaadmin”用户组或者单独给用户授予对应操作权限，否则将会鉴权失败。
- 非安全模式下，KafkaUI对所有操作不作鉴权处理。
- 本章节内容仅适用于MRS 3.1.2及之后版本。

## 查看生产消费详情

### 步骤1 进入KafkaUI界面。

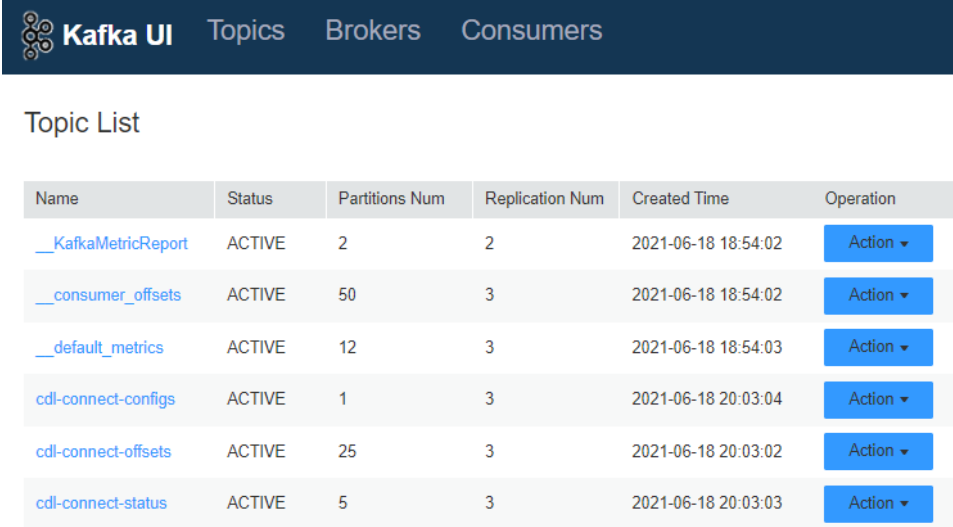
- 使用具有KafkaUI页面访问权限的用户登录FusionInsight Manager，选择“集群 > 服务 > Kafka”。

如需在页面上进行相关操作，例如创建Topic，需同时授予用户相关权限，请参考[Kafka用户权限说明](#)。

- 在“KafkaManager WebUI”右侧，单击URL链接，访问KafkaUI的页面。

### 步骤2 单击“Topics”，进入Topic管理页面。可在当前页面进行如下操作：

- 在“Topic List”栏可查看当前集群已创建的Topic的名称、状态、分区数量、创建时间和副本个数等信息。



The screenshot shows the Kafka UI interface with a navigation bar containing 'Kafka UI', 'Topics', 'Brokers', and 'Consumers'. Below the navigation bar is the 'Topic List' section, which contains a table with the following data:

Name	Status	Partitions Num	Replication Num	Created Time	Operation
<a href="#">_KafkaMetricReport</a>	ACTIVE	2	2	2021-06-18 18:54:02	Action ▾
<a href="#">__consumer_offsets</a>	ACTIVE	50	3	2021-06-18 18:54:02	Action ▾
<a href="#">__default_metrics</a>	ACTIVE	12	3	2021-06-18 18:54:03	Action ▾
<a href="#">cdl-connect-configs</a>	ACTIVE	1	3	2021-06-18 20:03:04	Action ▾
<a href="#">cdl-connect-offsets</a>	ACTIVE	25	3	2021-06-18 20:03:02	Action ▾
<a href="#">cdl-connect-status</a>	ACTIVE	5	3	2021-06-18 20:03:03	Action ▾

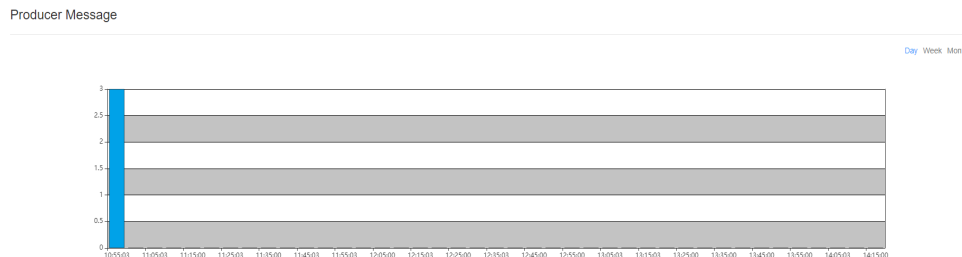
### Producer Message

- 单击Topic名称可进入Topic详情页面。在该页面可查看Topic与分区的详细信息。

Partition Summary

Partition Id	Leader	Replicas	In Sync Replicas	Logsize (B)	Start Offset	End Offset
0	1	[1, 2, 3]	[1, 2, 3]	0.0B	0	0
1	2	[2, 3, 1]	[2, 3, 1]	0.0B	0	0
2	3	[3, 1, 2]	[3, 1, 2]	0.0B	0	0
3	1	[1, 3, 2]	[1, 3, 2]	0.0B	0	0
4	2	[2, 1, 3]	[2, 1, 3]	0.0B	0	0
5	3	[3, 2, 1]	[3, 2, 1]	3.0MB	0	14583
6	1	[1, 2, 3]	[1, 2, 3]	0.0B	0	0
7	2	[2, 3, 1]	[2, 3, 1]	0.0B	0	0
8	3	[3, 1, 2]	[3, 1, 2]	0.0B	0	0
9	1	[1, 3, 2]	[1, 3, 2]	0.0B	0	0

- 在“Producer Message”栏可根据业务需求选择“Day”、“Week”、“Month”不同时段查看此Topic生产数据条数。

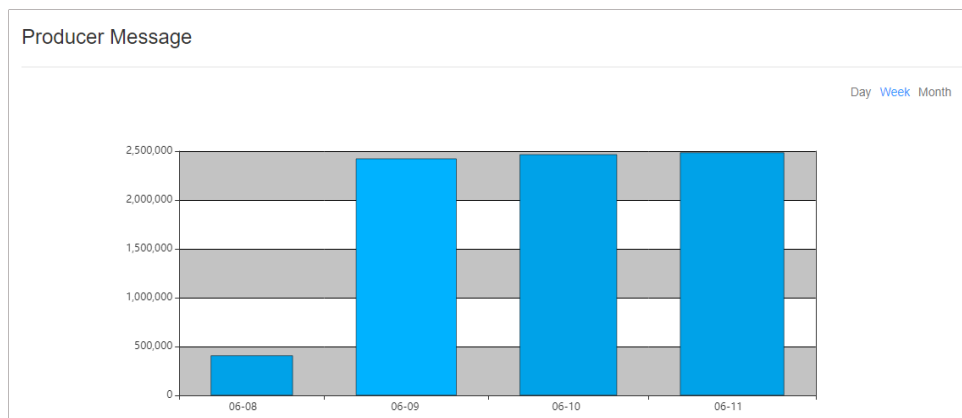


- 修改Topic配置。  
在待修改项的“Operation”列单击“Action > Config”，弹出的页面中可修改Topic的“Key”和“Value”值，如需要添加多条，可单击+添加。单击“OK”完成修改。
- 搜索Topic。  
在页面右上角，用户可以输入Topic名称搜索查看该Topic信息。
- 删除Topic。  
在待修改项的“Operation”列单击“Action > Delete”。在弹出的确认信息页面中单击“OK”即可完成删除。

#### 📖 说明

系统默认内置的Topic不支持删除操作。

- 查看生产数据条数。  
在“Producer Message”栏可选择“Day”、“Week”、“Month”不同时段查看当前集群所有集群生产数据条数。



----结束

## 15.8 Kafka 企业级能力增强

## 15.8.1 配置 Kafka 高可用和高可靠

### 操作场景

Kafka消息传输保障机制，可以通过配置不同的参数来保障消息传输，进而满足不同的性能和可靠性要求。本章节介绍如何配置Kafka高可用和高可靠参数。

本章节内容适用于MRS 3.x及后续版本。

### 对系统的影响

- 配置高可用、高性能的影响：

#### 须知

配置高可用、高性能模式后，数据可靠性会降低。在磁盘故障、节点故障等场景下存在数据丢失风险。

- 配置高可靠性的影响：
  - 性能降低：

在生产数据时，配置了高可靠参数ack=-1之后，需要多个副本均写入成功之后才认为是写入成功。这样会导致单条消息时延增加，客户端处理能力下降。具体性能以现场实际测试数据为准。
  - 可用性降低：

不允许不在ISR中的副本被选举为Leader。如果Leader下线时，其他副本均不在ISR列表中，那么该分区将保持不可用，直到Leader节点恢复。当分区的一个副本所在节点故障时，无法满足最小写入成功的副本数，那么将会导致业务写入失败。
- 参数配置项为服务级配置需要重启Kafka，建议在变更窗口做服务级配置修改。

### 参数描述

- 如果业务需要保证高可用和高性能。

在服务端配置如表15-3中参数，参数配置入口请参考[修改集群服务配置参数](#)。

表 15-3 服务端高可用性和高性能参数说明

参数	默认值	说明
unclean.leader.election.enable	true	是否允许不在ISR中的副本被选举为Leader，若设置为true，可能会造成数据丢失。
auto.leader.rebalance.enable	true	是否使用Leader自动均衡功能。 如果设为true，Controller会周期性的为所有节点的每个分区均衡Leader，将Leader分配给更优先的副本。
min.insync.replicas	1	当Producer设置acks为-1时，指定需要写入成功的副本的最小数目。

在客户端配置文件producer.properties中配置如表15-4中参数，producer.properties存放路径为：/opt/client/Kafka/kafka/config/producer.properties，其中/opt/client为Kafka客户端安装目录。

表 15-4 客户端高可用性和高性能参数说明

参数	默认值	说明
acks	1	<p>需要Leader确认消息是否已经接收并认为已经处理完成。该参数会影响消息的可靠性和性能。</p> <ul style="list-style-type: none"> <li>acks=0：Producer将不会等待服务端任何响应。消息将会被认为成功。</li> <li>acks=1：当副本所在Leader确认数据已写入，但是其不会等待所有的副本完全写入即返回响应。在这种情况下，如果Leader确认后但是副本未同步完成时Leader异常，那么数据就会丢失。</li> <li>acks=-1：意味着等待所有的同步副本确认后才认为成功，配合“min.insync.replicas”可以确保多副本写入成功，只要有一个副本保持活跃状态，记录将不会丢失。</li> </ul>

- 如果业务需要保证数据高可靠性。  
在服务端配置如表15-5参数，参数配置入口请参考[修改集群服务配置参数](#)。

表 15-5 服务端高可靠性参数说明

参数	建议值	说明
unclean.leader.election.enable	false	不允许不在ISR中的副本被选举为Leader。

参数	建议值	说明
min.insync.replicas	2	当Producer设置acks为-1时，指定需要写入成功的副本的最小数目。 需要满足min.insync.replicas <= replication.factor。

在客户端配置文件producer.properties中配置如表15-6中参数，producer.properties存放路径为：/opt/client/Kafka/kafka/config/producer.properties，其中/opt/client为Kafka客户端安装目录。

表 15-6 客户端高可靠性参数说明

参数	建议值	说明
acks	-1	Producer需要Leader确认消息是否已经接收并认为已经处理完成。 acks=-1需要等待在ISR列表的副本都确认接收到消息并处理完成才表示消息成功。配合“min.insync.replicas”可以确保多副本写入成功，只要有一个副本保持活跃状态，记录将不会丢失，此参数配置为-1时，会降低生产性能，请权衡后配置。

## 配置建议

请根据以下业务场景对可靠性和性能要求进行评估，采用合理参数配置。

- 对于价值数据，这两种场景下建议Kafka数据目录磁盘配置raid1或者raid5，从而提高单个磁盘故障情况下数据可靠性。
- 参数配置项均为Topic级别可修改的参数，默认采用服务级配置。

可针对不同Topic可靠性要求对Topic进行单独配置。以root用户登录Kafka客户端节点，在客户端安装目录下配置Topic名称为test的可靠性参数命令：

```
cd Kafka/kafka/bin
```

```
kafka-topics.sh --zookeeper 192.168.1.205:2181/kafka --alter --topic test --config unclean.leader.election.enable=false --config min.insync.replicas=2
```

其中192.168.1.205为ZooKeeper业务IP地址。

- 参数配置项为服务级配置需要重启Kafka，建议在变更窗口做服务级配置修改。



## 15.8.2 配置 Kafka 数据安全传输协议

本章节内容适用于MRS 3.x及后续版本。

### Kafka API 简单说明

- Producer API  
指`org.apache.kafka.clients.producer.KafkaProducer`中定义的接口，在使用“`kafka-console-producer.sh`”时，默认使用此API。
- Consumer API  
指`org.apache.kafka.clients.consumer.KafkaConsumer`中定义的接口，在使用“`kafka-console-consumer.sh`”时，默认会调用此API。

#### 📖 说明

MRS 3.x后，Kafka不支持旧Producer API和旧Consumer API。

### Kafka 访问协议说明

请参考[修改集群服务配置参数](#)查看或配置参数。

Kafka当前支持四种协议类型的访问：PLAINTEXT、SSL、SASL\_PLAINTEXT、SASL\_SSL。

Kafka服务启动时，默认会启动PLAINTEXT和SASL\_PLAINTEXT两种协议类型的安全认证。可通过设置Kafka服务配置“`ssl.mode.enable`”为“`true`”，来启动SSL和SASL\_SSL两种协议类型的安全认证。下表是四种协议类型的简单说明：

协议类型	说明	默认端口
PLAINTEXT	支持无认证的明文访问	获取参数“ <code>port</code> ”的值，默认为9092
SASL_PLAINTEXT	支持Kerberos认证的明文访问	获取参数“ <code>sasl.port</code> ”的值，默认为21007
SSL	支持无认证的SSL加密访问	获取参数“ <code>ssl.port</code> ”的值，默认为9093
SASL_SSL	支持Kerberos认证的SSL加密访问	获取参数“ <code>sasl-ssl.port</code> ”的值，默认为21009

## Topic 的 ACL 设置

Topic的权限信息，需要在Linux客户端上，使用“kafka-acls.sh”脚本进行查看和设置，具体可参考[Kafka用户权限说明](#)。

### 针对不同的 Topic 访问场景，Kafka 中 API 使用说明

- 场景一：访问设置了ACL的Topic

使用的API	用户属组	客户端参数	服务端参数	访问的端口
API	用户需满足以下条件之一即可： <ul style="list-style-type: none"> <li>加入 System_administrator角色</li> <li>属于 kafkaadmin组</li> <li>属于 kafkasuperuser组</li> <li>被授权的 kafka组的用户</li> </ul>	security.inter.broker.protocol=SASL_PLAINTEXT sasl.kerberos.service.name = kafka	-	sasl.port（默认21007）
		security.protocol=SASL_SSL sasl.kerberos.service.name = kafka	“ssl.mode.enabled”配置为true	sasl-ssl.port（默认21009）

- 场景二：访问未设置ACL的Topic

使用的API	用户属组	客户端参数	服务端参数	访问的端口
API	用户需满足以下条件之一： <ul style="list-style-type: none"> <li>加入 System_administrator角色</li> <li>属于 kafkaadmin组</li> <li>属于 kafkasuperuser组</li> </ul>	security.protocol=SASL_PLAINTEXT sasl.kerberos.service.name = kafka	-	sasl.port（默认21007）

使用的 API	用户属组	客户端参数	服务端参数	访问的端口
	用户属于kafka组		“allow.everyone.if.no.acl.found”配置为 true <b>说明</b> 普通集群下不涉及服务端参数“allow.everyone.if.no.acl.found”的修改	sasl.port （默认 21007）
	用户需满足以下条件之一： <ul style="list-style-type: none"> <li>加入 System_administrator角色</li> <li>属于kafkaadmin组</li> <li>kafkasuperuser组用户</li> </ul>	security.protocol=SASL_SSL sasl.kerberos.service.name = kafka	“ssl.mode.enable”配置为“true”	sasl-ssl.port （默认 21009）
	用户属于kafka组		1. “allow.everyone.if.no.acl.found”配置为“true” 2. “ssl.mode.enable”配置为“true”	sasl-ssl.port （默认 21009）
	-	security.protocol=PLAINTEXT	“allow.everyone.if.no.acl.found”配置为“true”	port（默认 9092）
	-	security.protocol=SSL	1. “allow.everyone.if.no.acl.found”配置为“true” 2. “ssl.mode.enable”配置为“true”	ssl.port（默认9063）

## 15.8.3 配置 Kafka 数据均衡工具

### 操作场景

该任务指导管理员根据业务需求，在客户端中执行Kafka均衡工具来均衡Kafka集群的负载，一般用于节点的退服、入服以及负载均衡的场景。

本章节内容适用于MRS 3.x及后续版本。3.x之前版本请参考[均衡Kafka扩容节点后数据](#)

### 前提条件

- MRS集群管理员已明确业务需求，并准备一个Kafka管理员用户（属于kafkaadmin组，普通模式不需要）。
- 已安装Kafka客户端。

### 操作步骤

**步骤1** 以客户端安装用户，登录已安装Kafka客户端的节点。

**步骤2** 切换到Kafka客户端安装目录，例如“/opt/client”。

```
cd /opt/client
```

**步骤3** 执行以下命令，配置环境变量。

```
source bigdata_env
```

**步骤4** 执行以下命令，进行用户认证（普通模式跳过此步骤）。

```
kinit 组件业务用户
```

**步骤5** 执行以下命令，切换到Kafka客户端安装目录。

```
cd Kafka/kafka
```

**步骤6** 使用“kafka-balancer.sh”进行用户集群均衡，常用命令如下：

- 使用--run命令执行集群均衡：

```
./bin/kafka-balancer.sh --run --zookeeper <ZooKeeper的任意一个节点的业务IP:zkPort/kafka>--bootstrap-server <Kafka集群IP: port> --throttle 10000000 --consumer-config config/consumer.properties --enable-az-aware --show-details
```

该命令包含均衡方案的生成和执行两部分，其中--show-details为可选参数，表示是否打印方案明细，--throttle表示均衡方案执行时的带宽限制，单位:bytes/sec，--enable-az-aware为可选参数，表明生成均衡方案时，开启跨AZ特性，使用此参数时，请务必保证集群已开启跨AZ特性。

- 使用--run命令执行节点退服：

```
./bin/kafka-balancer.sh --run --zookeeper <ZooKeeper的任意一个节点的业务IP:zkPort/kafka>--bootstrap-server <Kafka集群IP: port>--throttle 10000000 --consumer-config config/consumer.properties --remove-brokers<BrokerId列表> --enable-az-aware --force
```

其中--remove-brokers表示要删除的BrokerId列表，多个间用逗号分隔，--force参数为可选参数，表示忽略磁盘使用率告警，强制生成迁移方案，-enable-az-aware为可选参数，表明生成均衡方案时，开启跨AZ特性，使用此参数时，请务必保证集群已开启跨AZ特性。

### 📖 说明

此退服命令会将待退服Broker节点上的数据迁移至其他Broker节点。

- 查看执行状态：

```
./bin/kafka-balancer.sh --status --zookeeper <ZooKeeper的任意一个节点的业务IP:zkPort/kafka>
```

- 生成均衡方案：

```
./bin/kafka-balancer.sh --generate --zookeeper <ZooKeeper的任意一个节点的业务IP:zkPort/kafka>--bootstrap-server<Kafka集群IP:port> --consumer-config config/consumer.properties --enable-az-aware
```

该命令仅根据集群当前状态生成迁移方案，并打印到控制台，其中--enable-az-aware为可选参数，表明生成迁移方案时，开启跨AZ特性，使用此参数时，请务必保证集群已开启跨AZ特性。

- 清理中间状态

```
./bin/kafka-balancer.sh --clean --zookeeper <ZooKeeper的任意一个节点的业务IP:zkPort/kafka>
```

一般在迁移没有正常执行完成时用来清理ZooKeeper上的中间状态信息。

### 须知

Kafka集群IP端口号安全模式下是21007，普通模式下是9092。

---结束

## 异常情况处理

在使用Kafka均衡工具进行Partition迁移的过程中，如果出现集群中Broker故障导致均衡工具的执行进度阻塞，这时需要人工介入来恢复，分为以下几种场景：

- 存在Broker因为磁盘占有率达到100%导致Broker故障的情况。
  - a. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Kafka > 实例”，将运行状态为“正在恢复”的Broker实例停止并记录实例所在节点的管理IP地址以及对应的“broker.id”，该值可通过单击角色名称，在“实例配置”页面中选择“全部配置”，搜索“broker.id”参数获取。
  - b. 以root用户登录记录的管理IP地址，并执行df -lh命令，查看磁盘占用率为100%的挂载目录，例如“\${BIGDATA\_DATA\_HOME}/kafka/data1”。
  - c. 进入该目录，执行du -sh \*命令，查看该目录下各文件夹的大小。查看是否存在除“kafka-logs”目录外的其他文件，并判断是否可以删除或者迁移。
    - 是，删除或者迁移相关数据，然后执行8。
    - 否，执行4。
  - d. 进入“kafka-logs”目录，执行du -sh \*命令，选择一个待移动的Partition文件夹，其名称命名规则为“Topic名称-Partition标识”，记录Topic及Partition。
  - e. 修改“kafka-logs”目录下的“recovery-point-offset-checkpoint”和“replication-offset-checkpoint”文件（两个文件做同样的修改）。

- i. 减少文件中第二行的数字（若移出多个目录，则减少的数字为移出的目录个数）。
- ii. 删除待移出的Partition所在的行（行结构为“Topic名称 Partition标识 Offset”，删除前先将该行数据保存，后续此内容还要添加到目的目录下的同名文件中）。
- f. 修改目的数据目录下（例如：“\${BIGDATA\_DATA\_HOME}/kafka/data2/kafka-logs”）的“recovery-point-offset-checkpoint”和“replication-offset-checkpoint”文件（两个文件做同样的修改）。
  - 增加文件中第二行的数字（若移入多个Partition目录，则增加的数字为移入的Partition目录个数）。
  - 添加待移入的Partition行到文件末尾（行结构为“Topic名称 Partition标识 Offset”，直接复制5中保存的行数据即可）。
- g. 移动数据，将待移动的Partition文件夹移动到目的目录下，移动完成后执行 **chown omm:wheel -R Partition目录** 命令修改Partition目录属组。
- h. 登录FusionInsight Manager，选择“集群 > 服务 > Kafka > 实例”，启动停止的Broker实例。
- i. 等待5至10分钟后查看Broker实例的运行状态是否为“良好”。
  - 是，修复完成后按照“ALM-38001 Kafka磁盘容量不足”告警指导彻底解决磁盘容量不足问题。
  - 否，联系运维人员。

按照上述步骤将故障Broker进行恢复后，阻塞的均衡任务会继续执行，可使用--status命令来查看任务的执行进度。

- 存在由其他原因导致的Broker故障，且问题场景单一明确，短时间内可以恢复Broker的情况。
  - a. 根据问题根因指定恢复方案，恢复故障Broker。
  - b. 故障Broker恢复后，阻塞的均衡任务会继续执行，可使用--status命令来查看任务的执行进度。
- 存在由其他原因导致的Broker故障，且问题场景复杂，短时间内无法恢复Broker的情况。
  - a. 执行 **kinit Kafka管理员用户**。（普通模式跳过此步骤）
  - b. 使用 **zkCli.sh -server <ZooKeeper集群业务IP:zkPort/kafka>** 登录ZooKeeper Shell。
  - c. 执行 **addauth krbgroup**。（普通模式跳过此步骤）
  - d. 删除“/admin/reassign\_partitions”目录和“/controller”目录。
  - e. 通过以上步骤强行终止迁移，待集群恢复后使用 **kafka-reassign-partitions.sh** 命令手动将中间过程中导致的多余的副本删除。

## 15.9 Kafka 性能调优

### 操作场景

通过调整Kafka服务端参数，可以提升特定业务场景下Kafka的处理能力。

## 参数调优

修改服务配置参数，请参考[修改集群服务配置参数](#)。调优参数请参考[表15-7](#)。

表 15-7 调优参数

配置参数	缺省值	调优场景
num.recovery.threads.per.data.dir	10	在Kafka启动过程中，数据量较大情况下，可调大此参数，可以提升启动速度。
background.threads	10	Broker后台任务处理的线程数目。数据量较大的情况下，可适当调大此参数，以提升Broker处理能力。
num.replica.fetchers	1	副本向Leader请求同步数据的线程数，增大这个数值会增加副本的I/O并发度。
num.io.threads	8	Broker用来处理磁盘I/O的线程数目，这个线程数目建议至少等于硬盘的个数。
KAFKA_HEAP_OPTS	-Xmx6G -Xms6G	Kafka JVM堆内存设置。当Broker上数据量较大时，应适当调整堆内存大小。

## 15.10 Kafka 运维管理

### 15.10.1 Kafka 常用配置参数

本章节内容适用于MRS 3.x及后续版本。

#### 参数入口

请参考[修改集群服务配置参数](#)进入Kafka服务参数“全部配置”页面。

#### 常用参数

表 15-8 参数说明

配置参数	说明	缺省值
log.dirs	Kafka数据存储目录列表，以逗号分隔多个目录。	% {@auto.detect.datapart.bk.log.logs}
KAFKA_HEAP_OPTS	Kafka启动Broker时使用的jvm选项。建议根据业务需要进行设置。	-Xmx6G -Xms6G

配置参数	说明	缺省值
auto.create.topics.enable	是否自动创建Topic，若参数设置为false，发消息前需要通过命令创建Topic。	true
default.replication.factor	自动创建Topic时的默认副本数。	2
monitor.preInitDelay	服务启动后，第一次健康检查的延迟时间。如果启动需要较长时间，可以通过调大参数，来完成启动。单位为毫秒。	600000

## 超时参数

表 15-9 Broker 相关超时参数

参数名称	参数说明	默认值	影响分析
controller.socket.timeout.ms	Controller连接Broker的超时时间。单位：毫秒。	30000	Controller连接Broker的超时时间，一般不需要调整。
group.max.session.timeout.ms	Consumer注册时允许的最大会话超时时间。单位：毫秒。	180000	允许Consumer配置的session.timeout.ms的最大值（不包含此值）。
group.min.session.timeout.ms	Consumer注册时允许的最小会话超时时间。单位：毫秒。	6000	允许Consumer配置的session.timeout.ms的最小值（不包含此值）。
offsets.commit.timeout.ms	Offset提交请求的超时时间。单位：毫秒。	5000	Offset提交时被延迟处理的最大超时时间。
replica.socket.timeout.ms	副本数据同步请求的超时时间，配置值不得小于replica.fetch.wait.max.ms。单位：毫秒。	30000	同步线程在发送同步请求之前等待通道建立的最大超时时间，要求配置大于replica.fetch.wait.max.ms。



参数名称	参数说明	默认值	影响分析
request.timeout.ms	设置客户端发送连接请求后，等待响应的超时时间。单位：毫秒。	30000	Broker节点上的Controller、Replica线程中传入networkclient连接的超时参数，如果在超时时间内没有接收到响应，那么客户端重新发送，并在达到重试次数后返回请求失败。
transaction.max.timeout.ms	事务允许的最大超时。单位：毫秒。	900000	事务最大超时时间，如果客户端的请求时间超过该值，则Broker将在InitProducerIdRequest中返回一个错误。这样可以防止客户端超时时间过长，而导致消费者无法接收topic。
user.group.cache.timeout.seconds	指定缓存中保存用户对应组信息的时间。单位：秒。	300	缓存中用户和组对应关系缓存时间，超过此时间用户信息才会再次通过id -Gn命令查询，在此期间，仅使用缓存中的用户和组对应关系。
zookeeper.connection.timeout.ms	连接ZooKeeper的超时时间。单位：毫秒。	45000	ZooKeeper连接超时时间，这个时间决定了zkclient中初次连接建立过程时允许消耗的时间，超过该时间，zkclient会主动断开。

参数名称	参数说明	默认值	影响分析
zookeeper.session.timeout.ms	ZooKeeper会话超时时间。如果Broker在此时间内未向ZooKeeper上报心跳，则被认为失效。单位：毫秒。	45000	<p>ZooKeeper会话超时时间。</p> <p>作用一：这个时间结合传入的ZKURL中ZooKeeper的地址个数，ZooKeeper客户端以（sessionTimeout/传入ZooKeeper地址个数）为连接一个节点的超时时间，超过此时间未连接成功，则尝试连接下一个节点。</p> <p>作用二：连接建立后，一个会话的超时时间，如ZooKeeper上注册的临时节点BrokerId，当Broker被停止，则该BrokerId，会经过一个sessionTimeout才会被ZooKeeper清理。</p>

表 15-10 Producer 相关超时参数

配置名称	说明	默认值	影响分析
request.timeout.ms	指定发送消息请求的请求超时时间。单位：毫秒。	30000	请求超时时间，出现网络问题时，需调大此参数；配置过小，则容易出现Batch Expire异常。

表 15-11 Consumer 相关超时参数

配置名称	说明	默认值	影响分析
connections.max.idle.ms	空闲连接的保留时间。单位：毫秒	60000	空闲连接的保留时间，连接空闲时间大于此时间，则会销毁该连接，有需要时重新创建连接。

配置名称	说明	默认值	影响分析
request.timeout.ms	消费请求的超时时间。单位：毫秒。	30000	请求超时时间，请求超时时会失败然后不断重试。

## 15.10.2 Kafka 日志介绍

本章节内容适用于MRS 3.x及后续版本。

### 日志描述

**日志路径：**Kafka相关日志的默认存储路径为“/var/log/Bigdata/kafka”，审计日志的默认存储路径为“/var/log/Bigdata/audit/kafka”。

- Broker：“/var/log/Bigdata/kafka/broker”（运行日志）
- KafkaUI：“/var/log/Bigdata/kafka/ui”（运行日志）

**日志归档规则：**Kafka的日志启动了自动压缩归档功能，默认情况下，当日志大小超过30MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd\_hh-mm-ss>.[编号].log.zip”。默认最多保留最近的20个压缩文件，压缩文件保留个数和压缩文件阈值可以配置。

表 15-12 Kafka 日志列表

日志类型	日志文件名	描述
运行日志	server.log	Broker进程的server运行日志。
	controller.log	Broker进程的controller运行日志。
	kafka-request.log	Broker进程的request运行日志。
	log-cleaner.log	Broker进程的cleaner运行日志。
	state-change.log	Broker进程的state-change运行日志。
	kafkaServer-<SSH_USER>-<DATE>-<PID>-gc.log	Broker进程的GC日志。
	postinstall.log	Broker安装后的工作日志。
	prestart.log	Broker启动前的工作日志。
	checkService.log	Broker启动是否成功的检查日志。

日志类型	日志文件名	描述
	start.log	Broker进程启动日志。
	stop.log	Broker进程停止日志。
	checkavailable.log	Kafka服务健康状态检查日志。
	checkInstanceHealth.log	Broker实例健康状态检测日志。
	kafka-authorizer.log	Broker鉴权日志。
	kafka-root.log	Broker基础日志。
	cleanup.log	Broker卸载的清理日志。
	metadata-backup-recovery.log	Broker备份恢复日志。
	ranger-kafka-plugin-enable.log	Broker启动Ranger插件日志。
	server.out	Broker jvm日志。
	audit.log	Ranger鉴权插件鉴权日志。此日志统一归档在“/var/log/Bigdata/audit/kafka”目录下。

表 15-13 KafkaUI 日志列表

日志类型	日志文件名	描述
运行日志	kafka-ui.log	KafkaUI进程的运行日志。
	postinstall.log	KafkaUI安装后的工作日志。
	cleanup.log	KafkaUI卸载的清理日志。
	prestart.log	KafkaUI启动前的工作日志。
	ranger-kafka-plugin-enable.log	KafkaUI启动Ranger插件日志。
	start.log	KafkaUI进程启动日志。
	stop.log	KafkaUI进程停止日志。
	start.out	KafkaUI进程启动信息。
审计日志	audit.log	KafkaUI服务审计日志。

日志类型	日志文件名	描述
鉴权日志	kafka-authorizer.log	Kafka开源自带鉴权插件运行日志。 此日志统一归档在“/var/log/Bigdata/audit/kafka/kafkai”目录下。
	ranger-authorizer.log	Ranger鉴权插件运行日志。 此日志统一归档在“/var/log/Bigdata/audit/kafka/kafkai”目录下。

## 日志级别

Kafka提供了如表15-14所示的日志级别。

运行日志的级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 15-14 日志级别

级别	描述
ERROR	ERROR表示系统运行的错误信息。
WARN	WARN表示当前事件处理存在异常信息。
INFO	INFO表示记录系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 请参考[修改集群服务配置参数](#)，进入Kafka的“全部配置”页面。
- 步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤3** 选择所需修改的日志级别。
- 步骤4** 保存配置，在弹出窗口中单击“确定”使配置生效。

----结束

## 日志格式

Kafka的日志格式如下所示

表 15-15 日志格式

日志类型	格式	示例
运行日志	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线 程名字> <log中的 message> <日志事件调用 类全名>(<日志打印文件 >:<行号>)	2015-08-08 11:09:53,483   INFO   [main]   Loading logs.   kafka.log.LogManager (Logging.scala:68)
	<yyyy-MM-dd HH:mm:ss><HostName> <组件名 ><logLevel><Message>	2015-08-08 11:09:51 10-165-0-83 Kafka INFO Running kafka-start.sh.

## 15.10.3 更改 Broker 的存储目录

### 操作场景

本章节内容适用于MRS 3.x及后续版本。

增加Broker的存储目录时，MRS集群管理员需要在FusionInsight Manager中修改Broker的存储目录，以保证Kafka正常工作，新创建的主题分区将在分区最少的目录中生成。适用于以下场景：

#### 说明

由于Kafka不感知磁盘容量，建议各Broker实例配置的磁盘个数和容量保持一致。

- 更改Broker角色的存储目录，所有Broker实例的存储目录将同步修改。
- 更改Broker单个实例的存储目录，只对单个实例生效，其他节点Broker实例存储目录不变。

### 对系统的影响

- 更改Broker角色的存储目录需要重新启动服务，服务重启时无法访问。
- 更改Broker单个实例的存储目录需要重新启动实例，该节点Broker实例重启时无法提供服务。
- 服务参数配置如果使用旧的存储目录，需要更新为新目录。

### 前提条件

- 在各个数据节点准备并安装好新磁盘，并格式化磁盘。
- 已安装好Kafka客户端。
- 更改Broker单个实例的存储目录时，保持活动的Broker实例数必须大于创建主题时指定的备份数。

### 操作步骤

#### 更改Kafka角色的存储目录

**步骤1** 以root用户登录到安装Kafka服务的各个数据节点中，执行如下操作。

1. 创建目标目录。  
例如目标目录为“`${BIGDATA_DATA_HOME}/kafka/data2`”：  
执行`mkdir ${BIGDATA_DATA_HOME}/kafka/data2`。
2. 挂载目录到新磁盘。例如挂载“`${BIGDATA_DATA_HOME}/kafka/data2`”到新磁盘。
3. 修改新目录的权限。  
例如新目录路径为“`${BIGDATA_DATA_HOME}/kafka/data2`”：  
执行`chmod 700 ${BIGDATA_DATA_HOME}/kafka/data2 -R`和`chown omm:wheel ${BIGDATA_DATA_HOME}/kafka/data2 -R`。

**步骤2** MRS 3.x及后续版本，登录FusionInsight Manager，然后选择“集群 > 服务 > Kafka > 配置”。

**步骤3** 添加新目录到“log.dirs”的默认值后面。

在搜索框中输入“log.dirs”进行搜索，将新目录添加到配置项“log.dirs”的默认值后面，多个目录使用逗号分隔。例如“

`${BIGDATA_DATA_HOME}/kafka/data1/kafka-logs,${BIGDATA_DATA_HOME}/kafka/data2/kafka-logs`”。

**步骤4** 单击“保存”，并单击“确定”。界面提示“操作成功”，单击“完成”。

**步骤5** 选择“集群 > 服务 > Kafka”，右上角选择“更多 > 重启服务”，重启Kafka服务。

#### 更改Kafka单个实例的存储目录

**步骤6** 以root用户登录到Broker节点，执行如下操作。

1. 创建目标目录。  
例如目标目录为“`${BIGDATA_DATA_HOME}/kafka/data2`”：  
执行`mkdir ${BIGDATA_DATA_HOME}/kafka/data2`。
2. 挂载目录到新磁盘。例如挂载“`${BIGDATA_DATA_HOME}/kafka/data2`”到新磁盘。
3. 修改新目录的权限。  
例如新目录路径为“`${BIGDATA_DATA_HOME}/kafka/data2`”：  
执行`chmod 700 ${BIGDATA_DATA_HOME}/kafka/data2 -R`和`chown omm:wheel ${BIGDATA_DATA_HOME}/kafka/data2 -R`。

**步骤7** MRS 3.x及后续版本，登录FusionInsight Manager，然后选择“集群 > 服务 > Kafka > 实例”。

**步骤8** 单击指定的Broker实例并切换到“实例配置”。

在搜索框中输入“log.dirs”进行搜索，将新目录添加到配置项“log.dirs”的默认值后面，多个目录使用逗号分隔。例如“`${BIGDATA_DATA_HOME}/kafka/data1/kafka-logs,${BIGDATA_DATA_HOME}/kafka/data2/kafka-logs`”。

**步骤9** 单击“保存”，并单击“确定”，界面提示“操作成功”，单击“完成”。

**步骤10** 在Broker实例页面选择“更多 > 重启实例”，重启Broker实例。

----结束

## 15.10.4 迁移 Kafka 节点内数据

### 操作场景

用户可以根据业务需求，通过Kafka客户端命令，在不停止服务的情况下，进行节点内磁盘间的分区数据迁移。也可以通过KafkaUI进行分区迁移。

### 前提条件

- MRS集群管理员已明确业务需求，并准备一个Kafka用户（属于kafkaadmin组，普通模式不需要）。
- 已安装Kafka客户端。
- Kafka实例状态和磁盘状态均正常。
- 根据待迁移分区当前的磁盘空间占用情况，评估迁移后，不会导致新迁移后的磁盘空间不足。

### 使用 Kafka 客户端迁移数据

**步骤1** 以客户端安装用户，登录已安装Kafka客户端的节点。

**步骤2** 执行以下命令，切换到Kafka客户端安装目录，例如“/opt/kafkaclient”。

```
cd /opt/kafkaclient
```

**步骤3** 执行以下命令，配置环境变量。

```
source bigdata_env
```

**步骤4** 执行以下命令，进行用户认证（普通模式跳过此步骤）。

```
kinit 组件业务用户
```

**步骤5** 执行以下命令，切换到Kafka客户端目录。

```
cd Kafka/kafka/bin
```

**步骤6** 执行以下命令，查看待迁移的Partition对应的Topic的详细信息。

**安全模式：**

```
./kafka-topics.sh --describe --bootstrap-server Kafka集群IP:21007 --command-config ../config/client.properties --topic 主题名称
```

**普通模式：**

```
./kafka-topics.sh --describe --bootstrap-server Kafka集群IP:21005 --command-config ../config/client.properties --topic 主题名称
```



```

Topic:testws PartitionCount:24 ReplicationFactor:2 Configs:
Topic: testws Partition: 0 Leader: 4 Replicas: 4,3 Isr: 4,3
Topic: testws Partition: 1 Leader: 5 Replicas: 5,4 Isr: 5,4
Topic: testws Partition: 2 Leader: 6 Replicas: 6,5 Isr: 6,5
Topic: testws Partition: 3 Leader: 3 Replicas: 3,6 Isr: 3,6
Topic: testws Partition: 4 Leader: 4 Replicas: 4,5 Isr: 4,5
Topic: testws Partition: 5 Leader: 5 Replicas: 5,4 Isr: 5,4
Topic: testws Partition: 6 Leader: 6 Replicas: 6,3 Isr: 6,3
Topic: testws Partition: 7 Leader: 3 Replicas: 3,4 Isr: 3,4
Topic: testws Partition: 8 Leader: 4 Replicas: 4,6 Isr: 4,6
Topic: testws Partition: 9 Leader: 5 Replicas: 5,3 Isr: 5,3
Topic: testws Partition: 10 Leader: 6 Replicas: 6,4 Isr: 6,4
Topic: testws Partition: 11 Leader: 3 Replicas: 3,5 Isr: 3,5
Topic: testws Partition: 12 Leader: 4 Replicas: 4,3 Isr: 4,3
Topic: testws Partition: 13 Leader: 5 Replicas: 5,4 Isr: 5,4
Topic: testws Partition: 14 Leader: 6 Replicas: 6,5 Isr: 6,5
Topic: testws Partition: 15 Leader: 3 Replicas: 3,6 Isr: 3,6
Topic: testws Partition: 16 Leader: 4 Replicas: 4,5 Isr: 4,5
Topic: testws Partition: 17 Leader: 5 Replicas: 5,6 Isr: 5,6
Topic: testws Partition: 18 Leader: 6 Replicas: 6,3 Isr: 6,3
Topic: testws Partition: 19 Leader: 3 Replicas: 3,4 Isr: 3,4
Topic: testws Partition: 20 Leader: 4 Replicas: 4,6 Isr: 4,6
Topic: testws Partition: 21 Leader: 5 Replicas: 5,3 Isr: 5,3
Topic: testws Partition: 22 Leader: 6 Replicas: 6,4 Isr: 6,4

```

**步骤7** 执行以下命令，查询Broker\_ID和IP对应关系。

```
./kafka-broker-info.sh --zookeeper ZooKeeper的quorumpeer实例业务IP.ZooKeeper客户端端口号/kafka
```

```

Broker_ID IP_Address

4 192.168.0.100
5 192.168.0.101
6 192.168.0.102

```

#### 📖 说明

- ZooKeeper的quorumpeer实例业务IP：  
ZooKeeper服务所有quorumpeer实例业务IP。登录FusionInsight Manager，选择“集群 > 服务 > ZooKeeper > 实例”，可查看所有quorumpeer实例所在主机业务IP地址。
- ZooKeeper客户端端口号：  
登录FusionInsight Manager，选择“集群 > 服务 > ZooKeeper”，在“配置”页签查看“clientPort”的值。默认为24002。

**步骤8** 从**步骤6**和**步骤7**回显中获取分区的分布信息和节点信息，在当前目录下创建执行重新分配的json文件。

以迁移的是Broker\_ID为6的节点的分区为例，迁移到“/srv/BigData/hadoop/data1/kafka-logs”，完成迁移所需的json配置文件，内容如下。

```
{"partitions":[{"topic": "testws","partition": 2,"replicas": [6,5],"log_dirs": ["/srv/BigData/hadoop/data1/kafka-logs","any"]}],"version":1}
```

#### 📖 说明

- topic为Topic名称，此处以testws为例，具体以实际为准。
- partition为Topic分区。
- replicas中的数字对应Broker\_ID。replicas必须与分区的副本数相对应，不然会造成副本缺少的情况。在本案例中分区所在的replicas对应6和5，只迁移Broker\_ID为6的节点的分区中的数据时，也必须把Broker\_ID为5的节点的分区带上。
- log\_dirs为需要迁移的磁盘路径。此样例迁移的是Broker\_ID为6的节点，Broker\_ID为5的节点对应的log\_dirs可设置为“any”，Broker\_ID为6的节点对应的log\_dirs设置为“/srv/BigData/hadoop/data1/kafka-logs”。**注意路径需与节点对应。**

**步骤9** 使用如下命令，执行重分配操作。

**安全模式：**

```
./kafka-reassign-partitions.sh --bootstrap-server Broker业务IP:21007 --
command-config ../config/client.properties --zookeeper {zk_host}:{port}/kafka
--reassignment-json-file 步骤8中编写的json文件路径 --execute
```

普通模式：

```
./kafka-reassign-partitions.sh --bootstrap-server Broker业务IP:21005 --
command-config ../config/client.properties --zookeeper {zk_host}:{port}/kafka
--reassignment-json-file 步骤8中编写的json文件路径 --execute
```

提示 “Successfully started reassignment of partitions” 表示执行成功。

----结束

## 使用 KafkaUI 迁移分区（MRS 3.1.2 及之后版本）

**步骤1** 进入KafkaUI界面。

1. 使用具有KafkaUI页面访问权限的用户登录FusionInsight Manager，选择“集群 > 服务 > Kafka”。  
如需在页面上进行相关操作，例如创建Topic，需同时授予用户相关权限，请参考 [Kafka用户权限说明](#)。
2. 在“KafkaManager WebUI”右侧，单击URL链接，访问KafkaUI的页面。

**步骤2** 单击“Generate assignment”进入分区迁移页面。

**步骤3** 在“Brokers”处选择要将主题重新分配的Broker。

**步骤4** 单击“Generate Partition Assignments”生成分区迁移方案。

Generate Partition Assignments

Choose brokers to reassign topic to:

\* Brokers:

Select All  
 1       2       3

Current Assignments

Partition	Replicas
__KafkaMetricReport-0	[3, 2]
__KafkaMetricReport-1	[1, 3]
cdl-connect-configs-0	[3, 1, 2]
cdl-connect-status-0	[1, 3, 2]
cdl-connect-status-1	[2, 1, 3]
cdl-connect-status-2	[3, 2, 1]
cdl-connect-status-3	[1, 2, 3]
cdl-connect-status-4	[2, 3, 1]
cdl-connect-offsets-0	[1, 3, 2]

**步骤5** 继续单击“Run assignment”执行分区迁移方案，完成分区迁移。

----结束

## 15.10.5 均衡 Kafka 扩容节点后数据

### 操作场景

用户可以在Kafka扩容节点后，在客户端中执行Kafka均衡工具来均衡Kafka集群的负载。

本章节内容适用于MRS 3.x之前版本。3.x及之后版本请参考[配置Kafka数据均衡工具](#)。

### 前提条件

- MRS集群管理员已明确业务需求，并准备一个Kafka管理员用户（属于kafkaadmin组，普通模式不需要）。
- 已安装Kafka客户端，客户端安装目录如“/opt/client”。
- 本示例需创建两个Topic，可参考[步骤7](#)，分别命名为“test\_2”和“test\_3”，并创建“move-kafka-topic.json”文件，创建路径如“/opt/client/Kafka/kafka”，Topic格式内容如下：

```
{
 "topics":
 [{"topic":"test_2"}, {"topic":"test_3"}],
 "version":1
}
```

### 操作步骤

**步骤1** 以客户端安装用户，登录安装Kafka客户端的节点。

**步骤2** 切换到Kafka客户端安装目录。

```
cd /opt/client
```

**步骤3** 执行以下命令，配置环境变量。

```
source bigdata_env
```

**步骤4** 执行以下命令，进行用户认证。（普通模式跳过此步骤）

```
kinit 组件业务用户
```

**步骤5** 执行以下命令进入Kafka客户端的bin目录。

```
cd Kafka/kafka/bin
```

**步骤6** 执行以下命令生成执行计划。

```
./kafka-reassign-partitions.sh --zookeeper 172.16.0.119:2181/kafka --topics-to-move-json-file ../move-kafka-topic.json --broker-list "1,2,3" --generate
```

#### 📖 说明

- 172.16.0.119: ZooKeeper实例的业务IP。
- --broker-list "1,2,3": 参数中的“1,2,3”为扩容后的所有broker\_id。

```
[root@node-master1SPXC bin]# ./kafka-reassign-partitions.sh --zookeeper 172.16.0.119:2181/kafka --topics-to-move-json-file ../reassignment.json --execute --throttle 50000000
Current partition replica assignment
{"version":1,"partitions":[{"topic":"test_2","partition":3,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":4,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":5,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":3,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_2","partition":2,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":0,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_3","partition":2,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":6,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":4,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":0,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":1,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_2","partition":1,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":5,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_3","partition":6,"replicas":[1,2],"log_dirs":["any","any"]}]}

Proposed partition reassignment configuration
{"version":1,"partitions":[{"topic":"test_3","partition":0,"replicas":[2,3],"log_dirs":["any","any"]}, {"topic":"test_2","partition":1,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":6,"replicas":[3,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":2,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":3,"replicas":[3,2],"log_dirs":["any","any"]}, {"topic":"test_3","partition":5,"replicas":[1,3],"log_dirs":["any","any"]}, {"topic":"test_2","partition":0,"replicas":[3,1],"log_dirs":["any","any"]}, {"topic":"test_2","partition":5,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":4,"replicas":[3,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":2,"replicas":[2,3],"log_dirs":["any","any"]}, {"topic":"test_3","partition":1,"replicas":[3,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":6,"replicas":[2,3],"log_dirs":["any","any"]}, {"topic":"test_3","partition":3,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_2","partition":4,"replicas":[1,3],"log_dirs":["any","any"]}]}
[root@node-master1SPXC bin]#
```

**步骤7** 执行vim ../reassignment.json创建“reassignment.json”文件并保存，保存路径为“/opt/kafkaclient/Kafka/kafka”。

拷贝步骤6中生成的“Proposed partition reassignment configuration”下的内容至“reassignment.json”文件，如下所示：

```
{"version":1,"partitions":[{"topic":"test","partition":4,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test","partition":1,"replicas":[1,3],"log_dirs":["any","any"]}, {"topic":"test","partition":3,"replicas":[3,1],"log_dirs":["any","any"]}, {"topic":"test","partition":0,"replicas":[3,2],"log_dirs":["any","any"]}, {"topic":"test","partition":2,"replicas":[2,1],"log_dirs":["any","any"]}]}

```

**步骤8** 执行以下命令进行分区重分布。

```
./kafka-reassign-partitions.sh --zookeeper 172.16.0.119:2181/kafka --reassignment-json-file ../reassignment.json --execute --throttle 50000000
```

#### 说明

--throttle 50000000：限制网络带宽为50MB。带宽可根据数据量大小及客户对均衡时间的要求进行调整，5TB数据量，使用50MB带宽，均衡时长约8小时。

```
[root@node-master1SPXC bin]# vim ../reassignment.json
[root@node-master1SPXC bin]# ./kafka-reassign-partitions.sh --zookeeper 172.16.0.119:2181/kafka --reassignment-json-file ../reassignment.json --execute --throttle 50000000
Current partition replica assignment

{"version":1,"partitions":[{"topic":"test_2","partition":3,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":4,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":5,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":3,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_2","partition":2,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":0,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_3","partition":2,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":6,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":4,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":0,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_3","partition":1,"replicas":[2,1],"log_dirs":["any","any"]}, {"topic":"test_2","partition":1,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_2","partition":5,"replicas":[1,2],"log_dirs":["any","any"]}, {"topic":"test_3","partition":6,"replicas":[1,2],"log_dirs":["any","any"]}]}

Save this to use as the --reassignment-json-file option during rollback
Warning: You must run Verify periodically, until the reassignment completes, to ensure the throttle is removed. You can also lift the throttle by rerunning the Execute command passing a new value.
The inter-broker throttle limit was set to 50000000 B/s
Successfully started reassignment of partitions.
[root@node-master1SPXC bin]#
```

**步骤9** 执行以下命令查看迁移状态。

```
./kafka-reassign-partitions.sh --zookeeper 172.16.0.119:2181/kafka --reassignment-json-file ../reassignment.json --verify
```

```
Terminal 1: [root@node-str-coreRuzk0001 kafka-logs]# ll
total 56
-rw-r----- 1 omm wheel 4 Sep 14 21:30 cleaner-offset-check
-rw-r----- 1 omm wheel 4 Sep 14 21:31 log-start-offset-check
-rw-r----- 1 omm wheel 54 Sep 14 19:39 meta.properties
-rw-r----- 1 omm wheel 103 Sep 14 21:31 recovery-point-offset
-rw-r----- 1 omm wheel 103 Sep 14 21:32 replication-offset-check
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-0
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-1
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-4
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-5
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-6
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-1
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-2
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-3
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-5
[root@node-str-coreRuzk0001 kafka-logs]#

Terminal 2: [root@node-str-coreRuzk0002 kafka-logs]# ll
total 56
-rw-r----- 1 omm wheel 4 Sep 14 21:30 cleaner-offset-check
-rw-r----- 1 omm wheel 4 Sep 14 21:31 log-start-offset-check
-rw-r----- 1 omm wheel 54 Sep 14 19:39 meta.properties
-rw-r----- 1 omm wheel 103 Sep 14 21:31 recovery-point-offset
-rw-r----- 1 omm wheel 103 Sep 14 21:32 replication-offset-check
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-1
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-2
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-3
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:11 test_2-5
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-0
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-2
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-3
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-4
drwxr-xr-x 2 omm wheel 4096 Sep 14 21:12 test_3-6
[root@node-str-coreRuzk0002 kafka-logs]#

Terminal 3: [root@node-master1SPXC bin]# ./kafka-reassign-partitions.sh --zookeeper.connect zoo1:2181,zoo2:2181,zoo3:2181 --command.config kafka-reassignment.json --verify
Status of partition reassignment:
Reassignment of partition test_2-3 completed successfully
Reassignment of partition test_2-4 completed successfully
Reassignment of partition test_3-5 completed successfully
Reassignment of partition test_3-3 completed successfully
Reassignment of partition test_2-2 completed successfully
Reassignment of partition test_3-0 completed successfully
Reassignment of partition test_3-2 completed successfully
Reassignment of partition test_2-6 completed successfully
Reassignment of partition test_3-4 completed successfully
Reassignment of partition test_2-0 completed successfully
Reassignment of partition test_3-1 completed successfully
Reassignment of partition test_2-1 completed successfully
Reassignment of partition test_2-5 completed successfully
Reassignment of partition test_3-6 completed successfully
Throttle was removed.
[root@node-master1SPXC bin]#
```

---结束

## 15.11 Kafka 常见问题

### 15.11.1 Kafka 业务规格说明

本章节内容适用于MRS 3.x及后续版本。

#### 支持的 Topic 上限

支持Topic的个数，受限于进程整体打开的文件句柄数（现场环境一般主要是数据文件和索引文件占用比较多）。

1. 可通过- 2. 执行lsof -p <Kafka PID>命令，查看当前单节点上Kafka进程打开的文件句柄（会继续增加）；
- 3. 权衡当前需要创建的Topic创建完成后，会不会达到文件句柄上限，每个Partition文件夹下会最多保存多大的数据，会产生多少个数据文件（\*.log文件，默认配置为1GB，可通过修改log.segment.bytes来调整大小）和索引文件（\*.index文件，默认配置为10MB，可通过修改log.index.size.max.bytes来调整大小），是否会影响Kafka正常运行。

#### Consumer 的并发量

在一个应用中，同一个Group的Consumer并发量建议与Topic的Partition个数保持一致，保证每个Consumer对应消费一个Partition上的数据。若Consumer的并发量多于Partition个数，那么多余的Consumer将消费不到数据。

#### Topic 和 Partition 的划分关系说明

- 假设集群中部署了K个Kafka节点，每个节点上配置的磁盘个数为N，每块磁盘大小为M，集群共有n个Topic（T1,T2...Tn），并且其中第m个Topic的每秒输入数据总流量为X(Tm) MB/s，配置的副本数为R(Tm)，配置数据保存时间为Y(Tm)小时，那么整体必须满足：

$$M \times N \times K > \sum_{i=1}^{Tn} (X(i)R(i)Y(i) \times 3600)$$

- 假设单个磁盘大小为M，该磁盘上有n个Partition ( P0,P1……Pn )，并且其中第m个Partition的每秒写入数据流量为Q(Pm) MB/s（计算方法：所属Topic的数据流量除以Partition数）、数据保存时间为T(Pm)小时，那么单个磁盘必须满足：

$$M > \sum_{i=P0}^{Pn} (Q(i)T(i) \times 3600)$$

- 根据吞吐量粗略计算，假设生产者可以达到的吞吐量为P，消费者可以达到的吞吐量为C，预期Kafka吞吐量为T，那么建议该Topic的Partition数目设置为Max(T/P, T/C)。

#### 📖 说明

- 在Kafka集群中，分区越多吞吐量越高，但是分区过多也存在潜在影响，例如文件句柄增加、不可用性增加（如：某个节点故障后，部分Partition重选Leader后时间窗口会比较大）及端到端时延增加等。
- 建议：单个Partition的磁盘占用最大不超过100GB；单节点上Partition数目不超过3000；整个集群的分区总数不超过10000。

## 15.11.2 Kafka 相关特性说明

### Kafka Idempotent 特性

特性说明：Kafka从0.11.0.0版本引入了创建幂等性Producer的功能，开启此特性后，Producer自动升级成幂等性Producer，当Producer发送了相同字段值的消息后，Broker会自动感知消息是否重复，继而避免数据重复。需要注意的是，这个特性只能保证单分区上的幂等性，即一个幂等性Producer能够保证某个主题的一个分区内不出现重复消息；只能实现单会话上的幂等性，这里的会话指的是Producer进程的一次运行，即重启Producer进程后，幂等性不保证。

开启方法：

1. 二次开发代码中添加 “props.put(“enable.idempotence”, true)”。
2. 客户端配置文件中添加 “enable.idempotence = true”。

### Kafka Transaction 特性

特性说明：Kafka在0.11版本中，引入了事务特性，Kafka事务特性指的是一系列的生产者生产消息和消费者提交偏移量的操作在一个事务中，或者说是一个原子操作，生产消息和提交偏移量同时成功或者失败，此特性提供的是read committed隔离级别的事务，保证多条消息原子性的写入到目标分区，同时也能保证Consumer只能看到成功提交的事务消息。Kafka中的事务特性主要用于以下两种场景：

1. 生产者发送多条数据可以封装在一个事务中，形成一个原子操作。多条消息要么都发送成功，要么都发送失败。
2. read-process-write模式：将消息消费和生产封装在一个事务中，形成一个原子操作。在一个流式处理的应用中，常常一个服务需要从上游接收消息，然后经过处理后送达到下游，这就对应着消息的消费和生产。

二次开发代码样例如下：

```
// 初始化配置,开启事务特性
Properties props = new Properties();
props.put("enable.idempotence", true);
props.put("transactional.id", "transaction1");
...

KafkaProducer producer = new KafkaProducer<String, String>(props);

// init 事务
producer.initTransactions();
try {
 // 开启事务
 producer.beginTransaction();
 producer.send(record1);
 producer.send(record2);
 // 结束事务
 producer.commitTransaction();
} catch (KafkaException e) {
 // 事务 abort
 producer.abortTransaction();
}
```

## 就近消费特性

特性说明：Kafka 2.4.0之前版本，客户端的生产、消费都是面向各个partition的leader副本，follower副本仅用来进行数据冗余，不对外提供服务，常会导致leader副本压力较大，且在跨机房、机架的消费场景下，常会导致大量的机房、机架间的数据传输；Kafka 2.4.0及之后版本，Kafka内核支持从follower副本消费数据，在跨机房、机架的场景中，会大大降低数据传输量，减轻网络带宽压力。社区开放了ReplicaSelector接口来支持此特性，MRS Kafka中默认提供两种实现此接口的方式。

1. RackAwareReplicaSelector：优先从相同机架的副本进行消费（机架内就近消费特性）。
2. AzAwareReplicaSelector：优先从相同AZ内的节点上的副本进行消费（AZ内就近消费特性）。

以RackAwareReplicaSelector为例，描述实现就近消费副本的选取：

```
public class RackAwareReplicaSelector implements ReplicaSelector {

 @Override
 public Optional<ReplicaView> select(TopicPartition topicPartition,
 ClientMetadata clientMetadata,
 PartitionView partitionView) {
 if (clientMetadata.rackId() != null && !clientMetadata.rackId().isEmpty()) {
 Set<ReplicaView> sameRackReplicas = partitionView.replicas().stream()
 // 过滤与客户端处于相同Rack的副本
 .filter(replicaInfo -> clientMetadata.rackId().equals(replicaInfo.endpoint().rack()))
 .collect(Collectors.toSet());
 if (sameRackReplicas.isEmpty()) {
 // 如果没有副本与客户端处于相同Rack，则返回leader副本
 return Optional.of(partitionView.leader());
 } else {
 // 到这里说明存在与客户端位于同一Rack的副本
 if (sameRackReplicas.contains(partitionView.leader())) {
 // 如果客户端和leader在同一个机架，则优先返回leader副本
 return Optional.of(partitionView.leader());
 } else {
 // 否则，返回和leader同步最新的副本
 return sameRackReplicas.stream().max(ReplicaView.comparator());
 }
 }
 }
 } else {
 // 如果客户端请求中不包含机架信息，则默认返回leader副本
 return Optional.of(partitionView.leader());
 }
 }
}
```

```
}
}
```

开启方法：

1. 服务端：根据不同特性更新“replica.selector.class”配置项：
  - 开启“机架内就近消费特性”，配置为“org.apache.kafka.common.replica.RackAwareReplicaSelector”。
  - 开启“AZ内就近消费特性”，配置为“org.apache.kafka.common.replica.AzAwareReplicaSelector”。
2. 客户端：在“{客户端安装目录}/Kafka/kafka/config”目录中的“consumer.properties”消费配置文件里添加“client.rack”配置项：
  - 若服务端开启“机架内就近消费特性”，添加客户端所处的机架信息，如 client.rack = /default0/rack1。
  - 若服务端开启“AZ内就近消费特性”，添加客户端所处的机架信息，如 client.rack = /AZ1/rack1。

## Ranger 统一鉴权特性

特性说明：在Kafka 2.4.0之前版本，Kafka组件仅支持社区自带的SimpleAclAuthorizer鉴权插件，Kafka 2.4.0及之后版本，MRS Kafka同时支持Ranger鉴权插件和社区自带鉴权插件。默认使用Ranger鉴权，基于Ranger鉴权插件，可进行细粒度的Kafka Acl管理。

### 📖 说明

服务端使用Ranger鉴权插件时，若“allow.everyone.if.no.acl.found”配置为“true”，使用非安全端口访问时，所有行为将直接放行。建议使用Ranger鉴权插件的安全集群，不要开启“allow.everyone.if.no.acl.found”。

## 15.11.3 基于 binlog 的 MySQL 数据同步到 MRS 集群中

本章节为您介绍使用Maxwell同步工具将线下基于binlog的数据迁移到MRS Kafka集群中的指导。

Maxwell是一个开源程序（<https://maxwells-daemon.io>），通过读取MySQL的binlog日志，将增删改等操作转为JSON格式发送到输出端（如控制台/文件/Kafka等）。Maxwell可部署在MySQL机器上，也可独立部署在其他与MySQL网络可通的机器上。

Maxwell运行在Linux服务器上，常见的有EulerOS、Ubuntu、Debian、CentOS、OpenSUSE等，且需要Java 1.8+支持。

同步数据具体内容如下。

1. [配置MySQL](#)
2. [安装Maxwell](#)
3. [配置Maxwell](#)
4. [启动Maxwell](#)
5. [验证Maxwell](#)
6. [停止Maxwell](#)
7. [Maxwell生成的数据格式及常见字段含义](#)



## 配置 MySQL

**步骤1** 开启binlog，在MySQL中打开my.cnf文件，在[mysqld] 区块检查是否配置server\_id, log-bin与binlog\_format，若没有配置请执行如下命令添加配置项并重启MySQL，若已经配置则忽略此步骤。

```
$ vi my.cnf

[mysqld]
server_id=1
log-bin=master
binlog_format=row
```

**步骤2** Maxwell需要连接MySQL，并创建一个名称为maxwell的数据库存储元数据，且需要能访问需要同步的数据库，所以建议新建一个MySQL用户专门用来给Maxwell使用。使用root登录MySQL之后，执行如下命令创建maxwell用户（其中XXXXXX是密码，请修改为实际值）。

- 若Maxwell程序部署在非MySQL机器上，则创建maxwell用户需要有远程登录数据库的权限，此时创建命令为

```
mysql> GRANT ALL on maxwell.* to 'maxwell'@'%' identified by 'XXXXXX';
```

```
mysql> GRANT SELECT, REPLICATION CLIENT, REPLICATION SLAVE on *.* to 'maxwell'@'%';
```

- 若Maxwell部署在MySQL机器上，则创建maxwell用户可以设置为只能在本机登录数据库，此时创建命令为

```
mysql> GRANT SELECT, REPLICATION CLIENT, REPLICATION SLAVE on *.* to 'maxwell'@'localhost' identified by 'XXXXXX';
```

```
mysql> GRANT ALL on maxwell.* to 'maxwell'@'localhost';
```

----结束

## 安装 Maxwell

**步骤1** 下载安装包，下载路径为<https://github.com/zendesk/maxwell/releases>，选择名为maxwell-XXX.tar.gz的二进制文件下载，其中XXX为版本号。

**步骤2** 将tar.gz包上传到任意目录下（本示例路径为Master节点的/opt）。

**步骤3** 登录部署Maxwell的服务器，并执行如下命令进入tar.gz包所在目录。

```
cd /opt
```

**步骤4** 执行如下命令解压“maxwell-XXX.tar.gz”压缩包，并进入“maxwell-XXX”文件夹。

```
tar -zxvf maxwell-XXX.tar.gz
```

```
cd maxwell-XXX
```

----结束

## 配置 Maxwell

在maxwell-XXX文件夹下若有conf目录则配置config.properties文件，配置项说明请参见表15-16。若没有conf目录，则是在maxwell-XXX文件夹下将config.properties.example修改成config.properties。

表 15-16 Maxwell 配置项说明

配置项	是否必填	说明	默认值
user	是	连接MySQL的用户名，即步骤2中新创建的用户	-
password	是	连接MySQL的密码，配置文件中包含认证密码信息可能存在安全风险，建议当前场景执行完毕后删除相关配置文件或加强安全管理。	-
host	否	MySQL地址	localhost
port	否	MySQL端口	3306
log_level	否	日志打印级别，可选值为 <ul style="list-style-type: none"> <li>• debug</li> <li>• info</li> <li>• warn</li> <li>• error</li> </ul>	info
output_ddl	否	是否发送DDL(数据库与数据表的定义修改)事件 <ul style="list-style-type: none"> <li>• true: 发送DDL事件</li> <li>• false: 不发送DDL事件</li> </ul>	false
producer	是	生产者类型，配置为kafka <ul style="list-style-type: none"> <li>• stdout: 将生成的事件打印在日志中</li> <li>• kafka: 将生成的事件发送到kafka</li> </ul>	stdout
producer_partition_by	否	分区策略，用来确保相同一类的数据写入到kafka同一分区 <ul style="list-style-type: none"> <li>• database: 使用数据库名称做分区，保证同一个数据库的事件写入到kafka同一个分区中</li> <li>• table: 使用表名称做分区，保证同一个表的事件写入到kafka同一个分区中</li> </ul>	database
ignore_producer_error	否	是否忽略生产者发送数据失败的错误 <ul style="list-style-type: none"> <li>• true: 在日志中打印错误信息并跳过错误的数据，程序继续运行</li> <li>• false: 在日志中打印错误信息并终止程序</li> </ul>	true
metrics_slf4j_interval	否	在日志中输出上传kafka成功与失败数据的数量统计的时间间隔，单位为秒	60
kafka.bootstrap.servers	是	kafka代理节点地址，配置形式为HOST:PORT[,HOST:PORT]	-
kafka_topic	否	写入kafka的topic名称	maxwell

配置项	是否必填	说明	默认值
dead_letter_topic	否	当发送某条记录出错时，记录该条出错记录主键的kafka topic	-
kafka_version	否	Maxwell使用的kafka producer版本号，不能在config.properties中配置，需要在启动命令时用-- kafka_version xxx参数传入	-
kafka_partition_hash	否	划分kafka topic partition的算法，支持default或murmur3	default
kafka_key_format	否	Kafka record的key生成方式，支持array或Hash	Hash
ddl_kafka_topic	否	当output_ddl配置为true时，DDL操作写入的topic	{kafka_topic}
filter	否	过滤数据库或表。 <ul style="list-style-type: none"> <li>若只想采集mydatabase的库，可以配置为exclude: *.*;include: mydatabase.*</li> <li>若只想采集mydatabase.mytable的表，可以配置为exclude: *.*;include: mydatabase.mytable</li> <li>若只想采集mydatabase库下的mytable, mydate_123, mydate_456表，可以配置为exclude: *.*;include: mydatabase.mytable, include: mydatabase./mydate_\\d*/</li> </ul>	-

## 启动 Maxwell

**步骤1** 登录Maxwell所在的服务器。

**步骤2** 执行如下命令进入Maxwell安装目录。

```
cd /opt/maxwell-1.21.0/
```

### 📖 说明

如果是初次使用Maxwell，建议将conf/config.properties中的log\_level改为debug(调试级别)，以便观察启动之后是否能正常从MySQL获取数据并发送到kafka，当整个流程调试通过之后，再把log\_level修改为info，然后先停止再启动Maxwell生效。

```
log level [debug | info | warn | error]
```

```
log_level=debug
```

**步骤3** 执行如下命令启动Maxwell。

```
source /opt/client/bigdata_env
```

```
bin/Maxwell
```

```
bin/maxwell --user='maxwell' --password='XXXXXX' --host='127.0.0.1' \
```

```
--producer=kafka --kafka.bootstrap.servers=kafkahost:9092 --
kafka_topic=Maxwell
```

其中，user，password和host分别表示MySQL的用户名，密码和IP地址，这三个参数可以通过修改配置项配置也可以通过上述命令配置，kafkaHost为流式集群的Core节点的IP地址。

命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。

显示类似如下信息，表示Maxwell启动成功。

```
Success to start Maxwell [78092].
```

----结束

## 验证 Maxwell

**步骤1** 登录Maxwell所在的服务器。

**步骤2** 查看日志。如果日志里面没有ERROR日志，且有打印如下日志，表示与MySQL连接正常。

```
BinlogConnectorLifecycleListener - Binlog connected.
```

**步骤3** 登录MySQL数据库，对测试数据进行更新/创建/删除等操作。操作语句可以参考如下示例。

```
-- 创建库
create database test;
-- 创建表
create table test.e (
 id int(10) not null primary key auto_increment,
 m double,
 c timestamp(6),
 comment varchar(255) charset 'latin1'
);
-- 增加记录
insert into test.e set m = 4.2341, c = now(3), comment = 'I am a creature of light.';
-- 更新记录
update test.e set m = 5.444, c = now(3) where id = 1;
-- 删除记录
delete from test.e where id = 1;
-- 修改表
alter table test.e add column torvalds bigint unsigned after m;
-- 删除表
drop table test.e;
-- 删除库
drop database test;
```

**步骤4** 观察Maxwell的日志输出，如果没有WARN/ERROR打印，则表示Maxwell安装配置正常。

若要确定数据是否成功上传，可设置config.properties中的log\_level为debug，则数据上传成功时会立刻打印如下JSON格式数据，具体字段含义请参考[Maxwell生成的数据格式及常见字段含义](#)。

```
{"database":"test","table":"e","type":"insert","ts":1541150929,"xid":60556,"commit":true,"data":
{"id":1,"m":4.2341,"c":"2018-11-02 09:28:49.297000","comment":"I am a creature of light."}}
.....
```

### 说明

当整个流程调试通过之后，可以把config.properties文件中的配置项log\_level修改为info，减少日志打印量，并重启Maxwell。

```
log level [debug | info | warn | error]
log_level=info
```

----结束

## 停止 Maxwell

**步骤1** 登录Maxwell所在的服务器。

**步骤2** 执行如下命令，获取Maxwell的进程标识（PID）。输出的第二个字段即为PID。

```
ps -ef | grep Maxwell | grep -v grep
```

**步骤3** 执行如下命令，强制停止Maxwell进程。

```
kill -9 PID
```

----结束

## Maxwell 生成的数据格式及常见字段含义

Maxwell生成的数据格式为JSON，常见字段含义如下：

- type：操作类型，包含database-create, database-drop, table-create, table-drop, table-alter, insert, update, delete
- database：操作的数据库名称
- ts：操作时间，13位时间戳
- table：操作的表名
- data：数据增加/删除/修改之后的内容
- old：数据修改前的内容或者表修改前的结构定义
- sql：DDL操作的SQL语句
- def：表创建与表修改的结构定义
- xid：事务唯一ID
- commit：数据增加/删除/修改操作是否已提交

## 15.11.4 如何解决 Kafka topic 无法删除的问题

### 问题

删除Kafka topic后发现未成功删除，如何正常删除？

### 回答

- 可能原因一：配置项“delete.topic.enable”未配置为“true”，只有配置为“true”才能执行真正删除。
- 可能原因二：“auto.create.topics.enable”配置为“true”，其他应用程序有使用该Topic，并且一直在后台运行。

解决方法：

- 针对原因一：配置页面上将“delete.topic.enable”设置为“true”。
- 针对原因二：先停掉后台使用该Topic的应用程序，或者“auto.create.topics.enable”配置为“false”（需要重启Kafka服务），然后再做删除操作。

# 16 使用 KafkaManager

## 16.1 KafkaManager 介绍

KafkaManager是Apache Kafka的管理工具，提供Kafka集群界面化的Metric监控和集群管理。

通过KafkaManager可以：

- 支持管理多个Kafka集群
- 支持界面检查集群状态（主题，消费者，偏移量，分区，副本，节点）
- 支持界面执行副本的leader选举
- 使用选择生成分区分配以选择要使用的分区方案
- 支持界面执行分区重新分配（基于生成的分区方案）
- 支持界面选择配置创建主题（支持多种Kafka版本集群）
- 支持界面删除主题（仅支持0.8.2+并设置了delete.topic.enable = true）
- 支持批量生成多个主题的分区分配，并可选择要使用的分区方案
- 支持批量运行重新分配多个主题的分区
- 支持为已有主题增加分区
- 支持更新现有主题的配置
- 可以为分区级别和主题级别度量标准启用JMX查询
- 可以过滤掉zookeeper中没有ids / owner / & offsets /目录的使用者。

## 16.2 访问 KafkaManager 的 WebUI

用户可以通过KafkaManager的WebUI，在图形化界面监控管理Kafka集群。

### 前提条件

- 已安装KafkaManager服务的集群。
- 获取用户“admin”账号密码。“admin”密码在创建MRS集群时由用户指定。

## 访问 KafkaManager 的 WebUI

**步骤1** 登录集群详情页面，选择“组件管理 > KafkaManager”。

### 📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

**步骤2** 在KafkaManager概述的“KafkaManager WebUI”中单击任意一个UI链接，打开KafkaManager的WebUI页面。

KafkaManager的WebUI支持查看以下信息：

- Kafka集群列表
- Kafka集群Broker节点列表和Metric监控
- Kafka集群副本监控
- Kafka集群Consumer监控

### 📖 说明

在KafkaManager的任何子页面单击左上角KafkaManager的Logo都可以回到KafkaManager的WebUI主界面，显示集群列表信息。

----结束

## 16.3 管理 Kafka 集群

管理Kafka集群包含以下内容：

- [添加集群到KafkaManager的WebUI界面](#)
- [更新集群参数](#)
- [删除KafkaManager的WebUI界面的集群](#)

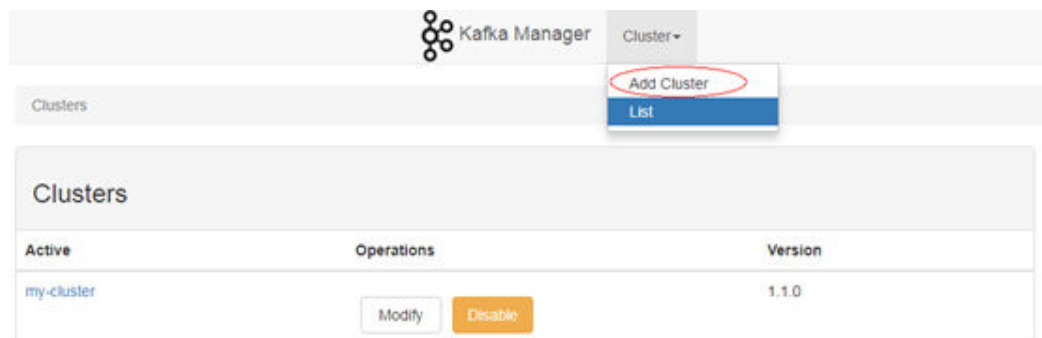
### 添加集群到 KafkaManager 的 WebUI 界面

首次创建Kafka集群后会在KafkaManager的WebUI界面创建名为my-cluster的默认Kafka集群，用户也可以在KafkaManager的WebUI界面自行添加已经通过MRS控制台创建的Kafka集群，用于管理多个Kafka集群。

**步骤1** 登录KafkaManager的WebUI界面。

**步骤2** 在页面上方选择“Cluster > Add Cluster”。

图 16-1 添加集群



**步骤3** 设置待添加集群的参数，如下参数请参考样例，其他参数默认不需要修改。

**表 16-1** 需修改的集群参数

参数名称	取值样例	说明
Cluster Name	mrs-demo	待添加集群在 KafkaManager 的 WebUI 界面中显示的名称。
Cluster Zookeeper Hosts	zk1_ip:zk1_port, zk2_ip:zk2_port/kafka	待添加集群的 Zookeeper 地址。
Kafka Version	1.1.0	待添加集群的 Kafka 版本，默认 1.1.0。
Enable JMX Polling (Set JMX_PORT env variable before starting kafka server)	勾选	-
Poll consumer information (Not recommended for large # of consumers)	勾选	-
Enable Active OffsetCache (Not recommended for large # of consumers)	勾选	-
Display Broker and Topic Size (only works after applying this patch)	勾选	-
Security Protocol	PLAINTEXT	<ul style="list-style-type: none"> <li>• 开启 Kerberos 的 Kafka 集群选择 SASL_PLAINTEXT</li> <li>• 未开启 Kerberos 集群选择 PLAINTEXT</li> </ul>

**步骤4** 单击“Save”完成添加集群。

----结束

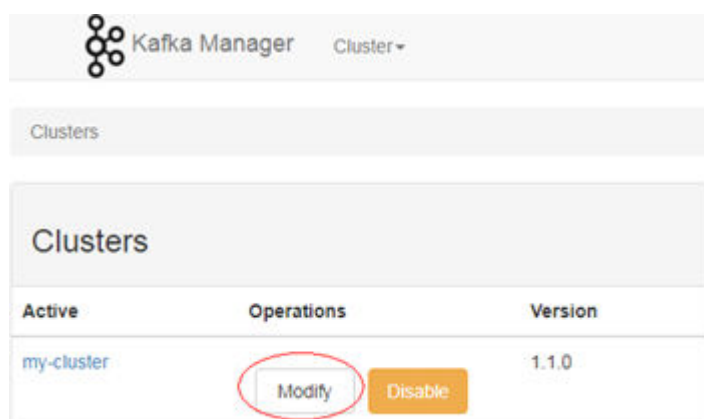
## 更新集群参数

**步骤1** 登录 KafkaManager 的 WebUI 界面。

**步骤2** 在对应集群的“Operations”列单击“Modify”。



图 16-2 更新集群参数



步骤3 进入集群配置参数页面，修改集群参数。

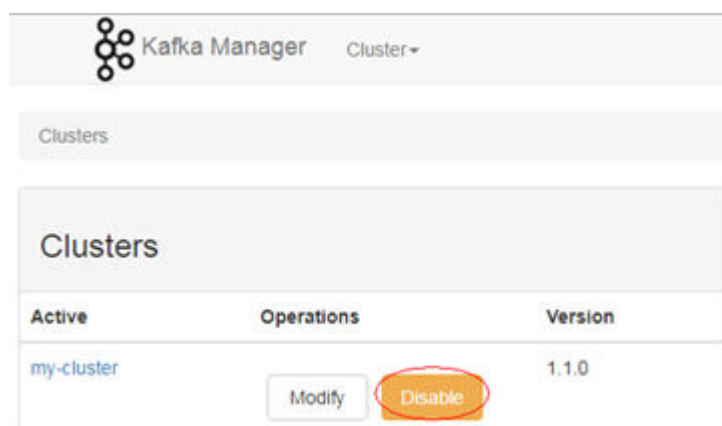
----结束

## 删除 KafkaManager 的 WebUI 界面的集群

步骤1 登录KafkaManager的WebUI界面。

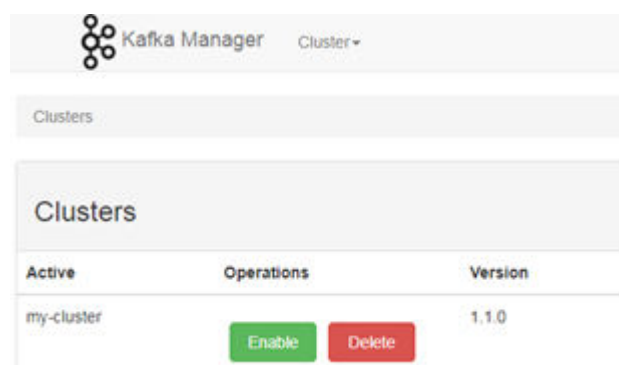
步骤2 在对应集群的“Operations”列单击“Disable”。

图 16-3 停用集群



步骤3 等待集群列表页面的“Operations”列出现“Delete”或“Enable”时，单击“Delete”删除集群。也可以单击“Enable”启用集群。

图 16-4 启用或删除集群



----结束

## 16.4 Kafka 集群监控管理

Kafka集群监控管理包含以下内容：

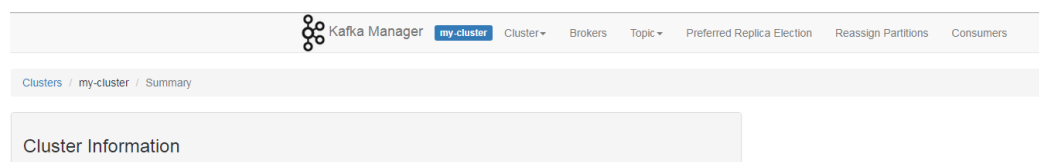
- [查看Broker信息](#)
- [查看Topic信息](#)
- [查看Consumers信息](#)
- [通过KafkaManager修改Topic的partition](#)

### 查看 Broker 信息

**步骤1** 登录KafkaManager的WebUI界面。

**步骤2** 在集群列表页面单击对应集群名称进入集群Summary页面。

图 16-5 集群 Summary 页面



**步骤3** 单击“Brokers”进入Broker监控页面，该页面包括Broker列表和Broker节点的IO统计信息。

图 16-6 Broker 监控页面

Brokers							Combined Metrics				
Id	Host	Port	JMX Port	Bytes In	Bytes Out	Size	Rate	Mean	1 min	5 min	15 min
1	...	SSL:9093,PLAINTEXT:9092	21006	0.00	0.00	0 B	Messages in /sec	0.00	0.00	0.00	0.00
2	...	SSL:9093,PLAINTEXT:9092	21006	0.00	0.00	0 B	Bytes in /sec	0.05	0.00	0.00	0.00
3	...	SSL:9093,PLAINTEXT:9092	21006	0.00	0.00	0 B	Bytes out /sec	0.02	0.00	0.00	0.00
							Bytes rejected /sec	0.00	0.00	0.00	0.00
							Failed fetch request /sec	0.00	0.00	0.00	0.00
							Failed produce request /sec	0.00	0.00	0.00	0.00

----结束

## 查看 Topic 信息

- 步骤1 登录KafkaManager的WebUI界面。
- 步骤2 在集群列表页面单击对应集群名称进入集群Summary页面。
- 步骤3 单击“Topic > List”查看当前集群的Topic列表及每个Topic的相关信息。

图 16-7 Topic 列表

Topic	# Partitions	# Brokers	Brokers Spread %	Brokers Skew %	Brokers Leader Skew %	# Replicas	Under Replicated %	Leader Size	Producer Message/Sec
_consumer_offsets	50	1	100	0	0	1	0		0.00
test1	2	1	100	0	0	1	0		0.00

- 步骤4 单击具体的Topic名称查看该Topic的详细信息。

图 16-8 Topic 的详细信息

The screenshot shows the Kafka Manager web interface for a topic named 'test1'. The interface is organized into several sections:

- Topic Summary:** A table with the following data:

Replication	1
Number of Partitions	2
Sum of partition offsets	3,000
Total number of Brokers	1
Number of Brokers for Topic	1
Preferred Replicas %	100
Brokers Skewed %	0
Brokers Leader Skewed %	0
Brokers Spread %	100
Under-replicated %	0
Leader Size	
- Operations:** A set of buttons for managing the topic: Delete Topic, Reassign Partitions, Generate Partition Assignments, Add Partitions, Update Config, and Manual Partition Assignments.
- Metrics:** A table showing various metrics over time (Mean, 1 min, 5 min, 15 min):

Rate	Mean	1 min	5 min	15 min
Messages in /sec	0.00	0.00	0.00	0.00
Bytes in /sec	0.00	0.00	0.00	0.00
Bytes out /sec	0.00	0.00	0.00	0.00
Bytes rejected /sec	0.00	0.00	0.00	0.00
Failed fetch request /sec	0.00	0.00	0.00	0.00
Failed produce request /sec	0.00	0.00	0.00	0.00
- Partitions by Broker:** A table showing the distribution of partitions across brokers:

Broker	# of Partitions	# as Leader	Partitions	Skewed?	Leader Skewed?
1	2	2	(0,1)	false	false
- Consumers consuming from this topic:** A list showing the consumer group 'group1' with a Kafka version of 'KF'.
- Partition Information:** A table showing details for each partition:

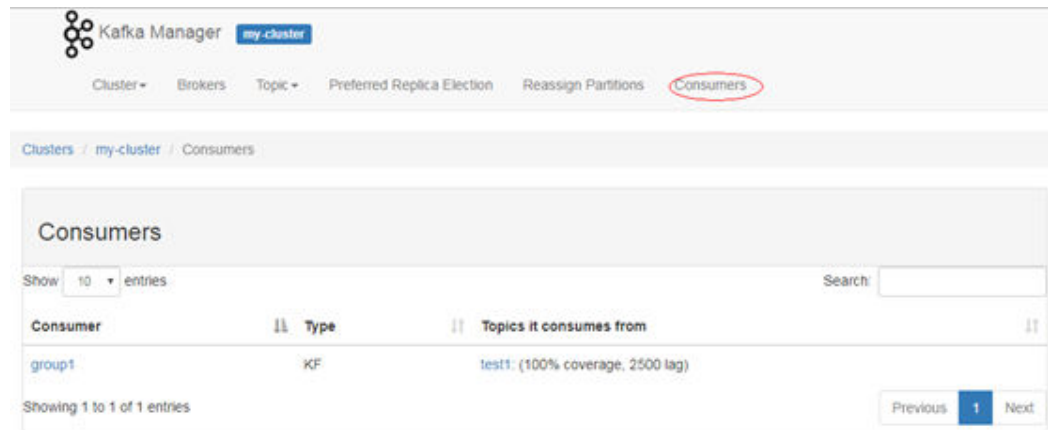
Partition	Latest Offset	Leader	Replicas	In Sync Replicas	Preferred Leader?	Under Replicated?	Leader Size
0	1,500	1	(1)	(1)	true	false	
1	1,500	1	(1)	(1)	true	false	

---结束

## 查看 Consumers 信息

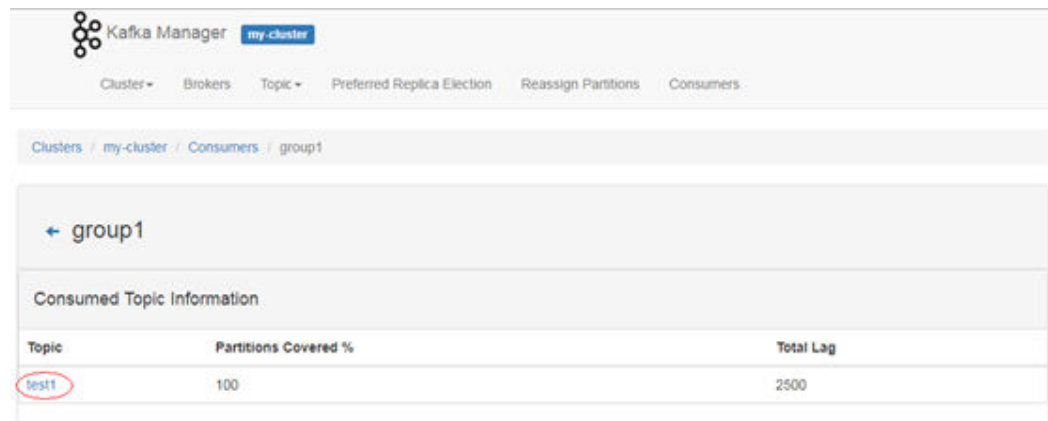
- 步骤1** 登录KafkaManager的WebUI界面。
- 步骤2** 在集群列表页面单击对应集群名称进入集群Summary页面。
- 步骤3** 单击“Consumers”查看当前集群的Consumers列表及每个Consumer的消费信息。

图 16-9 Consumers 列表



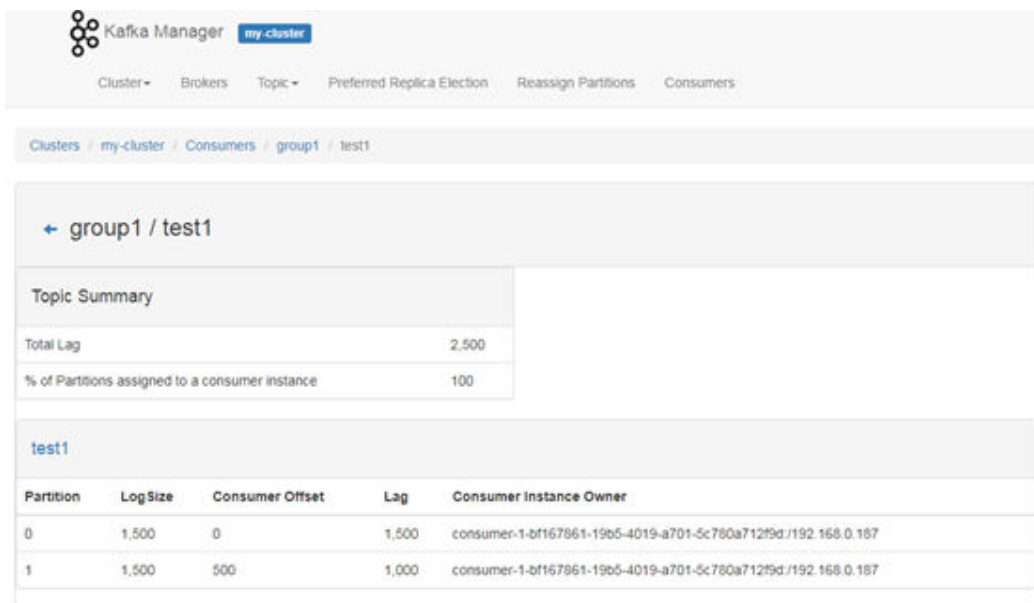
步骤4 单击Consumer的名称查看消费的Topic列表。

图 16-10 Consumer 消费的 Topic 列表



步骤5 单击Consumer下Topic列表中的Topic名称，查看该Consumer对Topic的具体消费情况。

图 16-11 Consumer 对 Topic 的具体消费情况

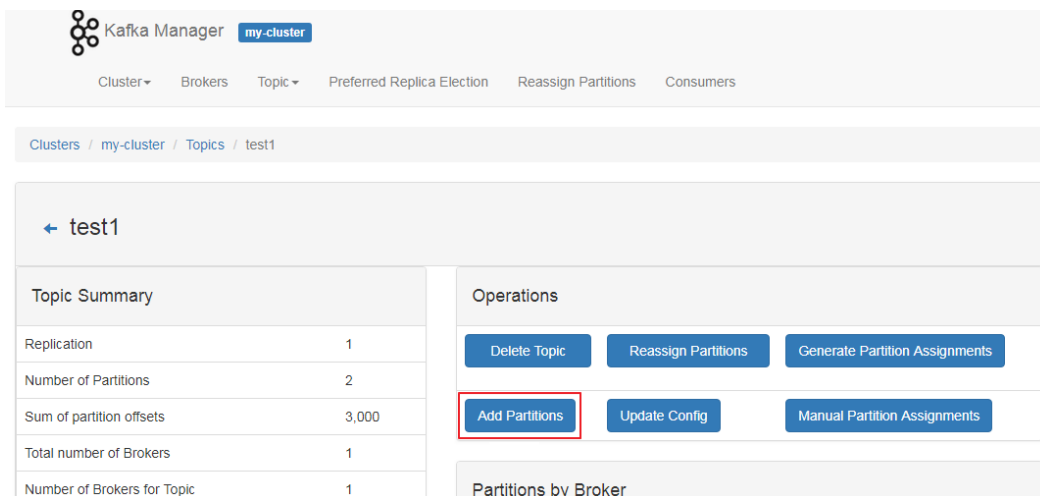


----结束

## 通过 KafkaManager 修改 Topic 的 partition

- 步骤1 登录KafkaManager的WebUI界面。
- 步骤2 在集群列表页面单击对应集群名称进入集群Summary页面。
- 步骤3 单击“Topic > List”进入当前集群的Topic列表页面。
- 步骤4 单击具体的Topic名称进入Topic Summary页面。
- 步骤5 单击“add partitions”，进入添加分区页面。

图 16-12 添加分区



- 步骤6 确认Topic名称并修改“Partitions”数量，单击“Add Partitions”进行分区添加。

图 16-13 修改 Partitions 数量

Clusters / my-cluster / Topics / test1 / Add Partitions

### ← Add Partitions

Add Partitions	Brokers
Topic test1	Select All Select None
Partitions 2	<input checked="" type="checkbox"/> 1 - 192.168.0.112

Add Partitions Cancel

**步骤7** 分区添加成功后，单击“Go to topic view.”返回Topic Summary页面。

**步骤8** 在Topic Summary页面的下方“Partition Information”中确认partition数量。

图 16-14 Partition Information

The screenshot displays the Kafka Manager interface for a cluster named 'my-cluster'. The breadcrumb trail is 'Clusters / my-cluster / Topics / test1'. The main content area is titled '+ test1' and is divided into several sections:

- Topic Summary:** A table with the following data:

Replication	1
Number of Partitions	5
Sum of partition offsets	3,000
Total number of Brokers	1
Number of Brokers for Topic	1
Preferred Replicas %	100
Brokers Skewed %	0
Brokers Leader Skewed %	0
Brokers Spread %	100
Under-replicated %	0
Leader Size	
- Operations:** A set of buttons including 'Delete Topic', 'Reassign Partitions', 'Generate Partition Assignments', 'Add Partitions', 'Update Config', and 'Manual Partition Assignments'.
- Partitions by Broker:** A table with the following data:

Broker	# of Partitions	# as Leader	Partitions	Skewed?	Leader Skewed?
1	5	5	(0,1,2,3,4)	false	false
- Consumers consuming from this topic:** A table with the following data:

Consumer Group	Consumer
group1	KF
- Metrics:** A table showing various metrics over time (Mean, 1 min, 5 min, 15 min). All values are 0.00.

Rate	Mean	1 min	5 min	15 min
Messages in /sec	0.00	0.00	0.00	0.00
Bytes in /sec	0.00	0.00	0.00	0.00
Bytes out /sec	0.00	0.00	0.00	0.00
Bytes rejected /sec	0.00	0.00	0.00	0.00
Failed fetch request /sec	0.00	0.00	0.00	0.00
Failed produce request /sec	0.00	0.00	0.00	0.00
- Partition Information:** A table with the following data:

Partition	Latest Offset	Leader	Replicas	In Sync Replicas	Preferred Leader?	Under Replicated?	Leader Size
0	1,500	1	(1)	(1)	true	false	
1	1,500	1	(1)	(1)	true	false	
2	0	1	(1)	(1)	true	false	
3	0	1	(1)	(1)	true	false	
4	0	1	(1)	(1)	true	false	

**步骤9**（可选）若对分配的分区不满意，可以执行Partition的重新分配功能来重新自动分配分区。

1. 在Topic Summary页面单击“Generate Partition Assignments”。
2. 勾选broker实例，单击“Generate Partition Assignments”生成分区。
3. 分区生成完成，单击“Go to topic view.”返回Topic Summary页面。
4. 在Topic Summary页面单击“Reassign Partitions”可以在集群的broker实例上重新自动分配分区。



5. 单击“Go to reassign partitions.”查看重新分配的分区详情。

**步骤10**（可选）若对自动分配的分区不满意，可以执行手动分配来重新分配分区。

1. 在Topic Summary页面单击“Manual Partition Assignments”进入手动分配分区页面。
2. 手动为每个分区的副本分配Broker id，然后单击“Save Partition Assignment”保存修改。
3. 单击“Go to topic view.”返回Topic Summary页面，查看分区详情。

----**结束**

# 17 使用 Loader

## 17.1 从零开始使用 Loader

用户可以使用Loader将数据从SFTP服务器导入到HDFS。

本章节适用于MRS 3.x之前版本。

### 前提条件

- 已准备业务数据。
- 已创建分析集群。

### 操作步骤

**步骤1** 访问Loader页面。

1. 登录集群详情页面，选择“服务管理”。
  2. 选择“Hue”，在“Hue概述”的“Hue WebUI”，单击“Hue (主)”，打开Hue的WebUI。
  3. 选择“Data Browsers > Sqoop”。
- 默认显示Loader页面中的作业管理界面。

**步骤2** 在Loader页面，单击“管理连接”。

**步骤3** 单击“新建连接”，参考[文件服务器连接](#)，创建sftp-connector。

**步骤4** 单击“新建连接”，输入连接名称，选择连接器为hdfs-connector，创建hdfs-connector。

**步骤5** 访问Loader页面，单击“管理作业”。

**步骤6** 单击“新建作业”。

**步骤7** 在“基本信息”填写参数。

1. 在“名称”填写一个作业的名称。
2. 选择[步骤3](#)创建的“源连接”和[步骤4](#)创建的“目的连接”。

**步骤8** 在“自”填写源连接的作业配置。

具体请参见[ftp-connector](#)或[sftp-connector](#)。

**步骤9** 在“至”填写目的连接的作业配置。

具体请参见[hdfs-connector](#)。

**步骤10** 在“任务配置”填写作业的运行参数。

表 17-1 Loader 作业运行属性

参数	说明
抽取并发数	设置map任务的个数。
加载(写入)并发数	设置reduce任务的个数。 该参数只有在目的字段为Hbase和Hive时才会显示。
单个分片的最大错误记录数	设置一个错误阈值，如果单个map任务的错误记录超过设置阈值则任务自动结束，已经获取的数据不回退。 <b>说明</b> “generic-jdbc-connector”的“MYSQL”和“MPPDB”默认批量读写数据，每一批次数据最多只记录一次错误记录。
脏数据目录	设置一个脏数据目录，在出现脏数据的场景中在该目录保存脏数据。如果不设置则不保存。

**步骤11** 单击“保存”。

----结束

## 17.2 Loader 使用简介

本章节适用于MRS 3.x之前版本。

### 使用流程

通过Loader迁移用户数据时，基本流程如下所示。

1. 访问Hue WebUI的Loader页面。
2. 管理Loader连接。
3. 创建作业，选择数据源的连接以及保存数据的连接。
4. 运行作业，完成数据迁移。

### Loader 页面介绍

Loader页面是基于开源Sqoop WebUI的图形化数据迁移管理工具，该页面托管在Hue的WebUI中。进入Loader页面请执行以下操作：

1. 访问Hue WebUI，参见[访问Hue WebUI界面](#)。
2. 选择“Data Browsers > Sqoop”。  
默认显示Loader页面中的作业管理界面。



## Loader 连接介绍

Loader连接保存了数据具体位置的相关信息，Loader使用连接来访问数据，或将数据保存到指定的位置。进入Loader连接管理页面请执行以下操作：

1. 进入Loader页面。
2. 单击“管理连接”。  
显示Loader连接管理页面。  
可单击“管理作业”回到作业管理页面。
3. 单击“新建连接”，进入配置页面，并填写参数创建一个Loader连接。

## Loader 作业介绍

Loader作业用于管理数据迁移任务，每个作业包含一个源数据的连接，和一个目的数据的连接，通过从源连接读取数据，再将数据保存到目的连接，完成数据迁移任务。

# 17.3 Loader 常用参数

本章节适用于MRS 3.x及后续版本。

## 参数入口

参数入口，请参考[修改集群服务配置参数](#)。

## 参数说明

表 17-2 Loader 常用参数

配置参数	说明	默认值	范围
mapreduce.client.submit.file.replication	MapReduce任务在运行时依赖的相关job文件在HDFS上的副本数。当集群中DataNode个数小于该参数值时，副本数等于DataNode的个数。当DataNode个数大于或等于该参数值，副本数为该参数值。	10	3 ~ 256

配置参数	说明	默认值	范围
loader.fault.tolerance.rate	容错率。 值大于0时使能容错机制。使能容错机制时建议将作业的Map数设置为大于等于3，推荐在作业数据量大的场景下使用。	0	0 ~ 1.0
loader.input.field.separator	默认的输入字段分割符，需要配置输入与输出转换步骤才生效，转换步骤的内容可以为空；如果作业的转换步骤中没有配置分割符，则以此处的默认分割符为准。	,	-
loader.input.line.separator	默认的输入行分割符，需要配置输入与输出转换步骤才生效，转换步骤的内容可以为空；如果作业的转换步骤中没有配置分割符，则以此处的默认分割符为准。	-	-
loader.output.field.separator	默认的输出字段分割符，需要配置输入与输出转换步骤才生效，转换步骤的内容可以为空；如果作业的转换步骤中没有配置分割符，则以此处的默认分割符为准。	,	-
loader.output.line.separator	Loader输出数据的行分隔符。	-	-

#### 📖 说明

- 由于容错率的统计需要时间，为保证使用效果，建议在作业运行时间在2分钟以上时使用“loader.fault.tolerance.rate”参数。
- 此处参数设置的为Loader全局的默认分割符，如果作业的转换步骤中配置了分割符，则以转换步骤为准，转换步骤中没有配置分割符则以此处的默认分割符为准。

## 17.4 创建 Loader 角色

### 操作场景

该任务指导MRS集群管理员在FusionInsight Manager创建并设置Loader的角色。Loader角色可设置Loader管理员权限、作业连接、作业分组以及Loader作业的操作和调度权限。

本章节适用于MRS 3.x及后续版本。

## 前提条件

- MRS集群管理员已明确业务需求。
- 已登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。

## 操作步骤

**步骤1** 选择“系统 > 权限 > 角色”。

**步骤2** 单击“添加角色”，然后“角色名称”和“描述”输入角色名字与描述。

**步骤3** 设置角色“权限”请参见[表17-3](#)。

### 说明

设置角色的权限时，不能同时选择跨资源权限，如果需要设置多个资源的相关权限，请依次逐一设置。

Loader权限：

- “管理员”：Loader管理员权限。
- “作业连接器”：Loader的连接权限。
- “作业分组”：Loader的作业分组操作权限。用户可以在指定作业分组下设置具体作业的操作权限，包括作业的编辑“编辑”与执行“执行”权限。
- “作业调度”：Loader的作业调度权限。

表 17-3 设置 Loader 角色

任务场景	角色授权操作
设置Loader管理员权限	在“配置资源权限”的表格中选择“待操作集群的名称 > Loader”，勾选“管理员”。
设置Loader的连接权限 (包括Job Connection的编辑、删除和引用权限)	<ol style="list-style-type: none"> <li>1. 在“配置资源权限”的表格中选择“待操作集群的名称 &gt; Loader &gt; 作业连接器”。</li> <li>2. 在指定作业连接的“权限”列，勾选“编辑”。</li> </ol>
设置Loader作业分组的编辑权限 (包括修改作业分组的名称、删除指定分组、在指定分组下创建作业的权限、从外部将作业批量导入到指定分组的权限、将其他分组的作业迁移到指定分组的权限)	<ol style="list-style-type: none"> <li>1. 在“权限”的表格中选择“Loader &gt; 作业分组”。</li> <li>2. 在指定作业分组的“权限”列，勾选“分组编辑”。</li> </ol>
设置Loader作业分组下所有作业的编辑权限 (包括对分组下现有或后续新增所有作业的编辑权限)	<ol style="list-style-type: none"> <li>1. 在“配置资源权限”的表格中选择“待操作集群的名称 &gt; Loader &gt; 作业分组”。</li> <li>2. 在指定作业分组的“权限”列，勾选“作业编辑”。</li> </ol>

任务场景	角色授权操作
设置Loader作业分组下所有作业的执行权限 (包括对分组下现有或后续新增所有作业的执行权限)	<ol style="list-style-type: none"> <li>在“配置资源权限”的表格中选择“待操作集群的名称 &gt; Loader &gt; 作业分组”。</li> <li>在指定作业分组的“权限”列，勾选“作业执行”。</li> </ol>
设置Loader作业的编辑权限 (包括作业的编辑、删除、复制和导出权限)	<ol style="list-style-type: none"> <li>在“配置资源权限”的表格中选择“待操作集群的名称 &gt; Loader &gt; 作业分组”。</li> <li>选择某个作业分组。</li> <li>在指定作业的“权限”列，勾选“编辑”。</li> </ol>
设置Loader作业的执行权限 (包括作业的启动、停止和查看历史记录权限)	<ol style="list-style-type: none"> <li>在“配置资源权限”的表格中选择“待操作集群的名称 &gt; Loader &gt; 作业分组”。</li> <li>选择某个作业分组。</li> <li>在指定作业的“权限”列，勾选“执行”。</li> </ol>
设置Loader作业调度的操作权限 (包括Scheduler的编辑、删除、是否生效权限)	<ol style="list-style-type: none"> <li>在“配置资源权限”的表格中选择“待操作集群的名称 &gt; Loader &gt; 作业调度”。</li> <li>在指定作业调度行的“权限”列，勾选“编辑”。</li> </ol>

### 📖 说明

- 除了“管理员”权限，以上权限只针对存量的资源信息进行权限配置。
- 未设置以上角色的用户也可以创建任务、分组、连接器，但是无法对存量的资源进行操作。

**步骤4** 单击“确定”完成，返回“角色”。

----结束

## 17.5 Loader 连接配置说明

本章节适用于MRS 3.x之前版本。

### 基本介绍

Loader支持以下多种连接，每种连接的配置介绍可根据本章节内容了解。

- obs-connector
- generic-jdbc-connector
- ftp-connector或sftp-connector

- hbase-connector、hdfs-connector或hive-connector

## OBS 连接

OBS连接是Loader与OBS进行数据交换的通道，配置参数如表17-4所示。

表 17-4 obs-connector 配置

参数	说明
名称	指定一个Loader连接的名称。
OBS服务器	输入OBS endpoint地址，一般格式为 <b>OBS.Region.DomainName</b> 。 例如执行如下命令查看OBS endpoint地址： <b>cat /opt/Bigdata/apache-tomcat-7.0.78/webapps/web/WEB-INF/classes/cloud-obs.properties</b>
端口	访问OBS数据的端口。默认值为“443”。
访问标识(AK)	表示访问OBS的用户的访问密钥AK。
密钥(SK)	表示访问密钥对应的SK。

## 关系型数据库连接

关系型数据库连接是Loader与关系型数据库进行数据交换的通道，配置参数如表17-5所示。

### 说明

部分参数需要单击“显示高级属性”后展开，否则默认隐藏。

表 17-5 generic-jdbc-connector 配置

参数	说明
名称	指定一个Loader连接的名称。
数据库类型	表示Loader连接支持的数据，可以选择“ORACLE”、“MYSQL”和“MPPDB”。
数据库服务器	表示数据库的访问地址，可以是IP地址或者域名。
端口	表示数据库的访问端口。
数据库名称	表示保存数据的具体数据库名。
用户名	表示连接数据库使用的用户名称。
密码	表示此用户对应的密码。需要与实际密码保持一致。



表 17-6 高级属性配置

参数	说明
一次请求行数	表示每次连接数据库时，最多可获取的数据量。
连接属性	不同类型数据库支持该数据库连接特有的驱动属性，例如 MySQL 的“autoReconnect”。如果需要定义驱动属性，单击“添加”。
引用符号	表示数据库的 SQL 中保留关键字的定界符，不同类型数据库定义的定界符不完全相同。

## 文件服务器连接

文件服务器连接包含 FTP 连接和 SFTP 连接，是 Loader 与文件服务器进行数据交换的通道，配置参数如表 17-7 所示。

表 17-7 ftp-connector 或 sftp-connector 配置

参数	说明
名称	指定一个 Loader 连接的名称。
主机名或 IP	输入文件服务器的访问地址，可以是服务器的主机名或者 IP 地址。
端口	访问文件服务器的端口。 <ul style="list-style-type: none"><li>FTP 协议请使用端口“21”。</li><li>SFTP 协议请使用端口“22”。</li></ul>
用户名	表示文件服务器的用户名称。
密码	表示此用户对应的密码。

## MRS 集群连接

MRS 集群连接包含 HBase 连接、HDFS 连接和 Hive 连接，是 Loader 与对应各数据进行数据交换的通道。

配置 MRS 集群连接时，需要设置名称、选择对应的连接器“hbase-connector”、“hdfs-connector”或“hive-connector”，然后保存即可。

# 17.6 管理 Loader 连接（MRS 3.x 之前版本）

## 操作场景

Loader 页面支持创建、查看、编辑和删除连接。

本章节适用于 MRS 3.x 之前版本。

## 前提条件

已访问Loader页面，参见[Loader页面介绍](#)。

## 创建连接

**步骤1** 在Loader页面，单击“管理连接”。

**步骤2** 单击“新建连接”，配置连接参数。

参数介绍具体可参见[Loader连接配置说明](#)。

**步骤3** 单击“保存”。

如果连接配置，例如IP地址、端口、访问用户等信息不正确，将导致验证连接失败无法保存。另外VPC相关设置，也可能影响网络连通性。

### 说明

用户可以直接单击“测试”立即检测连接是否可用。

----结束

## 查看连接

**步骤1** 在Loader页面，单击“管理连接”。

- 如果集群启用了Kerberos认证，则默认显示所有当前用户创建的连接，不支持显示其他用户创建的连接。
- 如果集群未启用Kerberos认证，则显示集群中全部的Loader连接。

**步骤2** 在“Sqoop连接”中输入指定连接的名称，可以筛选该连接。

----结束

## 编辑连接

**步骤1** 在Loader页面，单击“管理连接”。

**步骤2** 单击指定连接的名称，进入编辑页面。

**步骤3** 根据业务需要，修改连接配置参数。

**步骤4** 单击“测试”。

如果显示测试成功，则执行**步骤5**；如果显示不能连接至OBS Server，则需要重复**步骤3**。

**步骤5** 单击“保存”。

如果某个Loader作业已集成一个Loader连接，那么编辑连接参数后可能导致Loader作业运行效果也产生变化。

----结束

## 删除连接

**步骤1** 在Loader页面，单击“管理连接”。

**步骤2** 在指定连接所在行，单击“删除”。

**步骤3** 在弹出的对话框窗口，单击“是，将其删除”。

如果某个Loader作业已集成一个Loader连接，那么该连接不可以被删除。

----结束

## 17.7 管理 Loader 连接（MRS 3.x 及之后版本）

### 操作场景

Loader页面支持创建、查看、编辑和删除连接。

本章节适用于MRS 3.x及之后版本。

### 创建连接

**步骤1** 登录服务页面：

MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)，选择“集群 > 服务”。

**步骤2** 选择“Loader”，在“Loader WebUI”右侧，单击链接，打开Loader的WebUI。

**步骤3** 在Loader页面，单击“新建作业”。

**步骤4** 在“连接”后单击“添加”，配置连接参数。

参数介绍具体可参见[Loader连接配置说明](#)。

**步骤5** 单击“确定”。

如果连接配置，例如IP地址、端口、访问用户等信息不正确，将导致验证连接失败无法保存。

#### 说明

用户可以直接单击“测试”立即检测连接是否可用。

----结束

### 查看连接

**步骤1** 在Loader页面，单击“新建作业”。

**步骤2** 单击“连接”后的下拉列表框，可以查看已创建的连接。

----结束

### 编辑连接

**步骤1** 在Loader页面，单击“新建作业”。

**步骤2** 单击“连接”后的下拉列表框，选择待编辑的连接名称。

**步骤3** 在“连接”后单击“编辑”，进入编辑页面。

**步骤4** 根据业务需要，修改连接配置参数。

**步骤5** 单击“测试”。

- 如果显示测试成功，则执行**步骤6**。
- 如果显示测试失败，则需要重复**步骤4**。

**步骤6** 单击“保存”。

如果某个Loader作业已集成一个Loader连接，那么编辑连接参数后可能导致Loader作业运行效果也产生变化。

----结束

## 删除连接

**步骤1** 在Loader页面，单击“新建作业”。

**步骤2** 单击“连接”后的下拉列表框，选择待删除的连接名称。

**步骤3** 单击“删除”。

**步骤4** 在弹出的对话框窗口，单击“确定”。

如果某个Loader作业已集成一个Loader连接，那么该连接不可以被删除。

----结束

## Loader 连接配置说明

Loader支持以下多种连接：

- generic-jdbc-connector：参数配置请参见**表17-8**。
- ftp-connector：参数配置请参见**表17-9**。
- sftp-connector：参数配置请参见**表17-10**。
- hdfs-connector：参数配置请参见**表17-11**。
- oracle-connector：参数配置请参见**表17-12**。
- mysql-fastpath-connector：参数配置请参见**表17-14**。
- oracle-partition-connector：参数配置请参见**表17-13**。

表 17-8 generic-jdbc-connector 配置

参数	说明
名称	给定一个Loader连接的名称。
连接器	选择“generic-jdbc-connector”。

参数	说明
JDBC驱动程序类	JDBC驱动类如下： <ul style="list-style-type: none"> <li>• oracle: oracle.jdbc.driver.OracleDriver</li> <li>• SQLServer: com.microsoft.sqlserver.jdbc.SQLServerDriver</li> <li>• mysql: com.mysql.jdbc.Driver</li> <li>• postgresql: org.postgresql.Driver</li> <li>• gaussdb200: com.huawei.gauss200.jdbc.Driver</li> </ul>
JDBC连接字符串	表示数据库的访问地址，可以是IP地址或者域名。 输入数据库连接字符串（以下以IP为10.10.10.10，样例数据库为“test”为例）： <ul style="list-style-type: none"> <li>• oracle: jdbc:oracle:thin:@10.10.10.10:1521:orcl</li> <li>• SQLServer: jdbc:sqlserver://10.10.10.10:1433;DatabaseName=test</li> <li>• mysql: jdbc:mysql://10.10.10.10/test?&amp;useUnicode=true&amp;characterEncoding=GBK</li> <li>• postgresql: jdbc:postgresql://10.10.10.10:5432/test</li> <li>• gaussdb200: jdbc:gaussdb://10.10.10.10:15400/test（15400为样例端口）</li> </ul>
用户名	表示连接数据库使用的用户名称。
密码	表示此用户对应的密码。需要与实际密码保持一致。

表 17-9 ftp-connector 配置

参数	说明
名称	指定一个Loader连接的名称。
连接器	选择“ftp-connector”。
FTP模式	选择“ACTIVE”或者“PASSIVE”。
FTP协议	选择： <ul style="list-style-type: none"> <li>• FTP</li> <li>• SSL_EXPLICIT</li> <li>• SSL_IMPLICIT</li> <li>• TLS_EXPLICIT</li> <li>• TLS_IMPLICIT</li> </ul>
文件名编码类型	文件名或者文件路径名的编码类型。

表 17-10 sftp-connector 配置

参数	说明
名称	指定一个Loader连接的名称。
连接器	选择“sftp-connector”。

表 17-11 hdfs-connector 配置

参数	说明
名称	指定一个Loader连接的名称。
连接器	选择“hdfs-connector”。

表 17-12 oracle-connector 配置

参数	说明
名称	指定一个Loader连接的名称。
连接器	选择“oracle-connector”。
JDBC连接字符串	输入用于连接数据库的连接串，例如“jdbc:oracle:thin:@IP.port.database”。
用户名	表示连接数据库使用的用户名称。
密码	表示此用户对应的密码。需要与实际密码保持一致。

表 17-13 oracle-partition-connector 配置

参数	说明
名称	指定一个Loader连接的名称。
连接器	选择“oracle-partition-connector”。
JDBC驱动程序类	输入“oracle.jdbc.driver.OracleDriver”。
JDBC连接字符串	输入用于连接数据库的连接串，例如“jdbc:oracle:thin:@IP.port.database”。
用户名	表示连接数据库使用的用户名称。
密码	表示此用户对应的密码。需要与实际密码保持一致。

表 17-14 mysql-fastpath-connector 配置

参数	说明
名称	指定一个Loader连接的名称。
连接器	<p>选择“mysql-fastpath-connector”。</p> <p><b>须知</b> 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlimport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<a href="http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html">http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html</a>，安装MySQL client applications and tools。</p> <p>例如：在RHEL-x86系统上需要安装如下RPM包（请根据实际情况选择版本）</p> <ul style="list-style-type: none"> <li>mysql-community-client-5.7.23-1.el7.x86_64.rpm</li> <li>mysql-community-common-5.7.23-1.el7.x86_64.rpm</li> <li>mysql-community-devel-5.7.23-1.el7.x86_64.rpm</li> <li>mysql-community-embedded-5.7.23-1.el7.x86_64.rpm</li> <li>mysql-community-libs-5.7.23-1.el7.x86_64.rpm</li> <li>mysql-community-libs-compat-5.7.23-1.el7.x86_64.rpm</li> </ul>
JDBC连接字符串	输入用于连接数据库的连接串，例如“jdbc:mysql://IP/database?&useUnicode=true&characterEncoding=GBK”。
用户名	表示连接数据库使用的用户名称。
密码	表示此用户对应的密码。需要与实际密码保持一致。

## 17.8 Loader 作业源连接配置说明

### 基本介绍

Loader作业需要从不同数据源获取数据时，应该选择对应类型的连接，每种连接在该场景中需要配置连接的属性。

本章节适用于MRS 3.x之前版本。

### obs-connector

表 17-15 obs-connector 数据源连接属性

参数	说明
桶名	保存源数据的OBS文件系统。
源目录或文件	源数据实际存储的形态，可能是文件系统包含一个目录中的全部数据文件，或者是文件系统包含的单个数据文件。

参数	说明
文件格式	Loader支持OBS中存储数据的文件格式，默认支持以下两种： <ul style="list-style-type: none"> <li>• CSV_FILE：表示文本格式文件。目的连接为数据库型连接时，只支持文本格式。</li> <li>• BINARY_FILE：表示文本格式以外的二进制文件。</li> </ul>
换行符	源数据的每行结束标识字符。
字段分割符	源数据的每个字段分割标识字符。
编码类型	源数据的文本编码类型。只对文本类型文件有效。
文件分割方式	支持以下两种： <ul style="list-style-type: none"> <li>• File：按总文件个数分配map任务处理的文件数量，计算规则为“文件总个数/抽取并发数”。</li> <li>• Size：按文件总大小分配map任务处理的文件大小，计算规则为“文件总大小/抽取并发数”。</li> </ul>

## generic-jdbc-connector

表 17-16 generic-jdbc-connector 数据源连接属性

参数	说明
模式或表空间	表示源数据对应的数据库名称，支持通过界面查询并选择。
表名	存储源数据的数据表，支持通过界面查询并选择。
抽取分区字段	分区字段，如果需读取多个字段，使用该字段分割结果并获取数据。
Where子句	表示读取数据库时使用的查询语句。

## ftp-connector 或 sftp-connector

表 17-17 ftp-connector 或 sftp-connector 数据源连接属性

参数	说明
源目录或文件	源数据实际存储的形态，可能是文件服务器包含一个目录中的全部数据文件，或者是单个数据文件。
文件格式	Loader支持文件服务器中存储数据的文件格式，默认支持以下两种： <ul style="list-style-type: none"> <li>• CSV_FILE：表示文本格式文件。目的连接为数据库型连接时，只支持文本格式。</li> <li>• BINARY_FILE：表示文本格式以外的二进制文件。</li> </ul>



参数	说明
换行符	源数据的每行结束标识字符。 <b>说明</b> ftp或sftp作为源连接时，当“文件格式”配置为BINARY_FILE时，高级属性中的“换行符”配置无效。
字段分割符	源数据的每个字段分割标识字符。 <b>说明</b> ftp或sftp作为源连接时，当“文件格式”配置为BINARY_FILE时，高级属性中的“字段分割符”配置无效
编码类型	源数据的文本编码类型。只对文本类型文件有效。
文件分割方式	支持以下两种： <ul style="list-style-type: none"> <li>• File：按总文件个数分配map任务处理的文件数量，计算规则为“文件总个数/抽取并发数”。</li> <li>• Size：按文件总大小分配map任务处理的文件大小，计算规则为“文件总大小/抽取并发数”。</li> </ul>

## hbase-connector

表 17-18 hbase-connector 数据源连接属性

参数	说明
表名	源数据实际存储的HBase表。

## hdfs-connector

表 17-19 hdfs-connector 数据源连接属性

参数	说明
源目录或文件	源数据实际存储的形态，可能是HDFS包含一个目录中的全部数据文件，或者是单个数据文件。
文件格式	Loader支持HDFS中存储数据的文件格式，默认支持以下两种： <ul style="list-style-type: none"> <li>• CSV_FILE：表示文本格式文件。目的连接为数据库型连接时，只支持文本格式。</li> <li>• BINARY_FILE：表示文本格式以外的二进制文件。</li> </ul>
换行符	源数据的每行结束标识字符。 <b>说明</b> hdfs作为源连接时，当“文件格式”配置为BINARY_FILE时，高级属性中的“换行符”配置无效。

参数	说明
字段分割符	源数据的每个字段分割标识字符。 <b>说明</b> hdfs作为源连接时，当“文件格式”配置为BINARY_FILE时，高级属性中的“字段分割符”配置无效。
文件分割方式	支持以下两种： <ul style="list-style-type: none"><li>• File：按总文件个数分配map任务处理的文件数量，计算规则为“文件总个数/抽取并发数”。</li><li>• Size：按文件总大小分配map任务处理的文件大小，计算规则为“文件总大小/抽取并发数”。</li></ul>

## hive-connector

表 17-20 hive-connector 数据源连接属性

参数	说明
数据库名称	数据源的Hive数据库名称，支持通过界面查询并选择。
表名	数据源的Hive表名称，支持通过界面查询并选择。

## 17.9 Loader 作业目的连接配置说明

### 基本介绍

Loader作业需要将数据保存到不同目的存储位置时，应该选择对应类型的目的连接，每种连接在该场景中需要配置连接的属性。

### obs-connector

表 17-21 obs-connector 目的连接属性

参数	说明
桶名	保存最终数据的OBS文件系统。
写入目录	最终数据在文件系统保存时的具体目录。必须指定一个目录。
文件格式	Loader支持OBS中存储数据的文件格式，默认支持以下两种： <ul style="list-style-type: none"><li>• CSV_FILE：表示文本格式文件。目的连接为数据库型连接时，只支持文本格式。</li><li>• BINARY_FILE：表示文本格式以外的二进制文件。</li></ul>
换行符	最终数据的每行结束标识字符。
字段分割符	最终数据的每个字段分割标识字符。

参数	说明
编码类型	最终数据的文本编码类型。只对文本类型文件有效。

## generic-jdbc-connector

表 17-22 generic-jdbc-connector 目的连接属性

参数	说明
模式名称	保存最终数据的数据库名称。
表名	保存最终数据的数据表名称。

## ftp-connector 或 sftp-connector

表 17-23 ftp-connector 或 sftp-connector 目的连接属性

参数	说明
写入目录	最终数据在文件服务器保存时的具体目录。必须指定一个目录。
文件格式	Loader支持文件服务器中存储数据的文件格式，默认支持以下两种： <ul style="list-style-type: none"> <li>• CSV_FILE：表示文本格式文件。目的连接为数据库型连接时，只支持文本格式。</li> <li>• BINARY_FILE：表示文本格式以外的二进制文件。</li> </ul>
换行符	最终数据的每行结束标识字符。 <b>说明</b> ftp或sftp作为目的连接时，当“文件格式”配置为BINARY_FILE时，高级属性中的“换行符”配置无效。
字段分割符	最终数据的每个字段分割标识字符。 <b>说明</b> ftp或sftp作为目的连接时，当“文件格式”配置为BINARY_FILE时，高级属性中的“字段分割符”配置无效
编码类型	最终数据的文本编码类型。只对文本类型文件有效。

## hbase-connector

表 17-24 hbase-connector 目的连接属性

参数	说明
表名	保存最终数据的HBase表名称，支持通过界面查询并选择。
导入方式	支持BULKLOAD、PUTLIST两种方式导入数据到HBase表。
导入前清空数据	标识是否需要清空目标HBase表中的数据，支持以下两种类型： <ul style="list-style-type: none"><li>• True：清空表中的数据。</li><li>• False：不清空表中的数据，选择False时如果表中存在数据，则作业运行会报错。</li></ul>

## hdfs-connector

表 17-25 hdfs-connector 目的连接属性

参数	说明
写入目录	最终数据在HDFS保存时的具体目录。必须指定一个目录。
文件格式	Loader支持HDFS中存储数据的文件格式，默认支持以下两种： <ul style="list-style-type: none"><li>• CSV_FILE：表示文本格式文件。目的连接为数据库型连接时，只支持文本格式。</li><li>• BINARY_FILE：表示文本格式以外的二进制文件。</li></ul>
压缩格式	文件在HDFS保存时的压缩行为。支持NONE、DEFLATE、GZIP、BZIP2、LZ4和SNAPPY。
是否覆盖	文件在导入HDFS时对写入目录中原有文件的处理行为，支持以下两种： <ul style="list-style-type: none"><li>• True：默认清空目录中的文件并导入新文件。</li><li>• False：不清空文件。如果写入目录中有文件，则作业运行失败。</li></ul>
换行符	最终数据的每行结束标识字符。 <b>说明</b> hdfs作为目的连接时，当“文件格式”配置为BINARY_FILE时，高级属性中的“换行符”配置无效。
字段分割符	最终数据的每个字段分割标识字符。 <b>说明</b> hdfs作为目的连接时，当“文件格式”配置为BINARY_FILE时，高级属性中的“字段分割符”配置无效。

## hive-connector

表 17-26 hive-connector 目的连接属性

参数	说明
数据库名称	保存最终数据的Hive数据库名称，支持通过界面查询并选择。
表名	保存最终数据的Hive表名称，支持通过界面查询并选择。

## 17.10 管理 Loader 作业

### 操作场景

Loader页面支持创建、查看、编辑和删除作业。

本章节适用于MRS 3.x之前版本。

### 前提条件

已访问Loader页面，参见[Loader页面介绍](#)。

### 创建作业

**步骤1** 访问Loader页面，单击“新建作业”。

**步骤2** 在“基本信息”填写参数。

1. 在“名称”填写一个作业的名称。
2. 在“源连接”和“目的连接”选择对应的连接。  
选择某个类型的连接，表示从指定的源获取数据，并保存到目的位置。

#### 说明

如果没有需要的连接，可单击“添加新连接”。

**步骤3** 在“自”填写源连接的作业配置。

具体请参见[Loader作业源连接配置说明](#)。

**步骤4** 在“至”填写目的连接的作业配置。

具体请参见[Loader作业目的连接配置说明](#)。

**步骤5** 在“目的连接”是否选择了数据库类型的连接。

数据库类型的连接包含以下几种：

- generic-jdbc-connector
- hbase-connector
- hive-connector

“目的连接”选择数据库类型的连接时，还需要配置业务数据与数据库表字段的对应关系：

- 是，请执行**步骤6**。
- 否，请执行**步骤7**。

**步骤6** 在“字段映射”填写字段对应关系。然后执行**步骤7**。

“字段映射”的对应关系，表示用户数据中每一列与数据库的表字段的匹配关系。

**表 17-27 “字段映射”属性**

参数	说明
列号	表示业务数据的字段顺序。
样本	表示业务数据的第一行值样例。
列族	“目的连接”为hbase-connector类型时，支持定义保存数据的具体列族。
目的字段	配置保存数据的具体字段。
类型	显示用户选择字的类型。
行键	“目的连接”为hbase-connector类型时，需要勾选作为行键的“目的字段”。

#### 说明

如果From是sftp/ftp/obs/hdfs等文件类型连接器，Field Mapping 样值取自文件第一行数据，需要保证第一行数据是完整的，Loader作业不会抽取没有Mapping上的列。

**步骤7** 在“任务配置”填写作业的运行参数。

**表 17-28 Loader 作业运行属性**

参数	说明
抽取并发数	设置map任务的个数。
加载(写入)并发数	设置reduce任务的个数。 该参数只有在目的字段为Hbase和Hive时才会显示。
单个分片的最大错误记录数	设置一个错误阈值，如果单个map任务的错误记录超过设置阈值则任务自动结束，已经获取的数据不回退。 <b>说明</b> “generic-jdbc-connector”的“MYSQL”和“MPPDB”默认批量读写数据，每一批次数据最多只记录一次错误记录。
脏数据目录	设置一个脏数据目录，在出现脏数据的场景中在该目录保存脏数据。如果不设置则不保存。

**步骤8** 单击“保存”。

----结束

## 查看作业

- 步骤1** 访问Loader页面，默认显示Loader作业管理页面。
- 如果集群启用了Kerberos认证，则默认显示所有当前用户创建的作业，不支持显示其他用户的作业。
  - 如果集群未启用Kerberos认证，则显示集群中全部的作业。
- 步骤2** 在“Sqoop作业”中输入指定作业的名称或连接类型，可以筛选该作业。
- 步骤3** 单击“刷新列表”，可以获取作业的最新状态。
- 结束

## 编辑作业

- 步骤1** 访问Loader页面，默认显示Loader作业管理页面。
- 步骤2** 单击指定作业的名称，进入编辑页面。
- 步骤3** 根据业务需要，修改作业配置参数。
- 步骤4** 单击“保存”。

### 说明

左侧导航栏支持作业的基本操作，包含“运行”、“复制”、“删除”、“激活”、“历史记录”和“显示作业JSON定义”。

----结束

## 删除作业

- 步骤1** 访问Loader页面。
- 步骤2** 在指定作业所在行，单击✕。
- 您还可以勾选一个或多个作业，单击作业列表右上方的“删除作业”。
- 步骤3** 在弹出的对话框窗口，单击“是，将其删除”。
- 如果某个Loader作业正处于“运行中”的状态，则无法删除作业。
- 结束

# 17.11 准备 MySQL 数据库连接的驱动

## 操作场景

Loader作为批量数据导出的组件，可以通过关系型数据库导入、导出数据。

## 前提条件

已准备业务数据。

## 操作步骤

### MRS 3.x之前版本：

**步骤1** 从MySQL官网下载MySQL jdbc驱动程序“mysql-connector-java-5.1.21.jar”，具体MySQL jdbc驱动程序选择参见下表。

表 17-29 版本信息

jdbc驱动程序版本	MySQL版本
Connector/J 5.1	MySQL 4.1、MySQL 5.0、MySQL 5.1、MySQL 6.0 alpha
Connector/J 5.0	MySQL 4.1、MySQL 5.0 servers、distributed transaction (XA)
Connector/J 3.1	MySQL 4.1、MySQL 5.0 servers、MySQL 5.0 except distributed transaction (XA)
Connector/J 3.0	MySQL 3.x、MySQL 4.1

**步骤2** 将“mysql-connector-java-5.1.21.jar”上传至MRS master 主备节点loader安装目录。

- 针对MRS 3.x之前版本，上传至“/opt/Bigdata/MRS\_XXX/install/FusionInsight-Sqoop-1.99.7/FusionInsight-Sqoop-1.99.7/server/jdbc/”  
其中“XXX”为MRS版本号，请根据实际情况修改。

**步骤3** 修改“mysql-connector-java-5.1.21.jar”包属主为“omm:wheel”。

**步骤4** 修改配置文件“jdbc.properties”。

将“MYSQL”的键值修改为上传的jdbc驱动包名“mysql-connector-java-5.1.21.jar”，例如：MYSQL=mysql-connector-java-5.1.21.jar。

**步骤5** 重启Loader服务。

----结束

### MRS 3.x及之后版本：

修改关系型数据库对应的驱动jar包文件权限。

**步骤1** 登录Loader服务的主备管理节点，获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：“\${BIGDATA\_HOME}/FusionInsight\_Porter\_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib”。

#### 说明

此处版本号8.1.0.1为示例，具体以实际环境的版本号为准。

**步骤2** 使用root用户在Loader服务主备节点分别执行以下命令修改权限：

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib
```



`chown omm:wheel jar包文件名`

`chmod 600 jar包文件名`

**步骤3** 登录FusionInsight Manager系统，选择“集群 > 待操作集群名称 > 服务 > Loader > 更多 > 重启服务”输入管理员密码重启Loader服务。

----结束

## 17.12 数据导入

### 17.12.1 概述

“数据导入”章节适用于MRS 3.x及后续版本。

#### 简介

Loader是实现MRS与外部数据源如关系型数据库、SFTP服务器、FTP服务器之间交换数据和文件的ETL工具，支持将数据或文件从关系型数据库或文件系统导入到MRS服务中。

Loader支持如下数据导入方式：

- 从关系型数据库导入数据到HDFS/OBS。
- 从关系型数据库导入数据到HBase。
- 从关系型数据库导入数据到Phoenix表。
- 从关系型数据库导入数据到Hive表。
- 从SFTP服务器导入数据到HDFS/OBS。
- 从SFTP服务器导入数据到HBase。
- 从SFTP服务器导入数据到Phoenix表。
- 从SFTP服务器导入数据到Hive表。
- 从FTP服务器导入数据到HDFS/OBS。
- 从FTP服务器导入数据到HBase。
- 从FTP服务器导入数据到Phoenix表。
- 从FTP服务器导入数据到Hive表。
- 从同一集群内HDFS/OBS导入数据到HBase。

MRS与外部数据源交换数据和文件时需要连接数据源。系统提供以下连接器，用于配置不同类型数据源的连接参数：

- `generic-jdbc-connector`：关系型数据库连接器。
- `ftp-connector`：FTP数据源连接器。
- `hdfs-connector`：HDFS数据源连接器。
- `oracle-connector`：Oracle数据库专用连接器，使用`row_id`作为分区列，相对`generic-jdbc-connector`来说，Map任务分区更均匀，并且不依赖分区列是否有创建索引。

- `mysql-fastpath-connector`: MySQL数据库专用连接器，使用MySQL的`mysqldump`和`mysqlimport`工具进行数据的导入导出，相对`generic-jdbc-connector`来说，导入导出速度更快。
- `sftp-connector`: SFTP数据源连接器。
- `oracle-partition-connector`: 支持Oracle分区特性的连接器，专门对Oracle分区表的导入导出进行优化。

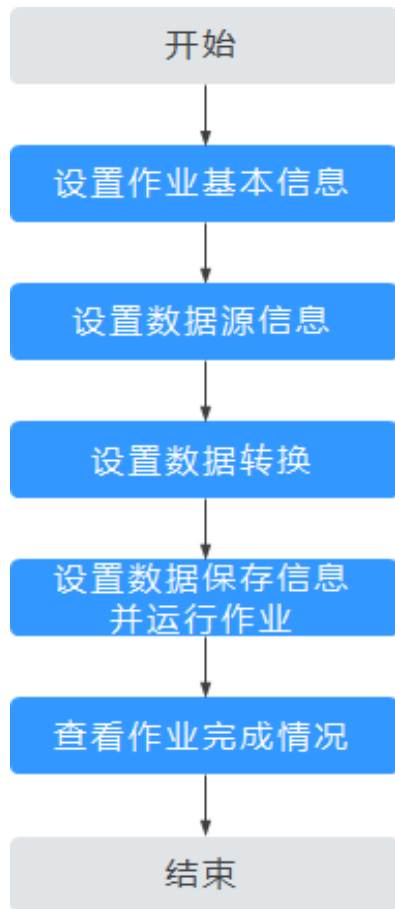
#### 📖 说明

- 使用FTP数据源连接器时不加密数据，可能存在安全风险，建议使用SFTP数据源连接器。
- 建议将SFTP服务器、FTP服务器和数据库服务器与Loader部署在独立的子网中，以保障数据安全地导入。
- 与关系数据库连接时，可以选择通用数据库连接器（`generic-jdbc-connector`）或者专用数据库连接器（`oracle-connector`、`oracle-partition-connector`、`mysql-fastpath-connector`），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用`mysql-fastpath-connector`时，要求在NodeManager节点上有MySQL的`mysqldump`和`mysqlimport`命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。
- 使用`oracle-connector`时，要求给连接用户赋予如下系统表或者视图的select权限：  
`dba_tab_partitions`、`dba_constraints`、`dba_tables`、`dba_segments`、`v$instance`、`dba_objects`、`v$instance`、`SYS_CONTEXT`函数、`dba_extents`、`dba_tab_subpartitions`。
- 使用`oracle-partition-connector`时，要求给连接用户赋予如下系统表的select权限：  
`dba_objects`、`dba_extents`。

## 导入流程

用户通过Loader界面进行数据导入作业，导入流程如[图17-1](#)所示。

图 17-1 导入流程示意



用户也可以通过shell脚本来更新与运行Loader作业，该方式需要对已安装的Loader客户端进行配置。

## 17.12.2 使用 Loader 导入数据

### 操作场景

该任务指导用户完成将数据从外部的数据源导入到MRS的工作。

一般情况下，用户可以手工在Loader界面管理数据导入导出作业。当用户需要通过shell脚本来更新与运行Loader作业时，必须对已安装的Loader客户端进行配置。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HDFS/OBS目录、HBase表和数据。
- 获取外部数据源（SFTP服务器或关系型数据库）使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 使用Loader从SFTP、FTP和HDFS/OBS导入数据时，确保外部数据源的输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符/"/":;,中的任意字符。
- 如果设置的任务需要使用指定Yarn队列功能，该用户需要已授权有相关Yarn队列的权限。

- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

## 操作步骤

**步骤1** 是否第一次从MRS导入数据到关系型数据库：

- 是，执行**步骤2**。
- 否，执行**步骤3**。

**步骤2** 修改关系型数据库对应的驱动jar包文件权限。

1. 登录Loader服务的主备管理节点，获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：“\${BIGDATA\_HOME}/FusionInsight\_Porter\_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib”。

### 📖 说明

此处版本号8.1.0.1为示例，具体以实际环境的版本号为准。

2. 使用root用户在Loader服务主备节点分别执行以下命令修改权限：

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib
```

```
chown omm:wheel jar包文件名
```

```
chmod 600 jar包文件名
```

3. 登录FusionInsight Manager系统，选择“集群 > 待操作集群名称 > 服务 > Loader > 更多 > 重启服务”输入管理员密码重启Loader服务。

**步骤3** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-2 Loader WebUI 界面



**步骤4** 创建Loader数据导入作业，单击“新建作业”，在“1.基本信息”选择所需要的作业类型，然后单击“下一步”。

1. “名称”输入作业的名称，“类型”选择“导入”。
2. “连接”选择一个连接。默认没有已创建的连接，单击“添加”创建一个新的连接，完成后单击“测试”，测试是否可用，待提示成功后单击“确定”。

MRS与外部数据源交换数据和文件时需要连接数据源，“连接”表示连接数据源时的连接参数集合。

表 17-30 连接配置参数一览表

连接器类型	参数名	说明
generic-jdbc-connector	JDBC驱动程序类	JDBC驱动类名。
	JDBC连接字符串	JDBC连接字符串。
	用户名	连接数据库使用的用户名。
	密码	连接数据库使用的密码。
	JDBC连接属性	JDBC连接属性，单击“添加”手动添加。 - 名称：连接属性名。 - 值：连接属性值。
ftp-connector	FTP服务器的IP	FTP服务器的IP地址。
	FTP服务器端口	FTP服务器的端口号。
	FTP用户名	访问FTP服务器的用户名。
	FTP密码	访问FTP服务器的密码。
	FTP模式	设置FTP访问模式，“ACTIVE”表示主动模式，“PASSIVE”表示被动模式。不指定参数值，默认为被动模式。
	FTP协议	设置FTP传输协议： - “FTP”：FTP协议。 - “SSL_EXPLICIT”：显式SSL协议。 - “SSL_IMPLICIT”：隐式SSL协议。 - “TLS_EXPLICIT”：显式TLS协议。 - “TLS_IMPLICIT”：隐式TLS协议。 不指定参数值，默认为FTP协议。
	文件名编码类型	填写FTP服务器支持的文件名、文件路径编码格式，不填写时使用系统默认格式“UTF-8”。
hdfs-connector	-	-
oracle-connector	JDBC连接字符串	用户连接数据库的连接字符串。
	用户名	连接数据库使用的用户名。
	密码	连接数据库使用的密码。
	连接属性	连接属性，单击“添加”手动添加。 - 名称：连接属性名。 - 值：连接属性值。

连接器类型	参数名	说明
mysql-fastpath-connector	JDBC连接字符串	JDBC连接字符串。
	用户名	连接数据库使用的用户名。
	密码	连接数据库使用的密码。
	连接属性	连接属性，单击“添加”手动添加。 - 名称：连接属性名。 - 值：连接属性值。
sftp-connector	Sftp服务器的IP	SFTP服务器的IP地址。
	Sftp服务器端口	SFTP服务器的端口号。
	Sftp用户名	访问SFTP服务器的用户名。
	Sftp密码	访问SFTP服务器的密码。
	Sftp公钥	Sftp服务器公钥。
oracle-partition-connector	JDBC驱动程序类	JDBC驱动类名。
	JDBC连接字符串	JDBC连接字符串。
	用户名	连接数据库使用的用户名。
	密码	连接数据库使用的密码。
	连接属性	连接属性，单击“添加”手动添加。 - 名称：连接属性名。 - 值：连接属性值。

- “组”设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，单击“确定”保存。
- “队列”设置Loader的任务在指定的Yarn队列中执行。默认值“root.default”表示任务在“default”队列中执行。
- “优先级”设置Loader的任务在指定的Yarn队列中的优先级。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。默认值为“NORMAL”。

**步骤5** 在“2.输入设置”，设置数据来源，然后单击“下一步”。

#### 说明

- 创建或者编辑Loader作业时，在配置SFTP路径、HDFS/OBS路径、SQL的Where条件等参数时，可以使用宏定义，具体请参考[配置项中使用宏定义](#)章节。
- Loader支持常见的字段数据类型，如Char、VarChar、Boolean、Binary、SmallInt、Int、BigInt、Decimal、Float、Double、Date、Time、TimeStamp、String等，具体支持类型根据数据来源的不同可能会有所变化，具体支持的类型可以参考Loader界面中相应输入算子（如表输入等）的字段数据类型下拉框中的内容。一些数据库的特有字段可能不被支持，例如Loader不支持oracle中的CLOB和XMLType、BLOB字段。

表 17-31 输入配置参数一览表

源文件类型	参数名	解释说明
sftp-connector或ftp-connector	输入路径	SFTP服务器中源文件的输入路径，如果连接器配置多个地址此处可对应使用分号分隔多个输入路径，数量需要与连接器中服务器的数量一致。
	文件分割方式	选择按文件或大小分割源文件，作为数据导入的MapReduce任务中各个map的输入文件。 <ul style="list-style-type: none"> <li>选择“FILE”表示每个map处理1个或多个完整的源文件，同一个源文件不可分配至不同map，数据保存至输出目录时将保留输入路径的目录结构</li> <li>选择“SIZE”表示每个map处理一定大小的输入文件，同一个源文件可分割至多个map，数据保存至输出目录时保存的文件数与map数量相同，文件名格式为“import_part_xxxx”，“xxxx”为系统生成的随机数，具有唯一性。</li> </ul>
	过滤器类型	选择文件过滤的条件。“WILCARD”表示使用通配符过滤，“REGEX”表示使用正则表达式匹配。与“路径过滤器”和“文件过滤器”配合使用。不选择值时默认为通配符过滤。
	路径过滤器	与“过滤器类型”配合使用，配置通配符或正则表达式对源文件的输入路径包含的目录进行过滤。输入路径“输入路径”不参与过滤。配置多个过滤条件时使用逗号隔开，配置为空时表示不过滤目录。
	文件过滤器	与“过滤器类型”配合使用，配置通配符或正则表达式对源文件的输入文件名进行过滤。配置多个过滤条件时使用逗号隔开。不能配置为空。
	编码类型	源文件的编码格式，如UTF-8。导入文本文件时才能配置。
	后缀名	源文件导入成功后对输入文件增加的后缀值。该值为空，表示不加后缀。
	压缩	使用SFTP协议导出数据时，是否开启压缩传输功能以减小带宽使用。“true”为开启压缩，“false”为关闭压缩。
	hdfs-connector	输入路径
路径过滤器		配置通配符对源文件的输入路径包含的目录进行过滤。输入路径不参与过滤。配置多个过滤条件时使用逗号隔开，配置为空时表示不过滤目录。不支持正则表达式过滤。
文件过滤器		配置通配符对源文件的输入文件名进行过滤。配置多个过滤条件时使用逗号隔开。不能配置为空。不支持正则表达式过滤。

源文件类型	参数名	解释说明
	编码类型	源文件的编码格式，如UTF-8。导入文本文件时才能配置。
	后缀名	源文件导入成功后对输入文件增加的后缀值。该值为空，表示不加后缀。
generic-jdbc-connector	架构名称	“表方式”模式下存在，数据库模式名。
	表名	“表方式”模式下存在，数据库表名。
	SQL语句	“SQL方式”模式下存在，配置要查询的SQL语句，使Loader可通过SQL语句查询结果并作为导入的数据。SQL语句需要有查询条件“WHERE \$ {CONDITIONS}”，否则无法正常工作。例如，“select * from TABLE WHERE A>B and \$ {CONDITIONS}”。如果同时配置“表列名”，SQL语句中查询的列将被“表列名”配置的列代替。不能和“架构名称”、“表名”同时配置。
	表列名	配置要导入的列，使Loader将列的内容全部导入。配置多个字段时使用逗号分隔。
	分区列名	指定数据库表的一列，根据该列来划分要导入的数据，在map任务中用于分区。建议配置主键字段。 <b>说明</b> <ul style="list-style-type: none"> <li>分区列必须有索引，如果没有索引，请不要指定分区列，指定没有索引的分区列会导致数据库服务器磁盘I/O繁忙，影响其他业务访问数据库，并且导入时间长。</li> <li>在有索引的多个字段中，选择字段值最离散的字段作为分区列，不离散的分区列会导致多个导入MR任务负载不均衡。</li> <li>分区列的排序规则必须支持大小写敏感，否则在数据导入过程中，可能会出现数据丢失。</li> <li>不建议分区列选择类型为float或double的字段，因为精度问题，可能导致分区列字段的最小值、最大值所在记录无法导入。</li> </ul>
	分区列空值	配置对数据库列中为null值记录的处理方式。值为“true”时，分区列的值为null的数据会被处理；值为“false”时，分区列的值为null的数据不会被处理。
	是否指定分区列	是否指定分区列。
	oracle-connector	表名
	列名	列名。
	查询条件	SQL语句中的查询条件。
	切分方式	指定数据的切分方式，有“ROWID”和“PARTITION”两种。



源文件类型	参数名	解释说明
	表分区名	表分区名，使用逗号分隔不同的分区。
	数据块分配方式	指定数据切分后，如何分配。
	读取大小	指定每次读取多大的数据量。
mysql-fastpath-connector	架构名称	数据库模式名。
	表名	数据库表名。
	查询条件	指定表的查询条件。
	分区列名	指定数据库表的一列，根据该列来划分要导入的数据，在map任务中用于分区。建议配置主键字段。 <b>说明</b> <ul style="list-style-type: none"> <li>分区列必须有索引，如果没有索引，请不要指定分区列，指定没有索引的分区列会导致数据库服务器磁盘I/O繁忙，影响其他业务访问数据库，并且导入时间长。</li> <li>在有索引的多个字段中，选择字段值最离散的字段作为分区列，不离散的分区列会导致多个导入MR任务负载不均衡。</li> <li>不建议分区列选择类型为float或double的字段，因为精度问题，可能导致分区列字段的最小值、最大值所在记录无法导入。</li> </ul>
	分区列空值	配置对数据库列中为null值记录的处理方式。 <ul style="list-style-type: none"> <li>值为“true”时，分区列的值为null的数据会被处理。</li> <li>值为“false”时，分区列的值为null的数据不会被处理。</li> </ul>
	是否指定分区列	是否指定分区列。
oracle-partition-connector	架构名称	数据库模式名。
	表名	分区表名。
	查询条件	SQL语句中的查询条件。
	表列名	配置要导入的列，使Loader将列的内容全部导入。配置多个字段时使用逗号分隔。

**步骤6** 在“3.转换”设置数据传输过程中的转换操作。

确认Loader创建的数据操作作业中，源数据的值是否满足直接使用需求而不进行转换，例如大小写转换、截取、拼接和分隔。

- 满足需求，请单击“下一步”。
  - 不满足需求，请执行[步骤6.1](#) ~ [步骤6.4](#)。
1. 默认没有已创建的转换步骤，可拖动左侧样例到编辑框，添加一个新的转换步骤。

2. 完整的转换流程包含以下类型，每个类型请根据业务需要进行选择。
  - a. 输入类型，第一个转换步骤，仅添加一种，任务涉及HBase或关系型数据库必须添加。
  - b. 转换类型，中间转换步骤，可添加一种以上或不添加。
  - c. 输出类型，最后一个转换步骤，仅添加一种，任务涉及HBase或关系型数据库必须添加。

表 17-32 样例一览表

类型	描述
输入类型	<ul style="list-style-type: none"><li>▪ CSV文件输入：CSV文件输入步骤，配置分隔符以转换生成多个字段。</li><li>▪ 固定宽度文件输入：文本文件输入步骤，配置截取字符或字节的长度以转换生成多个字段。</li><li>▪ 表输入：关系型数据输入步骤，配置数据库的指定列为输入的字段。</li><li>▪ HBase输入：HBase表输入步骤，配置HBase表的列定义到指定字段。</li><li>▪ HTML输入：HTML网页数据输入步骤，配置获取HTML网页文件目标数据到指定字段。</li><li>▪ Hive输入：Hive表输入步骤，配置Hive表的列定义到指定字段。</li><li>▪ Spark输入：SparkSQL表输入步骤，配置SparkSQL表的列定义到指定字段。仅支持SparkSQL存取Hive数据。</li></ul>

类型	描述
转换类型	<ul style="list-style-type: none"><li>▪ 长整型时间转换：长整型日期转换步骤，配置长整型数值与日期的转换。</li><li>▪ 空值转换：空值转换步骤，配置指定值替换空值。</li><li>▪ 随机值转换：随机数据生成步骤，配置新增值为随机数据的字段。</li><li>▪ 增加常量字段：增加常量步骤，配置直接生成常量字段。</li><li>▪ 拼接转换：拼接字段步骤，配置已生成的字段通过连接符连接，转换出新的字段。</li><li>▪ 分隔转换：分隔字段步骤，配置已生成的字段通过分隔符分隔，转换出新的字段。</li><li>▪ 取模转换：取模运算步骤，配置已生成的字段通过取模，转换出新的字段。</li><li>▪ 剪切字符串：字符串截取步骤，配置已生成的字段通过指定位置截取，转换出新的字段。</li><li>▪ EL操作转换：计算器，可以对字段值进行运算，目前支持的算子有：md5sum、sha1sum、sha256sum和sha512sum等。</li><li>▪ 字符串大小写转换：字符串转换步骤，配置已生成的字段通过大小写变换，转换出新的字段。</li><li>▪ 字符串逆序转换：字符串逆序步骤，配置已生成的字段通过逆序，转换出新的字段。</li><li>▪ 字符串空格清除转换：字符串空格清除步骤，配置已生成的字段通过清除空格，转换出新的字段。</li><li>▪ 过滤行转换：过滤行步骤，配置逻辑条件过滤掉含触发条件的行。</li><li>▪ 更新域：更新域步骤，配置当满足某些条件时，更新指定字段的值。</li></ul>

类型	描述
输出类型	<ul style="list-style-type: none"> <li>文件输出：文本文件输出步骤，配置已生成的字段通过分隔符连接并输出到文件。</li> <li>表输出：关系型数据库输出步骤，配置输出的字段对应到数据库的指定列。</li> <li>HBase输出：HBase表输出步骤，配置已生成的字段输出到HBase表的列。</li> <li>Hive输出：Hive表输出步骤，配置已生成的字段输出到Hive表的列。</li> <li>Spark输出：SparkSQL表输出步骤，配置已生成的字段输出到SparkSQL表的列。仅支持SparkSQL存取Hive数据。</li> </ul>

编辑栏包括以下几种任务：

- 重命名：重命名样例。
- 编辑：编辑步骤转换，参考[步骤6.3](#)。
- 删除：删除样例。

#### 说明

也可使用快捷键“Del”删除。

- 单击“编辑”，编辑步骤转换信息，配置字段与数据。  
步骤转换信息中的具体参数设置请参考[算子帮助](#)。

#### 说明

- 使用sftp-connector或ftp-connector导入数据时，在数据转换步骤中，需要将原数据中时间类型数值对应的字段，设置为字符串类型，才能精确到毫秒并完成导入。数据中包含比毫秒更精确的部分不会被导入。
- 使用generic-jdbc-connector导入数据时，在数据转换步骤中，建议“CHAR”或“VARCHAR”类型字段设置数据长度为“-1”，使全部数据正常导入，避免实际数据字符太长时被部分截取，出现缺失。
- 使用generic-jdbc-connector导入数据时，在数据转换步骤中，需要将原数据中时间类型数值对应的字段，设置为时间类型，才能精确到秒并完成导入。数据中包含比秒更精确的部分不会被导入。
- 导入到Hive分区表内表时，Hive默认不会扫描新导入的数据，需要执行如下HQL修复表才可以查询到新导入数据：

**MSCK REPAIR TABLE** *table\_name*;

转换步骤配置不正确时，传输的数据将无法转换并成为脏数据，脏数据标记规则如下：

- 任意输入类型步骤中，原数据包含字段的个数小于配置字段的个数，或者原数据字段值与配置字段的类型不匹配时，全部数据成为脏数据。
- “CSV文件输入”步骤中，“验证输入字段”检验输入字段与值的类型匹配情况，检查不匹配时跳过该行，当前行成为脏数据。

- “固定宽度文件输入”步骤中，“固定长度”指定字段分割长度，长度大于原字段值的长度则数据分割失败，当前行成为脏数据。
- “HBase输入”步骤中，“HBase表名”指定HBase表名不正确，或者“主键”没有配置主键列，全部数据成为脏数据。
- 任意转换类型步骤中，转换失败的行成为脏数据。例如“分隔转换”步骤中，生成的字段个数小于配置字段的个数，或者原数据不能转换为String类型，当前行成为脏数据。
- “过滤行转换”步骤中，被筛选条件过滤的行成为脏数据。
- “取模转换”步骤中，原字段值为“NULL”，当前行成为脏数据。
- 对于导入数据到Hive/SparkSQL表的作业，必须配置Hive的转换步骤。

4. 单击“下一步”。

**步骤7** 在“4.输出设置”，设置数据保存目标位置，然后单击“保存”保存作业或“保存并运行”，保存作业并运行作业。

**表 17-33** 输出配置参数一览表

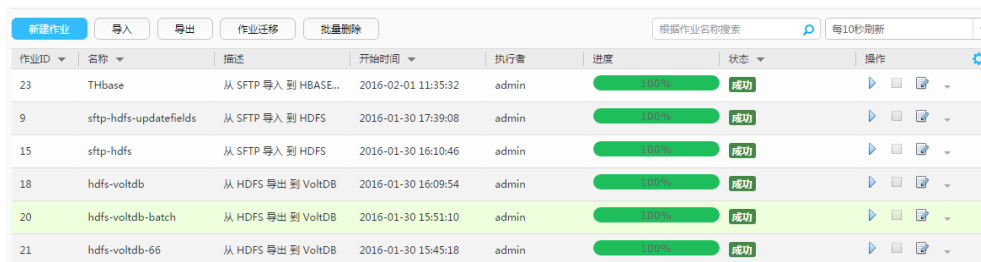
存储类型	参数名	解释说明
HDFS	文件类型	在下拉菜单中选择数据导入HDFS后保存文件的文件类型。 <ul style="list-style-type: none"> <li>• “TEXT_FILE”：导入文本文件并保存为文本文件。</li> <li>• “SEQUENCE_FILE”：导入文本文件并保存为sequence file文件格式。</li> <li>• “BINARY_FILE”：以二进制流的方式导入文件，可以导入任何格式的文件，不对文件做任何处理。</li> </ul> <b>说明</b> 文件类型选择“TEXT_FILE”或“SEQUENCE_FILE”导入时，Loader会自动根据文件的后缀选择对应的解压方法，对文件进行解压。
	压缩格式	在下拉菜单中选择数据导入HDFS后保存文件的压缩格式，未配置或选择NONE表示不压缩数据。
	输出目录	数据导入到HDFS里存储的保存目录。

存储类型	参数名	解释说明
	文件操作方式	<p>数据导入时的操作行为。全部数据从输入路径导入到目标路径时，先保存在临时目录，然后再从临时目录复制转移至目标路径，任务完成时删除临时路径的文件。转移临时文件存在同名文件时有以下行为：</p> <ul style="list-style-type: none"> <li>“OVERRIDE”：直接覆盖旧文件。</li> <li>“RENAME”：重命名新文件。无扩展名的文件直接增加字符串后缀，有扩展名的文件在文件名增加字符串后缀。字符串具有唯一性。</li> <li>“APPEND”：在旧文件尾部合并新文件内容。合并操作只是简单的追加，不保证追加文件是否可以使用。例如文本文件可合并，压缩文件合并后可能无法使用。</li> <li>“IGNORE”：保留旧文件，不复制新文件。</li> <li>“ERROR”：转移过程中出现同名文件时任务将停止执行并报错，已转移的文件导入成功，同名的文件及未转移的文档导入失败。</li> </ul>
	Map数	配置数据操作的MapReduce任务中同时启动的map数量。不可与“Map数据块大小”同时配置。参数值必须小于或等于“3000”。
	Map数据块大小	配置数据操作的MapReduce任务中启动map所处理的数据大小，单位为MB。参数值必须大于或等于“100”，建议配置值为“1000”。不可与“Map数”同时配置。当使用关系型数据库连接器时，不支持“Map数据块大小”，请配置“Map数”。
HBASE_BULKLOAD	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。
	导入前清理数据	导入前清空原表的数据。“true”为执行清空，“false”为不执行。不配置此参数则默认不执行清空。
	Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于“3000”。
	Map数据块大小	HBase不支持此参数，请配置“Map数”。

存储类型	参数名	解释说明
HBASE_PUTLIST	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。
	Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于“3000”。
	Map数据块大小	HBase不支持此参数，请配置“Map数”。
HIVE	输出目录	数据导入到Hive里存储的保存目录。
	Map数	配置数据操作的MapReduce任务中同时启动的map数量。不可与“Map数据块大小”同时配置。参数值必须小于或等于“3000”。
	Map数据块大小	配置数据操作的MapReduce任务中启动map所处理的数据大小，单位为MB。参数值必须大于或等于“100”，建议配置值为“1000”。不可与“Map数”同时配置。当使用关系型数据库连接器时，不支持“Map数据块大小”，请配置“Map数”。
SPARK	输出目录	仅支持SparkSQL存取Hive数据，指定数据导入到Hive里存储的保存目录。
	Map数	配置数据操作的MapReduce任务中同时启动的map数量。不可与“Map数据块大小”同时配置。参数值必须小于或等于“3000”。
	Map数据块大小	配置数据操作的MapReduce任务中启动map所处理的数据大小，单位为MB。参数值必须大于或等于“100”，建议配置值为“1000”。不可与“Map数”同时配置。当使用关系型数据库连接器时，不支持“Map数据块大小”，请配置“Map数”。

**步骤8** 已创建的作业可以在“Loader WebUI”界面上进行浏览，可进行启动、停止、复制、删除、编辑和查看历史信息操作。

图 17-3 查看 Loader 作业



作业ID	名称	描述	开始时间	执行者	进度	状态	操作
23	THbase	从 SFTP 导入 到 HBASE...	2016-02-01 11:35:32	admin	100%	成功	▶ □ 📄 ⚙
9	sftp-hdfs-updatefields	从 SFTP 导入 到 HDFS	2016-01-30 17:39:08	admin	100%	成功	▶ □ 📄 ⚙
15	sftp-hdfs	从 SFTP 导入 到 HDFS	2016-01-30 16:10:46	admin	100%	成功	▶ □ 📄 ⚙
18	hdfs-voltdb	从 HDFS 导出 到 VoltDB	2016-01-30 16:09:54	admin	100%	成功	▶ □ 📄 ⚙
20	hdfs-voltdb-batch	从 HDFS 导出 到 VoltDB	2016-01-30 15:51:10	admin	100%	成功	▶ □ 📄 ⚙
21	hdfs-voltdb-66	从 HDFS 导出 到 VoltDB	2016-01-30 15:45:18	admin	100%	成功	▶ □ 📄 ⚙

----结束

## 17.12.3 典型场景：从 SFTP 服务器导入数据到 HDFS/OBS

### 操作场景

该任务指导用户使用Loader将数据从SFTP服务器导入到HDFS/OBS。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HDFS/OBS目录和数据。
- 获取SFTP服务器使用的用户和密码，且该用户具备SFTP服务器上源文件的读取权限。若源文件在导入后文件名要增加后缀，则该用户还需具备源文件的写入权限。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 使用Loader从SFTP服务器导入数据时，确保SFTP服务器输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符“\”、“:”、“;”中的任意字符。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

### 操作步骤

#### 设置作业基本信息

步骤1 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-4 Loader WebUI 界面





**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-5 “基本信息”界面

1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“sftp-connector”，单击“添加”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。Loader支持配置多个SFTP服务器操作数据，单击“添加”可增加多行SFTP服务器的配置信息。

表 17-34 连接参数

参数名	说明	示例
名称	SFTP服务器连接的名称。	sftpName
Sftp服务器的IP	SFTP服务器的IP地址。	10.16.0.1
Sftp服务器端口	SFTP服务器的端口号。	22
Sftp用户名	访问SFTP服务器的用户名。	root
Sftp密码	访问SFTP服务器的密码。	xxxx
Sftp公钥	Sftp服务器公钥。	OdDt/yn...etM

#### 📖 说明

配置多个SFTP服务器时，多个SFTP服务器指定目录的数据导入到HDFS/OBS的同一个目录下。

#### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息。

表 17-35 输入设置参数

参数名	说明	示例
输入路径	<p>SFTP服务器中源文件的输入路径，如果连接器配置多个地址此处可对应使用“;”分隔多个输入路径，数量需要与连接器中服务器的数量一致。</p> <p><b>说明</b> 路径参数可以使用宏定义，具体请参考<a href="#">配置项中使用宏定义</a>。</p>	/opt/ tempfile/ opt
文件分割方式	<p>选择按文件或大小分割源文件，作为数据导入的MapReduce任务中各个map的输入文件。</p> <ul style="list-style-type: none"> <li>选择“FILE”，表示按文件分割源文件，即每个map处理一个或多个完整的源文件，同一个源文件不可分配至不同map，完成数据导入后保持源文件的目录结构。</li> <li>选择“SIZE”，表示按大小分割源文件，即每个map处理一定大小的输入文件，同一个源文件可分割至多个map，数据保存至输出目录时保存的文件数与map数量相同，文件名格式为“import_part_xxxx”，“xxxx”为系统生成的随机数，具有唯一性。</li> </ul>	FILE
过滤类型	<p>选择文件过滤的条件，与“路径过滤器”、“文件过滤器”配合使用。</p> <ul style="list-style-type: none"> <li>选择“WILDCARD”，表示使用通配符过滤。</li> <li>选择“REGEX”，表示使用正则表达式匹配。</li> <li>不选择，则默认为通配符过滤。</li> </ul>	WILDCARD
路径过滤器	<p>与“过滤类型”配合使用，配置通配符或正则表达式对源文件的输入路径包含的目录进行过滤。“输入路径”不参与过滤。使用分号“;”分隔多个服务器上的路径过滤器，每个服务器的多个过滤条件使用逗号“,”隔开。配置为空时表示不过滤目录。</p> <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul> <p>例如，当“过滤类型”选择“WILDCARD”时，将该参数设置为“*”;当“过滤类型”选择“REGEX”时，将该参数设置为“\\.*”。</p>	1*,2*;1*
文件过滤器	<p>与“过滤类型”配合使用，配置通配符或正则表达式对源文件的输入文件名进行过滤。使用分号“;”分隔多个服务器上的文件过滤器，每个服务器的多个过滤条件使用逗号“,”隔开。该参数不能配置为空。</p> <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul> <p>例如，当“过滤类型”选择“WILDCARD”时，将该参数设置为“*”;当“过滤类型”选择“REGEX”时，将该参数设置为“\\.*”。</p>	*.txt,*.csv; *.txt

参数名	说明	示例
编码类型	源文件的编码格式，如UTF-8、GBK。导入文本文件时才能配置。	UTF-8
后缀名	源文件导入成功后对输入文件增加的后缀值。该值为空，则表示不加后缀。数据源为文件系统，该参数才有效。用户若需增量导入数据建议设置该参数。 例如设置为“.txt”，源文件为“test-loader.csv”，则导出后源文件名为“test-loader.csv.txt”。	.log
压缩	使用SFTP协议导入数据时，是否开启压缩传输功能以减小带宽使用。 <ul style="list-style-type: none"> <li>选择“true”，表示开启压缩。</li> <li>选择“false”，表示关闭压缩。</li> </ul>	true

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-36](#)。

表 17-36 算子输入、输出参数设置

输入类型	输出类型
CSV文件输入	文件输出
HTML输入	文件输出
固定宽度文件输入	文件输出

图 17-6 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，在“存储类型”中选择“HDFS”，设置数据保存方式。

表 17-37 输出设置参数

参数名	说明	示例
文件类型	<p>文件导入后保存的类型：</p> <ul style="list-style-type: none"> <li>“TEXT_FILE”：导入文本文件并保存为文本文件</li> <li>“SEQUENCE_FILE”：导入文本文件并保存在“sequence file”文件格式</li> <li>“BINARY_FILE”：以二进制流的方式导入文件，可以导入任何格式的文件</li> </ul>	TEXT_FILE
压缩格式	在下拉菜单中选择数据导入HDFS/OBS后保存文件的压缩格式，未配置或选择NONE表示不压缩数据。	NONE
输出目录	<p>数据导入到HDFS/OBS里存储的保存目录。</p> <p><b>说明</b> 路径参数可以使用宏定义，具体请参考<a href="#">配置项中使用宏定义</a>。</p>	/user/test
文件操作方式	<p>数据导入时的操作行为。全部数据从输入路径导入到目标路径时，先保存在临时目录，然后再从临时目录复制转移至目标路径，任务完成时删除临时路径的文件。转移临时文件存在同名文件时有以下行为：</p> <ul style="list-style-type: none"> <li>“OVERRIDE”：直接覆盖旧文件。</li> <li>“RENAME”：重命名新文件。无扩展名的文件直接增加字符串后缀，有扩展名的文件在文件名增加字符串后缀。字符串具有唯一性。</li> <li>“APPEND”：在旧文件尾部合并新文件内容。合并操作只是简单的追加，不保证追加文件是否可以使用。例如文本文件可合并，压缩文件合并后可能无法使用。</li> <li>“IGNORE”：保留旧文件，不复制新文件。</li> <li>“ERROR”：转移过程中出现同名文件时任务将停止执行并报错，已转移的文件导入成功，同名的文件及未转移的文档导入失败。</li> </ul>	OVERRIDE

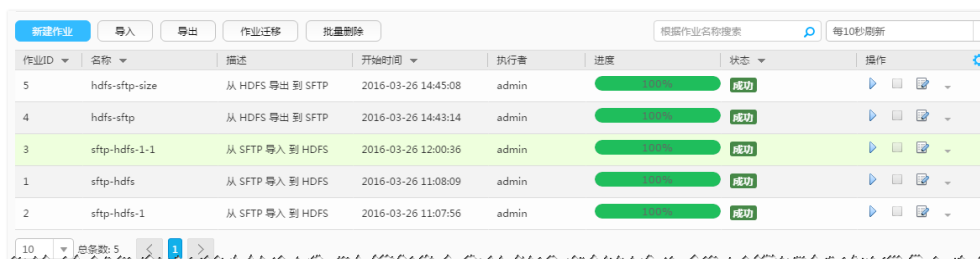
参数名	说明	示例
Map数	配置数据操作的MapReduce任务中同时启动的Map数量。不可与“Map数据块大小”同时配置。参数值必须小于或等于3000，建议以SFTP服务器的CPU的核数作为其取值。 <b>说明</b> 为了提高导入数据速度，需要确保以下条件： <ul style="list-style-type: none"> <li>每个Map连接时，相当于一个客户端连接，因此需要确保SFTP服务器最大连接数大于Map数量。</li> <li>确保SFTP服务器上的磁盘IO或是网络带宽都未达到上限。</li> </ul>	20
Map数据块大小	配置数据操作的MapReduce任务中启动map所处理的数据大小，单位为MB。参数值必须大于或等于100，建议配置值为1000。不可与“Map数”同时配置。	1000

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-7 查看作业



----结束

## 17.12.4 典型场景：从 SFTP 服务器导入数据到 HBase

### 操作场景

该任务指导用户使用Loader将数据从SFTP服务器导入到HBase。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HBase表或phoenix表。
- 获取SFTP服务器使用的用户和密码，且该用户具备SFTP服务器上源文件的读取权限。若源文件在导入后文件名要增加后缀，则该用户还需具备源文件的写入权限。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。

- 使用Loader从SFTP服务器导入数据时，确保SFTP服务器输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符“\”“;”中的任意字符。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

## 操作步骤

### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

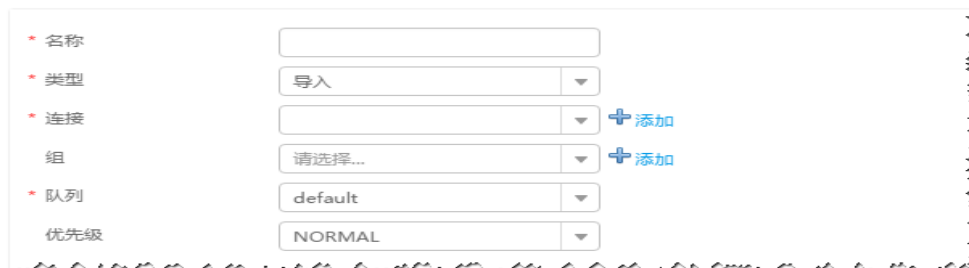
1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-8 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-9 “基本信息”界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“sftp-connector”，单击“添加”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。Loader支持配置多个SFTP服务器操作数据，单击“添加”可增加多行SFTP服务器的配置信息。

表 17-38 连接参数

参数名	说明	示例
名称	SFTP服务器连接的名称。	sftpName
Sftp服务器的IP	SFTP服务器的IP地址。	10.16.0.1
Sftp服务器端口	SFTP服务器的端口号。	22
Sftp用户名	访问SFTP服务器的用户名。	root
Sftp密码	访问SFTP服务器的密码。	xxxx
Sftp公钥	Sftp服务器公钥。	OdDt/yn...etM

**说明**

配置多个SFTP服务器，多个服务器指定目录的数据将导入到HBase。

**设置数据源信息**

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息。

表 17-39 输入设置参数

参数名	说明	示例
输入路径	SFTP服务器中源文件的输入路径，如果连接器配置多个地址此处可对应使用“;”分隔多个输入路径，数量需要与连接器中服务器的数量一致。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/opt/ tempfile;/opt
文件分割方式	选择按文件或大小分割源文件，作为数据导入的MapReduce任务中各个map的输入文件。 <ul style="list-style-type: none"> <li>选择“FILE”，表示按文件分割源文件，即每个map处理一个或多个完整的源文件，同一个源文件不可分配至不同map，完成数据导入后保持源文件的目录结构。</li> <li>选择“SIZE”，表示按大小分割源文件，即每个map处理一定大小的输入文件，同一个源文件可分割至多个map，数据保存至输出目录时保存的文件数与map数量相同，文件名格式为“import_part_xxxx”，“xxxx”为系统生成的随机数，具有唯一性。</li> </ul>	FILE

参数名	说明	示例
过滤类型	<p>选择文件过滤的条件，与“路径过滤器”、“文件过滤器”配合使用。</p> <ul style="list-style-type: none"> <li>选择“WILDCARD”，表示使用通配符过滤。</li> <li>选择“REGEX”，表示使用正则表达式匹配。</li> <li>不选择，则默认为通配符过滤。</li> </ul>	WILDCARD
路径过滤器	<p>与“过滤类型”配合使用，配置通配符或正则表达式对源文件的输入路径包含的目录进行过滤。“输入路径”不参与过滤。使用分号“;”分隔多个服务器上的路径过滤器，每个服务器的多个过滤条件使用逗号“,”隔开。配置为空时表示不过滤目录。</p> <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul> <p>例如，当“过滤类型”选择“WILDCARD”时，将该参数设置为“*”；当“过滤类型”选择“REGEX”时，将该参数设置为“\\.*”。</p>	1*,2*;1*
文件过滤器	<p>与“过滤类型”配合使用，配置通配符或正则表达式对源文件的输入文件名进行过滤。使用分号“;”分隔多个服务器上的文件过滤器，每个服务器的多个过滤条件使用逗号“,”隔开。该参数不能配置为空。</p> <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul> <p>例如，当“过滤类型”选择“WILDCARD”时，将该参数设置为“*”；当“过滤类型”选择“REGEX”时，将该参数设置为“\\.*”。</p>	*.txt,*.csv;*.txt
编码类型	<p>源文件的编码格式，如UTF-8、GBK。导入文本文件时才能配置。</p>	UTF-8
后缀名	<p>源文件导入成功后对输入文件增加的后缀值。该值为空，则表示不加后缀。数据源为文件系统，该参数才有效。用户若需增量导入数据建议设置该参数。</p> <p>例如设置为“.txt”，源文件为“test-loader.csv”，则导出后源文件名为“test-loader.csv.txt”。</p>	.log



参数名	说明	示例
压缩	使用SFTP协议导入数据时，是否开启压缩传输功能以减小带宽使用。 <ul style="list-style-type: none"> <li>选择“true”，表示开启压缩。</li> <li>选择“false”，表示关闭压缩。</li> </ul>	true

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-40](#)。

表 17-40 算子输入、输出参数设置

输入类型	输出类型
CSV文件输入	HBase输出
HTML输入	HBase输出
固定宽度文件输入	HBase输出

图 17-10 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，根据实际场景在“存储类型”选择“HBASE\_BULKLOAD”或“HBASE\_PUTLIST”，设置数据保存方式。

表 17-41 输出设置参数

存储类型	适用场景	参数名	说明	示例
HBASE_BULKLOAD	数据量大	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HBase
		导入前清理数据	导入前清空原表的数据。 “true”为执行清空，“false”为不执行。不配置此参数则默认不执行清空。	true
		Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000，建议以SFTP服务器当前最大连接数作为其取值。	20
		Map数据块大小	HBase不支持此参数，请配置“Map数”。	-
HBASE_PUTLIST	数据量小	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HBase
		Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000。	20

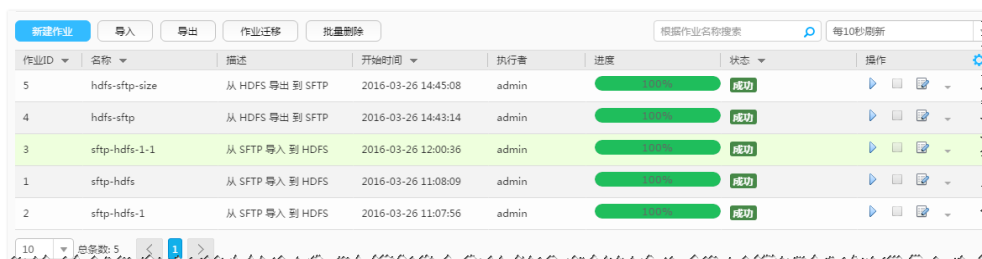
存储类型	适用场景	参数名	说明	示例
		Map数据块大小	HBase不支持此参数，请配置“Map数”。	-

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-11 查看作业



----结束

## 17.12.5 典型场景：从 SFTP 服务器导入数据到 Hive

### 操作场景

该任务指导用户使用Loader将数据从SFTP服务器导入到Hive。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业中指定的Hive表的权限。
- 获取SFTP服务器使用的用户和密码，且该用户具备SFTP服务器上源文件的读取权限。若源文件在导入后文件名要增加后缀，则该用户还需具备源文件的写入权限。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 使用Loader从SFTP服务器导入数据时，确保SFTP服务器输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符“\”、“/”、“:”、“;”中的任意字符。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

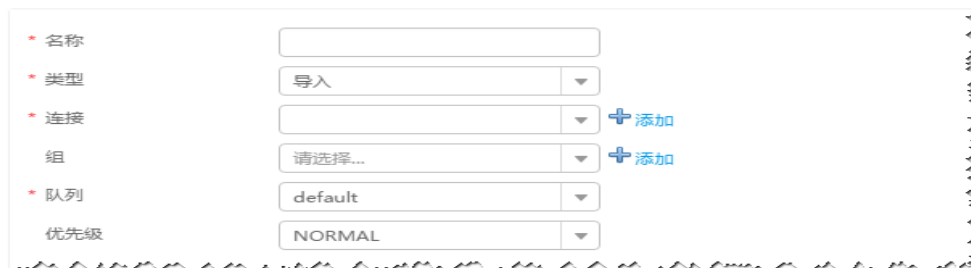
1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

**图 17-12** Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

**图 17-13** “基本信息”界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“sftp-connector”，单击“添加”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。Loader支持配置多个SFTP服务器操作数据，单击“添加”可增加多行SFTP服务器的配置信息。

**表 17-42** 连接参数

参数名	说明	示例
名称	SFTP服务器连接的名称。	sftpName
Sftp服务器的IP	SFTP服务器的IP地址。	10.16.0.1
Sftp服务器端口	SFTP服务器的端口号。	22

参数名	说明	示例
Sftp用户名	访问SFTP服务器的用户名。	root
Sftp密码	访问SFTP服务器的密码。	xxxx
Sftp公钥	Sftp服务器公钥。	OdDt/yn...etM

### 📖 说明

配置多个SFTP服务器，多个服务器指定目录的数据将导入到Hive。

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息。

表 17-43 输入设置参数

参数名	说明	示例
输入路径	SFTP服务器中源文件的输入路径，如果连接器配置多个地址此处可对应使用“;”分隔多个输入路径，数量需要与连接器中服务器的数量一致。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/opt/tem pfile; opt
文件分割方式	选择按文件或大小分割源文件，作为数据导入的MapReduce任务中各个map的输入文件。 <ul style="list-style-type: none"> <li>选择“FILE”，表示按文件分割源文件，即每个map处理一个或多个完整的源文件，同一个源文件不可分配至不同map，完成数据导入后保持源文件的目录结构。</li> <li>选择“SIZE”，表示按大小分割源文件，即每个map处理一定大小的输入文件，同一个源文件可分割至多个map，数据保存至输出目录时保存的文件数与map数量相同，文件名格式为“import_part_xxxx”，“xxxx”为系统生成的随机数，具有唯一性。</li> </ul>	FILE
过滤器类型	选择文件过滤的条件，与“路径过滤器”、“文件过滤器”配合使用。 <ul style="list-style-type: none"> <li>选择“WILDCARD”，表示使用通配符过滤。</li> <li>选择“REGEX”，表示使用正则表达式匹配。</li> <li>不选择，则默认为通配符过滤。</li> </ul>	WIL DC AR D

参数名	说明	示例
路径过滤器	<p>与“过滤类型”配合使用，配置通配符或正则表达式对源文件的输入路径包含的目录进行过滤。“输入路径”不参与过滤。使用分号“;”分隔多个服务器上的路径过滤器，每个服务器的多个过滤条件使用逗号“,”隔开。配置为空时表示不过滤目录。</p> <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul> <p>例如，当“过滤类型”选择“WILDCARD”时，将该参数设置为“*”;当“过滤类型”选择“REGEX”时，将该参数设置为“\\.*”。</p>	1*;2 *;1*
文件过滤器	<p>与“过滤类型”配合使用，配置通配符或正则表达式对源文件的输入文件名进行过滤。使用分号“;”分隔多个服务器上的文件过滤器，每个服务器的多个过滤条件使用逗号“,”隔开。该参数不能配置为空。</p> <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul> <p>例如，当“过滤类型”选择“WILDCARD”时，将该参数设置为“*”;当“过滤类型”选择“REGEX”时，将该参数设置为“\\.*”。</p>	*.txt ;.csv *.txt
编码类型	源文件的编码格式，如UTF-8、GBK。导入文本文件时才能配置。	UTF-8
后缀名	源文件导入成功后对输入文件增加的后缀值。该值为空，则表示不加后缀。数据源为文件系统，该参数才有效。用户若需增量导入数据建议设置该参数。 例如设置为“.txt”，源文件为“test-loader.csv”，则导出后源文件名为“test-loader.csv.txt”。	.log
压缩	使用SFTP协议导入数据时，是否开启压缩传输功能以减少带宽使用。 <ul style="list-style-type: none"> <li>选择“true”，表示开启压缩。</li> <li>选择“false”，表示关闭压缩。</li> </ul>	true

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-44](#)。

表 17-44 算子输入、输出参数设置

输入类型	输出类型
CSV文件输入	Hive输出
HTML输入	Hive输出
固定宽度文件输入	Hive输出

图 17-14 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，在“存储类型”选择“HIVE”，设置数据保存方式。

表 17-45 输出设置参数

参数名	说明	示例
输出目录	数据导入到Hive里存储的保存目录。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/opt/tempfile
Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000，建议以SFTP服务器当前最大连接数作为其取值。	20

参数名	说明	示例
Map数据块大小	Hive不支持此参数，请配置“Map数”。	-

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-15 查看作业

作业ID	名称	描述	开始时间	执行者	进度	状态	操作
5	hdfs-sftp-size	从 HDFS 导出到 SFTP	2016-03-26 14:45:08	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️ -
4	hdfs-sftp	从 HDFS 导出到 SFTP	2016-03-26 14:43:14	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️ -
3	sftp-hdfs-1-1	从 SFTP 导入到 HDFS	2016-03-26 12:00:36	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️ -
1	sftp-hdfs	从 SFTP 导入到 HDFS	2016-03-26 11:08:09	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️ -
2	sftp-hdfs-1	从 SFTP 导入到 HDFS	2016-03-26 11:07:56	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️ -

----结束

## 17.12.6 典型场景：从 FTP 服务器导入数据到 HBase

### 操作场景

该任务指导用户使用Loader将数据从FTP服务器导入到HBase。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 获取FTP服务器使用的用户和密码，且该用户具备FTP服务器上源文件的读取权限。若源文件在导入后文件名要增加后缀，则该用户还需具备源文件的写入权限。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 使用Loader从FTP服务器导入数据时，确保FTP服务器输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符/"/":;,中的任意字符。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。



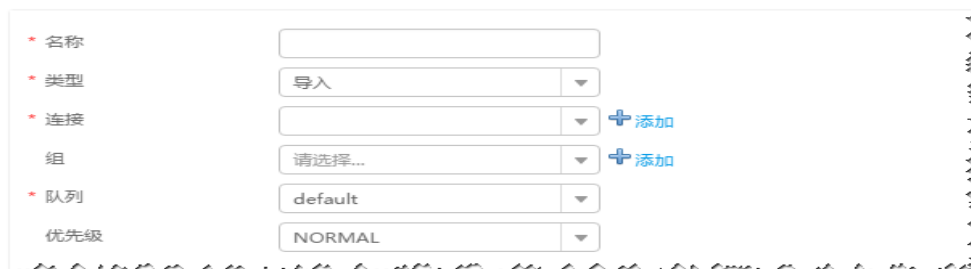
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-16 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-17 “基本信息”界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“ftp-connector”，单击“添加”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。Loader支持配置多个FTP服务器操作数据，单击“添加”可增加多行FTP服务器的配置信息。

表 17-46 连接参数

参数名	说明	示例
FTP服务器的IP	FTP服务器的IP地址。	ftpName
FTP服务器端口	FTP服务器的端口号。	22
FTP用户名	访问FTP服务器的用户名。	root
FTP密码	访问FTP服务器的密码。	xxxx

参数名	说明	示例
FTP模式	设置FTP访问模式，“ACTIVE”表示主动模式，“PASSIVE”表示被动模式。不指定参数值，默认为被动模式。	PASSIVE
FTP协议	设置FTP传输协议： <ul style="list-style-type: none"> <li>“FTP”：FTP协议。</li> <li>“SSL_EXPLICIT”：显式SSL协议。</li> <li>“SSL_IMPLICIT”：隐式SSL协议。</li> <li>“TLS_EXPLICIT”：显式TLS协议。</li> <li>“TLS_IMPLICIT”：隐式TLS协议。</li> </ul> 不指定参数值，默认为FTP协议。	FTP
文件名编码类型	填写FTP服务器支持的文件名、文件路径编码格式，不填写时使用系统默认格式UTF-8。	UTF-8

#### 说明

配置多个FTP服务器，多个服务器指定目录的数据将导入到HBase。

#### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息。

表 17-47 输入设置参数

参数名	说明	示例
输入路径	FTP服务器中源文件的输入路径，如果连接器配置多个地址此处可对应使用“;”分隔多个输入路径，数量需要与连接器中服务器的数量一致。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/opt/tempfile;/opt
文件分割方式	选择按文件或大小分割源文件，作为数据导入的MapReduce任务中各个map的输入文件。 <ul style="list-style-type: none"> <li>选择“FILE”，表示按文件分割源文件，即每个map处理一个或多个完整的源文件，同一个源文件不可分配至不同map，完成数据导入后保持源文件的目录结构。</li> <li>选择“SIZE”，表示按大小分割源文件，即每个map处理一定大小的输入文件，同一个源文件可分割至多个map，数据保存至输出目录时保存的文件数与map数量相同，文件名格式为“import_part_xxxx”，“xxxx”为系统生成的随机数，具有唯一性。</li> </ul>	FILE

参数名	说明	示例
过滤类型	<p>选择文件过滤的条件，与“路径过滤器”、“文件过滤器”配合使用。</p> <ul style="list-style-type: none"> <li>选择“WILDCARD”，表示使用通配符过滤。</li> <li>选择“REGEX”，表示使用正则表达式匹配。</li> <li>不选择，则默认为通配符过滤。</li> </ul>	WILDCARD
路径过滤器	<p>与“过滤类型”配合使用，配置通配符或正则表达式对源文件的输入路径包含的目录进行过滤。“输入路径”不参与过滤。使用分号“;”分隔多个服务器上的路径过滤器，每个服务器的多个过滤条件使用逗号“,”隔开。配置为空时表示不过滤目录。</p> <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul> <p>例如，当“过滤类型”选择“WILDCARD”时，将该参数设置为“*”;当“过滤类型”选择“REGEX”时，将该参数设置为“\\.*”。</p>	1*,2*;1*
文件过滤器	<p>与“过滤类型”配合使用，配置通配符或正则表达式对源文件的输入文件名进行过滤。使用分号“;”分隔多个服务器上的文件过滤器，每个服务器的多个过滤条件使用逗号“,”隔开。该参数不能配置为空。</p> <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul> <p>例如，当“过滤类型”选择“WILDCARD”时，将该参数设置为“*”;当“过滤类型”选择“REGEX”时，将该参数设置为“\\.*”。</p>	*.txt,*.csv;*.txt
编码类型	<p>源文件的编码格式，如UTF-8、GBK。导入文本文件时才能配置。</p>	UTF-8
后缀名	<p>源文件导入成功后对输入文件增加的后缀值。该值为空，则表示不加后缀。数据源为文件系统，该参数才有效。用户若需增量导入数据建议设置该参数。</p> <p>例如设置为“.txt”，源文件为“test-loader.csv”，则导出后源文件名为“test-loader.csv.txt”。</p>	.log

参数名	说明	示例
压缩	使用FTP协议导入数据时，是否开启压缩传输功能以减小带宽使用。 <ul style="list-style-type: none"> <li>选择“true”，表示开启压缩。</li> <li>选择“false”，表示关闭压缩。</li> </ul>	true

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-48](#)。

表 17-48 算子输入、输出参数设置

输入类型	输出类型
CSV文件输入	HBase输出
HTML输入	HBase输出
固定宽度文件输入	HBase输出

图 17-18 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，根据实际场景在“存储类型”选择“HBASE\_BULKLOAD”或“HBASE\_PUTLIST”，设置数据保存方式。

表 17-49 输出设置参数

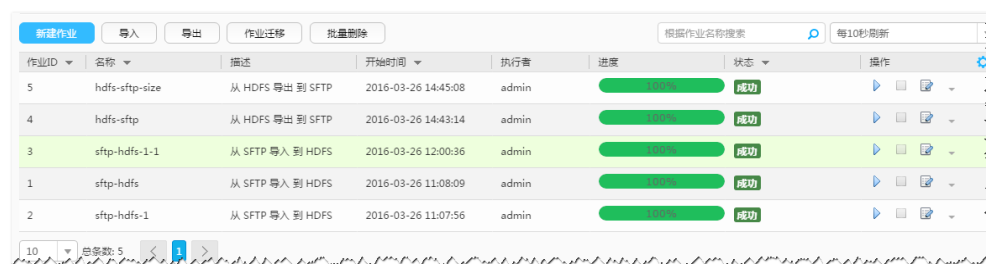
存储类型	适用场景	参数名	说明	示例
HBASE_B ULKLOAD	数据量大	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HBase
		导入前清理数据	导入前清空原表的数据。“True”为执行清空，“False”为不执行。不配置此参数则默认不执行清空。	true
		Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000，建议以FTP服务器当前最大连接数作为其取值。	20
		Map数据块大小	HBase不支持此参数，请配置“Map数”。	-
HBASE_P UTLIST	数据量小	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HBase
		Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000。	20
		Map数据块大小	HBase不支持此参数，请配置“Map数”。	-

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-19 查看作业



----结束

## 17.12.7 典型场景：从关系型数据库导入数据到 HDFS/OBS

### 操作场景

该任务指导用户使用Loader将数据从关系型数据库导入到HDFS/OBS。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HDFS/OBS目录和数据。
- 获取关系型数据库使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。
- 操作前需要进行如下配置：
  - a. 获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：  
“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib`”。
  - b. 使用root用户在主备节点分别执行以下命令修改权限：  

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/
FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/
webapps/loader/WEB-INF/ext-lib
chown omm:wheel jar包文件名
chmod 600 jar包文件名
```
  - c. 登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Loader > 更多 > 重启服务”，输入管理员密码重启Loader服务。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-20 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-21 “基本信息”界面

1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“generic-jdbc-connector”或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

**说明**

- 与关系数据库连接时，可以选择通用数据库连接器（generic-jdbc-connector）或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlimport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。

表 17-50 “generic-jdbc-connector” 连接参数

参数名	说明	示例
名称	关系型数据库连接的名称。	dbName
JDBC驱动程序类	JDBC驱动类名。	oracle.jdbc.driver.OracleDriver
JDBC连接字符串	JDBC连接字符串。	jdbc:oracle:thin:@//10.16.0.1:1521/oradb
用户名	连接数据库使用的用户名。	omm
密码	连接数据库使用的密码。	xxxx

参数名	说明	示例
JDBC连接属性	JDBC连接属性，单击“添加”手动添加。 <ul style="list-style-type: none"> <li>名称：连接属性名</li> <li>值：连接属性值</li> </ul>	<ul style="list-style-type: none"> <li>名称：socketTimeout</li> <li>值：20</li> </ul>

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息。

表 17-51 输入设置参数

参数名	说明	示例
架构名称	“表方式”模式下存在，数据库模式名。	public
表名	“表方式”模式下存在，数据库表名。	test
SQL语句	“SQL方式”模式下存在，配置要查询的SQL语句，使Loader可通过SQL语句查询结果并作为导入的数据。SQL语句需要有查询条件“WHERE \$ {CONDITIONS}”，否则无法正常工作。例如，“select * from TABLE WHERE A>B and \$ {CONDITIONS}”。如果同时配置“表列名”，SQL语句中查询的列将被“表列名”配置的列代替。不能和“架构名称”、“表名”同时配置。 <b>说明</b> SQL Where语句可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	select * from TABLE WHERE A>B and \$ {CONDITIONS}
表列名	配置要导入的列，使Loader将列的内容全部导入。配置多个字段时使用“,”分隔。如果不配置，则导入所有列，同时“Select *”的顺序作为列的位置。	id,name
分区列名	指定数据库表的一列，根据该列来划分要导入的数据，在Map任务中用于分区。建议配置主键字段。 <b>说明</b> <ul style="list-style-type: none"> <li>分区列必须有索引，如果没有索引，请不要指定分区列，指定没有索引的分区列会导致数据库服务器磁盘I/O繁忙，影响其他业务访问数据库，并且导入时间长。</li> <li>在有索引的多个字段中，选择字段值最离散的字段作为分区列，不离散的分区列会导致多个导入MR任务负载不均衡。</li> <li>分区列的排序规则必须支持大小写敏感，否则在数据导入过程中，可能会出现数据丢失。</li> <li>不建议分区列选择类型为float或double的字段，因为精度问题，可能导致分区列字段的最小值、最大值所在记录无法导入。</li> </ul>	id



参数名	说明	示例
分区列空值	配置对数据库列中为null值记录的处理方式。 <ul style="list-style-type: none"> <li>值为“true”时，分区列的值为null的数据会被处理；</li> <li>值为“false”时，分区列的值为null的数据不会被处理。</li> </ul>	true
是否指定分区列	是否指定分区列。	true

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-52](#)。

表 17-52 算子输入、输出参数设置

输入类型	输出类型
表输入	文件输出

图 17-22 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，在“存储类型”中选择“HDFS”，设置数据保存方式。

表 17-53 输出设置参数

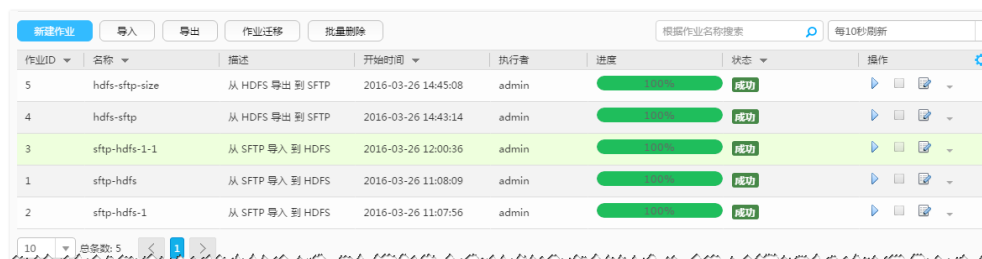
参数名	说明	示例
文件类型	<p>文件导入后保存的类型：</p> <ul style="list-style-type: none"> <li>“TEXT_FILE”：导入文本文件并保存为文本文件</li> <li>“SEQUENCE_FILE”：导入文本文件并保存在“sequence file”文件格式</li> <li>“BINARY_FILE”：以二进制流的方式导入文件，可以导入任何格式的文件</li> </ul>	TEXT_FILE
压缩格式	在下拉菜单中选择数据导入HDFS/OBS后保存文件的压缩格式，未配置或选择“NONE”表示不压缩数据。	NONE
输出目录	<p>数据导入到HDFS/OBS里存储的保存目录。</p> <p><b>说明</b> 路径参数可以使用宏定义，具体请参考<a href="#">配置项中使用宏定义</a>。</p>	/user/test
文件操作方式	<p>数据导入时的操作行为。全部数据从输入路径导入到目标路径时，先保存在临时目录，然后再从临时目录复制转移至目标路径，任务完成时删除临时路径的文件。转移临时文件存在同名文件时有以下行为：</p> <ul style="list-style-type: none"> <li>“OVERRIDE”：直接覆盖旧文件。</li> <li>“RENAME”：重命名新文件。无扩展名的文件直接增加字符串后缀，有扩展名的文件在文件名增加字符串后缀。字符串具有唯一性。</li> <li>“APPEND”：在旧文件尾部合并新文件内容。合并操作只是简单的追加，不保证追加文件是否可以使用。例如文本文件可合并，压缩文件合并后可能无法使用。</li> <li>“IGNORE”：保留旧文件，不复制新文件。</li> <li>“ERROR”：转移过程中出现同名文件时任务将停止执行并报错，已转移的文件导入成功，同名的文件及未转移的文档导入失败。</li> </ul>	OVERRIDE
Map数	配置数据操作的MapReduce任务中同时启动的Map数量。不可与“Map数据块大小”同时配置。参数值必须小于或等于3000。	-
Map数据块大小	配置数据操作的MapReduce任务中启动map所处理的数据大小，单位为MB。参数值必须大于或等于100，建议配置值为1000。不可与“Map数”同时配置。当使用关系型数据库连接器时，不支持“Map数据块大小”，请配置“Map数”。	1000

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

步骤8 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-23 查看作业



作业ID	名称	描述	开始时间	执行者	进度	状态	操作
5	hdfs-sftp-size	从 HDFS 导出 到 SFTP	2016-03-26 14:45:08	admin	100%	成功	▶ □ 🗑️
4	hdfs-sftp	从 HDFS 导出 到 SFTP	2016-03-26 14:43:14	admin	100%	成功	▶ □ 🗑️
3	sftp-hdfs-1-1	从 SFTP 导入 到 HDFS	2016-03-26 12:00:36	admin	100%	成功	▶ □ 🗑️
1	sftp-hdfs	从 SFTP 导入 到 HDFS	2016-03-26 11:08:09	admin	100%	成功	▶ □ 🗑️
2	sftp-hdfs-1	从 SFTP 导入 到 HDFS	2016-03-26 11:07:56	admin	100%	成功	▶ □ 🗑️

----结束

## 17.12.8 典型场景：从关系型数据库导入数据到 HBase

### 操作场景

该任务指导用户使用Loader将数据从关系型数据库导入到HBase。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HBase表或phoenix表。
- 获取关系型数据库使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。
- 操作前需要进行如下配置：
  - a. 获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：  
“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib`”。
  - b. 使用root用户在主备节点分别执行以下命令修改权限：  

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/
FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/
webapps/loader/WEB-INF/ext-lib
chown omm:wheel jar包文件名
chmod 600 jar包文件名
```
  - c. 登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Loader > 更多 > 重启服务”，输入管理员密码重启Loader服务。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

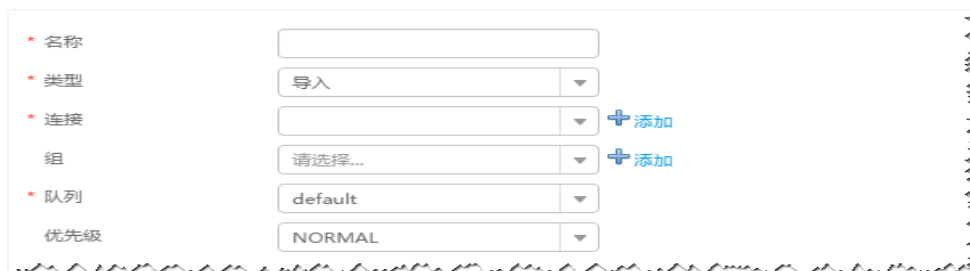
1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-24 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-25 “基本信息”界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“generic-jdbc-connector”或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

**说明**

- 与关系数据库连接时，可以选择通用数据库连接器（generic-jdbc-connector）或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlimport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。

表 17-54 “generic-jdbc-connector” 连接参数

参数名	说明	示例
名称	关系型数据库连接的名称。	dbName
JDBC驱动程序类	JDBC驱动类名。	oracle.jdbc.driver.OracleDriver
JDBC连接字符串	JDBC连接字符串。	jdbc:oracle:thin:@//10.16.0.1:1521/oradb
用户名	连接数据库使用的用户名。	omm
密码	连接数据库使用的密码。	xxxx
JDBC连接属性	JDBC连接属性，单击“添加”手动添加。 <ul style="list-style-type: none"> <li>名称：连接属性名</li> <li>值：连接属性值</li> </ul>	<ul style="list-style-type: none"> <li>名称：socketTimeout</li> <li>值：20</li> </ul>

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息。

表 17-55 输入设置参数

参数名	说明	示例
架构名称	“表方式”模式下存在，数据库模式名。	dbo
表名	“表方式”模式下存在，数据库表名。	test
SQL语句	“SQL方式”模式下存在，配置要查询的SQL语句，使Loader可通过SQL语句查询结果并作为导入的数据。SQL语句需要有查询条件“WHERE \${CONDITIONS}”，否则无法正常工作。例如，“select * from TABLE WHERE A>B and \${CONDITIONS}”。如果同时配置“表列名”，SQL语句中查询的列将被“表列名”配置的列代替。不能和“架构名称”、“表名”同时配置。 <b>说明</b> SQL Where语句可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	select * from test where \${CONDITIONS}
表列名	配置要导入的列，使Loader将列的内容全部导入。配置多个字段时使用“,”分隔。 如果不配置，则导入所有列，同时“Select *”的顺序作为列的位置。	-

参数名	说明	示例
分区列名	<p>指定数据库表的一列，根据该列来划分要导入的数据，在Map任务中用于分区。建议配置主键字段。</p> <p><b>说明</b></p> <ul style="list-style-type: none"> <li>分区列必须有索引，如果没有索引，请不要指定分区列，指定没有索引的分区列会导致数据库服务器磁盘I/O繁忙，影响其他业务访问数据库，并且导入时间长。</li> <li>在有索引的多个字段中，选择字段值最离散的字​​段作为分区列，不离散的分​​区列会导致多个导入MR任务负载不均衡。</li> <li>分区列的排序规则必须支持大小写敏感，否则在数据导入过程中，可能会出现数据丢失。</li> <li>不建议分区列选择类型为float或double的字段，因为精度问题，可能导致分区列字段的最小值、最大值所在记录无法导入。</li> </ul>	id
分区列空值	<p>配置对数据库列中为null值记录的处理方式。</p> <ul style="list-style-type: none"> <li>值为“true”时，分区列的值为null的数据会被处理；</li> <li>值为“false”时，分区列的值为null的数据不会被处理。</li> </ul>	true
是否指定分区列	是否指定分区列。	true

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-56](#)。

**表 17-56** 算子输入、输出参数设置

输入类型	输出类型
表输入	HBase输出

图 17-26 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，根据实际场景在“存储类型”选择“HBASE\_BULKLOAD”或“HBASE\_PUTLIST”，设置数据保存方式。

表 17-57 输出设置参数

存储类型	适用场景	参数名	说明	示例
HBASE_BULKLOAD	数据量大	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HB ase
		导入前清理数据	导入前清空原表的数据。“True”为执行清空，“False”为不执行。不配置此参数则默认不执行清空。	tru e
		Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000。	20
		Map数据块大小	HBase不支持此参数，请配置“Map数”。	-
HBASE_PUTLIST	数据量小	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HB ase

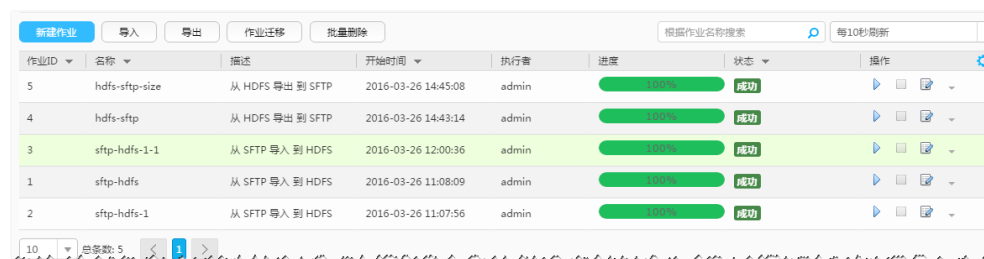
存储类型	适用场景	参数名	说明	示例
		Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000。	1000
		Map数据块大小	HBase不支持此参数，请配置“Map数”。	-

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-27 查看作业



----结束

## 17.12.9 典型场景：从关系型数据库导入数据到 Hive

### 操作场景

该任务指导用户使用Loader将数据从关系型数据库导入到Hive。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的Hive表。
- 获取关系型数据库使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。
- 操作前需要进行如下配置：
  - a. 获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：“\${BIGDATA\_HOME}/FusionInsight\_Porter\_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib”。



- b. 使用root用户在主备节点分别执行以下命令修改权限：  

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/
FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/
webapps/loader/WEB-INF/ext-lib
chown omm:wheel jar包文件名
chmod 600 jar包文件名
```
- c. 登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Loader > 更多 > 重启服务”，输入管理员密码重启Loader服务。

## 操作步骤

### 设置作业基本信息

#### 步骤1 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-28 Loader WebUI 界面



#### 步骤2 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-29 “基本信息”界面

1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“generic-jdbc-connector”或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

**说明**

- 与关系数据库连接时，可以选择通用数据库连接器（generic-jdbc-connector）或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlimport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。

**表 17-58 “generic-jdbc-connector” 连接参数**

参数名	说明	示例
名称	关系型数据库连接的名称。	dbName
JDBC驱动程序类	JDBC驱动类名。	oracle.jdbc.driver.OracleDriver
JDBC连接字符串	JDBC连接字符串。	jdbc:oracle:thin:@//10.16.0.1:1521/oradb
用户名	连接数据库使用的用户名。	omm
密码	连接数据库使用的密码。	xxxx
JDBC连接属性	JDBC连接属性，单击“添加”手动添加。 <ul style="list-style-type: none"> <li>名称：连接属性名</li> <li>值：连接属性值</li> </ul> <b>须知</b> 使用通用连接器连接MySQL时，在大数据量场景下，需要在MySQL的JDBC连接串中设置连接属性“useCursorFetch=true”。	<ul style="list-style-type: none"> <li>名称：socketTimeout</li> <li>值：20</li> </ul>

**设置数据源信息**

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息。

**表 17-59 输入设置参数**

参数名	说明	示例
架构名称	“表方式”模式下存在，数据库模式名。	dbo

参数名	说明	示例
表名	“表方式”模式下存在，数据库表名。	test
SQL语句	<p>“SQL方式”模式下存在，配置要查询的SQL语句，使Loader可通过SQL语句查询结果并作为导入的数据。SQL语句需要有查询条件“WHERE \${CONDITIONS}”，否则无法正常工作。例如，“select * from TABLE WHERE A&gt;B and \${CONDITIONS}”。如果同时配置“表列名”，SQL语句中查询的列将被“表列名”配置的列代替。不能和“架构名称”、“表名”同时配置。</p> <p><b>说明</b> SQL Where语句可以使用宏定义，具体请参考<a href="#">配置项中使用宏定义</a>。</p>	select * from test where \$ {CONDITIONS}
表列名	<p>配置要导入的列，使Loader将列的内容全部导入。配置多个字段时使用“,”分隔。</p> <p>如果不配置，则导入所有列，同时“Select *”的顺序作为列的位置。</p>	-
分区列名	<p>指定数据库表的一列，根据该列来划分要导入的数据，在Map任务中用于分区。建议配置主键字段。</p> <p><b>说明</b></p> <ul style="list-style-type: none"> <li>分区列必须有索引，如果没有索引，请不要指定分区列，指定没有索引的分区列会导致数据库服务器磁盘I/O繁忙，影响其他业务访问数据库，并且导入时间长。</li> <li>在有索引的多个字段中，选择字段值最离散的字段作为分区列，不分散的分区列会导致多个导入MR任务负载不均衡。</li> <li>分区列的排序规则必须支持大小写敏感，否则在数据导入过程中，可能会出现数据丢失。</li> <li>不建议分区列选择类型为float或double的字段，因为精度问题，可能导致分区列字段的最小值、最大值所在记录无法导入。</li> </ul>	id
分区列空值	<p>配置对数据库列中为null值记录的处理方式。</p> <ul style="list-style-type: none"> <li>值为“true”时，分区列的值为null的数据会被处理；</li> <li>值为“false”时，分区列的值为null的数据不会被处理。</li> </ul>	true
是否指定分区列	是否指定分区列。	true

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-60](#)。

**表 17-60** 算子输入、输出参数设置

输入类型	输出类型
表输入	Hive输出

**图 17-30** 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，在“存储类型”选择“HIVE”，设置数据保存方式。

**表 17-61** 输出设置参数

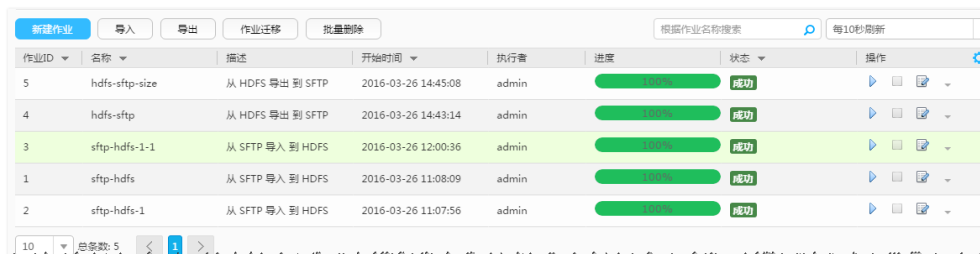
参数名	说明	示例
输出目录	数据导入到Hive里存储的保存目录。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/opt/ tempfile
Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000，建议以SFTP服务器当前最大连接数作为其取值。	20
Map数据块大小	Hive不支持此参数，请配置“Map数”。	-

**步骤7** 单击“保存并运行”，开始保存并运行作业。

#### 查看作业完成情况

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-31 查看作业



作业ID	名称	描述	开始时间	执行者	进度	状态	操作
5	hdfs-sftp-size	从 HDFS 导出到 SFTP	2016-03-26 14:45:08	admin	<div style="width: 100%;"></div>	成功	▶ □ 🔄 -
4	hdfs-sftp	从 HDFS 导出到 SFTP	2016-03-26 14:43:14	admin	<div style="width: 100%;"></div>	成功	▶ □ 🔄 -
3	sftp-hdfs-1-1	从 SFTP 导入到 HDFS	2016-03-26 12:00:36	admin	<div style="width: 100%;"></div>	成功	▶ □ 🔄 -
1	sftp-hdfs	从 SFTP 导入到 HDFS	2016-03-26 11:08:09	admin	<div style="width: 100%;"></div>	成功	▶ □ 🔄 -
2	sftp-hdfs-1	从 SFTP 导入到 HDFS	2016-03-26 11:07:56	admin	<div style="width: 100%;"></div>	成功	▶ □ 🔄 -

----结束

## 17.12.10 典型场景：从 HDFS/OBS 导入数据到 HBase

### 操作场景

该任务指导用户使用Loader将文件从HDFS/OBS导入到HBase。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HDFS/OBS目录和数据。
- 确保用户已授权访问作业执行时操作的HBase表或phoenix表。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 使用Loader从HDFS/OBS导入数据时，确保HDFS/OBS输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符“\”、“/”、“:”、“;”中的任意字符。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

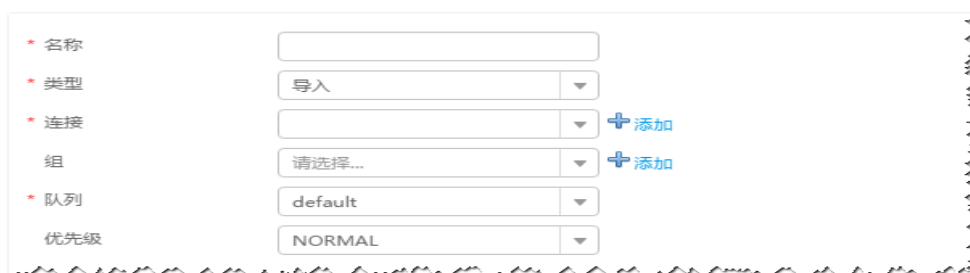
1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-32 Loader WebUI 界面



步骤2 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-33 “基本信息”界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

步骤3 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“hdfs-connector”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

### 设置数据源信息

步骤4 单击“下一步”，进入“输入设置”界面，设置数据源信息。

表 17-62 输入设置参数

参数名	说明	示例
输入路径	HDFS/OBS中源文件的输入路径。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/user/test
路径过滤器	配置通配符对源文件的输入路径包含的目录进行过滤。“输入路径”不参与过滤。配置多个过滤条件时使用“,”隔开，配置为空时表示不过滤目录。不支持正则表达式过滤。	*

参数名	说明	示例
文件过滤器	配置通配符对源文件的输入文件名进行过滤。配置多个过滤条件时使用“,” 隔开。不能配置为空。不支持正则表达式过滤。	*
编码类型	源文件的编码格式，如UTF-8。导入文本文件时才能配置。	UTF-8
后缀名	源文件导入成功后对输入文件增加的后缀值。该值为空，表示不加后缀。	.log

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-63](#)。

**表 17-63** 算子输入、输出参数设置

输入类型	输出类型
CSV文件输入	HBase输出
HTML输入	HBase输出
固定宽度文件输入	HBase输出

**图 17-34** 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，根据实际场景在“存储类型”选择“HBASE\_BULKLOAD”或“HBASE\_PUTLIST”，设置数据保存方式。

表 17-64 输出设置参数

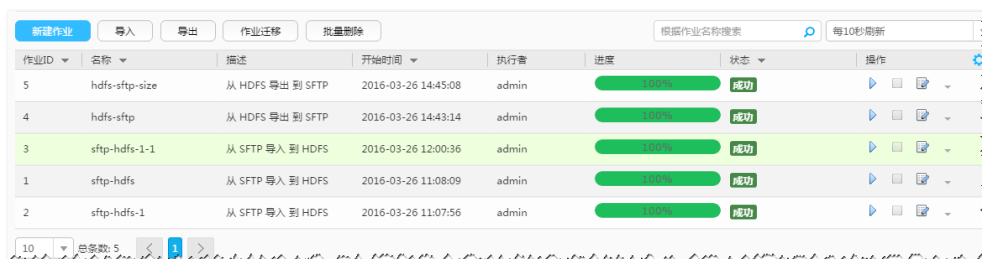
存储类型	适用场景	参数名	说明	示例
HBASE_BULKLOAD	数据量大	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HBase
		导入前清理数据	导入前清空原表的数据。“True”为执行清空，“False”为不执行。不配置此参数则默认不执行清空。	true
		Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000。	20
		Map数据块大小	HBase不支持此参数，请配置“Map数”。	-
HBASE_PUTLIST	数据量小	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HBase
		Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000。	20
		Map数据块大小	HBase不支持此参数，请配置“Map数”。	-

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-35 查看作业



----结束



## 17.12.11 典型场景：从关系型数据库导入数据到 ClickHouse

### 操作场景

该任务指导用户使用Loader将数据从关系型数据库导入到ClickHouse，本章节以MySQL为例进行操作。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- ClickHouse表已创建，确保用户已授权访问作业执行时操作该表的权限。
- 获取MySQL数据库使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。
- 操作前需要进行如下配置：
  - a. 从MySQL数据库安装路径下获取MySQL客户端jar包（如mysqlclient-5.8.1.jar），将其保存在Loader服务主备节点的lib路径：“\${BIGDATA\_HOME}/FusionInsight\_Porter\_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib”。
  - b. 在ClickHouse的安装目录获取clickhouse-jdbc-\*.jar包，将其保存在Loader服务主备节点的lib路径：“\${BIGDATA\_HOME}/FusionInsight\_Porter\_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib”。
  - c. 使用root用户在主备节点分别执行以下命令修改权限：

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/
FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/
webapps/loader/WEB-INF/ext-lib
chown omm:wheel jar包文件名
chmod 600 jar包文件名
```
  - d. 登录FusionInsight Manager系统，选择“集群 > 服务 > Loader > 更多 > 重启服务”，输入管理员密码重启Loader服务。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

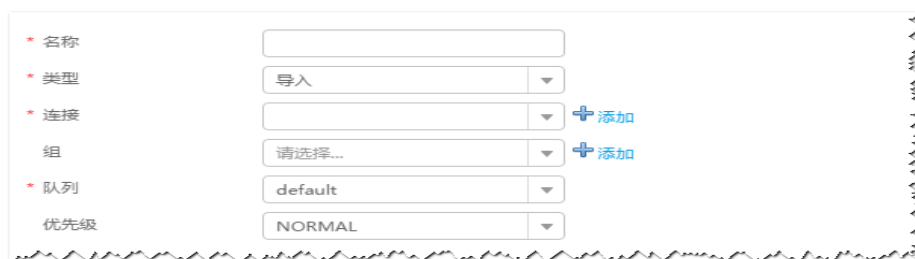
1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-36 Loader WebUI 界面



步骤2 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-37 “基本信息”界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

步骤3 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“generic-jdbc-connector”或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

说明

- 与关系数据库连接时，可以选择通用数据库连接器（generic-jdbc-connector）或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlimport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。

表 17-65 “generic-jdbc-connector” 连接参数

参数名	说明	示例
名称	关系型数据库连接的名称。	mysql_test
JDBC驱动程序类	JDBC驱动类名。	com.mysql.jdbc.Driver

参数名	说明	示例
JDBC连接字符串	JDBC连接字符串。	jdbc:mysql://10.254.144.102:3306/test?useUnicode=true&characterEncoding=UTF-8
用户名	连接数据库使用的用户名。	root
密码	连接数据库使用的密码。	xxxx

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息，暂只支持选择“表方式”。

表 17-66 输入设置参数

参数名	说明	示例
架构名称	用户指定数据库的模式名。	public
表名	表名称。	test
表列名	指定要输入的列名。	id,name
是否指定分区列	暂只支持不指定分区模式。	false

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-67](#)。

表 17-67 算子输入、输出参数设置

输入类型	输出类型
MySQL输入	ClickHouse输出

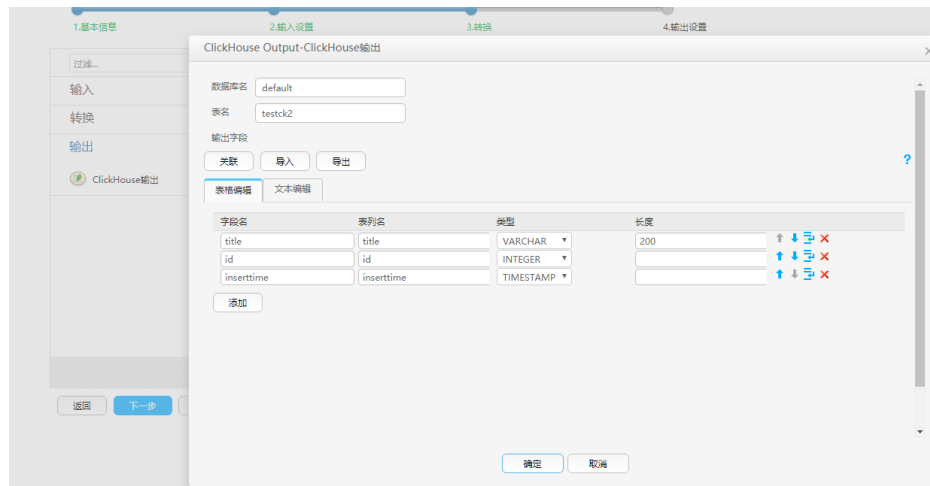
在输入中把“表输入”拖拽到网格中，双击“表输入”，选择“自动识别”如[图17-38](#)所示。

图 17-38 算子输入



在输出中把“ClickHouse输出”拖拽到网格中，双击“表输出”，选择“关联”或者手动编辑表格，与输入的表格对应，如图17-39所示。

图 17-39 算子输出



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，在“存储类型”中选择“CLICKHOUSE”，设置数据保存方式。

表 17-68 输出设置参数

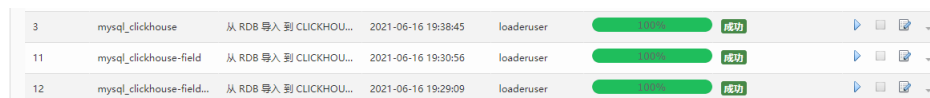
参数名	说明	示例
存储类型	选择CLICKHOUSE。	-
ClickHouse实例	选择ClickHouse。	-
导入前清理数据	选择“true”或“false”。 <b>说明</b> 如果导入的表为ClickHouse分布式表，且需要清理数据时，请在导入前手动删除ClickHouse分布式表对应的本地表中的数据。	true

**步骤7** 单击“保存并运行”，开始保存并运行作业。

### 查看作业完成情况

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-40 查看作业



**步骤9** 使用ClickHouse客户端，查询ClickHouse表数据是否和MySQL的表数据一致。

----结束

## 17.12.12 典型场景：从 HDFS 导入数据到 ClickHouse

### 操作场景

该任务指导用户使用Loader将文件从HDFS导入到ClickHouse。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HDFS目录和数据。
- ClickHouse相关表已创建，并确保用户已授权访问作业执行时操作该表的权限。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 使用Loader从HDFS导入数据时，确保HDFS输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符“\”“:”“;”中的任意字符。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

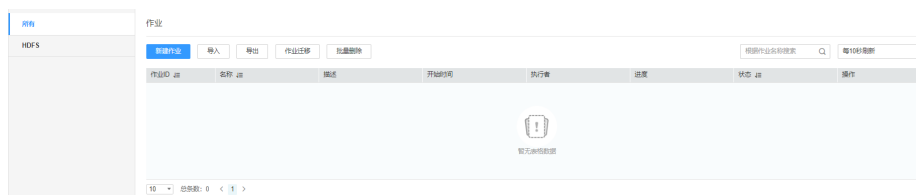
### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

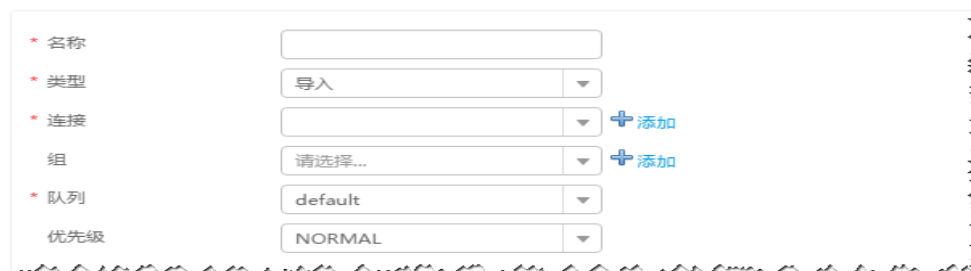
1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-41 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-42 “基本信息”界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导入”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“hdfs-connector”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，设置数据源信息。

**表 17-69** 输入设置参数

参数名	说明	示例
输入路径	HDFS中源文件的输入路径。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/use r/ test
路径过滤器	配置通配符对源文件的输入路径包含的目录进行过滤。“输入路径”不参与过滤。配置多个过滤条件时使用“,”隔开，配置为空时表示不过滤目录。不支持正则表达式过滤。	*
文件过滤器	配置通配符对源文件的输入文件名进行过滤。配置多个过滤条件时使用“,”隔开。不能配置为空。不支持正则表达式过滤。	*
编码类型	源文件的编码格式，如UTF-8。导入文本文件时才能配置。	UT F-8
后缀名	源文件导入成功后对输入文件增加的后缀值。该值为空，表示不加后缀。	.log

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-70](#)。

**表 17-70** 算子输入、输出参数设置

输入类型	输出类型
CSV文件输入	ClickHouse输出

图 17-43 算子操作方法示意



设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，根据实际场景在“存储类型”选择“CLICKHOUSE”，设置数据保存方式。

表 17-71 输出设置参数

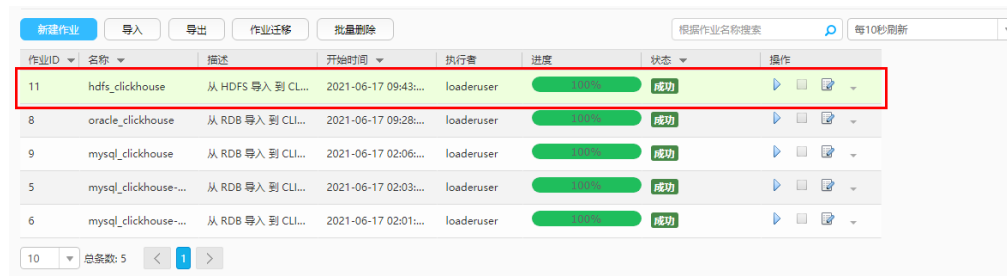
存储类型	参数名	说明	示例
CLICKHOUSE	ClickHouse实例	在ClickHouse作业中，Loader支持从集群可添加的所有ClickHouse服务实例中选择一个。如果选定的ClickHouse服务实例在集群中未添加，则此作业无法正常运行。	ClickHouse
	导入前清理数据	导入前清空原表的数据。“True”为执行清空，“False”为不执行。不配置此参数则默认不执行清空。 <b>说明</b> 如果导入的表为ClickHouse分布式表，且需要清理数据时，请在导入前手动删除ClickHouse分布式表对应的本地表中的数据。	false
	Map数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000。	20
	Map数据块大小	ClickHouse不支持此参数，请配置“Map数”。	-
	个数	Map任务的个数。	-

**步骤7** 单击“保存并运行”，开始保存并运行作业。

查看作业完成情况

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-44 查看作业



作业ID	名称	描述	开始时间	执行者	进度	状态	操作
11	hdfs_clickhouse	从 HDFS 导入到 CL...	2021-06-17 09:43:...	loaderuser	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">🗑️</a>
8	oracle_clickhouse	从 RDB 导入到 CLI...	2021-06-17 09:28:...	loaderuser	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">🗑️</a>
9	mysql_clickhouse	从 RDB 导入到 CLI...	2021-06-17 02:06:...	loaderuser	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">🗑️</a>
5	mysql_clickhouse...	从 RDB 导入到 CLI...	2021-06-17 02:03:...	loaderuser	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">🗑️</a>
6	mysql_clickhouse...	从 RDB 导入到 CLI...	2021-06-17 02:01:...	loaderuser	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">🗑️</a>

步骤9 使用ClickHouse客户端，查询ClickHouse表数据是否和HDFS导入的数据一致。

----结束

## 17.13 数据导出

### 17.13.1 概述

“数据导出”章节适用于MRS 3.x及后续版本。

#### 简介

Loader是实现MRS与关系型数据库、文件系统之间交换数据和文件的ETL工具，支持将数据或者文件从MRS系统中导出到关系型数据库或文件系统中。

Loader支持如下数据导出方式：

- 从HDFS/OBS中导出数据到SFTP服务器
- 从HDFS/OBS中导出数据到关系型数据库
- 从HBase中导出数据到SFTP服务器
- 从HBase中导出数据到关系型数据库
- 从Phoenix表导出数据到SFTP服务器
- 从Phoenix表导出数据到关系型数据库
- 从Hive中导出数据到SFTP服务器
- 从Hive中导出数据到关系数据库
- 从同一集群内HBase导出数据到HDFS/OBS

MRS与外部数据源交换数据和文件时需要连接数据源。系统提供以下连接器，用于配置不同类型数据源的连接参数：

- generic-jdbc-connector：关系型数据库连接器。
- hdfs-connector：HDFS数据源连接器。
- oracle-connector：Oracle数据库专用连接器，使用row\_id作为分区列，相对generic-jdbc-connector来说，Map任务分区更均匀，并且不依赖区分列是否有创建索引。
- mysql-fastpath-connector：MySQL数据库专用连接器，使用MySQL的mysqldump和mysqlimport工具进行数据的导入导出，相对generic-jdbc-connector来说，导入导出速度更快。



- sftp-connector: SFTP数据源连接器。
- oracle-partition-connector: 支持Oracle分区特性的连接器，专门对Oracle分区表的导入导出进行优化。

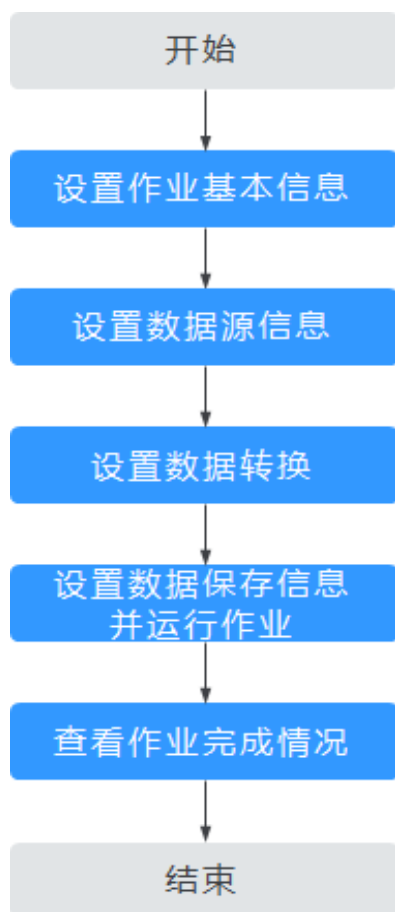
#### 📖 说明

- 建议将SFTP服务器和数据库服务器与Loader部署在独立的子网中，以保障数据安全地导出。
- 与关系数据库连接时，可以选择通用数据库连接器（generic-jdbc-connector）或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlimport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。
- 使用oracle-connector时，要求给连接用户赋予如下系统表或者视图的select权限：  
dba\_tab\_partitions、dba\_constraints、dba\_tables、dba\_segments、v\$instance、dba\_objects、v\$instance、dba\_extents、dba\_tab\_partitions、dba\_tab\_subpartitions。
- 使用oracle-partition-connector时，要求给连接用户赋予如下系统表的select权限：  
dba\_objects、dba\_extents。

## 导出流程

用户通过Loader界面进行数据导出作业，导出流程如图17-45所示。

图 17-45 导出流程示意



用户也可以通过Shell脚本来更新与运行Loader作业。该方式需要对已安装的Loader客户端进行配置。

## 17.13.2 使用 Loader 导出数据

### 操作场景

该任务指导用户完成将数据从MRS导出到外部的数据源的工作。

一般情况下，用户可以手工在Loader界面管理数据导入导出作业。当用户需要通过shell脚本来更新与运行Loader作业时，必须对已安装的Loader客户端进行配置。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的目录、HBase表和数据。
- 获取外部数据源（SFTP服务器或关系型数据库）使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 使用Loader从HDFS/OBS导出数据时，确保HDFS/OBS数据源的输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符“\":;”中的任意字符。
- 如果设置的任务需要使用指定Yarn队列功能，该用户需要已授权有相关Yarn队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

### 操作步骤

**步骤1** 是否第一次从Loader导出数据到关系型数据库？

- 是，执行**步骤2**。
- 否，执行**步骤3**。

**步骤2** 修改关系型数据库对应的驱动jar包文件权限。

1. 获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：“\${BIGDATA\_HOME}/FusionInsight\_Porter\_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib”。
2. 使用root用户在主备节点分别执行以下命令修改权限：  

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib
chown omm:wheel jar包文件名
chmod 600 jar包文件名
```
3. 登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Loader > 更多 > 重启服务”输入管理员密码重启Loader服务。

**步骤3** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。

2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-46 Loader WebUI 界面



**步骤4** 创建Loader数据导出作业，单击“新建作业”，在“1.基本信息”选择所需要的作业类型，然后单击“下一步”。

1. “名称”输入作业的名称，“类型”选择“导出”即导出。
2. “连接”选择一个连接。默认没有已创建的连接，单击“添加”创建一个新的连接，完成后单击“测试”，测试是否可用，待提示成功后单击“确定”。

表 17-72 连接配置参数一览表

连接器类型	参数名	说明
generic-jdbc-connector	JDBC驱动程序类	JDBC驱动类名。
	JDBC连接字符串	JDBC连接字符串。
	用户名	连接数据库使用的用户名。
	密码	连接数据库使用的密码。
	JDBC连接属性	JDBC连接属性，单击“添加”手动添加。 - 名称：连接属性名 - 值：连接属性值
hdfs-connector	-	-
oracle-connector	JDBC连接字符串	用户连接数据库的连接字符串。
	用户名	连接数据库使用的用户名。
	密码	连接数据库使用的密码。
	连接属性	连接属性，单击“添加”手动添加。 - 名称：连接属性名 - 值：连接属性值
mysql-fastpath-connector	JDBC连接字符串	JDBC连接字符串。
	用户名	连接数据库使用的用户名。

连接器类型	参数名	说明
	密码	连接数据库使用的密码。
	连接属性	连接属性，单击“添加”手动添加。 - 名称：连接属性名 - 值：连接属性值
sftp-connector	Sftp服务器的IP	SFTP服务器的IP地址。
	Sftp服务器端口	SFTP服务器的端口号。
	Sftp用户名	访问SFTP服务器的用户名。
	Sftp密码	访问SFTP服务器的密码。
	Sftp公钥	Sftp服务器公钥。
oracle-partition-connector	JDBC驱动程序类	JDBC驱动类名。
	JDBC连接字符串	JDBC连接字符串。
	用户名	连接数据库使用的用户名。
	密码	连接数据库使用的密码。
	连接属性	连接属性，单击“添加”手动添加。 - 名称：连接属性名 - 值：连接属性值

- “组”设置“作业”所属组，默认没有已创建的组，单击“添加”创建一个新的组，单击“确定”保存。
- “队列”设置Loader的任务在指定的Yarn队列中执行。默认值“root.default”表示任务在“default”队列中执行。
- “优先级”设置Loader的任务在指定的Yarn队列中的优先级。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。默认值为“NORMAL”。

**步骤5** 在“2.输入设置”，设置数据来源，然后单击“下一步”。

#### 📖 说明

创建或者编辑Loader作业时，在配置SFTP路径、HDFS/OBS路径、SQL的Where条件等参数时，可以使用宏定义，具体请参考[配置项中使用宏定义](#)章节。

**表 17-73** 输入配置参数一览表

源文件类型	参数名	解释说明
HDFS/OBS	输入目录	从HDFS/OBS导出时的输入路径。

源文件类型	参数名	解释说明
	路径过滤器	配置通配符对源文件的输入路径包含的目录进行过滤。输入路径“输入目录”不参与过滤。配置多个过滤条件时使用逗号隔开，配置为空时表示不过滤目录。不支持正则表达式过滤。
	文件过滤器	配置通配符对源文件的输入文件名进行过滤。配置多个过滤条件时使用逗号隔开。不能配置为空。不支持正则表达式过滤。
	文件类型	文件导入类型： <ul style="list-style-type: none"> <li>TEXT_FILE：导入文本文件并保存为文本文件。</li> <li>SEQUENCE_FILE：导入文本文件并保存在 sequence file 文件格式。</li> <li>BINARY_FILE：以二进制流的方式导入文件，可以导入任何格式的文件。</li> </ul>
	文件分割方式	选择按FILE文件或SIZE大小分割源文件成多份，作为数据导出的MapReduce任务中各个map的输入文件。
	Map数	配置数据操作的MapReduce任务中同时启动的map数量。不可与“Map数据块大小”同时配置。参数值必须小于或等于“3000”。
	Map数据块大小	配置数据操作的MapReduce任务中启动map所处理的数据大小，单位为MB。参数值必须大于或等于“100”，建议配置值为“1000”。不可与“Map数”同时配置。当使用关系型数据库连接器时，不支持“Map数据块大小”，请配置“Map数”。
HBASE	HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。
	个数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于“3000”。
HIVE	Hive实例	在Hive作业中，Loader支持从集群可添加的所有Hive服务实例中选择任意一个。如果选定的Hive服务实例在集群中未添加，则此作业无法正常运行。
	个数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于“3000”。
SPARK	Spark实例	仅支持SparkSQL存取Hive数据。在SparkSQL作业中，Loader支持从集群可添加的所有Spark服务实例中选择任意一个。如果选定的Spark服务实例在集群中未添加，则此作业无法正常运行。

源文件类型	参数名	解释说明
	个数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于“3000”。

**步骤6** 在“3.转换”设置数据传输过程中的转换操作。

确认Loader创建的数据操作作业中，源数据的值是否满足直接使用需求而不进行转换，例如大小写转换、截取、拼接和分隔。

- 满足需求，请单击“下一步”。
  - 不满足需求，请执行[步骤6.1](#) ~ [步骤6.4](#)。
1. 默认没有已创建的转换步骤，可拖动左侧样例到编辑框，添加一个新的转换步骤。
  2. 完整的转换流程包含以下类型，每个类型请根据业务需要进行选择。
    - a. 输入类型，第一个转换步骤，仅添加一种，任务涉及HBase或关系型数据库必须添加。
    - b. 转换类型，中间转换步骤，可添加一种以上或不添加。
    - c. 输出类型，最后一个转换步骤，仅添加一种，任务涉及HBase或关系型数据库必须添加。

**表 17-74** 样例一览表

类型	描述
输入类型	<ul style="list-style-type: none"> <li>▪ CSV文件输入：CSV文件输入步骤，配置分隔符以转换生成多个字段。</li> <li>▪ 固定宽度文件输入：文本文件输入步骤，配置截取字符或字节的长度以转换生成多个字段。</li> <li>▪ 表输入：关系型数据输入步骤，配置数据库的指定列为输入的字段。</li> <li>▪ HBase输入：HBase表输入步骤，配置HBase表的列定义到指定字段。</li> <li>▪ HTML输入：HTML网页数据输入步骤，配置获取HTML网页文件目标数据到指定字段。</li> <li>▪ Hive输入：Hive表输入步骤，配置Hive表的列定义到指定字段。</li> <li>▪ Spark输入：SparkSQL表输入步骤，配置SparkSQL表的列定义到指定字段。仅支持存取Hive数据。</li> </ul>

类型	描述
转换类型	<ul style="list-style-type: none"><li>▪ 长整型时间转换：长整型日期转换步骤，配置长整型数值与日期的转换。</li><li>▪ 空值转换：空值转换步骤，配置指定值替换空值。</li><li>▪ 随机值转换：随机数据生成步骤，配置新增值为随机数据的字段。</li><li>▪ 增加常量字段：增加常量步骤，配置直接生成常量字段。</li><li>▪ 拼接转换：拼接字段步骤，配置已生成的字段通过连接符连接，转换出新的字段。</li><li>▪ 分隔转换：分隔字段步骤，配置已生成的字段通过分隔符分隔，转换出新的字段。</li><li>▪ 取模转换：取模运算步骤，配置已生成的字段通过取模，转换出新的字段。</li><li>▪ 剪切字符串：字符串截取步骤，配置已生成的字段通过指定位置截取，转换出新的字段。</li><li>▪ EL操作转换：计算器，可以对字段值进行运算，目前支持的算子有：md5sum、sha1sum、sha256sum和sha512sum等。</li><li>▪ 字符串大小写转换：字符串转换步骤，配置已生成的字段通过大小写变换，转换出新的字段。</li><li>▪ 字符串逆序转换：字符串逆序步骤，配置已生成的字段通过逆序，转换出新的字段。</li><li>▪ 字符串空格清除转换：字符串空格清除步骤，配置已生成的字段通过清除空格，转换出新的字段。</li><li>▪ 过滤行转换：过滤行步骤，配置逻辑条件过滤掉含触发条件的行。</li><li>▪ 更新域：更新域步骤，配置当满足某些条件时，更新指定字段的值。</li></ul>

类型	描述
输出类型	<ul style="list-style-type: none"> <li>▪ 文件输出：文本文件输出步骤，配置已生成的字段通过分隔符连接并输出到文件。</li> <li>▪ 表输出：关系型数据库输出步骤，配置输出的字段对应到数据库的指定列。</li> <li>▪ HBase输出：HBase表输出步骤，配置已生成的字段输出到HBase表的列。</li> <li>▪ Hive输出：Hive表输出步骤，配置已生成的字段输出到Hive表的列。</li> <li>▪ Spark输出：SparkSQL表输出步骤，配置已生成的字段输出到SparkSQL表的列。仅支持存取Hive数据。</li> </ul>

编辑栏包括以下几种任务：

- 重命令：重命名样例。
- 编辑：编辑步骤转换，参考[步骤6.3](#)。
- 删除：删除样例。

#### 说明

也可使用快捷键“Del”删除。

3. 单击“编辑”，编辑步骤转换信息，配置字段与数据。

步骤转换信息中的具体参数设置请参考[算子帮助](#)。

转换步骤配置不正确时，传输的数据将无法转换并成为脏数据，脏数据标记规则如下：

- 任意输入类型步骤中，原数据包含字段的个数小于配置字段的个数，或者原数据字段值与配置字段的类型不匹配时，全部数据成为脏数据。
- “CSV文件输入”步骤中，“验证输入字段”检验输入字段与值的类型匹配情况，检查不匹配时跳过该行，当前行成为脏数据。
- “固定宽度文件输入”步骤中，“固定长度”指定字段分割长度，长度大于原字段值的长度则数据分割失败，当前行成为脏数据。
- “HBase输入”步骤中，“HBase表名”指定HBase表名不正确，或者“主键”没有配置主键列，全部数据成为脏数据。
- 任意转换类型步骤中，转换失败的行成为脏数据。例如“分隔转换”步骤中，生成的字段个数小于配置字段的个数，或者原数据不能转换为String类型，当前行成为脏数据。
- “过滤行转换”步骤中，被筛选条件过滤的行成为脏数据。
- “取模转换”步骤中，原字段值为“NULL”，当前行成为脏数据。

4. 单击“下一步”。

**步骤7** 在“4.输出设置”，设置数据保存目标位置，然后单击“保存”保存作业或“保存并运行”，保存作业并运行作业。



表 17-75 输出配置参数一览表

连接类型	参数名	解释说明
sftp-connector	输出路径	SFTP服务器中导出文件的路径或者文件名，如果连接器配置多个地址此处可对应使用分号分隔多个路径或者文件名，数量需要与连接器中服务器的数量一致。
	文件操作方式	数据导入时的操作行为。全部数据从输入路径导入到目标路径时，先保存在临时目录，然后再从临时目录复制转移至目标路径，任务完成时删除临时路径的文件。转移临时文件存在同名文件时有以下行为： <ul style="list-style-type: none"> <li>“OVERRIDE”：直接覆盖旧文件。</li> <li>“RENAME”：重命名新文件。无扩展名的文件直接增加字符串后缀，有扩展名的文件在文件名增加字符串后缀。字符串具有唯一性。</li> <li>“APPEND”：在旧文件尾部合并新文件内容。合并操作只是简单的追加，不保证追加文件是否可以使用。例如文本文件可合并，压缩文件合并后可能无法使用。</li> <li>“IGNORE”：保留旧文件，不复制新文件。</li> <li>“ERROR”：转移过程中出现同名文件时任务将停止执行并报错，已转移的文件导入成功，同名的文件及未转移的文档导入失败。</li> </ul>
	编码类型	导出文件的编码格式，如UTF-8。导出文本文件时才能配置。
	压缩	使用SFTP协议导出数据时，是否开启压缩传输功能以减小带宽使用。“true”为开启压缩，“false”为关闭压缩。
hdfs-connector	输出路径	导出文件在HDFS/OBS的输出目录或者文件名。
	文件格式	文件导出类型： <ul style="list-style-type: none"> <li>“TEXT_FILE”：导入文本文件并保存为文本文件。</li> <li>“SEQUENCE_FILE”：导入文本文件并保存在sequence file文件格式。</li> <li>“BINARY_FILE”：以二进制流的方式导入文件，可以导入任何格式的文件。</li> </ul>
	压缩格式	在下拉菜单中选择数据导出到HDFS/OBS后保存文件的压缩格式，未配置或选择NONE表示不压缩数据。

连接类型	参数名	解释说明
	自定义压缩格式	自定义压缩格式类型的名称。
generic-jdbc-connector	架构名称	数据库模式名。
	表名	数据库表名，用于最终保存传输的数据。
	临时表	数据库临时表的表名，用于临时保存传输过程中的数据，字段需要和“表名”配置的表一致。
oracle-partition-connector	架构名称	数据库模式名。
	表名	数据库表名，用于最终保存传输的数据。
	临时表	数据库临时表的表名，用于临时保存传输过程中的数据，字段需要和“表名”配置的表一致。
oracle-connector	表名	目标表，用于存储数据。
	列名	指定要写入的列名，没有指定的列允许被设置为空值或者默认值。
mysql-fastpath-connector	架构名称	数据库模式名。
	表名	数据库表名，用于最终保存传输的数据。
	临时表名	临时表名称，用于预存数据，作业执行成功后，再将数据转移到正式表。

**步骤8** 已创建的作业可以在“Loader WebUI”界面上进行浏览，可进行启动、停止、复制、删除、编辑和查看历史信息操作。

图 17-47 查看 Loader 作业

作业ID	名称	描述	开始时间	执行者	进度	状态	操作
23	THbase	从 SFTP 导入到 HBASE...	2016-02-01 11:35:32	admin	100%	成功	▶ □ 📄 ⚙️
9	sftp-hdfs-updatefields	从 SFTP 导入到 HDFS	2016-01-30 17:39:08	admin	100%	成功	▶ □ 📄 ⚙️
15	sftp-hdfs	从 SFTP 导入到 HDFS	2016-01-30 16:10:46	admin	100%	成功	▶ □ 📄 ⚙️
18	hdfs-voltdb	从 HDFS 导出到 VoltDB	2016-01-30 16:09:54	admin	100%	成功	▶ □ 📄 ⚙️
20	hdfs-voltdb-batch	从 HDFS 导出到 VoltDB	2016-01-30 15:51:10	admin	100%	成功	▶ □ 📄 ⚙️
21	hdfs-voltdb-66	从 HDFS 导出到 VoltDB	2016-01-30 15:45:18	admin	100%	成功	▶ □ 📄 ⚙️

----结束

### 17.13.3 典型场景：从 HDFS/OBS 导出数据到 SFTP 服务器

#### 操作场景

该任务指导用户使用Loader将数据从HDFS/OBS导出到SFTP服务器。

## 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HDFS/OBS目录和数据。
- 获取SFTP服务器使用的用户和密码，且该用户具备SFTP服务器数据导出目录的写入权限。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 使用Loader从HDFS/OBS导出数据时，确保HDFS/OBS数据源的输入路径目录名、输入路径的子目录名及子文件名不能包含特殊字符“\”“:”“;”中的任意字符。
- 如果设置的任务需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

## 操作步骤

### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-48 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-49 基本信息界面

The screenshot shows the 'Basic Information' form for creating a job. The form has the following fields:

- \* 名称 (Name): A text input field.
- \* 类型 (Type): A dropdown menu with '导出' (Export) selected.
- \* 连接 (Connection): A dropdown menu with a '+ 添加' (Add) button next to it.
- 组 (Group): A dropdown menu with '请选择...' (Please select...) and a '+ 添加' (Add) button next to it.
- \* 队列 (Queue): A dropdown menu with 'default' selected.
- 优先级 (Priority): A dropdown menu with 'NORMAL' selected.

1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导出”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。

4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“sftp-connector”，单击“添加”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。Loader支持配置多个SFTP服务器操作数据，单击“添加”可增加多行SFTP服务器的配置信息。

表 17-76 连接参数

参数名	说明	示例
名称	SFTP服务器连接的名称。	sftpName
Sftp服务器的IP	SFTP服务器的IP地址。	10.16.0.1
Sftp服务器端口	SFTP服务器的端口号。	22
Sftp用户名	访问SFTP服务器的用户名。	root
Sftp密码	访问SFTP服务器的密码。	xxxx
Sftp公钥	Sftp服务器公钥。	OdDt/yn...etM

### 说明

配置多个SFTP服务器时，HDFS/OBS的数据将分为多份随机导出到各个SFTP服务器。

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，在“源文件类型”中选择“HDFS”，设置数据源信息。

表 17-77 数据来源配置参数

参数名	解释说明	示例
输入目录	从HDFS/OBS导出时的输入路径。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/use r/ test
路径过滤器	配置通配符对源文件的输入路径包含的目录进行过滤。“输入目录”不参与过滤。配置多个过滤条件时使用“,”隔开，配置为空时表示不过滤目录。不支持正则表达式过滤。 <ul style="list-style-type: none"> <li>● “?” 匹配单个字符。</li> <li>● “*” 配置多个字符。</li> <li>● 在匹配条件前加“^”表示取反，即文件过滤。</li> </ul>	*

参数名	解释说明	示例
文件过滤器	<p>配置通配符对源文件的输入文件名进行过滤。配置多个过滤条件时使用“,” 隔开。不能配置为空。不支持正则表达式过滤。</p> <ul style="list-style-type: none"> <li>“?” 匹配单个字符。</li> <li>“*” 配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul>	*
文件类型	<p>文件导入类型：</p> <ul style="list-style-type: none"> <li>“TEXT_FILE”：导入文本文件并保存为文本文件</li> <li>“SEQUENCE_FILE”：导入文本文件并保存在“sequence file”文件格式</li> <li>“BINARY_FILE”：以二进制流的方式导入文件，可以导入任何格式的文件，不对文件做任何处理。</li> </ul> <p><b>说明</b> 文件类型选择“TEXT_FILE”或“SEQUENCE_FILE”导入时，Loader会自动根据文件的后缀选择对应的解压方法，对文件进行解压。</p>	TEXT_FILE
文件分割方式	<p>选择按文件或大小分割源文件，作为数据导出的MapReduce任务中各个map的输入文件。</p> <ul style="list-style-type: none"> <li>选择“FILE”，表示按文件分割源文件，即每个map处理一个或多个完整的源文件，同一个源文件不可分配至不同map，完成数据导入后保持源文件的目录结构。</li> <li>选择“SIZE”，表示按大小分割源文件，即每个map处理一定大小的输入文件，同一个源文件可分割至多个map，数据保存至输出目录时保存的文件数与map数量相同，文件名格式为“import_part_xxxx”，“xxxx”为系统生成的随机数，具有唯一性。</li> </ul>	FILE
Map数	<p>配置数据操作的MapReduce任务中同时启动的Map数量。不可与“Map数据块大小”同时配置。参数值必须小于或等于3000，建议以SFTP服务器的CPU的核数作为其取值。</p> <p><b>说明</b> 为了提高导入数据速度，需要确保以下条件：</p> <ul style="list-style-type: none"> <li>每个Map连接时，相当于一个客户端连接，因此需要确保SFTP服务器最大连接数大于Map数量。</li> <li>确保SFTP服务器上的磁盘IO或网络带宽都未达到上限。</li> </ul>	20
Map数据块大小	<p>配置数据操作的MapReduce任务中启动map所处理的数据大小，单位为MB。参数值必须大于或等于100，建议配置值为1000。不可与“Map数”同时配置。</p>	-

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-78](#)。

**表 17-78** 算子输入、输出参数设置

输入类型	输出类型
CSV文件输入	文件输出
HTML输入	文件输出
固定宽度文件输入	文件输出

**图 17-50** 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，设置数据保存方式。

**表 17-79** 输出设置参数

参数名	解释说明	示例
输出路径	SFTP服务器中导出文件的路径或者文件名，如果连接器配置多个地址此处可对应使用“;”分隔多个路径或者文件名，数量需要与连接器中服务器的数量一致。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/opt/ tempfile

参数名	解释说明	示例
文件操作方式	<p>数据导入时的操作行为。全部数据从输入路径导入到目标路径时，先保存在临时目录，然后再从临时目录复制转移至目标路径，任务完成时删除临时路径的文件。转移临时文件存在同名文件时有以下行为：</p> <ul style="list-style-type: none"> <li>“OVERRIDE”：直接覆盖旧文件。</li> <li>“RENAME”：重命名新文件。无扩展名的文件直接增加字符串后缀，有扩展名的文件在文件名增加字符串后缀。字符串具有唯一性。</li> <li>“APPEND”：在旧文件尾部合并新文件内容。合并操作只是简单的追加，不保证追加文件是否可以使用。例如文本文件可合并，压缩文件合并后可能无法使用。</li> <li>“IGNORE”：保留旧文件，不复制新文件。</li> <li>“ERROR”：转移过程中出现同名文件时任务将停止执行并报错，已转移的文件导入成功，同名的文件及未转移的文档导入失败。</li> </ul>	OVERRIDE
编码类型	导出文件的编码格式，如UTF-8。导出文本文件时才能配置。	UTF-8
压缩	<p>使用SFTP协议导入数据时，是否开启压缩传输功能以减小带宽使用。</p> <ul style="list-style-type: none"> <li>选择“true”，表示开启压缩。</li> <li>选择“false”，表示关闭压缩。</li> </ul>	true

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-51 查看作业

作业ID	名称	描述	开始时间	执行者	进度	状态	操作
5	hdfs-sftp-size	从 HDFS 导出到 SFTP	2016-03-26 14:45:08	admin	<div style="width: 100%;"></div>	成功	<span>▶</span> <span>📄</span> <span>🗑️</span>
4	hdfs-sftp	从 HDFS 导出到 SFTP	2016-03-26 14:43:14	admin	<div style="width: 100%;"></div>	成功	<span>▶</span> <span>📄</span> <span>🗑️</span>
3	sftp-hdfs-1-1	从 SFTP 导入到 HDFS	2016-03-26 12:00:36	admin	<div style="width: 100%;"></div>	成功	<span>▶</span> <span>📄</span> <span>🗑️</span>
1	sftp-hdfs	从 SFTP 导入到 HDFS	2016-03-26 11:08:09	admin	<div style="width: 100%;"></div>	成功	<span>▶</span> <span>📄</span> <span>🗑️</span>
2	sftp-hdfs-1	从 SFTP 导入到 HDFS	2016-03-26 11:07:56	admin	<div style="width: 100%;"></div>	成功	<span>▶</span> <span>📄</span> <span>🗑️</span>

----结束

## 17.13.4 典型场景：从 HBase 导出数据到 SFTP 服务器

### 操作场景

该任务指导用户使用Loader将数据从HBase导出到SFTP服务器。

## 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HBase表或phoenix表。
- 获取SFTP服务器使用的用户和密码，且该用户具备SFTP服务器数据导出目录的写入权限。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的任务需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

## 操作步骤

### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-52 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-53 基本信息界面

The screenshot shows the 'Basic Information' form for creating a job. It includes the following fields:

- \* 名称 (Name): A text input field.
- \* 类型 (Type): A dropdown menu with '导出' (Export) selected.
- \* 连接 (Connection): A dropdown menu with a '+ 添加' (Add) button next to it.
- 组 (Group): A dropdown menu with '请选择...' (Please select...) and a '+ 添加' (Add) button next to it.
- \* 队列 (Queue): A dropdown menu with 'default' selected.
- 优先级 (Priority): A dropdown menu with 'NORMAL' selected.

1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导出”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。



- 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“sftp-connector”，单击“添加”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。Loader支持配置多个SFTP服务器操作数据，单击“添加”可增加多行SFTP服务器的配置信息。

**表 17-80** 连接参数

参数名	说明	示例
名称	SFTP服务器连接的名称。	sftpName
Sftp服务器的IP	SFTP服务器的IP地址。	10.16.0.1
Sftp服务器端口	SFTP服务器的端口号。	22
Sftp用户名	访问SFTP服务器的用户名。	root
Sftp密码	访问SFTP服务器的密码。	xxxx
Sftp公钥	Sftp服务器公钥。	OdDt/yn...etM

#### 说明

配置多个SFTP服务器时，HBase表或phoenix表将分成多份随机保存到各个SFTP服务器。

#### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，在“源文件类型”中选择“HBASE”，设置数据源信息。

**表 17-81** 数据源配置参数说明

参数名	解释说明	示例
HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HB ase
个数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000，建议以SFTP服务器当前最大连接数作为其取值。	20

#### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-82](#)。

表 17-82 算子输入、输出参数设置

输入类型	输出类型
HBase输入	文件输出

图 17-54 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，设置数据保存方式。

表 17-83 输出设置参数

参数名	解释说明	示例
输出路径	SFTP服务器中导出文件的路径或者文件名，如果连接器配置多个地址此处可对应使用“;”分隔多个路径或者文件名，数量需要与连接器中服务器的数量一致。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/opt/ tem pfil e

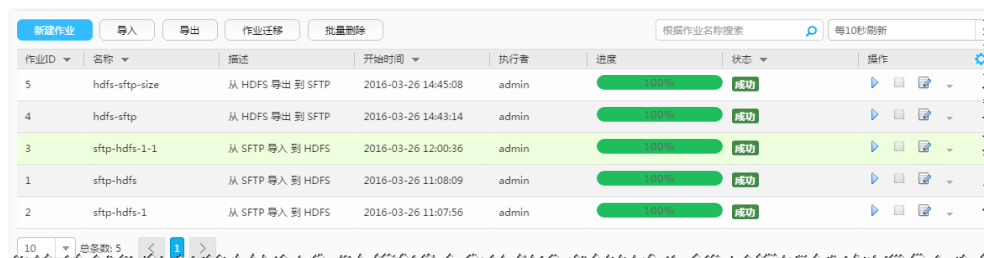
参数名	解释说明	示例
文件操作方式	<p>数据导入时的操作行为。全部数据从输入路径导入到目标路径时，先保存在临时目录，然后再从临时目录复制转移至目标路径，任务完成时删除临时路径的文件。转移临时文件存在同名文件时有以下行为：</p> <ul style="list-style-type: none"> <li>“OVERRIDE”：直接覆盖旧文件。</li> <li>“RENAME”：重命名新文件。无扩展名的文件直接增加字符串后缀，有扩展名的文件在文件名增加字符串后缀。字符串具有唯一性。</li> <li>“APPEND”：在旧文件尾部合并新文件内容。合并操作只是简单的追加，不保证追加文件是否可以使用。例如文本文件可合并，压缩文件合并后可能无法使用。</li> <li>“IGNORE”：保留旧文件，不复制新文件。</li> <li>“ERROR”：转移过程中出现同名文件时任务将停止执行并报错，已转移的文件导入成功，同名的文件及未转移的文档导入失败。</li> </ul>	OVERRIDE
编码类型	导出文件的编码格式，如UTF-8。导出文本文件时才能配置。	UTF-8
压缩	<p>使用SFTP协议导入数据时，是否开启压缩传输功能以减少带宽使用。</p> <ul style="list-style-type: none"> <li>选择“true”，表示开启压缩。</li> <li>选择“false”，表示关闭压缩。</li> </ul>	true

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

**图 17-55 查看作业**



----结束

## 17.13.5 典型场景：从 Hive 导出数据到 SFTP 服务器

### 操作场景

该任务指导用户使用Loader将数据从Hive导出到SFTP服务器。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业中指定的Hive表的权限。
- 获取SFTP服务器使用的用户和密码，且该用户具备SFTP服务器数据导出目录的写入权限。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的任务需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-56 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-57 基本信息界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导出”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“sftp-connector”，单击“添加”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。Loader支持配置多个SFTP服务器操作数据，单击“添加”可增加多行SFTP服务器的配置信息。

表 17-84 连接参数

参数名	说明	示例
名称	SFTP服务器连接的名称。	sftpName
Sftp服务器的IP	SFTP服务器的IP地址。	10.16.0.1
Sftp服务器端口	SFTP服务器的端口号。	22
Sftp用户名	访问SFTP服务器的用户名。	root
Sftp密码	访问SFTP服务器的密码。	xxxx
Sftp公钥	Sftp服务器公钥。	OdDt/yn...etM

#### 说明

配置多个SFTP服务器时，Hive表将分成多份随机保存到各个SFTP服务器。

#### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，在“源文件类型”中选择“HIVE”，设置数据源信息。

表 17-85 数据源配置参数说明

参数名	解释说明	示例
Hive实例	在Hive作业中，Loader支持从集群可添加的所有Hive服务实例中选择任意一个。如果选定的Hive服务实例在集群中未添加，则此作业无法正常运行。	hive
个数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000，建议以SFTP服务器当前最大连接数作为其取值。	20

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-86](#)。

**表 17-86** 算子输入、输出参数设置

输入类型	输出类型
Hive输入	文件输出

**图 17-58** 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，设置数据保存方式。

**表 17-87** 输出设置参数

参数名	解释说明	示例
输出路径	SFTP服务器中导出文件的路径或者文件名，如果连接器配置多个地址此处可对应使用“;”分隔多个路径或者文件名，数量需要与连接器中服务器的数量一致。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/opt/ tempfile

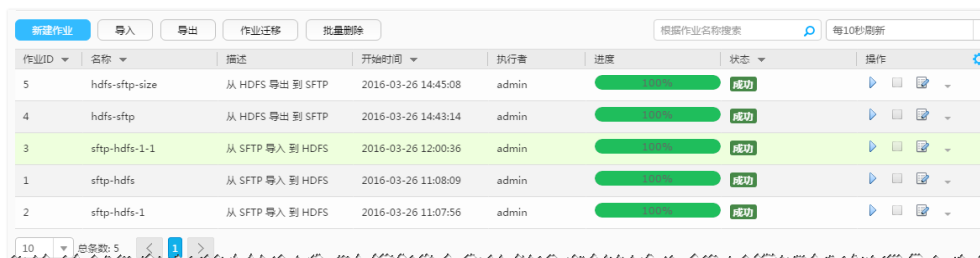
参数名	解释说明	示例
文件操作方式	<p>数据导入时的操作行为。全部数据从输入路径导入到目标路径时，先保存在临时目录，然后再从临时目录复制转移至目标路径，任务完成时删除临时路径的文件。转移临时文件存在同名文件时有以下行为：</p> <ul style="list-style-type: none"> <li>“OVERRIDE”：直接覆盖旧文件。</li> <li>“RENAME”：重命名新文件。无扩展名的文件直接增加字符串后缀，有扩展名的文件在文件名增加字符串后缀。字符串具有唯一性。</li> <li>“APPEND”：在旧文件尾部合并新文件内容。合并操作只是简单的追加，不保证追加文件是否可以使用。例如文本文件可合并，压缩文件合并后可能无法使用。</li> <li>“IGNORE”：保留旧文件，不复制新文件。</li> <li>“ERROR”：转移过程中出现同名文件时任务将停止执行并报错，已转移的文件导入成功，同名的文件及未转移的文档导入失败。</li> </ul>	OVERRIDE
编码类型	导出文件的编码格式，如UTF-8。导出文本文件时才能配置。	UTF-8
压缩	<p>使用SFTP协议导入数据时，是否开启压缩传输功能以减少带宽使用。</p> <ul style="list-style-type: none"> <li>选择“true”，表示开启压缩。</li> <li>选择“false”，表示关闭压缩。</li> </ul>	true

**步骤7** 单击“保存并运行”，开始保存并运行作业。

**查看作业完成情况**

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

**图 17-59 查看作业**



----结束

## 17.13.6 典型场景：从 HDFS/OBS 导出数据到关系型数据库

### 操作场景

该任务指导用户使用Loader将数据从HDFS/OBS导出到关系型数据库。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HDFS/OBS目录和数据。
- 获取关系型数据库使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。
- 操作前需要进行如下配置：
  - a. 获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：  
“`${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib`”。
  - b. 使用root用户在主备节点分别执行以下命令修改权限：  

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/
FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/
webapps/loader/WEB-INF/ext-lib
chown omm:wheel jar包文件名
chmod 600 jar包文件名
```
  - c. 登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Loader > 更多 > 重启服务”，输入管理员密码重启Loader服务。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-60 Loader WebUI 界面





**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

**图 17-61** 基本信息界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导出”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“generic-jdbc-connector”或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

**说明**

- 与关系数据库连接时，可以选择通用数据库连接器（generic-jdbc-connector）或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlexport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。

**表 17-88** “generic-jdbc-connector” 连接参数

参数名	说明	示例
名称	关系型数据库连接的名称。	dbName
JDBC驱动程序类	JDBC驱动类名。	oracle.jdbc.driver.OracleDriver
JDBC连接字符串	JDBC连接字符串。	jdbc:oracle:thin:@//10.16.0.1:1521/oradb
用户名	连接数据库使用的用户名。	omm
密码	连接数据库使用的密码。	xxxx

参数名	说明	示例
JDBC连接属性	JDBC连接属性，单击“添加”手动添加。 <ul style="list-style-type: none"> <li>名称：连接属性名</li> <li>值：连接属性值</li> </ul>	<ul style="list-style-type: none"> <li>名称：socketTimeout</li> <li>值：20</li> </ul>

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，在“源文件类型”中选择“HDFS”，设置数据源信息。

表 17-89 数据来源配置参数

参数名	解释说明	示例
输入目录	从HDFS/OBS导出时的输入路径。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/use r/ test
路径过滤器	配置通配符对源文件的输入路径包含的目录进行过滤。“输入目录”不参与过滤。配置多个过滤条件时使用“,”隔开，配置为空时表示不过滤目录。不支持正则表达式过滤。 <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul>	*
文件过滤器	配置通配符对源文件的输入文件名进行过滤。配置多个过滤条件时使用“,”隔开。不能配置为空。不支持正则表达式过滤。 <ul style="list-style-type: none"> <li>“?”匹配单个字符。</li> <li>“*”配置多个字符。</li> <li>在匹配条件前加“^”表示取反，即文件过滤。</li> </ul>	*
文件类型	文件导入类型： <ul style="list-style-type: none"> <li>“TEXT_FILE”：导入文本文件并保存为文本文件。</li> <li>“SEQUENCE_FILE”：导入文本文件并保存在sequence file文件格式。</li> <li>“BINARY_FILE”：以二进制流的方式导入文件，可以导入任何格式的文件，不对文件做任何处理。</li> </ul> <b>说明</b> 文件类型选择“TEXT_FILE”或“SEQUENCE_FILE”导入时，Loader会自动根据文件的后缀选择对应的解压方法，对文件进行解压。	TEXT_F ILE

参数名	解释说明	示例
文件分割方式	<p>选择按文件或大小分割源文件，作为数据导出的 MapReduce 任务中各个 map 的输入文件。</p> <ul style="list-style-type: none"> <li>选择“FILE”，表示按文件分割源文件，即每个 map 处理一个或多个完整的源文件，同一个源文件不可分配至不同 map，完成数据导入后保持源文件的目录结构。</li> <li>选择“SIZE”，表示按大小分割源文件，即每个 map 处理一定大小的输入文件，同一个源文件可分割至多个 map，数据保存至输出目录时保存的文件数与 map 数量相同，文件名格式为“import_part_xxxx”，“xxxx”为系统生成的随机数，具有唯一性。</li> </ul>	FILE
Map数	配置数据操作的 MapReduce 任务中同时启动的 Map 数量。不可与“Map 数据块大小”同时配置。参数值必须小于或等于 3000。	20
Map 数据块大小	配置数据操作的 MapReduce 任务中启动 map 所处理的数据大小，单位为 MB。参数值必须大于或等于 100，建议配置值为 1000。不可与“Map 数”同时配置。当使用关系型数据库连接器时，不支持“Map 数据块大小”，请配置“Map 数”。	-

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-90](#)。

**表 17-90** 算子输入、输出参数设置

输入类型	输出类型
CSV文件输入	表输出
HTML输入	表输出
固定宽度文件输入	表输出

图 17-62 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，设置数据保存方式。

表 17-91 输出设置参数

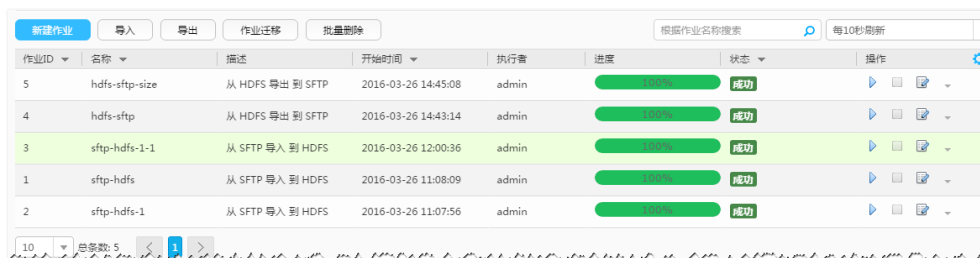
参数名	说明	示例
架构名称	数据库模式名。	dbo
表名	数据库表名，用于最终保存传输的数据。 <b>说明</b> 表名可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	test
临时表	数据库临时表表名，用于临时保存传输过程中的数据，字段需要和“表名”配置的表一致。 <b>说明</b> 使用临时表是为了使得导出数据到数据库时，不会在目的表中产生脏数据。只有在所有数据成功写入临时表后，才会将数据从临时表迁移到目的表。使用临时表会增加作业的执行时间。	tm p_ est

**步骤7** 单击“保存并运行”，开始保存并运行作业。

### 查看作业完成情况

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-63 查看作业



作业ID	名称	描述	开始时间	执行者	进度	状态	操作
5	hdfs-sftp-size	从 HDFS 导出到 SFTP	2016-03-26 14:45:08	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">⌵</a>
4	hdfs-sftp	从 HDFS 导出到 SFTP	2016-03-26 14:43:14	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">⌵</a>
3	sftp-hdfs-1-1	从 SFTP 导入到 HDFS	2016-03-26 12:00:36	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">⌵</a>
1	sftp-hdfs	从 SFTP 导入到 HDFS	2016-03-26 11:08:09	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">⌵</a>
2	sftp-hdfs-1	从 SFTP 导入到 HDFS	2016-03-26 11:07:56	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">📄</a> <a href="#">⌵</a>

----结束

## 17.13.7 典型场景：从 HBase 导出数据到关系型数据库

### 操作场景

该任务指导用户使用Loader将数据从HBase导出到关系型数据库。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HBase表或phoenix表。
- 获取关系型数据库使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。
- 操作前需要进行如下配置：
  - a. 获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：“\${BIGDATA\_HOME}/FusionInsight\_Porter\_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib”。
  - b. 使用root用户在主备节点分别执行以下命令修改权限：

```
cd ${BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/
FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/
webapps/loader/WEB-INF/ext-lib
chown omm:wheel jar包文件名
chmod 600 jar包文件名
```
  - c. 登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Loader > 更多 > 重启服务”，输入管理员密码重启Loader服务。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-64 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-65 基本信息界面



1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导出”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“generic-jdbc-connector”或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

#### 说明

- 与关系数据库连接时，可以选择通用数据库连接器（generic-jdbc-connector）或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlimport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。

表 17-92 “generic-jdbc-connector” 连接参数

参数名	说明	示例
名称	关系型数据库连接的名称。	dbName
JDBC驱动程序类	JDBC驱动类名。	oracle.jdbc.driver.OracleDriver
JDBC连接字符串	JDBC连接字符串。	jdbc:oracle:thin:@//10.16.0.1:1521/oradb
用户名	连接数据库使用的用户名。	omm
密码	连接数据库使用的密码。	xxxx
JDBC连接属性	JDBC连接属性，单击“添加”手动添加。 <ul style="list-style-type: none"> <li>名称：连接属性名</li> <li>值：连接属性值</li> </ul>	<ul style="list-style-type: none"> <li>名称：socketTimeout</li> <li>值：20</li> </ul>

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，在“源文件类型”中选择“HBASE”，设置数据源信息。

表 17-93 数据源配置参数说明

参数名	解释说明	示例
HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HB ase
个数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于“3000”。	20

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-94](#)。

表 17-94 算子输入、输出参数设置

输入类型	输出类型
HBase输入	表输出

图 17-66 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，设置数据保存方式。

表 17-95 输出设置参数

参数名	说明	示例
架构名称	数据库模式名。	dbo
表名	数据库表名，用于最终保存传输的数据。 <b>说明</b> 表名可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	test
临时表	数据库临时表表名，用于临时保存传输过程中的数据，字段需要和“表名”配置的表一致。 <b>说明</b> 使用临时表是为了使得导出数据到数据库时，不会在目的表中产生脏数据。只有在所有数据成功写入临时表后，才会将数据从临时表迁移到目的表。使用临时表会增加作业的执行时间。	tm p_ est

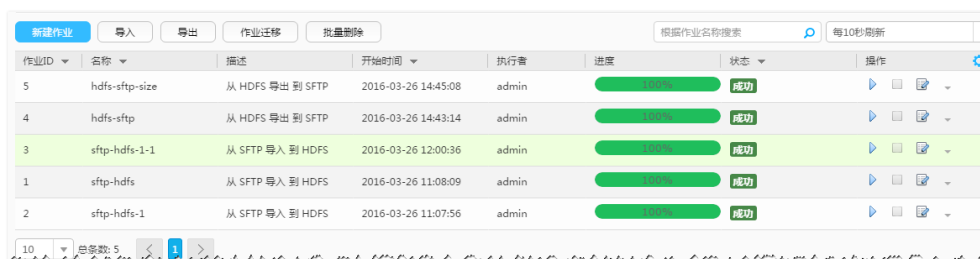
**步骤7** 单击“保存并运行”，开始保存并运行作业。

### 查看作业完成情况

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。



图 17-67 查看作业



作业ID	名称	描述	开始时间	执行者	进度	状态	操作
5	hdfs-sftp-size	从 HDFS 导出到 SFTP	2016-03-26 14:45:08	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">🔍</a> <a href="#">-</a>
4	hdfs-sftp	从 HDFS 导出到 SFTP	2016-03-26 14:43:14	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">🔍</a> <a href="#">-</a>
3	sftp-hdfs-1-1	从 SFTP 导入到 HDFS	2016-03-26 12:00:36	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">🔍</a> <a href="#">-</a>
1	sftp-hdfs	从 SFTP 导入到 HDFS	2016-03-26 11:08:09	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">🔍</a> <a href="#">-</a>
2	sftp-hdfs-1	从 SFTP 导入到 HDFS	2016-03-26 11:07:56	admin	<div style="width: 100%;"></div>	成功	<a href="#">▶</a> <a href="#">🔍</a> <a href="#">-</a>

----结束

## 17.13.8 典型场景：从 Hive 导出数据到关系型数据库

### 操作场景

该任务指导用户使用Loader将数据从Hive导出到关系型数据库。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的Hive表。
- 获取关系型数据库使用的用户和密码。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。
- 操作前需要进行如下配置：
  - a. 获取关系型数据库对应的驱动jar包保存在Loader服务主备节点的lib路径：  
“`#{BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/webapps/loader/WEB-INF/ext-lib`”。
  - b. 使用root用户在主备节点分别执行以下命令修改权限：  

```
cd #{BIGDATA_HOME}/FusionInsight_Porter_8.1.0.1/install/
FusionInsight-Sqoop-1.99.3/FusionInsight-Sqoop-1.99.3/server/
webapps/loader/WEB-INF/ext-lib
chown omm:wheel jar包文件名
chmod 600 jar包文件名
```
  - c. 登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Loader > 更多 > 重启服务”，输入管理员密码重启Loader服务。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。

2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-68 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-69 基本信息界面

1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导出”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“generic-jdbc-connector”或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

#### 📖 说明

- 与关系数据库连接时，可以选择通用数据库连接器（generic-jdbc-connector）或者专用数据库连接器（oracle-connector、oracle-partition-connector、mysql-fastpath-connector），专用数据库连接器特别针对具体数据库类型进行优化，相对通用数据库连接器来说，导出、导入速度更快。
- 使用mysql-fastpath-connector时，要求在NodeManager节点上有MySQL的mysqldump和mysqlimport命令，并且此两个命令所属MySQL客户端版本与MySQL服务器版本兼容，如果没有这两个命令或版本不兼容，请参考<http://dev.mysql.com/doc/refman/5.7/en/linux-installation-rpm.html>，安装MySQL client applications and tools。

表 17-96 “generic-jdbc-connector” 连接参数

参数名	说明	示例
名称	关系型数据库连接的名称。	dbName
JDBC驱动程序类	JDBC驱动类名。	oracle.jdbc.driver.OracleDriver
JDBC连接字符串	JDBC连接字符串。	jdbc:oracle:thin:@//10.16.0.1:1521/oradb
用户名	连接数据库使用的用户名。	omm
密码	连接数据库使用的密码。	xxxx
JDBC连接属性	JDBC连接属性，单击“添加”手动添加。 <ul style="list-style-type: none"> <li>名称：连接属性名</li> <li>值：连接属性值</li> </ul>	<ul style="list-style-type: none"> <li>名称：socketTimeout</li> <li>值：20</li> </ul>

### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，在“源文件类型”中选择“HIVE”，设置数据源信息。

表 17-97 数据源配置参数说明

参数名	解释说明	示例
Hive实例	在Hive作业中，Loader支持从集群可添加的所有Hive服务实例中选择一个。如果选定的Hive服务实例在集群中未添加，则此作业无法正常运行。	hive
个数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于“3000”。	20

### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-98](#)。

表 17-98 算子输入、输出参数设置

输入类型	输出类型
Hive输入	表输出

图 17-70 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，设置数据保存方式。

表 17-99 输出设置参数

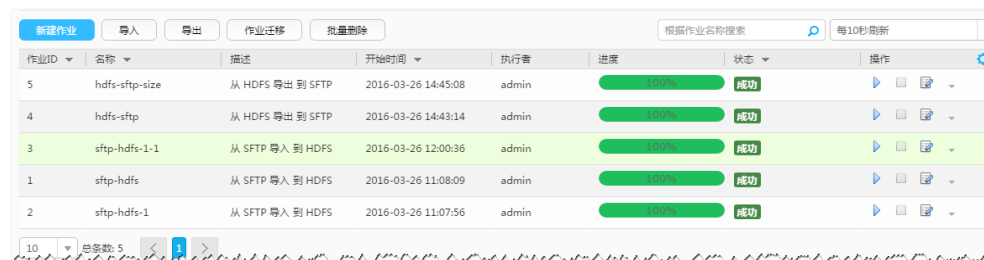
参数名	说明	示例
架构名称	数据库模式名。	dbo
表名	数据库表名，用于最终保存传输的数据。 <b>说明</b> 表名可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	test
临时表	数据库临时表表名，用于临时保存传输过程中的数据，字段需要和“表名”配置的表一致。 <b>说明</b> 使用临时表是为了使得导出数据到数据库时，不会在目的表中产生脏数据。只有在所有数据成功写入临时表后，才会将数据从临时表迁移到目的表。使用临时表会增加作业的执行时间。	tm p_ est

**步骤7** 单击“保存并运行”，开始保存并运行作业。

### 查看作业完成情况

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-71 查看作业



作业ID	名称	描述	开始时间	执行者	进度	状态	操作
5	hdfs-sftp-size	从 HDFS 导出到 SFTP	2016-03-26 14:45:08	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️
4	hdfs-sftp	从 HDFS 导出到 SFTP	2016-03-26 14:43:14	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️
3	sftp-hdfs-1-1	从 SFTP 导入到 HDFS	2016-03-26 12:00:36	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️
1	sftp-hdfs	从 SFTP 导入到 HDFS	2016-03-26 11:08:09	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️
2	sftp-hdfs-1	从 SFTP 导入到 HDFS	2016-03-26 11:07:56	admin	<div style="width: 100%;"></div>	成功	▶ □ 🗑️

----结束

## 17.13.9 典型场景：从 HBase 导出数据到 HDFS/OBS

### 操作场景

该任务指导用户使用Loader将数据从HBase导出到HDFS/OBS。

### 前提条件

- 创建或获取该任务中创建Loader作业的业务用户和密码。
- 确保用户已授权访问作业执行时操作的HDFS/OBS目录和数据。
- 确保用户已授权访问作业执行时操作的HBase表或phoenix表。
- 检查磁盘空间，确保没有出现告警且余量满足导入、导出数据的大小。
- 如果设置的作业需要使用指定YARN队列功能，该用户需要已授权有相关YARN队列的权限。
- 设置任务的用户需要获取该任务的执行权限，并获取该任务对应的连接的使用权限。

### 操作步骤

#### 设置作业基本信息

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-72 Loader WebUI 界面



**步骤2** 单击“新建作业”，进入“基本信息”界面，创建作业基本信息。

图 17-73 基本信息界面

1. 在“名称”中输入作业的名称。
2. 在“类型”中选择“导出”。
3. 在“组”中设置作业所属组，默认没有已创建的组，单击“添加”创建一个新的组，输入组的名称，单击“确定”保存。
4. 在“队列”中选择执行该作业的YARN队列。默认值“root.default”。
5. 在“优先级”中选择执行该作业的YARN队列优先级。默认值为“NORMAL”。可选值为“VERY\_LOW”、“LOW”、“NORMAL”、“HIGH”和“VERY\_HIGH”。

**步骤3** 在“连接”区域，单击“添加”新建一个的连接，在“连接器”中选择“hdfs-connector”，输入配置连接参数，单击“测试”验证连接是否可用，待提示“测试成功”后单击“确定”。

#### 设置数据源信息

**步骤4** 单击“下一步”，进入“输入设置”界面，在“源文件类型”中选择“HBASE”，设置数据源信息。

表 17-100 输入设置参数

参数名	解释说明	示例
HBase实例	在HBase作业中，Loader支持从集群可添加的所有HBase服务实例中选择任意一个。如果选定的HBase服务实例在集群中未添加，则此作业无法正常运行。	HB ase
个数	配置数据操作的MapReduce任务中同时启动的map数量。参数值必须小于或等于3000。	20

#### 设置数据转换

**步骤5** 单击“下一步”，进入“转换”界面，设置数据传输过程中的转换操作。算子的选择和参数设置具体请参考[算子帮助](#)及[表17-101](#)。

表 17-101 算子输入、输出参数设置

输入类型	输出类型
HBase输入	文件输出

图 17-74 算子操作方法示意



### 设置数据保存信息并运行作业

**步骤6** 单击“下一步”，进入“输出设置”界面，设置数据保存方式。

表 17-102 输出设置参数

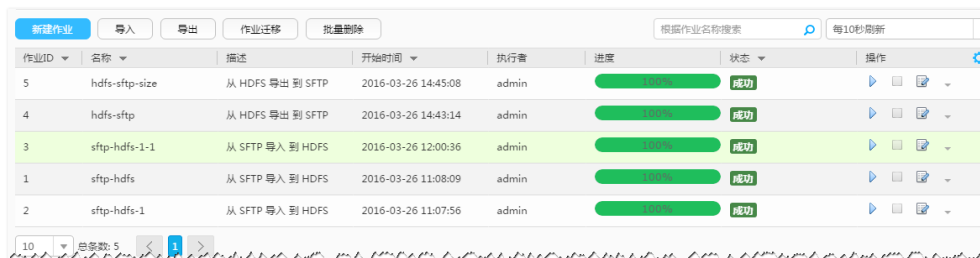
参数名	解释说明	示例
输出路径	导出文件在HDFS/OBS的输出目录或者文件名。 <b>说明</b> 路径参数可以使用宏定义，具体请参考 <a href="#">配置项中使用宏定义</a> 。	/user/test
文件格式	文件导出类型： <ul style="list-style-type: none"> <li>“TEXT_FILE”：导入文本文件并保存为文本文件。</li> <li>“SEQUENCE_FILE”：导入文本文件并保存在“sequence file”文件格式。</li> <li>“BINARY_FILE”：以二进制流的方式导入文件，可以导入任何格式的文件。</li> </ul>	TEXT_FILE
压缩格式	在下拉菜单中选择数据导出到HDFS/OBS后保存文件的压缩格式，未配置或选择“NONE”表示不压缩数据。	NONE

**步骤7** 单击“保存并运行”，开始保存并运行作业。

### 查看作业完成情况

**步骤8** 进入“Loader WebUI”界面，待“状态”显示“成功”则说明作业完成。

图 17-75 查看作业



作业ID	名称	描述	开始时间	执行者	进度	状态	操作
5	hdfs-sftp-size	从 HDFS 导出 到 SFTP	2016-03-26 14:45:08	admin	100%	成功	▶ □ 🗑️
4	hdfs-sftp	从 HDFS 导出 到 SFTP	2016-03-26 14:43:14	admin	100%	成功	▶ □ 🗑️
3	sftp-hdfs-1-1	从 SFTP 导入 到 HDFS	2016-03-26 12:00:36	admin	100%	成功	▶ □ 🗑️
1	sftp-hdfs	从 SFTP 导入 到 HDFS	2016-03-26 11:08:09	admin	100%	成功	▶ □ 🗑️
2	sftp-hdfs-1	从 SFTP 导入 到 HDFS	2016-03-26 11:07:56	admin	100%	成功	▶ □ 🗑️

----结束

## 17.14 作业管理

### 17.14.1 批量迁移 Loader 作业

#### 操作场景

Loader支持将作业批量从一个分组（源分组）迁移到另一个分组（目标分组）。  
本章节适用于MRS 3.x及后续版本。

#### 前提条件

- 源分组和目标分组均存在。
- 当前用户具备源分组和目标分组的编辑“Group Edit”权限。
- 当前用户具备源分组的作业编辑“Jobs Edit”权限或待迁移作业的作业编辑“Edit”权限。

#### 操作步骤

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-76 Loader WebUI 界面



**步骤2** 单击“作业迁移”，进入作业迁移界面。

**步骤3** 在“源分组”中选择待迁移作业当前所属分组，在“目标分组”中选择待迁移作业的目标分组。



**步骤4** 在“选择迁移类型”中选择迁移类型。

- “所有”：将源分组所有作业迁移到目标分组。
- “指定作业”：将源分组中指定的作业迁移到目标分组。选择“指定作业”，在作业列表中勾选需要迁移的作业。

**步骤5** 单击“确定”，开始作业迁移。当弹出框中进度条显示100%，则说明作业迁移完成。

----结束

## 17.14.2 批量删除 Loader 作业

### 操作场景

Loader支持批量删除已有作业。

本章节适用于MRS 3.x及后续版本。

### 前提条件

当前用户具备待删除作业的编辑“Edit”权限或作业所在分组的编辑“Jobs Edit”权限。

### 操作步骤

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-77 Loader WebUI 界面



**步骤2** 单击“批量删除”，进入作业批量删除界面。

**步骤3** 在“批量删除”中选择删除作业类型。

- “所有”，表示删除当前所有的作业。
- “指定作业”，表示指定需要删除的作业。选择“指定作业”，在作业列表中勾选需要删除的作业。

**步骤4** 单击“确定”，开始删除作业。当弹出框中进度条显示100%，则说明作业删除完成。

----结束

## 17.14.3 批量导入 Loader 作业

### 操作场景

Loader支持批量导入某个配置文件中的所有作业。

本章节适用于MRS 3.x及后续版本。

### 前提条件

当前用户具备待导入作业所在分组的编辑“Jobs Edit”权限。

#### 说明

如果作业所在的分组不存在，则会先创建该分组。当前用户就是该分组的创建者，拥有该分组的编辑“Jobs Edit”权限。

### 操作步骤

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-78 Loader WebUI 界面



**步骤2** 单击“导入”，进入作业导出界面。

**步骤3** 在“导入”界面中选择要导入的配置文件的路径。

**步骤4** 单击“上传”，开始导入作业。当弹出框中进度条显示100%，则说明作业导出完成。

----结束

## 17.14.4 批量导出 Loader 作业

### 操作场景

Loader支持批量导出已有作业。

本章节适用于MRS 3.x及后续版本。

### 前提条件

当前用户具备待导出作业的编辑“Edit”权限或作业所在分组的编辑“Jobs Edit”权限。

## 操作步骤

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-79 Loader WebUI 界面



**步骤2** 单击“导出”，进入作业导出界面。

**步骤3** 在“选择导出类型”中选择删除作业类型。

- “所有”：表示导出当前所有的作业。
- “指定作业”：表示指定需要导出的作业。选择“指定作业”，在作业列表中勾选需要导出的作业。
- “指定组别”：表示导出某个指定分组中的所有作业。选择“指定分组”，在分组列表中勾选需要导出的作业分组。
- “是否导出密码”：导出时是否导出连接器密码，勾选时，导出加密后的密码串。

**步骤4** 单击“确定”，开始导出作业。当弹出框中进度条显示100%，则说明作业导出完成。

----结束

## 17.14.5 查看作业历史信息

### 操作场景

该任务指导您在日常运维中，查看某个Loader作业的历史执行状态以及每次执行时长，同时提供该作业两种操作：

- 脏数据：查看作业执行过程中处理失败的数据、或者被清洗过滤掉的数据，针对该数据可以查看源数据中哪些数据不符合转换、清洗规则。
- 日志：查看作业在MapReduce执行的日志信息。

本章节适用于MRS 3.x及后续版本。

### 前提条件

获取登录“Loader WebUI”的账户和密码。

### 操作步骤

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-80 Loader WebUI 界面



**步骤2** 查看Loader作业的历史记录。

1. 选择待查看的作业所在行。
2. 如图所示，选择“更多>历史记录”查看作业执行的历史记录。

图 17-81 查看历史记录

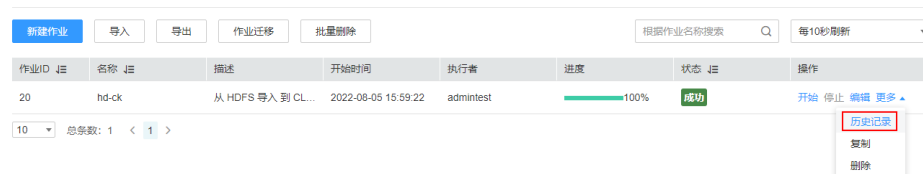


表 17-103 参数说明

名称	说明
行/文件 读取数	从输入源中读取的行数（文件数）。
行/文件 写入数	写入到输出源的行数（文件数）。
行/文件 跳过数	<ul style="list-style-type: none"> <li>- 转换过程中记录的坏行数（文件数）：输入格式不正确，无法进行转换。</li> <li>- 转换过程中配置过滤条件后跳过的行数。</li> </ul>

----结束

## 17.15 算子帮助

### 17.15.1 概述

“算子帮助”章节适用于MRS 3.x及后续版本。

## 转换流程

Loader读取源端数据，通过输入算子将数据按规则逐一转换成字段，再通过转换算子，对这些字段做清洗或转换，最后通过输出算子将处理后的字段，输出到目标端。

- 每个作业，如果进行数据转换操作，有且只能有一个输入算子，有且只能有一个输出算子。
- 不符合转换规则的数据，将成为脏数据跳过。

### 说明

- 从关系型数据库导入数据到HDFS/OBS，可以不用配置数据转换，数据将按“,”分隔保存到HDFS/OBS。
- 从HDFS/OBS导出数据到关系型数据库，可以不用配置数据转换，数据将按“,”分隔保存到关系型数据库。

## 算子简介

Loader算子包括以下类型：

- **输入算子**  
数据转换的第一步，负责将数据转换成字段，每次转换有且只能有一种输入算子，涉及HBase或Hive导入导出时，必须填写。
- **转换算子**  
数据转换的中间转换步骤，属于可选类型，各个转换算子可任意搭配使用。转换算子是针对字段而言，必须先使用输入算子，将数据转换成字段。
- **输出算子**  
数据转换的最后一步，每次转换有且只能有一种输出算子，用于输出处理后的字段。涉及HBase或Hive导入导出时，必须填写。

表 17-104 算子分类一览表

类型	描述
输入	<ul style="list-style-type: none"> <li>• CSV文件输入：将文件的每一行按指定分隔符转换成多个输入字段。</li> <li>• 固定宽度文件输入：将文件的每一行，按可配置长度的字符或字节，转换成多个输入字段。</li> <li>• 表输入：将关系型数据库表的指定列按顺序转换成同等数量的输入字段。</li> <li>• HBase输入：将HBase表的指定列转换成同等数量的输入字段。</li> <li>• HTML输入：将HTML文件中的元素转换成输入字段。</li> <li>• Hive输入：将Hive表的指定列转换成同等数量的输入字段。</li> </ul>

类型	描述
转换	<ul style="list-style-type: none"> <li>长整型时间转换：实现长整型数值与日期类型的互换。</li> <li>空值转换：将空值替换成指定值。</li> <li>增加常量字段：生成常量字段。</li> <li>随机值转换：生成随机数字段。</li> <li>拼接转换：拼接已有字段，生成新字段。</li> <li>分隔转换：将已有字段，按指定分隔符，分隔出新字段。</li> <li>取模转换：对已有字段取模，生成新字段。</li> <li>剪切字符串：通过指定起始位置，截取已有字符串类型的字段，生成新字段。</li> <li>EL操作转换：指定算法，对字段值进行运算，目前支持的算法有：md5sum、sha1sum、sha256sum和sha512sum等。</li> <li>字符串大小写转换：对已有的字符串类型字段，切换大小写，生成新字段。</li> <li>字符串逆序转换：对已有的字符串类型字段，做逆序变换，生成新字段。</li> <li>字符串空格清除转换：对已有的字符串类型字段，清除左右空格，生成新字段。</li> <li>过滤行转换：配置逻辑条件过滤掉含触发条件的行。</li> <li>更新域：当满足某些条件时，更新字段的值。</li> </ul>
输出	<ul style="list-style-type: none"> <li>Hive输出：将已生成的字段输出到Hive表。</li> <li>表输出：将已生成的字段输出到关系型数据库表。</li> <li>文件输出：将已生成的字段通过分隔符连接并输出到文件。</li> <li>HBase输出：将已生成的字段输出到HBase表。</li> </ul>

## 字段简介

作业配置中的字段是Loader按业务需要定义的与用户数据对应的一种数据项，它拥有具体类型，必须与用户实际数据类型保持一致。

## 17.15.2 输入算子

### 17.15.2.1 CSV 文件输入

#### 概述

“CSV文件输入”算子，用于导入所有能用文本编辑器打开的文件。

## 输入与输出

- 输入：文本文件。
- 输出：多个字段。

## 参数说明

表 17-105 算子参数说明

参数	含义	类型	是否必填	默认值
分隔符	CSV文件的列分隔符，用于分隔每行的数据。	string	是	,
换行符	用户根据数据实际情况，填写字符串作为换行符。支持任何字符串。默认使用操作系统的换行符。	string	否	\n
文件名是否作为字段	自定义一个字段，以当前数据所在的文件名作为该字段值。	string	否	无
绝对路径	配置“文件名是否作为字段”引用文件名环境，选中单选框时是带绝对路径的文件名；不选中单选框时是不带路径的文件名。	boolean	否	不选中
验证输入字段	是否检验输入字段与值的类型匹配情况，值为“NO”，不检查；值为“YES”，检查。若不匹配则跳过该行。	enum	是	YES
输入字段	配置输入字段的相关信息： <ul style="list-style-type: none"><li>• 位置：源文件每行被列分隔符分隔后，目标字段对应的位置，从1开始编号。</li><li>• 字段名：配置字段名。</li><li>• 类型：配置字段类型。</li><li>• 数据格式：字段类型为“DATE”或“TIM”E或“TIMESTAMP”时，需指定特定时间格式，其他字段类型指定无效。时间格式如：“yyyyMMdd HH:mm:ss”。</li><li>• 长度：配置字段长度，字段值太长则按配置的长度截取，类型为“CHAR”时实际长度不足则空格补齐，类型为“VARCHAR”时实际长度不足则不补齐。</li></ul>	map	是	无

## 数据处理规则

- 将每行数据按照指定的分隔符，分隔成多个字段，供之后的转换算子使用。
- 当字段的值与实际的类型不匹配时，该行数据会成为脏数据。

- 输入字段列数不等于原始数据实际包含字段列数，该行数据会保存为脏数据。

## 样例

源文件如下图：

```
2016,year
year,2016
```

配置“CSV文件输入”算子，分隔符为“,”，生成两个字段A、B。

分隔符: ,

换行符:

文件名是否作为字段:

绝对路径:

验证输入字段: YES

输入字段

导入 导出

表格编辑 文本编辑

位置	字段名	类型	数据格式	长度	
1	A	VARCHAR			↑ ↓ ↕ ×
2	B	VARCHAR			↑ ↓ ↕ ×

添加

将A、B输出，结果如下：

```
2016,year
year,2016
```

### 17.15.2.2 固定宽度文件输入

#### 概述

“固定宽度文件输入”算子，将文件的每一行，按可配置长度的字符或字节，转换成多个输入字段。

#### 输入与输出

- 输入：文本文件。
- 输出：多个字段。



## 参数说明

表 17-106 算子参数说明

参数	含义	类型	是否必填	默认值
换行符	用户根据数据实际情况，填写字符串作为换行符。支持任何字符串。默认使用操作系统的换行符。	string	否	\n
分割长度单位	长度单位，可选择“char”字符或“byte”字节。	enum	是	char
输入字段	配置输入字段相关信息： <ul style="list-style-type: none"> <li>固定长度：设置字段长度，第2个字段起点从第1个字段终点开始，以此类推。</li> <li>字段名：配置输入字段名。</li> <li>类型：配置字段类型。</li> <li>数据格式：字段类型为“DATE”或“TIME”或“TIMESTAMP”时，需指定特定时间格式，其他字段类型指定无效。时间格式如：“yyyyMMdd HH:mm:ss”。</li> <li>长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> </ul>	map	是	无

## 数据处理规则

- 按照输入字段的长度依次截取源文件，生成字段。
- 当字段的值与实际的类型不匹配时，该行数据会成为脏数据。
- 配置字段分割长度，大于原字段值的长度，则数据分割失败，当前行成为脏数据。

## 样例

源文件如下图：

```
fusionInsightbigdataprodu
```

配置“固定宽度文件输入”算子，生成三个字段A、B和C。

换行符

分割长度单位

输入字段

固定长度	字段名	类型	时间格式
<input type="text" value="13"/>	<input type="text" value="A"/>	<input type="text" value="VARCHAR"/>	<input type="text"/>
<input type="text" value="7"/>	<input type="text" value="B"/>	<input type="text" value="VARCHAR"/>	<input type="text"/>
<input type="text" value="7"/>	<input type="text" value="C"/>	<input type="text" value="VARCHAR"/>	<input type="text"/>

将三个字段依次输出，结果如下：

```
fusionInsight,bigdata,product
```

### 17.15.2.3 表输入

#### 概述

“表输入”算子，将关系型数据库表的指定列按顺序转换成同等数量的输入字段。

#### 输入与输出

- 输入：表列
- 输出：字段

#### 参数说明

表 17-107 算子参数说明

参数	含义	类型	是否必填	默认值
输入字段	配置关系型数据库输入字段的相关信息： <ul style="list-style-type: none"> <li>• 位置：配置输入字段的位置。</li> <li>• 字段名：配置输入字段名。</li> <li>• 类型：配置字段类型。</li> <li>• 长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> </ul>	map	是	无

## 数据处理规则

- 将指定的列按顺序生成字段。具体的表列是在作业配置的第二步“输入设置”中指定，当配置了“表列名”时，就是配置的值；当没配置“表列名”时，默认该表的所有列或者是“SQL语句”配置项里配置的查询条件中指定的列。
- 配置的输入字段个数不能大于实际指定的列数，否则全部数据成为脏数据。
- 当字段的值与实际的类型不匹配时，该行数据会成为脏数据。

## 样例

以sqlserver 2014为例，创建测试表test：

```
create table test (id int, name text, value text);
```

往测试表中插入三条数据：

```
insert into test values (1,'zhangshan','zhang');
```

```
insert into test values (2,'lisi','li');
```

```
insert into test values (3,'wangwu','wang');
```

查询表：

	id	name	value
1	1	zhangshan	zhang
2	2	lisi	li
3	3	wangwu	wang

配置“表输入”算子，生成三个字段：

设置了数据连接器后，可以单击“自动识别”，系统将自动读取数据库中的字段，可根据需要选择添加，然后根据业务场景手动进行完善或者修正即可，无需逐一手动添加。

### 说明

- 此操作会覆盖表格内已有数据。
- 单击“自动识别”后，建议手动检查系统自动识别出的字段类型，确保与表中实际的字段类型相符合。

例如Oracle数据库中的“date”类型，系统会自动识别为“timestamp”类型，若不手动处理会导致后续Hive表在查询数据时报错。

输入字段

导入 导出 自动识别

表格编辑 文本编辑

位置	字段名	类型	长度	
1	A	VARCHAR		↑ ↓ ↺ ×
2	B	VARCHAR		↑ ↓ ↺ ×
3	C	VARCHAR		↑ ↓ ↺ ×

添加

配置输出算子，输出到HDFS/OBS，结果如下：

```
1,zhangshan,zhang
2,lisi,li
3,wangwu,wang
```

## 17.15.2.4 HBase 输入

### 概述

“HBase输入”算子，将HBase表的指定列转换成同等数量的输入字段。

### 输入与输出

- 输入：HBase表列
- 输出：字段

### 参数说明

表 17-108 算子参数说明

参数	含义	类型	是否必填	默认值
HBase表类型	配置HBase表类型，可选项为normal（普通表）和phoenix表。	enum	是	normal
HBase表名	配置HBase表名。仅支持一个HBase表。	string	是	无
HBase输入字段	配置HBase输入信息： <ul style="list-style-type: none"> <li>• 列族：配置HBase列族名。</li> <li>• 列名：配置HBase列名。</li> <li>• 字段名：配置输入字段名。</li> <li>• 类型：配置字段类型。</li> <li>• 长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> <li>• 主键：配置是否为主键列。普通HBase表主键只能指定一个；phoenix表主键可以指定多个，配置多个列为主键时，会按照配置列的先后顺序对其进行拼接。必须配置一个主键列。</li> </ul>	map	是	无

### 数据处理规则

- 当配置HBase表名不存在时，作业提交失败。

- 当配置的列名与HBase表列名不匹配时，读取不到数据，导入数据条数会为0。
- 配置输入字段列数，大于原始数据实际包含字段列数，全部数据成为脏数据。
- 当字段的值与实际的类型不匹配时，该行数据会成为脏数据。

## 样例

以HBase导出到sqlserver2014数据库为例。

在sqlserver2014上创建一张空表test\_1用于存储HBase数据。执行以下语句：

```
create table test_1 (id int, name text, value text);
```

配置“HBase输入”算子，生成三个字段A、B和C：

设置了数据连接器后，可以单击“自动识别”，系统将自动读取数据库中的字段，可根据需要选择添加，然后根据业务场景手动进行完善或者修正即可，无需逐一手动添加。

### 说明

此操作会覆盖表格内已有数据。

列族名	列名	字段名	类型	长度	主键
f1	A	A	VARCHAR		<input checked="" type="checkbox"/>
f1	B	B	VARCHAR		<input type="checkbox"/>
f1	C	C	VARCHAR		<input type="checkbox"/>

通过“表输出”算子，将A、B和C输出到test\_1表中：

```
select * from test_1;
```

	id	name	value
1	1	zhangshan	zhang
2	2	lisi	li
3	3	wangwu	wang

## 17.15.2.5 HTML 输入

### 概述

“HTML输入”算子，导入有规则的HTML文件，并将HTML文件中的元素转换成输入字段。

## 输入与输出

输入：HTML文件

输出：多个字段

## 参数说明

表 17-109 算子参数说明

参数	含义	类型	是否必填	默认值
父标签	所有字段的上层HTML标签，用于限定搜索范围。	string	是	无
文件名	自定义一个字段，以当前数据所在的文件名作为该字段值。	string	否	无
绝对文件名	配置“文件名”引用文件名环境，选中单选框时是带绝对路径的文件名；不选中单选框时是不带路径的文件名。	boolean	否	否
验证输入字段	检验输入字段与值的类型匹配情况，值为“NO”，不检查；值为“YES”，检查。若不匹配则跳过该行。	enum	是	YES
输入字段	配置输入字段的相关信息： <ul style="list-style-type: none"> <li>位置：目标字段对应的位置，从1开始编号。</li> <li>字段名：配置字段名。</li> <li>字段所在的标签：字段的标签。</li> <li>关键字：配置关键字，能够匹配标签所在的内容，支持通配符，例如标签内容为“姓名”，可配置关键字“*姓名*”。</li> <li>类型：配置字段类型。</li> <li>数据格式：字段类型为“DATE”或“TIME”或“TIMESTAMP”时，需指定特定时间格式，其他字段类型指定无效。时间格式如：“yyyyMMdd HH:mm:ss”。</li> <li>长度：配置字段长度，字段值太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> </ul>	map	是	无

## 数据处理规则

- 首先配置父标签，限定搜索范围，父标签要存在，否则取到的内容为空。
- 配置输入字段，子标签用于精确定位字段所在的标签，相同的标签再通过关键字来精确匹配。
- 关键字用于匹配字段的内容，配置方法类似于“输入设置”中的“文件过滤器”字段，支持“\*”通配符，提供三种标记用于辅助定位，分别为：
  - a. “#PART”标记，表示取被通配符“\*”所匹配的值，如果存在多个“\*”号，可以指定一个序号，按从左到右的顺序，取得对应序号的“\*”所配置的内容。例如“#PART1”，表示取第1个“\*”号匹配的值；“#PART8”，表示取第8个“\*”号匹配的值。
  - b. “#NEXT”标记，表示取当前匹配的标签的下一个标签的值。
  - c. “#ALL”标记，表示取当前匹配的标签的所有内容作为值。
- 配置的标签有误时，取到的值为空，不会报错。

## 样例

源文件如下：

```
<html>
<body>
<table>
<tr>
<td>name:zhangshan</td>
<td>department:FusionInght</td>
<td>age:25</td>
</tr>
</table>
</body>
</html>
```

配置“HTML输入”算子，生成三个字段A、B和C：

父标签: tr

文件名:

绝对文件名:

验证输入字段: NO

输入字段

导入 导出

表格编辑 文本编辑

位置	字段名	字段所在的标签	关键字	类型	数据格式	长度
1	A	td	name:*#PART1	VARCHAR		
2	B	td	department:*#P,	VARCHAR		
3	C	td	age:*#PART1	VARCHAR		

添加

依次输出这三个字段，结果如下：

zhangshan,FusionInght,25

## 17.15.2.6 Hive 输入

### 概述

“Hive输入”算子，将Hive表的指定列转换成同等数量的输入字段。

### 输入与输出

- 输入：Hive表列
- 输出：字段

### 参数说明

表 17-110 算子参数说明

参数	含义	类型	是否必填	默认值
Hive数据库	Hive的数据库名称。	String	否	default
Hive表名	配置Hive表名。 仅支持一个Hive表。	String	是	无
分区过滤器	配置分区过滤器可以导出指定分区数据，默认为空，导出整个表数据。 例如导出分区字段locale的值为“CN”或“US”的表数据，输入如下： <b>locale = "CN" or locale = "US"</b>	String	否	-
Hive输入字段	配置Hive输入信息： <ul style="list-style-type: none"> <li>• 列名：配置Hive列名。</li> <li>• 字段名：配置输入字段名。</li> <li>• 类型：配置字段类型。</li> <li>• 长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> </ul>	map	是	-

### 数据处理规则

- 当配置Hive表名不存在时，作业提交失败。
- 当配置的列名与Hive表列名不匹配时，读取不到数据，导入数据条数会为0。



- 当字段的值与实际类型不匹配时，该行数据会成为脏数据。

## 样例

以Hive导出到sqlserver2014数据库为例。

在sqlserver2014上创建一张空表“test\_1”用于存储Hive数据。执行以下语句：

```
create table test_1 (id int, name text, value text);
```

配置“Hive输入”算子，生成三个字段A、B和C：

设置了数据连接器后，单击“自动识别”，系统将自动读取数据库中的字段，可根据需要选择添加，然后根据业务场景手动进行完善或者修正即可，无需逐一手动添加。

### 说明

此操作会覆盖表格内已有数据。

列名	字段名	类型	长度	
A	A	VARCHAR		↑ ↓ ↻ ✕
B	B	VARCHAR		↑ ↓ ↻ ✕
C	C	VARCHAR		↑ ↓ ↻ ✕

通过“表输出”算子，将A、B和C输出到“test\_1”表中：

```
select * from test_1;
```

	id	name	value
1	1	zhangshan	zhang
2	2	lisi	li
3	3	wangwu	wang

## 17.15.2.7 Spark 输入

### 概述

“Spark输入”算子，将SparkSQL表的指定列转换成同等数量的输入字段。

## 输入与输出

- 输入：SparkSQL表列
- 输出：字段

## 参数说明

表 17-111 算子参数说明

参数	含义	类型	是否必填	默认值
Spark数据库	SparkSQL的数据库名称。	String	否	default
Spark表名	配置SparkSQL表名。 仅支持一个SparkSQL表。	String	是	无
分区过滤器	配置分区过滤器可以导出指定分区数据，默认为空，导出整个表数据。 例如导出分区字段locale的值为“CN”或“US”的表数据，输入如下： <b>locale = "CN" or locale = "US"</b>	String	否	-
Spark输入字段	配置SparkSQL输入信息： <ul style="list-style-type: none"> <li>• 列名：配置SparkSQL列名。</li> <li>• 字段名：配置输入字段名。</li> <li>• 类型：配置字段类型。</li> <li>• 长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> </ul>	map	是	-

## 数据处理规则

- 当配置SparkSQL表名不存在时，作业提交失败。
- 当配置的列名与SparkSQL表列名不匹配时，读取不到数据，导入数据条数会为0。
- 当字段的值与实际的类型不匹配时，该行数据会成为脏数据。

## 样例

以SPARK导出到sqlserver2014数据库为例。

在sqlserver2014上创建一张空表“test\_1”用于存储SparkSQL数据。执行以下语句：

**create table test\_1 (id int, name text, value text);**

配置“Spark输入”算子，生成三个字段A、B和C：

设置了数据连接器后，单击“自动识别”，系统将自动读取数据库中的字段，可根据需要选择添加，然后根据业务场景手动进行完善或者修正即可，无需逐一手动添加。

### 说明

此操作会覆盖表格内已有数据。

Hive Input-Hive输入

Hive数据库

Hive表名

分区过滤器

Hive输入字段

列名	字段名	类型	长度				
A	A	VARCHAR	<input type="text"/>	↑	↓	↕	✖
B	B	VARCHAR	<input type="text"/>	↑	↓	↕	✖
C	C	VARCHAR	<input type="text"/>	↑	↓	↕	✖

通过“表输出”算子，将A、B和C输出到“test\_1”表中：

```
select * from test_1;
```

	id	name	value
1	1	zhangshan	zhang
2	2	lisi	li
3	3	wangwu	wang

## 17.15.3 转换算子

### 17.15.3.1 长整型时间转换

#### 概述

“长整型时间转换”算子，用于配置长整型数值与日期的转换。

#### 输入与输出

- 输入：需要转换的字段
- 输出：转换后的新字段

## 参数说明

表 17-112 算子参数说明

参数	含义	类型	是否必填	默认值
转换类型	配置长整型时间转换类型： <ul style="list-style-type: none"> <li>long to date：长整型数值转换为DATE类型。</li> <li>long to time：长整型数值转换为TIME类型。</li> <li>long to timestamp：长整型数值转换为TIMESTAMP类型。</li> <li>date to long：DATE类型转换为长整型数值。</li> <li>time to long：TIME类型转换为长整型数值。</li> <li>timestamp to long：TIMESTAMP类型转换为长整型数值。</li> </ul>	enum	是	long to date
输入字段名	配置输入的待转换字段名称，需填写上一个转换步骤生成的字段名。	string	是	无
输出字段名	配置输出字段的字段名。	string	是	无
字段单位	配置长整型数值字段的单位，根据“转换类型”长整型数据可以是输入字段或生成字段，可选值为“second”和“milisecond”。	enum	是	second
输出字段类型	配置输出字段的类型，可选值为“BIGINT”，“DATE”，“TIME”和“TIMESTAMP”。	enum	是	BIGINT
时间格式	配置时间字段格式，时间格式如：“yyyyMMdd HH:mm:ss”。	string	否	无

## 数据处理规则

- 原始数据包含null值，不做转换处理。
- 配置输入字段列数，大于原始数据实际包含字段列数，全部数据成为脏数据。
- 遇到类型转换错误，当前数据保存为脏数据。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下图：

```
1453431755874,2016-01-22 10:40:00
```

配置“长整型时间转换”算子，生成四个新字段C、D、E和F，类型分别为DATE、TIME、TIMESTAMP、BIGINT。

转换类型	输入字段名	输出字段名	字段单位	输出字段类型	时间格式
long to date ▼	A	C	millisecond ▼	DATE ▼	yyyy-MM-dd
long to time ▼	A	D	millisecond ▼	TIME ▼	HH:mm:ss
long to time ▼	A	E	millisecond ▼	TIMESTAMP ▼	yyyyMMdd HH:m
date to long ▼	B	F	millisecond ▼	BIGINT ▼	

添加

转换后，依次输出，结果如下：

```
1453431755874,2016-01-22,2016-01-22,11:02:35,20160122 11:02:35,1453430400000
```

### 17.15.3.2 空值转换

#### 概述

“空值转换”算子，用于将空值替换成指定值。

#### 输入与输出

- 输入：空值字段
- 输出：原字段，但值已经被替换

#### 参数说明

表 17-113 算子参数说明

参数	含义	类型	是否必填	默认值
输入字段名	配置可能出现空值的字段名，需填写已生成的字段名。	string	是	无
替换值	配置替换空值的指定值。	string	是	无

#### 数据处理规则

字段原值为null时，替换成指定的值。

#### 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下图：

```
,value1
key2,value2
key3,
```

配置“空值转换”算子，如下图：

输入字段名	替换值
<input type="text" value="A"/>	<input type="text" value="newKey"/>
<input type="text" value="B"/>	<input type="text" value="newValue"/>

转换后，将A和B的值输出后的结果如下：

```
newKey,value1
key2,value2
key3,newValue
```

### 17.15.3.3 增加常量字段

#### 概述

“增加常量字段”算子，用于直接生成常量字段。

#### 输入与输出

- 输入：无
- 输出：常量字段

## 参数说明

表 17-114 算子参数说明

参数	含义	类型	是否必填	默认值
配置字段	配置常量字段相关信息： <ul style="list-style-type: none"> <li>输出字段名：配置字段名。</li> <li>类型：配置字段类型。</li> <li>时间格式：字段类型为“DATE”或“TIME”或“TIMESTAMP”时，需指定特定时间格式，其他类型指定无效。时间格式如：“yyyyMMdd HH:mm:ss”。</li> <li>长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> <li>常量值：配置符合类型的常量值。</li> </ul>	map	是	无

## 数据处理规则

生成指定类型的常量字段。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下图：

```
,value1
key2,value2
key3,
```

配置“增加常量字段”算子，增加两个字段C和D：

输出字段名	类型	时间格式	长度	常量值
C	VARCHAR			constantsvalue1
D	INTEGER			2016

添加

转换后，将A、B、C和D按顺序输出，结果如下：

```
,value1,constantsvalue1,2016
key2,value2,constantsvalue1,2016
key3,,constantsvalue1,2016
```

### 17.15.3.4 随机值转换

#### 概述

“随机值转换”算子，用于配置新增值为随机数的字段。

#### 输入与输出

- 输入：无
- 输出：随机值字段

#### 参数说明

表 17-115 算子参数说明

参数	含义	类型	是否必填	默认值
输出字段名	配置生成随机值的字段名。	string	是	无
长度	配置字段长度。	map	是	无
类型	配置字段的类型，可选值为“VARCHAR”，“INTEGER”和“BIGINT”。	enum	是	VARCHAR

#### 数据处理规则

生成指定类型的随机值。

#### 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下图：

```
,value1
key2,value2
key3,
```

配置“随机值转换”算子，生成C、D、E三个字段：

输出字段名	类型
<input type="text" value="C"/>	<input type="text" value="VARCHAR"/>
<input type="text" value="D"/>	<input type="text" value="VARCHAR"/>
<input type="text" value="E"/>	<input type="text" value="VARCHAR"/>

转换后，按顺序输入这五个字段：



```
,value1,2druceak69ril,769974975,8452014577467885098
key2,value2,7oq2dku93q9cg,1631427868,867914116689501757
key3,,2jg5e7b1m17kq,654806209,2477823020516316030
```

可以发现，每次生成的随机值都不一样。

### 17.15.3.5 拼接转换

#### 概述

“拼接转换”算子，将已有字段的值通过连接符拼接，生成新的字段。

#### 输入与输出

- 输入：需要拼接的字段
- 输出：拼接后的字段

#### 参数说明

表 17-116 算子参数说明

参数	含义	类型	是否必填	默认值
输出字段名	配置拼接后的字段名。	string	是	无
分隔符	配置拼接符，可为空。	string	否	空字符串
被拼接字段名	配置需要被拼接字段名。 字段名：需填写上一个转换步骤生成的字段名，可添加多个。	map	是	无

#### 数据处理规则

- 按顺序将“被拼接字段名”中配置的字段的值，通过连接符拼接后，赋给“输出字段名”。
- 当有字段的值为null时，会转化为空字符串，再与其它字段值拼接。

#### 样例

通过“CSV文件输入”算子，生成三个字段A、B和C。

源文件如下图：

```
happy,new,year
welcome,to,2016
```

配置“拼接转换”算子，“分隔符”为空格，生成新字段D：

转换后，依次输出A、B、C和D，结果如下：

```
happy,new,year,happy new year
welcome,to,2016,welcome to 2016
```

### 17.15.3.6 分隔转换

#### 概述

“分隔转换”算子，将已有字段的值按指定的分隔符分隔后生成新字段。

#### 输入与输出

- 输入：需要分隔的字段
- 输出：分隔后的字段

#### 参数说明

表 17-117 算子参数说明

参数	含义	类型	是否必填	默认值
输入字段名	被分隔的字段名，需填写上一个转换步骤生成的字段名。	string	是	无
分隔符	配置分隔符。	string	是	无

参数	含义	类型	是否必填	默认值
分割后的字段	配置分隔后的字段，可为多个： <ul style="list-style-type: none"><li>位置：分隔后字段的位置。</li><li>输出字段名：分隔后的字段名。</li></ul>	map	是	无

## 数据处理规则

- 将输入字段的值按指定的分隔符分隔后，依次赋给配置的新字段。
- 配置分割后字段列数，大于原始数据实际可分割出来的字段列数，当前行成为脏数据。

## 样例

通过“CSV文件输入”算子，生成一个字段A。

源文件如下：

```
happy new year
welcome to 2016
```

配置“分隔转换”算子，“分隔符”为空格，生成三个字段B、C和D：

输入字段名

分隔符

分割后的字段

位置	输出字段名
<input type="text" value="1"/>	<input type="text" value="B"/>
<input type="text" value="2"/>	<input type="text" value="C"/>
<input type="text" value="3"/>	<input type="text" value="D"/>

转换后，依次输出A、B、C和D，结果如下：

```
happy new year,happy,new,year
welcome to 2016,welcome,to,2016
```

### 17.15.3.7 取模转换

#### 概述

“取模转换”算子，对整数字段取模，生成新字段。

#### 输入与输出

- 输入：整数字段
- 输出：模数字段

#### 参数说明

表 17-118 算子参数说明

参数	含义	类型	是否必填	默认值
取模字段名	配置取模运算信息： <ul style="list-style-type: none"><li>• 输入字段名：配置输入字段名，需填写上一个转换步骤生成的字段名。</li><li>• 输出字段名：配置输出字段名。</li><li>• 系数：指定取模的数值。</li></ul>	map	是	无

#### 数据处理规则

- 生成新字段，值为取模后的值。
- 字段的值须为整数，否则当前行会成为脏数据。

#### 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下图：

```
10,12
2015,2016
```

配置“取模转换”算子，生成两个新字段C和D：

输入字段名	输出字段名	系数
<input type="text" value="A"/>	<input type="text" value="C"/>	<input type="text" value="3"/>
<input type="text" value="B"/>	<input type="text" value="D"/>	<input type="text" value="3"/>
<input type="button" value="添加"/>		

转换后，依次输出A、B、C和D，结果如下：

```
10,12,1,0
2015,2016,2,0
```

### 17.15.3.8 剪切字符串

#### 概述

“剪切字符串”算子，截取已有字段的值，生成新的字段。

#### 输入与输出

- 输入：需要截取的字段
- 输出：截取后生成的新字段

#### 参数说明

表 17-119 算子参数说明

参数	含义	类型	是否必填	默认值
被截取的字段	配置被截取字段相关信息： <ul style="list-style-type: none"><li>• 输入字段名：配置输入字段名，需填写上一个转换步骤生成的字段名。</li><li>• 输出字段名：配置输出字段名。</li><li>• 开始位置：截取开始位置，从序号1开始。</li><li>• 结束位置：截取结束位置，不确定字符串长度时，可指定为-1表示被截取字段的末尾。</li><li>• 输出字段类型：输出字段的类型。</li><li>• 输出字段长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“输出字段类型”为“CHAR”时实际长度不足则空格补齐，“输出字段类型”为“VARCHAR”时实际长度不足则不补齐。</li></ul>	map	是	无

#### 数据处理规则

- 用开始位置和结束位置去截取原字段的值，生成新字段。
- 结束位置为“-1”时，表示字段的末尾。其它情况下，结束位置不能小于开始位置。
- 字符截取的开始位置或结束位置，大于输入字段的长度时，当前行成为脏数据。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下：

```
abcd,product
FusionInsight,Bigdata
```

配置“剪切字符串”算子后，生成两个新字段C和D：

输入字段名	输出字段名	开始位置	结束位置	输出字段类型	输出字段长度
A	C	1	3	VARCHAR	
B	D	1	4	VARCHAR	

添加

转换后，分别输出这三个字段：

```
abcd,product,abc,prod
FusionInsight,Bigdata,Fus,Bigd
```

### 17.15.3.9 EL 操作转换

#### 概述

“EL操作转换”算子，对字段值进行运算后生成新的字段，目前支持的算子有：md5sum、sha1sum、sha256sum和sha512sum等。

#### 输入与输出

- 输入：需要转换的字段
- 输出：经过EL表达式转换后的字段

## 参数说明

表 17-120 算子参数说明

参数	含义	类型	是否必填	默认值
el操作之后生成的字段	配置EL表达式： <ul style="list-style-type: none"> <li>名称：表达式输出结果的名称。</li> <li>el表达式：表达式，格式为：表达式名称（输入字段名,是否用小写字母表示输出结果）。例如，md5sum(fieldname,true)。                             <ul style="list-style-type: none"> <li>md5sum：生成md5校验值。</li> <li>sha1sum：生成sha1校验值。</li> <li>sha256sum：生成sha256校验值。</li> <li>sha512sum：生成sha512校验值。</li> </ul> </li> <li>类型：表达式输出结果类型，建议选择“VARCHAR”。</li> <li>时间格式：表达式输出结果格式。</li> <li>长度：表达式输出结果长度。</li> </ul>	map	是	无

## 数据处理规则

- 对字段值进行运算后生成新的字段。
- 当前新字段的类型只能为VARCHAR。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件见下图：

```
2016,year
year,2016
```

配置“EL操作转换”算子，生成C、D、E和F四个字段：

名称	el表达式	类型	时间格式
C	md5sum(A,false)	VARCHAR ▼	
D	sha1sum(A,true)	VARCHAR ▼	
E	sha256sum(B,false)	VARCHAR ▼	
F	sha512sum(B,true)	VARCHAR ▼	

依次输出这六个字段，结果如下图：

```
2016,year,95192C987323871658F8E396C0F2DAD2,ab39c54239118a4b086b878b7878100f769dd1
97,4CB4EA25583C25647247AE96FC90225D99AD7A6FABC3E2C2FD13C502E323CD9E,779edfe0463b2
596e7a83e4c59083e19242e8c51eace8e2ec57704643be5e15ba80f79af227cf3ea2e2362b4081377
96a1d82cb0535652b99844bb9a62019563
year,2016,84CDC76CABF418D7C961F6AB12F117D8,4ff0b1538469338a0073e2cdaab6a517801b6a
b4,DA6E2F539726FABD1F8CD7C9469A22B36769137975828ABC65FE2DC29E659B77,da0ae9104086a
1c58f89f82766ac55a02c8ab44277ce39f959ec0e73391bef651c6f9793657396ce47fbd846068465
ccbf3056764424bed9be7789bd1101ace7
```

### 17.15.3.10 字符串大小写转换

#### 概述

“字符串大小写转换”算子，用于配置已生成的字段通过大小写变换，转换出新的字段。

#### 输入与输出

- 输入：需要转换大小写的字段
- 输出：转换后的字段

#### 参数说明

表 17-121 算子参数说明

参数	含义	类型	是否必填	默认值
转换后的字段	配置字符串大小写转换的字段相关信息： <ul style="list-style-type: none"> <li>• 输入字段名：配置输入字段名，需填写上一个转换步骤生成的字段名。</li> <li>• 输出字段名：配置输出字段名。</li> <li>• 小写/大写：指定进行大写转换或小写转换。</li> </ul>	map	是	无

#### 数据处理规则

- 对字符串值做大小写转换。
- 传入数据为NULL值，不做转换处理。

#### 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下：

```
abcd,product
FusionInsight,Bigdata
```



配置“字符串大小写转换”算子后，生成两个新字段C和D：

输入字段名	输出字段名	小写/大写
<input type="text" value="A"/>	<input type="text" value="C"/>	Upper ▼
<input type="text" value="B"/>	<input type="text" value="D"/>	Lower ▼
<input type="button" value="添加"/>		

转换后，依次输出四个字段，结果如下：

```
abcd,product,ABCD,product
FusionInsight,Bigdata,FUSIONINSIGHT,bigdata
```

### 17.15.3.11 字符串逆序转换

#### 概述

“字符串逆序转换”算子，用于配置已生成的字段通过逆序，转换出新的字段。

#### 输入与输出

- 输入：需要逆序的字段
- 输出：逆序转换后的字段

#### 参数说明

表 17-122 算子参数说明

参数	含义	类型	是否必填	默认值
逆序转换的字段	<p>配置字符串逆序转换的字段相关信息：</p> <ul style="list-style-type: none"> <li>• 输入字段名：配置输入字段名，需填写上一个转换步骤生成的字段名。</li> <li>• 输出字段名：配置输出字段名。</li> <li>• 类型：配置字段类型。</li> <li>• 输出字段长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> </ul>	map	是	无

#### 数据处理规则

- 对字段的值做逆序操作。
- 传入数据为NULL值，不做转换处理。
- 配置输入字段列数，大于原始数据实际包含字段列数，全部数据成为脏数据。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下：

```
abcd,product
FusionInsight,Bigdata
```

配置“字符串逆序转换”算子后，生成两个新字段C和D：

输入字段名	输出字段名	类型	输出字段长度
A	C	VARCHAR ▼	
B	D	VARCHAR ▼	
添加			

转换后，依次输出四个字段，结果如下：

```
abcd,product,dcba,tcudorp
FusionInsight,Bigdata,thgislnoisuF,atadgiB
```

### 17.15.3.12 字符串空格清除转换

#### 概述

“字符串空格清除转换”算子，用于配置已生成的字段通过清除空格，转换出新的字段。

#### 输入与输出

- 输入：需要清除空格的字段
- 输出：转换后的字段

#### 参数说明

表 17-123 算子参数说明

参数	含义	类型	是否必填	默认值
清除空格的字段	配置字符串空格清除的字段相关信息： <ul style="list-style-type: none"> <li>• 输入字段名：配置输入字段名，需填写上一个转换步骤生成的字段名。</li> <li>• 输出字段名：配置输出字段名。</li> <li>• 对齐类型：配置清除方式（前空格、后空格、前后空格）。</li> </ul>	map	是	无

#### 数据处理规则

- 清空值两边的空格，支持只清除左边、只清除右边和同时清除左右空格。

- 传入数据为NULL值，不做转换处理。
- 配置输入字段列数，大于原始数据实际包含字段列数，全部数据成为脏数据。

## 样例

通过“CSV文件输入”算子，生成三个字段A、B和C。

源文件如下：

```
welcome ,to , 2016
happy ,new , year
```

配置“字符串空格清除转换”算子，生成三个新字段D、E和F。

输入字段名	输出字段名	对齐类型
<input type="text" value="A"/>	<input type="text" value="D"/>	both ▼
<input type="text" value="B"/>	<input type="text" value="E"/>	right ▼
<input type="text" value="C"/>	<input type="text" value="F"/>	left ▼

转换后，依次输出这六个字段，结果如下：

```
welcome ,to , 2016,welcome,to,2016
happy ,new , year,happy,new,year
```

### 17.15.3.13 过滤行转换

#### 概述

“过滤行转换”算子，用于配置逻辑条件过滤掉含触发条件的行。

#### 输入与输出

- 输入：用来做过滤条件的字段
- 输出：无

#### 参数说明

表 17-124 算子参数说明

参数	含义	类型	是否必填	默认值
条件逻辑连接符	配置条件逻辑连接符，可配置“AND”或“OR”。	enum	是	AND

参数	含义	类型	是否必填	默认值
条件	配置过滤条件相关信息： <ul style="list-style-type: none"> <li>输入字段名：配置输入字段名，需填写上一个转换步骤生成的字段名。</li> <li>操作：配置操作符。</li> <li>比较值：配置比较值，可直接输入值或输入“#{已存在的字段名}”格式引用字段的具体值。</li> </ul>	map	是	无

## 数据处理规则

- 条件逻辑为“AND”，如果未添加过滤条件，全部数据成为脏数据；或者原始数据满足添加的全部过滤条件，当前行成为脏数据。
- 条件逻辑为“OR”，如果未添加过滤条件，全部数据成为脏数据；或者原始数据满足任意添加的过滤条件，当前行成为脏数据。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下：

```
test, product
FusionInsight,Bigdata
```

配置“过滤行转换”算子，过滤掉含有test的行。

条件逻辑连接符 AND

条件

导入 导出

表格编辑 文本编辑

输入字段名	操作	比较值
A	==	test

添加

转换后，输入原字段，结果如下：

```
FusionInsight,Bigdata
```

### 17.15.3.14 更新域

#### 概述

“更新域”算子，当满足某些条件时，更新字段的值。

目前支持的类型有“BIGINT”、“DECIMAL”、“DOUBLE”、“FLOAT”、“INTEGER”、“SMALLINT”、“VARCHAR”。当类型为“VARCHAR”时，运算符为“+”时，表示在字符串后追加串，不支持“-”，当为其它类型时，“+”、“-”分别表示值的加和减。针对支持的所有类型，运算符“=”都表示直接赋新值。

#### 输入与输出

输入：字段

输出：输入字段

#### 参数说明

表 17-125 算子参数说明

参数	含义	类型	是否必填	默认值
更新字段名	需要更新的字段	string	是	无
操作符	操作符，支持“+”、“-”和“=”	enum	是	+
更新值	用来更新的值	与字段类型相匹配	否	无
条件逻辑连接符	配置条件逻辑连接符，可配置“AND”或“OR”。	enum	是	AND
条件	配置过滤条件相关信息： <ul style="list-style-type: none"> <li>输入字段名：配置输入字段名，需填写上一个转换步骤生成的字段名。</li> <li>操作：配置操作符。</li> <li>比较值：配置比较值，可直接输入值或输入“#{已存在的字段名}”格式引用字段的具体值。</li> </ul>	map	是	无

## 数据处理规则

- 首先判断条件是否成立。如果成立，更新字段的值；如果不成立，则不更新。
- 当更新字段为数值类型时，更新值需要为数值。
- 当更新字段为字符串类型时，更新操作不能为“-”。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下：

```
test, product
FusionInsight,Bigdata
```

配置“更新域”算子，当发现值为test时，更新值，在test后面加上good。

更新字段名

操作符

更新值

条件逻辑连接符

条件

输入字段名	操作	比较值
<input type="text" value="A"/>	<input style="border: none; border-bottom: 1px solid #ccc;" type="text" value="=="/>	<input type="text" value="test"/>

转换后，输出A和B，结果如下：

```
testgood ,product
FusionInsight,Bigdata
```

## 17.15.4 输出算子

### 17.15.4.1 Hive 输出

#### 概述

“Hive输出”算子，用于配置已生成的字段输出到Hive表的列。

## 输入与输出

- 输入：需要输出的字段
- 输出：Hive表

## 参数说明

表 17-126 算子参数说明

参数	含义	类型	是否必填	默认值
Hive文件存储格式	配置Hive表文件的存储格式（目前支持四种格式：CSV、ORC、RC和PARQUET）。 <b>说明</b> <ul style="list-style-type: none"> <li>• PARQUET格式是一种列式存储格式，PARQUET要求Loader的输出字段名和Hive表中的字段名保持一致。</li> <li>• Hive 1.2.0版本之后，Hive使用字段名称替代字段序号对ORC文件进行解析，因此，Loader的输出字段名和Hive表中的字段名需要保持一致。</li> </ul>	enum	是	CSV
Hive文件压缩格式	在下拉菜单中选择Hive表文件的压缩格式，未配置或选择“NONE”表示不压缩数据。	enum	是	NONE
Hive ORC文件版本	通过该字段配置ORC文件的版本（当Hive表文件的存储格式是ORC时）。	enum	是	0.12
输出分隔符	配置分隔符。	string	是	无
输出字段	配置输出信息： <ul style="list-style-type: none"> <li>• 位置：配置输出字段的位置。</li> <li>• 字段名：配置输出字段的字段名。</li> <li>• 类型：配置字段类型，字段类型为“DATE”或“TIME”或“TIMESTAMP”时，需指定特定时间格式，其他类型指定无效。时间格式如：“yyyyMMdd HH:mm:ss”。</li> <li>• 十进制格式：配置小数的刻度和精度。</li> <li>• 长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> <li>• 分区键：配置是否为分区列。可以不指定分区列，也可以指定多个分区列。配置多个列为分区列时，会按照配置列的先后顺序对其进行拼接。</li> </ul>	map	是	无

## 数据处理规则

- 将字段值输出到Hive表中。
- 如果指定了一个或多个列为分区列，则在作业配置第四步“输出设置”页面上，会显示“分割程序”属性，该属性表示使用多少个处理器去对分区数据进行处理。
- 如果没有指定任何列为分区列，则表示不需要对输入数据进行分区处理，“分割程序”属性默认隐藏。

## 样例

通过“CSV文件输入”算子，生成两个字段a\_str和b\_str。

源文件如下：

```
2016,year
year,2016
```

配置“Hive输出”算子，将a\_str和b\_str输出到Hive的表中。

位置	字段名	类型	十进制格式	长度	分区键
1	a_str	STRING			<input type="checkbox"/>
2	b_str	STRING			<input type="checkbox"/>

执行成功后，查看表数据：

```
0: jdbc:hive2://10.52.0.97:21066/> select * from hive_test;
+-----+-----+-----+
| hive_test.a_str | hive_test.b_str |
+-----+-----+-----+
| 2016 | year |
| year | 2016 |
+-----+-----+-----+
2 rows selected (1.6 seconds)
```

### 17.15.4.2 Spark 输出

#### 概述

“Spark输出”算子，用于配置已生成的字段输出到SparkSQL表的列。



## 输入与输出

- 输入：需要输出的字段
- 输出：SparkSQL表

## 参数说明

表 17-127 算子参数说明

参数	含义	类型	是否必填	默认值
Spark文件存储格式	配置SparkSQL表文件的存储格式（目前支持四种格式：CSV、ORC、RC和PARQUET）。 <b>说明</b> <ul style="list-style-type: none"><li>• PARQUET格式是一种列式存储格式，PARQUET要求Loader的输出字段名和SparkSQL表中的字段名保持一致。</li><li>• Hive 1.2.0版本之后，Hive使用字段名称替代字段序号对ORC文件进行解析，因此，Loader的输出字段名和SparkSQL表中的字段名需要保持一致。</li></ul>	enum	是	CSV
Spark文件压缩格式	在下拉菜单中选择SparkSQL表文件的压缩格式，未配置或选择“NONE”表示不压缩数据。	enum	是	NONE
Spark ORC文件版本	通过该字段配置ORC文件的版本（当SparkSQL表文件的存储格式是ORC时）。	enum	是	0.12
输出分隔符	配置分隔符。	string	是	无
输出字段	配置输出信息： <ul style="list-style-type: none"><li>• 位置：配置输出字段的位置。</li><li>• 字段名：配置输出字段的字段名。</li><li>• 类型：配置字段类型，字段类型为“DATE”或“TIME”或“TIMESTAMP”时，需指定特定时间格式，其他类型指定无效。时间格式如：“yyyyMMdd HH:mm:ss”。</li><li>• 十进制格式：配置小数的刻度和精度。</li><li>• 长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li><li>• 分区键：配置是否为分区列。可以不指定分区列，也可以指定多个分区列。配置多个列为分区列时，会按照配置列的先后顺序对其进行拼接。</li></ul>	map	是	无

## 数据处理规则

- 将字段值输出到SparkSQL表中。
- 如果指定了一个或多个列为分区列，则在作业配置第四步“输出设置”页面上，会显示“分割程序”属性，该属性表示使用多少个处理器去对分区数据进行处理。
- 如果没有指定任何列为分区列，则表示不需要对输入数据进行分区处理，“分割程序”属性默认隐藏。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下：

```
2016,year
year,2016
```

配置“Spark输出”算子，将A和B输出到SparkSQL的表中。

位置	字段名	类型	十进制格式	长度	分区键
1	A	STRING			<input type="checkbox"/>
2	B	STRING			<input type="checkbox"/>

### 17.15.4.3 表输出

#### 概述

“表输出”算子，用于配置输出的字段对应到关系型数据库的指定列。

#### 输入与输出

- 输入：需要输出的字段
- 输出：关系型数据库表

## 参数说明

表 17-128 算子参数说明

参数	含义	类型	是否必填	默认值
输出分隔符	配置分隔符。 <b>说明</b> 该配置仅用于MySQL专用连接器，当数据列内容中包含默认分隔符时，需要设置自定义分隔符，否则会出现数据错乱。	string	否	,
换行分隔符	用户根据数据实际情况，填写字符串作为换行符。支持任何字符串。默认使用操作系统的换行符。 <b>说明</b> 该配置仅用于MySQL专用连接器，当数据列内容中包含默认分隔符时，需要设置自定义分隔符，否则会出现数据错乱。	string	否	\n
输出字段	配置关系型数据库输出字段的相关信息： <ul style="list-style-type: none"> <li>• 字段名：配置输出字段的字段名。</li> <li>• 表列名：配置数据库表的列名。</li> <li>• 类型：配置字段类型，需要和数据库的字段类型一致。</li> <li>• 长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> </ul>	map	是	无

## 数据处理规则

将字段值输出到表中。

## 样例

以HBase导出到sqlserver2014数据库为例。

在sqlserver2014上创建一张空表test\_1用于存储HBase数据。执行以下语句：

```
create table test_1 (id int, name text, value text);
```

通过HBase输入步骤，生成三个字段A、B和C。

配置“表输出”算子，将A、B和C输出到test\_1表中：

输出结果如下：

	id	name	value
1	1	zhangshan	zhang
2	2	lisi	li
3	3	wangwu	wang

#### 17.15.4.4 文件输出

##### 概述

“文件输出”算子，用于配置已生成的字段通过分隔符连接并输出到文件。

##### 输入与输出

- 输入：需要输出的字段
- 输出：文件

##### 参数说明

表 17-129 算子参数说明

参数	含义	类型	是否必填	默认值
输出分隔符	配置分隔符。	string	是	无

参数	含义	类型	是否必填	默认值
换行符	用户根据数据实际情况，填写字符串作为换行符。支持任何字符串。默认使用操作系统的换行符。	string	否	\n
输出字段	配置输出信息： <ul style="list-style-type: none"> <li>位置：配置输出字段的位置。</li> <li>字段名：配置输出字段的字段名。</li> <li>类型：配置字段类型，字段类型为“DATE”或“TIME”或“TimeStamp”时，需指定特定时间格式，其他类型指定无效。时间格式如：“yyyyMMdd HH:mm:ss”。</li> <li>长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> </ul>	map	否	无

## 数据处理规则

将字段值输出到文件。

## 样例

通过“CSV文件输入”算子，生成两个字段A和B。

源文件如下：

```
aaa,product
bbb,Bigdata
```

配置“文件输出”算子，分隔符为“，”，将A和B输出到文件中：

输出分隔符

换行符

输出字段

位置	字段名	类型
<input type="text" value="1"/>	<input type="text" value="A"/>	<input type="text" value="VARCHAR"/> ▼
<input type="text" value="2"/>	<input type="text" value="B"/>	<input type="text" value="VARCHAR"/> ▼

输出后的结果如下：

```
aaa,product
bbb,Bigdata
```

### 17.15.4.5 HBase 输出

#### 概述

“HBase输出”算子，用于配置已生成的字段输出到HBase表的列。

#### 输入与输出

- 输入：需要输出的字段
- 输出：HBase表

#### 参数说明

表 17-130 算子参数说明

参数	含义	类型	是否必填	默认值
HBase表类型	配置HBase表类型，可选项为normal（普通HBase表）和phoenix表。	enum	是	normal
NULL值处理方式	配置NULL值处理方式。选中单选框时是将转换为空字符串并保存。不选中单选框时是不保存数据。	boolean	否	不选中单选框

参数	含义	类型	是否必填	默认值
HBase输出字段	<p>配置HBase输出信息：</p> <ul style="list-style-type: none"> <li>• 字段名：配置输出字段的字段名。</li> <li>• 表名：配置HBase表名。</li> <li>• 列族名：配置HBase列族名，如果HBase/Phoenix建表时未配置列族名，默认列族名为 '0'。</li> <li>• 列名：配置HBase列名。</li> <li>• 类型：配置字段类型，字段类型为“DATE”或“TIME”或“TIMESTAMP”时，需指定特定时间格式，其他类型指定无效。时间格式如：“yyyyMMdd HH:mm:ss”。</li> <li>• 长度：配置字段长度，字段值实际长度太长则按配置的长度截取，“类型”为“CHAR”时实际长度不足则空格补齐，“类型”为“VARCHAR”时实际长度不足则不补齐。</li> <li>• 主键：配置是否为主键列。普通HBase表主键只能指定一个；phoenix表主键可以指定多个，配置多个列为主键时，会按照配置列的先后顺序对其进行拼接。必须配置一个主键列。</li> </ul>	map	是	无

## 数据处理规则

- 将字段值输出到HBase表中。
- 原始数据包含NULL值，如果“NULL值处理方式”选中单选框时，将转换为空字符串并保存。如果“NULL值处理方式”不选中单选框时，不保存数据。

## 样例

以表输入为例，生成字段后，由HBase输出到对应的HBase表中，数据存放于test表中，如下图：

	id	name	value
1	1	zhangshan	zhang
2	2	lisi	li
3	3	wangwu	wang

创建HBase表：

```
create 'hbase_test','f1','f2';
```

配置“HBase输出”算子，如下图：

HBase表类型

NULL值处理方式

HBase输出字段

字段名	表名	列族名	列名	类型	长度	主键
<input type="text" value="A"/>	<input type="text" value="hbase_test"/>	<input type="text" value="f1"/>	<input type="text" value="A"/>	<input type="text" value="VARCHAR"/>	<input type="text"/>	<input checked="" type="checkbox"/>
<input type="text" value="B"/>	<input type="text" value="hbase_test"/>	<input type="text" value="f1"/>	<input type="text" value="B"/>	<input type="text" value="VARCHAR"/>	<input type="text"/>	<input type="checkbox"/>
<input type="text" value="C"/>	<input type="text" value="hbase_test"/>	<input type="text" value="f1"/>	<input type="text" value="C"/>	<input type="text" value="VARCHAR"/>	<input type="text"/>	<input type="checkbox"/>

作业执行成功后，查看hbase\_test表中数据：

```
hbase(main):001:0> scan 'hbase_test'
ROW
1
1
2
2
3
3
3 row(s) in 0.2720 seconds

COLUMN+CELL
column=f1:B, timestamp=1455855645760, value=zhangshan
column=f1:C, timestamp=1455855645760, value=zhang
column=f1:B, timestamp=1455855645760, value=lisi
column=f1:C, timestamp=1455855645760, value=li
column=f1:B, timestamp=1455855645760, value=wangwu
column=f1:C, timestamp=1455855645760, value=wang
```

## 17.15.4.6 ClickHouse 输出

### 概述

“ClickHouse输出”算子，用于配置已生成的字段输出到ClickHouse表的列。

### 输入与输出

- 输入：需要输出的字段。
- 输出：ClickHouse表。

### 参数说明

表 17-131 算子参数说明

参数	含义	类型	是否必填	默认值
数据库名	配置ClickHouse表所在的数据库。	string	是	default
表名	配置数据写入ClickHouse对应的表名。	string	是	无

### 数据处理规则

将字段值输出到ClickHouse表中。



## 样例

通过“CSV文件输入”算子，生成十二个字段。

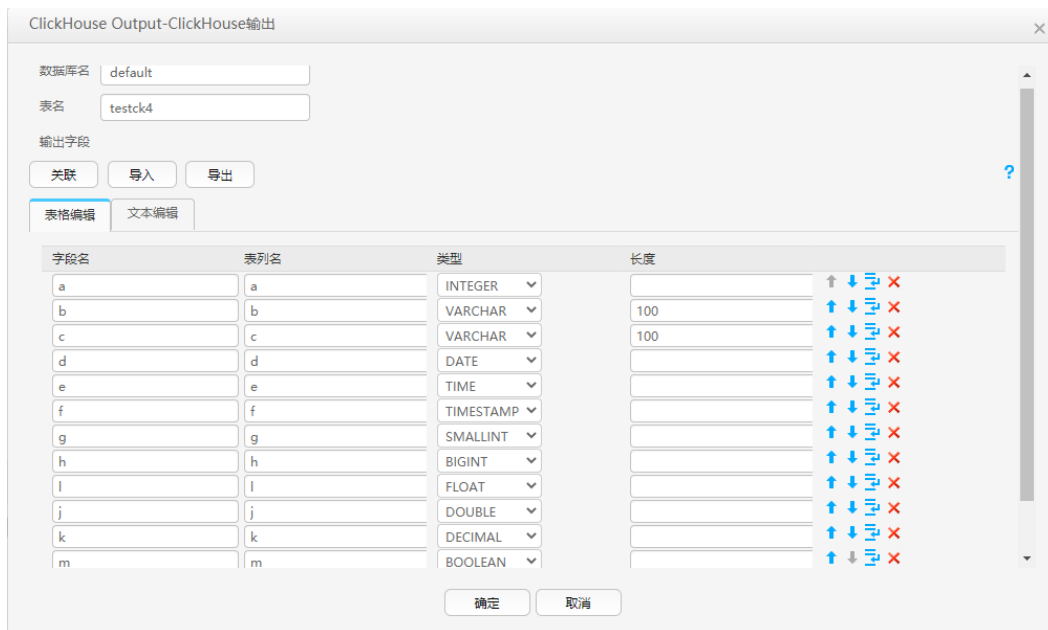
源文件如下：

```
1, 'b', 'abcd', '2021-06-15', '12:00:06', '2021-06-15 12:00:06', 1, 12, 6.8, 18.6, 12.8, true
2, 'abc', 'abcd', '2021-06-15', '12:00:06', '2021-06-15 12:00:06', 1, 12, 6.8, 18.6, 12.8, true
3, 'ab', 'abcd', '2021-06-15', '12:00:06', '2021-06-15 12:00:06', 1, 12, 6.8, 18.6, 12.8, true
4, 'abcdef', 'abcd', '2021-06-15', '12:00:06', '2021-06-15 12:00:06', 1, 12, 6.8, 18.6, 12.8, true
5, 'a', 'abcd', '2021-06-15', '12:00:06', '2021-06-15 12:00:06', 1, 12, 6.8, 18.6, 12.8, true
6, 'bg', 'cde', '2020-06-15', '13:00:06', '2021-06-15 12:00:06', 1, 12, 6.8, 18.6, 12.8, true
7, 'f', 'cde', '2020-06-15', '13:00:06', '2021-06-15 12:00:06', 1, 12, 6.8, 18.6, 12.8, true
8, 'h', 'cde', '2020-06-15', '13:00:06', '2021-06-15 12:00:06', 1, 12, 6.8, 18.6, 12.8, true
```

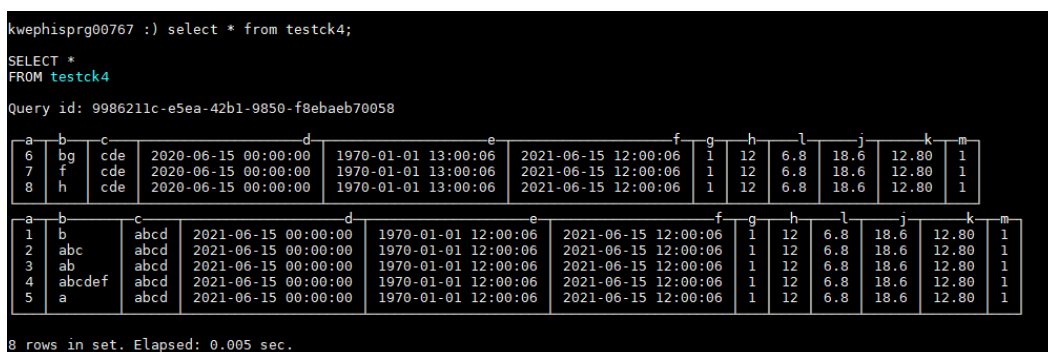
创建ClickHouse表的语句如下：

```
CREATE TABLE IF NOT EXISTS testck4 ON CLUSTER default_cluster(
 a Int32,
 b VARCHAR(100) NOT NULL,
 c char(100),
 d DateTime,
 e DateTime,
 f DateTime,
 g smallint,
 h bigint,
 l Float32,
 j Float64,
 k decimal(10,2),
 m boolean
)
ENGINE = ReplicatedMergeTree('/clickhouse/tables/{shard}/default/testck4',
'{replica}')
PARTITION BY toYYYYMM(d)ORDER BY a;
```

配置“ClickHouse输出”算子，如下图：



作业执行成功后，查看testck4表中数据：



## 17.15.5 关联、编辑、导入、导出算子的字段配置信息

### 操作场景

该任务指导用户在创建或编辑Loader作业时关联、导入或导出算子的字段配置信息。

- 关联操作  
将输入算子的字段配置信息关联到输出算子中。
- 编辑操作  
编辑算子配置参数中的字段信息。
- 导入操作  
通过算子导出文件或算子模板文件将字段配置信息导入到算子中。
- 导出操作  
将算子的字段配置信息以json文件导出保存到本地。

### 前提条件

获取登录“Loader WebUI”的账户和密码。

## 操作步骤

- 关联操作

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-82 Loader WebUI 界面



**步骤2** 编辑已有作业或者新建作业，进入“转换”界面。

**步骤3** 双击指定的输入算子（例如CSV文件输入）进入编辑页面，在输入字段的参数表格添加相应配置信息。

**步骤4** 双击指定的输出算子（例如文件输出）进入编辑页面，单击“关联”，并在弹出的“关联”对话框中勾选需要的字段信息。

### 说明

- 在输出算子的字段表格里已存在名称的字段信息，不会在“关联”窗口显示。
- 用户也可在“字段名”的列表中选择需要字段，相应配置信息会在输出字段的参数表格显示。

**步骤5** 单击“确定”，选中字段信息将会在输出字段的参数表格显示。

### ----结束

- 编辑操作

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-83 Loader WebUI 界面

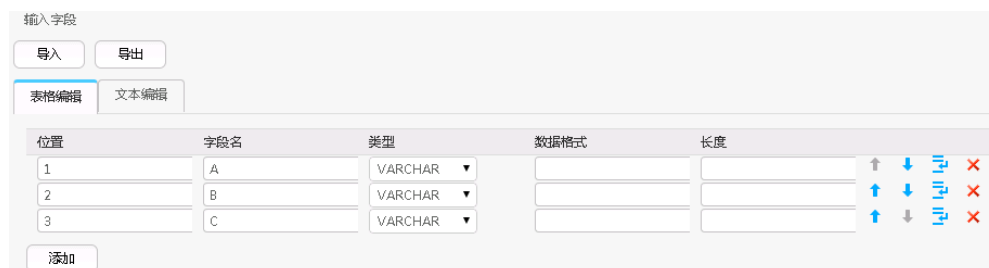


**步骤2** 编辑已有作业或者新建作业，进入“转换”界面。

**步骤3** 双击指定算子（例如CSV文件输入）进入编辑页面，在输入字段的“表格编辑”页签单击“添加”按钮，根据算子的参数格式要求填写相应字段信息。

**步骤4** 单击每行字段后的按钮可对字段进行上移、下移、下面插入一行以及删除等操作。

单击“文本编辑”，可以直接以文本形式对字段列表进行编辑，不同字段属性直接使用英文逗号“,”进行分隔。



**步骤5** 单击“确定”，保存字段信息。

----结束

- 导入操作

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-84 Loader WebUI 界面




**步骤2** 编辑已有作业或者新建作业，进入“转换”界面。

**步骤3** 双击指定的算子进入编辑页面，在输入或输出字段的参数表格添加相应配置信息。单击“导入”。

**步骤4** 选择导入的类型。

- 导出的文件  
通过算子导出的json文件导入字段的配置信息。
- 指导的模板  
通过根据算子模板手动编写txt文件，将字段配置信息导入到算子中。

**步骤5** 单击 ，选择上传文件对应路径。

**步骤6** 单击“上传”，字段的配置信息将会导入到算子。

----结束

- 导出操作

**步骤1** 登录“Loader WebUI”界面。

1. 登录FusionInsight Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群名称 > 服务 > Loader”。
3. 单击“LoaderServer(节点名称, 主)”打开“Loader WebUI”界面。

图 17-85 Loader WebUI 界面



**步骤2** 编辑已有作业或者新建作业，进入“转换”界面。

**步骤3** 双击指定的算子进入编辑页面，在输入或输出字段的参数表格添加相应配置信息，单击“导出”。

**步骤4** 选择导出的类型。

- 所有  
所有的字段信息将以json文件格式导出保存到本地。
- 指导字段  
在字段列表上勾选需要导出的字段以json文件格式导出保存到本地。

**步骤5** 单击“确定”，完成导出操作。

----结束

## 17.15.6 配置项中使用宏定义

用户在创建或者编辑Loader作业时，在配置参数时可以使用宏，在执行作业任务时会自动替换为宏对应的值。

### 说明

- 宏定义只在该作业范围内生效。
- 宏定义支持随作业导入导出，如果作业中有使用宏定义，则导出的作业包括宏定义。导入作业时默认也导入宏定义。
- 时间宏dataformat中的第一个参数的日期格式定义可参考“java.text.SimpleDateFormat.java”中的定义，但需要遵循目标系统的约束，例如HDFS/OBS目录不支持特殊符号等。

## Loader 宏定义

目前Loader默认支持以下时间宏定义：

表 17-132 Loader 常用宏定义

名称	替换后效果	说明
@{dateformat("yyyy-MM-dd")}@	2016-05-17	当前日期。
@{dateformat("yyyy-MM-dd HH:mm:ss")}@	2016-05-17 16:50:00	当前日期和时间。
@{timestamp()}@	1463476137557	从1970年到现在的毫秒数。
@{dateformat("yyyy-MM-dd HH:mm:ss",-7,DAYS)}@	2016-05-10 16:50:00	最近7天，即当前时间减7天。 第二个参数支持加减运算。 第三个参数为时间运算的单位，参考“java.util.concurrent.TimeUnit.java”定义，分为DAYS、HOURS、MINUTES、SECONDS。

在以下场景中，可以使用宏进行配置参数：

- 指定以当天时间命名的数据目录  
参数项配置为 “/user/data/inputdate\_@{dateformat("yyyy-MM-dd")}@”。
- 通过SQL语句查询最近7天的数据  
select \* from table where time between '@{dateformat("yyyy-MM-dd HH:mm:ss",-7,DAYS)}@' and '@{dateformat("yyyy-MM-dd HH:mm:ss")}@'
- 指定当天的表名  
参数项配置为 “table\_@{dateformat("yyyy-MM-dd")}@”。

## 17.15.7 算子数据处理规则

在Loader导入或导出数据的任务中，每个算子对于原始数据中NULL值、空字符串定义了不同的处理规则；在算子中无法正确处理的数据，将成为脏数据，无法导入或导出。

在转换步骤中，算子数据处理规则请参见下表。

表 17-133 数据处理规则一览表

转换步骤	规则描述
CSV文件输入	<ul style="list-style-type: none"> <li>分隔符在原始数据中连续出现两次，将生成空字符串字段。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>遇到类型转换错误，当前数据保存为脏数据。</li> </ul>

转换步骤	规则描述
固定宽度文件输入	<ul style="list-style-type: none"> <li>原始数据包含NULL值，不做转换处理。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>配置转换字段类型，与原始数据实际类型不同，全部数据成为脏数据。例如将字符串类型转换为数值类型。</li> <li>配置字段分割长度，大于原字段值的长度，则数据分割失败，当前行成为脏数据</li> </ul>
表输入	<ul style="list-style-type: none"> <li>原始数据包含NULL值，不做转换处理。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>配置转换字段类型，与原始数据实际类型不同，全部数据成为脏数据。例如将字符串类型转换为数值类型。</li> </ul>
HBase输入	<ul style="list-style-type: none"> <li>原始数据包含NULL值，不做转换处理。</li> <li>配置HBase表名错误，全部数据成为脏数据。</li> <li>“主键”没有配置主键列，全部数据成为脏数据。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>配置转换字段类型，与原始数据实际类型不同，全部数据成为脏数据。例如将字符串类型转换为数值类型。</li> </ul>
长整型时间转换	<ul style="list-style-type: none"> <li>原始数据包含NULL值，不做转换处理。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>遇到类型转换错误，当前数据保存为脏数据。</li> </ul>
空值转换	<ul style="list-style-type: none"> <li>原始数据包含NULL值，转换为用户指定的值。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> </ul>
随机值转换	不涉及处理NULL值、空字符串，不生成脏数据。
增加常量字段	不涉及处理NULL值、空字符串，不生成脏数据。
拼接转换	<ul style="list-style-type: none"> <li>原始数据包含NULL值，将转换为空字符串。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> </ul>
分隔转换	<ul style="list-style-type: none"> <li>原始数据包含NULL值，当前行成为脏数据。</li> <li>配置分割后字段列数，大于原始数据实际可分割出来的字段列数，当前行成为脏数据。</li> </ul>
取模转换	<ul style="list-style-type: none"> <li>原始数据包含NULL值，当前行成为脏数据。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>数据类型转换失败，当前行成为脏数据。</li> </ul>

转换步骤	规则描述
剪切字符串	<ul style="list-style-type: none"> <li>传入数据为NULL值，不做转换处理。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>字符截取的起点位置或终点位置，大于输入字段的长度时，当前行成为脏数据。</li> </ul>
EL操作转换	<ul style="list-style-type: none"> <li>传入数据为NULL值，不做转换处理。</li> <li>输入一个或多个字段的值，输出计算结果。</li> <li>输入类型和算子不兼容时，当前行为脏数据。</li> </ul>
字符串大小写转换	<ul style="list-style-type: none"> <li>传入数据为NULL值，不做转换处理。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> </ul>
字符串逆序转换	<ul style="list-style-type: none"> <li>传入数据为NULL值，不做转换处理。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> </ul>
字符串空格清除转换	<ul style="list-style-type: none"> <li>传入数据为NULL值，不做转换处理。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> </ul>
过滤行转换	<ul style="list-style-type: none"> <li>条件逻辑为“AND”，如果未添加过滤条件，全部数据成为脏数据；或者原始数据满足添加的全部过滤条件，当前行成为脏数据。</li> <li>条件逻辑为“OR”，如果未添加过滤条件，全部数据成为脏数据；或者原始数据满足任意添加的过滤条件，当前行成为脏数据。</li> </ul>
文件输出	<ul style="list-style-type: none"> <li>传入数据为NULL值，不做转换处理。</li> </ul>
表输出	<ul style="list-style-type: none"> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>数据类型转换失败，当前行成为脏数据。</li> </ul>
HBase输出	<ul style="list-style-type: none"> <li>原始数据包含NULL值，如果“NULL值处理方式”设置为“true”，将转换为空字符串并保存。如果“NULL值处理方式”设置为“false”，不保存数据。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>数据类型转换失败，当前行成为脏数据。</li> </ul>



转换步骤	规则描述
Hive输出	<ul style="list-style-type: none"> <li>如果指定了一个或多个列为分区列，则在“到”页面上，会显示“分割程序”属性，该属性表示使用多少个处理器去对分区数据进行处理。</li> <li>如果没有指定任何列为分区列，则表示不需要对输入数据进行分区处理，“分割程序”属性默认隐藏。</li> <li>配置输入字段列数，大于原始数据实际包含的字段列数，全部数据成为脏数据。</li> <li>数据类型转换失败，当前行成为脏数据。</li> </ul>

## 17.16 客户端工具说明

### 17.16.1 使用命令行运行 Loader 作业

#### 操作场景

一般情况下，用户可以手工在Loader界面管理数据导入导出作业。当用户需要通过shell脚本来更新与运行Loader作业时，必须对已安装的Loader客户端进行配置。

#### 说明

Loader不兼容旧版本客户端，如果重新安装集群或Loader服务，请重新下载并安装客户端，然后正常使用客户端。

本章节适用于MRS 3.x及后续版本。

#### 前提条件

- 完成Loader客户端的安装。使用非root用户安装Loader客户端时，如果其他用户也需要使用该客户端，则需要当前客户端的安装用户或者其他拥有更大权限的用户进行授权（将loader客户端的安装目录赋予“755”权限），请用户关注授权后的安全问题。
- 创建访问Loader服务的用户，如果是“机机”用户需要下载keytab文件。

#### 操作步骤

##### 步骤1 配置Loader shell客户端。

- 使用安装客户端的用户登录客户端所在节点。
- 执行以下命令，防止超时退出。

**TMOUT=0**

#### 说明

执行完本章节操作后，请及时恢复超时退出时间，执行命令**TMOUT=超时退出时间**。例如：**TMOUT=600**，表示用户无操作600秒后超时退出。

- 执行以下命令，进入Loader客户端安装目录。例如，Loader客户端安装目录为“/opt/client/Loader”。

**cd /opt/client/Loader**

4. 执行以下命令，配置环境变量。

**source/opt/client/bigdata\_env**

5. 如果集群为安全模式，执行以下命令进行用户认证。普通模式集群无需执行用户认证。

kinit *组件业务用户*

6. 执行以下命令修改工具授权配置文件“login-info.xml”，并保存退出。配置文件参数请参见表17-134。

**vi loader-tools-1.99.3/loader-tool/job-config/login-info.xml**

表 17-134 login-info.xml 参数

参数名称	描述
hadoop.config.path	填写MRS集群“core-site.xml”、“hdfs-site.xml”和“krb5.conf”三个配置文件的保存目录。默认保存在“Loader客户端安装目录/Loader/loader-tools-1.99.3/loader-tool/hadoop-config/”。
authentication.type	Loader服务的鉴权类型，请根据MRS集群认证模式填写： <ul style="list-style-type: none"> <li>- “kerberos”：表示安全模式。</li> <li>- “simple”：表示普通模式。</li> </ul>
user.keytab	是否使用keytab文件认证，参数值为“true”与“false”。
authentication.user	普通模式或者使用密码认证方式时，登录使用的用户。 keytab登录方式，则不需要设置该参数。
authentication.password	安全模式中若不使用keytab认证，配置访问Loader服务的用户密码加密字符串。 <b>说明</b> 使用安装客户端的用户执行以下命令加密密码。加密工具第一次执行时自动生成随机动态密钥并保存在“.loader-tools.key”中，加密工具每次加密密码时会使用此动态密钥。删除“.loader-tools.key”后加密工具执行时会重新生成新的随机密钥并保存在“.loader-tools.key”中。命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。 <b>sh Loader客户端安装目录/Loader/loader-tools-1.99.3/encrypt_tool password</b>

参数名称	描述
authentication.principal	安全模式中使用keytab认证，配置访问Loader服务的“机机”用户名。
authentication.keytab	安全模式中使用keytab认证，配置访问Loader服务的“机机”用户keytab文件目录，需包含绝对路径。
zookeeper.quorum	配置连接ZooKeeper节点的IP地址和端口，参数值格式为“IP1:port,IP2:port,IP3:port”，以此类推。默认端口号为“2181”。
sqoop.server.list	配置连接Loader的浮动IP和端口，参数值格式为“floatip:port”。默认端口号为“21351”。

## 步骤2 使用Loader shell客户端。

1. 执行以下命令，进入Loader shell客户端目录。例如，Loader客户端安装目录为“/opt/client/Loader”。

```
cd /opt/client/Loader/loader-tools-1.99.3/shell-client/
```

2. 执行以下命令，通过Loader shell客户端工具运行作业。

```
./submit_job.sh -n <arg> -u <arg> -jobType <arg> -connectorType <arg> -frameworkType <arg>
```

表 17-135 Loader shell 客户端工具参数一览表

参数名称	描述
“-n”	必配项，表示作业名称。
“-u”	必配项。 指定参数值为“y”表示更新作业参数并运行作业，此时需配置“-jobType”、“-connectorType”和“-frameworkType”。指定参数值为“n”表示不更新作业参数直接运行作业。
“-jobType”	表示作业类型，当“-u”的值为“y”时，必须配置。 指定参数值为“import”表示数据导入作业，指定参数值为“export”表示数据导出作业。

参数名称	描述
“-connectorType”	<p>表示连接器类型，当“-u”的值为“y”时，必须配置。根据业务需要可修改外部数据源的部分参数。</p> <p>指定参数值为“sftp”表示SFTP连接器。</p> <ul style="list-style-type: none"> <li>- 在导入作业中，支持修改源文件的输入路径“-inputPath”、源文件的编码格式“-encodeType”和源文件导入成功后对输入文件增加的后缀值“-suffixName”。</li> <li>- 在导出作业中，支持修改导出文件的路径或者文件名“-outputPath”。</li> </ul> <p>指定参数值为“rdb”表示关系型数据库连接器。</p> <ul style="list-style-type: none"> <li>- 在导入作业中，支持修改数据库模式名“-schemaName”、表名“-tableName”、SQL语句“-sql”、要导入的列名“-columns”和分区列“-partitionColumn”。</li> <li>- 在导出作业中，支持修改数据库模式名“-schemaName”、表名“-tableName”和临时表名称“-stageTableName”。</li> </ul>
“-frameworkType”	<p>表示MRS端数据保存的类型，当“-u”的值为“y”时，必须配置。根据业务需要可修改数据保存类型的部分参数。</p> <p>指定参数值为“hdfs”表示Hadoop端使用HDFS。</p> <ul style="list-style-type: none"> <li>- 在导入作业中，支持修改启动的map数量“-extractors”和数据导入到HDFS里存储的保存目录“-outputDirectory”。</li> <li>- 在导出作业中，支持修改启动的map数量“-extractors”、从HDFS导出时的输入路径“-inputDirectory”和导出作业的文件过滤条件“-fileFilter”。</li> </ul> <p>指定参数值为“hbase”表示MRS端使用HBase。在导入作业和导出作业中，支持修改启动的map数量“-extractors”。</p>

----结束

## 任务实例

- 不更新作业参数，直接运行名称为“sftp-hdfs”的作业。  
`./submit_job.sh -n sftp-hdfs -u n`
- 更新名称为“sftp-hdfs”导入作业的输入路径、编码类型、后缀、输出路径和启动的map数量参数，并运行作业。  
`./submit_job.sh -n sftp-hdfs -u y -jobType import -connectorType sftp -inputPath /opt/tempfile/1 -encodeType UTF-8 -suffixName " -frameworkType hdfs -outputDirectory /user/user1/tttest -extractors 10`
- 更新名称为“db-hdfs”导入作业的数据库模式、表名、输出路径参数，并运行作业。  
`./submit_job.sh -n db-hdfs -u y -jobType import -connectorType rdb -schemaName public -tableName sq_submission -sql " -partitionColumn sqs_id -frameworkType hdfs -outputDirectory /user/user1/dbdbt`

## 17.16.2 loader-tool 工具使用指导

### 概述

loader-tool工具是Loader客户端工具之一，包括“lt-ucc”、“lt-ucj”、“lt-ctl”三个工具。

Loader支持通过参数选项或作业模板这两种方式，对连接器进行创建、更新、查询和删除，以及对Loader作业进行创建、更新、查询、删除、启动和停止等操作。

本章节适用于MRS 3.x及后续版本。

#### 📖 说明

loader-tool工具是异步接口，命令提交后其结果不会实时返回到控制台，因此对连接器的创建、更新、查询和删除等操作，以及对Loader作业的创建、更新、查询、删除、启动和停止等操作，其成功与否需要在Loader WebUI确认或通过查询server端日志确认。

- 参数选项方式：

通过直接添加具体配置项的参数调用脚本。

- 作业模板方式：

修改作业模板中所有配置项的参数值，调用脚本时引用修改后的作业模板文件。

Loader客户端安装后，系统自动在“Loader客户端安装目录/loader-tools-1.99.3/loader-tool/job-config/”目录生成各种场景对应的作业模板，不同模板中配置项存在差异。作业模板中包含作业信息以及关联的连接器信息。

作业模板为xml文件，文件名格式为“数据原保存位置-to-数据新保存位置.xml”，例如“sftp-to-hdfs.xml”。如果此场景的作业支持转换步骤，则存在同名的转换步骤配置文件，文件类型为json，例如“sftp-to-hdfs.json”。

#### 📖 说明

作业模板中包含了连接器的配置信息。创建、更新连接器时，实际上仅调用到作业模板中的连接器的信息。

### 使用场景

不同的连接器或作业的配置项不同。

- 更新个别配置项时，使用参数选项方式。
- 创建连接器或作业时，使用作业模板方式。

**说明**

本工具目前支持FTP、HDFS、JDBC、MySQL、Oracle以及Oracle专用连接器，如果使用其他类型连接器，建议使用开源sqoop-shell工具。

**参数说明**

例如，Loader客户端的安装目录为：“/opt/client/Loader/”。

• **lt-ucc使用说明**

lt-ucc: loader-tool user-configuration-connection连接器配置工具，用于连接器的创建、更新和删除操作。

**表 17-136** lt-ucc 脚本“参数选项”说明

参数选项	说明	参数值示例
-help	获取帮助信息。	-
-a <arg>	执行的动作，有效值：create/update/delete，分别用于创建、更新和删除连接器。	create
-at <arg>	登录认证的类型，有效值kerberos、simple。	kerberos
-uk <arg>	是否使用keytab文件。	true
-au <arg>	登录认证的用户名。	bar
-ap <arg>	登录认证的密码，需要填写密文。 密码加密方法： <b>sh Loader客户端安装目录/Loader/loader-tools-1.99.3/encrypt_tool 用户非加密密码</b> <b>说明</b> 非加密密码中含有特殊字符时需要转义。例如，\$符号属于特殊字符，可使用单引号进行转义；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。	-
-c <arg>	登录认证的principal。	bar
-k <arg>	登录认证的keytab文件。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/hadoop-config/user.keytab

参数选项	说明	参数值示例
-h <arg>	MRS集群的配置文件路径。	-h /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/hadoop-config
-l <arg>	登录的模板文件。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml
-s <arg>	Loader服务的浮动IP和端口。 格式为： <i>浮动IP:端口</i> 端口默认值为21351	127.0.0.1:21351
-w <arg>	作业的模板文件路径，用于获取作业的详细信息。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.xml
-z <arg>	ZooKeeper quorum实例的IP地址和端口号，格式为 <i>IP地址:端口</i> ，配置多个用“,”分开。	127.0.0.0:2181, 127.0.0.1:2181
-n <arg>	连接器名称。	vt_sftp_test
-t <arg>	连接器类型。	sftp-connector
-P <arg>	更新某个属性的值，格式： -Pparam1=value1，param1为作业模板中连接器对应的属性名称。如果更新的是SFTP和FTP的连接器信息，还必须带上密码参数： -Pconnection.sftpPassword= <i>密码密文</i>	- Pconnection.sftpServerIp=10.6.26.11

完整示例如下：

```
./bin/lt-ucc -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -n vt_sftp_test -t sftp-connector -Pconnection.sftpPassword=密码密文 -Pconnection.sftpServerIp=10.6.26.111 -a update
```

lt-ucc脚本的作业模板配置说明：

以SFTP数据保存到HDFS为例，编辑“*loader客户端安装目录/loader-tools-1.99.3/loader-tool/job-config/*”目录下的“sftp-to-hdfs.xml”文件，连接器的配置如下：

```
<!-- 连接数据库的信息 -->
<sqoop.connection name="vt_sftp_test" type="sftp-connector">
<connection.sftpServerIp>10.96.26.111</connection.sftpServerIp>
<connection.sftpServerPort>22</connection.sftpServerPort>
```

```
<connection.sftpUser>root</connection.sftpUser>
<connection.sftpPassword>密码密文</connection.sftpPassword>
</sqoop.connection>
```

- 创建命令，如下：  
`./lt-ucc -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -w /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/ftp-to-hdfs.xml -a create`
- 更新命令，如下：  
`./lt-ucc -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -w /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/ftp-to-hdfs.xml -a update`
- 删除命令，如下：  
`./lt-ucc -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -w /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/ftp-to-hdfs.xml -a delete`

● **lt-ucj使用说明**

lt-ucj: loader-tool user-configuration-job作业配置工具，用于对作业的创作、更新、删除操作。

表 17-137 lt-ucj 脚本的“参数选项”配置说明

参数选项	说明	参数值示例
-help	获取帮助信息。	-
-a <arg>	执行的动作，有效值：create/update/delete，分别用于创建、更新和删除作业。	create
-at <arg>	登录认证的类型，有效值kerberos、simple。	kerberos
-uk <arg>	是否使用keytab文件。	true
-au <arg>	登录认证的用户名。	bar
-ap <arg>	登录认证的密码，需要填写密文。 密码加密方法： <b>sh Loader客户端安装目录/Loader/loader-tools-1.99.3/encrypt_tool 用户非加密密码</b> <b>说明</b> 非加密密码中含有特殊字符时需要转义。 例如，\$符号属于特殊字符，可使用单引号进行转义；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。	-
-c <arg>	登录认证的principal。	bar



参数选项	说明	参数值示例
-k <arg>	登录认证的keytab文件。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/hadoop-config/user.keytab
-h <arg>	MRS集群的配置文件路径。	-h /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/hadoop-config
-l <arg>	登录的模板文件。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml
-s <arg>	Loader服务的浮动IP和端口。 格式为： <i>浮动IP:端口</i> 端口默认值为21351。	127.0.0.1:21351
-w <arg>	作业的模板文件，用于获取作业的详细信息。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.xml
-z <arg>	ZooKeeper quorum实例的IP地址和端口号，格式为 <i>IP地址:端口</i> ，配置多个用“,”分开。	127.0.0.0:2181, 127.0.0.1:2181
-n <arg>	作业名称。	Sftp.to.Hdfs
-cn <arg>	连接器名称。	vt_sftp_test
-ct <arg>	连接器类型。	sftp-connector
-t <arg>	作业类型，有效值IMPORT、EXPORT。	IMPORT
-trans <arg>	作业关联的转换步骤文件。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.json
-priority <arg>	作业优先级，有效值：LOW/NORMAL/HIGH。	NORMAL
-queue <arg>	队列。	default

参数选项	说明	参数值示例
- storageType <arg>	存储类型。	HDFS
-P <arg>	更新某个属性的值，格式： - Pparam1=value1，param1为作业模板中连接器对应的属性名称。如果更新的是SFTP和FTP的连接器信息，还必须带上密码参数： -Pconnection.sftpPassword=密码密文	- Pconnection.sftpServerIp=10.6.26.11

完整示例如下：

```
./bin/lt-ucj -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -n Sftp.to.Hdfs -t IMPORT -ct sftp-connector -Poutput.outputDirectory=/user/loader/sftp-to-hdfs-test8888 -a update
```

lt-ucj 脚本的“作业模板”配置说明：

以SFTP数据保存到HDFS为例，编辑“loader客户端安装目录/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.xml”，作业的配置如下：

```
<!-- Job名称，全局唯一 -->
<sqoop.job name="Sftp.to.Hdfs" type="IMPORT" queue="default" priority="优先级NORMAL">

<!-- 外部数据源，参数配置 -->
<data.source connectionName="vt_sftp_test" connectionType="sftp-connector">
<file.inputPath>/opt/houjt/hive/all</file.inputPath>
<file.splitType>FILE</file.splitType>
<file.filterType>WILDCARD</file.filterType>
<file.pathFilter>*</file.pathFilter>
<file.fileFilter>*</file.fileFilter>
<file.encodeType>GBK</file.encodeType>
<file.suffixName></file.suffixName>
<file.isCompressive>FALSE</file.isCompressive>
</data.source>

<!-- MRS集群，参数配置 -->
<hadoop.source storageType="HDFS" >
<output.outputDirectory>/user/loader/sftp-to-hdfs</output.outputDirectory>
<output.fileOprType>OVERRIDE</output.fileOprType>
<throttling.extractors>3</throttling.extractors>
<output.fileType>TEXT_FILE</output.fileType>
</hadoop.source>

<!-- 作业关联的转换步骤文件 -->
<sqoop.job.trans.file>/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.json</sqoop.job.trans.file>
</sqoop.job>
```

- 创建命令，如下：

```
./bin/lt-ucj -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -w /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.xml -a create
```

- 更新命令，如下：

```
./bin/lt-ucj -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -w /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.xml -a update
```

- 删除命令，如下：  
`./bin/lt-ucj -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -w /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.xml -a delete`
- **lt-ctl使用说明**  
 lt-ctl: loader-tool controller作业管理工具，用于启停作业，查询作业状态与进度，查询作业是否运行中。

表 17-138 lt-ctl 脚本的“参数选项”配置说明

参数选项	说明	参数值示例
-help	获取帮助信息。	-
-a <arg>	执行的动作，有效值：status/start/stop/isrunning，分别用于查询作业状态、启动作业、停止作业以及判断作业是否在运行中。	create
-at <arg>	登录认证的类型，有效值kerberos、simple。	kerberos
-uk <arg>	是否使用keytab文件。	true
-au <arg>	登录认证的用户名。	bar
-ap <arg>	登录认证的密码，需要填写密文。 密码加密方法： <b>sh Loader客户端安装目录/Loader/loader-tools-1.99.3/encrypt_tool 用户非加密密码</b> <b>说明</b> 非加密密码中含有特殊字符时需要转义。例如，\$符号属于特殊字符，可使用单引号进行转义；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。	-
-c <arg>	登录认证的principal。	bar
-k <arg>	登录认证的keytab文件。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/hadoop-config/user.keytab
-h <arg>	MRS集群的配置文件路径。	-h /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/hadoop-config

参数选项	说明	参数值示例
-l <arg>	登录的模板文件。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml
-n <arg>	作业名称。	Sftp.to.Hdfs
-s <arg>	Loader服务的浮动IP和端口。 格式为： <i>浮动IP:端口</i> 端口默认值为21351。	127.0.0.1:21351
-w <arg>	作业的模板文件，用于获取作业的详细信息。	/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.xml
-z <arg>	ZooKeeper quorum实例的IP地址和端口号，格式为 <i>IP地址:端口</i> ，配置多个用“,”分开。	127.0.0.0:2181, 127.0.0.1:2181

- 启动作业：  
`./bin/lt-ctl -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -n Sftp.to.Hdfs -a start`
- 查看作业状态：  
`./bin/lt-ctl -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -n Sftp.to.Hdfs -a status`
- 判断作业是否运行中：  
`./bin/lt-ctl -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -n Sftp.to.Hdfs -a isrunning`
- 停止作业：  
`./bin/lt-ctl -l /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/login-info.xml -n Sftp.to.Hdfs -a stop`

### 17.16.3 loader-tool 工具使用示例

#### 操作场景

loader-tool工具支持通过作业模板或参数选项的方式，对连接器或者作业进行创建、更新、查询、删除等操作。

本文将“从SFTP服务器导入数据到HDFS”的作业为例，通过引用作业模板的方式，介绍loader-tool工具的使用方法。

本章节适用于MRS 3.x及后续版本。

## 前提条件

已安装并配置Loader客户端，具体操作请参见[使用命令行运行Loader作业](#)。

## 操作步骤

**步骤1** 使用安装客户端的用户登录客户端所在节点。

**步骤2** 执行以下命令，进入Loader客户端的loader-tool工具目录。例如，Loader客户端安装目录为“/opt/client/Loader/”。

```
cd /opt/client/Loader/loader-tools-1.99.3/loader-tool/
```

**步骤3** 执行以下命令，修改已有的作业模板。例如，“/opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/”目录下已有的作业模板“sftp-to-hdfs.xml”。

```
vi /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/sftp-to-hdfs.xml
```

```
<root>
<!-- 连接数据库的信息 -->
<sqoop.connection name="vt_sftp_test" type="sftp-connector">
<connection.sftpServerIp>10.96.26.111</connection.sftpServerIp>
<connection.sftpServerPort>22</connection.sftpServerPort>
<connection.sftpUser>root</connection.sftpUser>
<connection.sftpPassword>密码密文</connection.sftpPassword>
</sqoop.connection>

<!-- Job名称，全局唯一 -->
<sqoop.job name="Sftp.to.Hdfs" type="IMPORT" queue="default" priority="NORMAL">
<data.source connectionName="vt_sftp_test" connectionType="sftp-connector">
<file.inputPath>/opt/houjt/hive/all</file.inputPath>
<file.splitType>FILE</file.splitType>
<file.filterType>WILDCARD</file.filterType>
<file.pathFilter>*</file.pathFilter>
<file.fileFilter>*</file.fileFilter>
<file.encodeType>GBK</file.encodeType>
<file.suffixName></file.suffixName>
<file.isCompressive>FALSE</file.isCompressive>
</data.source>

<hadoop.source storageType="HDFS" >
<output.outputDirectory>/user/loader/sftp-to-hdfs</output.outputDirectory>
<output.fileOprType>OVERRIDE</output.fileOprType>
<throttling.extractors>3</throttling.extractors>
<output.fileType>TEXT_FILE</output.fileType>
</hadoop.source>

<sqoop.job.trans.file></sqoop.job.trans.file>
</sqoop.job>
</root>
```

### 说明

Loader每个作业都需要关联一个连接器，连接器主要作用：对于数据导入到集群的场景来说，就是从外部数据源读取数据；对于数据从集群导出出去的场景来说，就是将数据写入到外部数据源。上述示例配置的是一个SFTP数据源连接器。配置SFTP和FTP的数据源连接器需要设置密码并进行加密。密码加密方法如下：

1. 执行以下命令，进入到loader-tools-1.99.3目录。Loader客户端安装目录为“/opt/hadoopclient/Loader”。

```
cd /opt/hadoopclient/Loader/loader-tools-1.99.3
```

2. 执行以下命令，对非加密密码加密。

```
./encrypt_tool 未加密的密码
```

**步骤4** 执行以下命令，进入loader-tool工具目录。

```
cd /opt/client/Loader/loader-tools-1.99.3/loader-tool
```

**步骤5** 执行以下命令，使用lt-ucc工具创建连接器。

```
./bin/lt-ucc -l /opt/client/Loader/loader-tools-1.99.3/loader-tool/job-config/
login-info.xml -w /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/
job-config/sftp-to-hdfs.xml -a create
```

如无报错信息，且显示如下信息，则表示创建连接器的任务提交成功。

```
User login success. begin to execute task.
```

**步骤6** 执行以下命令，使用lt-ucj工具创建作业。

```
./bin/lt-ucj -l /opt/client/Loader/loader-tools-1.99.3/loader-tool/job-config/
login-info.xml -w /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/
job-config/sftp-to-hdfs.xml -a create
```

如无报错信息，且显示如下信息，则表示创建作业的任务提交成功。

```
User login success. begin to execute task.
```

**步骤7** 执行以下命令，使用lt-ctl工具提交作业。

```
./bin/lt-ctl -l /opt/client/Loader/loader-tools-1.99.3/loader-tool/job-config/
login-info.xml -n Sftp.to.Hdfs -a start
```

显示如下信息，表示作业提交成功。

```
Start job success.
```

**步骤8** 执行以下命令，查看作业状态。

```
./bin/lt-ctl -l /opt/client/Loader/loader-tools-1.99.3/loader-tool/job-config/
login-info.xml -n Sftp.to.Hdfs -a status
```

```
Job:Sftp.to.Hdfs
Status:RUNNING
Progress: 0.0
```

----结束

## 17.16.4 schedule-tool 工具使用指导

### 概述

schedule-tool工具，用于提交数据源为SFTP的作业。提交作业前可以修改输入路径、文件过滤条件，当目标源为HDFS时，可以修改输出路径。

本章节适用于MRS 3.x及后续版本。

## 参数说明

表 17-139 schedule.properties 配置参数说明

配置参数	说明	示例
server.url	Loader服务的浮动IP地址和端口。 端口默认为21351。 为了兼容性，此处支持配置多个IP地址和端口，并以“,”进行分隔。其中第一个必须是Loader服务的浮动IP地址和端口，其余的可根据业务需求配置。	10.96.26.111:213 51,127.0.0.2:2135 1
authentication.type	登录认证的方式。 <ul style="list-style-type: none"><li>“kerberos”，表示使用安全模式，进行Kerberos认证。Kerberos认证提供两种认证方式：密码和keytab文件。</li><li>“simple”，表示使用普通模式，不进行Kerberos认证。</li></ul>	kerberos
authentication.user	普通模式或者使用密码认证方式时，登录使用的用户。 keytab登录方式，则不需要设置该参数。	bar

配置参数	说明	示例
authentication.password	<p>使用密码认证方式时，登录使用的用户密码。普通模式或者keytab登录方式，则不需要设置该参数。</p> <p>用户需要对密码加密，加密方法如下：</p> <ol style="list-style-type: none"> <li>1. 进入“encrypt_tool”所在目录。例如，Loader客户端安装目录为“/opt/hadoopclient/Loader”，则执行如下命令。 <b>cd /opt/hadoopclient/Loader/loader-tools-1.99.3</b></li> <li>2. 执行以下命令，对非加密密码进行加密。 <b>./encrypt_tool 未加密的密码</b> 得到加密后的密文，作为“authentication.password”的取值。</li> </ol> <p><b>说明</b> 非加密密码中含有特殊字符时需要转义。例如，\$符号属于特殊字符，可使用单引号进行转义；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。</p> <p>命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。</p>	-
use.keytab	<p>是否使用keytab方式登录。</p> <ul style="list-style-type: none"> <li>• true，表示使用keytab文件登录</li> <li>• false，表示使用密码登录。</li> </ul>	true
client.principal	<p>使用keytab认证方式时，访问Loader服务的用户规则。</p> <p>普通模式或者密码登录方式，则不需要设置该参数。</p>	<p>loader/hadoop.&lt;系统域名&gt;</p> <p><b>说明</b> 用户可登录FusionInsight Manager，选择“系统 &gt; 权限 &gt; 域和互信”，查看“本端域”参数，即为当前系统域名。</p>
client.keytab	<p>使用keytab认证方式登录时，使用的keytab文件所在目录。</p> <p>普通模式或者密码登录方式，则不需要设置该参数。</p>	/opt/client/conf/loader.keytab



配置参数	说明	示例
krb5.conf.file	使用keytab认证方式登录时，使用的krb5.conf文件所在目录。 普通模式或者密码登录方式，则不需要设置该参数。	/opt/client/conf/ krb5.conf

表 17-140 job.properties 配置参数说明

配置参数	说明	示例
job.jobName	作业的名称。	job1
file.fileName.prefix	文件名的前缀。	table1
file.fileName.posfix	文件名的后缀。	.txt
file.filter	文件过滤器，通过匹配文件名来过滤文件。 <ul style="list-style-type: none"> <li>“true”，表示用上面的前缀/后缀，来匹配输入路径下的所有文件。详细使用，见最后示例。</li> <li>“false”，表示用上面的前缀/后缀，来匹配输入路径下的某一个文件。详细使用，见最后示例。</li> </ul>	true
date.day	顺延的天数，匹配导入文件的文件名中的日期。例如命令参数传入的日期是20160202，顺延天数是3，则匹配作业配置的输入路径中包含20160205日期字段的文件。详细使用见 <a href="#">schedule-tool工具使用示例</a> 。	3
file.date.format	待导入文件的文件名中所包含的日志格式。	yyyyMMdd
parameter.date.format	调用脚本时，所输入的日期格式，一般保持与“file.date.format”一致。	yyyyMMdd
file.format.iscompressed	待导入的文件是否为压缩文件。	false
storage.type	存储类型。待导入文件最终保存的类型，分别有HDFS、HBase、Hive等。	HDFS

## 说明

schedule-tool工具支持同时配置多个作业。配置多个作业时，表17-140中“job.jobName”、“file.fileName.prefix”、“file.fileName.posfix”参数需配置多个值，并且以“,”分隔。

## 注意事项

server.url属性必须需要配置两个IP地址和端口的格式串，用“,”分隔。

## 17.16.5 schedule-tool 工具使用示例

### 操作场景

通过Loader WebUI或客户端工具Loader-tool创建好作业后，可使用schedule-tool工具执行作业。

本章节适用于MRS 3.x及后续版本。

### 前提条件

完成了Loader客户端的安装与配置，具体操作请参见[使用命令行运行Loader作业](#)。

### 操作步骤

- 步骤1** 在SFTP服务器的“/opt/houjt/test03”路径中，创建多个以“table1”为前缀，“.txt”为后缀，中间为yyyyMMdd的日期格式的文件。

图 17-86 示例

```
[root@C12-RHEL64-ZYL111 test03]# ll
total 36
-rw-r--r--. 1 root root 54 Feb 29 19:11 table120160221.txt
-rw-r--r--. 1 root root 54 Feb 29 19:11 table120160222.txt
-rw-r--r--. 1 root root 54 Feb 29 19:11 table120160223.txt
-rw-r--r--. 1 root root 54 Feb 29 19:11 table120160224.txt
-rw-r--r--. 1 root root 54 Feb 29 19:11 table120160225.txt
-rw-r--r--. 1 root root 54 Feb 29 19:11 table120160226.txt
-rw-r--r--. 1 root root 54 Feb 29 18:43 table120160227.txt
-rw-r--r--. 1 root root 54 Feb 29 19:11 table120160228.txt
-rw-r--r--. 1 root root 54 Feb 29 19:11 table120160229.txt
```

- 步骤2** 创建一个从SFTP服务器导入数据到HDFS的Loader作业，具体操作请参见[典型场景：从SFTP服务器导入数据到HDFS/OBS](#)。
- 步骤3** 使用安装客户端的用户登录客户端所在节点。
- 步骤4** 执行以下命令，进入schedule-tool工具的conf目录。例如，Loader客户端安装目录为“/opt/client/Loader/”。
- ```
cd /opt/client/Loader/loader-tools-1.99.3/schedule-tool/conf
```
- 步骤5** 执行以下命令，编辑schedule.properties文件，配置登录方式。

```
vi schedule.properties
```

schedule-tool工具支持两种登录方式，两者只能选一。详细参数请参见[schedule-tool 工具使用指导](#)。

- 以密码方式登录，配置信息示例如下：
[server.url = 10.10.26.187:21351,127.0.0.2:21351]
[authentication.type = kerberos]
[use.keytab = false]
[authentication.user = admin]
密码明文存储存在安全风险，建议在配置文件或者环境变量中密文存放，使用时解密，确保安全
[authentication.password= xxx]

- 以keytab文件方式登录，配置信息示例如下：
[server.url = 10.10.26.187:21351,127.0.0.2:21351]
[authentication.type = kerberos]
[use.keytab = true]
[client.principal = bar]
[client.keytab = /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/hadoop-config/user.keytab]
[krb5.conf.file = /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/hadoop-config/krb5.conf]

步骤6 执行以下命令，编辑job.properties文件，配置作业信息。

vi job.properties

```
#job name
job.jobName = sftp2hdfs-schedule-tool

#Whether to update the loader configuration parameters(File filter)£?This parameter is used to match the
import file name.Values are true or false.
#false means update.the file name which is get by schedule tool will be updated to Loader configuration
parameters (File filter).
#false means no update.the file name which is get by schedule tool will be updated to Loader configuration
parameters (import path).
file.filter = false

#File name = prefix + date + suffix
#Need to import the file name prefix
file.fileName.prefix=table1

#Need to import the file name suffixes
file.fileName.posfix=.txt

#Date Days.Value is an integer.
#According to the date and number of days to get the date of the import file.
date.day = 1

#Date Format.Import file name contains the date format.Format Type£°yyyyMMdd,yyyyMMdd
HHmmss,yyyy-MM-dd,yyyy-MM-dd HH:mm:ss
file.date.format = yyyyMMdd

#Date Format.Scheduling script execution. Enter the date format.
parameter.date.format = yyyyMMdd

#Whether the import file is a compressed format.Values ??are true or false.
#true indicates that the file is a compressed format£?Execution scheduling tool will extract the files.false
indicates that the file is an uncompressed.Execution scheduling tool does not unpack.
file.format.iscompressed = false

#Hadoop storage type.Values are HDFS or HBase.
storage.type = HDFS
```

根据**步骤1**的所准备的数据，以文件table120160221.txt为例，过滤规则设置如下：

- 文件名的前缀
file.fileName.prefix=table1
- 文件名的后缀
file.fileName.posfix=.txt

- 文件名中包含的日期格式
file.date.format = yyyyMMdd
 - 调用脚本输入的日期参数
parameter.date.format = yyyyMMdd
 - 顺延的天数
date.day = 1
- 例如，脚本传入的日期参数是20160220，则通过加法计算，得到的结果是20160221。

📖 说明

如果执行的命令是 `./run.sh 20160220 /user/loader/schedule_01`时，以上过滤规则会拼凑出一个字符串：`"table1"+"20160221"+.txt = table120160221.txt`

步骤7 根据file.filter的值，选择过滤规则。

- 精确匹配某一个文件，请执行**步骤8**。
- 模糊匹配一系列文件，请执行**步骤9**。

步骤8 将job.properties文件中“file.filter”的值修改为“false”。

执行以下命令，运行作业，任务结束。

```
cd /opt/client/Loader/loader-tools-1.99.3/schedule-tool
```

```
./run.sh 20160220 /user/loader/schedule_01
```

其中20160220为输入的日期，/user/loader/schedule_01为输出的路径。

📖 说明

通过以上过滤规则，拼凑得到的字符串“table120160221.txt”，会直接作为文件名，追加到作业配置的输入路径中。所以，作业只会处理唯一匹配到的文件“table120160221.txt”。

步骤9 将job.properties文件中“file.filter”的值修改为“true”，“file.fileName.prefix”设置为“*”。

执行以下命令，运行作业，任务结束。

```
cd /opt/client/Loader/loader-tools-1.99.3/schedule-tool
```

```
./run.sh 20160220 /user/loader/schedule_01
```

其中20160220为输入的日期，/user/loader/schedule_01为输出的路径。

📖 说明

通过以上过滤规则，拼凑到的字符串“*20160221.txt”，会作为文件过滤器的模糊匹配模式，在作业配置的输入路径下，所有符合“*20160221.txt”这个模式的文件都将被作业处理。

----结束

17.16.6 使用 loader-backup 工具备份作业数据

操作场景

通过Loader WebUI或客户端工具loader-tool创建好作业后，可使用loader-backup工具进行数据备份。

说明

- 仅有数据导出的Loader作业才支持数据备份。
- 此工具为Loader的内部接口，供上层组件HBase调用，只支持HDFS到SFTP的数据备份。

本章节适用于MRS 3.x及后续版本。

前提条件

完成了Loader客户端的安装与配置，具体操作请参见[使用命令行运行Loader作业](#)。

操作步骤

步骤1 使用安装客户端的用户登录客户端所在节点，具体操作请参见[使用命令行运行Loader作业](#)。

步骤2 执行以下命令，进入“backup.properties”文件所在目录。例如，Loader客户端安装目录为“/opt/client/Loader/”。

```
cd /opt/client/Loader/loader-tools-1.99.3/loader-backup/conf
```

步骤3 执行以下命令，修改“backup.properties”文件的配置参数，参数具体说明如表17-141所示。

vi backup.properties

```
server.url = 10.0.0.1:21351,10.0.0.2:12000
authentication.type = kerberos
authentication.user =
authentication.password=
job.jobId = 1
use.keytab = true
client.principal = loader/hadoop
client.keytab = /opt/client/conf/loader.keytab
```

表 17-141 配置参数说明

| 配置参数 | 说明 | 示例 |
|---------------------|--|-------------------------------|
| server.url | Loader服务的浮动IP地址和端口（21351）。
为了兼容性，此处支持配置多个IP地址和端口，并以“,”进行分隔。其中第一个必须是Loader服务的浮动IP地址和端口（21351），其余的可根据业务需求配置。 | 10.0.0.1:21351,10.0.0.2:12000 |
| authentication.type | 登录认证的方式。
<ul style="list-style-type: none"> • “kerberos”，表示使用安全模式，进行Kerberos认证。Kerberos认证提供两种认证方式：密码和keytab文件。 • “simple”，表示使用普通模式，不进行Kerberos认证。 | kerberos |

| 配置参数 | 说明 | 示例 |
|-------------------------|--|---------------|
| authentication.user | 普通模式或者使用密码认证方式时，登录使用的用户。
keytab登录方式，则不需要设置该参数。 | bar |
| authentication.password | 使用密码认证方式时，登录使用的用户密码。
普通模式或者keytab登录方式，则不需要设置该参数。
用户需要对密码加密，加密方法：
1. 进入“encrypt_tool”所在目录。例如，Loader客户端安装目录为“/opt/hadoopclient/Loader”，则执行如下命令。
cd /opt/hadoopclient/Loader/loader-tools-1.99.3
2. 执行以下命令，对非加密密码进行加密。
./encrypt_tool 未加密的密码
得到加密后的密文，作为“authentication.password”的取值。
说明
非加密密码中含有特殊字符时需要转义。例如，\$符号属于特殊字符，可使用单引号进行转义；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。
命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。 | - |
| job.jobId | 需要执行数据备份的作业ID。
作业ID可通过登录Loader webUI在已创建的作业查看。 | 1 |
| use.keytab | 是否使用keytab方式登录。
<ul style="list-style-type: none"> • true，表示使用keytab文件登录 • false，表示使用密码登录。 | true |
| client.principal | 使用keytab认证方式时，访问Loader服务的用户规则。
普通模式或者密码登录方式，则不需要设置该参数。 | loader/hadoop |

| 配置参数 | 说明 | 示例 |
|---------------|---|--------------------------------|
| client.keytab | 使用keytab认证方式登录时，使用的keytab文件所在目录。
普通模式或者密码登录方式，则不需要设置该参数。 | /opt/client/conf/loader.keytab |

步骤4 执行以下命令，进入备份脚本“run.sh”所在目录。例如，Loader客户端安装目录为“/opt/hadoopclient/Loader”。

```
cd /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-backup
```

步骤5 执行以下命令，运行备份脚本“run.sh”，进行Loader作业数据备份。系统将数据备份到作业的输出路径同一层目录。

```
./run.sh 备份数据的输入目录
```

例如，备份数据的输入目录为“/user/hbase/”，作业的输出路径为/opt/client/sftp/sftp1，其中sftp1只起到一个占位符的作用。执行如下命令，数据将备份到/opt/client/sftp/hbase目录。

```
./run.sh /user/hbase/
```

----结束

17.16.7 开源 sqoop-shell 工具使用指导

概述

本章节适用于MRS 3.x及后续版本。

sqoop-shell是一个开源的shell工具，其所有功能都是通过执行脚本“sqoop2-shell”来实现的。

sqoop-shell工具提供了如下功能：

- 支持创建和更新连接器
- 支持创建和更新作业
- 支持删除连接器和作业
- 支持以同步或异步的方式启动作业
- 支持停止作业
- 支持查询作业状态
- 支持查询作业历史执行记录
- 支持复制连接器和作业
- 支持创建和更新转换步骤
- 支持指定行、列分隔符

sqoop-shell工具支持如下模式：

- 交互模式
通过执行不带参数的“sqoop2-shell”脚本，进入Loader特定的交互窗口，用户输入脚本后，工具会返回相应信息到交互窗口。

- 批量模式
通过执行“sqoop2-shell”脚本，带一个文件名作为参数，该文件中按行存储了多条命令，sqoop-shell工具将会按顺序执行文件中所有命令；或者在“sqoop2-shell”脚本后面通过“-c”参数附加一条命令，一次只执行一条命令。

sqoop-shell通过[表17-142](#)的命令来实现Loader各种功能。

表 17-142 命令一览表

| 命令 | 说明 |
|---------|--------------------------|
| exit | 表示退出交互模式。
该命令仅支持交互模式。 |
| history | 查看执行过的命令。
该命令仅支持交互模式。 |
| help | 查看工具帮助信息。 |
| set | 设置服务端属性。 |
| show | 显示服务属性和Loader所有元数据信息。 |
| create | 创建连接器和作业。 |
| update | 更新连接器和作业。 |
| delete | 删除连接器和作业。 |
| clone | 复制连接器和作业。 |
| start | 启动作业。 |
| stop | 停止作业。 |
| status | 查询作业状态。 |

命令参考

- sqoop2-shell有两种获取登录认证信息的方式，第一种通过配置文件获取，具体配置项请参考[开源sqoop-shell工具使用示例（SFTP - HDFS）](#)、[开源sqoop-shell工具使用示例（Oracle - HBase）](#)；第二种方式则使用参数直接提供认证信息，这个方式有两种模式：密码模式和Kerberos认证模式。
- 进入交互模式命令
通过执行不带参数的“sqoop2-shell”脚本，进入sqoop工具窗口，逐条执行命令。
通过读取配置文件获取认证信息：
./sqoop2-shell
通过密码模式认证：
./sqoop2-shell -uk false -u username -p encryptedPassword

命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的 history 命令记录功能，避免信息泄露。

通过 Kerberos 模式认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal
```

系统显示如下信息：

```
Welcome to sqoop client
Use the username and password authentication mode
Authentication success.
Sqoop Shell: Type 'help' or '\h' for help.

sqoop:000>
```

- 进入批量模式命令

进入批量模式有两种方式：

1. 通过执行“sqoop2-shell”脚本，带一个文本文件名作为参数，该文件中按行存储了多条命令，工具会按顺序执行该文件中的所有命令。使用这种方式有个限制条件，这个 sh 脚本必须放到当前用户的家目录下，如：`/root/batchCommand.sh`。

通过读取配置文件进行认证：

```
./sqoop2-shell /root/batchCommand.sh
```

通过密码模式认证：

```
./sqoop2-shell -uk false -u username -p encryptedPassword /root/  
batchCommand.sh
```

通过 Kerberos 模式认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal /root/  
batchCommand.sh
```

其中 `batchCommand.sh` 为用户自定义文本文件名称。

2. 通过执行“sqoop2-shell”脚本，在脚本后面通过“-c”参数附带一条命令，工具将执行该条命令。

通过取配置文件进行认证：

```
./sqoop2-shell -c expression
```

通过密码模式认证：

```
./sqoop2-shell -uk false -u username -p encryptedPassword -c expression
```

通过 Kerberos 模式认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal -c expression
```

其中 `expression` 为附带的语句，其格式和第一种方式中的文本内语句格式一致。

- exit 命令

该命令用于退出交互模式，仅在交互模式支持。

示例：

```
Welcome to sqoop client
Use the username and password authentication mode
Authentication success.
Sqoop Shell: Type 'help' or '\h' for help.

sqoop:000> exit
10-5-211-9:/opt/hadoopclient/Loader/loader-tools-1.99.3/sqoop-shell#
```

- history 命令

该命令用于查看已执行的命令，仅在交互模式支持。

示例：

```
sqoop:000> history
0 show connector
1 create connection -c 4
2 show connections;
3 show connection;
4 show connection -a;
5 show connections;
6 show connection;
7 show connection -x 53;
8 show connection -x 52;
9 show connection -x 2
10 show connection -x 53;
11 show connection
12 show connection -x 53
13 create job -x 53 -t import
14 show connector
15 create connection -c 5
16 show connection -x 54
17 exit
18 show connector
19 create connection -c 5
20 exit
21 show connector
22 create connection -c 6
23 create job -x 20 -t import
24 start job -j 85 -s
25 \x
26 exit
27 history
sqoop:000>
```

- help命令

该命令用于查看工具帮助信息。

示例：

```
sqoop:000> help
For information about Sqoop, visit: http://sqoop.apache.org/docs/1.99.3/index.html

Available commands:
exit (\x ) Exit the shell
history (\H ) Display, manage and recall edit-line history
help (\h ) Display this help message
set (\st ) Set server or option Info
show (\sh ) Show server, connector, framework, connection, job, submission or option Info
create (\cr ) Create connection or job Info
delete (\d ) Delete connection or job Info
update (\up ) Update connection or job Info
clone (\cl ) Clone connection or job Info
start (\sta) Start job
stop (\stp) Stop job
status (\stu) Status job

For help on a specific command type: help command

sqoop:000>
```

- set命令

set命令，用于设置客户端和服务端属性，支持如下属性：

- server表示设置服务端连接属性。

 说明

当设置了-u属性时，-h、-p、-w被会忽略。

- option表示设置客户端属性。

 说明

option通过键值对来赋值，例如：`set option --name verbose --value true`。

| 属性类别 | 子属性 | 含义 |
|--------|--------------|----------------|
| server | -h,--host | 服务IP地址 |
| | -p,--port | 服务端口 |
| | -w,--webapp | Tomcat应用名 |
| | -u,--url | Sqoop服务URL |
| option | verbose | 冗余模式，表示打印更多的信息 |
| | poll-timeout | 设置轮询超时时间 |

示例：

```
set option --name verbose --value false
set server --host 10.0.0.1 --port 21351 --webapp loader
```

- show命令

该命令用于显示变量信息、存储元数据信息等。

| 属性类别 | 子属性 | 含义 |
|------------|-------------|---------------|
| server | -a,--all | 显示所有server属性 |
| | -p,--port | 显示服务端口 |
| | -w,--webapp | 显示Tomcat应用名 |
| | -h,--host | 显示服务的IP地址 |
| option | -name | 显示指定名称的属性 |
| connector | -a,--all | 显示所有连接类型信息 |
| | -c,--cid | 显示指定ID的连接类型信息 |
| framework | 无 | 显示框架的元数据信息 |
| connection | -a,--all | 显示所有连接属性 |
| | -x,--xid | 显示指定ID的连接属性 |
| | -n,--name | 显示指定名称的连接属性 |
| job | -a,--all | 显示所有作业信息 |
| | -j,--jid | 显示指定ID的作业信息 |
| | -n,--name | 显示指定名称的作业信息 |
| submission | -j,--jid | 显示指定作业的提交记录 |

| 属性类别 | 子属性 | 含义 |
|------|-------------|--------|
| | -d,--detail | 显示详细信息 |

示例：

```
show server -all
show option --name verbose
show connector -all
show framework
show connection -all
show connection -n sftp-example
show job -all
show job -j 1
show submission --jid 1
show submission --jid 1 -d
```

- create命令

该命令用于创建连接器或作业。

| 属性类别 | 子属性 | 含义 |
|------------|-------------|--|
| connection | -c,--cid | 指定连接器类型的ID |
| | -cn,--cname | 指定连接器类型的名称 |
| job | -x,--xid | 指定连接器ID |
| | -xn,--xname | 指定连接器名称 |
| | -t,--type | 指定作业类型
可选值：
<ul style="list-style-type: none"> • import • export |

- 交互模式下，根据界面的提示逐一输入属性值。

创建连接器示例：

```
create connection -c 1
create connection -cn example
```

创建作业示例：

```
create job -x 1 -t import
create job -xn job_example -t export
```

- 批量模式下，需要先执行如下命令查看具体的属性，再对属性赋值。

create job -t import -x 1 --help

执行该命令有两种方式：

将命令保存到文本中，并在执行sqoop-shell脚本时将该文本作为附带参数：

```
./sqoop2-shell batchCommand.sh
```

使用-c参数，将需要执行的单条命令作为-c参数的输入：

```
./sqoop2-shell -c expression
```

可参考本节前文关于命令执行的描述。完整的命令语句可参考如下示例。

创建连接器示例：

```
create connection -c 4 --connector-connection-sftpPassword xxxxx --connector-connection-sftpServerIp 10.0.0.1 --connector-connection-sftpServerPort 22 --connector-connection-sftpUser root--name testConnection
```

创建作业示例：

```
create job -t import -x 1 --connector-file-inputPath /opt/tempfile --connector-file-fileFilter * --framework-output-outputDirectory /user/loader/1 --framework-output-storageType HDFS --framework-throttling-extractorSize 120 --framework-output-fileType TEXT_FILE --connector-file-splitType FILE -queue default -priority low -name newJob
```

- 批量模式下，可以使用“-c”参数附带一条语句。

创建连接器示例：

```
./sqoop2-shell -c "create connection -c 4 --connector-connection-sftpPassword xxxxx --connector-connection-sftpServerIp 10.0.0.1 --connector-connection-sftpServerPort 22 --connector-connection-sftpUser root--name testConnection"
```

- update命令

该命令用于更新连接器或作业。

| 属性类别 | 子属性 | 含义 |
|------------|----------|---|
| connection | -x,--xid | 指定连接器ID
说明
更新连接器一定要带上密码属性。 |
| job | -j,--jid | 指定作业ID |

- 交互模式

更新连接器示例：

```
update connection --xid 1
```

更新作业示例：

```
update job --jid 1
```

- 批量模式

更新连接器示例：

```
update connection -x 6 --connector-connection-sftpServerPort 21 - --name sfp_130--connector-connection-sftpPassword xxxx
```

更新作业示例：

示例1：

```
update job -jid 1 -name sftp2hdfs --connector-file-fileFilter *.txt
```

示例2：

```
./sqoop2-shell -uk true -k /opt/loader/user.keytab -s luser /opt/loader/testupdate.txt
./sqoop2-shell -uk true -k /opt/loader/user.keytab -s luser -c "update job --jid 24 --name oracle-hive --connector-table-sql 'SELECT * FROM range_example WHERE replace(datadt,\'-\',\'.\')=\'20240801\' and \${CONDITIONS}'"
```

 说明

更新作业可以将需要更新的命令写在文件中，例如“/opt/loader/testupdate.txt”（文件名自定义），也可以以--connector-table-sql来指定，后面跟随的sqlcmd需要用“'”单引号括起来，具体操作参考“更新作业示例-示例2”。涉及的命令还有connector-table-sql,connector-table-columns,connector-table-partitionColumn,connector-table-conditions,connector-table-queryCondition等。

- delete命令

该命令用于删除连接器或作业。

| 属性类别 | 子属性 | 含义 |
|------------|-----------|---------|
| connection | -x,--xid | 指定连接器ID |
| | -n,--name | 指定连接器名称 |
| job | -j,--jid | 指定作业ID |
| | -n,--name | 指定作业名称 |

示例：

```
delete connection -x 1
delete connection --name abc
delete job -j 1
delete job -n qwerty
```

- clone命令

该命令用于复制连接器或作业。

| 属性类别 | 子属性 | 含义 |
|------------|----------|--|
| connection | -x,--xid | 指定连接器ID
说明
复制连接器需要输入密码和连接器名称。 |
| job | -j,--jid | 指定作业ID |

示例如下：

```
clone job -j 1
```

- start命令

该命令用于启动作业。

| 属性类别 | 子属性 | 含义 |
|------|------------------|--------|
| job | -j,--jid | 指定作业ID |
| | -n,--name | 指定作业名称 |
| | -s,--synchronous | 是否同步 |

异步启动作业示例：

```
start job -j 1
start job -n abc
```

同步启动作业示例：

```
start job -j 1 -s
start job --name abc --synchronous
```

- stop命令
该命令用于停止作业。

| 属性类别 | 子属性 | 含义 |
|------|-----------|--------|
| job | -j,--jid | 指定作业ID |
| | -n,--name | 指定作业名称 |

示例：

```
stop job -j 1
stop job -n abc
```

- status命令
该命令用于查询作业状态。

| 属性类别 | 子属性 | 含义 |
|------|----------|--------|
| job | -j,--jid | 指定作业ID |

查询状态时，可以使用“-s”参数，只查询作业的状态枚举。

示例：

```
status job -j 1
status job -j 1 -s
```

create 命令扩展属性

针对HDFS与SFTP服务器或RDB进行数据交换场景，MRS在开源sqoop-shell工具的基础上对create命令属性进行扩展，以达到在创建作业时指定行、列分隔符及转换步骤的目的。

表 17-143 create 命令扩展属性

| 属性 | 说明 |
|-----------------------------|---|
| fields-terminated-by | 默认的列分隔符。 |
| lines-terminated-by | 默认的行分隔符。 |
| input-fields-terminated-by | 输入步骤的列分隔符，当不指定时，默认等于fields-terminated-by的值。 |
| input-lines-terminated-by | 输入步骤的行分隔符，当不指定时，默认等于lines-terminated-by的值。 |
| output-fields-terminated-by | 输出步骤的列分隔符，当不指定时，默认等于fields-terminated-by的值。 |
| output-lines-terminated-by | 输出步骤的行分隔符，当不指定时，默认等于lines-terminated-by的值。 |

| 属性 | 说明 |
|-------|--|
| trans | 指定转换步骤，值为转换步骤文件所在的路径。当指定文件的相对路径时，默认为“sqoop2-shell”脚本所在路径下的文件。当配置了该属性，其他扩展属性都被忽略。 |

sqoop1 对接 MRS 服务

步骤1 下载开源Sqoop，<http://www.apache.org/dyn/closer.lua/sqoop/1.4.7>。

步骤2 将下载好的sqoop-1.4.7.bin__hadoop-2.6.0.tar.gz 包放入MRS集群master节点的/opt/sqoop目录下并解压。

```
tar zxvf sqoop-1.4.7.bin__hadoop-2.6.0.tar.gz
```

步骤3 进入解压完成的目录，修改配置。

```
cd /opt/sqoop/sqoop-1.4.7.bin__hadoop-2.6.0/conf
```

```
cp sqoop-env-template.sh sqoop-env.sh
```

```
vi sqoop-env.sh
```

添加配置：

```
export HADOOP_COMMON_HOME=/opt/client/HDFS/hadoop
```

```
export HADOOP_MAPRED_HOME=/opt/client/HDFS/hadoop
```

```
export HIVE_HOME=/opt/Bigdata/MRS_1.9.X/install/FusionInsight-Hive-3.1.0/hive  
(请按照实际路径填写)
```

```
export HIVE_CONF_DIR=/opt/client/Hive/config
```

```
export HCAT_HOME=/opt/client/Hive/HCatalog
```

步骤4 添加系统变量，将“SQOOP_HOME”添加到PATH中。

```
vi /etc/profile
```

添加以下信息：

```
export SQOOP_HOME=/opt/sqoop/sqoop-1.4.7.bin__hadoop-2.6.0
```

```
export PATH=$PATH:$SQOOP_HOME/bin
```

步骤5 执行以下命令复制jline-2.12.jar文件到lib文件下。

```
cp /opt/share/jline-2.12/jline-2.12.jar /opt/sqoop/  
sqoop-1.4.7.bin__hadoop-2.6.0/lib
```

步骤6 执行以下命令，在文件中添加下列配置。

```
vim $JAVA_HOME/jre/lib/security/java.policy
```

```
permission javax.management.MBeanTrustPermission "register";
```


步骤7 执行以下命令，实现sqoop1对接MRS服务。

```
source /etc/profile  
----结束
```

17.16.8 开源 sqoop-shell 工具使用示例（SFTP - HDFS）

操作场景

本文将以“从SFTP服务器导入数据到HDFS”的作业为例，介绍如何分别在交互模式和批量模式下使用sqoop-shell工具进行创建和启动Loader作业。

本章节适用于MRS 3.x及后续版本。

前提条件

已安装并配置Loader客户端，具体操作请参见[使用命令行运行Loader作业](#)。

交互模式示例

步骤1 使用安装客户端的用户登录Loader客户端所在节点。

步骤2 执行以下命令，进入sqoop-shell工具的“conf”目录。例如，Loader客户端安装目录为“/opt/client/Loader/”。

```
cd /opt/client/Loader/loader-tools-1.99.3/sqoop-shell/conf
```

步骤3 执行以下命令，配置认证信息。

```
vi client.properties  
server.url=10.0.0.1:21351  
# simple or kerberos  
authentication.type=simple  
# true or false  
use.keytab=true  
  
authentication.user=  
authentication.password=  
  
client.principal=hdfs/hadoop@<系统域名>  
  
# keytab file  
client.keytab.file=./conf/login/hdfs.keytab
```

说明

登录FusionInsight Manager，选择“系统 > 权限 > 域和互信”，“本端域”参数即为当前系统域名。

表 17-144 配置参数说明

| 配置参数 | 说明 | 示例 |
|---------------------|--|----------------|
| server.url | Loader服务的浮动IP地址和端口（21351）。
为了兼容性，此处支持配置多个IP地址和端口，并以“,”进行分隔。其中第一个必须是Loader服务的浮动IP地址和端口（21351），其余的可根据业务需求配置。 | 10.0.0.1:21351 |
| authentication.type | 登录认证的方式。 <ul style="list-style-type: none">“kerberos”，表示使用安全模式，进行Kerberos认证。Kerberos认证提供两种认证方式：密码和keytab文件。“simple”，表示使用普通模式，不进行Kerberos认证。 | kerberos |
| authentication.user | 普通模式或者使用密码认证方式时，登录使用的用户。
keytab登录方式，则不需要设置该参数。 | bar |

| 配置参数 | 说明 | 示例 |
|-------------------------|--|--------------------------------|
| authentication.password | <p>使用密码认证方式时，登录使用的用户密码。</p> <p>普通模式或者keytab登录方式，则不需要设置该参数。</p> <p>用户需要对密码加密，加密方法：</p> <ol style="list-style-type: none"> 1. 进入“encrypt_tool”所在目录。例如，Loader客户端安装目录为“/opt/hadoopclient/Loader”，则执行如下命令。
cd /opt/hadoopclient/Loader/loader-tools-1.99.3 2. 执行以下命令，对非加密密码进行加密。
./encrypt_tool 未加密的密码
得到加密后的密文，作为“authentication.password”的取值。 <p>说明
非加密密码中含有特殊字符时需要转义。例如，\$符号属于特殊字符，可使用单引号进行转义；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。</p> <p>命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。</p> | - |
| use.keytab | <p>是否使用keytab方式登录。</p> <ul style="list-style-type: none"> • true，表示使用keytab文件登录 • false，表示使用密码登录。 | true |
| client.principal | <p>使用keytab认证方式时，访问Loader服务的用户规则。</p> <p>普通模式或者密码登录方式，则不需要设置该参数。</p> | loader/hadoop |
| client.keytab.file | <p>使用keytab认证方式登录时，使用的keytab文件所在目录。</p> <p>普通模式或者密码登录方式，则不需要设置该参数。</p> | /opt/client/conf/loader.keytab |

步骤4 执行以下命令，进入交互模式。

```
source /opt/client/bigdata_env
```

```
cd /opt/client/Loader/loader-tools-1.99.3/sqoop-shell
```

./sqoop2-shell

上述命令通过读取配置文件获取认证信息。

也可以直接通过密码或者Kerberos认证。

使用密码进行认证：

```
./sqoop2-shell -uk false -u username -p encryptedPassword
```

命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。

使用Kerberos认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal
```

```
Welcome to sqoop client
Use the username and password authentication mode
Authentication success.
Sqoop Shell: Type 'help' or '\h' for help.

sqoop:000>
```

步骤5 执行以下命令，查看当前连接器对应的ID。

show connector

显示如下信息：

```
+-----+-----+-----+-----+
| Id | Name | Version | Class |
+-----+-----+-----+-----+
1	generic-jdbc-connector	2.0.6-SNAPSHOT	org.apache.sqoop.connector.jdbc.GenericJdbcConnector
2	ftp-connector	2.0.5-SNAPSHOT	org.apache.sqoop.connector.ftp.FtpConnector
3	hdfs-connector	2.0.5-SNAPSHOT	org.apache.sqoop.connector.hdfs.HdfsConnector
4	oracle-connector	2.0.1-SNAPSHOT	org.apache.sqoop.connector.oracle.OracleConnector
5	mysql-fastpath-connector	2.0.1-SNAPSHOT	org.apache.sqoop.connector.mysql.MySqlConnector
6	sftp-connector	2.0.5-SNAPSHOT	org.apache.sqoop.connector.sftp.SftpConnector
7	oracle-partition-connector	2.0.6-SNAPSHOT	org.apache.sqoop.connector.oracle.partition.OraclePartitionConnector
+-----+-----+-----+-----+
```

根据如上信息，可知SFTP连接器类型ID为6。

步骤6 执行如下命令，创建连接器，根据提示输入具体的连接器信息。

```
create connection -c 连接器类型ID
```

例如，连接器类型的ID为6，则执行如下命令：

```
create connection -c 6
```

```
sqoop:000> create connection -c 6
Creating connection for connector with id 6
Please fill following values to create new connection object
Name: sftp14

Connection configuration

Sftp server IP: 10.0.0.1
Sftp server port: 22
```

```
Sftp user name: root
Sftp password: *****
Sftp public key:
New connection was successfully created with validation status FINE and persistent id 20
sqoop:000>
```

根据如上信息，可知连接器ID为20。

步骤7 根据连接器ID，执行如下命令，创建作业。

create job -x 连接器ID -t import

例如，连接器ID为20，则执行如下命令：

create job -x 20 -t import

显示如下信息：

```
Creating job for connection with id 20
Please fill following values to create new job object
Name: sftp-hdfs-test

File configuration

Input path: /opt/tempfile
File split type:
  0 : FILE
  1 : SIZE
Choose: 0
Filter type:
  0 : WILDCARD
  1 : REGEX
Choose: 0
Path filter: *
File filter: *
Encode type:
Suffix name:
Compression:

Output configuration

Storage type:
  0 : HDFS
  1 : HBASE_BULKLOAD
  2 : HBASE_PUTLIST
  3 : HIVE
Choose: 0
File type:
  0 : TEXT_FILE
  1 : SEQUENCE_FILE
  2 : BINARY_FILE
Choose: 0
Compression format:
  0 : NONE
  1 : DEFAULT
  2 : DEFLATE
  3 : GZIP
  4 : BZIP2
  5 : LZ4
  6 : SNAPPY
Choose:
Output directory: /user/loader/test
File operate type:
  0 : OVERRIDE
  1 : RENAME
  2 : APPEND
  3 : IGNORE
  4 : ERROR
Choose: 0
```

```
Throttling resources

Extractors: 2
Extractor size:
New job was successfully created with validation status FINE and persistent id 85
sqoop:000>
```

根据如上信息，可知作业ID为85。

步骤8 执行以下命令，启动作业。

```
start job -j 作业ID -s
```

例如，作业ID为85，则执行如下命令：

```
start job -j 85 -s
```

显示“SUCCEEDED”信息，则说明作业启动成功。

```
Submission details
Job ID: 85
Server URL: https://10.0.0.0:21351/loader/
Created by: admin
Creation date: 2016-07-20 16:25:38 GMT+08:00
Lastly updated by: admin
2016-07-20 16:25:38 GMT+08:00: BOOTING - Progress is not available
2016-07-20 16:25:46 GMT+08:00: BOOTING - 0.00 %
2016-07-20 16:25:53 GMT+08:00: BOOTING - 0.00 %
2016-07-20 16:26:08 GMT+08:00: RUNNING - 90.00 %
2016-07-20 16:26:08 GMT+08:00: RUNNING - 90.00 %
2016-07-20 16:26:27 GMT+08:00: SUCCEEDED
```

----结束

批量模式示例

步骤1 使用安装客户端的用户登录Loader客户端所在节点。

步骤2 执行以下命令，进入sqoop-shell工具的“conf”目录。例如，Loader客户端安装目录为“/opt/client/Loader/”。

```
cd /opt/client/Loader/loader-tools-1.99.3/sqoop-shell/conf
```

步骤3 执行以下命令，配置认证信息。

```
vi client.properties
server.url=10.0.0.1:21351
# simple or kerberos
authentication.type=simple
# true or false
use.keytab=true

authentication.user=
authentication.password=

client.principal=hdfs/hadoop@<系统域名>

# keytab file
client.keytab.file=./conf/login/hdfs.keytab
```

表 17-145 配置参数说明

| 配置参数 | 说明 | 示例 |
|-------------------------|---|----------------|
| server.url | Loader服务的浮动IP地址和端口（21351）。
为了兼容性，此处支持配置多个IP地址和端口，并以“,”进行分隔。其中第一个必须是Loader服务的浮动IP地址和端口（21351），其余的可根据业务需求配置。 | 10.0.0.1:21351 |
| authentication.type | 登录认证的方式。
<ul style="list-style-type: none"> “kerberos”，表示使用安全模式，进行Kerberos认证。Kerberos认证提供两种认证方式：密码和keytab文件。 “simple”，表示使用普通模式，不进行Kerberos认证。 | kerberos |
| authentication.user | 普通模式或者使用密码认证方式时，登录使用的用户。
keytab登录方式，则不需要设置该参数。 | bar |
| authentication.password | 使用密码认证方式时，登录使用的用户密码。
普通模式或者keytab登录方式，则不需要设置该参数。
用户需要对密码加密，加密方法：
<ol style="list-style-type: none"> 进入“encrypt_tool”所在目录。例如，Loader客户端安装目录为“/opt/hadoopclient/Loader”，则执行如下命令。
cd /opt/hadoopclient/Loader/loader-tools-1.99.3 执行以下命令，对非加密密码进行加密。
./encrypt_tool 未加密的密码
得到加密后的密文，作为“authentication.password”的取值。
说明
非加密密码中含有特殊字符时需要转义。例如，\$符号属于特殊字符，可使用单引号进行转义；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。 | - |

| 配置参数 | 说明 | 示例 |
|--------------------|---|--------------------------------|
| use.keytab | 是否使用keytab方式登录。
<ul style="list-style-type: none"> • true，表示使用keytab文件登录 • false，表示使用密码登录。 | true |
| client.principal | 使用keytab认证方式时，访问Loader服务的用户规则。
普通模式或者密码登录方式，则不需要设置该参数。 | loader/hadoop |
| client.keytab.file | 使用keytab认证方式登录时，使用的keytab文件所在目录。
普通模式或者密码登录方式，则不需要设置该参数。 | /opt/client/conf/loader.keytab |

步骤4 执行以下命令，进入“sqoop2-shell”脚本所在目录，并在该目录下创建一个文本文件，例如“batchCommand.sh”。

```
cd /opt/client/Loader/loader-tools-1.99.3/sqoop-shell
```

```
vi batchCommand.sh
```

“batchCommand.sh” 样例如下：

```
//查看参数
create connection -c 6 --help

//创建连接器
create connection -c 6 -name sftp-connection --connector-connection-sftpServerIp 10.0.0.1 --connector-connection-sftpServerPort 22 --connector-connection-sftpUser root --connector-connection-sftpPassword xxxxx

//创建作业
create job -t import -x 20 --connector-file-inputPath /opt/tempfile --connector-file-fileFilter * --framework-output-outputDirectory /user/loader/1 --framework-output-storageType HDFS --framework-throttling-extractorSize 120 --framework-output-fileType TEXT_FILE --connector-file-splitType FILE -name test

//启动作业
start job -j 85 -s
```

其中xxxxx为连接器密码。

步骤5 执行如下命令，sqoop-shell工具将依次执行上述命令。

```
./sqoop2-shell batchCommand.sh
```

也可以直接在命令里附带认证信息。

使用密码认证：

```
./sqoop2-shell -uk false -u username -p encryptedPassword batchCommand.sh
```

使用Kerberos认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal batchCommand.sh
```

显示“SUCCEEDED”信息，则说明作业启动成功。

```
Welcome to sqoop client
Use the username and password authentication mode
```



```
Authentication success.
sqoop:000> create connection -c 6 --help
usage: Show connection parameters:
  --connector-connection-sftpPassword <arg>
  --connector-connection-sftpServerIp <arg>
  --connector-connection-sftpServerPort <arg>
  --connector-connection-sftpUser <arg>
  --framework-security-maxConnections <arg>
  --name <arg>
====> FINE
sqoop:000> create connection -c 6 -name sftp-connection --connector-connection-sftpServerIp 10.0.0.1 --
connector-connection-sftpServerPort 22 --connector-connection-sftpUser root --connector-connection-
sftpPassword xxxxx
Creating connection for connector with id 6
New connection was successfully created with validation status FINE and persistent id 20
====> FINE
sqoop:000> create job -t import -x 20 --connector-file-inputPath /opt/tempfile --connector-file-fileFilter * --
framework-output-outputDirectory /user/loader/1 --framework-output-storageType HDFS --framework-
throttling-extractorSize 120 --framework-output-fileType TEXT_FILE --connector-file-splitType FILE -name
test
Creating job for connection with id 20
New job was successfully created with validation status FINE and persistent id 85
====> FINE

Submission details
Job ID: 85
Server URL: https://10.0.0.0:21351/loader/
Created by: admin
Creation date: 2016-07-20 16:25:38 GMT+08:00
Lastly updated by: admin
2016-07-20 16:25:38 GMT+08:00: BOOTING - Progress is not available
2016-07-20 16:25:46 GMT+08:00: BOOTING - 0.00 %
2016-07-20 16:25:53 GMT+08:00: BOOTING - 0.00 %
2016-07-20 16:26:08 GMT+08:00: RUNNING - 90.00 %
2016-07-20 16:26:08 GMT+08:00: RUNNING - 90.00 %
2016-07-20 16:26:27 GMT+08:00: SUCCEEDED
```

步骤6 批处理模式下，使用-c参数附带一条命令，sqoop-shell可以一次只执行附带的这一条命令。

执行如下命令将创建连接器。

```
./sqoop2-shell -c "create connection -c 6 -name sftp-connection --connector-
connection-sftpServerIp 10.0.0.1 --connector-connection-sftpServerPort 22 --
connector-connection-sftpUser root --connector-connection-sftpPassword
xxxxx"
```

可以在命令里直接附带认证信息。

使用密码认证：

```
./sqoop2-shell -uk false -u username -p encryptedPassword -c "create
connection -c 6 -name sftp-connection --connector-connection-sftpServerIp
10.0.0.1 --connector-connection-sftpServerPort 22 --connector-connection-
sftpUser root --connector-connection-sftpPassword xxxxx"
```

使用Kerberos认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal -c "create connection -c 6
-name sftp-connection --connector-connection-sftpServerIp 10.0.0.1 --
connector-connection-sftpServerPort 22 --connector-connection-sftpUser root
--connector-connection-sftpPassword xxxxx"
```

显示“FINE”信息，则说明连接创建成功。

```
Welcome to sqoop client
Use the username and password authentication mode
```

```
Authentication success.
sqoop:000> create connection -c 6 -name sftp-connection --connector-connection-sftpServerIp 10.0.0.1 --
connector-connection-sftpServerPort 22 --connector-connection-sftpUser root --connector-connection-
sftpPassword xxxx
Creating connection for connector with id 6
New connection was successfully created with validation status FINE and persistent id 20
====> FINE
```

----结束

17.16.9 开源 sqoop-shell 工具使用示例（Oracle - HBase）

操作场景

本文将以“从Oracle导入数据到HBase”的作业为例，介绍如何分别在交互模式和批量模式下使用sqoop-shell工具进行创建和启动Loader作业。

本章节适用于MRS 3.x及后续版本。

前提条件

已安装并配置Loader客户端，具体操作请参见[使用命令行运行Loader作业](#)。

交互模式示例

步骤1 使用安装客户端的用户登录Loader客户端所在节点。

步骤2 执行以下命令，进入sqoop-shell工具的“conf”目录。例如，Loader客户端安装目录为“/opt/client/Loader/”。

```
cd /opt/client/Loader/loader-tools-1.99.3/sqoop-shell/conf
```

步骤3 执行以下命令，配置认证信息。

```
vi client.properties
```

```
server.url=10.0.0.1:21351
# simple or kerberos
authentication.type=simple
# true or false
use.keytab=true

authentication.user=
authentication.password=

client.principal=oracle/hadoop@<系统域名>

# keytab file
client.keytab.file=./conf/login/oracle.keytab
```

说明

登录FusionInsight Manager，选择“系统 > 权限 > 域和互信”，“本端域”参数即为当前系统域名。

表 17-146 配置参数说明

| 配置参数 | 说明 | 示例 |
|---------------------|--|----------------|
| server.url | Loader服务的浮动IP地址和端口（21351）。
为了兼容性，此处支持配置多个IP地址和端口，并以“,”进行分隔。其中第一个必须是Loader服务的浮动IP地址和端口（21351），其余的可根据业务需求配置。 | 10.0.0.1:21351 |
| authentication.type | 登录认证的方式。 <ul style="list-style-type: none">“kerberos”，表示使用安全模式，进行Kerberos认证。Kerberos认证提供两种认证方式：密码和keytab文件。“simple”，表示使用普通模式，不进行Kerberos认证。 | kerberos |
| authentication.user | 普通模式或者使用密码认证方式时，登录使用的用户。
keytab登录方式，则不需要设置该参数。 | bar |

| 配置参数 | 说明 | 示例 |
|-------------------------|--|--------------------------------|
| authentication.password | <p>使用密码认证方式时，登录使用的用户密码。</p> <p>普通模式或者keytab登录方式，则不需要设置该参数。</p> <p>用户需要对密码加密，加密方法：</p> <ol style="list-style-type: none"> 1. 进入“encrypt_tool”所在目录。例如，Loader客户端安装目录为“/opt/hadoopclient/Loader”，则执行如下命令。
cd /opt/hadoopclient/Loader/loader-tools-1.99.3 2. 执行以下命令，对非加密密码进行加密。
./encrypt_tool 未加密的密码
得到加密后的密文，作为“authentication.password”的取值。 <p>说明
非加密密码中含有特殊字符时需要转义。例如，\$符号属于特殊字符，可使用单引号进行转义；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。</p> <p>命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。</p> | - |
| use.keytab | <p>是否使用keytab方式登录。</p> <ul style="list-style-type: none"> • true，表示使用keytab文件登录 • false，表示使用密码登录。 | true |
| client.principal | <p>使用keytab认证方式时，访问Loader服务的用户规则。</p> <p>普通模式或者密码登录方式，则不需要设置该参数。</p> | loader/hadoop |
| client.keytab.file | <p>使用keytab认证方式登录时，使用的keytab文件所在目录。</p> <p>普通模式或者密码登录方式，则不需要设置该参数。</p> | /opt/client/conf/loader.keytab |

步骤4 执行以下命令，进入交互模式。

```
source /opt/client/bigdata_env
```

```
cd /opt/client/Loader/loader-tools-1.99.3/sqoop-shell
```

./sqoop2-shell

上述命令通过读取配置文件获取认证信息。

也可以直接通过密码或者Kerberos认证。

使用密码进行认证：

```
./sqoop2-shell -uk false -u username -p encryptedPassword
```

命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。

使用Kerberos认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal
```

```
Welcome to sqoop client
Use the username and password authentication mode
Authentication success.
Sqoop Shell: Type 'help' or '\h' for help.

sqoop:000>
```

步骤5 执行以下命令，查看当前连接器对应的ID。

show connector

显示如下信息：

| Id | Name | Version | Class |
|----|----------------------------|----------------|--|
| 1 | generic-jdbc-connector | 2.0.7-SNAPSHOT | org.apache.sqoop.connector.jdbc.GenericJdbcConnector |
| 2 | ftp-connector | 2.0.5-SNAPSHOT | org.apache.sqoop.connector.ftp.FtpConnector |
| 3 | hdfs-connector | 2.0.5-SNAPSHOT | org.apache.sqoop.connector.hdfs.HdfsConnector |
| 4 | oracle-connector | 2.0.1-SNAPSHOT | org.apache.sqoop.connector.oracle.OracleConnector |
| 5 | mysql-fastpath-connector | 2.0.1-SNAPSHOT | org.apache.sqoop.connector.mysql.MySqlConnector |
| 6 | sftp-connector | 2.0.6-SNAPSHOT | org.apache.sqoop.connector.sftp.SftpConnector |
| 7 | oracle-partition-connector | 2.0.6-SNAPSHOT | org.apache.sqoop.connector.oracle.partition.OraclePartitionConnector |

根据如上信息，可知oracle连接器类型ID为4。

步骤6 执行如下命令，创建连接器，根据提示输入具体的连接器信息。

create connection -c 连接器类型ID

例如，连接器类型的ID为4，则执行如下命令：

create connection -c 4

```
sqoop:000> create connection -c 4
Creating connection for connector with id 4
Please fill following values to create new connection object
Name: oracle14

Oracle connection configuration

JDBC connection string: jdbc:oracle:thin:@189.120.84.106:1521:orcl
Username: oracledba
Password: *****
```

```
JDBC connection properties:  
There are currently 0 values in the map:  
entry#  
New connection was successfully created with validation status FINE and persistent id 3  
sqoop:000>
```

根据如上信息，可知连接器ID为3。

步骤7 根据连接器ID，执行如下命令，创建作业。

create job -x 连接器ID -t import --trans job-config目录的绝对路径/oracle-hbase.json

例如，连接器ID为3，则执行如下命令：

create job -x 3 -t import --trans /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/oracle-hbase.json

显示如下信息：

```
sqoop:000> create job -x 3 -t import --trans /opt/hadoopclient/Loader/loader-tools-1.99.3/loader-tool/job-config/oracle-to-hbase.json  
Creating job for connection with id 3  
Please fill following values to create new job object  
Name: run  
  
Database target  
  
Table name: test  
Columns:  
Conditions:  
Data split method:  
  0 : ROWID  
  1 : PARTITION  
Choose:  
Table Partitions:  
Data split allocation method:  
  0 : ROUNDROBIN  
  1 : SEQUENTIAL  
  2 : RANDOM  
Choose:  
JDBC fetch size:  
  
Output configuration  
  
Storage type:  
  0 : HDFS  
  1 : HBASE_BULKLOAD  
  2 : HBASE_PUTLIST  
  3 : HIVE  
  4 : SPARK  
Choose: 1  
HBase instance: HBase  
Clear data before import : false  
  
Throttling resources  
  
Extractors: 10  
Extractor size:  
New job was successfully created with validation status FINE and persistent id 7  
sqoop:000>
```

根据如信息，而知作业ID为7。

步骤8 执行以下命令，启动作业。

start job -j 作业ID -s

例如，作业ID为7，则执行如下命令：

start job -j 7 -s

显示“SUCCEEDED”信息，则说明作业启动成功。

```
Submission details
Job ID: 7
Server URL: https://10.0.0.0:21351/loader/
Created by: admintest
Creation date: 2019-12-04 16:37:34 CST
Lastly updated by: admintest
2019-12-04 16:37:34 CST: BOOTING - Progress is not available
2019-12-04 16:37:42 CST: BOOTING - 0.00 %
2019-12-04 16:37:42 CST: BOOTING - 0.00 %
2019-12-04 16:37:57 CST: RUNNING - 0.00 %
2019-12-04 16:38:12 CST: RUNNING - 45.00 %
2019-12-04 16:38:12 CST: RUNNING - 45.00 %
2019-12-04 16:38:27 CST: SUCCEEDED
```

----结束

批量模式示例

步骤1 使用安装客户端的用户登录Loader客户端所在节点。

步骤2 执行以下命令，进入sqoop-shell工具的“conf”目录。例如，Loader客户端安装目录为“/opt/client/Loader/”。

```
cd /opt/client/Loader/loader-tools-1.99.3/sqoop-shell/conf
```

步骤3 执行以下命令，配置认证信息。

```
vi client.properties
server.url=10.0.0.1:21351
# simple or kerberos
authentication.type=simple
# true or false
use.keytab=true

authentication.user=
authentication.password=

client.principal=hdfs/hadoop.<系统域名>@<系统域名>

# keytab file
client.keytab.file=./conf/login/hdfs.keytab
```

表 17-147 配置参数说明

| 配置参数 | 说明 | 示例 |
|------------|---|----------------|
| server.url | Loader服务的浮动IP地址和端口（21351）。
为了兼容性，此处支持配置多个IP地址和端口，并以“,”进行分隔。其中第一个必须是Loader服务的浮动IP地址和端口（21351），其余的可根据业务需求配置。 | 10.0.0.1:21351 |

| 配置参数 | 说明 | 示例 |
|-------------------------|---|---------------|
| authentication.type | <p>登录认证的方式。</p> <ul style="list-style-type: none"> “kerberos”，表示使用安全模式，进行Kerberos认证。Kerberos认证提供两种认证方式：密码和keytab文件。 “simple”，表示使用普通模式，不进行Kerberos认证。 | kerberos |
| authentication.user | <p>普通模式或者使用密码认证方式时，登录使用的用户。</p> <p>keytab登录方式，则不需要设置该参数。</p> | bar |
| authentication.password | <p>使用密码认证方式时，登录使用的用户密码。</p> <p>普通模式或者keytab登录方式，则不需要设置该参数。</p> <p>用户需要对密码加密，加密方法：</p> <ol style="list-style-type: none"> 进入“encrypt_tool”所在目录。例如，Loader客户端安装目录为“/opt/hadoopclient/Loader”，则执行如下命令。
cd /opt/hadoopclient/Loader/loader-tools-1.99.3 执行以下命令，对非加密密码进行加密。
./encrypt_tool 未加密的密码
得到加密后的密文，作为“authentication.password”的取值。 <p>说明
非加密密码中含有特殊字符时需要转义。例如，\$符号属于特殊字符，可使用单引号进行转义-；非加密密码中含有单引号时可用双引号进行转义，非加密密码中含有双引号应使用反斜杠\进行转义。可参考Shell的转义字符规则。</p> | - |
| use.keytab | <p>是否使用keytab方式登录。</p> <ul style="list-style-type: none"> true，表示使用keytab文件登录。 false，表示使用密码登录。 | true |
| client.principal | <p>使用keytab认证方式时，访问Loader服务的用户规则。</p> <p>普通模式或者密码登录方式，则不需要设置该参数。</p> | loader/hadoop |

| 配置参数 | 说明 | 示例 |
|--------------------|---|--------------------------------|
| client.keytab.file | 使用keytab认证方式登录时，使用的keytab文件所在目录。
普通模式或者密码登录方式，则不需要设置该参数。 | /opt/client/conf/loader.keytab |

步骤4 执行以下命令，进入“sqoop2-shell”脚本所在目录，并在该目录下创建一个文本文件，例如“batchCommand.sh”。

```
cd /opt/client/Loader/loader-tools-1.99.3/sqoop-shell
```

```
vi batchCommand.sh
```

“batchCommand.sh” 样例如下：

```
//查看参数
create connection -c 4 --help

//创建连接器
create connection -c 4 -name oracle-connection --connector-connection-oracleServerIp 10.0.0.1 --connector-connection-oracleServerPort 22 --connector-connection-oracleUser root --connector-connection-oraclePassword xxxxx

//创建作业
create job -t import -x 3 --connector-file-inputPath /opt/tempfile --connector-file-fileFilter * --framework-output-outputDirectory /user/loader/1 --framework-output-storageType HBase --framework-throttling-extractorSize 120 --framework-output-fileType TEXT_FILE --connector-file-splitType FILE -name test

//启动作业
start job -j 7 -s
```

其中xxxxxx为连接器密码。

步骤5 执行如下命令，sqoop-shell工具将依次执行上述命令。

```
./sqoop2-shell batchCommand.sh
```

也可以直接在命令里附带认证信息。

使用密码认证：

```
./sqoop2-shell -uk false -u username -p encryptedPassword batchCommand.sh
```

使用Kerberos认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal batchCommand.sh
```

显示“SUCCEEDED”信息，则说明作业启动成功。

```
Welcome to sqoop client
Use the username and password authentication mode
Authentication success.
sqoop:000> create connection -c 4 --help
usage: Show connection viparameters:
--connector-connection-oraclePassword <arg>
--connector-connection-oracleServerIp <arg>
--connector-connection-oracleServerPort <arg>
--connector-connection-oracleUser <arg>
--framework-security-maxConnections <arg>
--name <arg>
====> FINE
sqoop:000> create connection -c 4 -name oracle-connection --connector-connection-oracleServerIp 10.0.0.1 --connector-connection-oracleServerPort 22 --connector-connection-
```

```
oraclePassword xxxxx
Creating connection for connector with id 4
New connection was successfully created with validation status FINE and persistent id 3
====> FINE
sqoop:000> create job -t import -x 3 --connector-file-inputPath /opt/tempfile --connector-file-fileFilter * --
framework-output-outputDirectory /user/loader/1 --framework-output-storageType HDFS --framework-
throttling-extractorSize 120 --framework-output-fileType TEXT_FILE --connector-file-splitType FILE -name
test
Creating job for connection with id 3
New job was successfully created with validation status FINE and persistent id 7
====> FINE
Submission details
Job ID: 7
Server URL: https://10.0.0.0:21351/loader/
Created by: admintest
Creation date: 2019-12-04 16:37:34 CST
Lastly updated by: admintest
2019-12-04 16:37:34 CST: BOOTING - Progress is not available
2019-12-04 16:37:42 CST: BOOTING - 0.00 %
2019-12-04 16:37:42 CST: BOOTING - 0.00 %
2019-12-04 16:37:57 CST: RUNNING - 0.00 %
2019-12-04 16:38:12 CST: RUNNING - 45.00 %
2019-12-04 16:38:12 CST: RUNNING - 45.00 %
2019-12-04 16:38:27 CST: SUCCEEDED
```

步骤6 批处理模式下，使用-c参数附带一条命令，sqoop-shell可以一次只执行附带的这一条命令。

执行如下命令将创建连接器。

```
./sqoop2-shell -c "create connection -c 4 -name oracle-connection --
connector-connection-oracleServerIp 10.0.0.1 --connector-connection-
oracleServerPort 22 --connector-connection-oracleUser root --connector-
connection-oraclePassword xxxxx"
```

可以在命令里直接附带认证信息。

使用密码认证：

```
./sqoop2-shell -uk false -u username -p encryptedPassword -c "create
connection -c 4 -name oracle-connection --connector-connection-
oracleServerIp 10.0.0.1 --connector-connection-oracleServerPort 22 --
connector-connection-oracleUser root --connector-connection-oraclePassword
xxxxx"
```

使用Kerberos认证：

```
./sqoop2-shell -uk true -k user.keytab -s userPrincipal -c "create connection -c 4
-name oracle-connection --connector-connection-oracleServerIp 10.0.0.1 --
connector-connection-oracleServerPort 22 --connector-connection-oracleUser
root --connector-connection-oraclePassword xxxxx"
```

显示“FINE”信息，则说明连接创建成功。

```
Welcome to sqoop client
Use the username and password authentication mode
Authentication success.
sqoop:000> create connection -c 4 -name oracle-connection --connector-connection-oracleServerIp 10.0.0.1
--connector-connection-oracleServerPort 22 --connector-connection-oracleUser root --connector-connection-
oraclePassword xxxxx
Creating connection for connector with id 4
New connection was successfully created with validation status FINE and persistent id 3
====> FINE
```

----**结束**

17.17 Loader 日志介绍

日志描述

日志存储路径：Loader相关日志的默认存储路径为“/var/log/Bigdata/loader/日志分类”。

- runlog：“/var/log/Bigdata/loader/runlog”（运行日志）
- scriptlog：“/var/log/Bigdata/loader/scriptlog/”（脚本的执行日志）
- catalina：“/var/log/Bigdata/loader/catalina”（tomcat的启停日志）
- audit：“/var/log/Bigdata/loader/audit”（审计日志）

日志归档规则：

Loader的运行日志和审计日志，启动了自动压缩归档功能，默认情况下，当日志大小超过10MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd_hh-mm-ss>.[编号].log.zip”。最多保留最近的20个压缩文件，压缩文件保留个数可以在Manager界面中配置。

表 17-148 Loader 日志列表

| 日志类型 | 日志文件名 | 描述 |
|----------|--|-------------------------------------|
| 运行日志 | loader.log | Loader运行日志，记录Loader系统运行时候所产生的大部分日志。 |
| | loader-omm-***-pid***-gc.log.*.current | Loader进程gc日志 |
| | sqoopInstanceCheck.log | Loader实例健康检查日志 |
| 审计日志 | default.audit | Loader操作审计日志（例如：作业的增删改查、用户的登录）。 |
| tomcat日志 | catalina.out | tomcat的运行日志 |
| | catalina.<yyyy-mm-dd>.log | tomcat的运行日志 |
| | host-manager.<yyyy-mm-dd>.log | tomcat的运行日志 |
| | localhost_access_log.<yyyy-mm-dd>.txt | tomcat的运行日志 |
| | manager.<yyyy-mm-dd>.log | tomcat的运行日志 |
| | localhost.<yyyy-mm-dd>.log | tomcat的运行日志 |

| 日志类型 | 日志文件名 | 描述 |
|------|-----------------|---|
| 脚本日志 | postInstall.log | Loader安装脚本日志。
执行loader安装脚本（postInstall.sh）时产生的日志。 |
| | preStart.log | Loader服务的预启动脚本日志。Loader服务启动时，需要先执行一系列的准备操作（preStart.sh），例如生成keytab文件等，该日志正是记录了这些操作信息。 |
| | loader_ctl.log | Loader执行服务启停脚本（sqoop.sh）的日志。 |

日志级别

Loader中提供了如表17-149所示的日志级别，日志级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 17-149 日志级别

| 级别 | 描述 |
|-------|-----------------------|
| ERROR | ERROR表示错误日志输出。 |
| WARN | WARN表示当前事件处理存在异常信息。 |
| INFO | INFO表示系统及各事件正常运行状态信息。 |
| DEBUG | DEBUG表示系统及系统调试信息。 |

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 请参考[修改集群服务配置参数](#)，进入Loader的“全部配置”页面。
- 步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤3** 选择所需修改的日志级别。
- 步骤4** 保存配置，在弹出窗口中单击“确定”，完成后重启服务使配置生效。

----结束

日志格式

Loader的日志格式如下所示：

表 17-150 日志格式

| 日志类型 | 格式 | 示例 |
|------|--|---|
| 运行日志 | <yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的 message> <日志事件的发生位置> | 2015-06-29 14:54:35,553 INFO [localhost-startStop-1] ConnectionRequestHandler initialized org.apache.sqoop.handler.ConnectionRequestHandler.<init>(ConnectionRequestHandler.java:100) |
| 审计日志 | <yyyy-MM-dd HH:mm:ss,SSS> <Log Level> default <log中的 message> <日志事件的发生位置> | 2015-06-29 15:35:40,969 INFO default: UserName=admin, UserIP=10.52.0.111, Time=2015-06-29 15:35:40,969, Operation=submit, Resource=submission@21, Result=Failure, Detail={[reason:GET_SFTP_SESSION_FAILED:Failed to get sftp session - 10.162.0.35 (caused by: Auth cancel)]; [config:null]} |

17.18 样例：通过 Loader 将数据从 OBS 导入 HDFS

操作场景

用户需要将大量数据从集群外导入集群内的时候，可以选择从OBS导入到HDFS的方式。

前提条件

- 已准备业务数据。
- 已创建分析集群。

操作步骤

步骤1 将业务数据上传到用户的OBS文件系统。

步骤2 获取用户的AK/SK信息，然后创建一个OBS连接和一个HDFS连接。

具体可参见[Loader连接配置说明](#)。

步骤3 访问Loader页面。

如果是启用了Kerberos认证的分析集群，可参见[访问Hue WebUI界面](#)。

步骤4 单击“新建作业”。

步骤5 在“基本信息”填写参数。

1. 在“名称”填写一个作业的名称。例如“obs2hdfs”。
2. 在“源连接”选择已创建的OBS连接。
3. “目的连接”选择已创建的HDFS连接。

步骤6 在“自”填写源连接参数。

1. 在“桶名”填写业务数据所保存的OBS文件系统名称。
2. 在“源目录或文件”填写业务数据在文件系统的具体位置。
如果是单个文件，需要填写包含文件名的完整路径。如果是目录，填写目录的完整路径
3. “文件格式”填写业务数据文件的类型。

可参见[obs-connector](#)。

步骤7 在“至”填写目的连接参数。

1. 在“定入目录”填写业务数据在HDFS要保存的目录名称。
如果是启用Kerberos认证的集群，当前访问Loader的用户对保存数据的目录需要有写入权限。
2. 在“文件格式”填写业务数据文件的类型。
需要与[步骤6.3](#)的类型对应。
3. 在“压缩格式”填写一种压缩的算法。例如选择不压缩“NONE”。
4. 在“是否覆盖”选择已有文件的处理方式，选择“True”。
5. 单击“显示高级属性”，在“换行符”填写业务数据保存时，系统填充的换行字符。
6. 在“字段分割符”填写业务数据保存时，系统填充的分割字符。

可参见[hdfs-connector](#)。

步骤8 在“任务配置”填写作业的运行参数。

1. 在“抽取并发数”填写map任务的个数。
2. 在“加载(写入)并发数”填写reduce任务的个数。
目的连接为HDFS连接时，不显示“加载(写入)并发数”参数。
3. “单个分片的最大错误记录数”填写错误记录阈值。
4. 在“脏数据目录”填写一个脏数据的保存位置，例如“/user/sqoop/obs2hdfs-dd”。

步骤9 单击“保存并运行”。

在“管理作业界面”，查看作业运行结果。可以单击“刷新列表”获取作业的最新状态。

----结束

17.19 Loader 常见问题

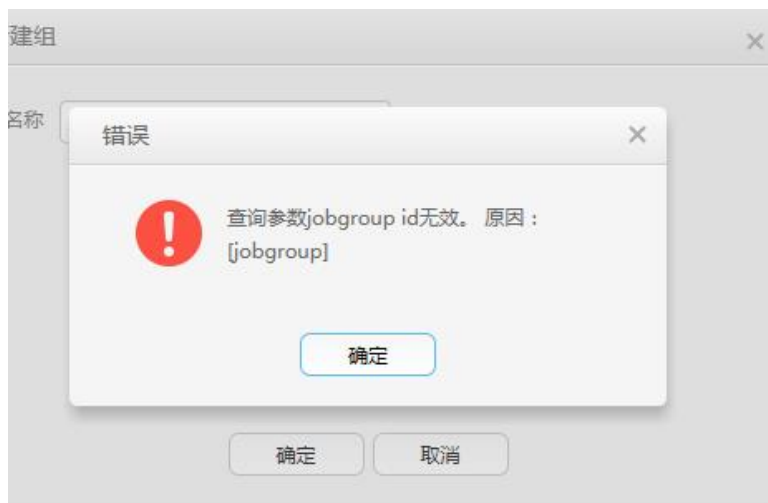
17.19.1 IE 10&IE 11 浏览器无法保存数据

问题

通过IE 10&IE 11浏览器访问Loader界面，提交数据后，会报错。

回答

- 现象
保存提交数据，出现类似报错：Invalid query parameter jobgroup id. cause: [jobgroup]。



- 原因
IE 11浏览器的某些版本在接收到HTTP 307响应时，会将POST请求转化为GET请求，从而使得POST数据无法下发到服务端。
- 解决建议
使用Google Chrome浏览器。

17.19.2 将 Oracle 数据库中的数据导入 HDFS 时各连接器的区别

问题

使用Loader将Oracle数据库中的数据导入到HDFS中时，可选择的连接器有generic-jdbc-connector、oracle-connector、oracle-partition-connector三种，要怎么选？有什么区别？

答案

- generic-jdbc-connector
使用JDBC方式从Oracle数据库读取数据，适用于支持JDBC的数据库。
在这种方式下，Loader加载数据的性能受限于分区列的数据分布是否均匀。当分区列的数据偏斜（数据集中在一个或者几个值）时，个别Map需要处理绝大部分数据，进而导致索引失效，造成SQL查询性能急剧下降。
generic-jdbc-connector支持视图的导入导出，而oracle-partition-connector和oracle-connector暂不支持，因此导入视图只能选择该连接器。

- oracle-partition-connector和oracle-connector

这两种连接器都支持按照Oracle的ROWID进行分区（oracle-partition-connector是自研，oracle-connector是社区开源版本），二者的性能较为接近。

oracle-connector需要的系统表权限较多，下面是各自需要的系统表，需要赋予读权限。

- oracle-connector: dba_tab_partitions、dba_constraints、dba_tables、dba_segments、v\$instance、dba_objects、v\$instance、SYS_CONTEXT函数、dba_extents、dba_tab_subpartitions。
- oracle-partition-connector: DBA_OBJECTS、DBA_EXTENTS。

相比于generic-jdbc-connector，oracle-partition-connector和oracle-connector具有以下优点：

- a. 负载均衡，数据分片的个数和范围与源表的数据无关，而是由源表的存储结构（数据块）确定，颗粒度可以达到“每个数据块一个分区”。
- b. 性能稳定，完全消除“数据偏斜”和“绑定变量窥探”导致的“索引失效”。
- c. 查询速度快，数据分片的查询速度比用索引快。
- d. 水平扩展性好，如果数据量越大，产生的分片就越多，所以只要增加任务的并发数，就可以获得较理想的性能；反之，减少任务并发数，就可以节省资源。
- e. 简化数据分片逻辑，不需要考虑“精度丢失”、“类型兼容”和“绑定变量”等问题。
- f. 易用性得到增强，用户不需要专门为Loader创建分区列、分区表。

18 使用 Kudu

18.1 从零开始使用 Kudu

Kudu是专为Apache Hadoop平台开发的列式存储管理器。Kudu具有Hadoop生态系统应用程序的共同技术特性：可水平扩展，并支持高可用性操作。

前提条件

已安装集群客户端，例如安装目录为“/opt/hadoopclient”，以下操作的客户端目录只是举例，请根据实际安装目录修改。

操作步骤

步骤1 以客户端安装用户，登录安装客户端的节点。

执行**su - omm**命令，切换到omm用户。

步骤2 执行以下命令，切换到客户端安装目录。

```
cd /opt/hadoopclient
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 运行Kudu命令行工具。

直接执行Kudu组件的命令行工具，查看帮助。

```
kudu -h
```

回显信息如下：

```
Usage: kudu <command> [<args>]

<command> can be one of the following:
  cluster  Operate on a Kudu cluster
  diagnose Diagnostic tools for Kudu servers and clusters
  fs       Operate on a local Kudu filesystem
  hms     Operate on remote Hive Metastores
  local_replica Operate on local tablet replicas via the local filesystem
  master  Operate on a Kudu Master
```

```
pbcc  Operate on PBC (protobuf container) files
perf  Measure the performance of a Kudu cluster
remote_replica  Operate on remote tablet replicas on a Kudu Tablet Server
table  Operate on Kudu tables
tablet  Operate on remote Kudu tablets
test  Various test actions
tserver  Operate on a Kudu Tablet Server
wal  Operate on WAL (write-ahead log) files
```

📖 说明

kudu命令行工具不提供DDL、DML等操作，但提供针对cluster、master、tserver、fs、table等的细化查询功能。

常用操作：

- 查看当前集群下有哪些表。

```
kudu table list KuduMaster实例IP1:7051, KuduMaster实例IP2:7051, KuduMaster实例IP3:7051
```

- 查询Kudu服务KuduMaster实例的配置信息。

```
kudu master get_flags KuduMaster实例IP:7051
```

- 查询表的schema。

```
kudu table describe KuduMaster实例IP1:7051, KuduMaster实例IP2:7051, KuduMaster实例IP3:7051 tablename
```

- 删除表。

```
kudu table delete KuduMaster实例IP1:7051, KuduMaster实例IP2:7051, KuduMaster实例IP3:7051 tablename
```

📖 说明

KuduMaster实例IP获取方式：在集群详情页面，选择“组件管理 > Kudu > 实例”，获取角色KuduMaster的IP地址。

----结束

18.2 访问 Kudu 的 WebUI

用户可以通过Kudu的WebUI，在图形化界面查看Kudu作业的相关信息。

前提条件

已安装Kudu服务的集群。

访问 KuduMaster WebUI（MRS 3.x 及之后版本）

步骤1 登录Manager页面，请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。

步骤2 选择“集群 > 服务 > Kudu”。

步骤3 在“Kudu 概览”的“KuduMaster WebUI”中单击“KuduMaster(KuduMaster)”，打开KuduMaster的WebUI页面。

图 18-1 KuduMaster WebUI

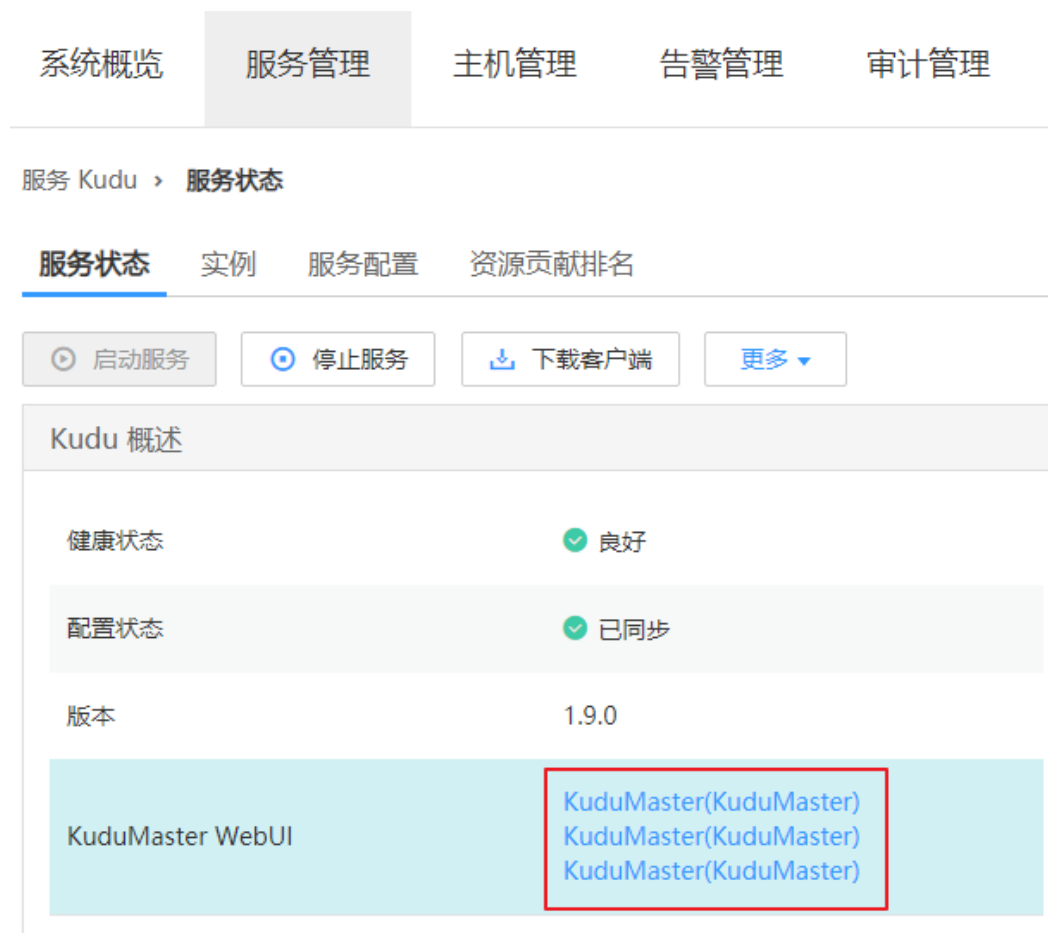


----结束

访问 KuduMaster WebUI（MRS 3.x 之前版本）

- 步骤1** 登录Manager页面，请参见[访问MRS Manager（MRS 3.x之前版本）](#)。
- 步骤2** 选择“服务管理 > Kudu”。
- 步骤3** 在“Kudu 概述”的“KuduMaster WebUI”中单击“KuduMaster(KuduMaster)”，打开KuduMaster的WebUI页面。

图 18-2 KuduMaster WebUI



---结束

19 使用 MapReduce

19.1 配置使用分布式缓存执行 MapReduce 任务

配置场景

📖 说明

本章节操作适用于MRS 3.x及之后版本。

分布式缓存在两种情况下非常有用。

- **滚动升级**

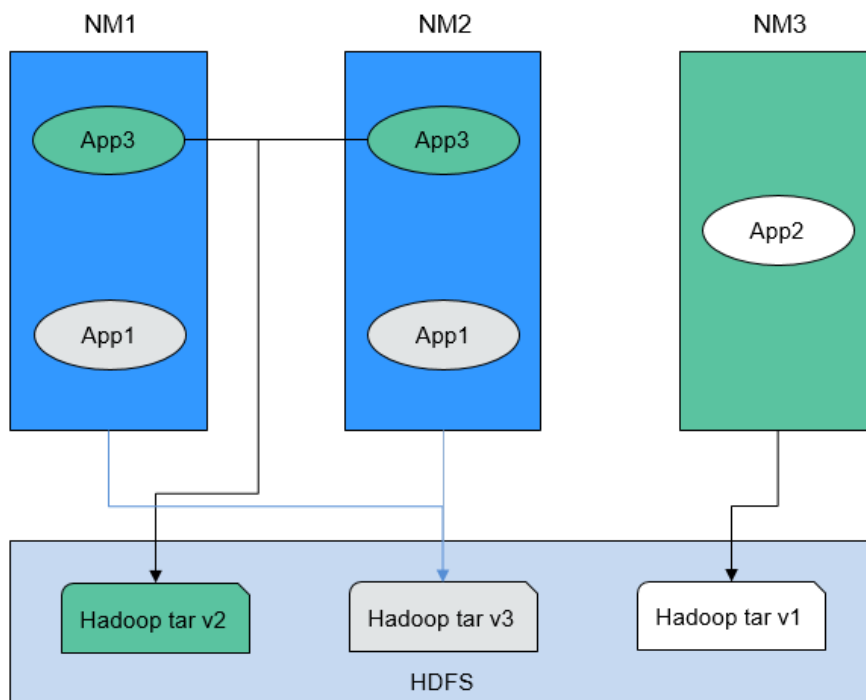
在升级过程中，应用程序必须保持文字内容（jar文件或配置文件）不变。而这些内容并非基于当前版本的Yarn，而是要基于其提交时的版本。一般情况下，应用程序（例如MapReduce、Hive、Tez等）需要进行完整的本地安装，将库安装至所有的集群机器（客户端及服务器端机器）中。当集群内开始进行滚动升级或降级时，本地安装的库的版本必然会在应用运行过程时发生改变。在滚动升级过程中，首先只会对少数NodeManager进行升级，这些NodeManager会获得新版本的软件。这导致了行为的不一致，并可能发生运行时错误。

- **同时存在多个Yarn版本**

集群管理员可能会在一个集群内运行使用多个版本Yarn及Hadoop jars的任务。这在当前很难实现，因为jars已被本地化且只有一个版本。

MapReduce应用框架可以通过分布式缓存进行部署，且无需依赖安装中复制的静态版本。因此，可以在HDFS中存放多版本的Hadoop，并通过配置“mapred-site.xml”文件指定任务默认使用的版本。只需设置适当的配置属性，用户就可以运行不同版本的MapReduce，而无需使用部署在集群中的版本。

图 19-1 具有多个版本 NodeManagers 及 Applications 的集群



在图19-1中：可以看出，应用程序可以使用HDFS中的Hadoop jars，而无需使用本地版本。因此在滚动升级中，即使NodeManager已经升级，应用程序仍然可以运行旧版本的Hadoop。

配置描述

步骤1 进入HDFS客户端。

1. 以客户端安装用户，登录安装客户端的节点。
2. 执行以下命令，切换到客户端安装目录。

```
cd 客户端安装路径
```

3. 执行以下命令配置环境变量。

```
source bigdata_env
```

4. 如果集群为安全模式，执行以下命令进行用户认证。普通模式集群无需执行用户认证。

```
kinit 组件业务用户
```

步骤2 将指定版本的MapReduce tar包存放至HDFS中应用程序可以访问的目录下，如下所示：

```
$HADOOP_HOME/bin/hdfs dfs -put hadoop-x.tar.gz /mapred/framework/
```

步骤3 执行以下命令，根据表19-1，对“客户端安装路径/Yarn/config/mapred-site.xml”文件中的参数进行设置。

```
vi 客户端安装路径/Yarn/config/mapred-site.xml
```

表 19-1 分布式缓存相关参数

| 参数 | 说明 | 默认值 |
|--------------------------------------|---|-----|
| mapreduce.application.framework.path | <p>此参数值为指向存档位置的URL。</p> <p>说明
如果对URL片段标示名称进行如下指定，该属性还可以为存档创建别名。作为示例，这里将别名设为了mr-framework。
<property> <name>mapreduce.application.framework.path</name>
<value>hdfs:/mapred/framework/hadoop-x.tar.gz#mr-framework</value> </property></p> | NA |
| mapreduce.application.classpath | <p>设定属性mapreduce.application.classpath，使其可以包含类目录中相关的MR jars。</p> <p>说明
例如，此处利用在框架路径中使用过的别名“mr-framework”对目录进行匹配。
<property> <name>mapreduce.application.classpath</name>
<value>\${PWD}/mr-framework/hadoop/share/hadoop/mapreduce/*:\${PWD}/mr-framework/hadoop/share/hadoop/mapreduce/lib/*:\${PWD}/mr-framework/hadoop/share/hadoop/common/*:\${PWD}/mr-framework/hadoop/share/hadoop/common/lib/*:\${PWD}/mr-framework/hadoop/share/hadoop/yarn/*:\${PWD}/mr-framework/hadoop/share/hadoop/yarn/lib/*:\${PWD}/mr-framework/hadoop/share/hadoop/hdfs/*:\${PWD}/mr-framework/hadoop/share/hadoop/hdfs/lib/*:/etc/hadoop/conf/secure</value></property></p> | NA |

可以将多个版本的MapReduce tar包上传至HDFS。不同的“mapred-site.xml”文件可以指向不同的位置。用户在此之后可以针对特定的“mapred-site.xml”文件运行任务。以下是一个针对x版本的MapReduce tar包运行MapReduce任务的样例：

```
hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-*.jar pi -conf etc/hadoop-x/mapred-site.xml 10 10
```

----结束

19.2 配置 MapReduce shuffle address

配置场景

当MapReduce shuffle服务启动时，它尝试基于localhost绑定IP。如果需要MapReduce shuffle服务连接特定IP，可以参考该章节进行配置。

配置描述

当需要MapReduce shuffle服务绑定特定IP时，需要在NodeManager实例所在节点的配置文件“mapred-site.xml”中（例如路径为：\${BIGDATA_HOME}/FusionInsight_HD_xxx/x_xx_NodeManager/etc/mapred-site.xml）设置如下参数。

表 19-2 参数描述

| 参数 | 描述 | 默认值 |
|---------------------------|---|-----|
| mapreduce.shuffle.address | 指定地址来运行shuffle服务，格式是IP:PORT，参数的默认值为空。当参数值为空时，将绑定localhost，默认端口为13562。
说明
如果涉及到的PORT值和配置的mapreduce.shuffle.port值不一样时，mapreduce.shuffle.port将不会生效。 | - |

19.3 配置 MapReduce 集群管理员列表

配置场景

该功能主要用于指定MapReduce集群管理员。

其中，集群管理员列表由参数“mapreduce.cluster.administrators”指定，集群管理员admin具有所有可以操作的权限。

配置描述

进入Mapreduce服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。

表 19-3 参数描述

| 参数 | 描述 | 默认值 |
|----------------------------------|---|--|
| mapreduce.cluster.acls.enabled | 是否开启对Job History Server权限控制的开关。 | true |
| mapreduce.cluster.administrators | 用于指定MapReduce集群管理员列表，可以配置用户和用户组，用户或者用户组之间用逗号间隔，用户和用户组之间用空格间隔，举例：userA,userB groupA,groupB。当配置为*时表示所有用户或用户组。 | MRS 3.x之前版本：mapred
MRS 3.x及之后版本：
mapred
supergroup,System_administrator_186 |

19.4 通过 Windows 系统提交 MapReduce 任务

配置场景

用户将MapReduce任务从Windows上提交到Linux上运行，则“mapreduce.app-submission.cross-platform”参数值需配置为“true”。若集群无此参数，或参数值为

“false”，则表示集群不支持此功能，需要按照如下操作增加该参数或修改参数值进行开启。

📖 说明

本章节操作适用于MRS 3.x及之后版本。

配置描述

在客户端的“mapred-site.xml”配置文件中进行如下配置。“mapred-site.xml”配置文件在客户端安装路径的config目录下，例如“/opt/client/Yarn/config”。

表 19-4 参数说明

| 参数 | 描述 | 默认值 |
|---|---|------|
| mapreduce.ap
p-
submission.cro
ss-platform | 支持在Windows上提交到Linux上运行MR任务的配置项。当该参数的值设为“true”时，表示支持。当该参数的值设为“false”时，表示不支持。 | true |

19.5 配置 MapReduce 任务日志归档和清理机制

配置场景

执行一个MapReduce应用会产生两种类型日志文件：作业日志和任务日志。

- 作业日志由MRApplicationMaster产生，详细记录了作业启动时间、运行时间，每个任务启动时间、运行时间、Counter值等信息。此日志内容被HistoryServer解析以后用于查看作业执行的详细信息。
- 任务日志记录了每个运行在Container中的任务输出的日志信息。默认情况下，任务日志只会存放在各NodeManager的本地磁盘上。打开日志聚合功能后，NodeManager会在作业运行完成后将本地的任务日志进行合并，写入到HDFS中。

由于MapReduce的作业日志和任务日志（聚合功能开启的情况下）都保存在HDFS上。对于计算任务量大的集群，如果不进行合理的配置对日志文件进行定期归档和删除，日志文件将占用HDFS大量内存空间，增加集群负载。

日志归档是通过Hadoop Archives功能实现的，Hadoop Archives启动的并行归档任务数（Map数）与待归档的日志文件总大小有关。计算公式为：并行归档任务数=待归档的日志文件总大小/归档文件大小。

配置描述

进入Mapreduce服务参数“全部配置”界面，具体操作请参考[修改集群服务配置参数](#)章节。

在搜索框中输入参数名称，修改并保存配置。然后在Mapreduce服务“概览”页面选择“更多 > 同步配置”。同步完成后重启Mapreduce服务。

- 作业日志参数：

表 19-5 参数说明

| 参数 | 描述 | 默认值 |
|--|--|-----------------|
| mapreduce.jobhistory.cleaner.enable | 是否开启作业日志文件清理功能。 | true |
| mapreduce.jobhistory.cleaner.interval-ms | 作业日志文件清理启动周期。只有保留时间比“mapreduce.jobhistory.max-age-ms”更长的日志文件才会被清除。 | 86400000（1天） |
| mapreduce.jobhistory.max-age-ms | 比此项设置的毫秒数保留时间更长的作业日志文件将被清理。 | 1296000000（15天） |

- 任务日志参数：

表 19-6 参数说明

| 参数 | 描述 | 默认值 |
|---|---|---------|
| yarn.log-aggregation.archive.files.minimum | 设置Mapreduce任务日志归档最小文件数。当“yarn.nodemanager.remote-app-log-dir”文件夹下文件数大于等于该设置的值时归档任务启动。
该参数适用于MRS 3.x版本。 | 5000 |
| yarn.log-aggregation.archive-check-interval-seconds | 设置Mapreduce任务日志归档任务启动周期（秒）。只有日志文件数达到“yarn.log-aggregation.archive.files.minimum”设置值时日志文件才会被归档。周期设置为“0”或“-1”时归档功能禁用。
该参数适用于MRS 3.x版本。 | -1 |
| yarn.log-aggregation.retain-seconds | 设置Mapreduce任务日志在HDFS上的保留时间。设置为“-1”时日志文件永久保存。 | 1296000 |
| yarn.log-aggregation.retain-check-interval-seconds | 设置Mapreduce任务日志清理任务的检查周期（秒）。设置为“-1”时检查周期为日志保留时间的十分之一。 | 86400 |

📖 说明

如果是任务日志将HDFS存储空间占用太多，主要修改“mapreduce.jobhistory.max-age-ms”和“yarn.log-aggregation.retain-check-interval-seconds”配置项来控制任务日志保存时间。

19.6 MapReduce 性能调优

19.6.1 多 CPU 内核下的 MapReduce 调优配置

操作场景

当CPU内核数很多时，如CPU内核为磁盘数的3倍时的调优配置。

操作步骤

以下参数有如下两个配置入口：

- 服务器端配置
进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。
- 客户端配置
直接在客户端中修改相应的配置文件。

说明

- HDFS客户端配置文件路径：*客户端安装目录*/HDFS/hadoop/etc/hadoop/hdfs-site.xml。
- Yarn客户端配置文件路径：*客户端安装目录*/HDFS/hadoop/etc/hadoop/yarn-site.xml。
- MapReduce客户端配置文件路径：*客户端安装目录*/HDFS/hadoop/etc/hadoop/mapred-site.xml。

表 19-7 多 CPU 内核设置

| 配置 | 参数 | 配置描述 |
|---------|--|---|
| 节点容器槽位数 | <p>yarn.nodemanager.resource.memory-mb</p> <ul style="list-style-type: none"> 参数解释：节点上YARN可使用的物理内存总量。单位：M。 默认值：
MRS 3.x之前版本：
8192
MRS 3.x及之后版本：
16384 参数入口：
MRS 3.x之前版本：需要在MRS控制台上进行配置。
MRS 3.x及之后版本：需要在 FusionInsight Manager系统进行配置。 | <ul style="list-style-type: none"> 参数配置组合决定了每节点任务 (map、reduce)的并发数。 如果所有的任务 (map/reduce) 需要读写数据至磁盘，多个进程将会同时访问一个磁盘。这将会导致磁盘的IO性能非常低下。为了改善磁盘的性能，请确保客户端并发访问磁盘的数不大于3。 最大并发的 container数量应该为$[2.5 * \text{Hadoop中磁盘配置数}]$。 |
| | <p>mapreduce.map.memory.mb</p> <ul style="list-style-type: none"> 参数解释：map任务的内存限制。单位：MB。 默认值：4096 参数入口：需要在客户端进行配置，配置文件路径：<i>客户端安装目录</i>/HDFS/hadoop/etc/hadoop/mapred-site.xml。 | |
| | <p>mapreduce.reduce.memory.mb</p> <ul style="list-style-type: none"> 参数解释：Reduce任务的内存限制。单位：MB。 默认值：4096 参数入口：需要在客户端进行配置，配置文件路径：<i>客户端安装目录</i>/HDFS/hadoop/etc/hadoop/mapred-site.xml。 | |

| 配置 | 参数 | 配置描述 |
|----------|--|---|
| Map输出与压缩 | <p>mapreduce.map.output.compress</p> <ul style="list-style-type: none"> 参数解释：指定了Map任务输出结果可以在网络传输前被压缩。这是一个per-job的配置。 默认值：true 参数入口：需要在客户端进行配置，配置文件路径：<i>客户端安装目录</i>/HDFS/hadoop/etc/hadoop/mapred-site.xml。 | <ul style="list-style-type: none"> Map任务所产生的输出可以在写入磁盘之前被压缩，这样可以节约磁盘空间并得到更快的写盘速度，同时可以减少至Reducer的数据传输量。需要在客户端进行配置。 在这种情况下，磁盘的IO是主要瓶颈。所以可以选择一种压缩率非常高的压缩算法。 编解码器可配置为Snappy，Benchmark测试结果显示Snappy是非常平衡以及高效的编码器。 |
| | <p>mapreduce.map.output.compress.codec</p> <ul style="list-style-type: none"> 参数解释：指定用于压缩的编解码器。 默认值：
org.apache.hadoop.io.compress.Lz4Codec 参数入口：需要在客户端进行配置，配置文件路径：<i>客户端安装目录</i>/HDFS/hadoop/etc/hadoop/mapred-site.xml。 | |
| Spills | <p>mapreduce.map.sort.spill.percent</p> <ul style="list-style-type: none"> 参数解释：序列化缓冲区中的软限制。一旦达到该限制，线程将在后台开始将内容溢出到磁盘。 默认值：0.8 参数入口：需要在客户端进行配置，配置文件路径：<i>客户端安装目录</i>/HDFS/hadoop/etc/hadoop/mapred-site.xml。 | <p>磁盘IO是主要瓶颈，合理配置“mapreduce.task.io.sort.mb”可以使溢出至磁盘的内容最小化。</p> |

| 配置 | 参数 | 配置描述 |
|-------|---|--|
| 数据包大小 | <p>dfs.client-write-packet-size</p> <ul style="list-style-type: none"> 参数解释：配置项可以指定该数据包的大小。可以通过每个job进行指定。 默认值：262144 参数入口：需要在客户端进行配置，配置文件路径：<i>客户端安装目录</i>/HDFS/hadoop/etc/hadoop/hdfs-site.xml。 | <ul style="list-style-type: none"> 当HDFS客户端写数据至数据节点时，数据会被累积，直到形成一个包。这个数据包会通过网络传输。 数据节点从HDFS客户端接收数据包，然后将数据包里的数据单线程写入磁盘。当磁盘处于并发写入状态时，增加数据包的大小可以减少磁盘寻道时间，从而提升IO性能。 dfs.client-write-packet-size = 262144 |

19.6.2 配置 MapReduce Job 基线

操作场景

确定Job基线是调优的基础，一切调优项效果的检查，都是通过和基线数据做对比来获得。

Job基线的确定有如下三个原则：

- 充分利用集群资源
- Reduce阶段尽量放在一轮
- 每个Task的执行时间要合理

操作步骤

- 原则一：充分利用集群资源。**

Job运行时，会让所有的节点都有任务处理，且处于繁忙状态，这样才能保证资源充分利用，任务的并发度达到最大。可以通过调整处理的数据量大小，以及调整map和reduce个数来实现。

reduce个数的控制使用“mapreduce.job.reduces”。

map个数取决于使用了哪种InputFormat，以及待处理的数据文件是否可分割。默认的TextFileInputFormat将根据block的个数来分配map数(一个block一个map)。通过如下配置参数进行调整。

参数入口：

进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。

表 19-8 参数配置-1

| 参数 | 描述 | 默认值 |
|---|---|-----|
| mapreduce.input.fileinputformat.split.maxsize | map输入信息应被拆分成的数据块的最大大小。
由用户定义的分片大小的设置及每个文件block大小的设置，可以计算分片的大小。计算公式如下：
splitSize = Math.max(minSize, Math.min(maxSize, blockSize))
如果maxSize设置大于blockSize，那么每个block就是一个分片，否则就会将一个block文件分隔为多个分片，如果block中剩下的一小段数据量小于splitSize，还是认为它是独立的分片。 | - |
| mapreduce.input.fileinputformat.split.minsize | 可以设置数据分片的数据最小值。 | 0 |

- **原则二：控制reduce阶段在一轮中完成。**

避免以下两种场景：

- 大部分的reduce在第一轮运行完后，剩下唯一一个reduce继续运行。这种情况下，这个reduce的执行时间将极大影响这个job的运行时间。因此需要将reduce个数减少。
- 所有的map运行完后，只有个别节点有reduce在运行。这时候集群资源没有得到充分利用，需要增加reduce的个数以便每个节点都有任务处理。

- **原则三：每个task的执行时间要合理。**

如果一个job，每个map或reduce的执行时间只有几秒钟，就意味着这个job的大部分时间都消耗在task的调度和进程启停阶段，因此需要增加每个task处理的数据大小。建议一个task处理时间为1分钟。

控制单个task处理时间的大小，可以通过如下配置来调整。

参数入口：

进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。

表 19-9 参数配置-2

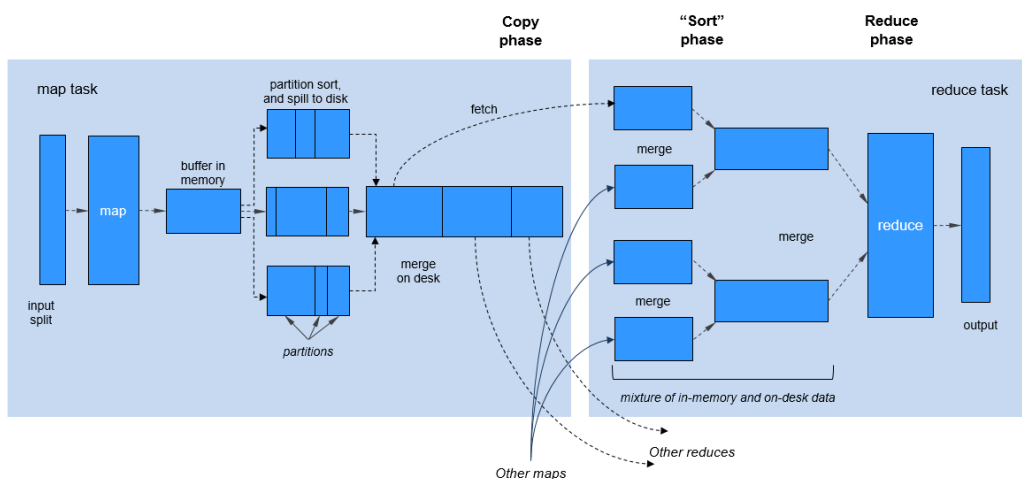
| 参数 | 描述 | 默认值 |
|---|--|-----|
| mapreduce.input.fileinputformat.split.maxsize | map输入信息应被拆分成的数据块的最大大小。
由用户定义的分片大小的设置及每个文件block大小的设置，可以计算分片的大小。计算公式如下：
$splitSize = \text{Math.max}(\text{minSize}, \text{Math.min}(\text{maxSize}, \text{blockSize}))$
如果maxSize设置大于blockSize，那么每个block就是一个分片，否则就会将一个block文件分隔为多个分片，如果block中剩下的一小段数据量小于splitSize，还是认为它是独立的分片。 | - |
| mapreduce.input.fileinputformat.split.minsize | 可以设置数据分片的数据最小值。 | 0 |

19.6.3 MapReduce Shuffle 调优

操作场景

Shuffle阶段是MapReduce性能的关键部分，包括了从Map task将中间数据写到磁盘一直到Reduce task拷贝数据并最终放到reduce函数的全部过程。这部分Hadoop提供了大量的调优参数。

图 19-2 Shuffle 过程



操作步骤

1. Map阶段的调优

- 判断Map使用的内存大小
判断Map分配的内存是否足够，一个简单的办法是查看运行完成的job的Counters中，对应的task是否发生过多次GC，以及GC时间占总task运行时间

之比。通常，GC时间不应超过task运行时间的10%，即GC time elapsed (ms)/CPU time spent (ms)<10%。

主要通过如下参数进行调整。

参数入口：

进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。

建议配置“mapreduce.map.java.opts”参数中“-Xmx”值为“mapreduce.map.memory.mb”参数值的0.8倍。

表 19-10 参数说明

| 参数 | 描述 | 默认值 |
|-------------------------|-------------|------|
| mapreduce.map.memory.mb | map任务的内存限制。 | 4096 |

| 参数 | 描述 | 默认值 |
|-------------------------|--|--|
| mapreduce.map.java.opts | map子任务的JVM参数。如果设置，会替代mapred.child.java.opts参数；如果未设置-Xmx，Xmx值从mapreduce.map.memory.mb*mapreduce.job.heap.memory-mb.ratio计算获取。 | <p>MRS 3.x之前版本： - Xmx2048M - Djava.net.preferIPv4Stack=true</p> <p>MRS 3.x及之后版本：</p> <ul style="list-style-type: none"> ● 集群已开启Kerberos认证： - Djava.net.preferIPv4Stack=true - Djava.net.preferIPv6Addresses=false - Djava.security.krb5.conf=\${BIGDATA_HOME}/common/runtime/krb5.conf - Dbeetle.application.home.path=\${BIGDATA_HOME}/common/runtime/security/config ● 集群未开启Kerberos认证： - Djava.net.preferIPv4Stack=true - Djava.net.preferIPv6Addresses=false - Dbeetle.application.home.path=\${BIGDATA_HOME}/common/runtime/security/config |

- 使用Combiner

在Map阶段，有一个可选过程，将同一个key值的中间结果合并，叫做Combiner。一般将reduce类设置为Combiner即可。通过Combiner，一般情况下可以显著减少Map输出的中间结果，从而减少shuffle过程的网络带宽占用。可通过如下接口为一个任务设置Combiner类。

表 19-11 Combiner 设置接口

| 类名 | 接口名 | 描述 |
|---------------------------------|--|--------------------|
| org.apache.hadoop.mapreduce.Job | public void
setCombinerClass(Class<?
extends Reducer> cls) | 为Job设置一个Combiner类。 |

2. Copy阶段的调优

数据是否压缩：

对Map的中间结果进行压缩，当数据量大时，会显著减少网络传输的数据量，但是也因为多了压缩和解压，带来了更多的CPU消耗。因此需要做好权衡。当任务属于网络瓶颈类型时，压缩Map中间结果效果明显。针对bulkload调优，压缩中间结果后性能提升60%左右。

配置方法：将“mapreduce.map.output.compress”参数值设置为“true”，将“mapreduce.map.output.compress.codec”参数值设置为“org.apache.hadoop.io.compress.SnappyCodec”。

3. Merge阶段的调优

通过调整如下参数减少reduce写磁盘的次数。

参数入口：

进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。

表 19-12 参数说明

| 参数 | 描述 | 默认值 |
|---|--|------|
| mapreduce.reduce.merge.inmem.threshold | 内存合并进程的文件数阈值。累计文件数达到阈值时会发起内存合并及溢出到磁盘。小于等于0的值表示该阈值不生效且仅基于ramfs的内存使用情况来触发合并。 | 1000 |
| mapreduce.reduce.shuffle.merge.percent | 发起内存合并的使用率阈值，表示为分配给映射输出信息的内存的比例（是由mapreduce.reduce.shuffle.input.buffer.percent设置的）。 | 0.66 |
| mapreduce.reduce.shuffle.input.buffer.percent | shuffle过程中分配给映射输出信息的内存占最大堆大小的比例。 | 0.70 |
| mapreduce.reduce.input.buffer.percent | Reduce过程中保存映射输出信息的内存相对于最大堆大小的比例。当shuffle结束时，需保证reduce开始前内存中所有剩余的映射输出信息所使用的内存小于该阈值。 | 0.0 |

19.6.4 MapReduce 大任务的 AM 调优

操作场景

当运行一个大任务（map总数达到了10万的规模），但是一直没有运行成功。经过查询发现是ApplicationMaster（以下简称AM）反应缓慢，最终超时失败。

此任务的问题是，task数量变多时，AM管理的对象也线性增长，因此就需要更多的内存来管理。AM默认分配的内存堆大小是1GB。

操作步骤

通过调大如下的参数来进行AM调优。

参数入口：

在Yarn客户端的“mapred-site.xml”配置文件中调整如下参数。“mapred-site.xml”配置文件在客户端安装路径的conf目录下，例如“/opt/client/Yarn/config”。

| 参数 | 描述 | 默认值 |
|------------------------------------|---|---|
| yarn.app.mapreduce.am.resource.mb | 该参数值必须大于下面参数的堆大小。单位：MB | 1536 |
| yarn.app.mapreduce.am.command-opts | 传递到MapReduce ApplicationMaster的JVM启动参数。 | <ul style="list-style-type: none"> MRS 3.x之前版本：-Xmx1024m -XX:CMSFullGCsBeforeCompaction=1 -XX:+UseConcMarkSweepGC -XX:+CMSParallelRemarkEnabled -XX:+UseCMSCompactAtFullCollection -verbose:gc MRS 3.x及之后版本：-Xmx1024m -XX:+UseConcMarkSweepGC -XX:+CMSParallelRemarkEnabled -verbose:gc -Djava.security.krb5.conf=\${KRB5_CONFIG} -Dhadoop.home.dir=\${BIGDATA_HOME}/FusionInsight_HD_xxx/install/FusionInsight-Hadoop-xxx/hadoop |

19.6.5 配置 MapReduce 任务推测执行

操作场景

当集群规模很大时（如几百上千台节点的集群），个别节点出现软硬件故障的概率会增大，并且会因此延长整个任务的执行时间（运行完成的任务会等待异常设备运行完成）。推测执行通过将一个task分给多台机器运行，取首先运行完成的节点。对于小集群，可以将该功能关闭。

操作步骤

参数入口：

进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。

| 参数 | 描述 | 默认值 |
|------------------------------|-----------------------------------|-------|
| mapreduce.map.speculative | 设置是否并行执行某些映射任务的多个实例。true表示开启。 | false |
| mapreduce.reduce.speculative | 设置是否并行执行某些reduce任务的多个实例。true表示开启。 | false |

19.6.6 通过 Slow Start 调优 MapReduce 任务

操作场景

Slow Start特性指定Map任务完成度为多少时Reduce任务可以启动，过早启动Reduce任务会导致资源占用，影响任务运行效率，但适当的提早启动Reduce任务会提高Shuffle阶段的资源利用率，提高任务运行效率。例如：某集群可启动10个Map任务，MapReduce作业共15个Map任务，那么在一轮Map任务执行完成后只剩5个Map任务，集群还有剩余资源，在这种场景下，配置Slow Start参数值小于1，比如0.8，则Reduce就可以利用集群剩余资源。

操作步骤

参数入口：

进入Mapreduce服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。

| 参数 | 描述 | 默认值 |
|--|---|-----|
| mapreduce.job.reduce.slowstart.completedmaps | 为job安排reduce前应完成的映射数的分数形式。默认100%的Map跑完后开始起Reduce。 | 1.0 |

19.6.7 MapReduce 任务 commit 阶段优化

操作场景

默认情况下，如果一个MR任务会产生大量的输出结果文件，那么该job在最后的commit阶段，会耗费较长的时间将每个task的临时输出结果commit到最终的结果输出目录。特别是在大集群中，大Job的commit过程会严重影响任务的性能表现。

针对以上情况，可以通过将以下参数

“mapreduce.fileoutputcommitter.algorithm.version”配置为“2”，来提升MR Job commit阶段的性能。

操作步骤

参数入口：

进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。具体操作请参考[修改集群服务配置参数](#)章节。

表 19-13 参数说明

| 参数 | 描述 | 默认值 |
|---|--|-----|
| mapreduce.fileoutputcommitter.algorithm.version | 用于指定Job的最终输出文件提交的算法版本，取值为“1”或“2”。
说明
版本2为建议的优化算法版本。该算法通过让任务直接将每个task的输出结果提交到最终的结果输出目录，从而减少大作业的输出提交时间。 | 2 |

19.6.8 降低 MapReduce 客户端运行任务失败率

配置场景

当网络不稳定或者集群IO、CPU负载过高的情况下，通过调整如下参数值，降低客户端应用的失败率，保证应用的正常运行。

配置描述

在客户端的“mapred-site.xml”配置文件中调整如下参数。

说明

“mapred-site.xml”配置文件在客户端安装路径的conf目录下，例如“/opt/client/Yarn/config”。

表 19-14 参数说明

| 参数 | 描述 | 默认值 |
|--|---|-----|
| mapreduce.reduce.shuffle.max-host-failures | MR任务在reduce过程中读取远端shuffle数据允许失败的次数。当设置次数大于5时，可以降低客户端应用的失败率。该参数适用于MRS 3.x版本。 | 5 |
| mapreduce.client.submit.file.replication | MR任务在运行时依赖的相关job文件在HDFS上的备份。当备份数大于10时，可以降低客户端应用的失败率。 | 10 |

19.7 MapReduce 日志介绍

日志描述

日志默认存储路径：

- JobhistoryServer：“/var/log/Bigdata/mapreduce/jobhistory”（运行日志），
“/var/log/Bigdata/audit/mapreduce/jobhistory”（审计日志）
- Container：“/srv/BigData/hadoop/data1/nm/containerlogs/application_{appid}/container_{\$contid}”

📖 说明

运行中的任务日志存储在以上路径中，运行结束后会基于YARN的配置是否汇聚到HDFS目录中，详情请参见[Yarn常用配置参数](#)。

日志归档规则：

MapReduce的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过50MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd_hh-mm-ss>.[编号].log.zip”。最多保留最近的100个压缩文件，压缩文件保留个数可以在参数配置界面中配置。

在MapReduce服务中，JobhistoryServer会定时去清理HDFS上存储的旧的日志文件（默认目录为HDFS文件系统中的“/mr-history/done”），具体清理的时间间隔参数配置为mapreduce.jobhistory.max-age-ms，默认值为1296000000，即15天。

表 19-15 MR 日志列表

| 日志类型 | 日志文件名 | 描述 |
|------|---|--------------------|
| 运行日志 | jhs-daemon-start-stop.log | 守护进程（Daemon）的启动日志。 |
| | hadoop-<SSH_USER>-jhshadaemon-<hostname>.log | 守护进程（Daemon）的运行日志。 |
| | hadoop-<SSH_USER>-<process_name>-<hostname>.out | MR运行环境信息日志。 |
| | historyserver-<SSH_USER>-<DATE>-<PID>-gc.log | MR服务垃圾回收日志。 |
| | jhs-haCheck.log | MR实例主备状态检查日志。 |
| | yarn-start-stop.log | MR服务启停操作日志。 |
| | yarn-prestart.log | MR服务启动前集群操作的记录日志。 |
| | yarn-postinstall.log | MR服务安装后启动前的工作日志。 |

| 日志类型 | 日志文件名 | 描述 |
|------|---|----------------|
| | yarn-cleanup.log | MR服务卸载时候的清理日志。 |
| | mapred-service-check.log | MR服务健康状态检测日志。 |
| | container_{\$contid} | Container日志。 |
| | hadoop-<SSH_USER>-<process_name>-<hostname>.log | MR运行日志。 |
| | mapred-switch-jhs.log | MR主备倒换日志。 |
| | env.log | 实例启停前的环境信息日志。 |
| 审计日志 | mapred-audit-jobhistory.log | MR操作审计日志。 |
| | SecurityAuth.audit | MR安全审计日志。 |

日志级别

MapReduce中提供了如表19-16所示的日志级别。其中日志级别优先级从高到低分别是FATAL、ERROR、WARN、INFO、DEBUG。程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 19-16 日志级别

| 级别 | 描述 |
|-------|-------------------------|
| FATAL | FATAL表示当前事件处理存在严重错误信息。 |
| ERROR | ERROR表示当前事件处理存在错误信息。 |
| WARN | WARN表示当前事件处理存在异常告警信息。 |
| INFO | INFO表示记录系统及各事件正常运行状态信息。 |
| DEBUG | DEBUG表示系统及系统的调试信息。 |

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 进入MapReduce服务参数“全部配置”界面，具体操作请参考[修改集群服务配置参数](#)。
- 步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤3** 选择所需修改的日志级别。
- 步骤4** 保存配置，在弹出窗口中单击“确定”使配置生效。

说明

配置完成后立即生效，不需要重启服务。

----结束

日志格式

MapReduce日志格式如下所示：

表 19-17 日志格式

| 日志类型 | 格式 | 示例 |
|------|---|--|
| 运行日志 | <yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置> | 2020-01-26 14:18:59,109 INFO main Client environment:java.compiler=<N/A> org.apache.zookeeper.Environment.logEnv(Environment.java:100) |
| 审计日志 | <yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事件的发生位置> | 2020-01-26 14:24:43,605 INFO main-EventThread USER=omm OPERATION=refreshAdminAcls TARGET=AdminService RESULT=SUCCESS org.apache.hadoop.yarn.server.resourcemanager.RMAuditLogger\$LogLevel \$6.printLog(RMAuditLogger.java:91) |

19.8 MapReduce 常见问题

19.8.1 ResourceManager 进行主备切换后，任务中断后运行时间过长

问题

在MapReduce任务运行过程中，ResourceManager发生主备切换，切换完成后，MapReduce任务继续执行，此时任务的运行时间过长。

回答

因为ResourceManager HA已启用，但是Work-preserving RM restart功能未启用。

如果Work-preserving RM restart功能未启用，ResourceManager切换时container会被kill，然后导致Application Master超时。Work-preserving RM restart功能介绍请参见：

MRS 3.2.0之前版本：<http://hadoop.apache.org/docs/r3.1.1/hadoop-yarn/hadoop-yarn-site/ResourceManagerRestart.html>

MRS 3.2.0及之后版本：<https://hadoop.apache.org/docs/r3.3.1/hadoop-yarn/hadoop-yarn-site/ResourceManagerRestart.html>

可以通过如下方式启用Work-preserving RM restart功能：

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中“yarn.resourcemanager.work-preserving-recovery.enabled”，设置参数值为“true”。保存配置后，在业务低峰期重启Yarn配置过期的实例。

19.8.2 MapReduce 任务长时间无进展

问题

MapReduce任务长时间无进展。

回答

一般是因为内存太少导致的。当内存较小时，任务中拷贝map输出的时间将显著增加。

为了减少等待时间，您可以适当增加堆内存空间。

任务的配置可根据mapper的数量和各mapper的数据大小来进行优化。根据输入数据的大小，优化“客户端安装路径/Yarn/config/mapred-site.xml”文件中的如下参数：

- “mapreduce.reduce.memory.mb”
- “mapreduce.reduce.java.opts”

例如：如果10个mapper的数据大小为5GB，那么理想的堆内存是1.5GB。随着数据大小的增加而增加堆内存大小。

19.8.3 为什么运行任务时客户端不可用

问题

当运行任务时，将MR ApplicationMaster或ResourceManager移动为D状态，为什么此时客户端会不可用？

回答

当运行任务时，将MR ApplicationMaster或ResourceManager移动为D状态（不间断睡眠状态）或T状态（停止状态），客户端会等待返回任务运行的状态，由于AM无返回，客户端会一直处于等待状态。

为避免出现上述场景，使用“core-site.xml”中的“ipc.client.rpc.timeout”配置项设置客户端超时时间。

该参数的参数值为毫秒。默认值为0，表示无超时。客户端超时的取值范围可以为0~2147483647毫秒。

📖 说明

- 如果Hadoop进程已处于D状态，重启该进程所处的节点。
- “core-site.xml”配置文件在客户端安装路径的conf目录下，例如“/opt/client/Yarn/config”。

19.8.4 在缓存中找不到 HDFS_DELEGATION_TOKEN 如何处理

问题

安全模式下，为什么在缓存中找不到HDFS_DELEGATION_TOKEN？

回答

在MapReduce中，默认情况下，任务完成之后，HDFS_DELEGATION_TOKEN将会被删除。因此如果在下一个任务中再次使用HDFS_DELEGATION_TOKEN，缓存中将会找不到HDFS_DELEGATION_TOKEN。

为了能够在随后的工作中再次使用同一个Token，为MapReduce任务配置参数。当参数为false时，用户能够再次使用同一个Token。

```
jobConf.setBoolean("mapreduce.job.complete.cancel.delegation.tokens", false);
```

19.8.5 如何在提交 MapReduce 任务时设置任务优先级

问题

如何在提交MapReduce任务时设置任务优先级？

回答

当您在客户端提交MapReduce任务时，可以在命令行中增加“-Dmapreduce.job.priority=<priority>”参数来设置任务优先级。格式如下：

```
yarn jar <jar> [mainClass] -Dmapreduce.job.priority=<priority> [path1] [path2]
```

命令行中参数含义为：

- <jar>：指定需要运行的jar包名称。
- [mainClass]：指jar包应用工程中的类的主方法。
- <priority>：指定任务的优先级，其取值可为：VERY_HIGH、HIGH、NORMAL、LOW、VERY_LOW。
- [path1]：指数据输入路径。
- [path2]：指数据输出路径。

例如，将“/opt/client/HDFS/hadoop/share/hadoop/mapreduce/hadoop-mapreduce-examples*.jar”包设置为高优先级任务。

```
yarn jar /opt/client/HDFS/hadoop/share/hadoop/mapreduce/hadoop-mapreduce-examples*.jar wordcount -Dmapreduce.job.priority=VERY_HIGH /DATA.txt /out/
```

19.8.6 MapReduce 任务运行失败，ApplicationMaster 出现物理内存溢出异常

问题

HBase bulkload任务有210000个map和10000个reduce，MapReduce任务运行失败，ApplicationMaster出现物理内存溢出异常。

```
For more detailed output, check the application tracking page:https://bigdata-55:8090/cluster/app/application_1449841777199_0003
Then click on links to logs of each attempt.
Diagnostics: Container [pid=21557,containerID=container_1449841777199_0003_02_000001] is running
beyond physical memory limits
Current usage: 1.0 GB of 1 GB physical memory used; 3.6 GB of 5 GB virtual memory used. Killing container.
Dump of the process-tree for container_1449841777199_0003_02_000001 :
|- PID PPID PGRPID SESSID CMD_NAME USER_MODE_TIME(MILLIS) SYSTEM_TIME(MILLIS)
  VMEM_USAGE(BYTES) RSSMEM_USAGE(PAGES) FULL_CMD_LINE
|- 21584 21557 21557 21557 (java) 12342 1627 3871748096 271331 ${BIGDATA_HOME}/jdk1.8.0_51//bin/
  java
  -Djava.io.tmpdir=/srv/BigData/hadoop/data1/nm/localdir/usercache/hbase/appcache/
  application_1449841777199_0003/container_1449841777199_0003_02_000001/tmp -
  Dlog4j.configuration=container-log4j.properties
  -Dyarn.app.container.log.dir=/srv/BigData/hadoop/data1/nm/containerlogs/
  application_1449841777199_0003/container_1449841777199_0003_02_000001 -
  Dyarn.app.container.log.filesize=0 -Dhadoop.root.logger=INFO,CLA
  -Dhadoop.root.logfile=syslog -Xmx784m org.apache.hadoop.mapreduce.v2.app.MRAppMaster
  |- 21557 21547 21557 21557 (bash) 0 0 13074432 368 /bin/bash -c ${BIGDATA_HOME}/jdk1.8.0_51//bin/
  java
  -Djava.io.tmpdir=/srv/BigData/hadoop/data1/nm/localdir/usercache/hbase/appcache/
  application_1449841777199_0003/container_1449841777199_0003_02_000001/tmp -
  Dlog4j.configuration=container-log4j.properties
  -Dyarn.app.container.log.dir=/srv/BigData/hadoop/data1/nm/containerlogs/
  application_1449841777199_0003/container_1449841777199_0003_02_000001 -
  Dyarn.app.container.log.filesize=0 -Dhadoop.root.logger=INFO,CLA
  -Dhadoop.root.logfile=syslog -Xmx784m org.apache.hadoop.mapreduce.v2.app.MRAppMaster 1>/srv/
  BigData/hadoop/data1/nm/containerlogs/application_1449841777199_0003/
  container_1449841777199_0003_02_000001/stdout
  2>/srv/BigData/hadoop/data1/nm/containerlogs/application_1449841777199_0003/
  container_1449841777199_0003_02_000001/stderr
  Container killed on request. Exit code is 143
  Container exited with a non-zero exit code 143
  Failing this attempt. Failing the application.
```

回答

这是性能规格的问题，MapReduce任务运行失败的根本原因是由于ApplicationMaster的内存溢出导致的，即物理内存溢出导致被NodeManager kill。

解决方案：

将ApplicationMaster的内存配置调大，在客户端“客户端安装路径/Yarn/config/mapred-site.xml”配置文件中优化如下参数：

- “yarn.app.mapreduce.am.resource.mb”
- “yarn.app.mapreduce.am.command-opts”，该参数中-Xmx值建议为0.8*“yarn.app.mapreduce.am.resource.mb”

参考规格：

ApplicationMaster配置如下时，可以同时支持并发Container数为2.4万个。

- “yarn.app.mapreduce.am.resource.mb” =2048

- “yarn.app.mapreduce.am.command-opts” 该参数中-Xmx=1638m

19.8.7 MapReduce 作业信息无法通过 ResourceManager Web UI 页面的 Tracking URL 打开

问题

MapReduce JobHistoryServer服务地址变更后，为什么运行完的MapReduce作业无法通过ResourceManager Web UI页面打开？

回答

JobHistoryServer地址（mapreduce.jobhistory.address / mapreduce.jobhistory.webapp.<https.>address）是MapReduce参数，MapReduce客户端提交作业时，会将此地址随任务一起提交给ResourceManager。ResourceManager在作业完成后，将此参数作为查看作业历史信息的跳转地址保存在RMStateStore中。

JobHistoryServer服务地址变更后，需要将新的服务地址及时更新到MapReduce客户端配置文件中，否则，新运行的作业在查看作业历史信息时，仍然会指向原JobHistoryServer地址，导致无法正常跳转到作业历史信息页面。服务地址变更前运行的MapReduce作业，由于其跳转信息已经保存在RMStateStore中，无法变更，因此从ResourceManager Web UI页面是无法进行正常跳转的，但可以直接访问新的JobHistoryServer服务地址进行查找，作业信息不会丢失。

19.8.8 多个 NameService 环境下运行 MapReduce 任务失败

问题

多个NameService环境下，运行使用viewFS功能的MapReduce或YARN任务失败。

回答

当使用viewFS时，只有在viewFS中挂载的目录才能被访问到。所以最可能的原因是配置的路径没有在viewFS的挂载点上。例如：

```
<property>
<name>fs.defaultFS</name>
<value>viewfs://ClusterX</value>
</property>
<property>
<name>fs.viewfs.mounttable.ClusterX.link./folder1</name>
<value>hdfs://NS1/folder1</value>
</property>
<property>
<name>fs.viewfs.mounttable.ClusterX.link./folder2</name>
<value>hdfs://NS2/folder2</value>
</property>
```

对于依赖HDFS的MR配置中，需要使用已挂载的目录。

错误示例：

```
<property>
<name>yarn.app.mapreduce.am.staging-dir</name>
<value>/tmp/hadoop-yarn/staging</value>
</property>
```

根目录 (/) 在viewFS中是无法访问的。

正确示例：

```
<property>  
<name>yarn.app.mapreduce.am.staging-dir</name>  
<value>/folder1/tmp/hadoop-yarn/staging</value>  
</property>
```

19.8.9 基于分区的任务黑名单异常如何处理

问题

Map&Reduce任务失败，并且故障节点数与集群总节点数的比值低于“yarn.resourceanager.am-scheduling.node-blacklisting-disable-threshold”配置的黑名单阈值，为什么Map&Reduce任务故障节点没有加入黑名单？

回答

当集群中有超过阈值的节点都被加入黑名单时，黑名单会释放这些节点，其中阈值为故障节点数与集群总节点数的比值。现在每个节点都有其标签表达式，黑名单阈值应根据有效节点标签表达式关联的节点数进行计算，其值为故障节点数与有效节点标签表达式关联的节点数的比值。

假设集群中有100个节点，其中有10个节点为有效节点标签表达式关联的节点（labelA）。其中所有有效节点标签表达式关联的节点都已经故障，黑名单节点释放阈值默认值为0.33，按照传统的计算方式， $10/100=0.1$ ，远小于该阈值。这就造成这10个节点永远无法得到释放，Map&Reduce任务一直无法获取节点，应用程序无法正常运行。实际需要根据与Map&Reduce任务的有效节点关联的节点总数进行计算，即 $10/10=1$ ，大于黑名单节点释放阈值，节点被释放。

因此即使故障节点数与集群总节点数的比值没有超过阈值，也存在黑名单将这些节点释放的情况。

20 使用 OpenTSDB

20.1 使用 MRS 客户端操作 OpenTSDB 指标数据

用户可以根据业务需要，在MRS集群的客户端中进行交互式操作。启用Kerberos认证的集群，需要操作的用户属于“opentsdb, hbase, opentsdbgroup和supergroup”组且拥有HBase权限。

前提条件

- 获取用户“admin”账号密码。“admin”密码在创建MRS集群时由用户指定。
- 已安装集群客户端，例如安装目录为“/opt/client”，以下操作的客户端目录只是举例，请根据实际安装目录修改。更新客户端，具体请参见[更新客户端（3.x之前版本）](#)。

使用客户端

步骤1 如果当前集群已启用Kerberos认证，登录MRS Manager页面，创建属于“opentsdb, hbase, opentsdbgroup和supergroup”组且拥有HBase权限的用户，例如创建用户为opentsdbuser，具体请参考[准备开发用户](#)。如果当前集群未启用Kerberos认证，则无需执行此步骤。

步骤2 根据业务情况，准备好客户端，并登录安装客户端的节点。

例如在Master2节点更新客户端，则登录该节点使用客户端，具体参见[更新客户端（3.x之前版本）](#)。

步骤3 执行以下命令切换用户。

```
sudo su - omm
```

步骤4 执行以下命令，切换到客户端目录，例如“/opt/client”。

```
cd /opt/client
```

步骤5 执行以下命令，配置环境变量。

```
source bigdata_env
```

步骤6 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

- 当用户为“人机”用户时：执行 `kinit opentsdbuser` 认证用户
- 当用户为“机机”用户时：下载用户认证凭据文件，保存并解压获取用户的 `user.keytab` 文件与 `krb5.conf` 文件，进入解压后的 `user.keytab` 目录下，执行 `kinit -kt user.keytab opentsdbuser` 认证用户

步骤7 操作 Opentsdb 数据，具体请参见 [操作数据](#)。

----结束

操作数据

- 查看帮助

执行 `tsdb` 命令打印出当前 opentsdb 所支持的所有命令。如，`fsck`, `import`, `mkmetric`, `query`, `tsd`, `scan`, `search`, `uid`, `version`。

回显信息：

```
tsdb: error: unknown command "  
usage: tsdb <command> [args]  
Valid commands: fsck, import, mkmetric, query, tsd, scan, search, uid, version
```

- 创建 OpenTSDB 指标

执行 `tsdb mkmetric` 命令创建指标。例如执行 `tsdb mkmetric sys.cpu.user` 命令创建名为 `sys.cpu.user` 的指标。

回显信息：

```
Start run net.opentsdb.tools.UidManager, args: assign metrics sys.cpu.user  
metrics sys.cpu.user: [0, 0, 6]
```

- 向 OpenTSDB 指标中导入数据

- a. 准备指标文件，例如包含如下信息的 `importData.txt` 文件。

```
sys.cpu.user 1356998400 41 host=web01 cpu=0  
sys.cpu.user 1356998401 42 host=web01 cpu=0  
sys.cpu.user 1356998402 44 host=web01 cpu=0  
sys.cpu.user 1356998403 47 host=web01 cpu=0  
sys.cpu.user 1356998404 42 host=web01 cpu=0  
sys.cpu.user 1356998405 42 host=web01 cpu=0
```

- b. 执行 `tsdb import` 命令导入指标数据。例如执行 `tsdb import importData.txt` 命令导入 `importData.txt` 文件。

```
Start run net.opentsdb.tools.TextImporter, args: importData.txt  
2019-06-26  
15:45:22,091 INFO [main] TextImporter:  
reading from file:importData.txt  
2019-06-26  
15:45:22,102 INFO [main] TextImporter:  
Processed importData.txt in 11 ms, 6 data points (545.5 points/s)  
2019-06-26  
15:45:22,102 INFO [main] TextImporter:  
Total: imported 6 data points in 0.012s (504.0 points/s)
```

- 查询 OpenTSDB 指标

执行 `tsdb uid metrics` 命令获取当前 OpenTSDB 中存入的指标。例如执行 `tsdb uid metrics sys.cpu.user` 命令查询 `sys.cpu.user` 的数据。

回显信息：

```
Start run net.opentsdb.tools.UidManager, args: metrics sys.cpu.user  
metrics sys.cpu.user: [0, 0, 6]
```

如需获得更多信息，请执行 `tsdb uid` 命令。


```
Start run net.opentsdb.tools.UidManager, args:
Not enough arguments
Usage: uid <subcommand> args
Sub commands:
  grep [kind] <RE>: Finds matching IDs.
  assign <kind> <name> [names]: Assign an ID for the given name(s).
  rename <kind> <name> <newname>: Renames this UID.
  delete <kind> <name>: Deletes this UID.
  fsck: [fix] [delete_unknown] Checks the consistency of UIDs.
      fix          - Fix errors. By default errors are logged.
      delete_unknown - Remove columns with unknown qualifiers.
                    The "fix" flag must be supplied as well.
  [kind] <name>: Lookup the ID of this name.
  [kind] <ID>: Lookup the name of this ID.
  metasync: Generates missing TSUID and UID meta entries, updates created timestamps
  metapurge: Removes meta data entries from the UID table
  treesync: Process all timeseries meta objects through tree rules
  treepurge <id> [definition]: Purge a tree and/or the branches from storage. Provide an integer Tree
  ID and                        optionally add "true" to delete the tree definition
Example values for [kind]: metrics, tagk (tag name), tagv (tag value).
--config=PATH  Path to a configuration file (default: Searches for file see docs).
--idwidth=N    Number of bytes on which the Uniqueld is encoded.
--ignore-case  Ignore case distinctions when matching a regexp.
--table=TABLE  Name of the HBase table where to store the time series (default: tsdb).
--uidtable=TABLE Name of the HBase table to use for Unique IDs (default: tsdb-uid).
--verbose      Print more logging messages and not just errors.
--zkbasedir=PATH Path under which is the znode for the -ROOT- region (default: /hbase).
--zkquorum=SPEC Specification of the ZooKeeper quorum to use (default: localhost).
-i            Short for --ignore-case.
-v            Short for --verbose.
```

- 扫描Opentsdb的指标数据

执行**tsdb query**命令批量查询导入的指标数据，命令格式如下：**tsdb query <START-DATE> <END-DATE> <aggregator> <metric> <tagk=tagv>**，例如执行**tsdb query 0 1h-ago sum sys.cpu.user host=web01**

```
Start run net.opentsdb.tools.CliQuery, args: 0 1h-ago sum sys.cpu.user host=web01
sys.cpu.user 1356998400000 41 {host=web01, cpu=0}
sys.cpu.user 1356998401000 42 {host=web01, cpu=0}
sys.cpu.user 1356998402000 44 {host=web01, cpu=0}
sys.cpu.user 1356998403000 47 {host=web01, cpu=0}
sys.cpu.user 1356998404000 42 {host=web01, cpu=0}
sys.cpu.user 1356998405000 42 {host=web01, cpu=0}
```

📖 说明

- **<START-DATE>**:要查询指标的起始时间点。
- **<END-DATE>**:要查询指标的结束时间点。
- **<aggregator>**:查询数据的聚合方式。
- **<metric>**:所需查询的指标名称。
- **<tagk=tagv>**:标签的key和value。
- 删除录入的Opentsdb指标
执行命令**tsdb uid delete**命令删除录入的指标及值。例如删除sys.cpu.user指标可执行命令**tsdb uid delete metrics sys.cpu.user**。

```
Start run net.opentsdb.tools.UidManager, args: delete metrics sys.cpu.user
```

20.2 使用 curl 命令操作 OpenTSDB

写入数据

例如，录入一个指标名称为testdata，时间戳为1524900185，值为true，标签为key，value的指标数据。

```
curl -ki -X POST -d '{"metric":"testdata", "timestamp":1524900185, "value":"true", "tags": {"key":"value"}}' https://<tsd_ip>:4242/api/put?sync
```

<tsd_ip>表示所需写入数据的Opentsdb服务的TSD实例的IP地址。

```
HTTP/1.1 204 No Content
Content-Type: application/json; charset=UTF-8
Content-Length:0
```

查询数据

例如，可查询指标testdata在过去三年的汇总信息。

```
curl -ks https://<tsd_ip>:4242/api/query?start=3y-ago&m=sum:testdata | python -m json.tool
```

- <tsd_ip>：所需访问Opentsdb服务的TSD实例IP或主机名。
- <start=3y-ago&m=sum:testdata>：在请求中可能无法识别“&”符号，需对其进行转义。
- <python -m json.tool>（可选）：把响应的请求转换为json格式。

```
[
  {
    "aggregateTags": [],
    "dps": {
      "1524900185": 1
    },
    "metric": "testdata",
    "tags": {
      "key": "value"
    }
  }
]
```

查询 tsd 状态信息

例如，可查询连接HBase的客户端信息。

```
curl -ks https://<tsd_ip>:4242/api/stats/region_clients | python -m json.tool
```

<tsd_ip>：所需访问Opentsdb服务的TSD实例IP地址。

```
[
  {
    "dead": false,
    "endpoint": "/xx.xx.xx.xx:16020",
    "inflightBreachd": 0,
    "pendingBatchedRPCs": 0,
    "pendingBreachd": 0,
    "pendingRPCs": 0,
    "rpcResponsesTimedout": 0,
    "rpcResponsesUnknown": 0,
    "rpcid": 78,
    "rpcsInFlight": 0,
    "rpcsSent": 79,
    "rpcsTimedout": 0,
    "writesBlocked": 0
  }
]
```

21 使用 Oozie

21.1 使用 Oozie 客户端提交作业

21.1.1 Oozie 客户端配置说明

操作场景

该任务指导用户在运维场景或业务场景中使用Oozie客户端。Oozie支持提交多种类型任务，例如Hive、Spark2x、Loader、Mapreduce、Java、DistCp、Shell、HDFS、SSH、SubWorkflow、Streaming、定时任务等。

前提条件

- 已安装客户端，具体请参考[安装客户端](#)。例如安装目录为“/opt/client”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- 各组件业务用户由MRS集群管理员根据业务需要创建。安全模式下，“机机”用户需要下载keytab文件。“人机”用户第一次登录时需修改密码。

使用 Oozie 客户端

步骤1 以客户端安装用户，登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录，该操作的客户端目录只是举例，请根据实际安装目录修改。

```
cd /opt/client
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 判断集群认证模式。

- 安全模式，执行以下命令进行用户认证。*exampleUser*为提交任务的用户名。

```
kinit exampleUser
```
- 普通模式，执行**步骤5**。

步骤5 配置Hue。

1. spark2x环境配置（如果不涉及spark2x任务，可以跳过此步骤）：

```
hdfs dfs -put /opt/client/Spark2x/spark/jars/*.jar /user/oozie/share/lib/spark2x/
```

当HDFS目录“/user/oozie/share”中的Jar包发生变化时，需要重启Oozie服务。

2. 上传Oozie配置文件以及Jar包至HDFS：

```
hdfs dfs -mkdir /user/exampleUser
```

```
hdfs dfs -put -f /opt/client/Oozie/oozie-client-*/examples /user/exampleUser/
```

📖 说明

- `exampleUser`为提交任务的用户名。
- 在提交任务的用户和非`job.properties`文件均无变更的前提下，客户端安装目录/Oozie/oozie-client-*/examples目录一经上传HDFS，后续可重复使用，无需多次提交。
- 解决Spark和Yarn关于jetty的jar冲突。

```
hdfs dfs -rm -f /user/oozie/share/lib/spark/jetty-all-9.2.22.v20170606.jar
```

- 普通模式下，上传过程如果遇到“Permission denied”的问题，可执行以下命令进行处理。

```
su - omm
```

```
source /opt/client/bigdata_env
```

```
hdfs dfs -chmod -R 777 /user/oozie
```

```
exit
```

步骤6 本操作以在Oozie客户端提交MapReduce任务为例进行演示。

1. 修改任务执行配置文件：

```
cd /opt/client/Oozie/oozie-client-*/examples/apps/map-reduce/  
vi job.properties
```

```
nameNode=hdfs://hacluster  
resourceManager=10.64.35.161:8032 ( 10.64.35.161为Yarn resourceManager ( Active ) 节点业务平面  
IP; 8032为yarn.resourcemanager.port )  
queueName=default  
examplesRoot=examples  
user.name=admin  
oozie.wf.application.path=${nameNode}/user/${user.name}/${examplesRoot}/apps/map-reduce #hdfs  
上传路径  
outputDir=map-reduce  
oozie.wf.rerun.failnodes=true
```

2. 运行Oozie任务：

```
oozie job -oozie https://oozie角色的主机名:21003/oozie/ -config  
job.properties -run
```

“21003”为Oozie HTTPS请求的运行端口，可在FusionInsight Manager，选择“集群 > 服务 > Oozie > 配置”，在搜索框中搜索“OOZIE_HTTPS_PORT”查看。

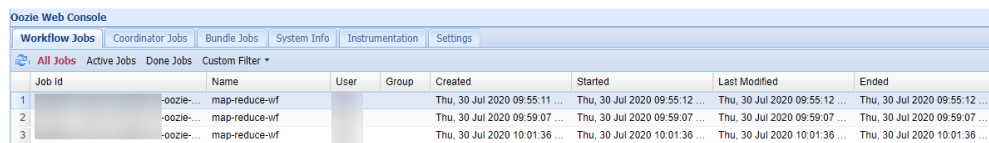
```
[root@kwephispra44947 map-reduce]# oozie job -oozie https://kwephispra44948:21003/oozie/ -  
config job.properties -run
```

```
.....
```

```
job: 0000000-200730163829770-oozie-omm-W
```

3. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Oozie”，单击“oozie WebUI”后的超链接进入Oozie页面，在Oozie的WebUI上查看任务运行结果。

图 21-1 任务运行结果



Job Id	Name	User	Group	Created	Started	Last Modified	Ended
1	-oozie-... map-reduce-wf			Thu, 30 Jul 2020 09:55:11 ...	Thu, 30 Jul 2020 09:55:12 ...	Thu, 30 Jul 2020 09:55:12 ...	Thu, 30 Jul 2020 09:55:12 ...
2	-oozie-... map-reduce-wf			Thu, 30 Jul 2020 09:59:07 ...	Thu, 30 Jul 2020 09:59:07 ...	Thu, 30 Jul 2020 09:59:07 ...	Thu, 30 Jul 2020 09:59:07 ...
3	-oozie-... map-reduce-wf			Thu, 30 Jul 2020 10:01:36 ...	Thu, 30 Jul 2020 10:01:36 ...	Thu, 30 Jul 2020 10:01:36 ...	Thu, 30 Jul 2020 10:01:36 ...

----结束

21.1.2 使用 Oozie 客户端提交 Hive 任务

操作场景

该任务指导用户在使用Oozie客户端提交Hive任务

Hive任务有如下类型：

- Hive作业
使用JDBC方式连接的Hive作业。
- Hive2作业
使用Beeline方式连接的Hive作业。

本文以使用Oozie客户端提交Hive作业为例介绍。

📖 说明

- 使用Oozie客户端提交Hive2作业与提交Hive作业操作步骤一致，只需将操作步骤中对应路径的“/Hive”改成“/Hive2”即可。
例如，Hive作业运行目录“/opt/client/Oozie/oozie-client-*/examples/apps/hive/”，则Hive2对应的运行目录为“/opt/client/Oozie/oozie-client-*/examples/apps/hive2/”。
- 建议下载使用最新版本的客户端。

前提条件

- Hive和Oozie组件及客户端已经安装，并且正常运行。
- 已创建或获取访问Oozie服务的人机用户账号及密码。

📖 说明

- 该用户需要从属于hadoop、supergroup、hive组，同时添加Oozie的角色操作权限。若使用Hive多实例，该用户还需要从属于具体的Hive实例组，如hive3。
- 用户同时还需要至少有manager_viewer权限的角色。
- 获取运行状态的Oozie服务器（任意实例）URL，如“https://10.1.130.10:21003/oozie”。
- 获取运行状态的Oozie服务器主机名，如“10-1-130-10”。
- 获取Yarn ResourceManager主节点IP，如10.1.130.11。

操作步骤

步骤1 以客户端安装用户，登录安装Oozie客户端的节点。

步骤2 执行以下命令，获取安装环境信息。其中“/opt/client”为客户端安装路径，该操作的客户端目录只是举例，请根据实际安装目录修改。

```
source /opt/client/bigdata_env
```

步骤3 判断集群认证模式。

- 安全模式，执行**kinit**命令进行用户认证。
例如，使用**oozieuser**用户进行认证。

```
kinit oozieuser
```

- 普通模式，执行**步骤4**。

步骤4 执行以下命令，进入样例目录。

```
cd /opt/client/Oozie/oozie-client-*/examples/apps/hive/
```

该目录下需关注文件如**表21-1**所示。

表 21-1 文件说明

文件名称	描述
hive-site.xml	Hive任务的配置文件。
job.properties	工作流的参数变量定义文件。
script.q	Hive任务的SQL脚本。
workflow.xml	工作流的规则定制文件。

步骤5 执行以下命令，编辑“job.properties”文件。

```
vi job.properties
```

修改如下内容：

更改“userName”的参数值为提交任务的人机用户名，例如“userName=oozieuser”。

步骤6 执行**oozie job**命令，运行工作流文件。

```
oozie job -oozie https://oozie角色的主机名:21003/oozie/ -config job.properties -run
```

说明

- 命令参数解释如下：
 - oozie：实际执行任务的Oozie服务器URL。
 - config：工作流属性文件。
 - run：运行工作流。
- 执行完工作流文件，显示job id表示提交成功，例如：job: 0000021-140222101051722-oozie-omm-W。登录Oozie管理页面，查看运行情况。
使用**oozieuser**用户，登录Oozie WebUI页面：<https://oozie角色的ip地址:21003/oozie>。
Oozie的WebUI界面中，可在页面表格根据jobid查看已提交的工作流信息。

----结束

21.1.3 使用 Oozie 客户端提交 Spark2x 任务

操作场景

该任务指导用户在使用Oozie客户端提交Spark2x任务。

📖 说明

请下载使用最新版本的客户端。

前提条件

- Spark2x和Oozie组件安装完成且运行正常，客户端安装成功。
如果当前客户端为旧版本，需要重新下载和安装客户端。
- 已创建或获取访问Oozie服务的人机用户账号及密码。

📖 说明

- 该用户需要从属于hadoop、supergroup、hive组，同时添加Oozie的角色操作权限。若使用Hive多实例，该用户还需要从属于具体的Hive实例组，如hive3。
- 用户同时还需要至少有manager_viewer权限的角色。
- 获取运行状态的Oozie服务器（任意实例）URL，如“https://10.1.130.10:21003/oozie”。
- 获取运行状态的Oozie服务器主机名，如“10-1-130-10”。
- 获取Yarn ResourceManager主节点IP，如“10.1.130.11”。

操作步骤

步骤1 以客户端安装用户登录安装Oozie客户端的节点。

步骤2 执行以下命令，获取安装环境信息。其中“/opt/client”为客户端安装路径，该操作的客户端目录只是举例，请根据实际安装目录修改。

```
source /opt/client/bigdata_env
```

步骤3 判断集群认证模式。

- 安全模式，执行kinit命令进行用户认证。
例如，使用oozieuser用户进行认证。

```
kinit oozieuser
```

- 普通模式，执行**步骤4**。

步骤4 执行以下命令，进入样例目录。

```
cd /opt/client/Oozie/oozie-client-*/examples/apps/spark2x/
```

该目录下需关注文件如表21-2所示。

表 21-2 文件说明

文件名称	描述
job.properties	工作流的参数变量定义文件。

文件名称	描述
workflow.xml	工作流的规则定制文件。
lib	工作流运行依赖的jar包目录。

步骤5 执行以下命令，编辑“job.properties”文件。

vi job.properties

修改如下内容：

更改“userName”的参数值为提交任务的人机用户名，例如“userName=oozieuser”。

步骤6 执行**oozie job**命令，运行工作流文件。

oozie job -oozie https://oozie角色的主机名:21003/oozie/ -config job.properties -run

说明

- 命令参数解释如下：
 - oozie：实际执行任务的Oozie服务器URL。
 - config：工作流属性文件。
 - run：运行工作流。
- 执行完工作流文件，显示“job id”表示提交成功，例如“job:0000021-140222101051722-oozie-omm-W”。登录Oozie管理页面，查看运行情况。使用**oozieuser**用户，登录Oozie WebUI页面：<https://oozie角色的ip地址:21003/oozie>。Oozie的WebUI界面中，可在页面表格根据“job id”查看已提交的工作流信息。

----结束

21.1.4 使用 Oozie 客户端提交 Loader 任务

操作场景

该任务指导用户在使用Oozie客户端提交Loader任务。

说明

请下载使用最新版本的客户端。

前提条件

- Loader和Oozie组件及客户端已经安装，并且正常运行。
- 已创建或获取访问Oozie服务的人机用户账号及密码。

说明

- 该用户需要从属于hadoop、supergroup、hive组，同时添加Oozie的角色操作权限。若使用Hive多实例，该用户还需要从属于具体的Hive实例组，如hive3。
- 用户同时还需要至少有manager_viewer权限的角色。

- 获取运行状态的Oozie服务器（任意实例）URL，如“https://10.1.130.10:21003/oozie”。
- 获取运行状态的Oozie服务器主机名，如“10-1-130-10”。
- 获取Yarn ResourceManager主节点IP，如10.1.130.11。
- 创建需要调度的Loader作业，并获取该作业ID。

操作步骤

步骤1 以客户端安装用户，登录安装Oozie客户端的节点。

步骤2 执行以下命令，获取安装环境信息。其中“/opt/client”为客户端安装路径，该操作的客户端目录只是举例，请根据实际安装目录修改。

```
source /opt/client/bigdata_env
```

步骤3 判断集群认证模式。

- 安全模式，执行kinit命令进行用户认证。
例如，使用oozieuser用户进行认证。

```
kinit oozieuser
```

- 普通模式，执行**步骤4**。

步骤4 执行以下命令，进入样例目录。

```
cd /opt/client/Oozie/oozie-client-*/examples/apps/sqoop/
```

该目录下需关注文件如表21-3所示。

表 21-3 文件说明

文件名称	描述
job.properties	工作流的参数变量定义文件。
workflow.xml	工作流的规则定制文件。

步骤5 执行以下命令，编辑“job.properties”文件。

```
vi job.properties
```

修改如下内容：

更改“userName”的参数值为提交任务的人机用户名，例如“userName=oozieuser”。

步骤6 执行以下命令，编辑“workflow.xml”文件。

```
vi workflow.xml
```

修改如下内容：

“command”的值修改为需要调度的已有Loader作业ID，例如1。

将“workflow.xml”文件上传至“job.properties”文件中的HDFS路径。

```
hdfs dfs -put -f workflow.xml /user/userName/examples/apps/sqoop
```

步骤7 执行`oozie job`命令，运行 workflow 文件。

```
oozie job -oozie https://oozie角色的主机名:21003/oozie/ -config job.properties -run
```

📖 说明

- 命令参数解释如下：
 - oozie：实际执行任务的Oozie服务器URL。
 - config：workflow属性文件。
 - run：运行workflow。
- 执行完workflow文件，显示job id表示提交成功，例如：job: 0000021-140222101051722-oozie-omm-W。登录Oozie管理页面，查看运行情况。
使用`oozieuser`用户，登录Oozie WebUI页面：`https://oozie角色的ip地址:21003/oozie`。
Oozie的WebUI界面中，可在页面表格根据jobid查看已提交的workflow信息。

---结束

21.1.5 使用 Oozie 客户端提交 DistCp 任务

操作场景

该任务指导用户在使用Oozie客户端提交DistCp任务。

📖 说明

请下载使用最新版本的客户端。

前提条件

- HDFS和Oozie组件安装完成且运行正常，客户端安装成功。
如果当前客户端为旧版本，需要重新下载和安装客户端。
- 已创建或获取访问Oozie服务的人机用户账号及密码。

📖 说明

- 该用户需要从属于hadoop、supergroup、hive组，同时添加Oozie的角色操作权限。若使用Hive多实例，该用户还需要从属于具体的Hive实例组，如hive3。
- 用户同时还需要至少有manager_viewer权限的角色。
- 已获取运行状态的Oozie服务器（任意实例）URL，如“`https://10.1.130.10:21003/oozie`”。
- 已获取运行状态的Oozie服务器主机名，如“`10-1-130-10`”。
- 已获取Yarn ResourceManager主节点IP，如“`10.1.130.11`”。

操作步骤

步骤1 以客户端安装用户登录安装Oozie客户端的节点。

步骤2 执行以下命令，获取安装环境信息。其中“`/opt/client`”为客户端安装路径，该操作的客户端目录只是举例，请根据实际安装目录修改。

```
source /opt/client/bigdata_env
```

步骤3 判断集群认证模式。

- 安全模式，执行**kinit**命令进行用户认证。
例如，使用**oozieuser**用户进行认证。

kinit oozieuser

- 普通模式，执行**步骤4**。

步骤4 执行以下命令，进入样例目录。

```
cd /opt/client/Oozie/oozie-client-*/examples/apps/distcp/
```

该目录下需关注文件如**表21-4**所示。

表 21-4 文件说明

文件名称	描述
job.properties	工作流的参数变量定义文件。
workflow.xml	工作流的规则定制文件。

步骤5 执行以下命令，编辑“job.properties”文件。

```
vi job.properties
```

修改如下内容：

更改“userName”的参数值为提交任务的人机用户名，例如“userName=oozieuser”。

步骤6 是否是跨安全集群的DistCp。

- 是，执行**步骤7**。
- 否，则执行**步骤9**。

步骤7 对两个集群进行跨Manager集群互信。

步骤8 备份并且修改workflow.xml的文件内容，命令如下：

```
cp workflow.xml workflow.xml.bak
```

```
vi workflow.xml
```

修改以下内容：

```
<workflow-app xmlns="uri:oozie:workflow:1.0" name="distcp-wf">
  <start to="distcp-node"/>
  <action name="distcp-node">
    <distcp xmlns="uri:oozie:distcp-action:1.0">
      <resource-manager>${resourceManager}</resource-manager>
      <name-node>${nameNode}</name-node>
      <prepare>
        <delete path="hdfs://target_ip:target_port/user/${userName}/${examplesRoot}/output-data/${outputDir}"/>
      </prepare>
      <configuration>
        <property>
          <name>mapred.job.queue.name</name>
          <value>${queueName}</value>
        </property>
      </configuration>
    </distcp>
  </action>
</workflow-app>
```

```
<name>oozie.launcher.mapreduce.job.hdfs-servers</name>
<value>hdfs://source_ip:source_port,hdfs://target_ip:target_port</value>
</property>
</configuration>
<arg>${nameNode}/user/${userName}/${examplesRoot}/input-data/text/data.txt</arg>
<arg>hdfs://target_ip:target_port/user/${userName}/${examplesRoot}/output-data/${outputDir}/
data.txt</arg>
</distcp>
<ok to="end"/>
<error to="fail"/>
</action>
<kill name="fail">
  <message>DistCP failed, error message[${wf.errorMessage(wf.lastErrorNode())}]</message>
</kill>
<end name="end"/>
</workflow-app>
```

其中“target_ip:target_port”为另一个互信集群的HDFS active namenode地址，例如：10.10.10.233:25000。

“source_ip:source_port”为源集群的HDFS active namenode地址，例如：10.10.10.223:25000。

两个IP地址和端口都需要根据自身的集群实际情况修改。

步骤9 执行`oozie job`命令，运行工作流文件。

```
oozie job -oozie https://oozie角色的主机名:21003/oozie/ -config job.properties -run
```

📖 说明

- 命令参数解释如下：
 - oozie：实际执行任务的Oozie服务器URL。
 - config：工作流属性文件。
 - run：运行工作流。
- 执行完工作流文件，显示“job id”表示提交成功，例如“job:0000021-140222101051722-oozie-omm-W”。登录Oozie管理页面，查看运行情况。使用`oozieuser`用户，登录Oozie WebUI页面：<https://oozie角色的ip地址:21003/oozie>。Oozie的WebUI界面中，可在页面表格根据“job id”查看已提交的工作流信息。

----结束

21.1.6 使用 Oozie 客户端提交其它任务

操作场景

除了Hive、Spark2x、Loader任务，也支持使用Oozie客户端提交MapReduce、Java、Shell、HDFS、SSH、SubWorkflow、Streaming、定时等任务。

📖 说明

请下载使用最新版本的客户端。

前提条件

- Oozie组件及客户端已经安装，并且正常运行。
- 已创建或获取访问Oozie服务的人机用户账号及密码。

📖 说明

- Shell任务：
该用户需要从属于hadoop、supergroup组，添加Oozie的角色操作权限，并确保Shell脚本在每个nodemanager节点都有执行权限。
- SSH任务：
该用户需要从属于hadoop、supergroup组，添加Oozie的角色操作权限，并完成互信配置。
- 其他任务：
该用户需要从属于hadoop、supergroup组，添加Oozie的角色操作权限，并具备对应任务类型所需的权限。
- 用户同时还需要至少manager_viewer权限的角色。
- 获取运行状态的Oozie服务器（任意实例）URL，如“https://10.1.130.10:21003/oozie”。
- 获取运行状态的Oozie服务器主机名，如“10-1-130-10”。
- 获取Yarn ResourceManager主节点IP，如10.1.130.11。

操作步骤

步骤1 以客户端安装用户，登录安装Oozie客户端的节点。

步骤2 执行以下命令，获取安装环境信息。其中“/opt/client”为客户端安装路径，该操作的客户端目录只是举例，请根据实际安装目录修改。

```
source /opt/client/bigdata_env
```

步骤3 判断集群认证模式。

- 安全模式，执行kinit命令进行用户认证。
例如，使用oozieuser用户进行认证。

```
kinit oozieuser
```

- 普通模式，执行**步骤4**。

步骤4 根据提交任务类型，进入对应样例目录。

表 21-5 样例目录列表

任务类型	样例目录
Mapreduce任务	客户端安装目录/Oozie/oozie-client-*/examples/apps/map-reduce
Java任务	客户端安装目录/Oozie/oozie-client-*/examples/apps/java-main
Shell任务	客户端安装目录/Oozie/oozie-client-*/examples/apps/shell
Streaming任务	客户端安装目录/Oozie/oozie-client-*/examples/apps/streaming
SubWorkflow任务	客户端安装目录/Oozie/oozie-client-*/examples/apps/subwf
SSH任务	客户端安装目录/Oozie/oozie-client-*/examples/apps/ssh

任务类型	样例目录
定时任务	客户端安装目录/Oozie/oozie-client-*/examples/apps/cron

📖 说明

其他任务样例中已包含HDFS任务样例。

样例目录下需关注文件如表21-6所示。

表 21-6 文件说明

文件名称	描述
job.properties	工作流的参数变量定义文件。
workflow.xml	工作流的规则定制文件。
lib	工作流运行依赖的jar包目录。
coordinator.xml	“cron”目录下存在，定时任务配置文件，用于设置定时策略。
oozie_shell.sh	“shell”目录下存在，提交Shell任务需要的Shell脚本文件。

步骤5 执行以下命令，编辑“job.properties”文件。

vi job.properties

修改如下内容：

更改“userName”的参数值为提交任务的人机用户名，例如
“userName=oozieuser”。

步骤6 执行**oozie job**命令，运行工作流文件。

oozie job -oozie https://oozie角色的主机名:21003/oozie -config job.properties文件所在路径 -run

例如：

oozie job -oozie https://10-1-130-10:21003/oozie -config

/opt/client/Oozie/oozie-client-*/examples/apps/map-reduce/job.properties -run

📖 说明

- 命令参数解释如下：
 - oozie: 实际执行任务的Oozie服务器URL。
 - config: 工作流属性文件。
 - run: 运行工作流。
- 执行完工作流文件，显示job id表示提交成功，例如：job: 0000021-140222101051722-oozie-omm-W。登录Oozie管理页面，查看运行情况。
使用oozieuser用户，登录Oozie WebUI页面：<https://oozie角色的ip地址:21003/oozie>。
Oozie的WebUI界面中，可在页面表格根据jobid查看已提交的工作流信息。

----结束

21.2 使用 Hue 提交 Oozie 作业

21.2.1 使用 Hue 创建工作流

操作场景

用户通过Hue管理界面可以进行提交Oozie作业，提交作业之前，首先需要创建一个工作流。

前提条件


使用Hue提交Oozie作业之前，需要提前配置好Oozie客户端，并上传样例配置文件和jar至HDFS指定目录，具体操作请参考[Oozie客户端配置说明](#)章节。

操作步骤

步骤1 准备一个具有对应组件操作权限的用户。

例如：使用admin用户登录FusionInsight Manager，选择“系统 > 用户 > 添加用户”，创建一个“人机”用户“hueuser”，并加入“hive”、“hadoop”、“supergroup”组和“System_administrator”角色，主组为“hive”。

步骤2 使用**步骤1**创建的用户登录FusionInsight Manager（首次登录需要修改密码），选择“集群 > 服务 > Hue”，单击“Hue WebUI”右侧的链接，进入Hue WebUI界面。

步骤3 在界面左侧导航栏单击，选择“Workflow”，打开Workflow编辑器。

步骤4 单击“文档”后的下拉框选择“操作”，在操作列表中选择需要创建的作业类型，将其拖到操作界面中即可。



不同类型作业提交请参考以下章节：

- [使用Hue提交Oozie Hive2作业](#)
- [使用Hue提交Oozie Spark2x作业](#)
- [使用Hue提交Oozie Java作业](#)
- [使用Hue提交Oozie Loader作业](#)
- [使用Hue提交Oozie Mapreduce作业](#)
- [使用Hue提交Oozie Sub workflow作业](#)
- [使用Hue提交Oozie Shell作业](#)
- [使用Hue提交Oozie HDFS作业](#)
- [使用Hue提交Oozie Streaming作业](#)
- [使用Hue提交Oozie Distcp作业](#)

----结束

21.2.2 使用 Hue 提交 Oozie Hive2 作业

操作场景

该任务指导用户通过Hue界面提交Hive2类型的Oozie作业。

操作步骤

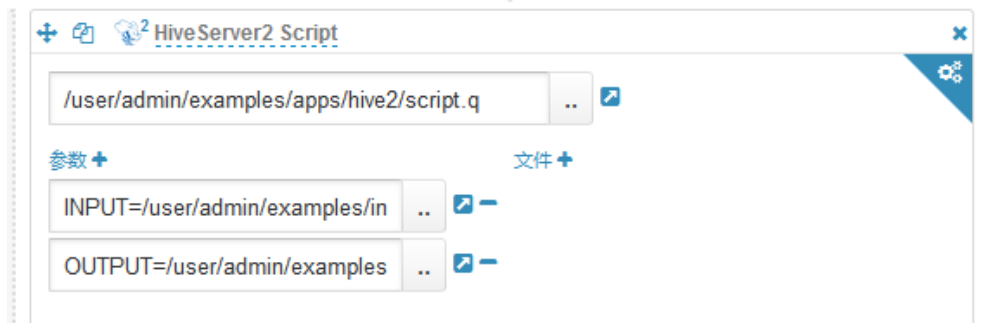
步骤1 创建工作流，请参考[使用Hue创建工作流](#)。


步骤2 在工作流编辑页面，选择“HiveServer2 脚本”按钮 ，将其拖到操作区中。

步骤3 在弹出的“HiveServer2 Script”窗口中配置HDFS上的脚本路径，例如“/user/admin/examples/apps/hive2/script.q”，然后单击“添加”。

步骤4 单击“参数+”，添加输入输出参数。

例如输入参数为“INPUT=/user/admin/examples/input-data/table”，输出参数为“OUTPUT=/user/admin/examples/output-data/hive2_workflow”。



步骤5 单击右上角的配置按钮 。在打开的配置界面中，单击“删除+”，添加删除目录，例如“/user/admin/examples/output-data/hive2_workflow”。

步骤6 配置“作业 XML”，值为“客户端安装目录/Oozie/oozie-client-*/examples/apps/hive/hive-site.xml”上传至HDFS目录中所在路径，例如“/user/admin/examples/apps/hive2/hive-site.xml”。“HiveServer2 URL”及其他参数无需配置。




说明

若以上的参数和值在使用过程中发生了修改，可在“Oozie客户端安装目录/oozie-client-*/conf/hive-site.xml”文件中查询。

步骤7 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“Hive2-Workflow”。

步骤8 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束

21.2.3 使用 Hue 提交 Oozie HQL 脚本

操作场景

该任务指导用户通过Hue界面提交Hive脚本作业。

操作步骤

步骤1 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

步骤2 在界面左侧导航栏选择“ > Workflow”，打开Workflow编辑器。


步骤3 单击“文档”，在操作列表中选择Hive脚本 ，将其拖到操作界面中。

步骤4 在弹出的“HiveServer2 Script”框中，选择之前保存的Hive脚本，关于保存Hive脚本参考[在Hue WebUI使用HiveQL编辑器](#)章节。选择脚本后单击“添加”。



步骤5 配置“作业 XML”，例如配置为hdfs路径“/user/admin/examples/apps/hive2/hive-site.xml”，配置方式参考[使用Hue提交Oozie Hive2作业](#)。

步骤6 单击Oozie编辑器右上角的 。

步骤7 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束


21.2.4 使用 Hue 提交 Oozie Spark2x 作业

操作场景

该任务指导用户通过Hue界面提交Spark2x类型的Oozie作业。

操作步骤

步骤1 创建工作流，请参考[使用Hue创建工作流](#)。

步骤2 在工作流编辑页面，选择“Spark 程序”按钮 ，将其拖到操作区中。

步骤3 在弹出的“Spark”窗口配置“Files”，例如“hdfs://hacluster/user/admin/examples/apps/spark2x/lib/oozie-examples.jar”。配置“jar/py name”，例如“oozie-examples.jar”，配置完成后单击“添加”。

步骤4 配置“Main class”的值。例如“org.apache.oozie.example.SparkFileCopy”。

步骤5 单击“参数+”，添加输入输出相关参数。


例如添加：

- “hdfs://hacluster/user/admin/examples/input-data/text/data.txt”
- “hdfs://hacluster/user/admin/examples/output-data/spark_workflow”

步骤6 在“Options list”文本框指定spark参数，例如“--conf spark.yarn.archive=hdfs://hacluster/user/spark2x/jars/8.1.0.1/spark-archive-2x.zip --conf spark.eventLog.enabled=true --conf spark.eventLog.dir=hdfs://hacluster/spark2xJobHistory2x”。

说明

此处版本号“8.1.0.1”为示例，可登录FusionInsight Manager界面，单击右上角的 ，在下拉框中单击“关于”，在弹框中查看Manager版本号。


步骤7 单击右上角的配置按钮 。配置“Spark Master”的值，例如“yarn-cluster”。配置“Mode”的值，例如“cluster”。

步骤8 在打开的配置界面中，单击“删除+”，添加删除目录，例如“hdfs://hacluster/user/admin/examples/output-data/spark_workflow”。

步骤9 单击“属性+”，添加oozie使用的sharelib，左边文本框填写属性名称“oozie.action.sharelib.for.spark”，右边文本框填写属性值“spark2x”。

步骤10 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“Spark-Workflow”。

步骤11 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束

21.2.5 使用 Hue 提交 Oozie Java 作业

操作场景

该任务指导用户通过Hue界面提交Java类型的Oozie作业。

操作步骤

步骤1 创建工作流，请参考[使用Hue创建工作流](#)。

步骤2 在工作流编辑页面，选择“Java 程序”按钮 ，将其拖到操作区中。

步骤3 在弹出的“Java program”窗口中配置“Jar name”的值，例如“/user/admin/examples/apps/java-main/lib/oozie-examples-5.1.0.jar”。配置“Main class”的值，例如“org.apache.oozie.example.DemoJavaMain”。然后单击“添加”。

步骤4 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“Java-Workflow”。

步骤5 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束

21.2.6 使用 Hue 提交 Oozie Loader 作业

操作场景

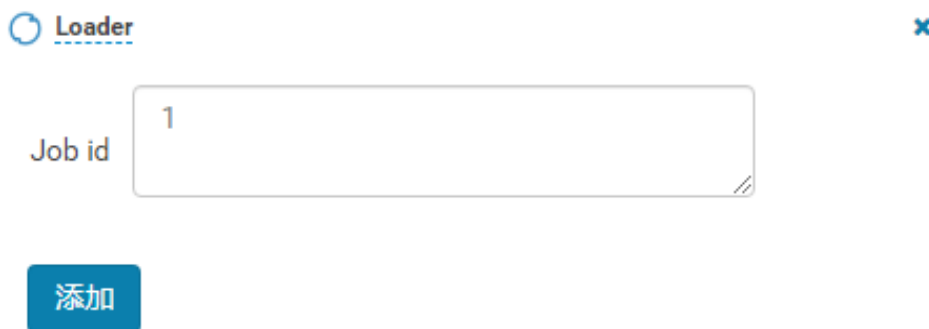
该任务指导用户通过Hue界面提交Loader类型的Oozie作业。

操作步骤

步骤1 创建工作流，请参考[使用Hue创建工作流](#)。

步骤2 在工作流编辑页面，选择“Loader”按钮 ，将其拖到操作区中。

步骤3 在弹出的“Loader”窗口中配置“Job id”的值，例如“1”。然后单击“添加”。



Loader x

Job id

添加

📖 说明

“Job id”是需要编排的Loader作业ID值，可从Loader页面获取。
创建需要调度的Loader作业，并获取该作业ID，具体操作请参见[使用Loader](#)相关章节。

步骤4 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“Loader-Workflow”。

步骤5 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束


21.2.7 使用 Hue 提交 Oozie Mapreduce 作业

操作场景

该任务指导用户通过Hue界面提交Mapreduce类型的Oozie作业。

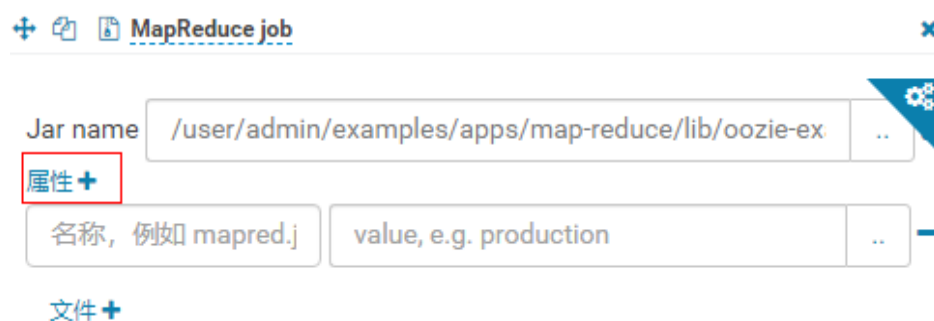
操作步骤

步骤1 创建工作流，请参考[使用Hue创建工作流](#)。


步骤2 在工作流编辑页面，选择“MapReduce 作业”按钮 ，将其拖到操作区中。

步骤3 在弹出的“MapReduce job”窗口中配置“Jar name”的值，例如“/user/admin/examples/apps/map-reduce/lib/oozie-examples-5.1.0.jar”。然后单击“添加”。

步骤4 单击“属性+”，添加输入输出相关属性。



例如配置“mapred.input.dir”的值为“/user/admin/examples/input-data/text”，配置“mapred.output.dir”的值为“/user/admin/examples/output-data/map-reduce_workflow”。

步骤5 单击右上角的配置按钮 。在打开的配置界面中，单击“删除+”，添加删除目录，例如“/user/admin/examples/output-data/map-reduce_workflow”。

步骤6 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“MapReduce-Workflow”。

步骤7 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束

21.2.8 使用 Hue 提交 Oozie Sub workflow 作业

操作场景

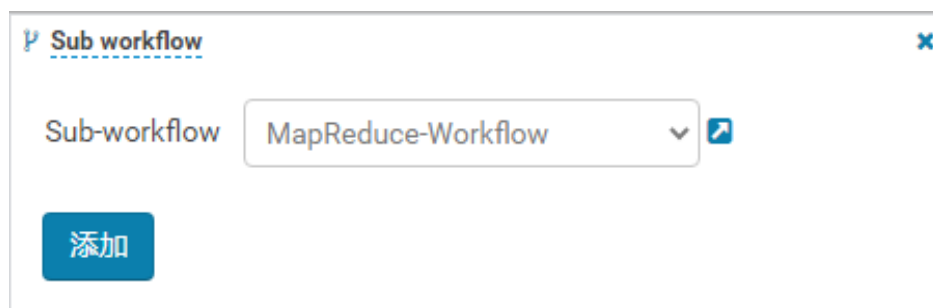
该任务指导用户通过Hue界面提交Sub Workflow类型的Oozie作业。

操作步骤

步骤1 创建工作流，请参考[使用Hue创建工作流](#)。

步骤2 在工作流编辑页面，选择“子Workflow”按钮 ，将其拖到操作区中。

步骤3 在弹出的“Sub workflow”窗口中配置“Sub-workflow”的值，例如从下拉列表中选择“Java-Workflow”（这个值是已经创建好的工作流之一），然后单击“添加”。



步骤4 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“Subworkflow-Workflow”。

步骤5 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束

21.2.9 使用 Hue 提交 Oozie Shell 作业

操作场景

该任务指导用户通过Hue界面提交Shell类型的Oozie作业。

操作步骤

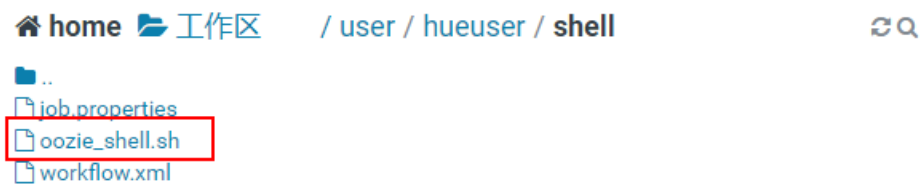
步骤1 创建工作流，请参考[使用Hue创建工作流](#)。

步骤2 在工作流编辑页面，选择“Shell”按钮 ，将其拖到操作区中。

步骤3 在弹出的“Shell”窗口中配置“Shell command”的值，例如“oozie_shell.sh”，然后单击“添加”。

步骤4 单击“文件+”，添加Shell命令执行文件或Oozie样例执行文件，可以选择存储在HDFS的文件或本地文件。

- 若文件存储在HDFS上，选择“.sh”文件所在路径即可，例如“user/hueuser/shell/oozie_shell.sh”。



上传文件

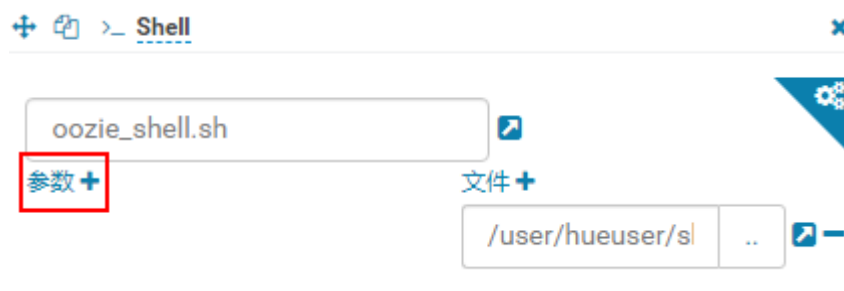
选择此文件夹

创建文件夹

- 若选择本地文件，则需在“选择文件”界面，单击“上传文件”，上传本地文件，文件上传成功后，选择该文件即可。



步骤5 如果执行的Shell文件需要传递参数，可单击“参数+”设置参数。



说明

传递参数的顺序需要和Shell脚本中保持一致。

步骤6 单击Oozie编辑器右上角的保存图标。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“Shell-Workflow”。

步骤7 保存完成后，单击提交按钮，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

说明

- 配置Shell命令为Linux指令时，请指定为原始指令，不要使用快捷键指令。例如：`ls -l`，不要配置成`ll`。可配置成Shell命令`ls`，参数添加一个“-l”。
- Windows上传Shell脚本到HDFS时，请保证Shell脚本的格式为Unix，格式不正确会导致Shell作业提交失败。

----结束


21.2.10 使用 Hue 提交 Oozie HDFS 作业

操作场景

该任务指导用户通过Hue界面提交HDFS类型的Oozie作业。

操作步骤

步骤1 创建工作流，请参考[使用Hue创建工作流](#)。

步骤2 在工作流编辑页面，选择“Fs”按钮 ，将其拖到操作区中。

步骤3 在弹出的“Fs”窗口中单击“添加”。

步骤4 单击“CREATE DIRECTORY+”，添加待创建的HDFS目录。例如“/user/admin/examples/output-data/mkdir_workflow”和“/user/admin/examples/output-data/mkdir_workflow1”。

注意

若单击了“DELETE PATH+”添加待删除的HDFS路径，该参数不能为空，否则会默认删除HDFS的“/user{提交用户名}”目录，可能会导致其他任务运行异常。

步骤5 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“HDFS-Workflow”。

步骤6 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束


21.2.11 使用 Hue 提交 Oozie Streaming 作业

操作场景

该任务指导用户通过Hue界面提交Streaming类型的Oozie作业。

操作步骤


步骤1 创建工作流，请参考[使用Hue创建工作流](#)。

步骤2 在工作流编辑页面，选择“数据流”按钮 ，将其拖到操作区中。

步骤3 在弹出的“Streaming”窗口中配置“Mapper”的值，例如“/bin/cat”。配置“Reducer”的值，例如“/usr/bin/wc”。然后单击“添加”。

步骤4 单击“文件+”，添加运行所需的文件。

例如“/user/oozie/share/lib/mapreduce-streaming/hadoop-streaming-xxx.jar”和“/user/oozie/share/lib/mapreduce-streaming/oozie-sharelib-streaming-5.1.0.jar”。


步骤5 单击右上角的配置按钮 。在打开的配置界面中，单击“删除+”，添加删除目录，例如“/user/admin/examples/output-data/streaming_workflow”。

步骤6 单击“属性+”，添加下列属性。

- 左边框填写属性名称“mapred.input.dir”，右边框填写属性值“/user/admin/examples/input-data/text”。
- 左边框填写属性名称“mapred.output.dir”，右边框填写属性值“/user/admin/examples/output-data/streaming_workflow”。

步骤7 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“Streaming-Workflow”。

步骤8 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束

21.2.12 使用 Hue 提交 Oozie Distcp 作业

操作场景

该任务指导用户通过Hue界面提交Distcp类型的Oozie作业。

操作步骤

步骤1 创建工作流，请参考[使用Hue创建工作流](#)。

步骤2 在工作流编辑页面，选择“DistCp”按钮 ，将其拖到操作区中。

步骤3 当前DistCp操作是否是跨集群操作。

- 是，执行[步骤4](#)。
- 否，执行[步骤7](#)。

步骤4 对两个集群进行跨Manager集群互信。

步骤5 在弹出的“Distcp”窗口中配置“源”的值，例如“hdfs://hacluster/user/admin/examples/input-data/text/data.txt”。配置“目标”的值，例如“hdfs://target_ip:target_port/user/admin/examples/output-data/distcp-workflow/data.txt”。然后单击“添加”。

步骤6 单击右上角的配置按钮 ，在打开的“属性”页签配置界面中，单击“属性+”，在左边文本框中填写属性名称“oozie.launcher.mapreduce.job.hdfs-servers”，在右边

文本框中填写属性值“`hdfs://source_ip:source_port,hdfs://target_ip:target_port`”，执行**步骤8**。

📖 说明


source_ip: 源集群的HDFS的NameNode的业务地址。

source_port: 源集群的HDFS的NameNode的端口号。

target_ip: 目标集群的HDFS的NameNode的业务地址。

target_port: 目标集群的HDFS的NameNode的端口号。


步骤7 在弹出的“Distcp”窗口中配置“源”的值，例如“`/user/admin/examples/input-data/text/data.txt`”。配置“目标”的值，例如“`/user/admin/examples/output-data/distcp-workflow/data.txt`”。然后单击“添加”。

步骤8 单击右上角的配置按钮 ，在打开的配置界面中，单击“删除+”，添加删除目录，例如“`/user/admin/examples/output-data/distcp-workflow`”。



步骤9 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“`My Workflow`”），可以直接单击该名称进行修改，例如“`Distcp-Workflow`”。

步骤10 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束

21.2.13 使用 Hue 提交 Oozie SSH 作业

操作场景

该任务指导用户通过Hue界面提交SSH类型的Oozie作业。

由于有安全攻击的隐患，所以默认是无法提交SSH作业的，如果想使用SSH功能，需要手动开启。

操作步骤

步骤1 开启SSH功能（若当前集群无“oozie.job.ssh.enable”参数，则跳过该操作）：


1. 在FusionInsight Manager界面，选择“集群 > 服务 > Oozie > 配置 > 全部配置 > oozie（角色） > 安全”，修改“oozie.job.ssh.enable”的值为“true”，单击“保存”，在弹出的“保存配置”界面单击“确定”，保存配置。



2. 在Oozie的“概览”界面，选择右上角“更多 > 重启服务”，重启Oozie服务。

步骤2 创建工作流，请参考[使用Hue创建工作流](#)。

步骤3 添加互信操作，请参考[配置Oozie节点间用户互信](#)。


步骤4 在工作流编辑页面，选择“Ssh”按钮 ，将其拖到操作区中。

步骤5 在弹出的“Ssh”窗口中配置以下参数并单击“添加”。

- User and Host: User为**步骤3**中配置互信的用户，参数配置格式为：*运行SSH任务的节点的用户@运行SSH任务的节点的IP地址*。例如该配置项的值可设置为：`root@x.x.x.x`。
- Ssh command: 提交作业的具体命令。

步骤6 单击Oozie编辑器右上角的 。

保存前如果需要修改作业名称（默认为“My Workflow”），可以直接单击该名称进行修改，例如“Ssh-Workflow”。

步骤7 保存完成后，单击 ，提交该作业。

作业提交后，可通过Hue界面查看作业的详细信息、日志、进度等相关内容。

----结束

21.2.14 使用 Hue 提交 Coordinator 定时调度作业

操作场景

该任务指导用户通过Hue界面提交定时调度类型的作业。

前提条件

提交Coordinator任务之前需要提前配置好相关的workflow作业。

操作步骤

步骤1 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

步骤2 在界面左侧导航栏单击，选择“计划”，打开Coordinator编辑器。

步骤3 在作业编辑界面中单击“My Schedule”修改作业的名称。


步骤4 单击“选择Workflow...”选择需要编排的Workflow。

My Schedule

[添加描述...](#)

要计划哪个 Workflow?

[选择 Workflow...](#)


步骤5 选择好Workflow，根据界面提示设置作业执行的频率，如果执行的Workflow需要传递参数，可单击“+添加参数”设置参数，然后单击右上角的保存作业。

说明

因时区转化的原因，此处时间有可能会与当地系统实际时间差异几个小时。比如在中国，此处的时间则会比当地时间晚8个小时。

如果需要统一同步配置为上海时间，操作如下：

1. 在Manager页面，选择“集群 > 服务 > Oozie > 配置 > 全部配置”，修改oozie的服务配置参数“oozie.processing.timezone”值为“GMT+0800”（修改配置需要重启服务生效）。
2. 在Oozie编辑器页面，提交Coordinator定时调度任务时，单击频率下方的“选项”按钮。弹出下拉选项，在“时区”中选择“Asia/Shanghai”。

步骤6 单击编辑器右上角的，设置定时任务执行的时间范围的起始值与结束值，然后单击“提交”提交作业。

说明

因时区转化的原因，此处时间有可能会与当地系统实际时间差异几个小时。

----结束

21.2.15 使用 Hue 提交提交 Bundle 批处理作业

操作场景

当同时存在多个定时任务的情况下，用户可以通过Bundle任务进行批量管理作业。该任务指导用户通过Hue界面提交批量类型的作业。

前提条件

提交Bundle批处理之前需要提前配置好相关的Workflow和Coordinator作业。


操作步骤



步骤1 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

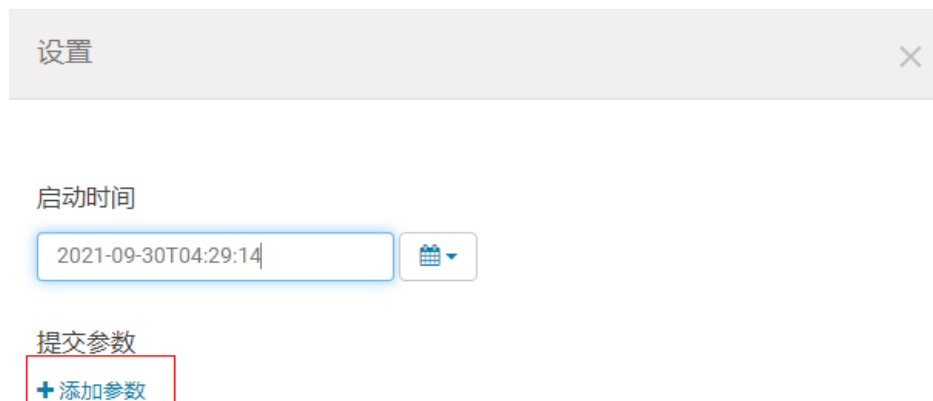
步骤2 在界面左侧导航栏单击，选择“Bundle”，打开Bundle编辑器。

步骤3 在作业编辑界面中单击“My Bundle”修改作业的名称。

步骤4 单击“+添加Coordinator”选择需要编排的Coordinator作业。


步骤5 根据界面提示设置Coordinator任务调度的开始、结束时间，然后单击右上角的保存作业。

步骤6 单击编辑器右上角的，在弹出菜单选择，设置Bundle任务的启动时间，根据实际需求单击“+添加参数”设置提交参数，然后关闭对话框保存设置。



说明

因时区转化的原因，此处时间有可能会与当地系统实际时间差异数个小时。比如在中国，此处的时间则会比当地时间晚8个小时。

步骤7 单击编辑器右上角的，在弹出的确认界面中单击“提交”提交作业。

----结束


21.2.16 在 Hue 界面中查询 Oozie 作业结果

操作场景

提交作业后，可以通过Hue界面查看具体作业的执行情况。

操作步骤

步骤1 访问Hue WebUI，请参考[访问Hue WebUI界面](#)。

步骤2 单击菜单左侧的，在打开的页面中可以查看Workflow、计划、Bundles任务的相关信息。

默认显示当前集群的所有作业。

说明

作业浏览器显示的数字表示集群中所有作业的总数。

“作业浏览器”将显示作业以下信息：

表 21-7 MRS 作业属性介绍

属性名	描述
名称	表示作业的名称。
用户	表示启动该作业的用户。
类型	表示作业的类型。
状态	表示作业的状态，包含“成功”、“正在运行”、“失败”。
进度	表示作业运行进度。
组	表示作业所属组。
开始	表示作业开始时间。
持续时间	表示作业运行使用的时间。
Id	表示作业的编号，由系统自动生成。

说明

如果MRS集群安装了Spark组件，则默认会启动一个作业“Spark-JDBCServer”，用于执行任务。

----结束

21.2.17 配置 Oozie 节点间用户互信

操作场景

在使用Oozie节点通过SSH作业执行外部节点的Shell，需要单向免密互信时，可以参考此示例。

前提条件

已经安装Oozie，而且能与外部节点（SSH连接的节点）通信。

操作步骤

步骤1 在外部节点上确保连接SSH时使用的用户存在，且该用户“~/ssh”目录存在。

步骤2 使用omm用户登录Oozie所在节点，查看“~/ssh/id_rsa.pub”文件是否存在。

- 是，执行**步骤3**。
- 否，执行以下命令生成公私钥：

```
ssh-keygen -t rsa
```

步骤3 以omm用户登录oozie实例所在节点，执行以下命令配置互信：

```
ssh-copy-id -i ~/.ssh/id_rsa.pub 运行SSH任务的用户@运行SSH任务的节点的IP地址
```

执行该命令后需要输入运行SSH任务的用户的密码。

说明

- Shell所在节点（外部节点）的账户需要有权限执行Shell脚本并对于所有Shell脚本里涉及到的所有目录文件有足够权限。
- 如果Oozie具有多个节点，需要在所有Oozie节点执行**步骤2~步骤3**。

步骤4 使用omm用户登录依次其他Oozie所在节点，重复执行**步骤2-步骤3**。

----结束

21.3 开启 Oozie HA 机制

操作场景

Oozie多个节点同时提供服务的时候，通过ZooKeeper来提供高可用（HA）功能，防止单节点故障以及多节点同时处理一个任务。

说明

本章节内容仅适用于MRS 3.1.2及之后版本。

对系统影响

操作过程中需要重启Oozie服务。重启过程中，Oozie服务无法提供服务。

前提条件

- 已安装Oozie、ZooKeeper服务，且服务正常运行。
- 没有任务正在运行。
- 如果当前集群不是安装最新的版本包，需要从“\$BIGDATA_HOME/FusionInsight_Porter_x.x.x/install/FusionInsight-Oozie-x.x.x/oozie-x.x.x/embedded-oozie-server/webapp/WEB-INF/lib”路径拷贝“curator-x-discovery-x.x.x.jar”包到“\$BIGDATA_HOME/FusionInsight_Porter_x.x.x/install/FusionInsight-Oozie-x.x.x/oozie-x.x.x/lib”目录下。

操作步骤

步骤1 在FusionInsight Manager界面选择“集群 > 服务 > Oozie > 配置 > 全部配置”，在“自定义”的“oozie.site.configs”参数中添加如下四个配置项。修改完成后单击“保存”，在弹框中单击“确定”保存配置。

名称	值	参数说明
oozie.services.ext	org.apache.oozie.service.ZKLocksService,org.apache.oozie.service.ZKXLogStreamingService,org.apache.oozie.service.ZKJobsConcurrencyService,org.apache.oozie.service.ZKUUIDService	HA启用的功能
oozie.zookeeper.connection.string	ZooKeeper实例的业务IP:端口（多个地址以逗号隔开）	ZooKeeper连接信息
oozie.zookeeper.namespace	oozie	Oozie在ZooKeeper的路径
oozie.zookeeper.secure	安全集群：true 普通集群：无需配置该参数	ZooKeeper是否启用kerberos

步骤2 在Oozie的“概览”界面，选择右上角“更多 > 重启服务”，重启Oozie集群。

----结束

21.4 Oozie 日志介绍

日志描述

日志路径： Oozie相关日志的默认存储路径为：

- 运行日志：“/var/log/Bigdata/oozie”。
- 审计日志：“/var/log/Bigdata/audit/oozie”。

日志归档规则： Oozie的日志分三类：运行日志、脚本日志和审计日志。运行日志每个文件最大20M，最多20个。审计日志每个文件最大20M，最多20个。

说明

“oozie.log”日志每小时生成一个日志压缩文件，默认保留720个（一个月的日志）。

表 21-8 Oozie 日志列表

日志类型	日志文件名	描述
运行日志	jetty.log	Oozie内置jetty服务器日志，处理OozieServlet的request/response信息
	jetty.out	Oozie进程启动日志
	oozie_db_temp.log	Oozie数据库连接日志
	oozie-instrumentation.log	Oozie仪表盘日志，主要记录Oozie运行状态，各组件的配置信息
	oozie-jpa.log	openJPa运行日志
	oozie.log	Oozie运行日志
	oozie-<SSH_USER>-<DATE>-<PID>-gc.log	Oozie服务垃圾回收日志
	oozie-ops.log	Oozie操作日志
	check-serviceDetail.log	Oozie健康检查日志
	oozie-error.log	Oozie运行错误日志
	threadDump-<DATE>.log	记录服务进程正常退出时堆栈信息的日志
脚本日志	postinstallDetail.log	安装后启动前的工作日志
	prestartDetail.log	预启动日志
	startDetail.log	服务启动日志
	stopDetail.log	服务停止日志
	upload-sharelib.log	sharelib上传操作日志
审计日志	oozie-audit.log	审计日志

日志级别

Oozie中提供了如表21-9所示的日志级别。

日志级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 21-9 日志级别

级别	描述
ERROR	ERROR表示错误日志，可能会导致进程异常。
WARN	WARN表示当前事件处理存在异常信息。

级别	描述
INFO	INFO表示系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及数据库底层数据传输的信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 登录FusionInsight Manager系统。
- 步骤2** 选择“集群 > 待操作集群的名称 > 服务 > Oozie > 配置”。
- 步骤3** 选择“全部配置”。
- 步骤4** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤5** 选择所需修改的日志级别。
- 步骤6** 单击“保存”，单击“确定”，处理结束后生效。

----结束

日志格式

Oozie的日志格式如下所示。

表 21-10 日志格式

日志类型	格式	示例
运行日志	<yyyy-MM-dd HH:mm:ss,SSS><Log Level><日志事件的发生位置> <log中的message>	2015-05-29 21:01:45,268 INFO StatusTransitService \$StatusTransitRunnable:539 - USER[-] GROUP[-] Released lock for [org.apache.oozie.service.StatusTransitService]
脚本日志	<yyyy-MM-dd HH:mm:ss,SSS><主机名>><Log Level><log中的message>	2015-06-01 17:18:03 001 suse11-192-168-0-111 oozie INFO Running oozie service check script
审计日志	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <线程名称> <log中的message> <日志事件的发生位置>	2015-06-01 22:38:41,323 INFO http-bio-21003-exec-8 IP [192.168.0.111] USER [null], GROUP [null], APP [null], JOBID [null], OPERATION [null], PARAMETER [null], RESULT [SUCCESS], HTTPCODE [200], ERRORCODE [null], ERRORMESSAGE [null] org.apache.oozie.util.XLog.log(XLog.java: 539)

21.5 Oozie 常见问题

21.5.1 Oozie 定时任务没有准时运行如何处理

问题

在Hue或者Oozie客户端设置执行Coordinator定时任务，但没有准时执行。

回答

设置任务时，需要使用UTC时间。

例如在“job.properties”中配置“start=2016-12-20T09:00Z”。

修改配置后重新启动定时任务即可。

21.5.2 HDFS 上更新了 oozie 的 share lib 目录但没有生效

问题

在HDFS的“/user/oozie/share/lib”目录上传了新的jar包，但执行任务时仍然报找不到类的错误。

回答

在客户端执行如下命令刷新目录：

```
oozie admin -oozie https://xxx.xxx.xxx.xxx:21003/oozie -sharelibupdate
```

21.5.3 Oozie 作业执行失败常用排查手段

1. 根据任务在Yarn上的任务日志排查，首先把实际的运行任务，比如Hive SQL通过beeline运行一遍，确认Hive无问题。
2. 出现“classnotfoundException”等报错，排查“/user/oozie/share/lib”路径下各组件有没有报错的类的Jar包，如果没有，添加Jar包并执行[HDFS上更新了oozie的share lib目录但没有生效](#)。如果执行了更新“share lib”目录依然报找不到类，那么可以查看执行更新“share lib”的命令打印出来的路径“sharelibDirNew”是否是“/user/oozie/share/lib”，一定不能是其它目录。

```
[root@host-... client]#
[root@host-... client]# oozie admin -oozie https://host-...:21003/oozie/ -sharelibupdate
INFO CMD-admin -oozie https://host-...:21003/oozie/ -sharelibupdate
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/client/Oozie/oozie-client-.../lib/slf4j-log4j12-1.7.30.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/client/Oozie/oozie-client-.../lib/slf4j-simple-1.7.30.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
[ShareLib update status]
sharelibDirOld = /user/oozie/share/lib
host = https://@:21003/oozie
sharelibDirNew = /user/oozie/share/lib
status = Successful
```

3. 出现NosuchMethodError，排查“/user/oozie/share/lib”路径下各组件的Jar包是不是有多个版本，注意业务本身上传的Jar包冲突，可通过Oozie在Yarn上的运行日志打印的加载的Jar包排查是否有Jar包冲突。
4. 自研代码运行异常，可以先运行Oozie的自带样例，排除Oozie自身的异常。

5. 寻求技术人员的支持，需要收集Yarn上Oozie任务运行日志、Oozie自身的日志及组件的运行的日志，例如使用Oozie运行Hive报异常，需收集Hive的日志。

22 使用 Presto

22.1 访问 Presto 的 WebUI

用户可以通过Presto的WebUI，在图形化界面查看Presto的统计信息。Presto的WebUI界面不支持使用IE浏览器访问，建议使用Google浏览器访问。

前提条件

- 已安装Presto服务的集群。
- 已安装集群客户端，例如安装目录为“/opt/client”。以下操作的客户端目录只是举例，请根据实际安装目录修改。

访问 Presto 的 WebUI

- 方法一（适用于MRS 3.x及之后版本）：
 - a. 登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务”。
 - b. 选择“Presto”并在“基本信息”的“Coordinator WebUI”中单击“Coordinator(Coordinator)”，打开Presto的WebUI页面。

图 22-1 Coordinator WebUI



- 方法二（适用于MRS 3.x之前版本）：
 - a. 登录MRS Manager页面，选择“服务管理”。

- b. 选择“Presto”并在“Presto 概述”的“Presto WebUI”中单击“Coordinator (主)”，打开Presto的WebUI页面。

图 22-2 Presto WebUI



说明

第一次访问Presto WebUI，需要在浏览器中添加站点信任以继续打开页面。

- 方法三（适用于MRS 1.9.2及之后版本）：
 - a. 在集群列表页面，单击集群名称，登录集群详情页面，选择“组件管理”。

说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- b. 选择“Presto”并在“Presto 概述”的“Presto WebUI”中单击“Coordinator (主)”，打开Presto的WebUI页面。

图 22-3 Presto WebUI



22.2 使用 Presto 客户端执行查询语句

用户可以根据业务需要，在MRS集群的客户端中进行交互式查询。启用Kerberos认证的集群，需要提交拓扑的用户属于“presto”组。

MRS 3.x版本Presto组件暂不支持开启Kerberos认证。

前提条件

- 获取用户“admin”账号密码。“admin”密码在创建MRS集群时由用户指定。
- 已刷新客户端。
- 3.x版本的集群需要手动安装Presto客户端。

操作步骤

步骤1 启用Kerberos认证的集群，登录MRS Manager页面，创建拥有“Hive Admin Privilege”权限的角色，创建角色请参考[创建角色](#)。

步骤2 创建属于“Presto”和“Hive”组的用户，同时为该用户绑定**步骤1**中创建的角色，然后下载用户认证文件，参见[创建用户](#)，[下载用户认证文件](#)。

步骤3 将下载的用户.keytab文件和krb5.conf上传到MRS客户端所在节点。

说明

步骤**步骤2-步骤3**仅启用Kerberos认证的集群执行，普通集群请直接从步骤**步骤4**开始执行。

步骤4 根据业务情况，准备好客户端，并登录安装客户端的节点。

例如在Master2节点更新客户端，则登录该节点使用客户端，具体参见[使用MRS客户端](#)。

步骤5 执行以下命令切换用户。


```
sudo su - omm
```

步骤6 执行以下命令，切换到客户端目录，例如“/opt/client”。

```
cd /opt/client
```

步骤7 执行以下命令，配置环境变量。

```
source bigdata_env
```

步骤8 连接Presto Server。根据客户端的不同，提供如下两种客户端的连接方式。

- 使用MRS提供的客户端。

- 未启用Kerberos认证的集群，执行以下命令连接本集群的Presto Server。

```
presto_cli.sh
```

- 未启用Kerberos认证的集群，执行以下命令连接其他集群的Presto Server，其中ip为对应集群的Presto的浮动IP（可通过在Presto配置项中搜索PRESTO_COORDINATOR_FLOAT_IP的值获得），port为Presto Server的端口号，默认为7520。

```
presto_cli.sh --server http://ip:port
```

- 启用Kerberos认证的集群，执行以下命令连接本集群的Presto Server。

```
presto_cli.sh --krb5-config-path krb5.conf文件路径 --krb5-principal 用户principal --krb5-keytab-path user.keytab文件路径 --user presto用户名
```

- 启用Kerberos认证的集群，执行以下命令连接其他集群的Presto Server，其中ip为对应集群的Presto的浮动IP（可通过在Presto配置项中搜索PRESTO_COORDINATOR_FLOAT_IP的值获得），port为Presto Server的端口号，默认为7521。

```
presto_cli.sh --krb5-config-path krb5.conf文件路径 --krb5-principal 用户principal --krb5-keytab-path user.keytab文件路径 --server https://ip:port --krb5-remote-service-name Presto Server name
```

- 使用原生客户端

Presto原生客户端为客户端目录下的Presto/presto/bin/presto，使用方式参见<https://prestodb.io/docs/current/installation/cli.html>和<https://prestodb.io/docs/current/security/cli.html>。

步骤9 执行Query语句，如“show catalogs;”，更多语句请参阅<https://prestodb.io/docs/current/sql.html>。

📖 说明

启用Kerberos认证的集群使用Presto查询Hive Catalog的数据时，运行Presto客户端的用户需要有Hive表的访问权限，并且需要在Hive beeline中执行命令**grant all on table [table_name] to group hive;**，给Hive组赋权限。

步骤10 查询结束后，执行以下命令退出客户端。

```
quit;
```

```
----结束
```

22.3 Presto 常见问题

22.3.1 Presto 配置多 Hive 连接

用户问题

Presto如何配置多Hive连接。

处理步骤

步骤1 将目标Hive集群的core-site.xml，hdfs-site.xml文件复制分发到Presto集群上，放置在omm用户有读权限的路径下（如/home/omm），将文件属主改为omm:wheel，文件权限改为750。

步骤2 进入Presto服务配置页面：

- MRS 1.8.10及之前版本，登录MRS Manager页面，具体请参见[访问MRS Manager](#)，然后选择“服务管理 > Presto > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。
- MRS 1.8.10之后及2.x版本，单击集群名称，登录集群详情页面，选择“组件管理 > Presto > 服务配置”，单击“基础配置”下拉菜单，选择“全部配置”。

说明

如果集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，然后选择“集群 > 待操作的集群名称 > 服务 > Presto > 配置 > 全部配置”。

步骤3 在搜索框中搜索“connector-customize”。

步骤4 添加名为myhive的connector。

在connector-customize中添加配置：

```
myhive.connector.name=hive-hadoop2
```

```
myhive.hive.metastore.uri=thrift://{目标集群的Hive Metastore IP}:{目标集群的Hive metastore端口} (多个MetaStore地址间用逗号分隔，如：thrift://192.0.2.1:9083,thrift://192.0.2.2:9083)。
```

```
myhive.hive.config.resources=/etc/hadoop/conf/core-site.xml,/etc/hadoop/conf/hdfs-site.xml。(步骤1中放置配置文件的绝对路径)
```

其他配置参考：<https://trino.io/docs/333/connector/hive.html#hive-configuration-properties>

添加配置时，在配置名前增加connector名作为前缀，此处为myhive。

步骤5 修改完成后，保存配置，重启Presto服务。

说明

使用presto客户端或jdbc、UI工具连接Presto服务，执行命令**show schemas from myhive**；可以看到目标集群上的database即为配置成功。

----结束

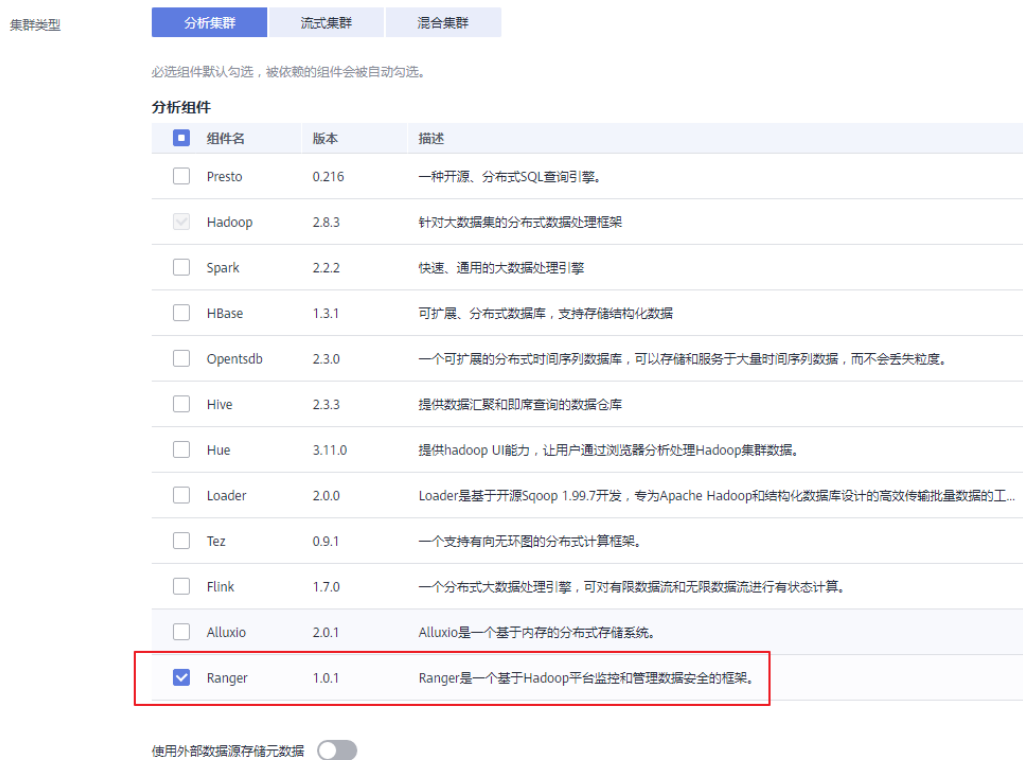
23 使用 Ranger（MRS 1.9.2）

23.1 创建 Ranger 集群

步骤1 参考[购买自定义集群](#)创建集群，组件选择时勾选Ranger组件。

目前MRS 1.9.2集群仅普通模式集群支持Ranger组件，开启Kerberos认证的安全集群不支持Ranger组件。

图 23-1 选择 Ranger 组件



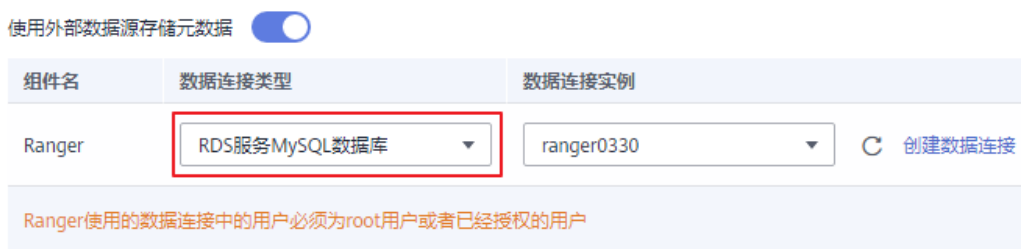
步骤2 选择是否开启“使用外部数据源存储元数据”功能。

- 开启：使用外置的MySQL数据库存储Ranger组件的User/Group/Policy等数据。

- 关闭：Ranger组件的User/Group/Policy等数据默认存放在当前集群本地数据库中。

步骤3 当“使用外部数据源存储元数据”开启时，选择数据连接类型为“RDS服务MySQL数据库”，数据连接实例选择已创建的数据连接实例，或单击“创建数据连接”新创建一个数据连接。

图 23-2 使用 RDS 服务 MySQL 数据库



说明

当用户选择的数据连接为“RDS服务MySQL数据库”时，请确保使用的数据库用户为root用户。如果为非root用户，需要先以root用户登录到数据库执行如下SQL命令为该数据库用户进行赋权，其中\${db_name}与\${db_user}为用户新建数据连接时输入的数据库名与用户名。

```
grant select on mysql.user to ${db_user};
grant all privileges on ${db_name}.* to '${db_user}'@'%' with grant option;
grant reload on *.* to '${db_user}'@'%' with grant option;
flush privileges;
```

步骤4 继续参考[购买自定义集群](#)配置其他参数并创建集群。

说明

- 在集群创建完成后，此时Ranger不会对用户访问Hive和HBase组件的权限进行控制。
- 使用Ranger管理各组件权限时，如管理hive表权限，在管理控制台或者客户端提交hive作业（操作hive数据表），可能会提示当前用户没有权限，需要在Ranger中给提交作业的用户配置具体数据库或者表权限，以免影响用户使用提交作业功能，具体请参考[在Ranger中配置Hive/Impala的访问权限](#)或[在Ranger中配置HBase的访问权限](#)页面的添加策略步骤。

----结束

23.2 访问 Ranger WebUI 及同步 Unix 用户到 Ranger WebUI

用户可以通过Ranger WebUI，在图形化界面上对Ranger进行管理。

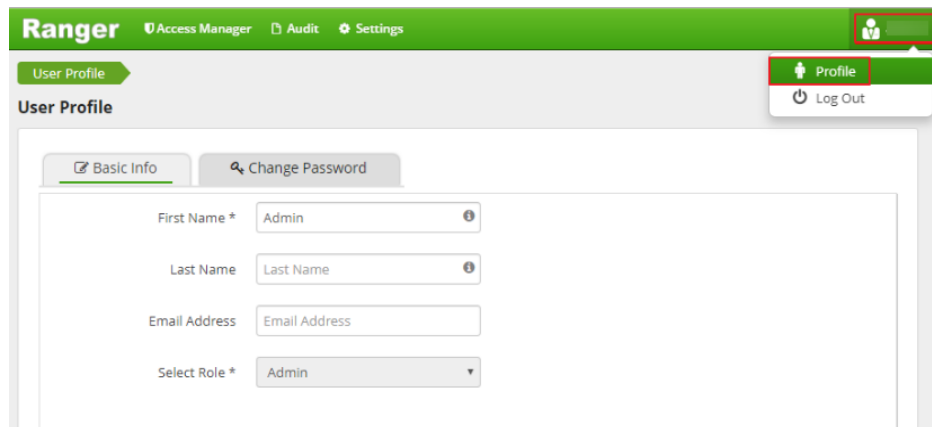
访问 Ranger Admin WebUI

- 步骤1** 在MRS控制台，单击集群名称进入集群详情页面。
- 步骤2** 选择“组件管理”。
- 步骤3** 选择“Ranger”，在“Ranger 概述”中单击“Ranger WebUI”对应的“RangerAdmin”。
- 步骤4** 进入Ranger WebUI登录界面，MRS 1.9.2版本集群默认用户名/默认密码为admin/admin@12345，MRS 1.9.3版本集群默认用户名/默认密码为admin/ranger@A1!。

首次登录Ranger WebUI界面后请修改用户密码并妥善保存。

步骤5 用户可以单击右上角的用户名，选择下拉菜单中的“Profile”，并选择“Change Password”修改用户密码。

图 23-3 修改 Ranger WebUI 登录密码



步骤6 修改完用户密码后，单击右上角用户名，选择下拉菜单中的“Log Out”，并使用新的密码重新进行登录。

----结束

使用 Ranger UserSync 同步集群节点上的 Unix 操作系统用户

Ranger UserSync是Ranger中一个重要的组件，它支持将Unix系统用户或LDAP用户同步到Ranger WebUI中，目前MRS服务只支持同步Ranger UserSync进程所在节点上的Unix用户。

步骤1 登录到UserSync进程所在的节点。

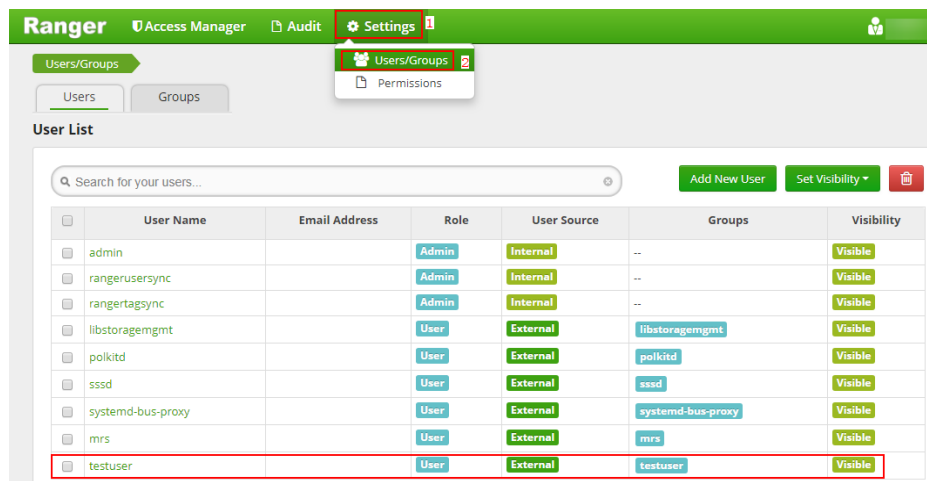
步骤2 执行useradd命令新增系统用户，例如“testuser”。

图 23-4 新增系统用户 testuser

```
[root@node-master1aHRf ~]# useradd testuser
[root@node-master1aHRf ~]# passwd testuser
Changing password for user testuser.
New password:
Retype new password:
passwd: all authentication tokens updated successfully.
```

步骤3 用户添加完成后等待1分钟左右，登录到Ranger WebUI，即可查看到该用户已经同步成功。

图 23-5 用户同步完成



---结束

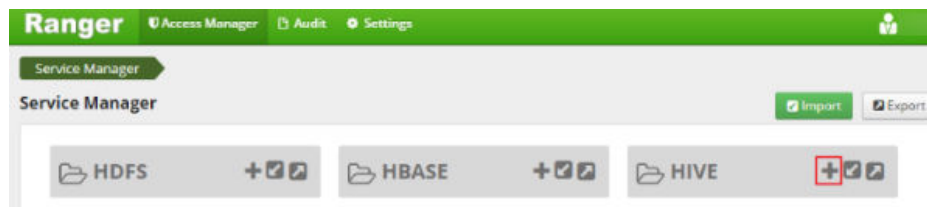
23.3 在 Ranger 中配置 Hive/Impala 的访问权限

在创建完安装了Ranger组件的MRS集群后，Hive/Impala的权限控制暂未集成在Ranger中，由于Hive与Impala配置方法一致，本章节主要介绍Hive组件如何集成在Ranger中。

步骤1 登录Ranger WebUI界面。

步骤2 在“Service Manager”中的HIVE处，单击+添加Hive Service。

图 23-6 添加 Hive Service



步骤3 请参考表23-1填写添加Hive Service的相关参数，未在表中列出的参数请保持默认值。

表 23-1 参数说明

参数	说明	示例值
Service Name	创建的服务名称，固定填写：hivedev。	hivedev
Username	根据实际需要填写。	admin
Password	根据实际需要填写。	-

参数	说明	示例值
jdbc.driverClassName	连接Hive的驱动类，固定填写： org.apache.hive.jdbc.HiveDriver。	org.apache.hive.jdbc.HiveDriver
jdbc.url	连接Hive的URL，格式为ZooKeeper Mode： jdbc:hive2://<host>:2181/;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2 其中<host>为ZooKeeper地址，ZooKeeper地址可通过登录MRS Manager然后选择“服务管理 > ZooKeeper > 实例”，查看ZooKeeper实例的“管理IP”地址获取。	jdbc:hive2://xx.xx.xx.xx:2181,xx.xx.xx.xx:2181,xx.xx.xx.xx:2181/;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2

图 23-7 Create hivedev

步骤4 单击“Add”添加服务。

步骤5 启动Ranger Hive Plugin，授权Ranger管理Hive。

1. 在MRS控制台，单击集群名称进入集群详情页面。
2. 选择“组件管理”。
3. 选择“Hive > 服务配置”，将“基础配置”切换为“全部配置”。
4. 搜索配置“hive.security.authorization”，修改如下两个配置：
 - hive.security.authorization.enabled = true
 - hive.security.authorization.manager = org.apache.ranger.authorization.hive.authorizer.RangerHiveAuthorizerFactory
5. 单击“保存配置”，并勾选“重新启动受影响的服务或实例。”重启Hive服务。

步骤6 添加访问控制策略，即Policy。

1. 登录Ranger WebUI界面。
2. 在HIVE区域单击已添加的服务名称“hivedev”。
3. 单击“Add New Policy”，新增访问控制策略。

4. 参考表23-2配置参数，未在表中列出的参数请保持默认值。

表 23-2 参数说明

参数	说明	示例值
Policy Name	策略名称。	Policy001
database	该策略允许访问的数据库名称。	test
table	该策略允许访问的数据库对应的表名称。	table1
Hive Column	该策略允许访问的数据库对应的表的列名。	name
Allow Conditions	<ul style="list-style-type: none"> - Select Group: 该策略允许访问的用户组。 - Select User: 该策略允许访问的用户组中的用户。 - Permissions: 该策略允许用户使用的权限。 	<ul style="list-style-type: none"> - Select Group: testuser - Select User: testuser - Permissions: Create 和select

图 23-8 新增 hivedev 的访问控制策略

The screenshot shows the 'Create Policy' interface in the Ranger web console. The breadcrumb navigation is 'Service Manager > hivedev Policies > Create Policy'. The main heading is 'Create Policy'. Under 'Policy Details', the 'Policy Type' is set to 'Access'. The 'Policy Name' is 'Policy001' with an 'enabled' toggle. The 'database' is 'test' with an 'include' toggle. The 'table' is 'table1' with an 'include' toggle. The 'Hive Column' is 'name' with an 'include' toggle. 'Audit Logging' is set to 'YES'. There is a 'Description' text area and a 'Policy Label' field. Under 'Allow Conditions', there is a table with columns: 'Select Group', 'Select User', 'Permissions', and 'Delegate Admin'. The first row contains 'testuser' in both 'Select Group' and 'Select User' fields, 'Create select' in the 'Permissions' field, and a 'Delegate Admin' checkbox. There are '+' and '-' buttons for adding and removing rows.

- 单击“Add”，完成策略添加，依据如上Policy示例，testuser用户组中的testuser用户将对Hive的“test”数据库中的表“table1”的“name”列有Create和select的权限，而对于其他列则没有任何的访问权限。

步骤7 参见[快速使用Hive进行数据分析](#)登录Hive客户端，验证Ranger是否已经完成集成Hive。

- 以客户端安装用户登录客户端安装节点，执行如下命令，进入hive beeline。

```
source /opt/client/bigdata_env
```

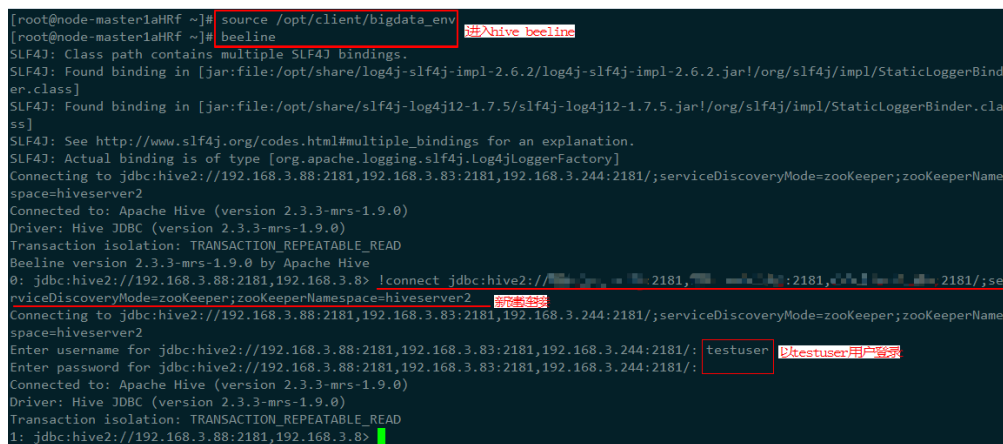
```
beeline
```

- 执行如下命令，建立连接并使用testuser登录。

```
!connect jdbc:hive2://
```

```
xx.xx.xx.xx:2181,xx.xx.3.81:2181,192.168.3.153:2181/;serviceDiscoveryMode  
=zooKeeper;zooKeeperNamespace=hiveserver2
```

图 23-9 登录 Hive



```
[root@node-master1aHRF ~]# source /opt/client/bigdata_env
[root@node-master1aHRF ~]# beeline
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/share/log4j-slf4j-impl-2.6.2/log4j-slf4j-impl-2.6.2.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/share/slf4j-log4j12-1.7.5/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Connecting to jdbc:hive2://192.168.3.88:2181,192.168.3.83:2181,192.168.3.244:2181/;serviceDiscoveryMode=zooKeeper;zooKeeperName
space=hiveserver2
Connected to: Apache Hive (version 2.3.3-mrs-1.9.0)
Driver: Hive JDBC (version 2.3.3-mrs-1.9.0)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 2.3.3-mrs-1.9.0 by Apache Hive
0: jdbc:hive2://192.168.3.88:2181,192.168.3.83:2181,192.168.3.244:2181/;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2
Connecting to jdbc:hive2://192.168.3.88:2181,192.168.3.83:2181,192.168.3.244:2181/;serviceDiscoveryMode=zooKeeper;zooKeeperName
space=hiveserver2
Enter username for jdbc:hive2://192.168.3.88:2181,192.168.3.83:2181,192.168.3.244:2181/: testuser
Enter password for jdbc:hive2://192.168.3.88:2181,192.168.3.83:2181,192.168.3.244:2181/:
```

- 查询数据，验证Ranger是否已经集成成功。

图 23-10 验证 Ranger 集成 Hive

```

1: jdbc:hive2://192.168.3.88:2181,192.168.3.88> select * from table1;
Error: Error while compiling statement: FAILED: HiveAccessControlException Permission denied: user [testuser] does not have [SELECT] privilege on [test/table1/*] (state=42000,code=40000)
1: jdbc:hive2://192.168.3.88:2181,192.168.3.88> select name from table1;
INFO : State: Compiling.
INFO : Compiling command(queryId=omm_20191204095459_5713dad6-1fff-4a98-96f8-0da351c63df9): select name from table1
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name=name, type:string, comment:null)], properties:null)
INFO : EXPLAIN output for queryId omm_20191204095459_5713dad6-1fff-4a98-96f8-0da351c63df9 : STAGE DEPENDENCIES:
  Stage-0 is a root stage [FETCH]

STAGE PLANS:
  Stage: Stage-0
    Fetch Operator
      limit: -1
      Processor Tree:
        TableScan
          alias: table1
          GatherStats: false
          Select Operator
            expressions: name (type: string)
            outputColumnNames: _col0
            ListSink

INFO : Completed compiling command(queryId=omm_20191204095459_5713dad6-1fff-4a98-96f8-0da351c63df9); Time taken: 0.12 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : State: Executing.
INFO : Executing command(queryId=omm_20191204095459_5713dad6-1fff-4a98-96f8-0da351c63df9): select name from table1
INFO : Completed executing command(queryId=omm_20191204095459_5713dad6-1fff-4a98-96f8-0da351c63df9); Time taken: 0.001 seconds
INFO : OK
+-----+
| name |
+-----+
+-----+
No rows selected (0.167 seconds)
    
```

----结束

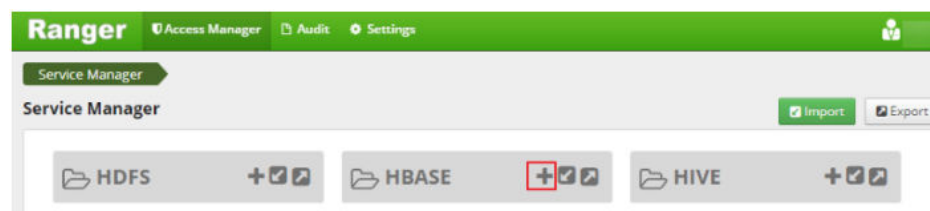
23.4 在 Ranger 中配置 HBase 的访问权限

在创建完安装了Ranger组件的MRS集群后，HBase的权限控制暂未集成在Ranger，本章节主要介绍HBase组件如何集成在Ranger中。

步骤1 登录Ranger WebUI界面。

步骤2 在“Service Manager”中的HBASE处，单击  添加HBase Service。

图 23-11 添加 HBase Service



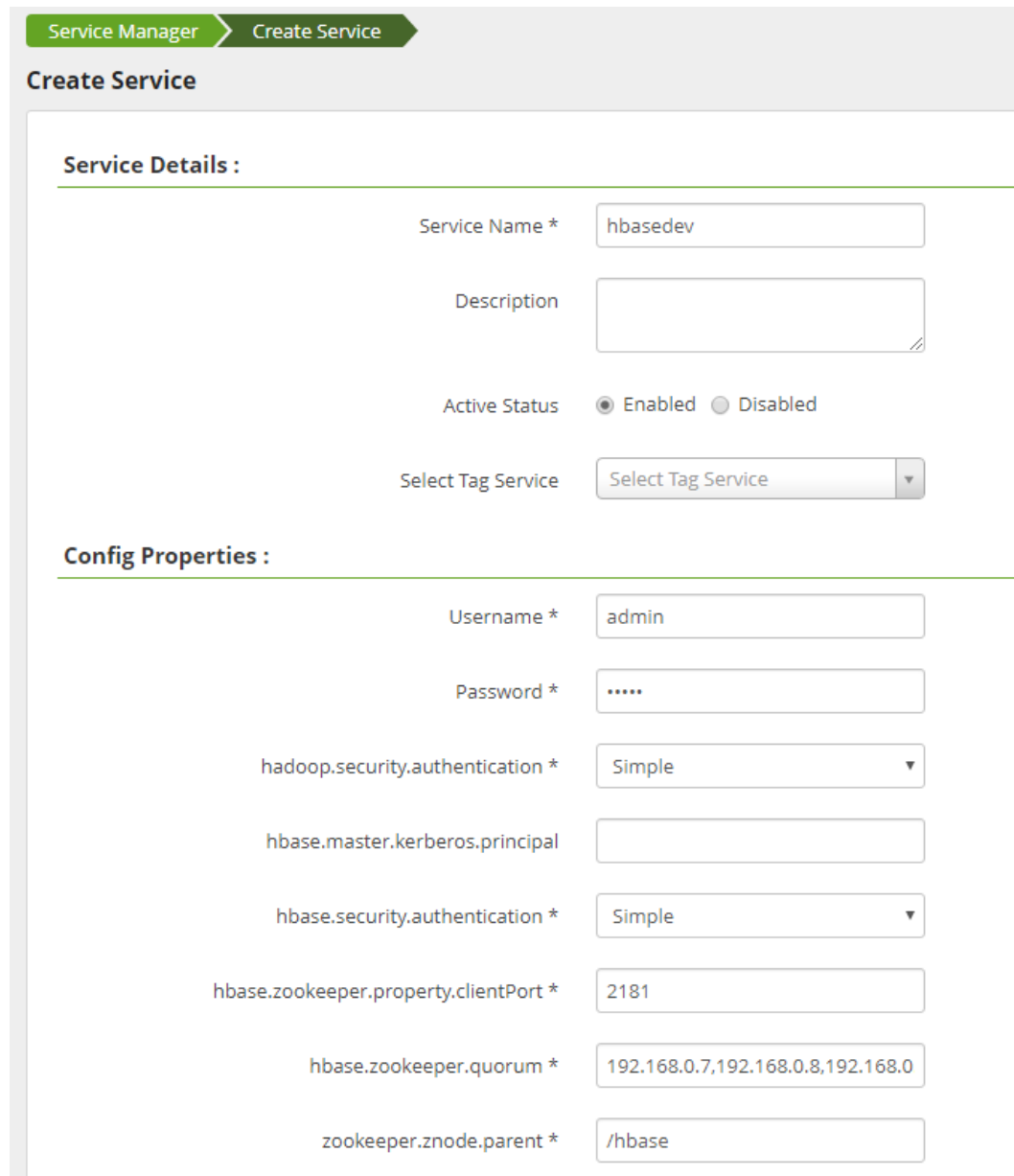
步骤3 请参考表23-3填写添加HBase Service的相关参数，未在表中列出的参数请保持默认值。

表 23-3 参数说明

名称	说明	示例值
Service Name	创建的服务名称，固定填写：hbasedev。	hbasedev
Username	根据实际需要填写。	admin

名称	说明	示例值
Password	根据实际需要填写。	-
hadoop.security.authentication	Hadoop的认证方式，固定填写：Simple。	Simple
hbase.security.authentication	HBase的认证方式，固定填写：Simple。	Simple
hbase.zookeeper.property.clientPort	HBase集群中ZooKeeper的端口号。	2181
hbase.zookeeper.quorum	HBase集群中ZooKeeper地址。	192.168.0.7,192.168.0.8,192.168.0.9
zookeeper.znode.parent	HBase存在ZooKeeper中的根节点路径，固定填写：/hbase。	/hbase

图 23-12 Create hbasedev



Service Manager > Create Service

Create Service

Service Details :

Service Name *

Description

Active Status Enabled Disabled

Select Tag Service

Config Properties :

Username *

Password *

hadoop.security.authentication *

hbase.master.kerberos.principal

hbase.security.authentication *

hbase.zookeeper.property.clientPort *

hbase.zookeeper.quorum *

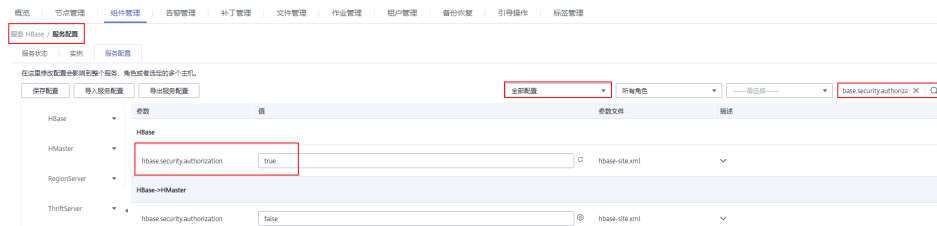
zookeeper.znode.parent *

步骤4 单击“Add”添加服务。

步骤5 启动Ranger HBase Plugin，授权Ranger管理HBase。

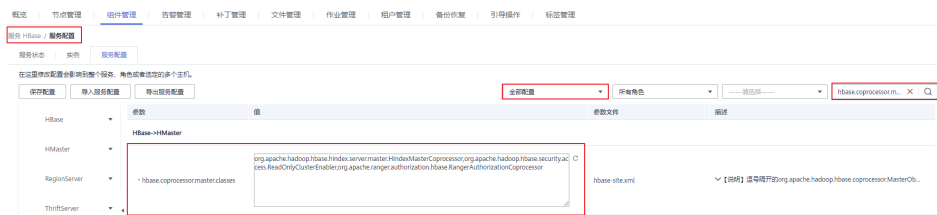
1. 在MRS控制台，单击集群名称进入集群详情页面。
2. 选择“组件管理”。
3. 选择“HBase > 服务配置”，将“基础配置”切换为“全部配置”。
4. 搜索并修改“hbase.security.authorization”为“true”（选择第一个HBase下的参数即可）。

图 23-13 修改 hbase.security.authorization



5. 搜索“hbase.coprocessor.master.classes”，并在原值后追加“org.apache.ranger.authorization.hbase.RangerAuthorizationCoprocesor”。

图 23-14 hbase.coprocessor.master.classes



6. 搜索“hbase.coprocessor.region.classes”，并在原值后追加“org.apache.ranger.authorization.hbase.RangerAuthorizationCoprocesor”。

图 23-15 hbase.coprocessor.region.classes



7. 单击“保存配置”，并勾选“重新启动受影响的服务或实例。”重启HMaster与RegionServer实例。

步骤6 在HBase Service hbasedev下创建对应的Policy。

1. 登录Ranger WebUI界面。
2. 在HBASE区域单击已添加的服务名称“hbasedev”。
3. 单击“Add New Policy”，新增访问控制策略。
4. 参考表23-4配置参数，未在表中列出的参数请保持默认值。

表 23-4 参数说明

参数	说明	示例值
Policy Name	策略名称。	Policy002
HBase Table	该策略允许访问的HBase表名称。	test1
HBase Column-family	该策略允许访问的HBase表对应的列族。	cf1

参数	说明	示例值
HBase Column	该策略允许访问的HBase表对应的表的列名。	name
Allow Conditions	<ul style="list-style-type: none"> - Select Group: 该策略允许访问的用户组。 - Select User: 该策略允许访问的用户组中的用户。 - Permissions: 该策略允许用户使用的权限。 	<ul style="list-style-type: none"> - Select Group: testuser - Select User: testuser - Permissions: Create和select

图 23-16 新增 hbasedev 的访问控制策略

5. 单击“Add”，完成策略添加，依据如上Policy，testuser用户组中的testuser用户拥有对HBase中“default” namespace下的“test1”表中“cf1:name”列有Create和select的权限，而对于其他列则没有任何的访问权限。

步骤7 参见[快速使用HBase进行离线数据分析](#)更新并登录HBase客户端，验证Ranger是否已经完成集成HBase。

1. 以客户端安装用户登录客户端安装节点，执行如下命令，进入hbase shell。
source /opt/client/bigdata_env
hbase shell

图 23-17 进入 hbase shell

```
[root@node-master1aHRf conf]# source /opt/client/bigdata_env
[root@node-master1aHRf conf]# hbase shell
INFO: Watching file:/opt/client/HBase/hbase/conf/log4j.properties for changes with interval : 60000
```

2. 添加数据，验证Ranger是否已经集成成功。
 - a. 为“test1”表中“cf1:name”列添加数据。
put 'test1','001','cf1:name','tom'
 - b. 为“test1”表中“cf1:age”列添加数据，该列用户无权限会添加数据失败。
put 'test1','001','cf1:age',10

图 23-18 验证 Ranger 集成 HBase

```
hbase(main):037:0> put 'test1','001','cf1:name','Tom'
0 row(s) in 0.1120 seconds
为001行添加cf1:name列数据

hbase(main):038:0> scan 'test1'
添加数据成功
ROW COLUMN+CELL
001 column=cf1:name, timestamp=1575426581007, value=Tom
1 row(s) in 0.0170 seconds

hbase(main):039:0> put 'test1','001','cf1:age',10
2019-12-04 10:30:15,637 WARN [hconnection-0x52ff99cd-shared--pool1-t26] client.AsyncProcess: #3, table=test1, attempt=1/7 failed-ops, last exception: org.apache.hadoop.hbase.security.AccessDeniedException: org.apache.hadoop.hbase.security.AccessDeniedException: Insufficient permissions for user 'testuser',action: put, tableName=test1, family:cf1, column: age
at org.apache.ranger.authorization.hbase.RangerAuthorizationCoprocessor.requirePermission(RangerAuthorizationCoprocessor.java:582)
为001行添加cf1:age列数据失败，由于没有权限
```

----结束

24 使用 Ranger（MRS 3.x）

24.1 登录 Ranger WebUI 界面

Ranger服务提供了集中式的权限管理框架，可以对HDFS、HBase、Hive、Yarn等组件进行细粒度的权限访问控制，并且提供了Web UI方便Ranger管理员进行操作。

Ranger 用户类型

Ranger中的用户可分为Admin、User、Auditor等类型，不同用户具有的Ranger管理界面查看和操作权限不同。

- Admin: Ranger安全管理员，可查看Ranger所有管理页面内容，进行服务权限管理插件及权限访问控制策略的管理操作，可查看审计信息内容，可进行用户类型设置。
- Auditor: Ranger审计管理员，可查看服务权限管理插件及权限访问控制策略的内容。
- User: 普通用户，可以被Ranger管理员赋予具体权限。

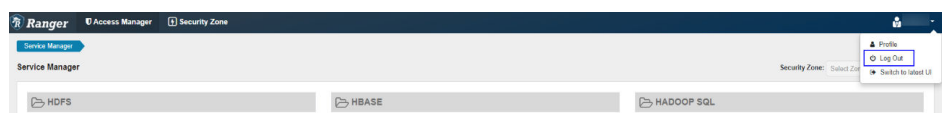
登录 Ranger 管理界面

安全模式（集群开启了Kerberos认证）

步骤1 使用admin用户登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。选择“集群 > 服务 > Ranger”，进入Ranger服务概览页面。

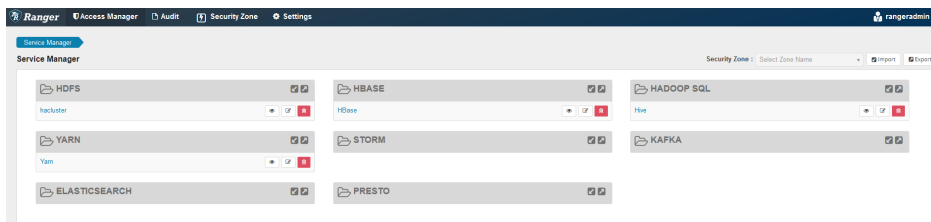
步骤2 单击“基本信息”区域中的“RangerAdmin”，进入Ranger WebUI界面。

- admin用户在Ranger中的用户类型为“User”，只能查看Access Manager和Security Zone页面。
- 如需查看所有管理页面，需要切换至rangeradmin用户或者其他具有Ranger管理员权限的用户：
 - a. 在Ranger WebUI界面，单击右上角用户名，选择“Log Out”，退出当前用户。



- b. 使用rangeradmin用户（默认密码为Rangeradmin@123）或者其他具有Ranger管理员权限用户重新登录。用户及默认密码请参考[用户信息一览表](#)。

图 24-1 Ranger WebUI



----结束

普通模式（集群关闭了Kerberos认证）：

步骤1 使用admin用户登录FusionInsight Manager，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。选择“集群 > 服务 > Ranger”，进入Ranger服务概览页面。

步骤2 单击“基本信息”区域中的“RangerAdmin”，进入Ranger WebUI界面。

普通模式下admin用户在Ranger中的用户类型为“Admin”，能查看Ranger所有管理页面，无需切换至rangeradmin用户。

说明

普通模式下使用rangeradmin用户登录Ranger WebUI界面，页面报错401。

----结束

在Ranger管理首页可查看当前Ranger已集成的各服务权限管理插件，用户可通过对应插件设置更细粒度的权限，具体主要操作页面功能描述参见[表24-1](#)。

表 24-1 Ranger 界面操作入口功能描述

入口	功能描述
Access Manager	查看当前Ranger已集成的各服务权限管理插件，用户可通过对应插件设置更细粒度的权限，具体操作请参考 添加Ranger权限策略 。
Audit	查看Ranger运行及权限管控相关审计日志信息，具体操作请参考 查看Ranger审计信息 。
Security Zone	配置安全区域，Ranger管理员可将各组件的资源切分为多个区域，由不同Ranger管理员为服务的指定资源设置安全策略，以便更好的管理，具体操作可参考 配置Ranger安全区信息 。
Settings	查看Ranger相关权限设置信息，例如查看用户、用户组、Role等，具体操作可参考 查看Ranger用户权限同步信息

24.2 MRS 集群服务启用 Ranger 鉴权

操作场景

该章节指导用户如何启用Ranger鉴权。安全模式默认开启Ranger鉴权，普通模式默认关闭Ranger鉴权。

操作步骤

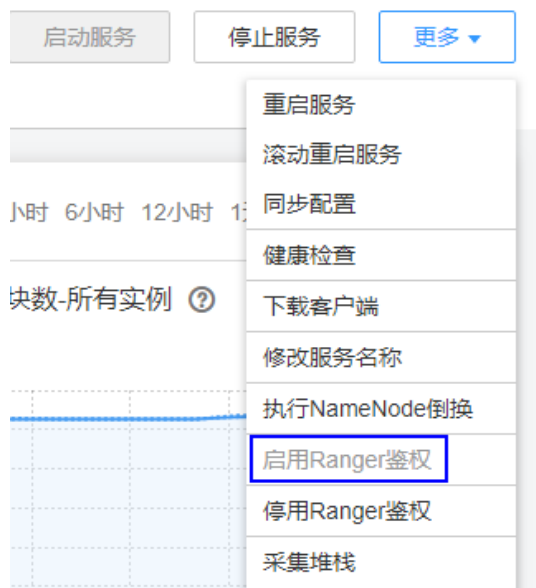
步骤1 登录FusionInsight Manager页面，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。选择“集群 > 服务 > 需要启用Ranger鉴权的服务名称”。

步骤2 在服务“概览”页面右上角单击“更多”，选择“启用Ranger鉴权”。在弹出的对话框中输入密码，单击“确定”，操作成功后单击“完成”。

说明

- 如果“启用Ranger鉴权”是灰色，表示已开启Ranger鉴权，如[图24-2](#)所示。
- 已启用Ranger授权的组件（HDFS与Yarn除外），Manager上非系统默认角色的权限将无法生效，需要通过配置Ranger策略为用户组赋权。

图 24-2 启用 Ranger 鉴权



步骤3 滚动重启服务或者重启服务。

----结束

24.3 添加 Ranger 权限策略

新安装的MRS集群默认安装Ranger服务并启用了Ranger鉴权模型，Ranger管理员可以通过组件权限插件对组件资源的访问设置细粒度的安全访问策略。

目前安全模式集群中支持Ranger的组件包括：HDFS、Yarn、HBase、Hive、Spark2x、Kafka、Storm。

通过 Ranger 配置用户权限策略

步骤1 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录 Ranger WebUI界面](#)。

步骤2 在Ranger首页的“Service Manager”区域内，单击组件名称下的权限插件名称，即可进入组件安全访问策略列表页面。

📖 说明

各组件的策略列表中，系统默认会生成若干条目，用于保证集群内的部分默认用户或用户组的权限（例如supergroup用户组），请勿删除，否则系统默认用户或用户组的权限会受影响。

步骤3 单击“Add New Policy”，根据业务场景规划配置相关用户或者用户组的资源访问策略。

不同组件的访问策略配置样例参考：

- [添加HDFS的Ranger访问权限策略](#)
- [添加HBase的Ranger访问权限策略](#)
- [添加Hive的Ranger访问权限策略](#)
- [添加Yarn的Ranger访问权限策略](#)
- [添加Spark2x的Ranger访问权限策略](#)
- [添加Kafka的Ranger访问权限策略](#)
- [添加Storm的Ranger访问权限策略](#)

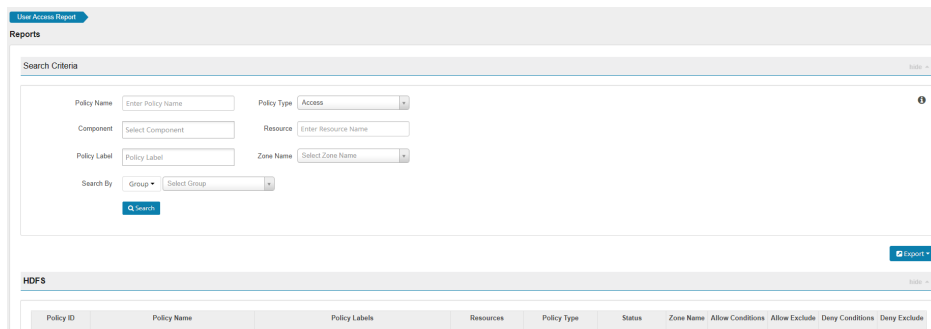
策略添加后，需等待30秒左右，待系统生效。

📖 说明

组件每次启动都会检查组件默认的Ranger Service是否存在，如果不存在则会创建以及为其添加默认Policy。如果用户在使用过程中误删了Service，可以重启或者滚动重启相应组件服务来恢复，若是误删了默认Policy，可先手动删除Service，再重启组件服务。

步骤4 单击“Access Manager > Reports”，可查看各组件所有的安全访问策略。

系统策略较多时，可通过策略名称、类型、组件、资源对象、策略标签、安全区域、用户或用户组等信息进行过滤搜索，也可以单击“Export”导出相关策略内容。



说明

- 对于同一个固定资源对象通常只能配置一条策略，多条策略针对的具体资源对象重复时将无法保存。
- 配置策略时，不同条件的优先级可参考[Ranger权限策略条件判断优先级](#)。

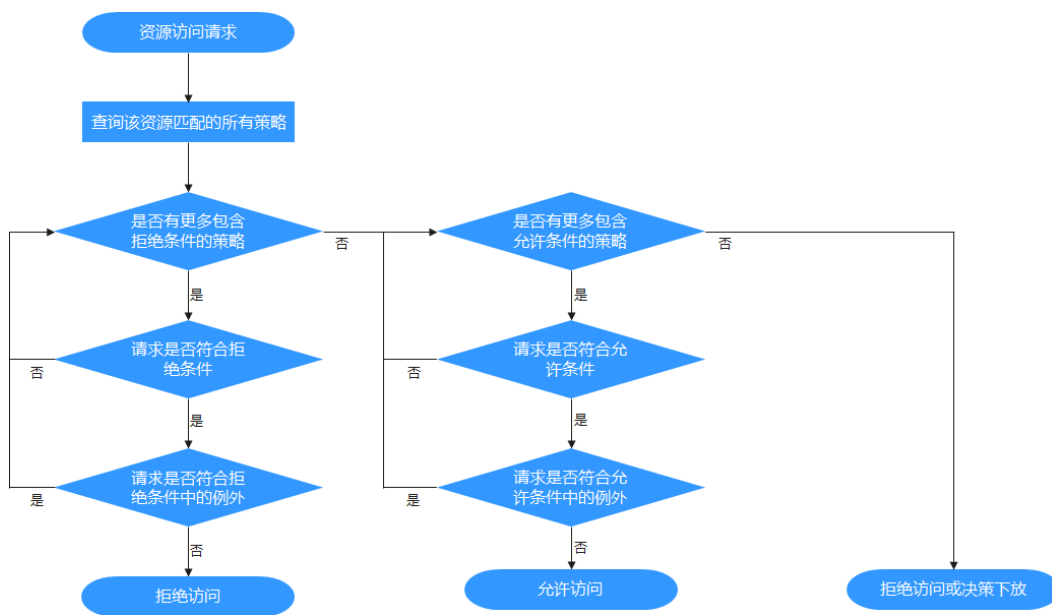
---结束

Ranger 权限策略条件判断优先级

配置资源的权限策略时，可配置针对该资源的允许条件（Allow Conditions）、允许例外条件（Exclude from Allow Conditions）、拒绝条件（Deny Conditions）以及拒绝例外条件（Exclude from Deny Conditions），以满足不同场景下的例外需求。

不同条件的优先级由高到低为：拒绝例外条件 > 拒绝条件 > 允许例外条件 > 允许条件。

系统判断流程可参考下图所示，如果组件资源请求未匹配到Ranger中的权限策略，系统默认将拒绝访问。但是对于HDFS和Yarn，系统会将决策下放给组件自身的访问控制层继续进行判断。



例如要将一个文件夹FileA的读写权限授权给用户组groupA，但是该用户组内某个用户UserA除外，这时可以增加一个允许条件及一个例外条件即可实现。

24.4 Ranger 权限策略配置示例

24.4.1 添加 HDFS 的 Ranger 访问权限策略

操作场景

Ranger管理员可通过Ranger为HDFS用户配置HDFS目录或文件的读、写和执行权限。

前提条件


- 已安装Ranger服务且服务运行正常。
- 已创建需要配置权限的用户、用户组或Role。

操作步骤

- 步骤1** 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录Ranger WebUI界面](#)。
- 步骤2** 在首页中单击“HDFS”区域的组件插件名称，例如“hacluster”。
- 步骤3** 单击“Add New Policy”，添加HDFS权限控制策略。
- 步骤4** 根据业务需求配置相关参数。

表 24-2 HDFS 权限参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，可以根据这些标签搜索报告和筛选策略。
Resource Path	资源路径，配置当前策略适用的HDFS路径文件夹或文件，可填写多个值，支持使用通配符“*”（例如“/test/*”）。 如需子目录继承上级目录权限，可打开递归开关按钮。 如果父目录开启递归，同时子目录也配置了策略，则子目录同时拥有父目录和子目录的策略；如果父目录与子目录的策略相悖，则以子目录策略为准。 <ul style="list-style-type: none">• non-recursive：关闭递归• recursive：打开递归
Description	策略描述信息。
Audit Logging	是否审计此策略。

参数名称	描述
Allow Conditions	<p>策略允许条件，配置本策略内允许的权限及例外，例外条件优先级高于正常条件。</p> <p>在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的Role、用户组或用户，单击“Add Conditions”，添加策略适用的IP地址范围，单击“Add Permissions”，添加对应权限。</p> <ul style="list-style-type: none"> • Read：读权限 • Write：写权限 • Execute：执行权限 • Select/Deselect All：全选/取消全选 <p>如需让当前条件中的用户或用户组管理本条策略，可勾选“Delegate Admin”使这些用户或用户组成为受委托的管理员。被委托的管理员可以更新、删除本策略，还可以基于原始策略创建子策略。</p> <p>如需添加多条权限控制规则，可单击  按钮添加。如需删除权限控制规则，可单击  按钮删除。</p> <p>Exclude from Allow Conditions：配置排除在允许条件之外的例外规则。</p>
Deny All Other Accesses	<p>是否拒绝其它所有访问。</p> <ul style="list-style-type: none"> • True：拒绝其它所有访问。 • False：设置为false，可配置Deny Conditions。
Deny Conditions	<p>策略拒绝条件，配置本策略内拒绝的权限及例外，配置方法与“Allow Conditions”类似，拒绝条件的优先级高于“Allow Conditions”中配置的允许条件。</p> <p>Exclude from Deny Conditions：配置排除在拒绝条件之外的例外规则。</p>

例如为用户“testuser”添加“/user/test”目录的写权限，配置如下：

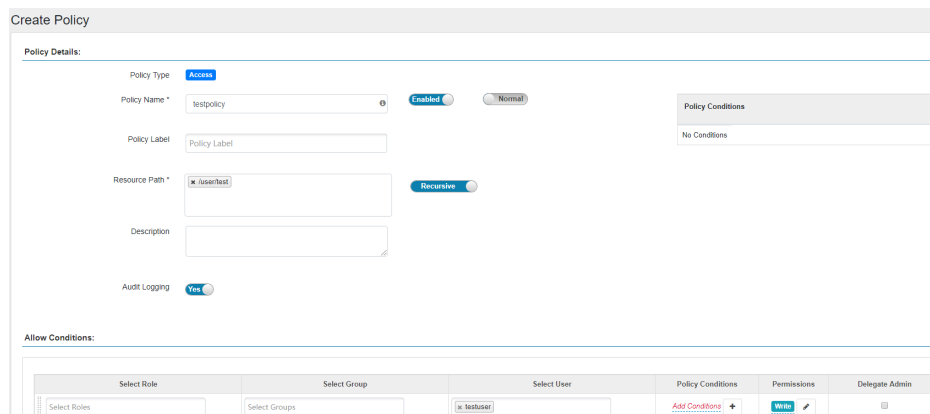






表 24-3 设置权限

任务场景	角色授权操作
设置HDFS管理员权限	<ol style="list-style-type: none"> 1. 在首页中单击“HDFS”区域的组件插件名称，例如“hacluster”。 2. 选择“Policy Name”为“all - path”的策略，单击  按钮编辑策略。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。
设置用户执行HDFS检查和HDFS修复的权限	<ol style="list-style-type: none"> 1. 在“Resource Path”配置文件夹或文件。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Read”和“Execute”。
设置用户读取其他用户的目录或文件的权限	<ol style="list-style-type: none"> 1. 在“Resource Path”配置文件夹或文件。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Read”和“Execute”。
设置用户在其他用户的文件写入数据的权限	<ol style="list-style-type: none"> 1. 在“Resource Path”配置文件夹或文件。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Write”和“Execute”。
设置用户在其他用户的目录新建或删除子文件、子目录的权限	<ol style="list-style-type: none"> 1. 在“Resource Path”配置文件夹或文件。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Write”和“Execute”。
设置用户在其他用户的目录或文件执行的权限	<ol style="list-style-type: none"> 1. 在“Resource Path”配置文件夹或文件。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Execute”。
设置子目录继承上级目录权限	<ol style="list-style-type: none"> 1. 在“Resource Path”配置文件夹或文件。 2. 打开递归开关按钮，“Recursive”即为打开递归。

步骤5 （可选）添加策略有效期。在页面右上角单击“Add Validity period”，设置“Start Time”和“End Time”，选择“Time Zone”。单击“Save”保存。如需添加多条策略有效期，可单击  按钮添加。如需删除策略有效期，可单击  按钮删除。

步骤6 单击“Add”，在策略列表可查看策略的基本信息。等待策略生效后，验证相关权限是否正常。

如需禁用某条策略，可单击  按钮编辑策略，设置策略开关为“Disabled”。

如果不再使用策略，可单击  按钮删除策略。

----结束

24.4.2 添加 HBase 的 Ranger 访问权限策略

操作场景

Ranger管理员可通过Ranger为HBase用户配置HBase表和列族，列的权限。

前提条件

- 已安装Ranger服务且服务运行正常。
- 已创建需要配置权限的用户、用户组或Role。

操作步骤

步骤1 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录Ranger WebUI界面](#)。

步骤2 在首页中单击“HBASE”区域的组件插件名称如“HBase”。

步骤3 单击“Add New Policy”，添加HBase权限控制策略。

步骤4 根据业务需求配置相关参数。

表 24-4 HBase 权限参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，可以根据这些标签搜索报告和筛选策略。
HBase Table	<p>将适用该策略的表。</p> <p>可支持通配符“*”，例如“table1:*”表示table1下的所有表。</p> <p>“Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。</p> <p>说明</p> <p>Ranger界面上HBase服务插件的“hbase.rpc.protection”参数值必须和HBase服务端的“hbase.rpc.protection”参数值保持一致。具体请参考Ranger界面添加或者修改HBase策略时，无法使用通配符搜索已存在的HBase表。</p>



参数名称	描述
HBase Column-family	将适用该策略的列族。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
HBase Column	将适用该策略的列。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
Description	策略描述信息。
Audit Logging	是否审计此策略。
Allow Conditions	<p>策略允许条件，配置本策略内允许的权限及例外。</p> <p>在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的Role、用户组或用户，单击“Add Conditions”，添加策略适用的IP地址范围，单击“Add Permissions”，添加对应权限。</p> <ul style="list-style-type: none"> • Read：读权限 • Write：写权限 • Create：创建权限 • Admin：管理权限 • Select/Deselect All：全选/取消全选 <p>如需让当前条件中的用户或用户组管理本条策略，可勾选“Delegate Admin”使这些用户或用户组成为受委托的管理员。被委托的管理员可以更新、删除本策略，还可以基于原始策略创建子策略。</p> <p>如需添加多条权限控制规则，可单击  按钮添加。如需删除权限控制规则，可单击  按钮删除。</p> <p>Exclude from Allow Conditions：配置策略例外条件。</p>
Deny All Other Accesses	<p>是否拒绝其它所有访问。</p> <ul style="list-style-type: none"> • True：拒绝其它所有访问 • False：设置为False，可配置Deny Conditions。
Deny Conditions	<p>策略拒绝条件，配置本策略内拒绝的权限及例外，配置方法与“Allow Conditions”类似。</p> <p>拒绝条件的优先级高于“Allow Conditions”中配置的允许条件。</p> <p>Exclude from Deny Conditions：配置排除在拒绝条件之外的例外规则。</p>



表 24-5 设置权限

任务场景	角色授权操作
设置HBase管理员权限	<ol style="list-style-type: none"> 1. 在首页中单击“HBase”区域的组件插件名称，例如“HBase”。 2. 选择“Policy Name”为“all - table, column-family, column”的策略，单击按钮编辑策略。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。
设置用户创建表的权限	<ol style="list-style-type: none"> 1. 在“HBase Table”配置表名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Create”。 4. 该用户具有以下操作权限： create table drop table truncate table alter table enable table flush table flush region compact disable enable desc
设置用户写入数据的权限	<ol style="list-style-type: none"> 1. 在“HBase Table”配置表名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Write”。 4. 该用户具有put, delete, append, incr, bulkload等操作权限。
设置用户读取数据的权限	<ol style="list-style-type: none"> 1. 在“HBase Table”配置表名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Read”。 4. 该用户具有get, scan操作权限。


任务场景	角色授权操作
设置用户管理命名空间或表的权限	<ol style="list-style-type: none"> 1. 在“HBase Table”配置表名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Admin”。 4. 该用户具有rsgroup, peer, assign, balance等操作权限。
设置列的读取或写入权限	<ol style="list-style-type: none"> 1. 在“HBase Table”配置表名。 2. 在“HBase Column-family”配置列族名。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“Read”或者“Write”。


说明

如果用户在hbase shell中执行desc操作，需要同时给该用户赋予hbase:quota表的读权限。

步骤5 （可选）添加策略有效期。在页面右上角单击“Add Validity period”，设置“Start Time”和“End Time”，选择“Time Zone”。单击“Save”保存。如需添加多条策略有效期，可单击  按钮添加。如需删除策略有效期，可单击  按钮删除。

步骤6 单击“Add”，在策略列表可查看策略的基本信息。等待策略生效后，验证相关权限是否正常。

如需禁用某条策略，可单击  按钮编辑策略，设置策略开关为“Disabled”。

如果不再使用策略，可单击  按钮删除策略。

----结束

24.4.3 添加 Hive 的 Ranger 访问权限策略

操作场景

Ranger管理员可通过Ranger为Hive用户进行相关的权限设置。Hive默认管理员账号为hive，初始密码为Hive@123。

前提条件

- 已安装Ranger服务且服务运行正常。
- 已创建用户需要配置权限的用户、用户组或Role。
- 用户加入hive组。

操作步骤

- 步骤1** 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录 Ranger WebUI界面](#)。
- 步骤2** 在首页中单击“HADOOP SQL”区域的组件插件名称如“Hive”。
- 步骤3** 在“Access”页签单击“Add New Policy”，添加Hive权限控制策略。
- 步骤4** 根据业务需求配置相关参数。

表 24-6 Hive 权限参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，您可以根据这些标签搜索报告和筛选策略。
database	将适用该策略的列Hive数据库名称。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
table	将适用该策略的Hive表名称。 如果需要添加基于UDF的策略，可切换为UDF，然后输入UDF的名称。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
Hive Column	将适用该策略的列名，填写*时表示所有列。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
Description	策略描述信息。
Audit Logging	是否审计此策略。


参数名称	描述
Allow Conditions	<p>策略允许条件，配置本策略内允许的权限及例外。</p> <p>在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的Role、用户组或用户，单击“Add Conditions”，添加策略适用的IP地址范围，然后再单击“Add Permissions”，添加对应权限。</p> <ul style="list-style-type: none"> ● select: 查询权限 ● update: 更新权限 ● Create: 创建权限 ● Drop: drop操作权限 ● Alter: alter操作权限 ● Index: 索引操作权限 ● All: 所有执行权限 ● Read: 可读权限 ● Write: 可写权限 ● Temporary UDF Admin: 临时UDF管理权限 ● Select/Deselect All: 全选/取消全选 <p>如需添加多条权限控制规则，可单击  按钮添加。</p> <p>如需当前条件中的用户或用户组管理本条策略，可勾选“Delegate Admin”，这些用户将成为受委托的管理员。被委托的管理员可以更新、删除本策略，它还可以基于原始策略创建子策略。</p>
Deny Conditions	<p>策略拒绝条件，配置本策略内拒绝的权限及例外，配置方法与“Allow Conditions”类型。</p>

表 24-7 设置权限

任务场景	角色授权操作
role admin操作	<ol style="list-style-type: none"> 1. 在首页中单击“Settings”，选择“Roles”。 2. 单击Role Name为admin的角色，在“Users”区域，单击“Select User”，选择对应用户名。 3. 单击Add Users按钮，在对应用户名所在行勾选“Is Role Admin”，单击“Save”保存配置。 <p>说明 Ranger页面的“Settings”选项只有rangeradmin用户有权限访问。用户绑定Hive管理员角色后，在每个维护操作会话中，还需要执行以下操作：</p> <ol style="list-style-type: none"> 1. 以客户端安装用户，登录安装Hive客户端的节点。 2. 执行以下命令配置环境变量。 例如，Hive客户端安装目录为“/opt/hiveclient”，执行 source /opt/hiveclient/bigdata_env 3. 执行以下命令认证用户。 kinit Hive业务用户 4. 执行以下命令登录客户端工具。 beeline 5. 执行以下命令更新用户的管理员权限。 set role admin;
创建库表操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写或选择对应的数据库(如果是创建表则在“table”右侧填写或选择对应的表)，在“column”右侧填写或选择“*”。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“Create”。
删除库表操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写或选择对应的数据库（如果是删除表则在“table”右侧填写或选择对应的表），在“column”右侧填写并选择“*”。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“Drop”。
查询操作(select、desc、show)	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写或选择对应的数据库(如果是表则在“table”右侧填写或选择对应的表)，在“column”右侧填写并选择对应的列（*代表所有列）。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“select”。


任务场景	角色授权操作
Alter操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写并选择对应的数据库(如果是表则在“table”右侧填写或选择对应的表), 在“column”右侧填写或选择“*”。 3. 在“Allow Conditions”区域, 单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”, 勾选“Alter”。
LOAD操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写或选择对应的数据库, 在“table”右侧填写或选择对应的表, 在“column”右侧填写并选择“*”。 3. 在“Allow Conditions”区域, 单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”, 勾选“update”。
INSERT、DELETE操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写或选择对应的数据库, 在“table”右侧填写或选择对应的表, 在“column”右侧填写并选择“*”。 3. 在“Allow Conditions”区域, 单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”, 勾选“update”。 5. 用户还需要具有Yarn任务队列的“submit”权限, 权限配置参考添加Yarn的Ranger访问权限策略。
GRANT、REVOKE操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写或选择对应的数据库, 在“table”右侧填写或选择对应的表, 在“column”右侧填写并选择“*”。 3. 在“Allow Conditions”区域, 单击“Select User”下选择框选择用户。 4. 勾选“Delegate Admin”。
ADD JAR操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. 单击“database”并在下拉菜单中选择“global”。在“global”右侧填写或选择“*”。 3. 在“Allow Conditions”区域, 单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”, 勾选“Temporary UDF Admin”。


任务场景	角色授权操作
UDF 操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写或选择对应的数据库，“udf”右侧填写对应的udf 函数名。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，根据需求，给用户勾选相应权限（udf支持 Create, select, Drop）。
VIEW操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写或选择对应的数据库，在“table”右侧填写或选择对应的VIEW名称，在“column”右侧填写并选择“*”。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，参照表格上述相关操作，根据需求，给用户勾选相应权限。
dfs命令操作	执行set role admin 操作才可使用。
其他用户库表操作	<ol style="list-style-type: none"> 1. 参照表格上述相关操作添加对应权限。 2. 给用户添加其他用户库表的HDFS路径的读、写、执行权限，详情请参考添加HDFS的Ranger访问权限策略。

📖 说明

- 如果用户在执行命令时指定了HDFS路径，需要给该用户添加HDFS路径的读、写、执行权限，详情请参考[添加HDFS的Ranger访问权限策略](#)。也可以不配置HDFS的Ranger策略，通过之前Hive权限插件的方式，给角色添加权限，然后把角色赋予对应用户。如果HDFS Ranger策略可以匹配到Hive库表的文件或目录权限，则优先使用HDFS Ranger策略。
- Ranger策略中的URL策略是hive表存储在obs上的场景涉及，URL填写对象在obs上的完整路径。与URL联合使用的Read, Write 权限，其他场景不涉及URL策略。
- Ranger策略中global策略仅用于和Temporary UDF Admin权限联合使用，控制UDF包的上传。
- Ranger策略中的hiveservice策略仅用于和服务 Admin权限联合使用，用于控制命令：**kill query <queryId>** 结束正在执行的任务的权限。
- lock、index、refresh、replAdmin 权限暂不支持。
- 使用**show grant**命令查看表权限，表owner的grantor列统一显示为hive用户，其他用户Ranger页面赋权或后台采用grant命令赋权，则grantor显示为对应用户；若用户需要查看之前Hive权限插件的结果，可设置hive-ext.ranger.previous.privileges.enable为true后采用**show grant**查看。

步骤5 单击“Add”，在策略列表可查看策略的基本信息。等待策略生效后，验证相关权限是否正常。

如需禁用某条策略，可单击  按钮编辑策略，设置策略开关为“Disabled”。

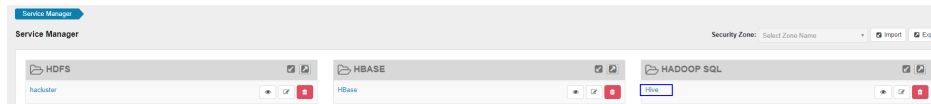
如果不再使用策略，可单击  按钮删除策略。

----结束

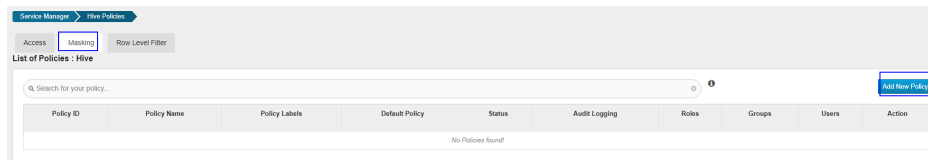
Hive 数据脱敏

Ranger支持对Hive数据进行脱敏处理（Data Masking），可对用户执行的select操作的返回结果进行处理，以屏蔽敏感信息。

步骤1 登录Ranger WebUI界面，在首页中单击“HADOOP SQL”区域的“Hive”




步骤2 在“Masking”页签单击“Add New Policy”，添加Hive权限控制策略。



步骤3 根据业务需求配置相关参数。

表 24-8 Hive 数据脱敏参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，您可以根据这些标签搜索报告和筛选策略。
Hive Database	配置当前策略适用的Hive中数据库名称。
Hive Table	配置当前策略适用的Hive中的表名称。
Hive Column	可添加列名。
Description	策略描述信息。
Audit Logging	是否审计此策略。

参数名称	描述
Mask Conditions	<p>在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的对象，单击“Add Conditions”，添加策略适用的IP地址范围，然后再单击“Add Permissions”，勾选“select”权限。</p> <p>单击“Select Masking Option”，选择数据脱敏时的处理策略：</p> <ul style="list-style-type: none"> • Redact: 用x屏蔽所有字母字符，用0屏蔽所有数字字符。 • Partial mask: show last 4: 只显示最后的4个字符，其他用x代替。 • Partial mask: show first 4: 只显示开始的4个字符，其他用x代替。 • Hash: 用值的哈希值替换原值，采用的是hive的内置mask_hash函数，只对string、char、varchar类型的字段生效，其他类型的字段会返回NULL值。 • Nullify: 用NULL值替换原值。 • Unmasked(retain original value): 原样显示。 • Date: show only year: 仅显示日期字符串的年份部分，并将月份和日期默认为01/01。 • Custom: 可使用任何有效返回与被屏蔽的列中的数据类型相同的数据类型来自定义策略。 <p>如需添加多列的脱敏策略，可单击  按钮添加。</p>

步骤4 单击“Add”，在策略列表可查看策略的基本信息。

步骤5 用户通过Hive客户端对配置了数据脱敏策略的表执行select操作，系统将对数据进行处理后进行展示。

说明

处理数据需要用户同时具有向Yarn队列提交任务的权限。

---结束

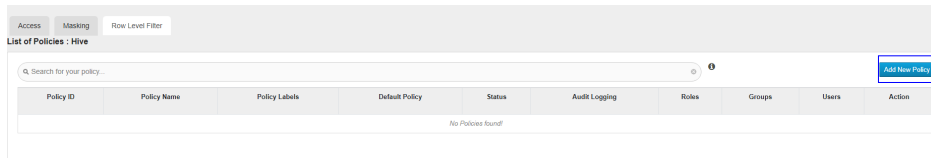
Hive 行级别数据过滤

Ranger支持用户对Hive数据表执行select操作时进行行级别的数据过滤。

步骤1 登录Ranger WebUI界面，在首页中单击“HADOOP SQL”区域的“Hive”。



步骤2 在“Row Level Filter”页签单击“Add New Policy”，添加行数据过滤策略。



步骤3 根据业务需求配置相关参数。

表 24-9 Hive 行数据过滤参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，您可以根据这些标签搜索报告和筛选策略。
Hive Database	配置当前策略适用的Hive中数据库名称。
Hive Table	配置当前策略适用的Hive中的表名称。
Description	策略描述信息。
Audit Logging	是否审计此策略。
Row Filter Conditions	<p>在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的对象，单击“Add Conditions”，添加策略适用的IP地址范围，然后再单击“Add Permissions”，勾选“select”权限。</p> <p>单击“Row Level Filter”，填写数据过滤规则。</p> <p>例如过滤表A中“name”列“zhangsan”行的数据，过滤规则为：name <> 'zhangsan'。更多信息可参考Ranger官方文档。</p> <p>如需添加更多规则，可单击  按钮添加。</p>

步骤4 单击“Add”，在策略列表可查看策略的基本信息。

步骤5 用户通过Hive客户端对配置了数据脱敏策略的表执行select操作，系统将对数据进行处理后进行展示。

 **说明**

处理数据需要用户同时具有向Yarn队列提交任务的权限。

----**结束**

24.4.4 添加 Impala 的 Ranger 访问权限策略

 **说明**

本章节仅适用于MRS 3.1.5及之后版本。

操作场景

在创建完安装了Ranger组件的MRS集群后，Hive/Impala的权限控制集成在Ranger中，由于Impala复用Hive的权限策略，本章节主要介绍Hive组件如何集成在Ranger中。

前提条件

- 已安装Ranger服务且服务运行正常。
- 已创建需要配置权限的用户、用户组或Role。

操作步骤

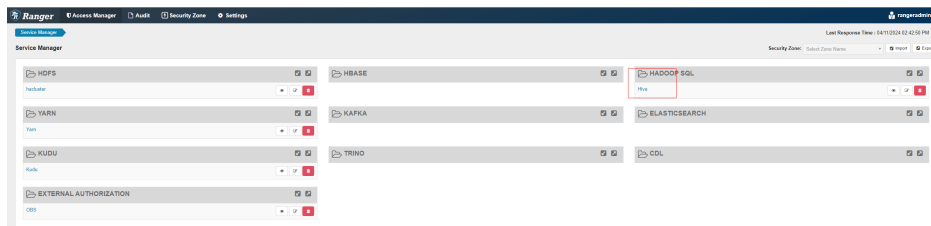
步骤1 普通模式（未开启Kerberos认证）下确认Hive，Impala是否启用Ranger鉴权。

如果未启用则需要先对Hive启用Ranger鉴权，并重启Hive服务，再重启Impala服务。然后对Impala启用Ranger鉴权，并重启Impala服务。**安全集群（开启Kerberos认证）默认启用Ranger，可跳过该步骤。**

步骤2 普通模式下，登录Manager页面，选择“集群 > 服务 > Ranger > 配置 > 全部配置 > UserSync（角色）”，在自定义配置“ranger.usersync.config.expandor”新增如下配置项，修改完成后重启Ranger服务。**安全集群（开启Kerberos认证）跳过该步骤。**

参数	值
ranger.usersync.sync.source	ldap

步骤3 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录Ranger WebUI界面](#)，选择Hive。



步骤4 添加访问控制策略，即Policy。

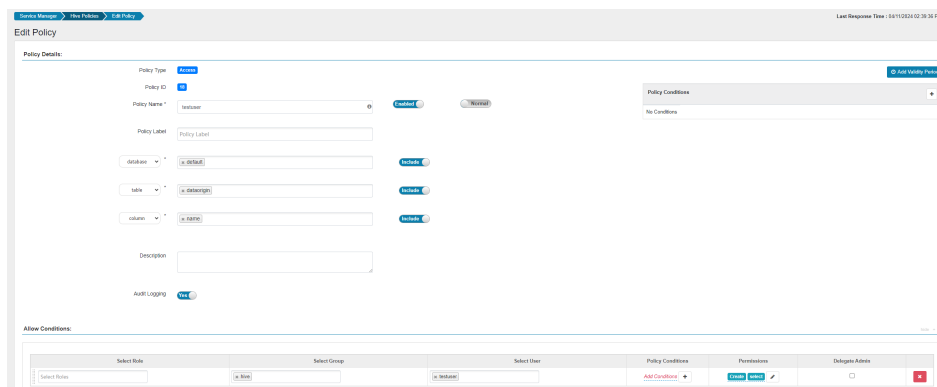
1. 在HIVE区域单击已添加的服务名称“Hive”。
2. 单击“Add New Policy”，新增访问控制策略。
3. 参考[表24-10](#)配置参数，未在表中列出的参数请保持默认值。

表 24-10 参数说明

参数	说明	示例值
Policy Name	策略名称。	testuser
database	该策略允许访问的数据库名称。	default

参数	说明	示例值
table	该策略允许访问的数据库对应的表名称。	dataorigin
Hive Column	该策略允许访问的数据库对应的表的列名。	name
Allow Conditions	<ul style="list-style-type: none"> - Select Group: 该策略允许访问的用户组。 - Select User: 该策略允许访问的用户组中的用户。 - Permissions: 该策略允许用户使用的权限。 	<ul style="list-style-type: none"> - Select Group: hive - Select User: testuser - Permissions: select,Create

图 24-3 新增 testuser 访问策略控制



- 单击“Add”，完成策略添加，依据如上Policy示例，hive用户组中的testuser用户将对Hive的“default”数据库中的表“dataorigin”的“name”列有Create和select的权限，而对于其他列则没有任何的访问权限。

步骤5 登录Impala客户端，验证Ranger是否已经完成集成Impala。

- 以客户端安装用户登录客户端安装节点，执行如下命令，初始化环境变量。

```
source /opt/client/bigdata_env
```

- 建立连接并使用testuser登录。

- 普通集群（未开启Kerberos认证）执行如下命令：

```
impala-shell -i <Impalad节点IP> -u testuser
```

- 安全集群（开启Kerberos认证）执行如下命令：

```
kinit testuser
```

并输入密码登录

```
impala-shell -i <Impalad节点IP>
```

- 查询数据，验证Ranger是否已经集成成功。

- 执行select * from dataorigin失败，报错显示权限不足。

- 执行select name from dataorigin成功，符合预设的权限。

```
[172.16.0.112:21050] default> select * from dataorigin;
Query: select * from dataorigin
Query submitted at: 2024-04-11 14:52:26 (Coordinator: http://172.16.0.112:25000)
ERROR: AuthorizationException: User 'testuser@523f6018_A35A_4301_B69E_EC38906755EE.COM' does not have privileges to execute 'SELECT' on: default.dataorigin

[172.16.0.112:21050] default> select name from dataorigin;
Query submitted at: 2024-04-11 14:52:53 (Coordinator: http://172.16.0.112:25000)
Query progress can be monitored at: http://172.16.0.112:25000/query_plan?query_id=6c41bf7b028e7973:d25c735208080800
-----+-----
| name |
-----+-----
| pgvpd0 |
| splited |
| kiewei |
| oftrkf |
| rknaql |
| kwhakb |
| fwoj |
-----+-----
Fetched 7 row(s) in 0.11s
```

注意

- 如果用户在执行命令时指定了HDFS路径，需要给该用户添加HDFS路径的读、写、执行权限，详情请参考[添加HDFS的Ranger访问权限策略](#)。也可以不配置HDFS的Ranger策略，通过之前Hive权限插件的方式，给角色添加权限，然后把角色赋予对应用户。如果HDFS Ranger策略可以匹配到Hive库表的文件或目录权限，则优先使用HDFS Ranger策略。
- 如果在Hive中创建的表，需要在Impala执行invalidate metadata刷新元数据。此时需要给该用户赋予refresh权限或者使用hive用户执行invalidate metadata，否则会遇到如下报错：

```
[192.168.0.24:21000] default> INVALIDATE METADATA;
Query: INVALIDATE METADATA
Query submitted at: 2024-04-03 11:45:47 (Coordinator: http://192.168.0.24:25000)
ERROR: AuthorizationException: User 'zhoujie' does not have privileges to execute 'INVALIDATE METADATA/REFRESH' on: impala_server

[192.168.0.24:21000] default> quit
> quit:
goodbye zhoujie
```

----结束

24.4.5 添加 Yarn 的 Ranger 访问权限策略

操作场景

Ranger管理员可通过Ranger为Yarn用户配置Yarn管理员权限以及Yarn队列资源管理权限。

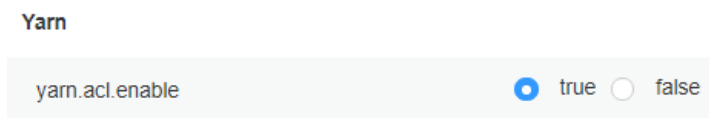
前提条件

- 已安装Ranger服务且服务运行正常。
- 已创建需要配置权限的用户、用户组或Role。

操作步骤



- 步骤1** 登录FusionInsight Manager界面，选择“集群 > 服务 > Yarn”。
- 步骤2** 选择“配置 > 全部配置”，搜索参数“yarn.acl.enable”，修改参数值为“true”。如果该参数值已经为“true”，则无需处理。

图 24-4 配置参数“yarn.acl.enable”




- 步骤3** 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录 Ranger WebUI界面](#)。
- 步骤4** 在首页中单击“YARN”区域的组件插件名称如“Yarn”。
- 步骤5** 单击“Add New Policy”，添加Yarn权限控制策略。
- 步骤6** 根据业务需求配置相关参数。



表 24-11 Yarn 权限参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，可以根据这些标签搜索报告和筛选策略。
Queue	队列名称，支持通配符“*”。 如需子队列继承上级队列权限，可打开递归开关按钮。 <ul style="list-style-type: none"> • Non-recursive: 关闭递归 • Recursive: 打开递归
Description	策略描述信息。
Audit Logging	是否审计此策略。
Allow Conditions	策略允许条件，配置本策略内允许的权限及例外。 在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的Role、用户组或用户，单击“Add Conditions”，添加策略适用的IP地址范围，单击“Add Permissions”，添加对应权限。 <ul style="list-style-type: none"> • submit-app: 提交队列任务权限 • admin-queue: 管理队列任务权限 • Select/Deselect All: 全选/取消全选 如需让当前条件中的用户或用户组管理本条策略，可勾选“Delegate Admin”使这些用户成为受委托的管理员。被委托的管理员可以更新、删除本策略，它还可以基于原始策略创建子策略。 如需添加多条权限控制规则，可单击  按钮添加。如需删除权限控制规则，可单击  按钮删除。 Exclude from Allow Conditions: 配置策略例外条件。

参数名称	描述
Deny All Other Accesses	是否拒绝其它所有访问。 <ul style="list-style-type: none"> • True: 拒绝其它所有访问 • False: 设置为False, 可配置Deny Conditions。
Deny Conditions	策略拒绝条件, 配置本策略内拒绝的权限及例外, 配置方法与“Allow Conditions”类似。拒绝条件的优先级高于“Allow Conditions”中配置的允许条件。 Exclude from Deny Conditions: 配置排除在拒绝条件之外的例外规则。


表 24-12 设置权限

任务场景	角色授权操作
设置Yarn管理员权限	<ol style="list-style-type: none"> 1. 在首页中单击“YARN”区域的组件插件名称, 例如“Yarn”。 2. 选择“Policy Name”为“all - queue”的策略, 单击  按钮编辑策略。 3. 在“Allow Conditions”区域, 单击“Select User”下选择框选择用户。
设置用户在指定Yarn队列提交任务的权限	<ol style="list-style-type: none"> 1. 在“Queue”配置队列名。 2. 在“Allow Conditions”区域, 单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”, 勾选“submit-app”。
设置用户在指定Yarn队列管理任务的权限	<ol style="list-style-type: none"> 1. 在“Queue”配置队列名。 2. 在“Allow Conditions”区域, 单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”, 勾选“admin-queue”。

步骤7 (可选) 添加策略有效期。在页面右上角单击“Add Validity period”, 设置“Start Time”和“End Time”, 选择“Time Zone”。单击“Save”保存。如需添加多条策略有效期, 可单击  按钮添加。如需删除策略有效期, 可单击  按钮删除。

步骤8 单击“Add”, 在策略列表可查看策略的基本信息。等待策略生效后, 验证相关权限是否正常。

如需禁用某条策略, 可单击  按钮编辑策略, 设置策略开关为“Disabled”。

如果不再使用策略, 可单击  按钮删除策略。

----结束

说明

Ranger Yarn上面各个权限之间相互独立，没有语义上的包含与被包含关系。当前支持下面两种权限：

- submit-app：提交队列任务权限
- admin-queue：管理队列任务权限

虽然admin-queue也有提交任务的权限，但和submit-app权限之间并没有包含关系。

24.4.6 添加 Spark2x 的 Ranger 访问权限策略

操作场景

Ranger管理员可通过Ranger为Spark2x用户进行相关的权限设置。

说明

1. Spark2x开启或关闭Ranger鉴权后，需要重启Spark2x服务。
2. 需要重新下载客户端，或手动刷新客户端配置文件“客户端安装目录/Spark2x/spark/conf/spark-defaults.conf”：
开启Ranger鉴权：spark.ranger.plugin.authorization.enable=true，同时需要修改参数“spark.sql.authorization.enabled”值为“true”。
关闭Ranger鉴权：spark.ranger.plugin.authorization.enable=false
3. Spark2x中，spark-beeline（即连接到JDBCServer的应用）支持Ranger的IP过滤策略（即Ranger权限策略中的**Policy Conditions**），spark-submit与spark-sql不支持。

前提条件

- 已安装Ranger服务且服务运行正常。
- 已启用Hive服务的Ranger鉴权功能，并且需要先重启Hive服务，再重启Spark服务，再启用Spark服务的Ranger鉴权。启用Spark服务的Ranger鉴权后再重启Spark服务。
- 已创建用户需要配置权限的用户、用户组或Role。
- 创建的用户已加入hive用户组。

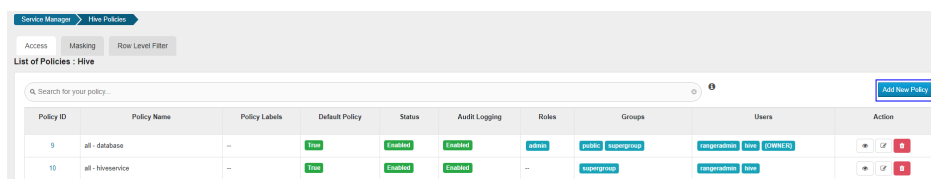
操作步骤

步骤1 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录Ranger WebUI界面](#)。

步骤2 在首页中单击“HADOOP SQL”区域的组件插件名称如“Hive”。



步骤3 在“Access”页签单击“Add New Policy”，添加Spark2x权限控制策略。



步骤4 根据业务需求配置相关参数。

表 24-13 Spark2x 权限参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如： 192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，您可以根据这些标签搜索报告和筛选策略。
database	适用该策略的Spark2x数据库名称。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
table	适用该策略的Spark2x表名称。 如果需要添加基于UDF的策略，可切换为UDF，然后输入UDF的名称。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
column	适用该策略的列名，填写*时表示所有列。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
Description	策略描述信息。
Audit Logging	是否审计此策略。


参数名称	描述
Allow Conditions	<p>策略允许条件，配置本策略内允许的权限及例外。</p> <p>在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的Role、用户组或用户，单击“Add Conditions”，添加策略适用的IP地址范围，然后再单击“Add Permissions”，添加对应权限。</p> <ul style="list-style-type: none"> • select: 查询权限 • update: 更新权限 • Create: 创建权限 • Drop: drop操作权限 • Alter: alter操作权限 • Index: 索引操作权限 • All: 所有执行权限 • Read: 可读权限 • Write: 可写权限 • Temporary UDF Admin: 临时UDF管理权限 • Select/Deselect All: 全选/取消全选 <p>如需添加多条权限控制规则，可单击  按钮添加。</p> <p>如需当前条件中的用户或用户组管理本条策略，可勾选“Delegate Admin”，这些用户将成为受委托的管理员。被委托的管理员可以更新、删除本策略，它还可以基于原始策略创建子策略。</p>
Deny Conditions	<p>策略拒绝条件，配置本策略内拒绝的权限及例外，配置方法与“Allow Conditions”类型。</p>

表 24-14 设置权限

任务场景	角色授权操作
role admin操作	<ol style="list-style-type: none"> 1. 在首页中单击“Settings”，选择“Roles > Add New Role”。 2. 设置“Role Name”为“admin”，在“Users”区域，单击“Select User”，选择对应用户名。 3. 单击Add Users按钮，在对应用户名所在行勾选“Is Role Admin”，单击“Save”保存配置。 <p>说明 用户绑定Hive管理员角色后，在每个维护操作会话中，还需要执行以下操作：</p> <ol style="list-style-type: none"> 1. 以客户端安装用户，登录安装Hive客户端的节点。 2. 执行以下命令配置环境变量。 例如，Spark2x客户端安装目录为“/opt/client”，执行 source /opt/client/bigdata_env 3. 执行以下命令认证用户。 kinit Spark2x业务用户 4. 执行以下命令登录客户端工具。 spark-beeline 5. 执行以下命令更新用户的管理员权限。 set role admin;
创建库表操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写并选择对应的数据库（如果是创建库，需填写将要创建的库名称，或填写“*”表示任意名称的数据库，然后选择所写名称），在“table”与“column”右侧填写并选择对应的表名称、列名称，均支持通配符（“*”）匹配。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“Create”。
删除库表操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写并选择对应的数据库（如果是删除库，需填写将要创建的库名称，或填写“*”表示任意名称的数据库，然后选择所写名称），在“table”与“column”右侧填写并选择对应的表名称、列名称，均支持通配符（“*”）匹配。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“Drop”。 <p>说明 对于CarbonData表，只有对应表的OWNER，才能执行“drop”操作。</p>

任务场景	角色授权操作
ALTER操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写并选择对应的数据库，在“table”右侧填写并选择对应的表，在“column”右侧填写并选择对应的列名称，支持通配符（“*”）匹配。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“Alter”。
LOAD操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写并选择对应的数据库，在“table”右侧填写并选择对应的表，在“column”右侧填写并选择对应的列名称，支持通配符（“*”）匹配。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“update”。
INSERT操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写并选择对应的数据库，在“table”右侧填写并选择对应的表，在“column”右侧填写并选择对应的列名称，支持通配符（“*”）匹配。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“update”。 5. 用户还需要具有Yarn任务队列的“submit-app”权限，默认情况下，hadoop用户组具有向所有Yarn任务队列“submit-app”权限。具体配置请参考添加Yarn的Ranger访问权限策略。
GRANT操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写并选择对应的数据库，在“table”右侧填写并选择对应的表，在“column”右侧填写并选择对应的列名称，支持通配符（“*”）匹配。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 勾选“Delegate Admin”。
ADD JAR操作	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. 单击“database”并在下拉菜单中选择“global”。在“global”右侧填写并选择“*”。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，勾选“Temporary UDF Admin”。

任务场景	角色授权操作
VIEW与INDEX权限	<ol style="list-style-type: none"> 1. 在“Policy Name”填写策略名称。 2. “database”右侧填写并选择对应的数据库，在“table”右侧填写并选择对应的VIEW或INDEX名称，在“column”右侧填写并选择“*”。 3. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 4. 单击“Add Permissions”，参照表格上述相关操作，根据需求，给用户勾选相应权限。
其他用户库表操作	<ol style="list-style-type: none"> 1. 参照表格上述操作添加对应权限。 2. 给当前用户添加其他用户库表的HDFS路径的读、写、执行权限，具体配置请参考添加HDFS的Ranger访问权限策略。

📖 说明

在Ranger上为用户添加Spark SQL的访问策略后，需要在HDFS的访问策略中添加相应的路径访问策略，否则无法访问数据文件，具体请参考[添加HDFS的Ranger访问权限策略](#)。

- Ranger策略中global策略仅用于联合Temporary UDF Admin权限，用来控制UDF包的上传。
- 通过Ranger对Spark SQL进行权限控制时，不支持empower语法。
- 开启Ranger鉴权后，对视图操作时，默认需要具备相关表的权限，如果需要对视图进行独立鉴权，不依赖相关表的权限，需要将参数spark.ranger.plugin.viewaccesscontrol.enable设置为true。
 - 使用非Spark-beeline方式提交作业时，需要在“客户端安装目录/Spark/spark/conf/spark-defaults.conf”中设置该参数。
 - 使用Spark-beeline方式提交作业时，需要在“客户端安装目录/Spark/spark/conf/spark-defaults.conf”和服务端“JDBCserver > 自定义配置”中设置该参数。

步骤5 单击“Add”，在策略列表可查看策略的基本信息。等待策略生效后，验证相关权限是否正常。

如果需要禁用某条策略，可单击  按钮编辑该策略，设置策略开关为“Disabled”。

如果不再使用某条策略，可单击  按钮删除该策略。

----结束

Spark2x 表数据脱敏

Ranger支持对Spark2x数据进行脱敏处理（Data Masking），可对用户执行的select操作的返回结果进行处理，以屏蔽敏感信息。

步骤1 修改服务端和客户端spark.ranger.plugin.masking.enable参数值为true。

- 服务端：登录FusionInsight Manage页面，选择“集群 > 服务 > Spark2x > 配置 > 全部配置”，搜索并修改所有的spark.ranger.plugin.masking.enable参数值为true，保存配置并重启服务。

- 客户端：登录Spark客户端节点，进入目录“{客户端安装目录}/Spark/spark/conf/spark-defaults.conf”，修改spark.ranger.plugin.masking.enable参数值为true。
- 步骤2** 登录Ranger WebUI界面，在首页单击“HADOOP SQL”区域的组件插件名称如“Hive”。
- 步骤3** 在“Masking”页签单击“Add New Policy”，添加Spark2x权限控制策略。
- 步骤4** 根据业务需求配置相关参数。

表 24-15 Spark2x 数据脱敏参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，您可以根据这些标签搜索报告和筛选策略。
Hive Database	配置当前策略适用的Spark2x中的数据库名称。
Hive Table	配置当前策略适用的Spark2x中的表名称。
Hive Column	配置当前策略适用的Spark2x中的列名称。
Description	策略描述信息。
Audit Logging	是否审计此策略。
Mask Conditions	<p>在“Select Group”、“Select User”列选择已创建好的需要授予权限的用户组或用户，单击“Add Conditions”，添加策略适用的IP地址范围，然后再单击“Add Permissions”，勾选“select”权限。</p> <p>单击“Select Masking Option”，选择数据脱敏时的处理策略：</p> <ul style="list-style-type: none"> • Redact：用x屏蔽所有字母字符，用0屏蔽所有数字字符。 • Partial mask: show last 4：只显示最后的4个字符。 • Partial mask: show first 4：只显示开始的4个字符。 • Hash：对数据进行Hash处理。 • Nullify：用NULL值替换原值。 • Unmasked(retain original value)：不脱敏，显示原数据。 • Date: show only year：日期格式数据只显示年份信息。 • Custom：可使用任何有效Hive UDF（返回与被屏蔽的列中的数据类型相同的数据类型）来自定义策略。 <p>如需添加多列的脱敏策略，可单击  按钮添加。</p>

参数名称	描述
Deny Conditions	策略拒绝条件，配置本策略内拒绝的权限及例外，配置方法与“Allow Conditions”类型。

----结束

Spark2x 行级别数据过滤

Ranger支持用户对Spark2x数据表执行select操作时进行行级别的数据过滤。

步骤1 修改服务端和客户端spark.ranger.plugin.rowfilter.enable参数值为true。

- 服务端：登录FusionInsight Manage页面，选择“集群 > 服务 > Spark2x > 配置 > 全部配置”，搜索并修改所有的spark.ranger.plugin.rowfilter.enable参数值为true，保存配置并重启服务。
- 客户端：登录Spark客户端节点，进入目录“{客户端安装目录}/Spark/spark/conf/spark-defaults.conf”，修改spark.ranger.plugin.rowfilter.enable参数值为true。


步骤2 登录Ranger WebUI界面，在首页单击“HADOOP SQL”区域的组件插件名称如“Hive”。

步骤3 在“Row Level Filter”页签单击“Add New Policy”，添加行数据过滤策略。

步骤4 根据业务需求配置相关参数。

表 24-16 Spark2x 行数据过滤参数

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Label	为当前策略指定一个标签，您可以根据这些标签搜索报告和筛选策略。
Hive Database	配置当前策略适用的Spark2x中的数据库名称。
Hive Table	配置当前策略适用的Spark2x中的表名称。
Description	策略描述信息。
Audit Logging	是否审计此策略。

参数名称	描述
Row Filter Conditions	<p>在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的对象，单击“Add Conditions”，添加策略适用的IP地址范围，然后再单击“Add Permissions”，勾选“select”权限。</p> <p>单击“Row Level Filter”，填写数据过滤规则。</p> <p>例如过滤表A中“name”列“zhangsan”行的数据，过滤规则为：name <> 'zhangsan'。更多信息可参考Ranger官方文档。</p> <p>如需添加更多规则，可单击  按钮添加。</p>

步骤5 单击“Add”，在策略列表可查看策略的基本信息。

步骤6 用户通过Spark2x客户端对配置了数据脱敏策略的表执行select操作，系统将对数据进行处理后进行展示。

----结束

24.4.7 添加 Kafka 的 Ranger 访问权限策略

操作场景

Ranger管理员可通过Ranger为Kafka用户配置Kafka主题的读、写、管理权限以及集群的管理权限，本章节以为用户“test”添加“test”主题的“生产”权限。

前提条件

- 已安装Ranger服务且服务运行正常。
- 已创建用户需要配置权限的用户、用户组或Role。

操作步骤

步骤1 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录 Ranger WebUI界面](#)。


步骤2 在首页中单击“KAFKA”区域的组件插件名称如“Kafka”。

步骤3 单击“Add New Policy”，添加Kafka权限控制策略。

步骤4 根据业务需求配置相关参数。

表 24-17 Kafka 权限参数

参数名称	描述
Policy Type	Access。
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。

参数名称	描述
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。
Policy Label	为当前策略指定一个标签，您可以根据这些标签搜索报告和筛选策略。
topic	配置当前策略适用的topic名，可以填写多个值。这里支持通配符，例如：test、test*、*。 “Include”策略适用于当前输入的对象，“Exclude”表示策略适用于除去当前输入内容之外的其他对象。
Description	策略描述信息。
Audit Logging	是否审计此策略。
Allow Conditions	<p>策略允许条件，配置本策略内允许的权限及例外，例外条件优先级高于正常条件。</p> <p>在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的Role、用户组或用户。</p> <p>单击“Add Conditions”，添加策略适用的IP地址范围，单击“Add Permissions”，添加对应权限。</p> <ul style="list-style-type: none"> ● Publish：生产权限。 ● Consume：消费权限。 ● Describe：查询权限。 ● Create：创建主题权限。 ● Delete：删除主题权限。 ● Describe Configs：查询配置权限。 ● Alter：修改topic的partition数量的权限。 ● Alter Configs：修改配置权限。 ● Select/Deselect All：全选/取消全选。 <p>如需添加多条权限控制规则，可单击  按钮添加。</p> <p>如需当前条件中的用户或用户组管理本条策略，可勾选“Delegate Admin”，这些用户将成为受委托的管理员。被委托的管理员可以更新、删除本策略，它还可以基于原始策略创建子策略。</p>
Deny Conditions	策略拒绝条件，配置本策略内拒绝的权限及例外，配置方法与“Allow Conditions”类型，拒绝条件的优先级高于“Allow Conditions”中配置的允许条件。

例如为用户“testuser”添加“test”主题的生产权限，配置如下：

图 24-5 Kafka 权限参数

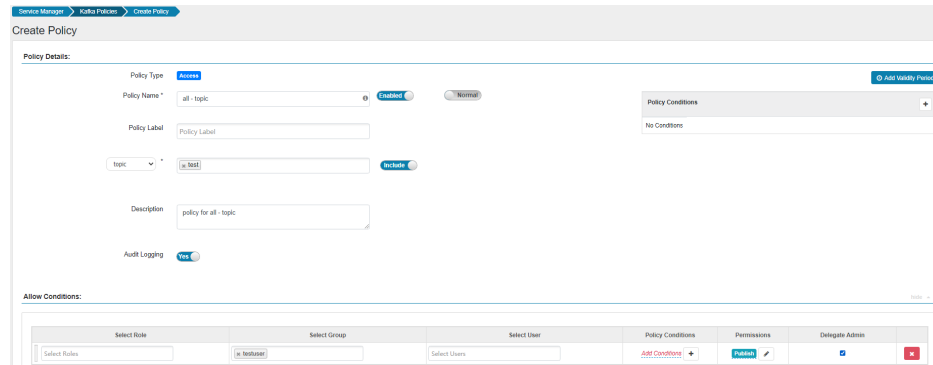






表 24-18 设置权限


任务场景	角色授权操作
设置Kafka管理员权限	<ol style="list-style-type: none"> 在首页中单击“KAFKA”区域的组件插件名称，例如“Kafka”。 选择“Policy Name”为“all - topic”的策略，单击  按钮编辑策略。 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 单击“Add Permissions”，勾选“Select/Deselect All”。
设置用户对Topic的创建权限	<ol style="list-style-type: none"> 在“topic”配置Topic名。 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 单击“Add Permissions”，勾选“Create”。 <p>说明 目前Kafka内核支持"--zookeeper"和"--bootstrap-server"两种方式创建Topic，社区将会在后续的版本中删掉对"--zookeeper"的支持，所以建议用户使用"--bootstrap-server"的方式创建Topic。 注意：目前Kafka只支持"--bootstrap-server"方式创建Topic行为的鉴权，不支持对"--zookeeper"方式的鉴权</p>
设置用户对Topic的删除权限	<ol style="list-style-type: none"> 在“topic”配置Topic名。 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 单击“Add Permissions”，勾选“Delete”。 <p>说明 目前Kafka内核支持"--zookeeper"和"--bootstrap-server"两种方式删除Topic，社区将会在后续的版本中删掉对"--zookeeper"的支持，所以建议用户使用"--bootstrap-server"的方式删除Topic。 注意：目前Kafka只支持对"--bootstrap-server"方式删除Topic行为的鉴权，不支持对"--zookeeper"方式的鉴权</p>



任务场景	角色授权操作
设置用户对Topic的查询权限	<ol style="list-style-type: none"> 1. 在“topic”配置Topic名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Describe”和“Describe Configs”。 <p>说明 目前Kafka内核支持"--zookeeper"和"--bootstrap-server"两种方式查询Topic，社区将会在后续的版本中删掉对"--zookeeper"的支持，所以建议用户使用"--bootstrap-server"的方式查询Topic。 注意：目前Kafka只支持对"--bootstrap-server"方式查询Topic行为的鉴权，不支持对"--zookeeper"方式的鉴权</p>
设置用户对Topic的生产权限	<ol style="list-style-type: none"> 1. 在“topic”配置Topic名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Publish”。
设置用户对Topic的消费权限	<ol style="list-style-type: none"> 1. 在“topic”配置Topic名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Consume”。 <p>说明 因为消费Topic时，涉及到Offset的管理操作，必须同时开启ConsumerGroup的“Consume”权限，详见“设置用户对ConsumerGroup Offsets 的提交权限”</p>
设置用户对Topic的扩容权限（增加分区）	<ol style="list-style-type: none"> 1. 在“topic”配置Topic名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Alter”。
设置用户对Topic的配置修改权限	当前Kafka内核暂不支持基于“--bootstrap-server”的Topic参数修改行为，故当前Ranger不支持对此行为的鉴权操作。
设置用户对Cluster的所有管理权限	<ol style="list-style-type: none"> 1. 在“cluster”右侧输入并选择集群名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Kafka Admin”。

任务场景	角色授权操作
设置用户对Cluster的创建权限	<ol style="list-style-type: none"> 1. 在首页中单击“KAFKA”区域的组件插件名称，例如“Kafka”。 2. 选择“Policy Name”为“all - cluster”的策略，单击  按钮编辑策略。 3. 在“cluster”右侧输入并选择集群名。 4. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 5. 单击“Add Permissions”，勾选“Create”。 <p>说明 对于Cluster的Create操作鉴权主要涉及以下两个场景：</p> <ol style="list-style-type: none"> 1. 集群开启了“auto.create.topics.enable”参数后，客户端向服务的还未创建的Topic发送数据的场景，此时会判断用户是否有集群的Create权限 2. 对于用户创建大量Topic的场景，如果授予用户Cluster Create权限，那么该用户可以在集群内部创建任意Topic
设置用户对Cluster的配置修改权限	<ol style="list-style-type: none"> 1. 在“cluster”右侧输入并选择集群名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Alter Configs”。 <p>说明 此处的配置修改权限，指的是Broker、Broker Logger的配置权限。 当授予用户配置修改权限后，即使不授予配置查询权限也可查询配置详情（配置修改权限高于且包含配置查询权限）。</p>
设置用户对Cluster的配置查询权限	<ol style="list-style-type: none"> 1. 在“cluster”右侧输入并选择集群名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Describe”和“Describe Configs”。 <p>说明 此处查询指的是查询集群内的Broker、Broker Logger信息。该查询不涉及Topic。</p>
设置用户对Cluster的Idempotent Write权限	<ol style="list-style-type: none"> 1. 在“cluster”右侧输入并选择集群名。 2. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 3. 单击“Add Permissions”，勾选“Idempotent Write”。 <p>说明 此权限会对用户客户端的Idempotent Produce行为进行鉴权。</p>


任务场景	角色授权操作
设置用户对Cluster的分区迁移权限管理	<ol style="list-style-type: none"> 在“cluster”右侧输入并选择集群名。 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 单击“Add Permissions”，勾选“Alter”。 <p>说明 Cluster的Alter权限可以对以下三种场景进行权限控制：</p> <ol style="list-style-type: none"> Partition Reassign场景下，迁移副本的存储目录。 集群里各分区内部leader选举。 Acl管理（添加或删除）。 <p>其中步骤4.1和步骤4.2都是集群内部Controller与Broker间、Broker与Broker间的操作，创建集群时，默认授予内置kafka用户此权限，普通用户授予此权限没有意义。</p> <p>步骤4.3涉及Acl的管理，Acl设计的就是用于鉴权，由于目前kafka鉴权已全部托管给Ranger，所以这个场景也基本不涉及（配置后亦不生效）。</p>
设置用户对Cluster的Cluster Action权限	<ol style="list-style-type: none"> 在“cluster”右侧输入并选择集群名。 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 单击“Add Permissions”，勾选“Cluster Action”。 <p>说明 此权限主要对集群内部副本主从同步、节点间通信进行控制，在集群创建时已经授权给内置kafka用户，普通用户授予此权限没有意义。</p>
设置用户对TransactionalId的权限	<ol style="list-style-type: none"> 在首页中单击“KAFKA”区域的组件插件名称，例如“Kafka”。 选择“Policy Name”为“all - transactionalid”的策略，单击按钮编辑策略。 在“transactionalid”配置事务ID。 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 单击“Add Permissions”，勾选“Publish”和“Describe”。 <p>说明 “Publish”权限主要对用户开启了事务特性的客户端请求进行鉴权，例如事务开启、结束、提交offset、事务性数据生产等行为。</p> <p>“Describe”权限主要对于开启事务特性的客户端与Coordinator的请求进行鉴权。</p> <p>建议在开启事务特性的场景下，给用户同时授予“Publish”和“Describe”权限。</p>


任务场景	角色授权操作
设置用户对 DelegationToken 的权限	<ol style="list-style-type: none"> 1. 在首页中单击“KAFKA”区域的组件插件名称，例如“Kafka”。 2. 选择“Policy Name”为“all - delegationtoken”的策略，单击  按钮编辑策略。 3. 在“delegationtoken”配置delegationtoken。 4. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 5. 单击“Add Permissions”，勾选“Describe”。 <p>说明 当前Ranger对DelegationToken的鉴权控制仅限于对查询的权限控制，不支持对DelegationToken的create、renew、expire操作的权限控制。</p>
设置用户对ConsumerGroup Offsets 的查询权限	<ol style="list-style-type: none"> 1. 在首页中单击“KAFKA”区域的组件插件名称，例如“Kafka”。 2. 选择“Policy Name”为“all - consumergroup”的策略，单击  按钮编辑策略。 3. 在“consumergroup”配置需要管理的consumergroup。 4. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 5. 单击“Add Permissions”，勾选“Describe”。
设置用户对ConsumerGroup Offsets 的提交权限	<ol style="list-style-type: none"> 1. 在首页中单击“KAFKA”区域的组件插件名称，例如“Kafka”。 2. 选择“Policy Name”为“all - consumergroup”的策略，单击  按钮编辑策略。 3. 在“consumergroup”配置需要管理的consumergroup。 4. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 5. 单击“Add Permissions”，勾选“Consume”。 <p>说明 当给用户授予了ConsumerGroup的“Consume”权限后，用户会同时被授予“Describe”权限。</p>

任务场景	角色授权操作
设置用户对ConsumerGroup Offsets 的删除权限	<ol style="list-style-type: none"> 1. 在首页中单击“KAFKA”区域的组件插件名称，例如“Kafka”。 2. 选择“Policy Name”为“all - consumergroup”的策略，单击  按钮编辑策略。 3. 在“consumergroup”配置需要管理的consumergroup。 4. 在“Allow Conditions”区域，单击“Select User”下选择框选择用户。 5. 单击“Add Permissions”，勾选“Delete”。 <p>说明 当给用户授予了ConsumerGroup的“Delete”权限后，用户会同时被授予“Describe”权限。</p>

步骤5 （可选）添加策略有效期。在页面右上角单击“Add Validity period”，设置“Start Time”和“End Time”，选择“Time Zone”。单击“Save”保存。如需添加多条策略有效期，可单击  按钮添加。如需删除策略有效期，可单击  按钮删除。

步骤6 单击“Add”，在策略列表可查看策略的基本信息。等待策略生效后，验证相关权限是否正常。

如需禁用某条策略，可单击  按钮编辑策略，设置策略开关为“Disabled”。

如果不再使用策略，可单击  按钮删除策略。

----结束

24.4.8 添加 Storm 的 Ranger 访问权限策略

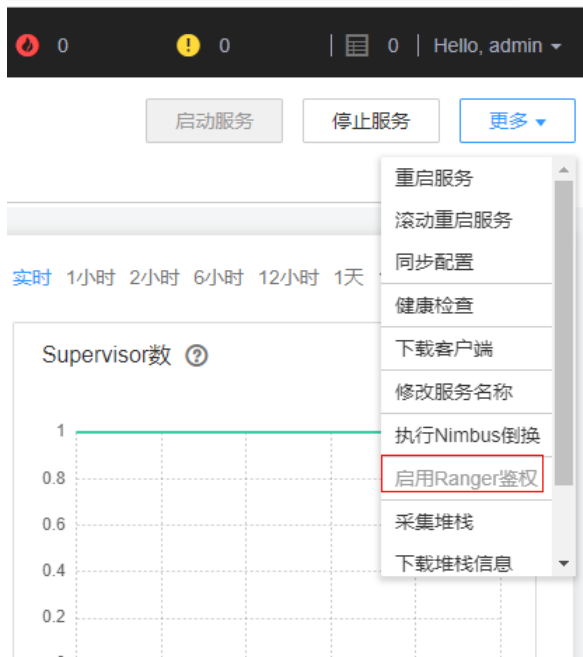
操作场景

Ranger管理员可通过Ranger为Storm用户进行相关的权限设置。

前提条件

- 已安装Ranger服务且服务运行正常。
- 已创建用户需要配置权限的用户、用户组或Role。
- 页面已启用Ranger鉴权开关，该按钮控制是否启用Ranger插件进行权限管控，启用则使用Ranger鉴权，否则使用组件自身鉴权机制。

图 24-6 启用 Ranger 鉴权






操作步骤

- 步骤1** 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录Ranger WebUI界面](#)。
- 步骤2** 在首页中单击“STORM”区域的“Storm”。
- 步骤3** 单击“Add New Policy”，添加Storm权限控制策略。
- 步骤4** 根据业务需求配置相关参数。

表 24-19 Storm 权限参数


参数名称	描述
Policy Conditions	IP过滤策略，可自定义，配置当前策略适用的主机节点，可填写一个或多个IP或IP段，并且IP填写支持“*”通配符，例如：192.168.1.10,192.168.1.20或者192.168.1.*。
Policy Name	策略名称，可自定义，不能与本服务内其他策略名称重复。 “include”策略适用于当前输入的对象， “exclude”表示策略适用于除去当前输入内容之外的其他对象。
Policy Label	为当前策略指定一个标签，您可以根据这些标签搜索报告和筛选策略。
Storm Topology	配置当前策略适用的拓扑名称。可以填写多个值。
Description	策略描述信息。

参数名称	描述
Audit Logging	是否审计此策略。
Allow Conditions	<p>策略允许条件，配置本策略内允许的权限及例外。在“Select Role”、“Select Group”、“Select User”列选择已创建好的需要授予权限的Role、用户组或用户，单击“Add Conditions”，添加策略适用的IP地址范围，单击“Add Permissions”，添加对应权限。</p> <ul style="list-style-type: none"> Submit Topology: 提交拓扑。 <p>说明 Submit Topology权限只有在Storm Topology为*的情况下可以赋权生效。</p> <ul style="list-style-type: none"> File Upload: 文件上传。 File Download: 文件下载。 Kill Topology: 删除拓扑。 Rebalance: Rebalance操作权限。 Activate: 激活权限。 Deactivate: 去激活权限。 Get Topology Conf: 获取拓扑配置。 Get Topology: 获取拓扑。 Get User Topology: 获取用户拓扑。 Get Topology Info: 获取拓扑信息。 Upload New Credential: 上传新的凭证。 Select/Deselect All: 全选/取消全选。 <p>如需添加多条权限控制规则，可单击  按钮添加。</p> <p>如需当前条件中的用户或用户组管理本条策略，可勾选“Delegate Admin”，这些用户将成为受委托的管理员。被委托的管理员可以更新、删除本策略，它还可以基于原始策略创建子策略。</p>
Deny Conditions	策略拒绝条件，配置本策略内拒绝的权限及例外，配置方法与“Allow Conditions”类似。

步骤5 （可选）添加策略有效期。在页面右上角单击“Add Validity period”，设置“Start Time”和“End Time”，选择“Time Zone”。单击“Save”保存。如需添加多条策略有效期，可单击  按钮添加。如需删除策略有效期，可单击  按钮删除。

步骤6 单击“Add”，在策略列表可查看策略的基本信息。等待策略生效后，验证相关权限是否正常。

如需禁用某条策略，可单击  按钮编辑策略，设置策略开关为“Disabled”。

如果不再使用策略，可单击  按钮删除策略。

----结束

24.5 查看 Ranger 审计信息

Ranger管理员可通过Ranger界面查看Ranger运行审计日志及组件使用Ranger鉴权后权限管控审计日志信息。

查看 Ranger 审计信息内容

- 步骤1** 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录 Ranger WebUI界面](#)。
- 步骤2** 单击“Audit”，查看相关审计信息，各页签内容说明请参考[表24-20](#)，条目较多时，单击搜索框可根据关键字字段进行筛选。

表 24-20 Audit 信息

页签	内容描述
Access	当前MRS不支持在线查看组件资源的审计日志信息，可登录组件安装节点，进入“/var/log/Bigdata/audit”目录下查看各组件的审计日志。
Admin	Ranger上操作审计信息，例如安全访问策略的创建/更新/删除、组件权限策略的创建/删除、role的创建/更新/删除等。
Login Sessions	登录Ranger的用户会话审计信息。
Plugins	Ranger内组件权限策略信息。
Plugin Status	各组件节点权限策略的同步审计信息。
User Sync	Ranger与LDAP用户同步审计信息。

----结束

24.6 配置 Ranger 安全区信息

Ranger支持配置安全区，Ranger管理员可将各组件的资源切分为多个安全区，由对应Ranger管理员用户为区域的指定资源设置安全策略，以便更好的细分资源管理。安全区中定义的策略仅适用于区域中的资源，服务的资源被划分到安全区后，非安全区针对该资源的访问权限策略将不再生效。安全区的管理员只能在其作为管理员的安全区中设置策略。

添加安全区

- 步骤1** 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录 Ranger WebUI界面](#)。


步骤2 单击“Security Zone”，在区域列表页面中单击 ，添加安全区。

表 24-21 安全区配置参数

参数名称	描述	示例
Zone Name	配置安全区的名称。	test
Zone Description	配置安全区的描述信息。	-
Admin Users/ Admin Usergroups	配置安全区的管理用户/用户组，可在安全区中添加及修改相关资源的权限策略。 必须至少配置一个用户或用户组。	zone_admin
Auditor Users/ Auditor Usergroups	添加审计用户/用户组，可在安全区中查看相关资源权限策略内容。 必须至少配置一个用户或用户组。	zone_user
Select Tag Services	选择服务的标签信息。	-
Select Resource Services	选择安全区内包含的服务及具体资源。 在“Select Resource Services”中选择服务后，需要在“Resource”列中添加具体的资源对象，例如HDFS服务器的文件目录、Yarn的队列、Hive的数据库及表、HBase的表及列。	/testzone

例如针对HDFS中的“/testzone”目录创建一个安全区，配置如下：

Zone Details :

Zone Name *

Zone Description

Zone Administration :

Admin Users

Admin Usergroups

Auditor Users

Auditor Usergroups

Services :

Select Tag Services

Select Resource Services *

Service Name	Service Type	Resource
hacluster	HDFS	path: /testzone <input type="checkbox"/> <input type="checkbox"/>

步骤3 单击“Save”，等待安全区添加成功。

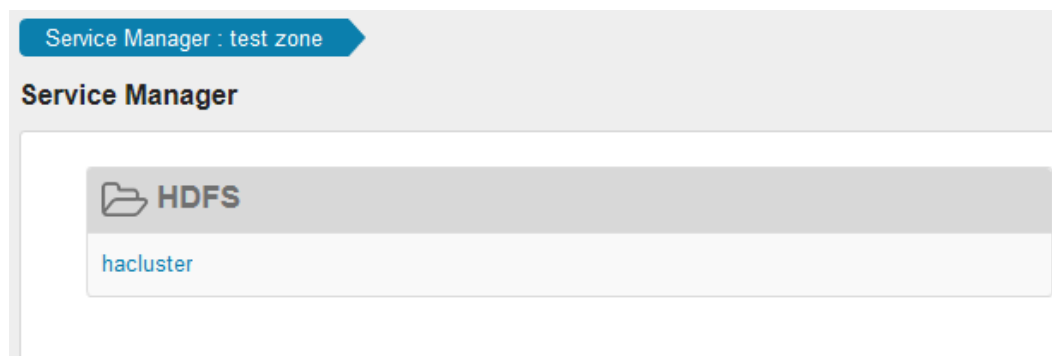
Ranger管理员可在“Security Zone”页面查看当前的所有安全区并单击“Edit”修改安全区的属性信息，当相关资源不需要在安全区中进行管理时，可单击“Delete”删除对应安全区。

----结束

在安全区中配置权限策略

步骤1 使用Ranger安全区管理员用户登录Ranger管理页面。

步骤2 在Ranger首页右上角的“Security Zone”选项的下拉列表中选择对应的安全区，即可切换至该安全区内的权限视图。



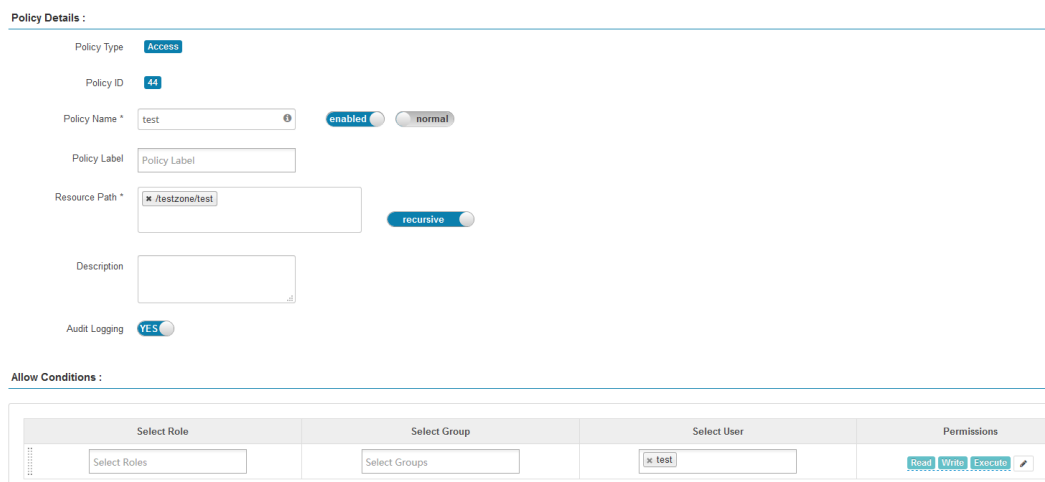
步骤3 单击组件名称下的权限插件名称，即可进入组件安全访问策略列表页面。

说明

各组件的策略列表中，系统默认生成的条目会自动继承至安全区内，用于保证集群内的部分系统默认用户或用户组的权限。

步骤4 单击“Add New Policy”，根据业务场景规划配置相关用户或者用户组的资源访问策略。

例如在本章节样例中，在安全区内配置一条允许“test”用户访问“/testzone/test”目录的策略：



其他不同组件的完整访问策略配置样例参考：

- [添加HDFS的Ranger访问权限策略](#)
- [添加HBase的Ranger访问权限策略](#)
- [添加Hive的Ranger访问权限策略](#)
- [添加Yarn的Ranger访问权限策略](#)
- [添加Spark2x的Ranger访问权限策略](#)
- [添加Kafka的Ranger访问权限策略](#)
- [添加Storm的Ranger访问权限策略](#)

策略添加后，需等待30秒左右，待系统生效。

📖 说明

- 安全区中定义的策略仅适用于区域中的资源，服务的资源被划分到安全区后，非安全区针对该资源的访问权限策略将不再生效。
- 如需配置针对当前安全区之外其他资源的访问策略，需在Ranger首页右上角的“Security Zone”选项中退出当前安全区后进行配置。

---结束

24.7 普通集群修改 Ranger 数据源为 Ldap

安全集群Ranger数据源默认为FusionInsight Manager Ldap用户。普通集群Ranger数据源默认为集群Unix用户。

前提条件

- 集群模式为普通模式。
- 已安装Ranger组件。

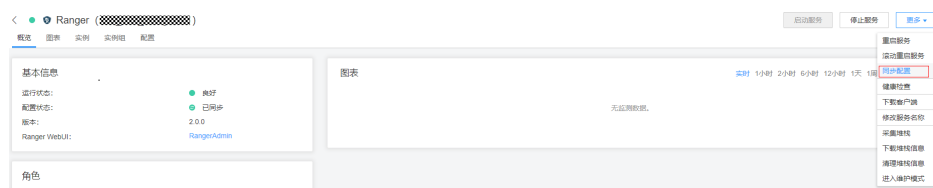
操作步骤

步骤1 登录FusionInsight Manager页面，具体请参见[访问FusionInsight Manager \(MRS 3.x及之后版本\)](#)。选择“集群 > 服务 > Ranger > 配置 > 全部配置 > UserSync (角色) > 自定义”。

步骤2 在“ranger.usersync.config.expandor”参数配置中添加“ranger.usersync.sync.source”值为“ldap”和“ranger.usersync.cookie.enabled”值为“false”，如下图所示：

名称	值
ranger.usersync.config.expandor	ranger.usersync.sync.source ldap
	ranger.usersync.cookie.enabled false

步骤3 在Ranger服务“概览”页面右上角单击“更多”，选择“同步配置”。



步骤4 在Ranger实例页面，勾选“UserSync”实例，选择“更多 > 重启实例”。



步骤5 在Ranger服务“概览”页面，单击“RangerAdmin”，查看“Settings > Users/Groups/Roles”页面是否有ldap用户。

----结束

24.8 查看 Ranger 用户权限同步信息

查看Ranger相关权限设置信息，例如查看用户、用户组、Role。

查看 Ranger 权限信息

步骤1 使用Ranger管理员用户rangeradmin登录Ranger管理页面，具体操作可参考[登录 Ranger WebUI界面](#)。

步骤2 选择“Settings > Users/Groups/Roles”，可查看系统中的用户、用户组、Roles信息。

- Users: 显示Ranger从LDAP或者OS同步的所有用户信息。
- Groups: 显示Ranger从LDAP或者OS同步的所有用户组、角色信息。
- Roles: 显示Ranger中创建的Role信息。

📖 说明

- 在FusionInsight Manager上创建的用户、角色、用户组会定期自动同步至Ranger，默认周期为300000毫秒（5分钟）。FusionInsight Manager中的角色和用户组在同步至Ranger后都变为用户组（Group）。只有被用户关联了的角色和用户组才会自动同步至Ranger。
- Ranger界面中创建的Role为用户或用户组的集合，用于灵活设置组件的权限访问策略，与FusionInsight Manager中的“角色”不同，请注意区分。

----结束

调整 Ranger 用户类型

步骤1 登录Ranger管理页面。

调整Ranger用户类型须使用Admin类型的用户（例如**admin**）进行操作，具体用户类型请参考[Ranger用户类型](#)。

步骤2 选择“Settings > Users/Groups/Roles”，在“Users”用户列表中，单击待修改类型的用户名。

步骤3 设置“Select Role”配置项为待修改的类型。

步骤4 单击“Save”。

----结束

创建 Ranger Role

Ranger管理员在设置组件的权限访问策略时，可基于用户、用户组或者Role灵活配置，其中用户与用户组信息从LDAP中自动同步，Role可手动添加。

步骤1 登录Ranger管理页面。

步骤2 选择“Settings > Users/Groups/Roles > Roles > Add New Role”。

步骤3 根据界面提示填写Role的名称与描述信息。

步骤4 添加Role内需要包含的用户、用户组、子Role信息。

- 在“Users”区域，选择系统中已创建的用户，然后单击“Add Users”。
- 在“Groups”区域，选择系统中已创建的用户组，然后单击“Add Group”。
- 在“Roles”区域，选择系统中已创建的Role，然后单击“Add Role”。

Users:

User Name	Is Role Admin	Action
test01	<input type="checkbox"/>	<input type="button" value="✖"/>

Select User

Groups:

Group Name	Is Role Admin	Action
hadoop	<input type="checkbox"/>	<input type="button" value="✖"/>

Select Group

Roles:

Role Name	Is Role Admin	Action
admin	<input type="checkbox"/>	<input type="button" value="✖"/>

Select Role

步骤5 单击“Save”，Role添加成功。

📖 说明

新创建的role，页面不提供删除操作，可以修改。

----结束

24.9 Ranger 日志介绍

日志描述

日志存储路径：Ranger相关日志的默认存储路径为“/var/log/Bigdata/ranger/角色名”

- RangerAdmin：“/var/log/Bigdata/ranger/rangeradmin”（运行日志）。
- TagSync：“/var/log/Bigdata/ranger/tagsync”（运行日志）。
- UserSync“/var/log/Bigdata/ranger/usersync”（运行日志）。

日志归档规则：Ranger的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过20MB的时，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd_hh-mm-ss>.[编号].log.zip”，最多保留最近的20个压缩文件。

表 24-22 Ranger 日志列表

日志类型	日志文件名	描述
RangerAdmin运行日志	access_log.<DATE>.log	Tomcat访问日志。
	catalina.out	Tomcat服务运行日志。
	gc-worker.log	RangerAdmin的GC日志。
	postinstallDetail.log	实例安装前启动后工作日志。

日志类型	日志文件名	描述
	prestartDetail.log	实例启动前准备工作日志。
	ranger-admin- <i><hostname></i> .log	RangerAdmin运行日志。
	ranger_admin_sql- <i><hostname></i> .log	RangerAdmin检索DBService的日志。
	startDetail.log	实例启动日志。
TagSync运行日志	cleanupDetail.log	实例清理日志。
	gc-worker.log	实例GC日志。
	postinstallDetail.log	实例安装前启动后工作日志。
	prestartDetail.log	实例启动前准备工作日志。
	ranger-tagsync- <i><hostname></i> .log	TagSync运行日志。
	startDetail.log	实例启动日志。
	tagsync.out	TagSync的运行日志。
UserSync运行日志	auth.log	unixauth服务运行日志。
	cleanupDetail.log	实例清理日志。
	gc-worker.log	实例GC日志。
	postinstallDetail.log	实例安装前启动后工作日志。
	prestartDetail.log	实例启动前准备工作日志。
	ranger-usersync- <i><hostname></i> .log	UserSync运行日志。
	startDetail.log	实例启动日志。

日志级别

HDFS中提供了如表24-23所示的日志级别，日志级别优先级从高到低分别是FATAL、ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 24-23 日志级别

级别	描述
FATAL	FATAL表示当前事件处理出现严重错误信息，可能导致系统崩溃。
ERROR	ERROR表示当前事件处理出现错误信息，系统运行出错。
WARN	WARN表示当前事件处理存在异常信息，但认为是正常范围，不会导致系统出错。
INFO	INFO表示记录系统及各事件正常运行状态信息
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 登录FusionInsight Manager。
- 步骤2** 选择“集群 > 服务 > Ranger > 配置”。
- 步骤3** 选择“全部配置”。
- 步骤4** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤5** 选择所需修改的日志级别。
- 步骤6** 单击“保存”，在弹出窗口中单击“确定”使配置生效。

 **说明**

配置完成后立即生效，不需要重启服务。

----**结束**

日志格式

Ranger的日志格式如下所示：

表 24-24 日志格式

日志类型	格式	示例
运行日志	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的 message> <日志事件的发生位置>	2020-04-29 20:09:28,543 INFO http-bio-21401-exec-56 Request comes from API call, skip cas filter. CasAuthenticationFilter Wrapper.java:25

24.10 Ranger 常见问题

24.10.1 安装集群过程中 Ranger 启动失败

问题

MRS集群创建过程中，Ranger启动失败，Manager进程任务列表里打印“ERROR: cannot drop sequence X_POLICY_REF_ACCESS_TYPE_SEQ”等关于数据库信息。

回答

该现象可能出现在安装两个RangerAmdin实例的场景下。

启动失败后，请先手动重启一个RangerAdmin，然后再逐步重启其他实例。

24.10.2 如何判断某个服务是否使用了 Ranger 鉴权

问题

如何判断某个支持使用Ranger鉴权的服务当前是否启用了Ranger鉴权？

回答

登录FusionInsight Manager，选择“集群 > 服务 > 服务名称”，在服务详情页上继续单击“更多”，查看“启用Ranger鉴权”是否为可单击？

- 是，表示当前本服务未启用Ranger鉴权插件，可单击“启用Ranger鉴权”启用该功能。
- 否，表示当前本服务已启用Ranger鉴权插件，可通过Ranger管理界面配置访问该服务资源的权限策略。

说明

若不存在该选项，则表示当前服务不支持Ranger鉴权插件，未开启Ranger鉴权。

24.10.3 新创建用户修改完密码后无法登录 Ranger

问题

使用新建用户登录Ranger页面，为什么在修改完密码后登录报401错误？

回答

由于UserSync同步用户数据有时间周期，默认是5分钟，因此在Manager上新创建的用户在用户同步成功前无法登录Ranger，因为Ranger的DB里暂时还没有该用户信息，需要等待同步周期所设置的时间后再尝试登录。

未开启Kerberos认证时，由于Ranger并不从Manager同步用户数据，因此，仅有admin用户可以登录Ranger，暂时不支持其他用户登录。

24.10.4 Ranger 界面添加或者修改 HBase 策略时，无法使用通配符搜索已存在的 HBase 表

问题


添加HBase的Ranger访问权限策略时，在策略中使用通配符搜索已存在的HBase表时，搜索不到已存在的表，并且在/var/log/Bigdata/ranger/rangeradmin/ranger-admin-*.log中报以下错误

```
Caused by: javax.security.sasl.SaslException: No common protection layer between client and server
at com.sun.security.sasl.gsskerb.GssKrb5Client.doFinalHandshake(GssKrb5Client.java:253)
at com.sun.security.sasl.gsskerb.GssKrb5Client.evaluateChallenge(GssKrb5Client.java:186)
at
org.apache.hadoop.hbase.security.AbstractHBaseSaslRpcClient.evaluateChallenge(AbstractHBaseSaslRpcClient.java:142)
at org.apache.hadoop.hbase.security.NettyHBaseSaslRpcClientHandler
$.run(NettyHBaseSaslRpcClientHandler.java:142)
at org.apache.hadoop.hbase.security.NettyHBaseSaslRpcClientHandler
$.run(NettyHBaseSaslRpcClientHandler.java:138)
at java.security.AccessController.doPrivileged(Native Method)
at javax.security.auth.Subject.doAs(Subject.java:422)
at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1761)
at
org.apache.hadoop.hbase.security.NettyHBaseSaslRpcClientHandler.channelRead0(NettyHBaseSaslRpcClientHandler.java:138)
at
org.apache.hadoop.hbase.security.NettyHBaseSaslRpcClientHandler.channelRead0(NettyHBaseSaslRpcClientHandler.java:42)
at
org.apache.hadoop.hbase.thirdparty.io.netty.channel.SimpleChannelInboundHandler.channelRead(SimpleChannelInboundHandler.java:105)
at
org.apache.hadoop.hbase.thirdparty.io.netty.channel.AbstractChannelHandlerContext.invokeChannelRead(AbstractChannelHandlerContext.java:362)
```

回答

Ranger界面上HBase服务插件的“hbase.rpc.protection”参数值和HBase服务端的“hbase.rpc.protection”参数值必须保持一致。

步骤1 参考[登录Ranger WebUI界面](#)章节，登录Ranger管理界面。

步骤2 在首页中“HBASE”区域，单击组件插件名称，如HBase的按钮

步骤3 搜索配置项“hbase.rpc.protection”，修改配置项的value值，与HBase服务端的“hbase.rpc.protection”的值保持一致。

步骤4 单击“保存”。

----结束

24.10.5 在 Ranger 管理界面查看不到创建的 MRS 用户

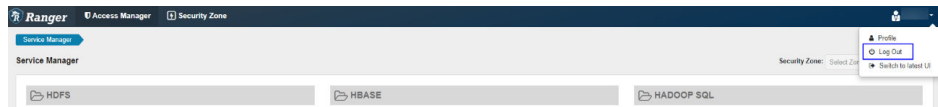
问题

在MRS Manager界面创建了一个账户，登录到Ranger管理界面后查看不到该用户。

回答

登录到Ranger管理界面的用户权限不够，需要切换至rangeradmin用户或者其他具有Ranger管理员权限的用户。

1. 在Ranger WebUI界面，单击右上角用户名，选择“Log Out”，退出当前用户。



2. 使用rangeradmin用户（默认密码为Rangeradmin@123）或者其他具有Ranger管理员权限用户重新登录。

24.10.6 MRS 用户无法同步至 Ranger 管理界面

问题

在MRS Manager界面创建了一个账户，在Ranger管理界面无法查看到该用户，重启UserSync后，可以查看到该用户。

回答

需要修改UserSync进程的GC内存默认为“-Xms1G -Xmx1G”，需要根据业务实际情况调整该参数值：

登录到MRS Manager界面，选择“集群 > 服务 > Ranger > 配置 > 全部配置 > UserSync（角色） > 系统”，修改参数“GC_OPTS”值。例如将内存修改为“-Xms2G -Xmx2G”。

25 使用 Spark（MRS 3.x 之前版本）

25.1 从零开始使用 Spark

本章节提供从零开始使用Spark提交sparkPi作业的操作指导，sparkPi是最经典的Spark作业，它用来计算Pi（ π ）值。

操作步骤

步骤1 准备sparkPi程序。

开源的Spark的样例程序包含多个例子，其中包含sparkPi。可以从<https://archive.apache.org/dist/spark/spark-2.1.0/spark-2.1.0-bin-hadoop2.7.tgz>中下载Spark的样例程序。

解压后在“spark-2.1.0-bin-hadoop2.7/examples/jars”路径下获取“spark-examples_2.11-2.1.0.jar”，即为Spark的样例程序。spark-examples_2.11-2.1.0.jar样例程序包含sparkPi程序。

步骤2 上传数据至OBS。

1. 登录OBS控制台。
2. 单击“并行文件系统 > 创建并行文件系统”，创建一个名称为sparkpi的文件系统。
sparkpi仅为示例，文件系统名称必须全局唯一，否则会创建并行文件系统失败。其他参数分别保持默认值。
3. 单击sparkpi文件系统名称，并选择“文件”。
4. 单击“新建文件夹”，分别创建program文件夹，创建完成后如[图25-1](#)所示。

图 25-1 文件夹列表



5. 进入program文件夹，单击上传文件，从本地选择步骤1中下载的程序包，“存储类别”选择“标准存储”。

步骤3 登录MRS控制台，在左侧导航栏选择“现有集群”，单击集群名称。

步骤4 提交sparkPi作业。

在MRS控制台选择“作业管理”，单击“添加”，进入“添加作业”页面，具体请参见[运行Spark作业](#)。

图 25-2 sparkPi 作业

The screenshot shows the '添加作业' (Add Job) configuration page. It includes the following fields and options:

- 作业类型** (Job Type): SparkSubmit (dropdown menu)
- 作业名称** (Job Name): sparkPi (text input)
- 执行程序路径** (Executable Path): obs://[bucket]/program/spark-examples_2.11-2.1.0.jar (text input). Storage options: HDFS, OBS.
- 运行程序参数** (Run Program Parameters): --class (text input), org.apache.spark.examples.SparkPi (dropdown menu).
- 执行程序参数** (Executable Parameters): 10 (text input). Storage options: HDFS, OBS.
- 服务配置参数** (Service Configuration Parameters): 参数 (Parameter) and 值 (Value) (text inputs).
- 命令参考** (Command Reference): spark-submit --class org.apache.spark.examples.SparkPi --master yarn-cluster obs://[bucket]/program/spark-examples_2.11-2.1.0.jar 10

At the bottom, there are two buttons: '确定' (Confirm) and '取消' (Cancel).

- 作业类型选择“SparkSubmit”。
- 作业名称为“sparkPi”。
- 执行程序路径配置为OBS上存放程序的地址。例如：obs://sparkpi/program/spark-examples_2.11-2.1.0.jar。
- 运行程序参数选择“--class”，值填写“org.apache.spark.examples.SparkPi”。
- 执行程序参数中填写的参数为：10。
- 服务配置参数无需填写。

只有集群处于“运行中”状态时才能提交作业。

作业提交成功后默认为“已接受”状态，不需要用户手动执行作业。

步骤5 查看作业执行结果。

1. 进入“作业管理”页面，查看作业是否执行完成。
作业运行需要时间，作业运行结束后，刷新作业列表。
作业执行成功或失败后都不能再次执行，只能新增作业，配置作业参数后重新提交作业。
2. 进入Yarn原生界面，查看作业输出信息。
 - a. 进入“作业管理”页面，单击对应作业所在行“操作”列的“查看详情”，获取“作业实际编号”。

图 25-3 作业实际编号

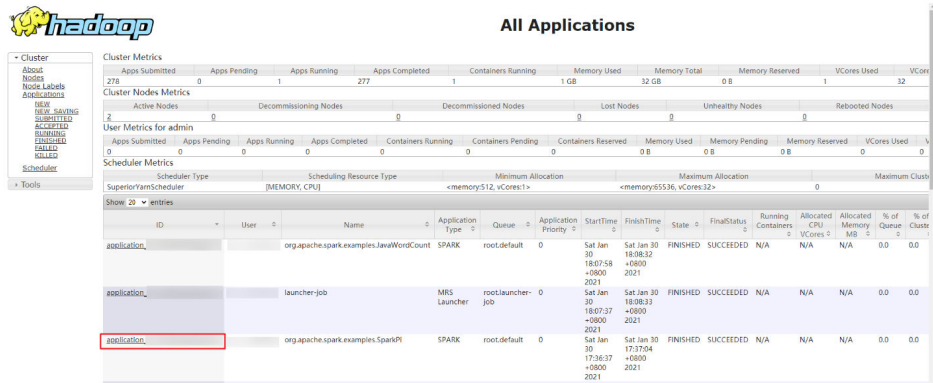
查看详情

作业类型	SparkSubmit
作业ID	5b7939ac-f4d4-4ed2-ac44-664baecdac74
Launcher作业编号	application_1611887005648_0274
实际作业编号	application_1611887005648_0275
作业提交时间	2021/01/30 17:35:56 GMT+08:00
作业开始时间	2021/01/30 17:36:37 GMT+08:00
作业结束时间	2021/01/30 17:37:04 GMT+08:00
作业进度	100%
作业执行时长	0.5 分钟
作业状态	已完成

确定

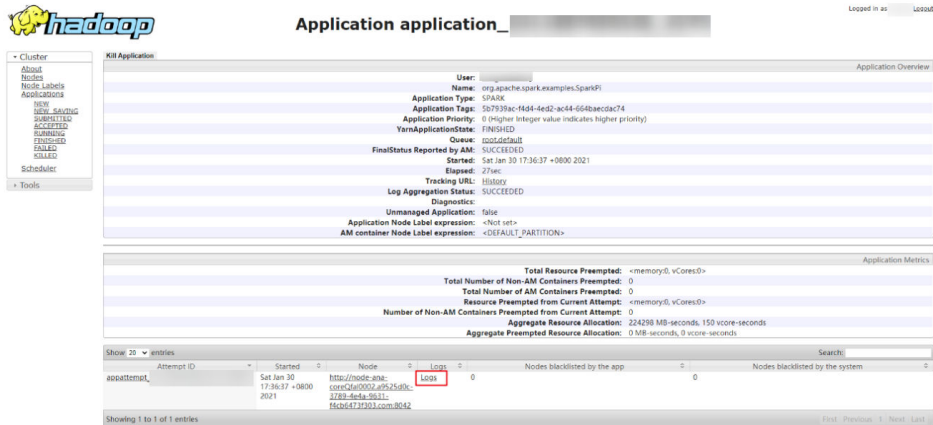
- b. 登录Manager页面，选择“服务管理 > Yarn > ResourceManager WebUI > ResourceManager (主)”进入Yarn界面。
- c. 单击“作业实际编号”对应ID。

图 25-4 Yarn 界面



d. 单击作业日志中的“Logs”。

图 25-5 sparkPi 作业日志



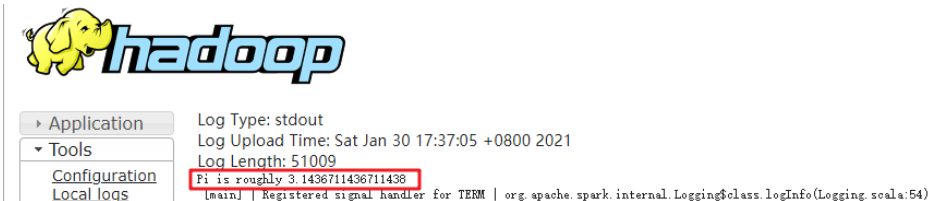
e. 单击“here”获取更详细日志。

图 25-6 sparkPi 作业更详细日志

Log Type: stdout
 Log Upload Time: Sat Jan 30 17:37:05 +0800 2021
 Log Length: 51009
 Showing 4096 bytes of 51009 total. Click [here](#) for the full log.

f. 获取作业执行结果。

图 25-7 sparkPi 作业执行结果



----结束

25.2 从零开始使用 Spark SQL

Spark提供类似SQL的Spark SQL语言操作结构化数据，本章节提供从零开始使用Spark SQL，创建一个名称为src_data的表，然后在src_data表中每行写入一条数据，最后将数据存储在“mrs_20160907”集群中。再使用SQL语句查询src_data表中的数据，最后可将src_data表删除。

前提条件

将OBS数据源中的数据写入Spark SQL表中时，需要先获取AK/SK。获取方法如下：

1. 登录管理控制台。
2. 单击用户名，在下拉列表中单击“我的凭证”。
3. 单击“访问密钥”。
4. 单击“新增访问密钥”，进入“新增访问密钥”页面。
5. 输入登录密码和短信验证码，单击“确定”，下载密钥，请妥善保管。

操作步骤

步骤1 准备使用Spark SQL分析的数据源。

样例txt文件如下：

```
abcd3ghji  
efgh658ko  
1234jjyu9  
7h8kodfg1  
kk99icxz3
```

步骤2 上传数据至OBS。

1. 登录OBS控制台。
2. 单击“并行文件系统 > 创建并行文件系统”，创建一个名称为sparksql的文件系统。
sparksql仅为示例，文件系统名称必须全局唯一，否则会创建并行文件系统失败。
3. 单击sparksql文件系统名称，并选择“文件”。
4. 单击“新建文件夹”，创建input文件夹。
5. 进入input文件夹，单击“上传文件 > 添加文件”，选择本地的txt文件，然后单击“上传”。

步骤3 登录MRS控制台，在左侧导航栏选择“现有集群”，单击集群名称。

步骤4 将OBS中的txt文件导入至HDFS中。

1. 选择“文件管理”。
2. 在“HDFS文件列表”页签中单击“新建”，创建一个名称为用户input的文件夹。
3. 进入userinput文件夹，单击“导入数据”。
4. 选择OBS和HDFS路径，单击“确定”。

OBS路径：obs://sparksql/input/sparksql-test.txt

HDFS路径：/user/userinput

图 25-8 从 OBS 导入数据至 HDFS



步骤5 提交Spark SQL语句。

1. 在MRS控制台选择“作业管理”，具体请参见[运行Spark作业](#)。
只有“mrs_20160907”集群处于“运行中”状态时才能提交Spark SQL语句。
2. 输入创建表的Spark SQL语句。
输入Spark SQL语句时，总字符数应当小于或等于10000字符，否则会提交语句失败。

语法格式：

```
CREATE [EXTERNAL] TABLE [IF NOT EXISTS] table_name [(col_name data_type [COMMENT col_comment], ...)] [COMMENT table_comment] [PARTITIONED BY (col_name data_type [COMMENT col_comment], ...)] [CLUSTERED BY (col_name, col_name, ...) [SORTED BY (col_name [ASC|DESC], ...)] INTO num_buckets BUCKETS] [ROW FORMAT row_format] [STORED AS file_format] [LOCATION hdfs_path];
```

创建表样例存在以下两种方式。

- 方式一：创建一个src_data表，将数据源中的数据一行一行写入src_data表中。
 - 数据源存储在HDFS的“/user/omm/userinput”文件夹下：**create external table src_data(line string) row format delimited fields terminated by '\n' stored as textfile location '/user/omm/userinput';**
 - 数据源存储在OBS的“/sparksql/input”文件夹下：**create external table src_data(line string) row format delimited fields terminated by '\n' stored as textfile location 'obs://AK:SK@sparksql/input';**
AK/SK获取方法，请参见[前提条件](#)。
- 方式二：创建一个表src_data1，将数据源中的数据批量load到src_data1表中。

```
create table src_data1 (line string) row format delimited fields
terminated by ',';
load data inpath '/user/omm/userinput/sparksql-test.txt' into table
src_data1;
```

📖 说明

采用方式二时，只能将HDFS上的数据load到新建的表中，OBS上的数据不支持直接load到新建的表中。

3. 输入查询表的Spark SQL语句。

语法格式：

```
SELECT col_name FROM table_name;
```

查询表样例，查询src_data表中的所有数据：

```
select * from src_data;
```

4. 输入删除表的Spark SQL语句。

语法格式：

```
DROP TABLE [IF EXISTS] table_name;
```

删除表样例：

```
drop table src_data;
```

5. 单击“检查”，检查输入语句的语法是否正确。

6. 单击“确定”。

Spark SQL语句提交后，是否执行成功会在“执行结果”列中展示。

步骤6 删除集群。

----结束

25.3 使用 Spark 客户端

MRS集群创建完成后，可以通过客户端去创建和提交作业。客户端可以安装在集群内部节点或集群外部节点上：

- 集群内部节点：MRS集群创建完成后，集群内的master和core节点默认已经安装好客户端，详情请参见[集群内节点使用MRS客户端](#)章节，登录安装客户端的节点。
- 集群外部节点：用户可以将客户端安装在集群外部节点上，详情请参见[集群外节点使用MRS客户端](#)章节，登录安装客户端的节点。

使用 Spark 客户端

步骤1 登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 如果集群为安全模式，执行以下命令进行用户认证。普通模式集群无需执行用户认证。

kinit 组件业务用户

步骤5 直接执行Spark Shell命令。例如：

```
spark-beeline
```

```
----结束
```

25.4 访问 Spark Web UI 界面

Spark Web UI界面主要用于查看Spark应用程序运行情况，推荐使用Google chrome浏览器以获得更好的体验。

Spark主要有两个Web页面。

- Spark UI页面，用于展示正在执行的应用的运行情况。
页面主要包括了Jobs、Stages、Storage、Environment、Executors、SQL、JDBC/ODBC Server等部分。Streaming应用会多一个Streaming标签页。
- History Server页面，用于展示已经完成的和未完成的Spark应用的运行情况。
页面包括了应用ID、应用名称、开始时间、结束时间、执行时间、所属用户等信息。

Spark UI

步骤1 进入组件管理页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理”。

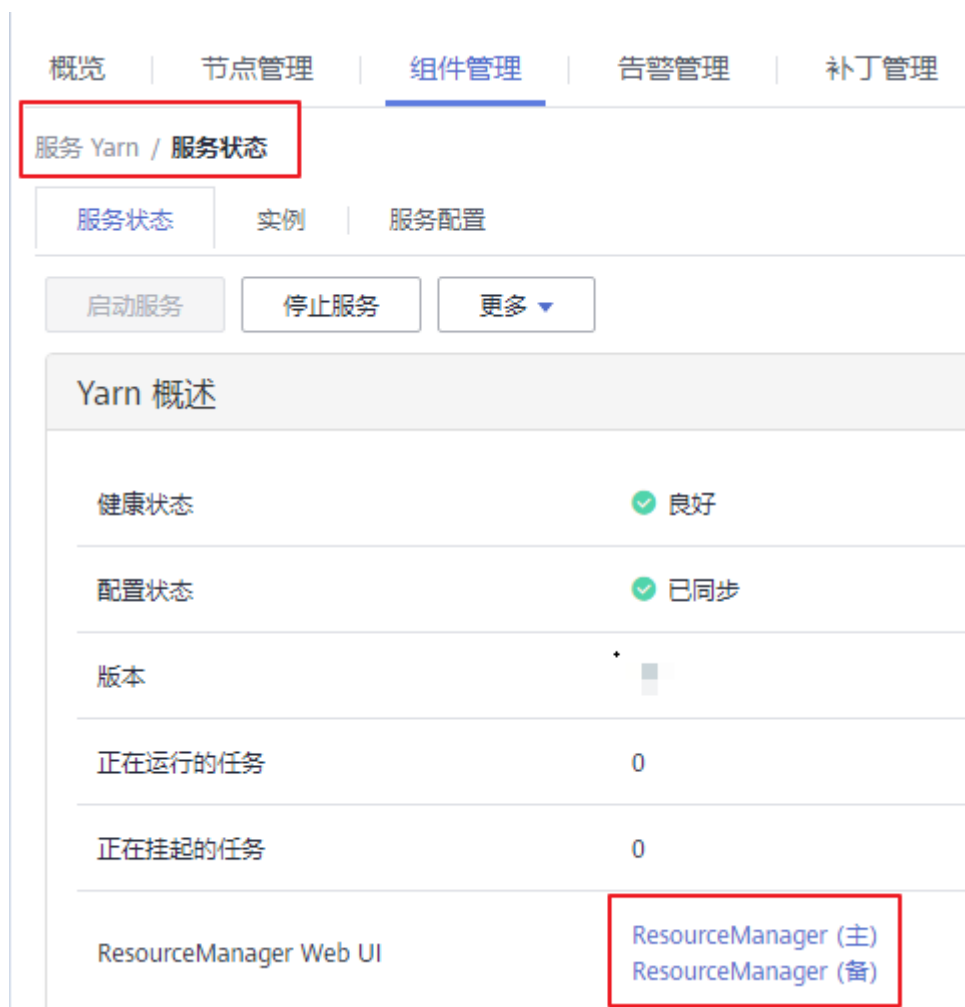
说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务”。

步骤2 选择“Yarn”并在“Yarn 概述”中“ResourceManager Web UI”中单击“ResourceManager Web UI”对应的“ResourceManager”进入Web界面。

图 25-9 ResourceManager Web UI



步骤3 查找到对应的Spark应用程序，单击应用信息的最后一列“ApplicationMaster”，即可进入Spark UI页面。

图 25-10 ApplicationMaster

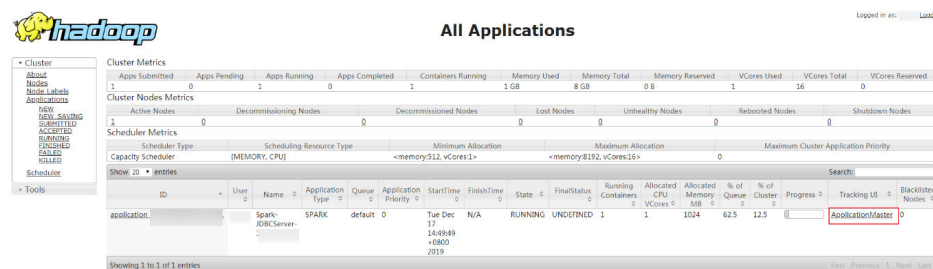


图 25-11 Spark UI 页面



----结束

History Server

步骤1 进入组件管理页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理”。

📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

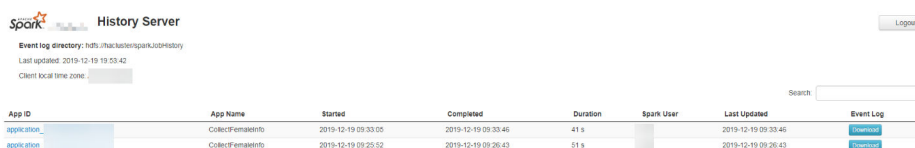
- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务”。

步骤2 选择“Spark”并在“Spark 概述”中“Spark Web UI”中单击“Spark Web UI”对应的“JobHistory”进入Web界面。

图 25-12 ResourceManager Web UI



图 25-13 Spark History Server



----结束

25.5 Spark 对接 OpenTSDB

25.5.1 创建表关联 OpenTSDB

功能描述

MRS的Spark实现了访问OpenTSDB的Datasource，能够在Spark中创建关联表，查询和插入OpenTSDB数据。

使用CREATE TABLE命令创建表并关联OpenTSDB上已有的metric。

📖 说明

若OpenTSDB上不存在metric，查询对应的表会报错。

语法格式

```
CREATE TABLE [IF NOT EXISTS] OPENTSDB_TABLE_NAME USING OPENTSDB OPTIONS (  
'metric' = 'METRIC_NAME',  
'tags' = 'TAG1,TAG2'  
);
```

关键字

参数	描述
metric	所创建的表对应的OpenTSDB中的指标名称。
tags	metric对应的标签，用于归类、过滤、快速检索等操作。可以是1个到8个，以“,”分隔，包括对应metric下所有tagk的值。

注意事项

创建表时，不需要指定timestamp和value字段，系统会根据指定的tags自动构建字段，包含以下字段，其中TAG1和TAG2由tags指定。

- TAG1 String
- TAG2 String
- timestamp Timestamp
- value double

示例

创建opentsdb_table表并关联到OpenTSDB组件的city.temp这个metric。

```
CREATE table opentsdb_table using opentsdb OPTIONS ('metric'='city.temp', 'tags'='city,location');
```

25.5.2 插入数据至 OpenTSDB 表

功能描述

使用INSERT INTO命令将表中的数据插入到已关联的OpenTSDB metric中。

语法格式

```
INSERT INTO TABLE_NAME SELECT * FROM SRC_TABLE;  
INSERT INTO TABLE_NAME VALUES(XXX);
```

关键字

参数	描述
TABLE_NAME	所关联的OpenTSDB表名。
SRC_TABLE	获取数据的表名，普通表即可。

注意事项

- 插入的数据不能为null；插入的数据相同，会覆盖原数据；插入的数据只有value值不同，也会覆盖原数据。
- 不支持INSERT OVERWRITE语法。
- 不建议对同一张表并发插入数据，因为有一定概率发生并发冲突，导致插入失败。
- 时间戳格式只支持yyyy-MM-dd hh:mm:ss。

示例

在opentsdb_table表中插入数据。

```
insert into opentsdb_table values('city1','city2','2018-05-03 00:00:00',21);
```

25.5.3 查询 OpenTSDB 表

SELECT命令用于查询OpenTSDB表中的数据。

语法格式

```
SELECT * FROM table_name WHERE tagk=tagv LIMIT number;
```

关键字

参数	描述
LIMIT	对查询结果进行限制。
number	参数仅支持INT类型。

注意事项

- 所查询的表必须是已经存在的表，否则会出错。
- 查询的tagv必须是已经有的值，否则会出错。

示例

查询表opentsdb_table中的数据。

```
SELECT * FROM opentsdb_table LIMIT 100;
SELECT * FROM opentsdb_table WHERE city='city1';
```

25.5.4 默认配置修改

默认会连接Spark的Executor所在节点本地的TSD进程，在MRS中一般使用默认配置即可，无需修改。

表 25-1 OpenTSDB 数据源相关配置

配置名	描述	样例值
spark.sql.datasource.opentsdb.host	连接的TSD进程地址	空（默认值） xx.xx.xx.xx，多个地址间用英文逗号间隔。
spark.sql.datasource.opentsdb.port	TSD进程端口号	4242（默认值）
spark.sql.datasource.opentsdb.randomSeed	当 spark.sql.datasource.opentsdb.host配置多个地址时，是否使用随机种子。配置为否时，所有在相同节点的executor会连接相同的host，这样可以配合 spark.blacklist.enabled=true来实现Task容错。	false（默认）

示例

在spark-sql，spark-beeline执行set语句后，再执行其他SQL：

```
set spark.sql.datasource.opentsdb.host = 192.168.2.143,192.168.2.158;
SELECT * FROM opentsdb_table ;
```

26 使用 Spark2x（MRS 3.x 及之后版本）

26.1 Spark 用户权限管理

26.1.1 SparkSQL 权限介绍

SparkSQL 权限

类似于Hive，SparkSQL也是建立在Hadoop上的数据仓库框架，提供类似SQL的结构化数据。

MRS提供用户、用户组和角色，集群中的各类权限需要先授予角色，然后将用户或者用户组与角色绑定。用户只有绑定角色或者加入绑定角色的用户组，才能获得权限。

📖 说明

- 如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加Spark2x的Ranger访问权限策略](#)。
- Spark2x开启或关闭Ranger鉴权后，需要重启Spark2x服务，并重新下载客户端，或刷新客户端配置文件spark/conf/spark-defaults.conf：
开启Ranger鉴权：spark.ranger.plugin.authorization.enable=true
关闭Ranger鉴权：spark.ranger.plugin.authorization.enable=false

权限管理介绍

SparkSQL的权限管理是指SparkSQL中管理用户操作数据库的权限系统，以保证不同用户之间操作数据库的独立性和安全性。如果一个用户想操作另一个用户的表、数据库等，需要获取相应的权限才能进行操作，否则会被拒绝。

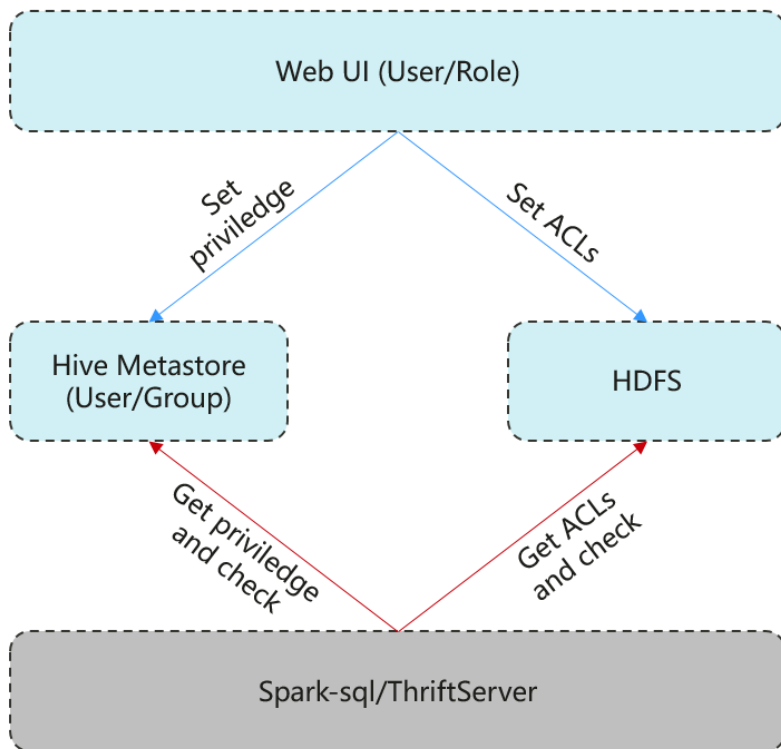
SparkSQL权限管理部分集成了Hive权限管理的功能。使用SparkSQL权限管理功能需要使用Hive的MetaStore服务和页面上的赋权功能。

图26-1展示了SparkSQL权限管理的基本架构。主要包含了两部分：页面赋权和服务获权并判断。

- 页面赋权：SparkSQL仅支持页面赋权的方式。在FusionInsight Manager的“系统 > 权限”中，可以进行用户、用户组和角色的添加/删除操作，可以对某个角色进行赋权/撤权。

- 服务获权并判断：当接收到客户端的DDL、DML的SQL命令时，SparkSQL服务会向MetaStore服务获取客户端用户对数据库信息的已有权限，并检查是否包含了所需的所有权限，如果是则继续执行，否则拒绝该用户的操作。当通过了MetaStore的权限检查后，还需进行HDFS的ACLs权限检查。

图 26-1 SparkSQL 权限管理架构图



SparkSQL还提供了列权限和视图权限，以满足用户不同场景的需求。

- 列权限介绍

SparkSQL权限控制由元数据权限控制和HDFS ACL权限控制两部分组成。Hive MetaStore会将表权限自动同步到HDFS ACL中时，不会同步列级别的权限。也就是说，当用户对表具有部分列权限或全部列权限时，不能通过HDFS Client访问HDFS文件。

- 在spark-sql模式下，用户仅具有列级别权限（即列权限用户）将不能访问HDFS文件，因此无法访问相应表的列。
- Beeline/JDBCServer模式下，用户间赋权，例如将A用户创建的表赋权给B用户时。
 - “hive.server2.enable.doAs” =true（在Spark服务端的“hive-site.xml”文件中配置）
此时用户B不可查询。需在HDFS上手动为文件赋读权限。
 - “hive.server2.enable.doAs” =false
 - 用户A和B均通过Beeline连接，用户B可查询。
 - A用户通过SQL方式建表，B用户可在Beeline进行查询。

而其他情况，如A用户使用Beeline建表，B用户通过SQL查询，或者A用户通过SQL方式建表，B用户使用SQL方式查询的情况均不支持。需在HDFS上手动为文件赋读权限。

📖 说明

由于“spark”用户在HDFS ACL的权限控制上为Spark管理员用户权限，Beeline客户端用户的权限控制仅取决于Spark侧的元数据权限。

- 视图权限介绍

视图权限是指仅对表的视图具有查询、修改等操作的权限，不再依赖于视图所在的表的相应权限。即用户拥有视图的查询权限时，不管是否有表权限都可以进行查询。视图的权限是针对整个表而言的，不支持对其中的部分列创建视图权限。

视图权限在SparkSQL权限上的限制与列权限相似，具体如下：

- 在spark-sql模式下，只有视图权限而没有表权限，且没有HDFS的读取权限时，用户不能访问HDFS上存储的表的数据，即该情况下不支持对该表的视图进行查询。
- Beeline/JDBCServer模式下，用户间赋权，例如将A用户创建的视图赋权给B用户时。

- “hive.server2.enable.doAs” =true（在Spark服务端的“hive-site.xml”文件中配置）

此时用户B不可查询。需在HDFS上手动为文件赋读权限。

- “hive.server2.enable.doAs” =false

- 用户A和B均通过Beeline连接，用户B可查询。

- A用户通过SQL方式创建视图，B用户可在Beeline进行查询。

而其他情况，如A用户使用Beeline创建视图，B用户通过SQL查询，或者A用户通过SQL方式创建视图，B用户使用SQL方式查询的情况均不支持。需在HDFS上手动为文件赋读权限。

对表的视图进行相应操作，分别需要具有以下权限。

- 创建视图时，需要数据库的CREATE权限、表的SELECT、SELECT_of_GRANT权限。
- 查询、描述视图时，只需要视图的SELECT权限，不需要视图所依赖的表或依赖的视图的SELECT权限。若同时查询视图和其他表，则仍然需要其他表的SELECT权限，例如：select * from v1 join t1时，需要有视图v1和表t1的SELECT权限，即使v1是基于t1的视图，也需要表t1的SELECT权限。

📖 说明

在Beeline/JDBCServer模式下，查询视图只需表的SELECT权限；而在spark-sql模式下，查询视图需要视图的SELECT权限和表的SELECT权限。

- 删除、修改视图时，必须要有视图的owner权限。

SparkSQL 权限模型

用户使用SparkSQL服务进行SQL操作，必须对SparkSQL数据库和表（含外表和视图）拥有相应的权限。完整的SparkSQL权限模型由元数据权限与HDFS文件权限组成。使用数据库或表时所需要的各种权限都是SparkSQL权限模型中的一种。

- 元数据权限

元数据权限即在元数据层上进行权限控制，与传统关系型数据库类似，SparkSQL数据库包含“创建”和“查询”权限，表和列包含“查询”、“插入”、

“UPDATE”和“删除”权限。SparkSQL中还包含拥有者权限“OWNERSHIP”和Spark管理员权限“管理”。

- 数据文件权限，即HDFS文件权限

SparkSQL的数据库、表对应的文件保存在HDFS中。默认创建的数据库或表保存在HDFS目录“/user/hive/warehouse”。系统自动以数据库名称和数据库中表的名称创建子目录。访问数据库或者表，需要在HDFS中拥有对应文件的权限，包含“读”、“写”和“执行”权限。

用户对SparkSQL数据库或表执行不同操作时，需要关联不同的元数据权限与HDFS文件权限。例如，对SparkSQL数据表执行查询操作，需要关联元数据权限“查询”，以及HDFS文件权限“读”和“执行”。

使用Manager界面图形化的角色管理功能来管理SparkSQL数据库和表的权限，只需要设置元数据权限，系统会自动关联HDFS文件权限，减少界面操作，提高效率。

SparkSQL 使用场景及对应权限

用户通过SparkSQL服务创建数据库需要加入Hive组，不需要角色授权。用户在Hive和HDFS中对自己创建的数据库或表拥有完整权限，可直接创建表、查询数据、删除数据、插入数据、更新数据以及授权他人访问表与对应HDFS目录与文件。

如果用户访问别人创建的表或数据库，需要授予权限。所以根据SparkSQL使用场景的不同，用户需要的权限可能也不相同。

表 26-1 SparkSQL 使用场景

主要场景	用户需要的权限
使用SparkSQL表、列或数据库	使用其他用户创建的表、列或数据库，不同的场景需要不同的权限，例如： <ul style="list-style-type: none">• 创建表，需要“创建”。• 查询数据，需要“查询”。• 插入数据，需要“插入”。
关联使用其他组件	部分场景除了SparkSQL权限，还可能需要组件的权限，例如：使用Spark on HBase，在SparkSQL中查询HBase表数据，需要设置HBase权限。

在一些特殊SparkSQL使用场景下，需要单独设置其他权限。

表 26-2 SparkSQL 授权注意事项

场景	用户需要的权限
创建SparkSQL数据库、表、外表，或者为已经创建的表或外表添加分区，且Hive用户指定数据文件保存在“/user/hive/warehouse”以外的HDFS目录。	<ul style="list-style-type: none"> 需要此目录已经存在，客户端用户是目录的属主，且用户对目录拥有“读”、“写”和“执行”权限。同时用户对此目录上层的每一级目录都拥有“读”和“执行”权限。 在Spark2x中，在创建HBase的外表时，需要拥有Hive端database的“创建”权限。而在Spark 1.5中，在创建HBase的外表时，需要拥有Hive端database的“创建”权限，也需要拥有HBase端Namespace的“创建”权限。
用户使用load将指定目录下所有文件或者指定文件，导入数据到表中。	<ul style="list-style-type: none"> 数据源为Linux本地磁盘，指定目录时需要此目录已经存在，系统用户“omm”对此目录以及此目录上层的每一级目录拥有“r”和“x”的权限。指定文件时需要此文件已经存在，“omm”对此文件拥有“r”的权限，同时对此文件上层的每一级目录拥有“r”和“x”的权限。 数据源为HDFS，指定目录时需要此目录已经存在，SparkSQL用户是目录属主，且用户对此目录及其子目录拥有“读”、“写”和“执行”权限，并且其上层的每一级目录拥有“读”和“执行”权限。指定文件时需要此文件已经存在，SparkSQL用户是文件属主，且用户对文件拥有“读”、“写”和“执行”权限，同时对此文件上层的每一级目录拥有“读”和“执行”权限。
创建函数、删除函数或者修改任意数据库。	需要授予“管理”权限。
操作Hive中所有的数据库和表。	需加入到supergroup用户组，并且授予“管理”权限。
对部分datasource表赋予insert权限后，执行insert analyze操作前需要单独对hdfs上的表目录赋予写权限。	当前对spark datasource表赋予Insert权限时，若表格式为：text csv json parquet orc,则不会修改表目录的权限。因此，对以上几种类型的datasource表赋予Insert权限后，还需要单独对hdfs上的表目录赋予写权限，用户才能成功对表执行insert analyze操作。

26.1.2 创建 SparkSQL 角色

操作场景

该任务指导MRS集群管理员在Manager创建并设置SparkSQL的角色。SparkSQL角色可设置Spark管理员权限以及数据表的数据操作权限。

用户使用Hive并创建数据库需要加入hive组，不需要角色授权。用户在Hive和HDFS中对自己创建的数据库或表拥有完整权限，可直接创建表、查询数据、删除数据、插入数据、更新数据以及授权他人访问表与对应HDFS目录与文件。默认创建的数据库或表保存在HDFS目录“/user/hive/warehouse”。

📖 说明

- 如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理，具体操作可参考[添加Spark2x的Ranger访问权限策略](#)。
- Spark2x开启或关闭Ranger鉴权后，需要重启Spark2x服务，并重新下载客户端，或刷新客户端配置文件spark/conf/spark-defaults.conf：
开启Ranger鉴权：spark.ranger.plugin.authorization.enable=true
关闭Ranger鉴权：spark.ranger.plugin.authorization.enable=false

操作步骤

1. 登录Manager页面，选择“系统 > 权限 > 角色”。
2. 单击“添加角色”，然后“角色名称”和“描述”输入角色名字与描述。
3. 设置角色“配置资源权限”请参见[表26-3](#)。
 - “Hive管理员权限”：Hive管理员权限。
 - “Hive读写权限”：Hive数据表管理权限，可设置与管理已创建的表的数据操作权限。

📖 说明

- Hive角色管理支持授予Hive管理员权限、访问表和视图的权限，不支持数据库的授权。
- Hive管理员权限不支持管理HDFS的权限。
- 如果数据库中的表或者表中的文件数量比较多，在授权时可能需要等待一段时间。例如表的文件数量为1万时，可能需要等待2分钟。

表 26-3 设置角色

任务场景	角色授权操作
设置Hive管理员权限	<p>在“配置资源权限”的表格中选择“待操作集群的名称 > Hive”，勾选“Hive管理权限”。</p> <p>用户绑定Hive管理员角色后，在每个维护操作会话中，还需要执行以下操作：</p> <ol style="list-style-type: none"> 1. 以客户端安装用户，登录安装Spark2x客户端的节点。 2. 执行以下命令配置环境变量。 例如，Spark2x客户端安装目录为“/opt/client”，执行source /opt/client/bigdata_env source /opt/client/Spark2x/component_env 3. 执行以下命令认证用户。 kinit Hive业务用户 4. 执行以下命令登录客户端工具。 /opt/client/Spark2x/spark/bin/beeline -u "jdbc:hive2://<zkNode1_IP>:<zkNode1_Port>,<zkNode2_IP>:<zkNode2_Port>,<zkNode3_IP>:<zkNode3_Port>/;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=sparkthriftserver2x;user.principal=spark2x/hadoop.<系统域名>@<系统域名>;sasLQop=auth-conf;auth=KERBEROS;principal=spark2x/hadoop.<系统域名>@<系统域名>"; 说明 <ul style="list-style-type: none"> • 其中 “<zkNode1_IP>:<zkNode1_Port>,<zkNode2_IP>:<zkNode2_Port>,<zkNode3_IP>:<zkNode3_Port>”是Zookeeper的URL。例如 “192.168.81.37:2181,192.168.195.232:2181,192.168.169.84:2181”。 • 其中“sparkthriftserver”是Zookeeper上的目录，表示客户端从该目录下随机选择Triftserver实例或proxyThriftServer进行连接。 • 用户可登录Manager，选择“系统 > 权限 > 域和互信”，查看“本端域”参数，即为当前系统域名。 “spark2x/hadoop.<系统域名>”为用户名，用户的用户名所包含的系统域名所有字母为小写。例如“本端域”参数为“9427068F-6EFA-4833-B43E-60CB641E5B6C.COM”，用户名为“spark2x/hadoo.9427068f-6efa-4833-b43e-60cb641e5b6c.com”。 5. 执行以下命令更新用户的管理员权限。 set role admin;

任务场景	角色授权操作
设置在默认数据库中，查询其他用户表的权限	<ol style="list-style-type: none"> 1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive读写权限”。 2. 在数据库列表中单击指定的数据库名称，显示数据库中的表。 3. 在指定表的“权限”列，勾选“查询”。
设置在默认数据库中，导入数据到其他用户表的权限	<ol style="list-style-type: none"> 1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive读写权限”。 2. 在数据库列表中单击指定的数据库名称，显示数据库中的表。 3. 在指定表的“权限”列，勾选“删除”和“插入”。

4. 单击“确定”完成。

26.1.3 配置 Spark 表、列和数据库的用户权限

操作场景

使用SparkSQL操作表或者数据库时，如果用户访问别人创建的表或数据库，需要授予对应的权限。为了实现更严格权限控制，SparkSQL也支持列级别的权限控制。如果要访问别人创建的表上某些列，需要授予列权限。以下介绍使用Manager角色管理功能在表授权、列授权和数据库授权三个场景下的操作。

操作步骤

SparkSQL表授权、列授权、数据库授权与Hive的操作相同，详情请参见[Hive用户权限管理](#)。

说明

- 在权限管理中，为了方便用户使用，授予数据库下表的任意权限将自动关联该数据库目录的HDFS权限。为了避免产生性能问题，取消表的任意权限，系统不会自动取消数据库目录的HDFS权限，但对应的用户只能登录数据库和查看表名。
- 若为角色添加或删除数据库的查询权限，数据库中的表也将自动添加或删除查询权限。此机制为Hive实现，SparkSQL与Hive保持一致。
- Spark不支持struct数据类型中列名称含有特殊字符（除字母、数字、下划线外的其他字符）。如果struct类型中列名称含有特殊字符，在FusionInsight Manager的“编辑角色”页面进行授权时，该列将无法正确显示。

相关概念

SparkSQL的语句在SparkSQL中进行处理，权限要求如[表26-4](#)所示。

表 26-4 使用 SparkSQL 表、列或数据库场景权限一览

操作场景	用户需要的权限
CREATE TABLE	“创建”，RWX+ownership（for create external table - the location） 说明 按照指定文件路径创建datasource表时，需要path后面文件的RWX+ownership权限。
DROP TABLE	“Ownership”（of table）
DROP TABLE PROPERTIES	“Ownership”
DESCRIBE TABLE	“查询”
SHOW PARTITIONS	“查询”
ALTER TABLE LOCATION	“Ownership”，RWX+ownership（for new location）
ALTER PARTITION LOCATION	“Ownership”，RWX+ownership（for new partition location）
ALTER TABLE ADD PARTITION	“插入”，RWX+ownership（for partition location）
ALTER TABLE DROP PARTITION	“删除”
ALTER TABLE(all of them except the ones above)	“Update”，“Ownership”
TRUNCATE TABLE	“Ownership”
CREATE VIEW	“查询”，“Grant Of Select”，“创建”
ALTER VIEW PROPERTIES	“Ownership”
ALTER VIEW RENAME	“Ownership”
ALTER VIEW ADD PARTS	“Ownership”
ALTER VIEW AS	“Ownership”
ALTER VIEW DROPPARTS	“Ownership”
ANALYZE TABLE	“查询”，“插入”
SHOW COLUMNS	“查询”
SHOW TABLE PROPERTIES	“查询”
CREATE TABLE AS SELECT	“查询”，“创建”
SELECT	“查询” 说明 与表一样，对视图进行SELECT操作的时候需要有该视图的“查询”权限。

操作场景	用户需要的权限
INSERT	“插入”，“删除 (for overwrite)”
LOAD	“插入”，“删除”，RWX+ownership(input location)
SHOW CREATE TABLE	“查询”，“Grant Of Select”
CREATE FUNCTION	“管理”
DROP FUNCTION	“管理”
DESC FUNCTION	-
SHOW FUNCTIONS	-
MSCK (metastore check)	“Ownership”
ALTER DATABASE	“管理”
CREATE DATABASE	-
SHOW DATABASES	-
EXPLAIN	“查询”
DROP DATABASE	“Ownership”
DESC DATABASE	-
CACHE TABLE	“查询”
UNCACHE TABLE	“查询”
CLEAR CACHE TABLE	“管理”
REFRESH TABLE	“查询”
ADD FILE	“管理”
ADD JAR	“管理”
HEALTHCHECK	-

26.1.4 配置 SparkSQL 业务用户权限

操作场景

SparkSQL业务还可能需关联使用其他组件，例如spark on HBase需要HBase权限。以下介绍SparkSQL关联HBase服务的操作。

前提条件

- 完成Spark客户端的安装，例如安装目录为“/opt/client”。
- 获取一个拥有MRS集群管理员权限的用户，例如“admin”。

操作步骤

• Spark on HBase授权

用户如果需要使用类似SQL语句的方式来操作HBase表，授予权限后可以使用SparkSQL访问HBase表。以授予用户在SparkSQL中查询HBase表的权限为例，操作步骤如下：

📖 说明

设置“spark.yarn.security.credentials.hbase.enabled”为“true”。

- a. 在Manager角色界面创建一个角色，例如“hive_hbase_create”，并授予创建HBase表的权限。

在“配置资源权限”的表格中选择“待操作集群的名称 > HBase > HBase Scope > global”，勾选命名空间“default”的“创建”，单击“确定”保存。

📖 说明

本例中建表是保存在Hive的“default”数据库中，默认具有“default”数据库的“建表”权限。如果Hive的数据库不是“default”，则还需要执行以下步骤：

在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive读写权限”，勾选所需指定的数据库的“建表”，单击“确定”保存。

- b. 在Manager角色界面创建一个角色，例如“hive_hbase_submit”，并授予提交任务到Yarn的队列的权限。

在“配置资源权限”的表格中选择“待操作集群的名称 > Yarn > 调度队列 > root”，勾选队列“default”的“提交”，单击“确定”保存。

- c. 在Manager用户界面创建一个“人机”用户，例如“hbase_creates_user”，加入“hive”组，绑定角色“hive_hbase_create”和“hive_hbase_submit”，用于创建SparkSQL表和HBase表。

- d. 以客户端安装用户登录安装客户端的节点。

- e. 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
source /opt/client/Spark2x/component_env
```

- f. 执行以下命令，认证用户。

```
kinit hbase_creates_user
```

- g. 执行以下命令，进入Spark JDBCServer客户端shell环境：

```
/opt/client/Spark2x/spark/bin/beeline -u "jdbc:hive2://
<zkNode1_IP>:<zkNode1_Port>,<zkNode2_IP>:<zkNode2_Port>,<zkNode3_IP>:<zkNode3_Port>";serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=sparkthriftserver2x;user.principal=spark2x/hadoop.<系统域名>@<系统域名>;saslQop=auth-conf;auth=KERBEROS;principal=spark2x/hadoop.<系统域名>@<系统域名>";
```

- h. 执行以下命令，同时在SparkSQL和HBase中创建表。例如创建表hbaseTable。

```
create table hbaseTable (id string, name string, age int) using
org.apache.spark.sql.hbase.HBaseSource options (hbaseTableName
"table1", keyCols "id", colsMapping = "", name=cf1.cq1, age=cf1.cq2);
```

- 创建好的SparkSQL表和HBase表分别保存在Hive的数据库“default”和HBase的命名空间“default”。
- i. 在Manager角色界面创建一个角色，例如“hive_hbase_select”，并授予查询SparkSQL on HBase表hbaseTable和HBase表hbaseTable的权限。
 - 在“配置资源权限”的表格中选择“待操作集群的名称 > HBase > HBase Scope > global > default”，勾选表hbaseTable的“读”，单击“确定”保存，授予HBase角色查询表的权限。
 - 编辑角色，在“配置资源权限”的表格中选择“待操作集群的名称 > HBase > HBase Scope > global > hbase”，勾选表“hbase:meta”的“执行”，单击“确定”保存。
 - 编辑角色，在“配置资源权限”的表格中选择“待操作集群的名称 > Hive > Hive读写权限 > default”，勾选表hbaseTable的“查询”，单击“确定”保存。
 - j. 在Manager用户界面创建一个“人机”用户，例如“hbase_select_user”，加入“hive”组，绑定角色“hive_hbase_select”，用于查询SparkSQL表和HBase表。
 - k. 执行以下命令，配置环境变量。

```
source /opt/client/bigdata_env
source /opt/client/Spark2x/component_env
```
 - l. 执行以下命令，认证用户。

```
kinit hbase_select_user
```
 - m. 执行以下命令，进入Spark JDBCServer客户端shell环境：

```
/opt/client/Spark2x/spark/bin/beeline -u "jdbc:hive2://
<zkNode1_IP>:<zkNode1_Port>,<zkNode2_IP>:<zkNode2_Port>,<zkNode3_IP>:<zkNode3_Port>";serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=sparkthriftserver2x;user.principal=spark2x/hadoop.<系统域名>@<系统域名>;sasLQop=auth-conf;auth=KERBEROS;principal=spark2x/hadoop.<系统域名>@<系统域名>;"
```
 - n. 执行以下命令，使用SparkSQL语句查询HBase表的数据。

```
select * from hbaseTable;
```

26.1.5 配置 Spark2x Web UI ACL

配置场景

当Spark2x Web UI中有一些不允许其他用户看到的数据时，用户可能想对UI进行安全防护。用户一旦登录，Spark2x可以比较与这个用户相对应的视图ACLs来确认是否授权用户访问 UI。

Spark2x存在两种类型的Web UI，一种为运行中任务的Web UI，可以通过Yarn原生页面的应用链接或者REST接口访问。一种为已结束任务的Web UI，可以通过Spark2x JobHistory服务或者REST接口访问。

说明

本章节仅支持安全模式（开启了Kerberos认证）集群。

- 运行中任务Web UI ACL配置。
运行中的任务，可通过服务端对如下参数进行配置。
 - “spark.admin.acls”：指定Web UI的管理员列表。
 - “spark.admin.acls.groups”：指定管理员组列表。
 - “spark.ui.view.acls”：指定yarn界面的访问者列表。
 - “spark.modify.acls.groups”：指定yarn界面的访问者组列表。
 - “spark.modify.acls”：指定Web UI的修改者列表。
 - “spark.ui.view.acls.groups”：指定Web UI的修改者组列表。
- 运行结束后Web UI ACL配置。
运行结束的任务通过客户端的参数“spark.history.ui.acls.enable”控制是否开启ACL访问权限。
如果开启了ACL控制，由客户端的“spark.admin.acls”和“spark.admin.acls.groups”配置指定Web UI的管理员列表和管理员组列表，由客户端的“spark.ui.view.acls”和“spark.modify.acls.groups”配置指定查看Web UI任务明细的访问者列表和组列表，由客户端的“spark.modify.acls”和“spark.ui.view.acls.groups”配置指定修改Web UI任务明细的访问者列表和组列表。

配置描述

登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索acl，在对应的JobHistory，JDBCServer，SparkResource和Spark界面修改以下参数。

表 26-5 参数说明

参数	说明	默认值
spark.history.ui.acls.enable	配置JobHistory是否支持单一任务的权限校验。	true
spark.acls.enable	配置是否开启spark权限管理。 如果开启，将会检查用户是否有权限访问和修改任务信息。	true
spark.admin.acls	配置spark管理员列表，列表中成员有权限管理所有spark任务，此处可以配置多个管理员用户，使用“，”分隔。	admin
spark.admin.acls.groups	配置spark管理组列表，列表中的组有权限管理所有spark任务，此处可以配置多个管理组，使用“，”分隔。	-
spark.modify.acls	配置有权限修改spark任务的成员列表。 启动任务的用户默认有此权限，此处可以配置多个用户，使用“，”分隔。	-
spark.modify.acls.groups	配置有权限修改spark任务的组列表，此处可以配置多个组，使用“，”分隔。	-

参数	说明	默认值
spark.ui.view.acls	配置有权限访问spark任务的成员列表。启动任务的用户默认有此权限，此处可以配置多个用户，使用“，”分隔。	-
spark.ui.view.acls.groups	配置有权限访问spark任务的组列表，此处可以配置多个组，使用“，”分隔。	-

说明

若使用客户端提交任务，“spark.admin.acls”、“spark.admin.acls.groups”、“spark.modify.acls”、“spark.modify.acls.groups”、“spark.ui.view.acls”和“spark.ui.view.acls.groups”参数修改后需要重新下载客户端。

26.1.6 Spark 客户端和服务端权限参数配置说明

SparkSQL权限管理功能相关的配置如下所示，客户端与服务端的配置相同。要使用表权限功能，需要在服务端和客户端添加如下配置。

- “spark-defaults.conf”配置文件

表 26-6 参数说明（1）

参数	描述	默认值
spark.sql.authorization.enabled	是否开启datasource语句的权限认证功能。建议将此参数修改为true，开启权限认证功能。	true

- “hive-site.xml”配置文件

表 26-7 参数说明（2）

参数	描述	默认值
hive.metastore.uris	Hive组件中MetaStore服务的地址，如“thrift://10.10.169.84:21088,thrift://10.10.81.37:21088”	-
hive.metastore.sasl.enabled	MetaStore服务是否使用SASL安全加固。表权限功能需要设置为“true”。	true
hive.metastore.kerberos.principal	Hive组件中MetaStore服务的Principal，如“hive/hadoop.<系统域名>@<系统域名>”。	hive-metastore/_HOST@EXAMPLE.COM

参数	描述	默认值
hive.metastore.thrift.sasl.qop	开启SparkSQL权限管理功能后，需将此参数设置为“auth-conf”。	auth-conf
hive.metastore.token.signature	MetaStore服务对应的token标识，设为“HiveServer2ImpersonationToken”。	HiveServer2ImpersonationToken
hive.security.authenticator.manager	Hive客户端授权的管理器，需设为“org.apache.hadoop.hive.ql.security.SessionStateUserGroupAuthenticator”。	org.apache.hadoop.hive.ql.security.SessionStateUserGroupAuthenticator
hive.security.authorization.enabled	是否开启客户端的授权，需设为“true”。	true
hive.security.authorization.createtable.owner.grants	将哪些权限赋给创建表的owner，建议设置为“ALL”。	ALL

- MetaStore服务的core-site.xml配置文件

表 26-8 参数说明 (3)

参数	描述	默认值
hadoop.proxyuser.spark.hosts	允许Spark用户伪装成来自哪些host的用户，需设为“*”，代表所有节点。	-
hadoop.proxyuser.spark.groups	允许Spark用户伪装成哪些用户组的用户，需设为“*”，代表所有用户组。	-

26.2 Spark 客户端使用实践

本章节提供从零开始使用Spark2x提交spark应用程序，包括Spark Core及Spark SQL。其中，Spark Core为Spark的内核模块，主要负责任务的执行，用于编写spark应用程序；Spark SQL为执行SQL的模块。

场景说明

假定用户有某个周末网民网购停留时间的日志文本，基于某些业务要求，要求开发Spark应用程序实现如下要求：

- 统计日志文件中本周末网购停留总时间超过2个小时的女性网民信息。
- 周末两天的日志文件第一列为姓名，第二列为性别，第三列为本次停留时间，单位为分钟，分隔符为“，”。

log1.txt：周六网民停留日志

```
LiuYang,female,20  
YuanJing,male,10  
GuoYijun,male,5  
CaiXuyu,female,50  
Liyuan,male,20  
FangBo,female,50  
LiuYang,female,20  
YuanJing,male,10  
GuoYijun,male,50  
CaiXuyu,female,50  
FangBo,female,60
```

log2.txt：周日网民停留日志

```
LiuYang,female,20  
YuanJing,male,10  
CaiXuyu,female,50  
FangBo,female,50  
GuoYijun,male,5  
CaiXuyu,female,50  
Liyuan,male,20  
CaiXuyu,female,50  
FangBo,female,50  
LiuYang,female,20  
YuanJing,male,10  
FangBo,female,50  
GuoYijun,male,50  
CaiXuyu,female,50  
FangBo,female,60
```

前提条件

- 在Manager界面创建用户并开通其HDFS、YARN、Kafka和Hive权限。
- 根据所用的开发语言安装并配置IntelliJ IDEA及JDK等工具。
- 已完成Spark2x客户端的安装及客户端网络连接的配置。
- 对于Spark SQL程序，需要先在客户端启动Spark SQL或Beeline以输入SQL语句。

操作步骤

步骤1 获取样例工程并将其导入IDEA，导入样例工程依赖jar包。通过IDEA配置并生成jar包。

步骤2 准备样例工程所需数据。

将场景说明中的原日志文件放置在HDFS系统中。

1. 本地新建两个文本文件，分别将log1.txt及log2.txt中的内容复制保存到input_data1.txt和input_data2.txt。
2. 在HDFS上建立一个文件夹“/tmp/input”，并上传input_data1.txt、input_data2.txt到此目录。

步骤3 将生成的jar包上传至Spark2x运行环境下（Spark2x客户端），如“/opt/female”。

步骤4 进入客户端目录，执行以下命令加载环境变量并登录。若安装了Spark2x多实例或者同时安装了Spark和Spark2x，在使用客户端连接具体实例时，请执行以下命令加载具体实例的环境变量。

```
source bigdata_env
```

```
source Spark2x/component_env
```

```
kinit <用于认证的业务用户>
```

步骤5 在bin目录下调用以下脚本提交Spark应用程序。

```
spark-submit --class com.huawei.bigdata.spark.examples.FemaleInfoCollection--  
master yarn-client/opt/female/FemaleInfoCollection.jar <inputPath>
```

📖 说明

- FemaleInfoCollection.jar为**步骤1**生成的jar包。
- <inputPath>是**步骤2.2**创建的目录。

步骤6（可选）在bin目录下调用**spark-sql**或**spark-beeline**脚本后便可直接输入SQL语句执行查询等操作。

如创建一个表，插入一条数据再对表进行查询。

```
spark-sql> CREATE TABLE TEST(NAME STRING, AGE INT);  
Time taken: 0.348 seconds  
spark-sql> INSERT INTO TEST VALUES('Jack', 20);  
Time taken: 1.13 seconds  
spark-sql> SELECT * FROM TEST;  
Jack    20  
Time taken: 0.18 seconds, Fetched 1 row(s)
```

步骤7 查看Spark应用运行结果。

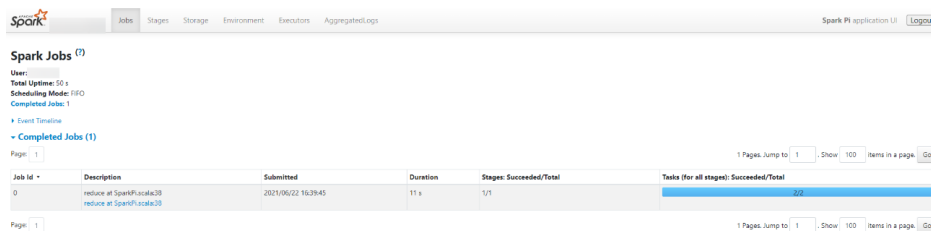
- 通过指定文件查看运行结果数据。
结果数据的存储路径和格式由Spark应用程序指定。
- 通过Web页面查看运行情况。
 - 登录Manager主页面。在服务中选择Spark2x。
 - 进入Spark2x概览页面，单击SparkWebUI任意一个实例，如JobHistory2x(host2)，登录History Server页面。
History Server页面用于展示已完成和未完成的应用的运行情况。

图 26-2 History Server 页面

Version	App ID	App Name	Started	Completed	Duration	Spark User	Last Updated	Event Log
	application_...	Spark Pi	2021-06-22 16:39:06	2021-06-22 16:39:56	50 s		2021-06-22 16:39:56	Download
	application_...	Spark Pi	2021-06-22 16:39:11	2021-06-22 16:39:55	45 s		2021-06-22 16:39:55	Download
	application_...	Spark Pi	2021-06-22 16:39:10	2021-06-22 16:39:55	44 s		2021-06-22 16:39:55	Download
	application_...	Spark Pi	2021-06-22 16:39:10	2021-06-22 16:39:46	35 s		2021-06-22 16:39:46	Download
	application_...	Spark Pi	2021-06-22 16:39:06	2021-06-22 16:39:44	38 s		2021-06-22 16:39:44	Download
	application_...	Spark Pi	2021-06-22 16:39:05	2021-06-22 16:39:26	21 s		2021-06-22 16:39:26	Download
	application_...	Spark Pi	2021-06-22 16:38:13	2021-06-22 16:39:05	52 s		2021-06-22 16:39:05	Download
	application_...	Spark Pi	2021-06-22 16:38:13	2021-06-22 16:38:57	45 s		2021-06-22 16:38:58	Download
	application_...	Spark Pi	2021-06-22 16:38:12	2021-06-22 16:38:57	45 s		2021-06-22 16:38:57	Download
	application_...	Spark Pi	2021-06-22 16:38:12	2021-06-22 16:38:54	42 s		2021-06-22 16:38:54	Download
	application_...	Spark Pi	2021-06-22 16:38:09	2021-06-22 16:38:47	38 s		2021-06-22 16:38:47	Download
	application_...	Spark Pi	2021-06-22 16:38:05	2021-06-22 16:38:46	41 s		2021-06-22 16:38:46	Download
	application_...	Spark Pi	2021-06-22 16:38:06	2021-06-22 16:38:27	21 s		2021-06-22 16:38:27	Download
	application_...	Spark Pi	2021-06-22 16:36:55	2021-06-22 16:38:06	1.2 min		2021-06-22 16:38:06	Download

- 选择一个应用ID，单击此页面将跳转到该应用的Spark UI页面。
Spark UI页面，用于展示正在执行的应用的运行情况。

图 26-3 Spark UI 页面



- 通过查看Spark日志获取应用运行情况。
可通过查看Spark2x日志了解应用运行情况，并根据日志信息调整应用程序，Spark2x相关日志请参见[Spark2x日志介绍](#)。

----结束

26.3 配置 Spark 读取 HBase 表数据

Spark On HBase

Spark on HBase为用户提供了在Spark SQL中查询HBase表，通过Beeline工具为HBase表进行存数据等操作。通过HBase接口可实现创建表、读取表、往表中插入数据等操作。

- 步骤1** 登录Manager界面，选择“集群 > 待操作集群的名称 > 集群属性”查看集群是否为安全模式。
- 是，执行**步骤2**。
 - 否，执行**步骤5**。
- 步骤2** 选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置 > 全部配置 > JDBCServer2x > 默认”，修改以下参数：

表 26-9 参数列表 1

参数	默认值	修改结果
spark.yarn.security.credentials.hbase.enabled	false	true

📖 说明

为了保证Spark2x可以长期访问HBase，建议不要修改HBase与HDFS服务的以下参数：

- dfs.namenode.delegation.token.renew-interval
- dfs.namenode.delegation.token.max-lifetime
- hbase.auth.key.update.interval
- hbase.auth.token.max.lifetime（不可修改，固定值为604800000毫秒，即7天）

如果必须要修改以上参数，请务必保证HDFS参数“dfs.namenode.delegation.token.renew-interval”的值不大于HBase参数“hbase.auth.key.update.interval”、“hbase.auth.token.max.lifetime”的值和HDFS参数“dfs.namenode.delegation.token.max-lifetime”的值。

步骤3 选择“SparkResource2x > 默认”，修改以下参数：

表 26-10 参数列表 2

参数	默认值	修改结果
spark.yarn.security.credentials.hbase.enabled	false	true

步骤4 重启Spark2x服务，配置生效。

说明

如果需要在Spark2x客户端用Spark on HBase功能，需要重新下载并安装Spark2x客户端。

步骤5 在Spark2x客户端使用spark-sql或者spark-beeline连接，可以查询由Hive on HBase所创建的表，支持通过SQL命令创建HBase表或创建外表关联HBase表。建表前，确认HBase中已存在对应 HBase表，下面以HBase表table1为例说明。

1. 通过Beeline工具创建HBase表，命令如下：

```
create table hbaseTable
(
  id string,
  name string,
  age int
)
using org.apache.spark.sql.hbase.HBaseSource
options(
  hbaseTableName "table1",
  keyCols "id",
  colsMapping "
  name=cf1.cq1,
  age=cf1.cq2
");
```

说明

- hbaseTable：是创建的spark表的表名。
- id string,name string, age int：是spark表的字段名和字段类型。
- table1：HBase表名。
- id：HBase表的rowkey列名。
- name=cf1.cq1, age=cf1.cq2：spark表的列和HBase表的列的映射关系。spark的name列映射HBase表的cf1列簇的cq1列，spark的age列映射HBase表的cf1列簇的cq2列。

2. 通过csv文件导入数据到HBase表，命令如下：

```
hbase org.apache.hadoop.hbase.mapreduce.ImportTsv -
Dimporttsv.separator=";" -
Dimporttsv.columns=HBASE_ROW_KEY,cf1:cq1,cf1:cq2,cf1:cq3,cf1:cq4,cf1:cq5
table1 /hperson
```

其中：table1为HBase表名，/hperson为csv文件存放的路径。

- 在spark-sql或spark-beeline中查询数据，*hbaseTable*为对应的spark表名。命令如下：

```
select * from hbaseTable;
```

----结束

Spark on HBaseV2

Spark on HBaseV2为用户提供了在Spark SQL中查询HBase表，通过Beeline工具为HBase表进行存数据等操作。通过HBase接口可实现创建表、读取表、往表中插入数据等操作。

- 步骤1** 登录Manager界面，选择“集群 > 待操作集群的名称 > 集群属性”查看集群是否为安全模式。
- 是，执行**步骤2**。
 - 否，执行**步骤5**。
- 步骤2** 选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置 > 全部配置 > JDBCServer2x > 默认”，修改以下参数：

表 26-11 参数列表 1

参数	默认值	修改结果
spark.yarn.security.credentials.hbase.enabled	false	true

说明

为了保证Spark2x可以长期访问HBase，建议不要修改HBase与HDFS服务的以下参数：

- dfs.namenode.delegation.token.renew-interval
- dfs.namenode.delegation.token.max-lifetime
- hbase.auth.key.update.interval
- hbase.auth.token.max.lifetime（不可修改，固定值为604800000毫秒，即7天）

如果必须要修改以上参数，请务必保证HDFS参数“dfs.namenode.delegation.token.renew-interval”的值不大于HBase参数“hbase.auth.key.update.interval”、“hbase.auth.token.max.lifetime”的值和HDFS参数“dfs.namenode.delegation.token.max-lifetime”的值。

- 步骤3** 选择“SparkResource2x > 默认”，修改以下参数：

表 26-12 参数列表 2

参数	默认值	修改结果
spark.yarn.security.credentials.hbase.enabled	false	true

- 步骤4** 重启Spark2x服务，配置生效。

📖 说明

如果需要在Spark2x客户端用Spark on HBase功能，需要重新下载并安装Spark2x客户端。

步骤5 在Spark2x客户端使用spark-sql或者spark-beeline连接，可以查询由Hive on HBase所创建的表，支持通过SQL命令创建HBase表或创建外表关联HBase表。具体见下面说明。下面以HBase表table1为例说明。

1. 通过spark-beeline工具创建表的语法命令如下：

```
create table hbaseTable1
(id string, name string, age int)
using org.apache.spark.sql.hbase.HBaseSourceV2
options(
hbaseTableName "table2",
keyCols "id",
colsMapping "name=cf1.cq1,age=cf1.cq2");
```

📖 说明

- hbaseTable1：是创建的spark表的表名。
 - id string,name string, age int：是spark表的字段名和字段类型。
 - table2：HBase表名。
 - id：HBase表的rowkey列名。
 - name=cf1.cq1, age=cf1.cq2：spark表的列和HBase表的列的映射关系。spark的name列映射HBase表的cf1列簇的cq1列，spark的age列映射HBase表的cf1列簇的cq2列。
2. 通过csv文件导入数据到HBase表，命令如下：
**hbase org.apache.hadoop.hbase.mapreduce.ImportTsv -
Dimporttsv.separator="," -
Dimporttsv.columns=HBASE_ROW_KEY,cf1:cq1,cf1:cq2,cf1:cq3,cf1:cq4,cf1:cq5
table2 /hperson**
其中：table2为HBase表名，/hperson为csv文件存放的路径。
 3. 在spark-sql或spark-beeline中查询数据，hbaseTable1为对应的spark表名，命令如下：
select * from hbaseTable1;

----结束

26.4 配置 Spark 任务不获取 HBase Token 信息

配置场景

使用Spark提交任务时，Driver默认会去HBase获取Token，访问HBase则需要配置文件“jaas.conf”进行安全认证。此时若用户未配置“jaas.conf”文件，会导致应用运行失败。

因此，根据应用是否涉及HBase进行以下处理：

- 当应用不涉及HBase时，即无需获取HBase Token。此时，将“spark.yarn.security.credentials.hbase.enabled”设置为“false”即可。

- 当应用涉及HBase时，将“spark.yarn.security.credentials.hbase.enabled”设置为“true”，且需要在Driver端配置“jaas.conf”文件，配置如下：
`{client}/spark/bin/spark-sql --master yarn-client --principal {principal} --keytab {keytab} --driver-java-options "-Djava.security.auth.login.config={LocalPath}/jaas.conf"`

在“jaas.conf”中指定Keytab和Prinical，示例如下：

```
Client {
  com.sun.security.auth.module.Krb5LoginModule required
  useKeyTab=true
  keyTab = "{LocalPath}/user.keytab"
  principal="super@<系统域名>"
  useTicketCache=false
  debug=false;
};
```

配置描述

在Spark客户端的“spark-defaults.conf”配置文件中设置。

表 26-13 参数说明

参数	说明	默认值
spark.yarn.security.credentials.hbase.enabled	HBase是否获取Token： <ul style="list-style-type: none"> true：获取 false：不获取 	false

26.5 Spark Core 企业级能力增强

26.5.1 配置 Spark HA 增强高可用

26.5.1.1 配置多主实例模式

配置场景

集群中支持同时共存多个ThriftServer服务，通过客户端可以随机连接其中的任意一个服务进行业务操作。即使集群中一个或多个ThriftServer服务停止工作，也不影响用户通过同一个客户端接口连接其他正常的ThriftServer服务。

配置描述

登录Manager，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索并修改以下参数。

表 26-14 多主实例参数说明

参数	说明	默认值
spark.thriftserver.zookeeper.connection.timeout	Zookeeper客户端连接超时时间，单位毫秒。	60000

参数	说明	默认值
spark.thriftserver.zookeeper.session.timeout	Zookeeper客户端会话超时时间，单位毫秒。	90000
spark.thriftserver.zookeeper.retry.times	Zookeeper客户端失联后，重试次数。	3
spark.yarn.queue	JDBCServer服务所在的Yarn队列。	default

26.5.1.2 配置 Spark 多租户模式

配置场景

多租户模式是将JDBCServer和租户绑定，每一个租户对应一个或多个JDBCServer，一个JDBCServer只给一个租户提供服务。不同的租户可以配置不同的Yarn队列，从而达到资源隔离。

说明

Yarn资源不足情况下，不建议开启多租户模式。

配置描述

登录Manager，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索并修改以下参数。

表 26-15 参数说明

参数	说明	默认值
spark.proxyserver.hash.enabled	是否使用Hash算法连接ProxyServer。 <ul style="list-style-type: none"> true为使用Hash算法，使用多租户模式时，该参数需配置为true。 false为使用随机连接，多主实例模式，配置为false。 	true 说明 该参数修改后需要重新下载客户端。
spark.thriftserver.proxy.enabled	是否使用多租户模式。 <ul style="list-style-type: none"> false表示使用多实例模式 true表示使用多租户模式 	true
spark.thriftserver.proxy.maxThriftServerPerTenancy	多租户模式下，一个租户可启动JDBCServer实例的最大个数。	1
spark.thriftserver.proxy.maxSessionPerThriftServer	多租户模式下，单个JDBCServer实例的session数量超过该值时，如果租户的JDBCServer最大实例数量没超过限制，则启动新的JDBCServer，否则输出警告日志。	50

参数	说明	默认值
spark.thriftserver.proxy.sessionWaitTime	多租户模式下，当JDBCServer的session连接数为0时，停止JDBCServer前的等待时间。	180000
spark.thriftserver.proxy.sessionThreshold	多租户模式下，当JDBCServer的session使用率（公式：当前session数 / (spark.thriftserver.proxy.maxSessionPerThriftServer * 当前JDBCServer个数)）达到阈值时，自动新增JDBCServer。	100
spark.thriftserver.proxy.healthcheck.period	多租户模式下，JDBCServer代理检查JDBCServer健康状态周期。	60000
spark.thriftserver.proxy.healthcheck.recheckTimes	多租户模式下，JDBCServer代理检查JDBCServer健康状态失败后重试次数。	3
spark.thriftserver.proxy.healthcheck.waitTime	多租户模式下，JDBCServer代理发送健康检查，等待JDBCServer响应的超时时间。	10000
spark.thriftserver.proxy.session.check.interval	多租户模式下，JDBCServer代理检查session的周期。	6h
spark.thriftserver.proxy.idle.session.timeout	多租户模式下，JDBCServer代理session的空闲超时时间。如果在这段时间内没有做任何操作，session会被关闭。	7d
spark.thriftserver.proxy.idle.session.check.operation	多租户模式下，JDBCServer代理session的过期是否要判断该session上还存在operation。	true
spark.thriftserver.proxy.idle.operation.timeout	多租户模式下，operation的超时时间。如果operation超时，operation会被关闭。	5d
hive.spark.client.server.connect.timeout	多租户模式下，客户端连接超时时间。	5min

26.5.1.3 配置多主实例与多租户模式切换

配置场景

在使用集群中，如果需要在多主实例模式与多租户模式之间切换，则还需要进行如下参数的设置。

- 多租户切换成多主实例模式
修改Spark2x服务的以下参数：
 - spark.thriftserver.proxy.enabled=false

- spark.scheduler.allocation.file=#{conf_dir}/fairscheduler.xml
- spark.proxyserver.hash.enabled=false
- 多主实例切换成多租户模式
修改Spark2x服务的以下参数：
 - spark.thriftserver.proxy.enabled=true
 - spark.scheduler.allocation.file=./__spark_conf__/__hadoop_conf__/fairscheduler.xml
 - spark.proxyserver.hash.enabled=true

配置描述

登录Manager，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索并修改以下参数。

表 26-16 参数说明

参数	说明	默认值
spark.thriftserver.proxy.enabled	是否使用多租户模式。 <ul style="list-style-type: none"> • false表示使用多实例模式 • true表示使用多租户模式 	true
spark.scheduler.allocation.file	公平调度文件路径。 <ul style="list-style-type: none"> • 多主实例配置为：#{conf_dir}/fairscheduler.xml • 多租户配置为：./__spark_conf__/__hadoop_conf__/fairscheduler.xml 	./__spark_conf__/__hadoop_conf__/fairscheduler.xml
spark.proxyserver.hash.enabled	是否使用Hash算法连接ProxyServer。 <ul style="list-style-type: none"> • true为使用Hash算法，使用多租户模式时，该参数需配置为true。 • false为使用随机连接，多主实例模式，配置为false。 	true 说明 该参数修改后需要重新下载客户端。

26.5.2 配置 Spark 事件队列大小

配置场景

Spark中见到的UI、EventLog、动态资源调度等功能都是通过事件传递实现的。事件有SparkListenerJobStart、SparkListenerJobEnd等，记录了每个重要的过程。

每个事件在发生后都会保存到一个队列中，Driver在创建SparkContext对象时，会启动一个线程循环的从该队列中依次拿出一个事件，然后发送给各个Listener，每个Listener感知到事件后就会做各自的处理。

因此当队列存放的速度大于获取的速度时，就会导致队列溢出，从而丢失了溢出的事件，影响了UI、EventLog、动态资源调度等功能。所以为了更灵活的使用，在这边添加一个配置项，用户可以根据Driver的内存大小设置合适的值。

配置描述

参数入口：

在执行应用之前，在Spark服务配置中修改。在Manager系统中，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”。在搜索框中输入参数名称。

表 26-17 参数说明

参数	描述	默认值
spark.scheduler.listenerbus.eventqueue.capacity	事件队列的大小，可以根据Driver的内存做适当的配置。	100000 0

说明

当Driver日志中出现如下的日志时，表示队列溢出了。

1. 普通应用：

Dropping SparkListenerEvent because no remaining room in event queue.
This likely means one of the SparkListeners is too slow and cannot keep up with the rate at which tasks are being started by the scheduler.

2. Spark Streaming应用：

Dropping StreamingListenerEvent because no remaining room in event queue.
This likely means one of the StreamingListeners is too slow and cannot keep up with the rate at which events are being started by the scheduler.

26.5.3 配置 parquet 表的压缩格式

配置场景

当前版本对于parquet表的压缩格式分以下两种情况进行配置：

- 对于分区表，需要通过parquet本身的配置项“parquet.compression”设置parquet表的数据压缩格式。如在建表语句中设置tblproperties：“parquet.compression”=“snappy”。
- 对于非分区表，需要通过“spark.sql.parquet.compression.codec”配置项来设置parquet类型的数据压缩格式。直接设置“parquet.compression”配置项是无效的，因为它会读取“spark.sql.parquet.compression.codec”配置项的值。当“spark.sql.parquet.compression.codec”未做设置时默认值为“snappy”，“parquet.compression”会读取该默认值。

因此，“spark.sql.parquet.compression.codec”配置项只适用于设置非分区表的parquet压缩格式。

配置参数

参数入口：

在Manager系统中，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，在搜索框中输入参数名称。

表 26-18 参数介绍

参数	描述	默认值
spark.sql.parquet.compression.codec	对于非分区parquet表，设置其存储文件的压缩格式。	snappy

26.5.4 使用 Ranger 时适配第三方 JDK

配置场景

当使用Ranger作为spark sql的权限管理服务时，访问RangerAdmin需要使用集群中的证书。若用户未使用集群中的JDK或者JRE，而是使用第三方JDK时，会出现访问RangerAdmin失败，进而spark应用程序启动失败的问题。

在这个场景下，需要进行以下操作，将集群中的证书导入第三方JDK或者JRE中。

配置方法

步骤1 导出集群中的证书：

1. 安装集群客户端，例如安装路径为“/opt/client”。
2. 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

3. 执行以下命令配置环境变量。

```
source bigdata_env
```

4. 生成证书文件

```
keytool -export -alias fusioninsightsubroot -storepass changeit -  
keystore /opt/client/JRE/jre/lib/security/cacerts -file  
fusioninsightsubroot.crt
```

步骤2 将集群中的证书导入第三方JDK或者JRE中

将**步骤1**中生成的fusioninsightsubroot.crt文件拷贝到第三方JRE节点上，设置好该节点的JAVA_HOME环境变量后，执行以下命令导入证书：

```
keytool -import -trustcacerts -alias fusioninsightsubroot -storepass changeit -  
file fusioninsightsubroot.crt -keystore MY_JRE/lib/security/cacerts
```

📖 说明

'MY_JRE'表示第三方JRE安装路径，请自行修改。

----结束

26.5.5 使用 Spark 小文件合并工具说明

工具介绍

在Hadoop大规模生产集群中，由于HDFS的元数据都保存在NameNode的内存中，集群规模受制于NameNode单点的内存限制。如果HDFS中有大量的小文件，会消耗NameNode大量内存，还会大幅降低读写性能，延长作业运行时间。因此，小文件问题是制约Hadoop集群规模扩展的关键问题。

本工具主要有如下两个功能：

1. 扫描表中有多少低于用户设定阈值的小文件，返回该表目录中所有数据文件的平均大小。
2. 对表文件提供合并功能，用户可设置合并后的平均文件大小。

支持的表类型

Spark：Parquet、ORC、CSV、Text、Json。

Hive：Parquet、ORC、CSV、Text、RCFile、Sequence、Bucket。

说明

1. 数据有压缩的表在执行合并后会采用Spark默认的压缩格式-Snappy。可以通过在客户端设置“spark.sql.parquet.compression.codec”（可选：uncompressed, gzip, snappy）和“spark.sql.orc.compression.codec”（可选：uncompressed, zlib, lzo, snappy）来选择Parquet和Orc表的压缩格式；由于Hive和Spark表在可选的压缩格式上有区别，除以上列出的压缩格式外，其他的压缩格式不支持。
2. 合并桶表数据，需要先在Spark2x客户端的hive-site.xml里加上配置：

```
<property>
<name>hive.enforce.bucketing</name>
<value>>false</value>
</property>
<property>
<name>hive.enforce.sorting</name>
<value>>false</value>
</property>
```
3. Spark暂不支持Hive的加密列特性。

工具使用

下载安装客户端，例如安装目录为“/opt/client”。进入“/opt/client/Spark2x/spark/bin”，执行mergetool.sh脚本。

加载环境变量

```
source /opt/client/bigdata_env
```

```
source /opt/client/Spark2x/component_env
```

扫描功能

命令形式：**sh mergetool.sh scan <db.table> <filesize>**

db.table的形式是“数据库名.表名”，filesize为用户自定义的小文件阈值（单位MB），返回结果为小于该阈值的文件个数，及整个表目录数据文件的平均大小。

例如：**sh mergetool.sh scan default.table1 128**

合并功能

命令形式：**sh mergetool.sh merge <db.table> <filesize> <shuffle>**

db.table的形式是“数据库名.表名”，filesize为用户自定义的合并后平均文件大小（单位MB），shuffle是一个boolean值，取值true/false，作用是设置合并过程中是否允许数据进行shuffle。

例如：**sh mergetool.sh merge default.table1 128 false**

提示如下，则操作成功：

```
SUCCESS: Merge succeeded
```

📖 说明

1. 请确保当前用户对合并的表具有owner权限。
2. 合并前请确保HDFS上有足够的存储空间，至少需要被合并表大小的一倍以上。
3. 合并表数据的操作需要单独进行，在此过程中读表，可能临时出现找不到文件的问题，合并完成后会恢复正常；另外在合并过程中请注意不要对相应的表进行写操作，否则可能会产生数据一致性问题。
4. 若合并完成后，在一直处于连接状态的spark-beeline/spark-sql session中查询分区表的数据，出现文件不存在的问题，根据提示可以执行“refresh table 表名”后再重新查询。
5. 请依据实际情况合理设置filesize值，例如可以在scan得到表中平均文件大小值average后，在merge时将filesize设置一个比average更大的值；否则，执行合并后可能出现文件数变得更多的情况。
6. 合并过程中，会将原表数据放入回收站，再填入已合并的数据。若在此过程中发生异常，根据工具提示，可将trash目录中的数据通过hdfs的mv命令恢复。
7. 在HDFS router联邦场景下，如果表的根路径与根路径“/user”的目标NameService不同，在二次合并时需要手动清理放入回收站的原表文件，否则会导致合并失败。
8. 此工具应用客户端配置，需要做性能调优可修改客户端配置文件的相关配置。

shuffle设置

对于合并功能，可粗略估计合并前后分区数的变化：

一般来说，旧分区数>新分区数，可设置shuffle为false；但如果旧分区远大于新分区数，例如高于100倍以上，可以考虑设置shuffle为true，增加并行度，提高合并的速度。

须知

- 设置shuffle为true（repartition），会有性能上的提升；但是由于Parquet和Orc存储方式的特殊性，repartition会使压缩率变小，直接表现是hdfs上表的总大小会增大到1.3倍。
- 设置shuffle为false（coalesce），合并后的大小不会非常平均，可能会分布在设置的filesize左右。

日志存放位置

默认日志存放位置为/tmp/SmallFilesLog.log4j，如需自定义日志存放位置，可在/opt/client/Spark2x/spark/tool/log4j.properties中配置log4j.appender.logfile.File。

26.5.6 配置流式读取 Spark Driver 执行结果

配置场景

在执行查询语句时，返回结果有可能会很大（10万数量以上），此时很容易导致 JDBCServer OOM（Out of Memory）。因此，提供数据汇聚功能特性，在基本不牺牲性能的情况下尽力避免OOM。

配置描述

提供两种不同的数据汇聚功能配置选项，两者在Spark JDBCServer服务端的tunning选项中进行设置，设置完后需要重启JDBCServer。

表 26-19 参数说明

参数	说明	默认值
spark.sql.bigdata.thriftServer.useHdfsCollect	<p>是否将结果数据保存到HDFS中而不是内存中。</p> <p>优点：由于查询结果保存在hdfs端，因此基本不会造成JDBCServer的OOM。</p> <p>缺点：速度慢。</p> <ul style="list-style-type: none"> • true：保存至HDFS中 • false：不使用该功能 <p>须知 spark.sql.bigdata.thriftServer.useHdfsCollect参数设置为true时，将结果数据保存到HDFS中，但JobHistory原生页面上Job的描述信息无法正常关联到对应的SQL语句，同时spark-beeline命令行中回显的Execution ID为null，为解决JDBCServer OOM问题，同时显示信息正确，建议选择 spark.sql.userlocalFileCollect参数进行配置。</p>	false
spark.sql.uselocalFileCollect	<p>是否将结果数据保存在本地磁盘中而不是内存里面。</p> <p>优点：结果数据小数据量情况下和原生内存的方式相比性能损失可以忽略，大数据情况下（亿级数据）性能远比使用hdfs，以及原生内存方式好。</p> <p>缺点：需要调优。大数据情况下建议JDBCServer driver端内存10G，executor端每个核心分配3G内存。</p> <ul style="list-style-type: none"> • true：使用该功能 • false：不使用该功能 	false

参数	说明	默认值
spark.sql.collect.Hive	<p>该参数在spark.sql.uselocalFileCollect开启的情况下生效。直接序列化的方式，还是间接序列化的方式保存结果数据到磁盘。</p> <p>优点：针对分区数特别多的表查询结果汇聚性能优于直接使用结果数据保证在磁盘的方式。</p> <p>缺点：和spark.sql.uselocalFileCollect开启时候的缺点一样。</p> <ul style="list-style-type: none"> • true：使用该功能 • false：不使用该功能 	false
spark.sql.collect.serialize	<p>该参数在spark.sql.uselocalFileCollect, spark.sql.collect.Hive同时开启的情况下生效。</p> <p>作用是进一步提升性能</p> <ul style="list-style-type: none"> • java：采用java序列化方式收集数据。 • kryo：采用kryo序列化方式收集数据，性能要比采用java好。 	java

说明

参数spark.sql.bigdata.thriftServer.useHdfsCollect和spark.sql.uselocalFileCollect不能同时设置为true。

26.6 Spark SQL 企业级能力增强

26.6.1 配置矢量化读取 ORC 数据

配置场景

ORC文件格式是一种Hadoop生态圈中的列式存储格式，它最初产生自Apache Hive，用于降低Hadoop数据存储空间和加速Hive查询速度。和Parquet文件格式类似，它并不是一个单纯的列式存储格式，仍然是首先根据行组分割整个表，在每一个行组内按列进行存储，并且文件中的数据尽可能的压缩来降低存储空间的消耗。矢量化读取ORC格式的数据能够大幅提升ORC数据读取性能。在Spark2.3版本中，SparkSQL支持矢量化读取ORC数据（这个特性在Hive的历史版本中已经得到支持）。矢量化读取ORC格式的数据能够获得比传统读取方式数倍的性能提升。

该特性可以通过下面的配置项开启：

- “spark.sql.orc.enableVectorizedReader”：指定是否支持矢量化方式读取ORC格式的数据，默认为true。
- “spark.sql.codegen.wholeStage”：指定是否需要将多个操作的所有stage编译为一个java方法，默认为true。
- “spark.sql.codegen.maxFields”：指定codegen的所有stage所支持的最大字段数（包括嵌套字段），默认为100。

- “spark.sql.orc.impl”：指定使用Hive还是Spark SQL native作为SQL执行引擎来读取ORC数据，默认为hive。

配置参数

登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索以下参数。

参数	说明	默认值	取值范围
spark.sql.orc.enableVectorizedReader	指定是否支持矢量化方式读取ORC格式的数据，默认为true。	true	[true,false]
spark.sql.codegen.wholeStage	指定是否需要将多个操作的所有stage编译为一个java方法，默认为true。	true	[true,false]
spark.sql.codegen.maxFields	指定codegen的所有stage所支持的最大字段数（包括嵌套字段），默认为100。	100	大于0
spark.sql.orc.impl	指定使用Hive还是Spark SQL native作为SQL执行引擎来读取ORC数据，默认为hive。	hive	[hive,native]

📖 说明

- 使用SparkSQL内置的矢量化方式读取ORC数据需要满足下面的条件：
 - spark.sql.orc.enableVectorizedReader：true，默认是true，一般不做修改。
 - spark.sql.codegen.wholeStage：true，默认为true，一般不做修改。
 - spark.sql.codegen.maxFields不小于scheme的列数。
 - 所有的数据类型均为AtomicType类型；所谓Atomic Type表示非NULL、UDTs、arrays、maps类型。如果列中存在这几种类型的任何一种，都无法获得预期的性能。
 - spark.sql.orc.impl：native，默认为hive。
- 若使用客户端提交任务，“spark.sql.orc.enableVectorizedReader”、“spark.sql.codegen.wholeStage”、“spark.sql.codegen.maxFields”、“spark.sql.orc.impl”、参数修改后需要重新下载客户端才能生效。

26.6.2 配置过滤掉分区表中路径不存在的分区

配置场景

当读取HIVE分区表时，如果指定的分区路径在HDFS上不存在，则执行select查询时会报FileNotFoundException异常。此时可以通过配置“spark.sql.hive.verifyPartitionPath”参数来过滤掉分区路径不存在的分区，来避免读取时报错。

配置描述

可以通过以下两种方式配置是否过滤掉分区表分区路径不存在的分区。

- 在Spark Driver端的“spark-defaults.conf”配置文件中设置。

表 26-20 参数说明

参数	说明	默认值
spark.sql.hive.verifyPartitionPath	配置读取HIVE分区表时，是否过滤掉分区表分区路径不存在的分区。 “true”：过滤掉分区路径不存在的分区； “false”：不进行过滤。	false

- 在spark-submit命令提交应用时，通过“--conf”参数配置是否过滤掉分区表分区路径不存在的分区。

示例：

```
spark-submit --class org.apache.spark.examples.SparkPi --conf spark.sql.hive.verifyPartitionPath=true $SPARK_HOME/lib/spark-examples_*.jar
```

26.6.3 配置 Hive 表分区动态覆盖

配置场景

在旧版本中，使用insert overwrite语法覆写分区表时，只支持对指定的分区表达式进行匹配，未指定表达式的分区将被全部删除。在spark2.3版本中，增加了对未指定表达式的分区动态匹配的支持，此种语法与Hive的动态分区匹配语法行为一致。

配置参数

登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索以下参数。

参数	说明	默认值	取值范围
spark.sql.sources.partitionOverwriteMode	当前执行insert overwrite 命令插入数据到分区表时，支持两种模式：STATIC模式和DYNAMIC模式。STATIC模式下，Spark会按照匹配条件删除所有分区。在DYNAMIC模式下，Spark按照匹配条件匹配分区，并动态匹配没有指定匹配条件的分区。	STATIC	[STATIC,DYNAMIC]

26.6.4 配置 Spark SQL 开启 Adaptive Execution 特性

配置场景

Spark SQL Adaptive Execution特性用于使Spark SQL在运行过程中，根据中间结果优化后续执行流程，提高整体执行效率。当前已实现的特性如下：

1. 自动设置shuffle partition数

在启用Adaptive Execution特性前，Spark SQL根据spark.sql.shuffle.partitions配置指定shuffle时的partition个数。此种方法在一个应用中执行多种SQL查询时缺乏灵活性，无法保证所有场景下的性能合适。开启Adaptive Execution后，Spark SQL将自动为每个shuffle过程动态设置partition个数，而不是使用通用配置，使每次shuffle过程自动使用最合理的partition数。

2. 动态调整执行计划

在启用Adaptive Execution特性前，Spark SQL根据RBO和CBO的优化结果创建执行计划，此种方法忽略了数据在运行过程中的结果集变化。比如基于某个大表创建的视图，与其他大表join时，即便视图的结果集很小，也无法将执行计划调整为BroadcastJoin。启用Adaptive Execution特性后，Spark SQL能够在运行过程中根据前面stage的运行结果动态调整后续的执行计划，从而获得更好的执行性能。

3. 自动处理数据倾斜

在执行SQL语句时，若存在数据倾斜，可能导致单个executor内存溢出、任务执行缓慢等问题。启动Adaptive Execution特性后，Spark SQL能自动处理数据倾斜场景，对倾斜的分区，启动多个task进行处理，每个task读取若干个shuffle输出文件，再对这部分任务的Join结果进行Union操作，以达到消除数据倾斜的效果

配置参数

登录FusionInsight Manager系统，选择“集群 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索以下参数。

参数	说明	默认值
spark.sql.adaptive.enabled	配置是否启用自适应执行功能。 注意：AQE特性与DPP（动态分区裁剪）特性同时开启时，SparkSQL任务执行中会优先执行DPP特性，从而使得AQE特性不生效。	false
spark.sql.optimizer.dynamicPartitionPruning.enabled	动态分区裁剪功能的开关。	true
spark.sql.adaptive.coalescePartitions.enabled	如果配置为true并且“spark.sql.adaptive.enabled”为true，Spark将根据目标大小（由spark.sql.adaptive.advisoryPartitionSizeInBytes指定）合并连续的随机播放分区，以避免执行过多的小任务。	true
spark.sql.adaptive.coalescePartitions.initialPartitionNum	合并之前的shuffle分区的初始数量，默认等于spark.sql.shuffle.partitions。只有当spark.sql.adaptive.enabled和spark.sql.adaptive.coalescePartitions.enabled都为true时，该配置才有效。创建时可选，初始分区数必须为正数。	200

参数	说明	默认值
spark.sql.adaptive.coalescePartitions.minPartitionNum	合并后的最小shuffle分区数。如果不设置，默认为Spark集群的默认并行度。只有当spark.sql.adaptive.enabled和spark.sql.adaptive.coalescePartitions.enabled都为true时，该配置才有效。创建时可选，最小分区数必须为正数。	1
spark.sql.adaptive.shuffle.targetPostShuffleInputSize	shuffle后单个分区的目标大小，从Spark3.0开始不再支持。	64MB
spark.sql.adaptive.advisoryPartitionSizeInBytes	自适应优化时（spark.sql.adaptive.enabled为true时）shuffle分区的咨询大小（单位：字节），在Spark聚合小shuffle分区或拆分倾斜的shuffle分区时生效。	64MB
spark.sql.adaptive.fetchShuffleBlocksInBatch	是否批量取连续的shuffle块。对于同一个map任务，批量读取连续的shuffle块可以减少IO，提高性能，而不是逐个读取块。注意，只有当spark.sql.adaptive.enabled和spark.sql.adaptive.coalescePartitions.enabled都为true时，单次读取请求中存在多个连续块。这个特性还依赖于一个可重定位的序列化器，使用的级联支持编解码器和新版本的shuffle提取协议。	true
spark.sql.adaptive.localShuffleReader.enabled	当“true”且spark.sql.adaptive.enabled为“true”时，Spark在不需要进行shuffle分区时，会尝试使用本地shuffle reader读取shuffle数据，例如：将sort-merge join转换为broadcast-hash join后。	true
spark.sql.adaptive.skewJoin.enabled	当此配置为true且spark.sql.adaptive.enabled设置为true时，启用运行时自动处理join运算中的数据倾斜功能	true
spark.sql.adaptive.skewJoin.skewedPartitionFactor	此配置为一个倍数因子，用于判定分区是否为数据倾斜分区。单个分区被判定为数据倾斜分区的条件为：当一个分区的数据大小超过除此分区外其他所有分区大小的中值与该配置的乘积，并且大小超过spark.sql.adaptive.skewJoin.skewedPartitionThresholdInBytes配置值时，此分区被判定为数据倾斜分区	5

参数	说明	默认值
spark.sql.adaptive.skewJoin.skewedPartitionThresholdInBytes	分区大小（单位：字节）大于该阈值且大于 spark.sql.adaptive.skewJoin.skewedPartitionFactor 与分区中值的乘积，则认为该分区存在倾斜。理想情况下，此配置应大于 spark.sql.adaptive.advisoryPartitionSizeInBytes。	256MB
spark.sql.adaptive.nonEmptyPartitionRatioForBroadcastJoin	两表进行 join 操作的时候，当非空分区比率低于此配置时，无论其大小如何，都不会被视为自适应执行中广播哈希连接的生成端。只有当 spark.sql.adaptive.enabled 为 true 时，此配置才有效。	0.2

26.6.5 配置 SparkSQL 的分块个数

配置场景

SparkSQL 在进行 shuffle 操作时默认的分块数为 200。在数据量特别大的场景下，使用默认的分块数就会造成单个数据块过大。如果一个任务产生的单个 shuffle 数据块大于 2G，该数据块在被 fetch 的时候还会报类似错误：

```
Adjusted frame length exceeds 2147483647: 2717729270 - discarded
```

例如，SparkSQL 运行 TPCDS 500G 的测试时，使用默认配置出现错误。所以当数据量较大时需要适当的调整该参数。

配置参数

参数入口：

在 Manager 系统中，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”。在搜索框中输入参数名称。

表 26-21 参数介绍

参数	描述	默认值
spark.sql.shuffle.partitions	SparkSQL 在进行 shuffle 操作时默认的分块数。	200

26.7 Spark Streaming 企业级能力增强

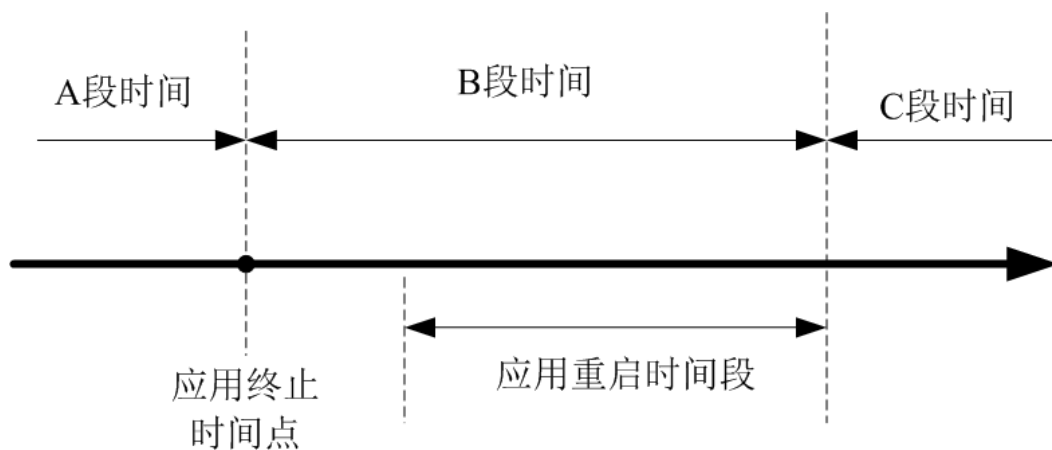
26.7.1 配置 Kafka 后进先出

配置场景

当Spark Streaming应用与Kafka对接，Spark Streaming应用异常终止并从checkpoint恢复重启后，对于进入Kafka数据的任务，系统默认优先处理应用终止前（A段时间）未完成的任務和应用终止到重启完成这段时间内（B段时间）进入Kafka数据生成的任务，最后再处理应用重启完成后（C段时间）进入Kafka数据生成的任务。并且对于B段时间进入Kafka的数据，Spark将按照终止时间（batch时间）生成相应个数的任务，其中第一个任务读取全部数据，其余任务可能不读取数据，造成任务处理压力不均匀。

若A段时间的任务和B段时间任务处理得较慢，则会影响C段时间任务的處理。针对上述场景，Spark提供Kafka后进先出功能。

图 26-4 Spark Streaming 应用重启时间轴



开启此功能后，Spark将优先调度C段时间内的任务，若存在多个C段任务，则按照任务产生的先后顺序调度执行，再执行A段时间和B段时间的任务。另外，对于B段时间进入Kafka的数据，Spark除了按照终止时间生成相应任务，还将这个期间进入Kafka的所有数据均匀分配到各个任务，避免任务处理压力不均匀。

约束条件：

- 目前该功能只适用于Spark Streaming中的Direct方式，且执行结果与上一个batch时间处理结果没有依赖关系（即无state操作，如updatestatebykey）。对多条数据输入流，需要相对独立无依赖的状态，否则可能导致数据切分后结果发生变化。
- Kafka后进先出功能的开启要求应用只能对接Kafka输入源。
- 若提交应用的同时开启Kafka后进先出和流控功能，对于B段时间进入Kafka的数据，将不启动流控功能，以确保读取这些数据的任务调度优先级最低。应用重新启动后C段时间的任务启用流控功能。

配置描述

在Spark Driver端的“spark-defaults.conf”配置文件中进行设置。

表 26-22 参数说明

参数	说明	默认值
spark.streaming.kafka.direct.lifo	配置是否开启Kafka后进先出功能。	false
spark.streaming.kafka010.inputstream.class	获取解耦在FusionInsight侧的类	org.apache.spark.streaming.kafka010.HWDirectKafkaInputDStream

26.7.2 配置对接 Kafka 可靠性

配置场景

Spark Streaming对接Kafka时，当Spark Streaming应用重启后，应用根据上一次读取的topic offset作为起始位置和当前topic最新的offset作为结束位置从Kafka上读取数据的。

Kafka服务的topic的leader异常后，若Kafka的leader和follower的offset相差太大，用户重启Kafka服务，Kafka的follower和leader相互切换，则Kafka服务重启后，topic的offset变小。

- 若Spark Streaming应用一直在运行，由于Kafka上topic的offset变小，会导致读取Kafka数据的起始位置比结束位置大，这样将无法从Kafka读取数据，应用报错。
- 若在重启Kafka服务前，先停止Spark Streaming应用，等Kafka重启后，再重启Spark Streaming应用使应用从checkpoint恢复。此时，Spark Streaming应用会记录终止前读取到的offset位置，以此为基准读取后面的数据，而Kafka offset变小（例如从10万变成1万），Spark Streaming会等待Kafka leader的offset增长至10万之后才会去消费，导致新发送的offset在1万至10万之间的数据丢失。

针对上述背景，提供配置Streaming对接Kafka更高级别的可靠性。对接Kafka可靠性功能开启后，上述场景处理方式如下。

- 若Spark Streaming应用在运行应用时Kafka上topic的offset变小，则会将Kafka上topic最新的offset作为读取Kafka数据的起始位置，继续读取后续的数据。
对于已经生成但未调度处理的任務，若读取的Kafka offset区间大于Kafka上topic的最新offset，则该任务会运行失败。

📖 说明

若任务失败过多，则会将executor加入黑名单，从而导致后续的任务无法部署运行。此时用户可以通过配置“spark.blacklist.enabled”参数关闭黑名单功能，黑名单功能默认为开启。

- 若Kafka上topic的offset变小后，Spark Streaming应用进行重启恢复终止前未处理完的任务若读取的Kafka offset区间大于Kafka上topic的最新offset，则该任务直接丢弃，不进行处理。

📖 说明

若Streaming应用中使用了state函数，则不允许开启对接Kafka可靠性功能。

配置描述

在Spark客户端的“spark-defaults.conf”配置文件中设置。

表 26-23 参数说明

参数	说明	默认值
spark.streaming.Kafka.reliability	Spark Streaming对接Kafka是否开启可靠性功能： <ul style="list-style-type: none">• true：开启可靠性功能• false：不开启可靠性功能	false

26.8 Spark Core 性能调优

26.8.1 Spark Core 数据序列化

操作场景

Spark支持两种方式的序列化：

- Java原生序列化JavaSerializer
- Kryo序列化KryoSerializer

序列化对于Spark应用的性能来说，具有很大的影响。在特定的数据格式的情况下，KryoSerializer的性能可以达到JavaSerializer的10倍以上，而对于一些Int之类的基本数据类型数据，性能的提升就几乎可以忽略。

KryoSerializer依赖Twitter的Chill库来实现，相对于JavaSerializer，主要的问题在于不是所有的Java Serializable对象都能支持，兼容性不好，所以需要手动注册类。

序列化功能用在两个地方：序列化任务和序列化数据。Spark任务序列化只支持JavaSerializer，数据序列化支持JavaSerializer和KryoSerializer。

操作步骤

Spark程序运行时，在shuffle和RDD Cache等过程中，会有大量的数据需要序列化，默认使用JavaSerializer，通过配置让KryoSerializer作为数据序列化器来提升序列化性能。

在开发应用程序时，添加如下代码来使用KryoSerializer作为数据序列化器。

- 实现类注册器并手动注册类。

```
package com.etl.common;

import com.esotericsoftware.kryo.Kryo;
import org.apache.spark.serializer.KryoRegistrator;

public class DemoRegistrator implements KryoRegistrator
{
    @Override
    public void registerClasses(Kryo kryo)
    {
        //以下为示例类，请注册自定义的类
        kryo.register(AggrateKey.class);
    }
}
```

```
kryo.register(AggrateValue.class);  
    }  
}
```

您可以在Spark客户端对spark.kryo.registrationRequired参数进行配置，设置是否需要Kryo注册序列化。

当参数设置为true时，如果工程中存在未被序列化的类，则会发生异常。如果设置为false（默认值），Kryo会自动将未注册的类名写到对应的对象中。此操作会对系统性能造成影响。设置为true时，用户需手动注册类，针对未序列化的类，系统不会自动写入类名，而是发生异常，相对比false，其性能较好。

- 配置KryoSerializer作为数据序列化器和类注册器。

```
val conf = new SparkConf()  
conf.set("spark.serializer", "org.apache.spark.serializer.KryoSerializer")  
.set("spark.kryo.registrator", "com.etl.common.DemoRegistrator")
```

26.8.2 Spark Core 内存调优

操作场景

Spark是内存计算框架，计算过程中内存不够对Spark的执行效率影响很大。可以通过监控GC（Garbage Collection），评估内存中RDD的大小来判断内存是否变成性能瓶颈，并根据情况优化。

监控节点进程的GC情况（在客户端的conf/spark-default.conf配置文件中，在spark.driver.extraJavaOptions和spark.executor.extraJavaOptions配置项中添加参数：“-verbose:gc -XX:+PrintGCDetails -XX:+PrintGCTimeStamps”

），如果频繁出现Full GC，需要优化GC。把RDD做Cache操作，通过日志查看RDD在内存中的大小，如果数据太大，需要改变RDD的存储级别来优化。

操作步骤

- 优化GC，调整老年代和新生代的大小和比例。在客户端的conf/spark-default.conf配置文件中，在spark.driver.extraJavaOptions和spark.executor.extraJavaOptions配置项中添加参数：-XX:NewRatio。如，“-XX:NewRatio=2”，则新生代占整个堆空间的1/3，老年代占2/3。
- 开发Spark应用程序时，优化RDD的数据结构。
 - 使用原始类型数组替代集合类，如可使用fastutil库。
 - 避免嵌套结构。
 - Key尽量不要使用String。
- 开发Spark应用程序时，建议序列化RDD。

RDD做cache时默认是不序列化数据的，可以通过设置存储级别来序列化RDD减小内存。例如：

```
testRDD.persist(StorageLevel.MEMORY_ONLY_SER)
```

26.8.3 Spark Core 内存调优

操作场景

并行度控制任务的数量，影响shuffle操作后数据被切分成的块数。调整并行度让任务的数量和每个任务处理的数据与机器的处理能力达到合适。

查看CPU使用情况和内存占用情况，当任务和数据不是平均分布在各节点，而是集中在个别节点时，可以增大并行度使任务和数据更均匀的分布在各个节点。增加任务的并行度，充分利用集群机器的计算能力，一般并行度设置为集群CPU总和的2-3倍。

操作步骤

并行度可以通过如下三种方式来设置，用户可以根据实际的内存、CPU、数据以及应用程序逻辑的情况调整并行度参数。

- 在会产生shuffle的操作函数内设置并行度参数，优先级最高。
`testRDD.groupByKey(24)`
- 在代码中配置“spark.default.parallelism”设置并行度，优先级次之。
`val conf = new SparkConf()
conf.set("spark.default.parallelism", 24)`
- 在“\$SPARK_HOME/conf/spark-defaults.conf”文件中配置“spark.default.parallelism”的值，优先级最低。
`spark.default.parallelism 24`

26.8.4 配置 Spark Core 广播变量

操作场景

Broadcast（广播）可以把数据集合分发到每一个节点上，Spark任务在执行过程中要使用这个数据集合时，就会在本地查找Broadcast过来的数据集合。如果不使用Broadcast，每次任务需要数据集合时，都会把数据序列化到任务里面，不但耗时，还使任务变得很大。

1. 每个任务分片在执行中都需要同一份数据集合时，就可以把公共数据集Broadcast到每个节点，让每个节点在本地都保存一份。
2. 大表和小表做join操作时可以把小表Broadcast到各个节点，从而就可以把join操作转变成普通的操作，减少了shuffle操作。

操作步骤

在开发应用程序时，添加如下代码，将“testArr”数据广播到各个节点。

```
def main(args: Array[String]) {  
  ...  
  val testArr: Array[Long] = new Array[Long](200)  
  val testBroadcast: Broadcast[Array[Long]] = sc.broadcast(testArr)  
  val resultRdd: RDD[Long] = inpputRdd.map(input => handleData(testBroadcast, input))  
  ...  
}  
  
def handleData(broadcast: Broadcast[Array[Long]], input: String) {  
  val value = broadcast.value  
  ...  
}
```

26.8.5 配置 Spark Executor 堆内存参数

配置场景

当分配的内存太小或者被更高优先级的进程抢占资源时，会出现物理内存超限的情况。调整如下参数，可以防止物理内存超限。

配置描述

参数入口:

在应用提交时通过“--conf”设置这些参数，或者在客户端的“spark-defaults.conf”配置文件中调整如下参数。

表 26-24 参数说明

参数	说明	默认值
spark.executor.memoryOverhead	用于指定每个executor的堆外内存大小(MB)，增大该参数值，可以防止物理内存超限。该值是通过 $\max(384, \text{executor-memory} * 0.1)$ 计算所得，最小值为384。	1024

26.8.6 使用 External Shuffle Service 提升 Spark Core 性能

操作场景

Spark系统在运行含shuffle过程的应用时，Executor进程除了运行task，还要负责写shuffle数据以及给其他Executor提供shuffle数据。当Executor进程任务过重，导致触发GC（Garbage Collection）而不能为其他Executor提供shuffle数据时，会影响任务运行。

External shuffle Service是长期存在于NodeManager进程中的一个辅助服务。通过该服务来抓取shuffle数据，减少了Executor的压力，在Executor GC的时候也不会影响其他Executor的任务运行。

操作步骤

步骤1 登录FusionInsight Manager系统。

步骤2 选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”。单击“全部配置”。

步骤3 选择“SparkResource2x > 默认”，修改以下参数：

表 26-25 参数列表

参数	默认值	修改结果
spark.shuffle.service.enabled	false	true

步骤4 重启Spark2x服务，配置生效。

📖 说明

如果需要在Spark2x客户端用External Shuffle Service功能，需要重新下载并安装Spark2x客户端。

----结束

26.8.7 配置 Yarn 模式下 Spark 动态资源调度

操作场景

对于Spark应用来说，资源是影响Spark应用执行效率的一个重要因素。当一个长期运行的服务（比如JDBCServer），若分配给它多个Executor，可是却没有任何任务分配给它，而此时有其他的应用却资源紧张，这就造成了很大的资源浪费和资源不合理的调度。

动态资源调度就是为了解决这种场景，根据当前应用任务的负载情况，实时的增减Executor个数，从而实现动态分配资源，使整个Spark系统更加健康。

操作步骤

步骤1 需要先配置External shuffle service。

步骤2 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置 > 全部配置”。在搜索框中输入“spark.dynamicAllocation.enabled”参数名称，将JDBCServer下的该参数值设置为“true”，表示开启动态资源调度功能。

----结束

下面是一些可选配置，如表26-26所示。

表 26-26 动态资源调度参数

配置项	说明	默认值
spark.dynamicAllocation.minExecutors	最小Executor个数。	0
spark.dynamicAllocation.initialExecutors	初始Executor个数。	0
spark.dynamicAllocation.maxExecutors	最大Executor个数。	2048
spark.dynamicAllocation.schedulerBacklogTimeout	调度第一次超时时间。	1s
spark.dynamicAllocation.sustainedSchedulerBacklogTimeout	调度第二次及之后超时时间。	1s
spark.dynamicAllocation.executorIdleTimeout	普通Executor空闲超时时间。	60s
spark.dynamicAllocation.cachedExecutorIdleTimeout	含有cached blocks的Executor空闲超时时间。	<ul style="list-style-type: none">JDBCServer2x: 2147483647sIndexServer2x: 2147483647sSparkResource2x: 120

📖 说明

使用动态资源调度功能，必须配置External Shuffle Service。

26.8.8 调整 Spark Core 进程参数

操作场景

Spark on Yarn模式下，有Driver、ApplicationMaster、Executor三种进程。在任务调度和运行的过程中，Driver和Executor承担了很大的责任，而ApplicationMaster主要负责container的启停。

因而Driver和Executor的参数配置对Spark应用的执行有着很大的影响意义。用户可通过如下操作对Spark集群性能做优化。

操作步骤

步骤1 配置Driver内存。

Driver负责任务的调度，和Executor、AM之间的消息通信。当任务数变多，任务平行度增大时，Driver内存都需要相应增大。

您可以根据实际任务数量的多少，为Driver设置一个合适的内存。

- 将“spark-defaults.conf”中的“spark.driver.memory”配置项设置为合适大小。
- 在使用spark-submit命令时，添加“--driver-memory MEM”参数设置内存。

步骤2 配置Executor个数。

每个Executor每个核同时能跑一个task，所以增加了Executor的个数相当于增大了任务的并发度。在资源充足的情况下，可以相应增加Executor的个数，以提高运行效率。

- 将“spark-defaults.conf”中的“spark.executor.instance”配置项或者“spark-env.sh”中的“SPARK_EXECUTOR_INSTANCES”配置项设置为合适大小。
- 在使用spark-submit命令时，添加“--num-executors NUM”参数设置Executor个数。

步骤3 配置Executor核数。

每个Executor多个核同时能跑多个task，相当于增大了任务的并发度。但是由于所有核共用Executor的内存，所以要在内存和核数之间做好平衡。

- 将“spark-defaults.conf”中的“spark.executor.cores”配置项或者“spark-env.sh”中的“SPARK_EXECUTOR_CORES”配置项设置为合适大小。
- 在使用spark-submit命令时，添加“--executor-cores NUM”参数设置核数。

步骤4 配置Executor内存。

Executor的内存主要用于任务执行、通信等。当一个任务很大的时候，可能需要较多资源，因而内存也可以做相应的增加；当一个任务较小运行较快时，就可以增大并发度减少内存。

- 将“spark-defaults.conf”中的“spark.executor.memory”配置项或者“spark-env.sh”中的“SPARK_EXECUTOR_MEMORY”配置项设置为合适大小。
- 在使用spark-submit命令时，添加“--executor-memory MEM”参数设置内存。

----结束

示例

- 在执行spark wordcount计算中。1.6T数据，250个executor。
在默认参数下执行失败，出现Futures timed out和OOM错误。
因为数据量大，task数多，而wordcount每个task都比较小，完成速度快。当task数多时driver端相应的一些对象就变大了，而且每个task完成时executor和driver都要通信，这就会导致由于内存不足，进程之间通信断连等问题。
当把Driver的内存设置到4g时，应用成功跑完。
- 使用JDBCServer执行TPC-DS测试套，默认参数配置下也报了很多错误：Executor Lost等。而当配置Driver内存为30g，executor核数为2，executor个数为125，executor内存为6g时，所有任务才执行成功。

26.8.9 Spark DAG 设计规范说明

操作场景

合理的设计程序结构，可以优化执行效率。在程序编写过程中要尽量减少shuffle操作，合并窄依赖操作。

操作步骤

以“同行车判断”例子讲解DAG设计的思路。

- **数据格式：**通过收费站时间、车牌号、收费站编号.....
- **逻辑：**以下两种情况下判定这两辆车是同行车：
 - 如果两辆车都通过相同序列的收费站，
 - 通过同一收费站之间的时间差小于一个特定的值。

该例子有两种实现模式，其中实现1的逻辑如图26-5所示，实现2的逻辑如图26-6所示。

图 26-5 实现 1 逻辑



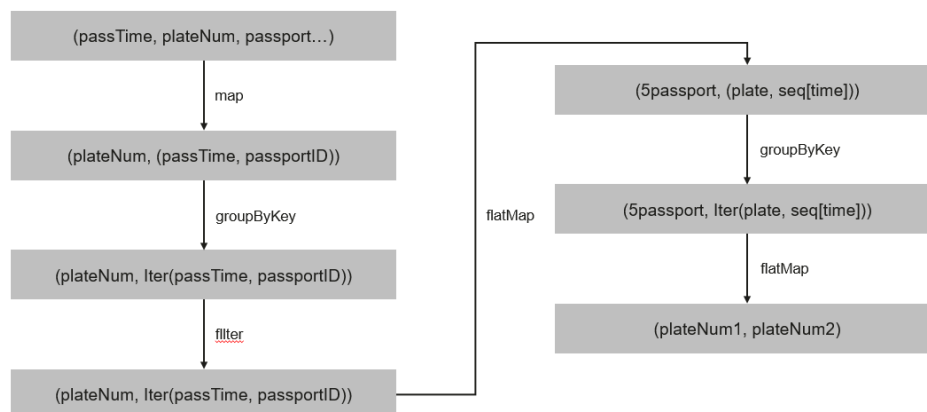
实现1的逻辑说明：

1. 根据车牌号聚合该车通过的所有收费站并排序，处理后数据如下：
车牌号1，[(通过时间, 收费站3), (通过时间, 收费站2), (通过时间, 收费站4), (通过时间, 收费站5)]
2. 标识该收费站是这辆车通过的第几个收费站。
(收费站3, (车牌号1, 通过时间, 通过的第1个收费站))
(收费站2, (车牌号1, 通过时间, 通过的第2个收费站))
(收费站4, (车牌号1, 通过时间, 通过的第3个收费站))
(收费站5, (车牌号1, 通过时间, 通过的第4个收费站))
3. 根据收费站聚合数据。
收费站1，[(车牌号1, 通过时间, 通过的第1个收费站), (车牌号2, 通过时间, 通过的第5个收费站), (车牌号3, 通过时间, 通过的第2个收费站)]
4. 判断两辆车通过该收费站的时间差是否满足同行车的要求，如果满足则取出这两辆车。
(车牌号1, 车牌号2), (通过的第1个收费站, 通过的第5个收费站)
(车牌号1, 车牌号3), (通过的第1个收费站, 通过的第2个收费站)
5. 根据通过相同收费站的两辆车的车牌号聚合数据，如下：
(车牌号1, 车牌号2), [(通过的第1个收费站, 通过的第5个收费站), (通过的第2个收费站, 通过的第6个收费站), (通过的第1个收费站, 通过的第7个收费站), (通过的第3个收费站, 通过的第8个收费站)]
6. 如果车牌号1和车牌号2通过相同收费站是顺序排列的（比如收费站3、4、5是车牌1通过的第1、2、3个收费站，是车牌2通过的第6、7、8个收费站）且数量大于同行车要求的数量则这两辆车是同行车。

实现1逻辑的缺点：

- 逻辑复杂
- 实现过程中shuffle操作过多，对性能影响较大。

图 26-6 实现 2 逻辑



实现2的逻辑说明：

1. 根据车牌号聚合该车通过的所有收费站并排序，处理后数据如下：

- 车牌号1, [(通过时间, 收费站3), (通过时间, 收费站2), (通过时间, 收费站4), (通过时间, 收费站5)]
- 根据同行车要通过的收费站数量（例子里为3）分段该车通过的收费站序列，如上面的数据被分解成：
收费站3->收费站2->收费站4, (车牌号1, [收费站3时间, 收费站2时间, 收费站4时间])
收费站2->收费站4->收费站5, (车牌号1, [收费站2时间, 收费站4时间, 收费站5时间])
 - 把通过相同收费站序列的车辆聚合，如下：
收费站3->收费站2->收费站4, [(车牌号1, [收费站3时间, 收费站2时间, 收费站4时间]), (车牌号2, [收费站3时间, 收费站2时间, 收费站4时间]), (车牌号3, [收费站3时间, 收费站2时间, 收费站4时间])]
 - 判断通过相同序列收费站的车辆通过相同收费站的时间差是不是满足同行车的要求，如果满足则说明是同行车。

实现2的优点如下：

- 简化了实现逻辑。
- 减少了一个groupByKey，也就减少了一次shuffle操作，提升了性能。

26.8.10 经验总结

使用 mapPartitions，按每个分区计算结果

如果每条记录的开销太大，例：

```
rdd.map{x=>conn=getDBConn;conn.write(x.toString);conn.close}
```

则可以使用MapPartitions，按每个分区计算结果，如

```
rdd.mapPartitions(records => conn.getDBConn;for(item <- records)  
write(item.toString); conn.close)
```

使用mapPartitions可以更灵活地操作数据，例如对一个很大的数据求TopN，当N不是很大时，可以先使用mapPartitions对每个partition求TopN，collect结果到本地之后再排序取TopN。这样相比直接对全量数据做排序取TopN效率要高很多。

使用 coalesce 调整分片的数量

coalesce可以调整分片的数量。coalesce函数有两个参数：

```
coalesce(numPartitions: Int, shuffle: Boolean = false)
```

当shuffle为true的时候，函数作用与repartition(numPartitions: Int)相同，会将数据通过Shuffle的方式重新分区；当shuffle为false的时候，则只是简单的将父RDD的多个partition合并到同一个task进行计算，shuffle为false时，如果numPartitions大于父RDD的切片数，那么分区不会重新调整。

遇到下列场景，可选择使用coalesce算子：

- 当之前的操作有很多filter时，使用coalesce减少空运行的任务数量。此时使用coalesce(numPartitions, false)，numPartitions小于父RDD切片数。
- 当输入切片个数太大，导致程序无法正常运行时使用。
- 当任务数过大时候Shuffle压力太大导致程序挂住不动，或者出现linux资源受限的问题。此时需要对数据重新进行分区，使用coalesce(numPartitions, true)。

localDir 配置

Spark的Shuffle过程需要写本地磁盘，Shuffle是Spark性能的瓶颈，I/O是Shuffle的瓶颈。配置多个磁盘则可以并行的把数据写入磁盘。如果节点中挂载多个磁盘，则在每个磁盘配置一个Spark的localDir，这将有效分散Shuffle文件的存放，提高磁盘I/O的效率。如果只有一个磁盘，配置了多个目录，性能提升效果不明显。

Collect 小数据

大数据量不适用collect操作。

collect操作会将Executor的数据发送到Driver端，因此使用collect前需要确保Driver端内存足够，以免Driver进程发生OutOfMemory异常。当不确定数据量大小，可使用saveAsTextFile等操作把数据写入HDFS中。只有在能够大致确定数据大小且driver内存充足的时候，才能使用collect。

使用 reduceByKey

reduceByKey会在Map端做本地聚合，使得Shuffle过程更加平缓，而groupByKey等Shuffle操作不会在Map端做聚合。因此能使用reduceByKey的地方尽量使用该算子，避免出现groupByKey().map(x=>(x._1,x._2.size))这类实现方式。

广播 map 代替数组

当每条记录需要查表，如果是Driver端用广播方式传递的数据，数据结构优先采用set/map而不是Iterator，因为Set/Map的查询速率接近O(1)，而Iterator是O(n)。

数据倾斜

当数据发生倾斜（某一部分数据量特别大），虽然没有GC（Garbage Collection，垃圾回收），但是task执行时间严重不一致。

- 需要重新设计key，以更小粒度的key使得task大小合理化。
- 修改并行度。

优化数据结构

- 把数据按列存放，读取数据时就可以只扫描需要的列。
- 使用Hash Shuffle时，通过设置spark.shuffle.consolidateFiles为true，来合并shuffle中间文件，减少shuffle文件的数量，减少文件IO操作以提升性能。最终文件数为reduce tasks数目。

26.9 Spark SQL 性能调优

26.9.1 Spark SQL join 优化

操作场景

Spark SQL中，当对两个表进行join操作时，利用Broadcast特性（见“使用广播变量”章节），将被广播的表Broadcast到各个节点上，从而转变成非shuffle操作，提高任务执行性能。

📖 说明

这里join操作，只指inner join。

操作步骤

在Spark SQL中进行Join操作时，可以按照以下步骤进行优化。为了方便说明，设表A和表B，且A、B表都有个名为name的列。对A、B表进行join操作。

1. 估计表的大小。

根据每次加载数据的大小，来估计表大小。

也可以在Hive的数据库存储路径下直接查看表的大小。首先在Spark的配置文件“hive-site.xml”中，查看Hive的数据库路径的配置，默认为“/user/hive/warehouse”。Spark服务多实例默认数据库路径为“/user/hive/warehouse”，例如“/user/hive1/warehouse”。

```
<property>
  <name>hive.metastore.warehouse.dir</name>
  <value>${test.warehouse.dir}</value>
  <description></description>
</property>
```

然后通过hadoop命令查看对应表的大小。如查看表A的大小命令为：

```
hadoop fs -du -s -h ${test.warehouse.dir}/a
```

📖 说明

进行广播操作，需要至少有一个表不是空表。

2. 配置自动广播的阈值。

Spark中，判断表是否广播的阈值为10485760（即10M）。如果两个表的大小至少有一个小于10M时，可以跳过该步骤。

自动广播阈值的配置参数介绍，见表26-27。

表 26-27 参数介绍

参数	默认值	描述
spark.sql.autoBroadcastJoinThreshold	10485760	当进行join操作时，配置广播的最大值。
	0	<ul style="list-style-type: none"> 当SQL语句中涉及的表中相应字段的大小小于该值时，进行广播。 配置为-1时，将不进行广播。

配置自动广播阈值的方法：

- 在Spark的配置文件“spark-defaults.conf”中，设置“spark.sql.autoBroadcastJoinThreshold”的值。

```
spark.sql.autoBroadcastJoinThreshold = <size>
```

- 利用Hive CLI命令，设置阈值。在运行Join操作时，提前运行下面语句：
SET spark.sql.autoBroadcastJoinThreshold=<size>;

3. 进行join操作。

- 两个表的大小都小于阈值。

- A表的字节数小于B表，则运行B join A，如
`SELECT A.name FROM B JOIN A ON A.name = B.name;`
 - 否则运行A join B。
`SELECT A.name FROM A JOIN B ON A.name = B.name;`
 - 一个表大于阈值一个表小于阈值。
将小表进行BroadCast操作。
 - 两个表的大小都大于阈值。
比较查询所涉及的字段大小与阈值的大小。
 - 若某表中涉及字段的大小小于阈值，将该表相应数据进行广播。
 - 若两表中涉及字段的大小都大于阈值，则不进行广播。
4. （可选）如下两种场景，需要执行Analyze命令（***ANALYZE TABLE tableName COMPUTE STATISTICS noscan;***）更新表元数据后进行广播。
- 需要广播的表是分区表，新建表且文件类型为非Parquet文件类型。
 - 需要广播的表是分区表，更新表数据后。

参考信息

被广播的表执行超时，导致任务结束。

默认情况下，BroadCastJoin只允许被广播的表计算5分钟，超过5分钟该任务会出现超时异常，而这个时候被广播的表的broadcast任务依然在执行，造成资源浪费。

这种情况下，有两种方式处理：

- 调整“spark.sql.broadcastTimeout”的数值，加大超时的时间限制。
- 降低“spark.sql.autoBroadcastJoinThreshold”的数值，不使用BroadCastJoin的优化。

26.9.2 优化数据倾斜场景下的 Spark SQL 性能

配置场景

在Spark SQL多表Join的场景下，会存在关联键严重倾斜的情况，导致Hash分桶后，部分桶中的数据远高于其它分桶。最终导致部分Task过重，跑得很慢；其它Task过轻，跑得很快。一方面，数据量大Task运行慢，使得计算性能低；另一方面，数据量少的Task在运行完成后，导致很多CPU空闲，造成CPU资源浪费。

通过如下配置项可开启自动进行数据倾斜处理功能，通过将Hash分桶后数据量很大的、且超过数据倾斜阈值的分桶拆散，变成多个task处理一个桶的数据机制，提高CPU资源利用率，提高系统性能。

📖 说明

未产生倾斜的数据，将采用原有方式进行分桶并运行。

使用约束：

- 只支持两表Join的场景。
- 不支持FULL OUTER JOIN的数据倾斜处理。
示例：执行下面SQL语句，a表倾斜或b表倾斜都无法触发该优化。

select aid FROM a FULL OUTER JOIN b ON aid=bid;

- 不支持LEFT OUTER JOIN的右表倾斜处理。
示例：执行下面SQL语句，b表倾斜无法触发该优化。

select aid FROM a LEFT OUTER JOIN b ON aid=bid;

- 不支持RIGHT OUTER JOIN的左表倾斜处理。
示例：执行下面SQL语句，a表倾斜无法触发该优化。

select aid FROM a RIGHT OUTER JOIN b ON aid=bid;

配置描述

在Spark Driver端的“spark-defaults.conf”配置文件中添加如下表格中的参数。

表 26-28 参数说明

参数	描述	默认值
spark.sql.adaptive.enabled	自适应执行特性的总开关。 注意：AQE特性与DPP（动态分区裁剪）特性同时开启时，SparkSQL任务执行中会优先执行DPP特性，从而使得AQE特性不生效。集群中DPP特性是默认开启的，因此开启AQE特性的同时，需要将DPP特性关闭。	false
spark.sql.optimize.r.dynamicPartitionPruning.enabled	动态分区裁剪功能的开关。	true
spark.sql.adaptive.skewJoin.enabled	当此配置为true且spark.sql.adaptive.enabled设置为true时，启用运行时自动处理join运算中的数据倾斜功能。	true
spark.sql.adaptive.skewJoin.skewedPartitionFactor	此配置为一个倍数因子，用于判定分区是否为数据倾斜分区。单个分区被判定为数据倾斜分区的条件为：当一个分区的数据大小超过除此分区外其他所有分区大小的中值与该配置的乘积，并且大小超过spark.sql.adaptive.skewJoin.skewedPartitionThresholdInBytes配置值时，此分区被判定为数据倾斜分区	5
spark.sql.adaptive.skewJoin.skewedPartitionThresholdInBytes	分区大小（单位：字节）大于该阈值且大于spark.sql.adaptive.skewJoin.skewedPartitionFactor与分区中值的乘积，则认为该分区存在倾斜。理想情况下，此配置应大于spark.sql.adaptive.advisoryPartitionSizeInBytes。	256MB
spark.sql.adaptive.shuffle.targetPostShuffleInputSize	每个task处理的shuffle数据的最小数据量。单位：Byte。	67108864

26.9.3 优化小文件场景下的 Spark SQL 性能

配置场景

Spark SQL的表中，经常会存在很多小文件（大小远小于HDFS块大小），每个小文件默认对应Spark中的一个Partition，也就是一个Task。在很多小文件场景下，Spark会起很多Task。当SQL逻辑中存在Shuffle操作时，会大大增加hash分桶数，严重影响性能。

在小文件场景下，您可以通过如下配置手动指定每个Task的数据量（Split Size），确保不会产生过多的Task，提高性能。

📖 说明

当SQL逻辑中不包含Shuffle操作时，设置此配置项，不会有明显的性能提升。

配置描述

要启动小文件优化，在Spark客户端的“spark-defaults.conf”配置文件中设置。

表 26-29 参数说明

参数	描述	默认值
spark.sql.files.maxPartitionBytes	在读取文件时，将单个分区打包的最大字节数。 单位：byte。	134217728 (即128M)
spark.files.openCostInBytes	打开文件的预估成本，按照同一时间能够扫描的字节数来测量。当一个分区写入多个文件时使用。高估更好，这样小文件分区将比大文件分区更先被调度。	4M

26.9.4 Spark INSERT SELECT 语句调优

操作场景

在以下几种情况下，执行INSERT...SELECT操作可以进行一定的调优操作。

- 查询的数据是大量的小文件。
- 查询的数据是较多的大文件。
- 在Beeline/JDBCServer模式下使用非Spark用户操作。

操作步骤

可对INSERT...SELECT操作做如下的调优操作。

- 如果建的是Hive表，将存储类型设为Parquet，从而减少执行INSERT...SELECT语句的时间。

- 建议使用spark-sql或者在Beeline/JDBCServer模式下使用spark用户来执行INSERT...SELECT操作，避免执行更改文件owner的操作，从而减少执行INSERT...SELECT语句的时间。

📖 说明

在Beeline/JDBCServer模式下，executor的用户跟driver是一致的，driver是JDBCServer服务的一部分，是由spark用户启动的，因此其用户也是spark用户，且当前无法实现在运行时将Beeline端的用户透传到executor，因此使用非spark用户时需要对文件进行更改owner为Beeline端的用户，即实际用户。

- 如果查询的数据是大量的小文件将会产生大量map操作，从而导致输出存在大量的小文件，在执行重命名文件操作时将会耗费较多时间，此时可以通过设置“spark.sql.files.maxPartitionBytes”与“spark.files.openCostInBytes”来设置一个partiton读取的最大字节，在一个partition中合并多个小文件来减少输出文件数及执行重命名文件操作的时间，从而减少执行INSERT...SELECT语句的时间。

📖 说明

上述优化操作并不能解决全部的性能问题，对于以下场景仍然需要较多时间：
对于动态分区表，如果其分区数非常多，那么也需要执行较长的时间。

26.9.5 动态分区插入场景内存优化

操作场景

SparkSQL在往动态分区表中插入数据时，分区数越多，单个Task生成的HDFS文件越多，则元数据占用的内存也越多。这就导致程序GC（Gabbage Collection）严重，甚至发生OOM（Out of Memory）。

经测试证明：10240个Task，2000个分区，在执行HDFS文件从临时目录rename到目标目录动作前，FileStatus元数据大小约29G。为避免以上问题，可修改SQL语句对数据进行重分区，以减少HDFS文件个数。

操作步骤

在动态分区语句中加入**distribute by**，by值为分区字段。

示例如下：

```
insert into table store_returns partition (sr_returned_date_sk) select
sr_return_time_sk,sr_item_sk,sr_customer_sk,sr_demo_sk,sr_hdemo_sk,sr_addr_sk,
sr_store_sk,sr_reason_sk,sr_ticket_number,sr_return_quantity,sr_return_amt,sr_return_tax,sr_return_amt_inc_tax,sr_fee,sr_return_ship_cost,sr_refunded_cash,sr_reversed_charge,sr_store_credit,sr_net_loss,sr_returned_date_sk from $
{SOURCE}.store_returns distribute by sr_returned_date_sk;
```

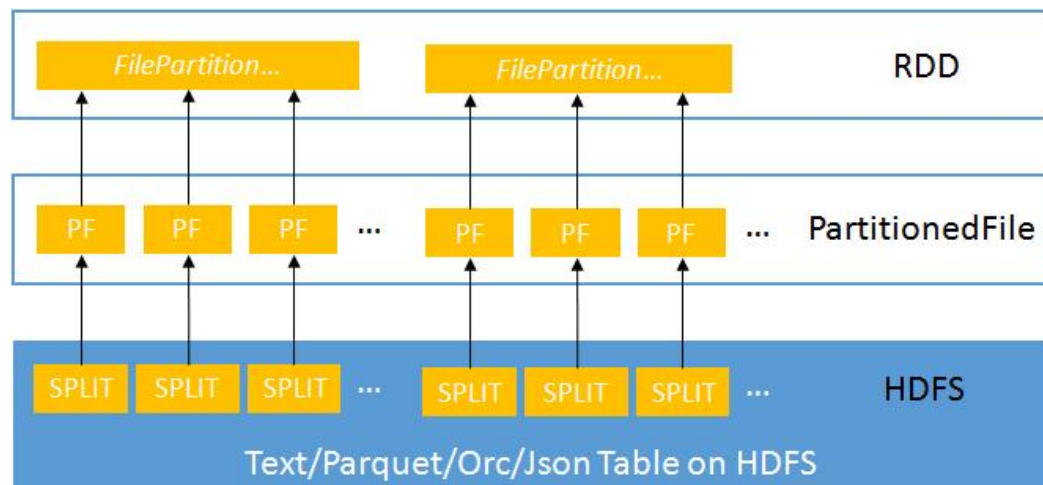
26.9.6 小文件优化

操作场景

Spark SQL表中，经常会存在很多小文件（大小远小于HDFS的块大小），每个小文件默认对应Spark中的一个Partition，即一个Task。在有很多小文件时，Spark会启动很多Task，此时当SQL逻辑中存在Shuffle操作时，会大大增加hash分桶数，严重影响系统性能。

针对小文件很多的场景，DataSource在创建RDD时，先将Table中的split生成PartitionedFile，再将这些PartitionedFile进行合并。即将多个PartitionedFile组成一个partition，从而减少partition数量，避免在Shuffle操作时生成过多的hash分桶，如图26-7所示。

图 26-7 小文件合并



操作步骤

要启动小文件优化，在Spark客户端的“spark-defaults.conf”配置文件中进行设置。

表 26-30 参数介绍

参数	描述	默认值
spark.sql.files.maxPartitionBytes	在读取文件时，将单个分区打包的最大字节数。 单位：byte。	134217728（即128M）
spark.files.openCostInBytes	打开文件的预估成本，按照同一时间能够扫描的字节数来测量。当一个分区写入多个文件时使用。高估更好，这样小文件分区将比大文件分区更先被调度。	4M

26.9.7 聚合算法优化

操作场景

在Spark SQL中支持基于行的哈希聚合算法，即使用快速聚合hashmap作为缓存，以提高聚合性能。hashmap替代了之前的ColumnarBatch支持，从而避免拥有聚合表的宽模式（大量key字段或value字段）时产生的性能问题。

操作步骤

要启动聚合算法优化，在Spark客户端的“spark-defaults.conf”配置文件中设置。

表 26-31 参数介绍

参数	描述	默认值
spark.sql.codegen.aggregate.map.twolevel.enabled	是否开启聚合算法优化： <ul style="list-style-type: none"> • true：开启 • false：不开启 	true

26.9.8 Datasource 表优化

操作场景

将datasource表的分区消息存储到Metastore中，并在Metastore中对分区消息进行处理。

- 优化datasource表，支持对表中分区执行增加、删除和修改等语法，从而增加与Hive的兼容性。
- 支持在查询语句中，把分区裁剪并下压到Metastore上，从而过滤掉不匹配的分区。

示例如下：

```
select count(*) from table where partCol=1; //partCol列为分区列
```

此时，在物理计划中执行TableScan操作时，只处理分区(partCol=1)对应的数据。

操作步骤

要启动Datasource表优化，在Spark客户端的“spark-defaults.conf”配置文件中设置。

表 26-32 参数介绍

参数	描述	默认值
spark.sql.hive.manageFilesourcePartitions	是否启用Metastore分区管理（包括数据源表和转换的Hive表）。 <ul style="list-style-type: none"> • true：启用Metastore分区管理，即数据源表存储分区在Hive中，并在查询语句中使用Metastore修剪分区。 • false：不启用Metastore分区管理。 	true
spark.sql.hive.metastorePartitionPruning	是否支持将predicate下压到Hive Metastore中。 <ul style="list-style-type: none"> • true：支持，目前仅支持Hive表的predicate下压。 • false：不支持 	true

参数	描述	默认值
spark.sql.hive.filesourcePartitionFileCacheSize	启用内存中分区文件元数据的缓存大小。所有表共享一个可以使用指定的num字节进行文件元数据的缓存。 只有当“spark.sql.hive.manageFilesourcePartitions”配置为“true”时，该配置项才会生效。	250 * 1024 * 1024
spark.sql.hive.convertMetastoreOrc	设置ORC表的处理方式： <ul style="list-style-type: none"> false: Spark SQL使用Hive SerDe处理ORC表。 true: Spark SQL使用Spark内置的机制处理ORC表。 	true

26.9.9 合并 CBO 优化

操作场景

Spark SQL默认支持基于规则的优化，但仅仅基于规则优化不能保证Spark选择合适的查询计划。CBO（Cost-Based Optimizer）是一种为SQL智能选择查询计划的技术。通过配置开启CBO后，CBO优化器可以基于表和列的统计信息，进行一系列的估算，最终选择出合适的查询计划。

操作步骤

要使用CBO优化，可以按照以下步骤进行优化。

1. 需要先执行特定的SQL语句来收集所需的表和列的统计信息。

SQL命令如下（根据具体情况选择需要执行的SQL命令）：

- 生成表级别统计信息（扫表）：

ANALYZE TABLE src COMPUTE STATISTICS

生成sizeInBytes和rowCount。

使用ANALYZE语句收集统计信息时，无法计算非HDFS数据源的表的文件大小。

- 生成表级别统计信息（不扫表）：

ANALYZE TABLE src COMPUTE STATISTICS NOSCAN

只生成sizeInBytes，如果原来已经生成过sizeInBytes和rowCount，而本次生成的sizeInBytes和原来的大小一样，则保留rowCount（若存在），否则清除rowCount。

- 生成列级别统计信息

ANALYZE TABLE src COMPUTE STATISTICS FOR COLUMNS a, b, c

生成列统计信息，为保证一致性，会同步更新表统计信息。目前不支持复杂数据类型（如Seq, Map等）和HiveStringType的统计信息生成。

- 显示统计信息

DESC FORMATTED src

在Statistics中会显示“xxx bytes, xxx rows”分别表示表级别的统计信息。也可以通过如下命令显示列统计信息：

DESC FORMATTED src a

使用限制：当前统计信息收集不支持针对分区表的分区级别的统计信息。

2. 在Spark客户端的“spark-defaults.conf”配置文件中[进行表26-33设置](#)。

表 26-33 参数介绍

参数	描述	默认值
spark.sql.cbo.enabled	CBO总开关。 <ul style="list-style-type: none"> • true表示打开， • false表示关闭。 要使用该功能，需确保相关表和列的统计信息已经生成。	false
spark.sql.cbo.joinReorder.enabled	使用CBO来自动调整连续的inner join的顺序。 <ul style="list-style-type: none"> • true: 表示打开 • false: 表示关闭 要使用该功能，需确保相关表和列的统计信息已经生成，且CBO总开关打开。	false
spark.sql.cbo.joinReorder.default.threshold	使用CBO来自动调整连续inner join的表的个数阈值。 如果超出该阈值，则不会调整join顺序。	12

26.9.10 多级嵌套子查询以及混合 Join 的 SQL 调优

操作场景

本章节介绍在多级嵌套以及混合Join SQL查询的调优建议。

前提条件

例如有一个复杂的查询样例如下：

```
select
s_name,
count(1) as numwait
from (
select s_name from (
select
s_name,
t2.l_orderkey,
l_suppkey,
count_suppkey,
max_suppkey
```

```
from
test2 t2 right outer join (
select
s_name,
l_orderkey,
l_suppkey from (
select
s_name,
t1.l_orderkey,
l_suppkey,
count_suppkey,
max_suppkey
from
test1 t1 join (
select
s_name,
l_orderkey,
l_suppkey
from
orders o join (
select
s_name,
l_orderkey,
l_suppkey
from
nation n join supplier s
on
s.s_nationkey = n.n_nationkey
and n.n_name = 'SAUDI ARABIA'
join lineitem l
on
s.s_suppkey = l.l_suppkey
where
l.l_receiptdate > l.l_commitdate
and l.l_orderkey is not null
) l1 on o.o_orderkey = l1.l_orderkey and o.o_orderstatus = 'F'
) l2 on l2.l_orderkey = t1.l_orderkey
) a
where
(count_suppkey > 1)
or ((count_suppkey=1)
and (l_suppkey <> max_suppkey))
) l3 on l3.l_orderkey = t2.l_orderkey
) b
where
(count_suppkey is null)
or ((count_suppkey=1)
and (l_suppkey = max_suppkey))
) c
group by
s_name
order by
numwait desc,
s_name
limit 100;
```

操作步骤

步骤1 分析业务。

从业务入手分析是否可以简化SQL，例如可以通过合并表去减少嵌套的层级和Join的次数。

步骤2 如果业务需求对应的SQL无法简化，则需要配置DRIVER内存：

- 使用spark-submit或者spark-sql运行SQL语句，执行[步骤3](#)。
- 使用spark-beeline运行SQL语句，执行[步骤4](#)。

步骤3 执行SQL语句时，需要添加参数“--driver-memory”，设置内存大小，例如：

```
/spark-sql --master=local[4] --driver-memory=512M -f /tpch.sql
```

步骤4 在执行SQL语句前，请使用MRS集群管理员用户修改内存大小配置。

1. 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”。
2. 单击“全部配置”，并搜索“SPARK_DRIVER_MEMORY”。
3. 修改参数值适当增加内存大小。仅支持整数数值，且需要输入单位M或者G。例如输入512M。

----结束

参考信息

DRIVER内存不足时，查询操作可能遇到以下错误提示信息：

```
2018-02-11 09:13:14,683 | WARN | Executor task launch worker for task 5 | Calling spill() on
RowBasedKeyValueBatch. Will not spill but return 0. |
org.apache.spark.sql.catalyst.expressions.RowBasedKeyValueBatch.spill(RowBasedKeyValueBatch.java:173)
2018-02-11 09:13:14,682 | WARN | Executor task launch worker for task 3 | Calling spill() on
RowBasedKeyValueBatch. Will not spill but return 0. |
org.apache.spark.sql.catalyst.expressions.RowBasedKeyValueBatch.spill(RowBasedKeyValueBatch.java:173)
2018-02-11 09:13:14,704 | ERROR | Executor task launch worker for task 2 | Exception in task 2.0 in stage
1.0 (TID 2) | org.apache.spark.internal.Logging$class.logError(Logging.scala:91)
java.lang.OutOfMemoryError: Unable to acquire 262144 bytes of memory, got 0
    at org.apache.spark.memory.MemoryConsumer.allocateArray(MemoryConsumer.java:100)
    at org.apache.spark.unsafe.map.BytesToBytesMap.allocate(BytesToBytesMap.java:791)
    at org.apache.spark.unsafe.map.BytesToBytesMap.<init>(BytesToBytesMap.java:208)
    at org.apache.spark.unsafe.map.BytesToBytesMap.<init>(BytesToBytesMap.java:223)
    at
org.apache.spark.sql.execution.UnsafeFixedWidthAggregationMap.<init>(UnsafeFixedWidthAggregationMap.j
ava:104)
    at
org.apache.spark.sql.execution.aggregate.HashAggregateExec.createHashMap(HashAggregateExec.scala:307)
    at org.apache.spark.sql.catalyst.expressions.GeneratedClass
$GeneratedIterator.agg_doAggregateWithKeys$(Unknown Source)
    at org.apache.spark.sql.catalyst.expressions.GeneratedClass$GeneratedIterator.processNext(Unknown
Source)
    at org.apache.spark.sql.execution.BufferedRowIterator.hasNext(BufferedRowIterator.java:43)
    at org.apache.spark.sql.execution.WholeStageCodegenExec$$anonfun$8$$anon
$1.hasNext(WholeStageCodegenExec.scala:381)
    at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:408)
    at
org.apache.spark.shuffle.sort.BypassMergeSortShuffleWriter.write(BypassMergeSortShuffleWriter.java:126)
    at org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask.scala:96)
    at org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask.scala:53)
    at org.apache.spark.scheduler.Task.run(Task.scala:99)
    at org.apache.spark.executor.Executor$TaskRunner.run(Executor.scala:325)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1149)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:624)
    at java.lang.Thread.run(Thread.java:748)
```

26.10 Spark Streaming 性能调优

操作场景

Streaming作为一种mini-batch方式的流式处理框架，它主要的特点是：秒级时延和高吞吐量。因此Streaming调优的目标：在秒级延迟的情景下，提高Streaming的吞吐能力，在单位时间处理尽可能多的数据。

📖 说明

本章节适用于输入数据源为Kafka的使用场景。

操作步骤

一个简单的流处理系统由以下三部分组件组成：数据源 + 接收器 + 处理器。数据源为Kafka，接收器为Streaming中的Kafka数据源接收器，处理器为Streaming。

对Streaming调优，就必须使该三个部件的性能都合理化。

• 数据源调优

在实际的应用场景中，数据源为了保证数据的容错性，会将数据保存在本地磁盘中，而Streaming的计算结果全部在内存中完成，数据源很有可能成为流式系统的最大瓶颈点。

对Kafka的性能调优，有以下几个点：

- 使用Kafka-0.8.2以后版本，可以使用异步模式的新Producer接口。
- 配置多个Broker的目录，设置多个IO线程，配置Topic合理的Partition个数。

详情请参见Kafka开源文档中的“性能调优”部分：<http://kafka.apache.org/documentation.html>。

• 接收器调优

Streaming中已有多种数据源的接收器，例如Kafka、Flume、MQTT、ZeroMQ等，其中Kafka的接收器类型最多，也是最成熟一套接收器。

Kafka包括三种模式的接收器API：

- KafkaReceiver：直接接收Kafka数据，进程异常后，可能出现数据丢失。
- ReliableKafkaReceiver：通过ZooKeeper记录接收数据位移。
- DirectKafka：直接通过RDD读取Kafka每个Partition中的数据，数据高可靠。

从实现上来看，DirectKafka的性能更好，实际测试上来看，DirectKafka也确实比其他两个API性能好了不少。因此推荐使用DirectKafka的API实现接收器。

数据接收器作为一个Kafka的消费者，对于它的配置优化，请参见Kafka开源文档：<http://kafka.apache.org/documentation.html>。

• 处理器调优

Spark Streaming的底层由Spark执行，因此大部分对于Spark的调优措施，都可以应用在Spark Streaming之中，例如：

- 数据序列化
- 配置内存
- 设置并行度
- 使用External Shuffle Service提升性能

📖 说明

在做Spark Streaming的性能优化时需注意一点，越追求性能上的优化，Spark Streaming整体的可靠性会越差。例如：

“spark.streaming.receiver.writeAheadLog.enable”配置为“false”的时候，会明显减少磁盘的操作，提高性能，但由于缺少WAL机制，会出现异常恢复时，数据丢失。

因此，在调优Spark Streaming的时候，这些保证数据可靠性的配置项，在生产环境中是不能关闭的。

- **日志归档调优**

参数“spark.eventLog.group.size”用来设置一个应用的JobHistory日志按照指定job个数分组，每个分组会单独创建一个文件记录日志，从而避免应用长期运行时形成单个过大日志造成JobHistory无法读取的问题，设置为“0”时表示不分组。

大部分Spark Streaming任务属于小型job，而且产生速度较快，会导致频繁的分组，产生大量日志小文件消耗磁盘I/O。建议增大此值，例如改为“1000”或更大值。

26.11 Spark 运维管理

26.11.1 快速配置参数

概述

本节介绍Spark2x使用过程中快速配置常用参数和不建议修改的配置参数。

快速配置常用参数

其他参数在安装集群时已进行了适配，以下参数需要根据使用场景进行调整。以下参数除特别指出外，一般在Spark2x客户端的“spark-defaults.conf”文件中配置。

表 26-34 快速配置常用参数

配置项	说明	默认值
spark.sql.parquet.compression.codec	对于非分区parquet表，设置其存储文件的压缩格式。 在JDBCServer服务端的“spark-defaults.conf”配置文件中设置。	snappy
spark.dynamicAllocation.enabled	是否使用动态资源调度，用于根据规模调整注册于该应用的executor的数量。目前仅在YARN模式下有效。 JDBCServer默认值为true，client默认值为false。	false
spark.executor.memory	每个Executor进程使用的内存数量，与JVM内存设置字符串的格式相同（例如：512m，2g）。	4G
spark.sql.autoBroadcastJoinThreshold	当进行join操作时，配置广播的最大值。 <ul style="list-style-type: none"> 当SQL语句中涉及的表中相应字段的大小小于该值时，进行广播。 配置为-1时，将不进行广播。 	10485760
spark.yarn.queue	JDBCServer服务所在的Yarn队列。 在JDBCServer服务端的“spark-defaults.conf”配置文件中设置。	default

配置项	说明	默认值
spark.driver.memory	大集群下推荐配置32~64g驱动程序进程使用的内存数量，即SparkContext初始化的进程（例如：512m, 2g）。	4G
spark.yarn.security.credentials.hbase.enabled	是否打开获取HBase token的功能。如果需要Spark-on-HBase功能，并且配置了安全集群，参数值设置为“true”。否则设置为“false”。	false
spark.serializer	用于序列化将通过网络发送或需要缓存的对象的类以序列化形式展现。 Java序列化的默认值适用于任何Serializable Java对象，但运行速度相当慢，所以建议使用org.apache.spark.serializer.KryoSerializer并配置Kryo序列化。可以是org.apache.spark.serializer.Serializer的任何子类。	org.apache.spark.serializer.JavaSerializer
spark.executor.cores	每个执行者使用的内核个数。 在独立模式和Mesos粗粒度模式下设置此参数。当有足够多的内核时，允许应用程序在同样的worker上执行多个执行程序；否则，在每个worker上，每个应用程序只能运行一个执行程序。	1
spark.shuffle.service.enabled	NodeManager中一个长期运行的辅助服务，用于提升Shuffle计算性能。	false
spark.sql.adaptive.enabled	是否开启自适应执行框架。	false
spark.executor.memoryOverhead	每个执行器要分配的堆内存量（单位为兆字节）。 这是占用虚拟机开销的内存，类似于内部字符串，其他内置开销等等。会随着执行器大小（通常为6-10%）而增长。	1GB
spark.streaming.kafka.direct.lifo	配置是否开启Kafka后进先出功能。	false

不建议修改的参数

以下参数在安装集群时已进行了适配，不建议用户进行修改。

表 26-35 不建议修改的参数说明

配置项	说明	默认值或配置示例
spark.password.factory	用于选择密钥解析方式。	org.apache.spark.om.util.FIPasswordFactory
spark.ssl.ui.protocol	配置ui的ssl协议。	TLSv1.2
spark.yarn.archive	Spark jars的存档，用于分发到YARN缓存。如果设置，此配置值将替换<code>spark.yarn.jars</code>，并存档在所有应用程序的容器中使用。存档应包含其根目录中的jar文件。与以前的选项一样，存档也可以在HDFS上托管，用来加快文件分发速度。	hdfs://hacluster/user/spark2x/jars/8.1.0.1/spark-archive-2x.zip 说明 此处版本号8.1.0.1为示例，具体以实际环境的版本号为准。
spark.yarn.am.extraJavaOptions	在Client模式下传递至YARN Application Master的一系列额外JVM选项。在Cluster模式下使用“spark.driver.extraJavaOptions”。	-Dlog4j.configuration=../__spark_conf__/_hadoop_conf__/log4j-executor.properties -Djava.security.auth.login.config=../__spark_conf__/_hadoop_conf__/jaas-zk.conf - Dzookeeper.server.principal=zookeeper/hadoop.<系统域名> - Djava.security.krb5.conf=../__spark_conf__/_hadoop_conf__/kdc.conf - Djdk.tls.ephemeralDHKeySize=2048
spark.shuffle.servicev2.port	Shuffle服务监测数据获取请求的端口。	27338
spark.ssl.historyServer.enabled	配置history server是否使用SSL。	true
spark.files.overwrite	当目标文件存在时，且其内容与源的文件不匹配。是否覆盖通过SparkContext.addFile()添加的文件。	false
spark.yarn.cluster.driver.extraClassPath	YARN-Cluster模式下，Driver使用的extraClassPath，配置为服务端的路径和参数。	\${BIGDATA_HOME}/common/runtime/security

配置项	说明	默认值或配置示例
spark.driver.extraClassPath	附加至driver的classpath的额外classpath条目。	\${BIGDATA_HOME}/common/runtime/security
spark.yarn.dist.innerfiles	配置YARN模式下Spark内部需要上传到HDFS的文件。	/Spark_path/spark/conf/s3p.file,/Spark_path/spark/conf/locals3.jceks <i>Spark_path</i> 为Spark客户端的安装路径。
spark.sql.bigdata.register.dialect	用于注册sql解析器。	org.apache.spark.sql.hbase.HBaseSQLParser
spark.shuffle.manager	处理数据的方式。有两种实现方式可用：sort和hash。sort shuffle对内存的使用率更高，是Spark 1.2及后续版本的默认选项。	SORT
spark.deploy.zookeeper.url	Zookeeper的地址，多个地址以逗号隔开。	For example: host1:2181,host2:2181,host3:2181
spark.broadcast.factory	使用的广播方式。	org.apache.spark.broadcast.TorrentBroadcastFactory
spark.sql.session.state.builder	指定会话状态构造器。	org.apache.spark.sql.hive.FIHiveACLSessionStateBuilder
spark.executor.extraLibraryPath	设置启动executor JVM时所使用的特殊的library path。	\${BIGDATA_HOME}/FusionInsight_HD_8.1.0.1/install/FusionInsight-Hadoop-*/hadoop/lib/native
spark.ui.customErrorMessagePage	配置网页有错误时是否允许显示自定义的错误信息页面。	true
spark.httpdProxy.enable	配置是否使用httpd代理。	true
spark.ssl.ui.enabledAlgorithms	配置ui ssl算法。	TLS_ECDHE_ECDSA_WITH_AES_256_GCM_SHA384,TLS_ECDHE_RSA_WITH_AES_256_GCM_SHA384,TLS_ECDHE_ECDSA_WITH_AES_128_GCM_SHA256,TLS_ECDHE_RSA_WITH_AES_128_GCM_SHA256,TLS_DHE_RSA_WITH_AES_256_GCM_SHA384,TLS_DHE_DSS_WITH_AES_256_GCM_SHA384,TLS_DHE_RSA_WITH_AES_128_GCM_SHA256,TLS_DHE_DSS_WITH_AES_128_GCM_SHA256

配置项	说明	默认值或配置示例
spark.ui.logout.enabled	针对Spark组件的WebUI，设置logout按钮。	true
spark.security.hideInfo.enabled	配置UI界面是否隐藏敏感信息。	true
spark.yarn.cluster.driver.extraLibraryPath	YARN-Cluster模式下driver的extraLibraryPath，配置成服务端的路径和参数。	\${BIGDATA_HOME}/ FusionInsight_HD_8.1.0.1/install/ FusionInsight-Hadoop-*/hadoop/lib/native
spark.driver.extraLibraryPath	设置一个特殊的library path在启动驱动程序JVM时使用。	\${DATA_NODE_INSTALL_HOME}/ hadoop/lib/native
spark.ui.killEnabled	允许停止Web UI中的stage和相应的job。	true
spark.yarn.access.hadoopFileSystems	Spark可以访问多个NameService。有多个NameService时，需要把所使用的NameService都配置进该配置项，之间以逗号分隔。	hdfs://hacluster,hdfs://hacluster

配置项	说明	默认值或配置示例
spark.yarn.cl uster.driver.e xtraJavaOpti ons	传递至Executor的额外JVM选项。例如，GC设置或其他日志记录。请注意不能通过此选项设置Spark属性或heap大小。Spark属性应该使用SparkConf对象或调用spark-submit脚本时指定的spark-defaults.conf文件来设置。Heap大小可以通过spark.executor.memory来设置。	-Xloggc:<LOG_DIR>/gc.log - XX:+PrintGCDetails -XX:- OmitStackTraceInFastThrow - XX:+PrintGCTimeStamps - XX:+PrintGCDateStamps - XX:+UseGCLogFileRotation - XX:NumberOfGCLogFiles=20 - XX:GCLogFileSize=10M - Dlog4j.configuration=./__spark_conf__/ __hadoop_conf__/log4j-executor.properties -Djava.security.auth.login.config=./ __spark_conf__/__hadoop_conf__/jaas- zk.conf - Dzookeeper.server.principal=zookeeper/ hadoop.<系统域名> - Djava.security.krb5.conf=./__spark_conf__/ __hadoop_conf__/kdc.conf - Djetty.version=x.y.z - Dorg.xerial.snappy.tmpdir=\$ {BIGDATA_HOME}/tmp/spark2x_app - Dcarbon.properties.filepath=./ __spark_conf__/__hadoop_conf__/ carbon.properties - Djdk.tls.ephemeralDHKeySize=2048
spark.driver.e xtraJavaOpti ons	传递至driver（驱动程序）的一系列额外JVM选项。	-Xloggc:\${SPARK_LOG_DIR}/indexserver- omm-%p-gc.log -XX:+PrintGCDetails -XX:- OmitStackTraceInFastThrow - XX:+PrintGCTimeStamps - XX:+PrintGCDateStamps - XX:MaxDirectMemorySize=512M - XX:MaxMetaspaceSize=512M - XX:+UseGCLogFileRotation - XX:NumberOfGCLogFiles=20 - XX:GCLogFileSize=10M - XX:OnOutOfMemoryError='kill -9 %p' - Djetty.version=x.y.z - Dorg.xerial.snappy.tmpdir=\$ {BIGDATA_HOME}/tmp/spark2x/ JDBCServer/snappy_tmp -Djava.io.tmpdir= \${BIGDATA_HOME}/tmp/spark2x/ JDBCServer/io_tmp - Dcarbon.properties.filepath=\$ {SPARK_CONF_DIR}/carbon.properties - Djdk.tls.ephemeralDHKeySize=2048 - Dspark.ssl.keyStore=\${SPARK_CONF_DIR}/ child.keystore #{java_stack_prefer}
spark.eventL og.overwrite	是否覆盖任何现有的文件。	false

配置项	说明	默认值或配置示例
spark.eventLog.dir	如果 spark.eventLog.enabled 为 true ，记录 Spark 事件的目录。在此目录下，Spark 为每个应用程序创建文件，并将应用程序的事件记录到文件中。用户也可设置为统一的与 HDFS 目录相似的地址，这样 History server 就可以读取历史文件。	hdfs://hacluster/spark2xJobHistory2x
spark.random.port.min	设置随机端口的最小值。	22600
spark.authenticate	是否 Spark 认证其内部连接。如果不是运行在 YARN 上，请参见 spark.authenticate.secret 的相关内容。	true
spark.random.port.max	设置随机端口的最大值。	22899
spark.eventLog.enabled	是否记录 Spark 事件，用于应用程序在完成后重构 webUI。	true

配置项	说明	默认值或配置示例
spark.executor.extraJavaOptions	传递至Executor的额外JVM选项。例如，GC设置或其他日志记录。请注意不能通过此选项设置Spark属性或heap大小。	<pre>-Xloggc:<LOG_DIR>/gc.log - XX:+PrintGCDetails -XX:- OmitStackTracelnFastThrow - XX:+PrintGCTimeStamps - XX:+PrintGCDateStamps - XX:+UseGCLogFileRotation - XX:NumberOfGCLogFiles=20 - XX:GCLogFileSize=10M - Dlog4j.configuration=./log4j- executor.properties - Djava.security.auth.login.config=./jaas- zk.conf - Dzookeeper.server.principal=zookeeper/ hadoop.<系统域名> - Djava.security.krb5.conf=./kdc.conf - Dcarbon.properties.filepath=./ carbon.properties -Xloggc:<LOG_DIR>/gc.log - XX:+PrintGCDetails -XX:- OmitStackTracelnFastThrow - XX:+PrintGCTimeStamps - XX:+PrintGCDateStamps - XX:+UseGCLogFileRotation - XX:NumberOfGCLogFiles=20 - XX:GCLogFileSize=10M - Dlog4j.configuration=./__spark_conf__/_ _hadoop_conf__/log4j-executor.properties -Djava.security.auth.login.config=./ __spark_conf__/_hadoop_conf__/jaas- zk.conf - Dzookeeper.server.principal=zookeeper/ hadoop.<系统域名> - Djava.security.krb5.conf=./__spark_conf__/_ _hadoop_conf__/kdc.conf - Dcarbon.properties.filepath=./ __spark_conf__/_hadoop_conf__/_ carbon.properties - Djdk.tls.ephemeralDHKeySize=2048</pre>
spark.sql.authorization.enabled	配置Hive client是否开启认证。	true

26.11.2 常用参数

概述

本节介绍Spark使用过程中的常用配置项。以特性为基础划分子章节，以使用户快速搜索到相应的配置项。如果用户使用MRS集群，本节介绍的参数大部分已经适配好，用户无需再进行配置。少数需要用户根据实际场景配置的参数，请参见[快速配置参数](#)。

配置 Stage 失败重试次数

Spark任务在遇到FetchFailedException时会触发Stage重试。为了防止Stage无限重试，对Stage重试次数进行限制。重试次数可以根据实际需要进行调整。

在Spark客户端的“spark-defaults.conf”文件中配置如下参数。

表 26-36 参数说明

参数	说明	默认值
spark.stage.maxConsecutiveAttempts	Stage失败重试最大次数。	4

配置是否使用笛卡尔积功能

要启动使用笛卡尔积功能，需要在Spark的“spark-defaults.conf”配置文件中进行如下设置。

表 26-37 笛卡尔积参数说明

参数	说明	默认值
spark.sql.crossJoin.enabled	是否允许隐性执行笛卡尔积。 <ul style="list-style-type: none">“true”表示允许“false”表示不允许，此时只允许query中显式包含CROSS JOIN语法。	true

说明

- JDBC应用在服务端的“spark-defaults.conf”配置文件中设置该参数。
- Spark客户端提交的任务在客户端配的“spark-defaults.conf”配置文件中设置该参数。

Spark 长时间任务安全认证配置

安全模式下，使用Spark CLI（如spark shell、spark sql、spark submit）时，如果使用kinit命令进行安全认证，当执行长时间运行任务时，会因为认证过期导致任务失败。

在客户端的“spark-defaults.conf”配置文件中设置如下参数，配置完成后，重新执行 Spark CLI即可。

说明

当参数值为“true”时，需要保证“spark-defaults.conf”和“hive-site.xml”中的Keytab和principal的值相同。

表 26-38 参数说明

参数名称	含义	默认值
spark.kerberos.principal	具有Spark操作权限的principal。请联系MRS集群管理员获取对应principal。	-
spark.kerberos.keytab	具有Spark操作权限的Keytab文件名称和文件路径。请联系MRS集群管理员获取对应Keytab文件。	-
spark.security.bigdata.loginOnce	Principal用户是否只登录一次。true为单次登录；false为多次登录。 单次登录与多次登录的区别在于：Spark社区使用多次Kerberos用户登录多次的方案，但容易出现TGT过期或者Token过期异常导致应用无法长时间运行。DataSight修改了Kerberos登录方式，只允许用户登录一次，可以有效的解决过期问题。限制在于，Hive相关的principal与keytab的配置项必须与Spark配置相同。 说明 当参数值为true时，需要保证“spark-defaults.conf”和“hive-site.xml”中的Keytab和principal的值相同。	true

Python Spark

Python Spark是Spark除了Scala、Java两种API之外的第三种编程语言。不同于Java和Scala都是在JVM平台上运行，Python Spark不仅会有JVM进程，还会有自身的Python进程。以下配置项只适用于Python Spark场景，而其他配置项也同样可以在Python Spark中生效。

表 26-39 参数说明

参数	描述	默认值
spark.python.profile	在Python worker中开启profiling。通过sc.show_profiles()展示分析结果。或者在driver退出前展示分析结果。可以通过sc.dump_profiles(path)将结果转储到磁盘中。如果一些分析结果已经手动展示，那么在Driver退出前，它们将不会再自动展示。 默认使用pyspark.profiler.BasicProfiler，可以在初始化SparkContext时传入指定的profiler来覆盖默认的profiler。	false

参数	描述	默认值
spark.python.worker.memory	聚合过程中每个python worker进程所能使用的内存大小，其值格式同指定JVM内存一致，如 512m，2g。如果进程在聚集期间所用的内存超过了该值，数据将会被写入磁盘。	512m
spark.python.worker.reuse	是否重用python worker。如是，它将使用固定数量的Python workers，那么下一批提交的task将重用这些Python workers，而不是为每个task重新fork一个Python进程。该功能在大型广播下非常有用，因为此时对下一批提交的task不需要将数据从JVM再一次传输至Python worker。	true

Dynamic Allocation

动态资源调度是On Yarn模式特有的特性，并且必须开启Yarn External Shuffle才能使用这个功能。在使用Spark作为一个常驻的服务时候，动态资源调度将大大的提高资源的利用率。例如JDBCServer服务，大多数时间该进程并不接受JDBC请求，因此将这段空闲时间的资源释放出来，将极大的节约集群的资源。

表 26-40 参数说明

参数	描述	默认值
spark.dynamicAllocation.enabled	是否使用动态资源调度，用于根据规模调整注册于该应用的executor的数量。注意目前仅在YARN模式下有效。 启用动态资源调度必须将spark.shuffle.service.enabled设置为true。以下配置也与此相关： spark.dynamicAllocation.minExecutors、 spark.dynamicAllocation.maxExecutors和 spark.dynamicAllocation.initialExecutors。	<ul style="list-style-type: none"> JDBCServer2x: true SparkResource2x: false
spark.dynamicAllocation.minExecutors	最小Executor个数。	0
spark.dynamicAllocation.initialExecutors	初始Executor个数。	spark.dynamicAllocation.minExecutors
spark.dynamicAllocation.maxExecutors	最大executor个数。	2048
spark.dynamicAllocation.schedulerBacklogTimeout	调度第一次超时时间。单位为秒。	1s

参数	描述	默认值
spark.dynamicAllocation.sustainedSchedulerBacklogTimeout	调度第二次及之后超时时间。	1s
spark.dynamicAllocation.executorIdleTimeout	普通Executor空闲超时时间。单位为秒。	60
spark.dynamicAllocation.cachedExecutorIdleTimeout	含有cached blocks的Executor空闲超时时间。	<ul style="list-style-type: none"> JDBCServer2x: 2147483647s IndexServer2x: 2147483647s SparkResource2x: 120

Spark Streaming

Spark Streaming是在Spark批处理平台提供的流式数据的处理能力，以“mini-batch”的方式处理从外部输入的数据。

在Spark客户端的“spark-defaults.conf”文件中配置如下参数。

表 26-41 参数说明

参数	描述	默认值
spark.streaming.receiver.writeAheadLog.enable	启用预写日志（WAL）功能。所有通过Receiver接收的输入数据将被保存至预写日志，预写日志可以保证Driver程序出错后数据可以恢复。	false
spark.streaming.unpersist	由Spark Streaming产生和保存的RDDs自动从Spark的内存中强制移除。Spark Streaming接收的原始输入数据也将自动清除。设置为false时原始输入数据和存留的RDDs不会自动清除，因此在streaming应用外部依然可以访问，但是这会占用更多的Spark内存。	true

Spark Streaming Kafka

Receiver是Spark Streaming一个重要的组成部分，它负责接收外部数据，并将数据封装为Block，提供给Streaming消费。最常见的数据源是Kafka，Spark Streaming对

Kafka的集成也是最完善的，不仅有可靠性的保障，而且也支持从Kafka直接作为RDD输入。

表 26-42 参数说明

参数	描述	默认值
spark.streaming.kafka.maxRatePerPartition	使用Kafka direct stream API时，从每个Kafka分区读取数据的最大速率（每秒记录数量）。	-
spark.streaming.blockInterval	在被存入Spark之前Spark Streaming Receiver接收数据累积成数据块的间隔（毫秒）。推荐最小值为50毫秒。	200ms
spark.streaming.receiver.maxRate	每个Receiver接收数据的最大速率（每秒记录数量）。配置设置为0或者负值将不会对速率设限。	-
spark.streaming.receiver.writeAheadLog.enabled	是否使用ReliableKafkaReceiver。该Receiver支持流式数据不丢失。	false

Netty/NIO 及 Hash/Sort 配置

Shuffle是大数据处理中最重要的一个性能点，网络是整个Shuffle过程的性能点。目前Spark支持两种Shuffle方式，一种是Hash，另外一种Sort。网络也有两种方式，Netty和NIO。

表 26-43 参数说明

参数	描述	默认值
spark.shuffle.manager	处理数据的方式。有两种实现方式可用：sort和hash。sort shuffle对内存的使用率更高，是Spark 1.2及后续版本的默认选项。	SORT
spark.shuffle.consolidateFiles	（仅hash方式）若要合并shuffle过程中创建的中间文件，需要将该值设置为“true”。文件创建的少可以提高文件系统处理性能，降低风险。使用ext4或者xfs文件系统时，建议设置为“true”。由于文件系统限制，在ext3上该设置可能会降低8核以上机器的处理性能。	false
spark.shuffle.sort.byPassMergeThreshold	该参数只适用于spark.shuffle.manager设置为sort时。在不做map端聚合并且reduce任务的partition数小于或等于该值时，避免对数据进行归并排序，防止系统处理不必要的排序引起性能下降。	200

参数	描述	默认值
spark.shuffle.io.maxRetries	（仅Netty方式）如果设为非零值，由于IO相关的异常导致的fetch失败会自动重试。该重试逻辑有助于大型shuffle在发生GC暂停或者网络闪断时保持稳定。	12
spark.shuffle.io.numConnectionsPerPeer	（仅Netty方式）为了减少大型集群的连接创建，主机间的连接会被重新使用。对于拥有较多硬盘和少数主机的集群，此操作可能会导致并发性不足以占用所有磁盘，所以用户可以考虑增加此值。	1
spark.shuffle.io.preferDirectBufs	（仅Netty方式）使用off-heap缓冲区减少shuffle和高速缓存块转移期间的垃圾回收。对于off-heap内存被严格限制的环境，用户可以将其关闭以强制所有来自Netty的申请使用堆内内存。	true
spark.shuffle.io.retryWait	（仅Netty方式）等待fetch重试期间的的时间（秒）。重试引起的最大延迟为maxRetries * retryWait，默认是15秒。	5

普通 Shuffle 配置

表 26-44 参数说明

参数	描述	默认值
spark.shuffle.spill	若设为“true”，通过将数据溢出至磁盘来限制reduce任务期间内存的使用量。	true
spark.shuffle.spill.compress	是否压缩shuffle期间溢出的数据。使用spark.io.compression.codec指定的算法进行数据压缩。	true
spark.shuffle.file.buffer	每个shuffle文件输出流的内存缓冲区大小（单位：KB）。这些缓冲区可以减少创建中间shuffle文件流过程中产生的磁盘寻道和系统调用次数。也可以通过配置项spark.shuffle.file.buffer.kb设置。	32KB
spark.shuffle.compress	是否压缩map任务输出文件。建议压缩。使用spark.io.compression.codec进行压缩。	true
spark.reducer.maxSizeInFlight	从每个reduce任务同时fetch的map任务输出最大值（单位：MB）。由于每个输出要求创建一个缓冲区进行接收，这代表了每个reduce任务固定的内存开销，所以除非拥有大量内存，否则保持低值。也可以通过配置项spark.reducer.maxMblnFlight设置。	48MB

Driver 配置

Spark Driver 可以理解为 Spark 提交应用的客户端，所有的代码解析工作都在这个进程中完成，因此该进程的参数尤其重要。下面将以如下顺序介绍 Spark 中进程的参数设置：

- JavaOptions: Java 命令中“-D”后面的参数，可以由 System.getProperty 获取。
- ClassPath: 包括 Java 类和 Native 的 Lib 加载路径。
- Java Memory and Cores: Java 进程的内存和 CPU 使用量。
- Spark Configuration: Spark 内部参数，与 Java 进程无关。

表 26-45 参数说明

参数	描述	默认值
spark.driver.extraJavaOptions	传递至 driver（驱动程序）的一系列额外 JVM 选项。例如，GC 设置或其他日志记录。 注意：在 Client 模式中，该配置禁止直接在应用程序中通过 SparkConf 设置，因为驱动程序 JVM 已经启动。请通过 --driver-java-options 命令行选项或默认 property 文件进行设置。	参考 快速配置参数
spark.driver.extraClassPath	附加至 driver 的 classpath 的额外 classpath 条目。 注意：在 Client 模式中，该配置禁止直接在应用程序中通过 SparkConf 设置，因为驱动程序 JVM 已经启动。请通过 --driver-java-options 命令行选项或默认 property 文件进行设置。	参考 快速配置参数
spark.driver.userClassPathFirst	（试验性）当在驱动程序中加载类时，是否授权用户添加的 jar 优先于 Spark 自身的 jar。这种特性可用于减缓 Spark 依赖和用户依赖之间的冲突。目前该特性仍处于试验阶段，仅用于 Cluster 模式中。	false
spark.driver.extraLibraryPath	设置一个特殊的 library path 在启动驱动程序 JVM 时使用。 注意：在 Client 模式中，该配置禁止直接在应用程序中通过 SparkConf 设置，因为驱动程序 JVM 已经启动。请通过 --driver-java-options 命令行选项或默认 property 文件进行设置。	<ul style="list-style-type: none"> • JDBCServer2x: \$ {SPARK_INSTALLED_HOME}/spark/native • SparkResource2x: \$ {DATA_NODE_INSTANCE_HOME}/hadoop/lib/native

参数	描述	默认值
spark.driver.cores	驱动程序进程使用的核数。仅适用于Cluster模式。	1
spark.driver.memory	驱动程序进程使用的内存数量，即SparkContext初始化的进程（例如：512M, 2G）。 注意：在Client模式中，该配置禁止直接在应用程序中通过SparkConf设置，因为驱动程序JVM已经启动。请通过--driver-java-options命令行选项或默认property文件进行设置。	4G
spark.driver.maxResultSize	对每个Spark action操作（例如“collect”）的所有分区序列化结果的总量限制，至少1M，设置成0表示不限制。如果总量超过该限制，工作任务会中止。限制值设置过高可能会引起驱动程序的内存不足错误（取决于spark.driver.memory和JVM的对象内存开销）。设置合理的限制可以避免驱动程序出现内存不足的错误。	1G
spark.driver.host	Driver监测的主机名或IP地址，用于Driver与Executor进行通信。	(local hostname)
spark.driver.port	Driver监测的端口，用于Driver与Executor进行通信。	(random)

ExecutorLauncher 配置

ExecutorLauncher只有在Yarn-Client模式下才会存在的角色，Yarn-Client模式下，ExecutorLauncher和Driver不在同一个进程中，需要对ExecutorLauncher的参数进行特殊的配置。

表 26-46 参数说明

参数	描述	默认值
spark.yarn.am.extraJavaOptions	在Client模式下传递至YARN Application Master的一系列额外JVM选项。在Cluster模式下使用spark.driver.extraJavaOptions。	参考 快速配置参数
spark.yarn.am.memory	针对Client模式下YARN Application Master使用的内存数量，与JVM内存设置字符串格式一致（例如：512m, 2g）。在集群模式下，使用spark.driver.memory。	1G
spark.yarn.am.memoryOverhead	和“spark.yarn.driver.memoryOverhead”一样，但只针对Client模式下的Application Master。	-

参数	描述	默认值
spark.yarn.am.cores	针对Client模式下YARN Application Master使用的核数。在Cluster模式下，使用spark.driver.cores。	1

Executor 配置

Executor也是单独一个Java进程，但不像Driver和AM只有一个，Executor可以有多个进程，而目前Spark只支持相同的配置，即所有Executor的进程参数都必然是一样的。

表 26-47 参数说明

参数	描述	默认值
spark.executor.extraJavaOptions	传递至Executor的额外JVM选项。例如，GC设置或其他日志记录。请注意不能通过此选项设置Spark属性或heap大小。Spark属性应该使用SparkConf对象或调用spark-submit脚本时指定的spark-defaults.conf文件来设置。Heap大小可以通过spark.executor.memory来设置。	参考 快速配置参数
spark.executor.extraClassPath	附加至Executor classpath的额外的classpath。这主要是为了向后兼容Spark的历史版本。用户一般不用设置此选项。	-
spark.executor.extraLibraryPath	设置启动executor JVM时所使用的特殊的library path。	参考 快速配置参数
spark.executor.userClassPathFirst	（试验性）与spark.driver.userClassPathFirst相同的功能，但应用于Executor实例。	false
spark.executor.memory	每个Executor进程使用的内存数量，与JVM内存设置字符串的格式相同（例如：512M，2G）。	4G
spark.executorEnv.[EnvironmentVariableName]	添加由EnvironmentVariableName指定的环境变量至executor进程。用户可以指定多个来设置多个环境变量。	-
spark.executor.logs.rolling.maxRetainedFiles	设置系统即将保留的最新滚动日志文件的数量。旧的日志文件将被删除。默认关闭。	-
spark.executor.logs.rolling.size.maxBytes	设置滚动Executor日志的文件的最大值。默认关闭。数值以字节为单位设置。若要自动清除旧日志，请查看spark.executor.logs.rolling.maxRetainedFiles。	-

参数	描述	默认值
spark.executor.logs.rolling.strategy	设置executor日志的滚动策略。默认滚动关闭。可以设置为“time”（基于时间的滚动）或“size”（基于大小的滚动）。当设置为“time”，使用spark.executor.logs.rolling.time.interval属性的值作为日志滚动的间隔。当设置为“size”，使用spark.executor.logs.rolling.size.maxBytes设置滚动的最大文件大小滚动。	-
spark.executor.logs.rolling.time.interval	设置executor日志滚动的时间间隔。默认关闭。合法值为“daily”、“hourly”、“minutely”或任意秒。若要自动清除旧日志，请查看spark.executor.logs.rolling.maxRetainedFiles。	daily

WebUI

WebUI展示了Spark应用运行的过程和状态。

表 26-48 参数说明

参数	描述	默认值
spark.ui.killEnabled	允许停止Web UI中的stage和相应的job。 说明 出于安全考虑，将此配置项的默认值设置成false，以避免用户发生误操作。如果需要开启此功能，则可以在spark-defaults.conf配置文件中将此配置项的值设为true。请谨慎操作。	true
spark.ui.port	应用程序dashboard的端口，显示内存和工作量数据。	<ul style="list-style-type: none"> JDBC Server2x: 4040 Spark Resource2x: 0 Index Server2x: 22901
spark.ui.retainedJobs	在垃圾回收之前Spark UI和状态API记住的job数。	1000
spark.ui.retainedStages	在垃圾回收之前Spark UI和状态API记住的stage数。	1000

HistoryServer

HistoryServer读取文件系统中的EventLog文件，展示已经运行完成的Spark应用在运行时的状态信息。

表 26-49 参数说明

参数	描述	默认值
spark.history.fs.logDirectory	History server的日志目录	-
spark.history.ui.port	JobHistory侦听连接的端口。	18080
spark.history.fs.updateInterval	History server所显示信息的更新周期，单位为秒。每次更新检查持久存储中针对事件日志进行的更改。	10s
spark.history.fs.updateInterval.seconds	每个事件日志更新检查的间隔。与spark.history.fs.updateInterval功能相同，推荐使用spark.history.fs.updateInterval。	10s
spark.history.updateInterval	该配置项与spark.history.fs.updateInterval.seconds和spark.history.fs.updateInterval功能相同，推荐使用spark.history.fs.updateInterval。	10s

HistoryServer UI 超时和最大访问数

表 26-50 参数说明

参数	描述	默认值
spark.session.maxAge	设置会话的超时时间，单位秒。此参数只适用于安全模式。普通模式下，无法设置此参数。	600
spark.connection.maxRequest	设置客户端访问Jobhistory的最大并发数量。	5000

EventLog

Spark应用在运行过程中，实时将运行状态以JSON格式写入文件系统，用于HistoryServer服务读取并重现应用运行时状态。

表 26-51 参数说明

参数	描述	默认值
spark.eventLog.enabled	是否记录Spark事件，用于应用程序在完成后重构webUI。	true

参数	描述	默认值
spark.eventLog.dir	如果spark.eventLog.enabled为true，记录Spark事件的目录。在此目录下，Spark为每个应用程序创建文件，并将应用程序的事件记录到文件中。用户也可设置为统一的与HDFS目录相似的地址，这样History server就可以读取历史文件。	hdfs://hacluster/spark2x/jobHistory2x
spark.eventLog.compress	spark.eventLog.enabled为true时，是否压缩记录的事件。	false

EventLog 的周期清理

JobHistory上的Event log是随每次任务的提交而累积的，任务提交的次数多了之后会造成太多文件的存放。Spark提供了周期清理Event log的功能，用户可以通过配置开关和相应的清理周期参数来进行控制。

表 26-52 参数说明

参数	描述	默认值
spark.history.fs.cleaner.enabled	是否打开清理功能。	true
spark.history.fs.cleaner.interval	清理功能的检查周期。	1d
spark.history.fs.cleaner.maxAge	日志的最长保留时间。	4d

Kryo

Kryo是一个非常高效的Java序列化框架，Spark中也默认集成了该框架。几乎所有的Spark性能调优都离不开将Spark默认的序列化器转化为Kryo序列化器的过程。目前Kryo序列化只支持Spark数据层面的序列化，还不支持闭包的序列化。设置Kryo序列化元，需要将配置项“spark.serializer”设置为“org.apache.spark.serializer.KryoSerializer”，同时也搭配设置以下的配置项，优化Kryo序列化的性能。

表 26-53 参数说明

参数	描述	默认值
spark.kryo.classesToRegister	使用Kryo序列化时，需要注册到Kryo的类名，多个类之间用逗号分隔。	-

参数	描述	默认值
spark.kryo.referenceTracking	当使用Kryo序列化数据时，是否跟踪对同一个对象的引用情况。适用于对象图有循环引用或同一对象有多个副本的情况。否则可以设置为关闭以提升性能。	true
spark.kryo.registrationRequired	是否需要使用Kryo来注册对象。当设为“true”时，如果序列化一个未使用Kryo注册的对象则会发生异常。当设为“false”（默认值）时，Kryo会将未注册的类名称一同写到序列化对象中。该操作会带来大量性能开销，所以在用户还没有从注册队列中删除相应的类时应该开启该选项。	false
spark.kryo.registrator	如果使用Kryo序列化，使用Kryo将该类注册至定制类。如果需要以定制方式注册类，例如指定一个自定义字段序列化器，可使用该属性。否则spark.kryo.classesToRegister会更简单。它应该设置为一个扩展KryoRegistrator的类。	-
spark.kryo.serializer.buffer.max	Kryo序列化缓冲区允许的最大值，单位为兆字节。这个值必须大于尝试序列化的对象。当在Kryo中遇到“buffer limit exceeded”异常时可以适当增大该值。也可以通过配置项spark.kryo.serializer.buffer.max配置。	64MB
spark.kryo.serializer.buffer	Kryo序列化缓冲区的初始值，单位为兆字节。每个worker的每个核心都会有一个缓冲区。如果有需要，缓冲区会增大到spark.kryo.serializer.buffer.max设置的值。也可以通过配置项spark.kryo.serializer.buffer配置。	64KB

Broadcast

Broadcast用于Spark进程间数据块的传输。Spark中无论Jar包、文件还是闭包以及返回的结果都会使用Broadcast。目前的Broadcast支持两种方式，Torrent与HTTP。前者将会把数据切成小片，分布到集群中，有需要时从远程获取；后者将文件存入到本地磁盘，有需要时通过HTTP方式将整个文件传输到远端。前者稳定性优于后者，因此Torrent为默认的Broadcast方式。

表 26-54 参数说明

参数	描述	默认值
spark.broadcast.factory	使用的广播方式。	org.apache.spark.broadcast.TorrentBroadcastFactory
spark.broadcast.blockSize	TorrentBroadcastFactory的块大小。该值过大会降低广播时的并行度（速度变慢），过小可能会影响BlockManager的性能。	4096

参数	描述	默认值
spark.broadcast.compress	在发送广播变量之前是否压缩。建议压缩。	true

Storage

内存计算是Spark的最大亮点，Spark的Storage主要管理内存资源。Storage中主要存储RDD在Cache过程中产生的数据块。JVM中堆内存是整体的，因此在Spark的Storage管理中，“Storage Memory Size”变成了一个非常重要的概念。

表 26-55 参数说明

参数	描述	默认值
spark.storage.memoryMapThreshold	超过该块大小的Block，Spark会对该磁盘文件进行内存映射。这可以防止Spark在内存映射时映射过小的块。一般情况下，对接近或低于操作系统的页大小的块进行内存映射会有高开销。	2m

PORT

表 26-56 参数说明

参数	描述	默认值
spark.ui.port	应用仪表盘的端口，显示内存和工作负载数据。	<ul style="list-style-type: none">JDBC Server2x : 4040SparkResource2x : 0
spark.blockManager.port	所有BlockManager监测的端口。这些同时存在于Driver和Executor上。	随机端口范围
spark.driver.port	Driver监测的端口，用于Driver与Executor进行通信。	随机端口范围

随机端口范围

所有随机端口必须在一定端口范围内。

表 26-57 参数说明

参数	描述	默认值
spark.random.port.min	设置随机端口的最小值。	22600
spark.random.port.max	设置随机端口的最大值。	22899

TIMEOUT

Spark默认配置能很好的处理中等数据规模的计算任务，但一旦数据量过大，会经常出现超时导致任务失败的场景。在大数据量场景下，需调大Spark中的超时参数。

表 26-58 参数说明

参数	描述	默认值
spark.files.fetchTimeout	获取通过驱动程序的SparkContext.addFile()添加的文件时的通信超时（秒）。	60s
spark.network.timeout	所有网络交互的默认超时（秒）。如未配置，则使用该配置代替 spark.core.connection.ack.wait.timeout, spark.akka.timeout, spark.storage.blockManagerSlaveTimeoutMs或 spark.shuffle.io.connectionTimeout。	360s
spark.core.connection.ack.wait.timeout	连接时应答的超时时间（单位：秒）。为了避免由于GC带来的长时间等待，可以设置更大的值。	60

加密

Spark支持Akka和HTTP（广播和文件服务器）协议的SSL，但WebUI和块转移服务仍不支持SSL。

SSL必须在每个节点上配置，并使用特殊协议为通信涉及到的每个组件进行配置。

表 26-59 参数说明

参数	描述	默认值
spark.ssl.enabled	是否在所有被支持协议上开启SSL连接。 与spark.ssl.xxx类似的所有SSL设置指示了所有被支持协议的全局配置。为了覆盖特殊协议的全局配置，在协议指定的命名空间中必须重写属性。 使用“spark.ssl.YYY.XXX”设置覆盖由YYY指示的特殊协议的全局配置。目前YYY可以是基于Akka连接的akka或广播与文件服务器的fs。	false

参数	描述	默认值
spark.ssl.enabledAlgorithms	以逗号分隔的密码列表。指定的密码必须被JVM支持。	-
spark.ssl.keyPassword	key-store的私人密钥密码。	-
spark.ssl.keyStore	key-store文件的路径。该路径可以绝对或相对于开启组件的目录。	-
spark.ssl.keyStorePassword	key-store的密码。	-
spark.ssl.protocol	协议名。该协议必须被JVM支持。本页所有协议的参考表。	-
spark.ssl.trustStore	trust-store文件的路径。该路径可以绝对或相对于开启组件的目录。	-
spark.ssl.trustStorePassword	trust-store的密码。	-

安全性

Spark目前支持通过共享密钥认证。可以通过spark.authenticate配置参数配置认证。该参数控制Spark通信协议是否使用共享密钥执行认证。该认证是确保双边都有相同的共享密钥并被允许通信的基本握手。如果共享密钥不同，通信将不被允许。共享密钥通过如下方式创建：

- 对于YARN部署的Spark，将spark.authenticate配置为真会自动处理生成和分发共享密钥。每个应用程序会独占一个共享密钥。
- 对于其他类型部署的Spark，应该在每个节点上配置Spark参数spark.authenticate.secret。所有Master/Workers和应用程序都将使用该密钥。

表 26-60 参数说明

参数	描述	默认值
spark.acls.enable	是否开启Spark acls。如果开启，它将检查用户是否有访问和修改job的权限。请注意这要求用户可以被识别。如果用户被识别为无效，检查将不被执行。UI可以使用过滤器认证和设置用户。	true
spark.admin.acls	逗号分隔的有权限访问和修改所有Spark job的用户/管理员列表。如果在共享集群上运行并且工作时有MRS集群管理员或开发人员帮助调试，可以使用该列表。	admin
spark.authenticate	是否Spark认证其内部连接。如果不是运行在YARN上，请参见spark.authenticate.secret。	true

参数	描述	默认值
spark.authenticate.secret	设置Spark各组件之间验证的密钥。如果不是运行在YARN上且认证未开启，需要设置该项。	-
spark.modify.acls	逗号分隔的有权限修改Spark job的用户列表。默认情况下只有开启Spark job的用户才有修改列表的权限（例如删除列表）。	-
spark.ui.view.acls	逗号分隔的有权限访问Spark web ui的用户列表。默认情况下只有开启Spark job的用户才有访问权限。	-

开启 Spark 进程间的认证机制

目前Spark进程间支持共享密钥方式的认证机制，通过配置spark.authenticate可以控制Spark在通信过程中是否做认证。这种认证方式只是通过简单的握手来确定通信双方享有共同的密钥。

在Spark客户端的“spark-defaults.conf”文件中配置如下参数。

表 26-61 参数说明

参数	描述	默认值
spark.authenticate	在Spark on YARN模式下，将该参数配置成true即可。密钥的生成和分发过程是自动完成的，并且每个应用独占一个密钥。	true

Compression

数据压缩是一个以CPU换内存的优化策略，因此当Spark内存严重不足的时候（由于内存计算的特质，这种情况非常常见），使用压缩可以大幅提高性能。目前Spark支持三种压缩算法：snappy, lz4, lzf。Snappy为默认压缩算法，并且调用native方法进行压缩与解压缩，在Yarn模式下需要注意堆外内存对Container进程的影响。

表 26-62 参数说明

参数	描述	默认值
spark.io.compression.codec	用于压缩内部数据的codec，例如RDD分区、广播变量和shuffle输出。默认情况下，Spark支持三种压缩算法：lz4, lzf和snappy。可以使用完全合格的类名称指定算法，例如org.apache.spark.io.LZ4CompressionCodec、org.apache.spark.io.LZFCompressionCodec及org.apache.spark.io.SnappyCompressionCodec。	lz4
spark.io.compression.lz4.block.size	当使用LZ4压缩算法时LZ4压缩中使用的块大小（字节）。当使用LZ4时降低块大小同样也会降低shuffle内存使用。	32768

参数	描述	默认值
spark.io.compression.snappy.block.size	当使用Snappy压缩算法时Snappy压缩中使用的块大小（字节）。当使用Snappy时降低块大小同样也会降低shuffle内存使用。	32768
spark.shuffle.compress	是否压缩map任务输出文件。建议压缩。使用spark.io.compression.codec进行压缩。	true
spark.shuffle.spill.compress	是否压缩在shuffle期间溢出的数据。使用spark.io.compression.codec进行压缩。	true
spark.eventLog.compress	设置当spark.eventLog.enabled设置为true时是否压缩记录的事件。	false
spark.broadcast.compress	在发送之前是否压缩广播变量。建议压缩。	true
spark.rdd.compress	是否压缩序列化的RDD分区（例如StorageLevel.MEMORY_ONLY_SER的分区）。牺牲部分额外CPU的时间可以节省大量空间。	false

在资源不足的情况下，降低客户端运行异常概率

在资源不足的情况下，Application Master会因等待资源出现超时，导致任务被删除。调整如下参数，降低客户端应用运行异常概率。

在客户端的“spark-defaults.conf”配置文件中调整如下参数。

表 26-63 参数说明

参数	说明	默认值
spark.yarn.applicationMaster.waitTries	设置Application Master等待Spark master的次数，同时也是等待SparkContext初始化的次数。增大该参数值，可以防止AM任务被删除，降低客户端应用运行异常的概率。	10
spark.yarn.am.memory	调整AM的内存。增大该参数值，可以防止AM因内存不足而被RM删除任务，降低客户端应用运行异常的概率。	1G

26.11.3 Spark2x 日志介绍

日志描述

日志存储路径：

- Executor运行日志：“\${BIGDATA_DATA_HOME}/hadoop/data\${i}/nm/containerlogs/application_\${appid}/container_{\$contid}”

 说明

运行中的任务日志存储在以上路径中，运行结束后会基于Yarn的配置确定是否汇聚到HDFS目录中，详情请参见[Yarn常用配置参数](#)。

- 其他日志：“/var/log/Bigdata/spark2x”

日志归档规则：

- 使用yarn-client或yarn-cluster模式提交任务时，Executor日志默认50MB滚动存储一次，最多保留10个文件，不压缩。
- JobHistory2x日志默认100MB滚动存储一次，最多保留100个文件，压缩存储。
- JDBCServer2x日志默认100MB滚动存储一次，最多保留100个文件，压缩存储。
- IndexServer2x日志默认100MB滚动存储一次，最多保留100个文件，压缩存储。
- JDBCServer2x审计日志默认20MB滚动存储一次，最多保留20个文件，压缩存储。
- 日志大小和压缩文件保留个数可以在FusionInsight Manager界面中配置。

表 26-64 Spark2x 日志列表

日志类型	日志文件名	描述
SparkResource2x 日志	spark.log	Spark2x服务初始化日志。
	prestart.log	prestart脚本日志。
	cleanup.log	安装卸载实例时的清理日志。
	spark-availability-check.log	Spark2x服务健康检查日志。
	spark-service-check.log	Spark2x服务检查日志
JDBCServer2x日志	JDBCServer-start.log	JDBCServer2x启动日志。
	JDBCServer-stop.log	JDBCServer2x停止日志。
	JDBCServer.log	JDBCServer2x运行时，Driver端日志。
	jdbc-state-check.log	JDBCServer2x健康检查日志。
	jdbcserver-omm-pid***-gc.log.*.current	JDBCServer2x进程gc日志。
	spark-omm-org.apache.spark.sql.hive.thriftserver.HiveThriftProxyServer2-***.out*	JDBCServer2x进程启动信息日志。若进程停止，会打印jstack信息。
JobHistory2x日志	jobHistory-start.log	JobHistory2x启动日志。
	jobHistory-stop.log	JobHistory2x停止日志。
	JobHistory.log	JobHistory2x运行过程日志。

日志类型	日志文件名	描述
	jobhistory-omm-pid***-gc.log.*.current	JobHistory2x进程gc日志。
	spark-omm-org.apache.spark.deploy.history.HistoryServer-***.out*	JobHistory2x进程启动信息日志。若进程停止，会打印jstack信息。
IndexServer2x日志	IndexServer-start.log	IndexServer2x启动日志。
	IndexServer-stop.log	IndexServer2x停止日志。
	IndexServer.log	IndexServer2x运行时，Driver端日志。
	indexserver-state-check.log	IndexServer2x健康检查日志。
	indexserver-omm-pid***-gc.log.*.current	IndexServer2x进程gc日志。
	spark-omm-org.apache.spark.sql.hive.thriftserver.IndexServerProxy-***.out*	IndexServer2x进程启动信息日志。若进程停止，会打印jstack信息。
审计日志	jdbcservice-audit.log ranger-audit.log	JDBCServer2x审计日志。

日志级别

Spark2x中提供了如表26-65所示的日志级别。日志级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG。程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 26-65 日志级别

级别	描述
ERROR	ERROR表示当前时间处理存在错误信息。
WARN	WARN表示当前事件处理存在异常信息。
INFO	INFO表示记录系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

说明

默认情况下配置Spark2x日志级别不需要重启服务。

- 步骤1 登录FusionInsight Manager系统。
- 步骤2 选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”。
- 步骤3 单击“全部配置”。
- 步骤4 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤5 选择所需修改的日志级别。
- 步骤6 单击“保存”，然后单击“确定”，成功后配置生效。

---结束

日志格式

表 26-66 日志格式

日志类型	格式	示例
运行日志	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生该日志的线程名字> <log中的message> <日志事 件的发生位置>	2014-09-22 11:16:23,980 INFO DAGScheduler: Final stage: Stage 0(reduce at SparkPi.scala:35)

26.11.4 调整 Spark 日志级别

配置场景

在某些场景下，当任务已经启动后，用户想要修改日志级别以定位问题或者查看想要的信息。

用户可以在进程启动前，在进程的JVM参数中增加参数“-Dlog4j.configuration.watch=true”来打开动态设置日志级别的功能。进程启动后，就可以通过修改进程对应的log4j配置文件，来调整日志打印级别。

目前支持动态设置日志级别功能的有：Driver日志、Executor日志、AM日志、JobHistory日志、JDBCServer日志。

允许设置的日志级别是：FATAL，ERROR，WARN，INFO，DEBUG，TRACE和ALL。

配置描述

在进程对应的JVM参数配置项中增加以下参数。

表 26-67 参数描述

参数	描述	默认值
-Dlog4j.configuration.watch	进程JVM参数，设置成“true”用于打开动态设置日志级别功能。	未配置，即为false。

Driver、Executor、AM进程的JVM参数如表26-68所示。在Spark客户端的配置文件“spark-defaults.conf”中进行配置。Driver、Executor、AM进程的日志级别在对应的JVM参数中的“-Dlog4j.configuration”参数指定的log4j配置文件中设置。

表 26-68 进程的 JVM 参数 1

参数	说明	默认日志级别
spark.driver.extraJavaOptions	Driver的JVM参数。	INFO
spark.executor.extraJavaOptions	Executor的JVM参数。	INFO
spark.yarn.am.extraJavaOptions	AM的JVM参数。	INFO

JobHistory Server和JDBCServer的JVM参数如表26-69所示。在服务端配置文件“ENV_VARS”中进行配置。JobHistory Server和JDBCServer的日志级别在服务端配置文件“log4j.properties”中设置。

表 26-69 进程的 JVM 参数 2

参数	说明	默认日志级别
GC_OPTS	JobHistory Server的JVM参数。	INFO
SPARK_SUBMIT_OPTS	JDBCServer的JVM参数。	INFO

示例：

为了动态修改Executor日志级别为DEBUG，在进程启动之前，修改“spark-defaults.conf”文件中的Executor的JVM参数“spark.executor.extraJavaOptions”，增加如下配置：

```
-Dlog4j.configuration.watch=true
```

提交用户应用后，修改“spark.executor.extraJavaOptions”中“-Dlog4j.configuration”参数指定的log4j日志配置文件（例如：“-Dlog4j.configuration=file:\${BIGDATA_HOME}/FusionInsight_Spark2x_8.1.0.1/install/FusionInsight-Spark2x-*/spark/conf/log4j-executor.properties”）中的日志级别为DEBUG，如下所示：

```
log4j.rootCategory=DEBUG, sparklog
```

DEBUG级别生效会有一定的时延。

26.11.5 配置 WebUI 上查看 Container 日志

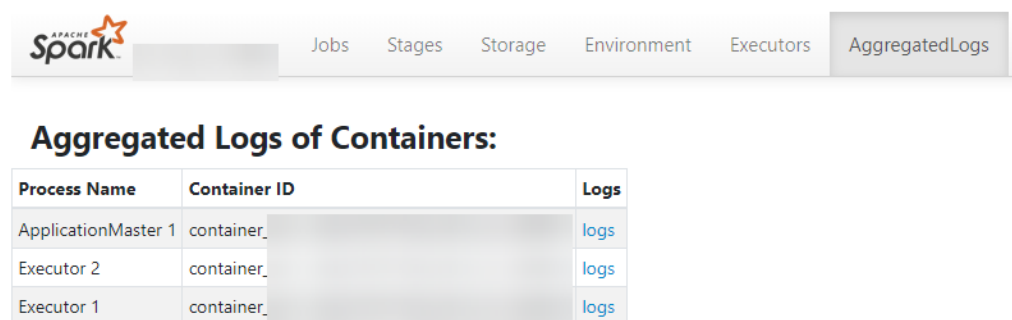
配置场景

当Yarn配置“yarn.log-aggregation-enable”为“true”时，就开启了container日志聚合功能。日志聚合功能是指：当应用在Yarn上执行完成后，NodeManager将本节点中所有container的日志聚合到HDFS中，并删除本地日志。详情请参见[配置Container日志聚合功能](#)。

然而，开启container日志聚合功能之后，其日志聚合至HDFS目录中，只能通过获取HDFS文件来查看日志。开源Spark和Yarn服务不支持通过WebUI查看聚合后的日志。

因此，Spark在此基础上进行了功能增强。如[图26-8](#)所示，在HistoryServer页面添加“AggregatedLogs”页签，可以通过“logs”链接查看聚合的日志。

图 26-8 聚合日志显示页面



配置描述

为了使WebUI页面显示日志，需要将聚合日志进行解析和展现。Spark是通过Hadoop的JobHistoryServer来解析聚合日志的，所以您可以通过“spark.jobhistory.address”参数，指定JobHistoryServer页面地址，即可完成解析和展现。

参数入口：

在应用提交时通过“--conf”设置这些参数，或者在客户端的“spark-defaults.conf”配置文件中调整如下参数。

📖 说明

- 此功能依赖Hadoop中的JobHistoryServer服务，所以使用聚合日志之前需要保证JobHistoryServer服务已经运行正常。
- 如果参数值为空，“AggregatedLogs”页签仍然存在，但是无法通过logs链接查看日志。
- 只有当App已经running，HDFS上已经有该App的事件日志文件时才能查看到聚合的container日志。
- 正在运行的任务的日志，用户可以通过“Executors”页面的日志链接进行查看，任务结束后日志会汇聚到HDFS上，“Executors”页面的日志链接就会失效，此时用户可以通过“AggregatedLogs”页面的logs链接查看聚合日志。

表 26-70 参数说明

参数	描述	默认值
spark.jobhistory.address	JobHistoryServer页面的地址，格式： <i>http(s)://ip:port/jobhistory</i> 。例如，将参数值设置为“ <i>https://10.92.115.1:26014/jobhistory</i> ”。 默认值为空，表示不能从WebUI查看container聚合日志。 修改参数后，需重启服务使得配置生效。	-

26.11.6 获取运行中 Spark 应用的 Container 日志

运行中Spark应用的Container日志分散在多个节点中，本章节用于说明如何快速获取Container日志。

场景说明

可以通过yarn logs命令获取运行在Yarn上的应用的日志，针对不同的场景，可以使用以下命令获取需要的日志：

1. 获取application的完整日志：**yarn logs --applicationId <appld> -out <outputDir>**

例如：**yarn logs --applicationId application_1574856994802_0016 -out /opt/test**

执行结果：

- a. 若该application处于运行状态，则无法获取dead状态的container日志
- b. 若该application处于结束状态，则可以获取全部归档的container日志

2. 获取指定Container日志：**yarn logs -applicationId <appld> -containerId <containerId>**

例如：**yarn logs -applicationId application_1574856994802_0018 -containerId container_e01_1574856994802_0018_01_000003**

执行结果：

- a. 若该application处于运行状态，则无法获取dead状态的Container日志
- b. 若该application处于结束状态，则可获取任意Container的日志

3. 获取任意状态的Container日志：**yarn logs -applicationId <appld> -containerId <containerId> -nodeAddress <nodeAddress>**

例如：**yarn logs -applicationId application_1574856994802_0019 -containerId container_e01_1574856994802_0019_01_000003 -nodeAddress 192-168-1-1:8041**

执行结果：可获取任意Container的日志

说明

此命令的参数中需要填入nodeAddress，可通过以下命令获取：

yarn node -list -all

26.11.7 配置 Spark Eventlog 日志回滚

配置场景

当Spark开启事件日志模式，即设置“spark.eventLog.enabled”为“true”时，就会往配置的一个日志文件中写事件，记录程序的运行过程。当程序运行很久，job很多，task很多时就会造成日志文件很大，如JDBCServer、Spark Streaming程序。

而日志回滚功能是指在写事件日志时，将元数据事件（EnvironmentUpdate, BlockManagerAdded, BlockManagerRemoved, UnpersistRDD, ExecutorAdded, ExecutorRemoved, MetricsUpdate, ApplicationStart, ApplicationEnd, LogStart）写入日志文件中，Job事件（StageSubmitted, StageCompleted, TaskResubmit, TaskStart, TaskEnd, TaskGettingResult, JobStart, JobEnd）按文件的大小进行决定是否写入新的日志文件。对于Spark SQL的应用，Job事件还包含ExecutionStart、ExecutionEnd。

Spark中有个HistoryServer服务，其UI页面就是通过读取解析这些日志文件获得的。在启动HistoryServer进程时，内存大小就已经定了。因此当日志文件很大时，加载解析这些文件就可能会造成内存不足，driver gc等问题。

所以为了在小内存模式下能加载较大日志文件，需要对大应用开启日志滚动功能。一般情况下，长时间运行的应用建议打开该功能。

配置参数

登录FusionInsight Manager系统，选择“集群 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索以下参数。

参数	说明	默认值
spark.eventLog.rolling.enabled	是否启用滚动event log文件。如果设置为true，则会将每个event log文件缩减到配置的大小。	true
spark.eventLog.rolling.maxFileSize	当spark.eventlog.rolling.enabled=true时，指定要滚动的event log文件的最大大小。	128M
spark.eventLog.compression.codec	用于压缩事件日志的编码解码器。默认情况下，spark提供四种编码解码器：lz4、lzf、snappy和zstd。如果没有给出，将使用spark.io.compression.codec。	无
spark.eventLog.logStageExecutorMetrics	是否将executor metrics的每个stage峰值（针对每个executor）写入event log。	false

26.11.8 配置 WebUI 上显示的 Lost Executor 信息的个数

配置场景

Spark WebUI 中“Executor”页面支持展示 Lost Executor 的信息，对于 JDBCServer 长任务来说，Executor 的动态回收是常态，Lost Executor 个数太多，会撑爆“Executor”页面，因此需要控制页面显示的 Lost Executor 个数。

配置描述

在 Spark 客户端的“spark-defaults.conf”配置文件中设置。

表 26-71 参数说明

参数	说明	默认值
spark.ui.retainedDeadExecutors	Spark UI 页面显示的 Lost Executor 的最大个数。	100

26.11.9 配置 JobHistory 本地磁盘缓存

配置场景

JobHistory 可使用本地磁盘缓存 spark 应用的历史数据，以防止 JobHistory 内存中加载大量应用数据，减少内存压力，同时该部分缓存数据可以复用以提高后续对相同应用的访问速度。

配置参数

登录 FusionInsight Manager 系统，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索以下参数。

参数	说明	默认值
spark.history.store.path	JobHistory 缓存历史信息的本地目录，如果设置了此配置，则 JobHistory 会将历史应用数据缓存在本地磁盘而不是内存中	\${BIGDATA_HOME}/tmp/spark2x_JobHistory
spark.history.store.maxDiskUsage	JobHistory 本地磁盘缓存的最大可用空间	10g

26.11.10 增强有限内存下的稳定性

配置场景

当前Spark SQL执行一个查询时需要使用大量的内存，尤其是在做聚合（Aggregate）和关联（Join）操作时，此时如果内存有限的情况下就容易出现OutOfMemoryError。有限内存下的稳定性就是确保在有限内存下依然能够正确执行相关的查询，而不出现OutOfMemoryError。

说明

有限内存并不意味着内存无限小，它只是在内存不足以放下大于内存可用总量几倍的数据时，通过利用磁盘来做辅助从而确保查询依然稳定执行，但依然有一些数据是必须留在内存的，如在做涉及到Join的查询时，对于当前用于Join的相同key的数据还是需要放在内存中，如果该数据量较大而内存较小依然会出现OutOfMemoryError。

有限内存下的稳定性涉及到3个子功能：

1. ExternalSort
外部排序功能，当执行排序时如果内存不足会将一部分数据溢出到磁盘中。
2. TungstenAggregate
新Hash聚合功能，默认对数据调用外部排序进行排序，然后再进行聚合，因此内存不足时在排序阶段会将数据溢出到磁盘，在聚合阶段因数据有序，在内存中只保留当前key的聚合结果，使用的内存较小。
3. SortMergeJoin、SortMergeOuterJoin
基于有序数据的等值连接。该功能默认对数据调用外部排序进行排序，然后再进行等值连接，因此内存不足时在排序阶段会将数据溢出到磁盘，在连接阶段因数据有序，在内存中只保留当前相同key的数据，使用的内存较小。

配置描述

参数入口：

在应用提交时通过“--conf”设置这些参数，或者在客户端的“spark-defaults.conf”配置文件中调整如下参数。

表 26-72 参数说明

参数	场景	描述	默认值
spark.sql.tungsten.enabled	/	类型为Boolean。 <ul style="list-style-type: none">当设置的值等于true时，表示开启tungsten功能，即逻辑计划等同于开启codegeneration，同时物理计划使用对应的tungsten执行计划。当设置的值等于false时，表示关闭tungsten功能。	true

参数	场景	描述	默认值
spark.sql.codegen.wholeStage		类型为Boolean。 ● 当设置的值等于true时，表示开启codegeneration功能，即运行时对于某些特定的查询将动态生成各逻辑计划代码。 ● 当设置的值等于false时，表示关闭codegeneration功能，运行时使用当前已有静态代码。	true

说明

1. 开启ExternalSort除配置spark.sql.planner.externalSort=true外，还需配置spark.sql.unsafe.enabled=false或者spark.sql.codegen.wholeStage =false。
2. 如果您需要开启TungstenAggregate，有如下几种方式：
 将spark.sql.codegen.wholeStage 和spark.sql.unsafe.enabled的值都设置为true（通过配置文件或命令行方式设置）。
 如果spark.sql.codegen.wholeStage 和spark.sql.unsafe.enabled都不为true或者其中一个不为true，只要spark.sql.tungsten.enabled的值设置为true时，TungstenAggregate会开启。

26.11.11 配置 YARN-Client 和 YARN-Cluster 不同模式下的环境变量

配置场景

当前，在YARN-Client和YARN-Cluster模式下，两种模式的客户端存在冲突的配置，即当客户端为一种模式的配置时，会导致在另一种模式下提交任务失败。

为避免出现如上情况，添加表26-73中的配置项，避免两种模式下来回切换参数，提升软件易用性。

- YARN-Cluster模式下，优先使用新增配置项的值，即服务端路径和参数。
- YARN-Client模式下，直接使用原有的三个配置项的值。
 原有的三个配置项为：“spark.driver.extraClassPath”、“spark.driver.extraJavaOptions”、“spark.driver.extraLibraryPath”。

说明

不添加表26-73中配置项时，使用方式与原有方式一致，程序可正常执行，只是在不同模式下需切换配置。

配置参数

参数入口：

在Manager系统中，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，在搜索框中输入参数名称。

表 26-73 参数介绍

参数	描述	默认值
spark.yarn.cluster.driver.extraClassPath	YARN-Cluster模式下，Driver使用的extraClassPath，配置为服务端的路径和参数。 同时，“spark.driver.extraClassPath”配置成Spark客户端路径，可以保证在YARN-Client模式下和YARN-Cluster模式下不需要切换配置。	<code>\${BIGDATA_HOME}/common/runtime/security</code>
spark.yarn.cluster.driver.extraJavaOptions	YARN-Cluster模式下Driver的extraJavaOptions，配置成服务端的路径和参数。 同时，“spark.driver.extraJavaOptions”配置成Spark客户端路径，可以保证YARN-Client模式和YARN-Cluster模式不需要切换配置。	<code>-Xloggc:<LOG_DIR>/indexserver-%p-gc.log -XX:+PrintGCDetails -XX:-OmitStackTraceInFastThrow -XX:+PrintGCTimeStamps -XX:+PrintGCDateStamps -XX:+UseGCLogFileRotation -XX:NumberOfGCLogFiles=20 -XX:GCLogFileSize=10M -Dlog4j.configuration=./__spark_conf__/__hadoop_conf__/log4j-executor.properties -Dlog4j.configuration.watch=true -Djava.security.auth.login.config=./__spark_conf__/__hadoop_conf__/jaas-zk.conf -Dzookeeper.server.principal=\${ZOOKEEPER_SERVER_PRINCIPAL} -Djava.security.krb5.conf=./__spark_conf__/__hadoop_conf__/kdc.conf -Djetty.version=x.y.z -Dorg.xerial.snappy.tmpdir=\${BIGDATA_HOME}/tmp -Dcarbon.properties.filepath=./__spark_conf__/__hadoop_conf__/carbon.properties -Djdk.tls.ephemeralDHKeySize=2048 -Dspark.ssl.keyStore=./child.keystore #{java_stack_prefer}</code>

26.11.12 Hive 分区修剪的谓词下推增强

配置场景

在旧版本中，对Hive表的分区修剪的谓词下推，只支持列名与整数或者字符串的比较表达式的下推，在2.3版本中，增加了对null、in、and、or表达式的下推支持。

配置参数

登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索以下参数。

参数	说明	默认值	取值范围
spark.sql.hive.advancedPartitionPredicatePushdown.enabled	用于配置是否开启Hive表的分区谓词下推增强功能。	true	[true,false]

26.11.13 配置列统计值直方图 Histogram 用以增强 CBO 准确度

配置场景

Spark优化sql的执行，一般的优化规则都是启发式的优化规则，启发式的优化规则，仅仅根据逻辑计划本身的特点给出优化，没有考虑数据本身的特点，也就是未考虑算子本身的执行代价。Spark在2.2中引入了基于代价的优化规则（CBO）。CBO会收集表和列的统计信息，结合算子的输入数据集来估计每个算子的输出条数以及字节大小，这些就是执行一个算子的代价。

CBO会调整执行计划，来最小化端到端的查询时间，中心思路2点：

- 尽早过滤不相关的数据。
- 最小化每个算子的代价。

CBO优化过程分为2步：

1. 收集统计信息。
2. 根据输入的数据集估算特定算子的输出数据集。

表级别统计信息包括：记录条数；表数据文件的总大小。

列级别统计信息包括：唯一值个数；最大值；最小值；空值个数；平均长度；最大长度；直方图。

有了统计信息后，就可以估算算子的执行代价了。常见的算子包括过滤条件Filter算子和Join算子。

直方图为列统计值的一种，可以直观的描述列数据的分布情况，将列的数据从最小值到最大值划分为事先指定数量的槽位（bin），计算各个槽位的上下界的值，使得全部数据都确定槽位后，所有槽位中的数据数量相同（等高直方图）。有了数据的详细分布后，各个算子的代价估计能更加准确，优化效果更好。

该特性可以通过下面的配置项开启：

spark.sql.statistics.histogram.enabled: 指定是否开启直方图功能，默认为false。

配置参数

登录FusionInsight Manager系统，选择“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，搜索以下参数。

参数	说明	默认值	取值范围
spark.sql.cbo.enabled	开启CBO来估计执行计划的统计值。	false	[true,false]
spark.sql.cbo.joinReorder.enabled	开启CBO连接重排序。	false	[true,false]
spark.sql.cbo.joinReorder.dp.threshold	动态规划算法中允许的最大的join节点数量。	12	>=1
spark.sql.cbo.joinReorder.card.weight	在重连接执行计划代价比较中维度（行数）所占的比重： $\text{行数} * \text{比重} + \text{文件大小} * (1 - \text{比重})$ 。	0.7	0-1
spark.sql.statistics.size.autoUpdate.enabled	开启当表的数据发生变化时，自动更新表的大小信息。注意如果表的数据文件总数量非常多时，这个操作会非常耗费资源，减慢对数据的操作速度。	false	[true,false]
spark.sql.statistics.histogram.enabled	开启后，当统计列信息时，会生成直方图。直方图可以提高估计准确度，但是收集直方图信息会有额外工作量。	false	[true,false]
spark.sql.statistics.histogram.numBins	生成的直方图的槽位数。	254	>=2
spark.sql.statistics.ndv.maxError	在生成列级别统计信息时，HyperLogLog++算法允许的最大估计误差。	0.05	0-1
spark.sql.statistics.percentile.accuracy	在生成等高直方图时百分位估计的准确率。该值越大意味着越准确。估计错误值可以通过 $(1.0 / \text{百分位估计的准确率})$ 来得到。	10000	>=1

📖 说明

- 如果希望直方图可以在CBO中生效，需要满足下面的条件：
 - spark.sql.statistics.histogram.enabled : true，默认是false，修改为true开启直方图功能。
 - spark.sql.cbo.enabled : true，默认为false，修改为true开启CBO。
 - spark.sql.cbo.joinReorder.enabled : true，默认为false，修改为true开启连接重排序。
- 若使用客户端提交任务，“spark.sql.cbo.enabled”、“spark.sql.cbo.joinReorder.enabled”、“spark.sql.cbo.joinReorder.dp.threshold”、“spark.sql.cbo.joinReorder.card.weight”、“spark.sql.statistics.size.autoUpdate.enabled”、“spark.sql.statistics.histogram.enabled”、“spark.sql.statistics.histogram.numBins”、“spark.sql.statistics.ndv.maxError”、“spark.sql.statistics.percentile.accuracy”参数修改后需要重新下载客户端才能生效。

26.11.14 CarbonData 首查优化工具

工具介绍

CarbonData 的首次查询较慢，对于实时性要求较高的节点可能会造成一定的时延。

本工具主要提供以下功能：

- 对查询时延要求较高的表进行首次查询预热。

工具使用

下载安装客户端，例如安装目录为“/opt/client”。进入目录“/opt/client/Spark2x/spark/bin”，执行start-prequery.sh。

参考表26-74，配置prequeryParams.properties。

表 26-74 参数列表

参数	说明	示例
spark.prequery.period.max.minute	预热的最大时长，单位分钟	60
spark.prequery.tables	表名配置 database.table:int，表名支持通配符*，int代表预热多长时间内有更新的表，单位为天。	default.test*:10
spark.prequery.maxThreads	预热时并发的最大线程数	50
spark.prequery.sslEnable	集群安全模式为true，非安全模式为false	true

参数	说明	示例
spark.prequery.driver	JDBCServer的地址 ip:port，如需要预热多个 Server则需填写多个 Server的IP,多个IP:port用 逗号隔开。	192.168.0.2:22550
spark.prequery.sql	预热的sql语句，不同语句 冒号隔开	SELECT COUNT(*) FROM %s;SELECT * FROM %s LIMIT 1
spark.security.url	安全模式下jdbc所需url	;sasLQop=auth- conf;auth=KERBEROS;pri ncipal=spark2x/ hadoop.hadoop.com@HA DOOP.COM;

说明

spark.prequery.sql 配置的语句在每个所预热的表中都会执行，表名用%s代替。

脚本使用

命令形式：**sh start-prequery.sh**

执行此条命令需要：将user.keytab或jaas.conf（二选一），krb5.conf（必须）放入conf目录中。

说明

- 此工具暂时只支持Carbon表。
- 此工具会初始化Carbon环境和预读取表的元数据到JDBCServer，所以更适合在多主实例、静态分配模式下使用。

26.12 Spark2x 常见问题

26.12.1 Spark Core

26.12.1.1 日志聚合下如何查看 Spark 已完成应用日志

问题

当YARN开启了日志聚合功能时，如何在页面看到聚合后的container日志？

回答

当Yarn配置“yarn.log-aggregation-enable”为“true”时，就开启了container日志聚合功能。

日志聚合功能是指：当应用在Yarn上执行完成后，NodeManager将本节点中所有container的日志聚合到HDFS中，并删除本地日志。详情请参见[配置Container日志聚合功能](#)。

然而，开启container日志聚合功能之后，其日志聚合至HDFS目录中，只能通过获取HDFS文件来查看日志。

开源Spark和Yarn服务不支持通过WebUI查看聚合后的日志。

26.12.1.2 Driver 返回码和 RM WebUI 上应用状态显示不一致

问题

ApplicationMaster与ResourceManager之间通信发生长时间异常时，为什么Driver返回码和RM WebUI上应用状态显示不一致？

回答

在yarn-client模式下，Spark的Driver和ApplicationMaster作为两个独立的进程在运行。当Driver完成任务退出时，会通知ApplicationMaster向ResourceManager注销自身，即调用unregister方法。

由于是远程调用，则存在发生网络故障的可能性。当发生网络故障时，ApplicationMaster会使用Yarn客户端的重试机制进行重试。在达到最大重试次数之前网络恢复正常，则ApplicationMaster会正常退出。

若超过重试次数和重试时长，则ApplicationMaster注销失败，ResourceManager会认为ApplicationMaster异常退出并尝试重新启动ApplicationMaster。新启动的ApplicationMaster在尝试连接已经退出的Driver失败后，会在ResourceManager页面上标记此次Application为FAILED状态。

这种情况为小概率事件且不影响Spark SQL对外展现的应用完成状态。也可以通过增大Yarn客户端连接次数和连接时长的方式减少此事件发生的概率。配置详情请参见：

<http://hadoop.apache.org/docs/r3.1.1/hadoop-yarn/hadoop-yarn-common/yarn-default.xml>

26.12.1.3 为什么 Driver 进程不能退出

问题

运行Spark Streaming任务，然后使用`yarn application -kill applicationID`命令停止任务，为什么Driver进程不能退出？

回答

使用`yarn application -kill applicationID`命令后Spark只会停掉任务对应的SparkContext，而不是退出当前进程。如果当前进程中存在其他常驻的线程（类似spark-shell需要不断检测命令输入，Spark Streaming不断在从数据源读取数据），SparkContext被停止并不会终止整个进程。

如果需要退出Driver进程，建议使用`kill -9 pid`命令手动退出当前Driver。

26.12.1.4 网络连接超时导致 FetchFailedException

问题

在380节点的大集群上，运行29T数据量的HiBench测试套中ScalaSort测试用例，使用以下关键配置（`--executor-cores 4`）出现如下异常：

```
org.apache.spark.shuffle.FetchFailedException: Failed to connect to /192.168.114.12:23242
    at
    org.apache.spark.storage.ShuffleBlockFetcherIterator.throwFetchFailedException(ShuffleBlockFetcherIterator.scala:321)
    at org.apache.spark.storage.ShuffleBlockFetcherIterator.next(ShuffleBlockFetcherIterator.scala:306)
    at org.apache.spark.storage.ShuffleBlockFetcherIterator.next(ShuffleBlockFetcherIterator.scala:51)
    at scala.collection.Iterator$$anon$11.next(Iterator.scala:328)
    at scala.collection.Iterator$$anon$13.hasNext(Iterator.scala:371)
    at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:327)
    at org.apache.spark.util.CompletionIterator.hasNext(CompletionIterator.scala:32)
    at org.apache.spark.InterruptibleIterator.hasNext(InterruptibleIterator.scala:39)
    at org.apache.spark.util.collection.ExternalSorter.insertAll(ExternalSorter.scala:217)
    at org.apache.spark.shuffle.hash.HashShuffleReader.read(HashShuffleReader.scala:102)
    at org.apache.spark.rdd.ShuffledRDD.compute(ShuffledRDD.scala:90)
    at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:301)
    at org.apache.spark.rdd.RDD.iterator(RDD.scala:265)
    at org.apache.spark.rdd.MapPartitionsRDD.compute(MapPartitionsRDD.scala:38)
    at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:301)
    at org.apache.spark.rdd.RDD.iterator(RDD.scala:265)
    at org.apache.spark.rdd.MapPartitionsRDD.compute(MapPartitionsRDD.scala:38)
    at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:301)
    at org.apache.spark.rdd.RDD.iterator(RDD.scala:265)
    at org.apache.spark.rdd.UnionRDD.compute(UnionRDD.scala:87)
    at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:301)
    at org.apache.spark.rdd.RDD.iterator(RDD.scala:265)
    at org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask.scala:73)
    at org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask.scala:41)
    at org.apache.spark.scheduler.Task.run(Task.scala:87)
    at org.apache.spark.executor.Executor$TaskRunner.run(Executor.scala:213)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1142)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:617)
    at java.lang.Thread.run(Thread.java:745)
Caused by: java.io.IOException: Failed to connect to /192.168.114.12:23242
    at org.apache.spark.network.client.TransportClientFactory.createClient(TransportClientFactory.java:214)
    at org.apache.spark.network.client.TransportClientFactory.createClient(TransportClientFactory.java:167)
    at org.apache.spark.network.netty.NettyBlockTransferService$$anon$1.createAndStart(NettyBlockTransferService.scala:91)
    at
    org.apache.spark.network.shuffle.RetryingBlockFetcher.fetchAllOutstanding(RetryingBlockFetcher.java:140)
    at org.apache.spark.network.shuffle.RetryingBlockFetcher.access$200(RetryingBlockFetcher.java:43)
    at org.apache.spark.network.shuffle.RetryingBlockFetcher$1.run(RetryingBlockFetcher.java:170)
    at java.util.concurrent.Executors$RunnableAdapter.call(Executors.java:511)
    at java.util.concurrent.FutureTask.run(FutureTask.java:266)
    ... 3 more
Caused by: java.net.ConnectException: Connection timed out: /192.168.114.12:23242
    at sun.nio.ch.SocketChannelImpl.checkConnect(Native Method)
    at sun.nio.ch.SocketChannelImpl.finishConnect(SocketChannelImpl.java:717)
    at io.netty.channel.socket.nio.NioSocketChannel.doFinishConnect(NioSocketChannel.java:224)
    at io.netty.channel.nio.AbstractNioChannel
    $AbstractNioUnsafe.finishConnect(AbstractNioChannel.java:289)
    at io.netty.channel.nio.NioEventLoop.processSelectedKey(NioEventLoop.java:528)
    at io.netty.channel.nio.NioEventLoop.processSelectedKeysOptimized(NioEventLoop.java:468)
    at io.netty.channel.nio.NioEventLoop.processSelectedKeys(NioEventLoop.java:382)
    at io.netty.channel.nio.NioEventLoop.run(NioEventLoop.java:354)
    at io.netty.util.concurrent.SingleThreadEventExecutor$2.run(SingleThreadEventExecutor.java:111)
    ... 1 more
```

回答

在运行应用程序时，使用Executor参数“--executor-cores 4”，单进程中并行度高导致IO非常繁忙，以至于任务运行缓慢。

```
16/02/26 10:04:53 INFO TaskSetManager: Finished task 2139.0 in stage 1.0 (TID 151149) in 376455 ms on 10-196-115-2 (694/153378)
```

单个任务运行时间超过6分钟，从而导致连接超时问题，最终使得任务失败。

将参数中的核数设置为1，“--executor-cores 1”，任务正常完成，单个任务处理时间在合理范围之内(15秒左右)。

```
16/02/29 02:24:46 INFO TaskSetManager: Finished task 59564.0 in stage 1.0 (TID 208574) in 15088 ms on 10-196-115-6 (59515/153378)
```

因此，处理这类网络超时任务，可以减少单个Executor的核数来规避该类问题。

26.12.1.5 当事件队列溢出时如何配置事件队列的大小

问题

当Driver日志中出现如下的日志时，表示事件队列溢出了。当事件队列溢出时如何配置事件队列的大小？

- 普通应用
Dropping SparkListenerEvent because no remaining room in event queue.
This likely means one of the SparkListeners is too slow and cannot keep up with the rate at which tasks are being started by the scheduler.
- Spark Streaming应用
Dropping StreamingListenerEvent because no remaining room in event queue.
This likely means one of the StreamingListeners is too slow and cannot keep up with the rate at which events are being started by the scheduler.

回答

1. 停止应用，在Spark的配置文件“spark-defaults.conf”中将配置项“spark.event.listener.logEnable”配置为“true”。并把配置项“spark.eventQueue.size”配置为1000W。如果需要控制打印频率（默认为1000毫秒打印1条日志），请根据需要修改配置项“spark.event.listener.logRate”，该配置项的单位为毫秒。
2. 启动应用，可以发现如下的日志信息（消费者速率、生产者速率、当前队列中的消息数量和队列中消息数量的最大值）。

```
INFO LiveListenerBus: [SparkListenerBus]:16044 events are consumed in 5000 ms.  
INFO LiveListenerBus: [SparkListenerBus]:51381 events are produced in 5000 ms, eventQueue still has 86417 events, MaxSize: 171764.
```
3. 用户可以根据日志信息【队列中消息数量的最大值MaxSize】，在配置文件“spark-defaults.conf”中将配置项“spark.eventQueue.size”配置成合适的队列大小。比如【队列中消息数量的最大值】为250000，那么配置合适的队列大小为300000。

26.12.1.6 Spark 应用执行过程中，日志中一直打印 getApplicationReport 异常且应用较长时间不退出

问题

Spark应用执行过程中，当driver连接RM失败时，会报下面的错误，且较长时间不退出。

```
16/04/23 15:31:44 INFO RetryInvocationHandler: Exception while invoking getApplicationReport of class
ApplicationClientProtocolPBClientImpl over 37 after 1 fail over attempts. Trying to fail over after sleeping
for 44160ms.
java.net.ConnectException: Call From vm1/192.168.39.30 to vm1:8032 failed on connection exception:
java.net.ConnectException: Connection refused; For more details see: http://wiki.apache.org/hadoop/
ConnectionRefused
```

回答

在Spark中有个定期线程，通过连接RM监测AM的状态。由于连接RM超时，就会报上面的错误，且一直重试。RM中对重试次数有限制，默认是30次，每次间隔默认为30秒左右，每次重试时都会报上面的错误。超过次数后，driver才会退出。

RM中关于重试相关的配置项如表26-75所示。

表 26-75 参数说明

参数	描述	默认值
yarn.resourcemanager.connect.max-wait.ms	连接RM的等待时间最大值。	900000
yarn.resourcemanager.connect.retry-interval.ms	重试连接RM的时间频率。	30000

重试次数=yarn.resourcemanager.connect.max-wait.ms/
yarn.resourcemanager.connect.retry-interval.ms，即重试次数=连接RM的等待时间最大值/重试连接RM的时间频率。

在Spark客户端机器中，通过修改“conf/yarn-site.xml”文件，添加并配置“yarn.resourcemanager.connect.max-wait.ms”和“yarn.resourcemanager.connect.retry-interval.ms”，这样可以更改重试次数，Spark应用可以提早退出。

26.12.1.7 Spark 执行应用时上报“Connection to ip:port has been quiet for xxx ms while there are outstanding requests”并导致应用结束

问题

Spark执行应用时上报如下类似错误并导致应用结束。

```
2016-04-20 10:42:00,557 | ERROR | [shuffle-server-2] | Connection to 10-91-8-208/10.18.0.115:57959 has
been quiet for 180000 ms while there are outstanding requests. Assuming connection is dead; please adju
st spark.network.timeout if this is wrong. |
org.apache.spark.network.server.TransportChannelHandler.userEventTriggered(TransportChannelHandler.java:
128)
2016-04-20 10:42:00,558 | ERROR | [shuffle-server-2] | Still have 1 requests outstanding when connection
```

```
from 10-91-8-208/10.18.0.115:57959 is closed | org.apache.spark.network.client.TransportResponseHandl
er.channelUnregistered(TransportResponseHandler.java:102)
2016-04-20 10:42:00,562 | WARN | [yarn-scheduler-ask-am-thread-pool-160] | Error sending message
 [message = DoShuffleClean(application_1459995017785_0108,319)] in 1 attempts |
org.apache.spark.Logging$class
s.logWarning(Logging.scala:92)
java.io.IOException: Connection from 10-91-8-208/10.18.0.115:57959 closed
    at
org.apache.spark.network.client.TransportResponseHandler.channelUnregistered(TransportResponseHandler.j
ava:104)
    at
org.apache.spark.network.server.TransportChannelHandler.channelUnregistered(TransportChannelHandler.jav
a:94)
    at
io.netty.channel.AbstractChannelHandlerContext.invokeChannelUnregistered(AbstractChannelHandlerContext
.java:158)
    at
io.netty.channel.AbstractChannelHandlerContext.fireChannelUnregistered(AbstractChannelHandlerContext.ja
va:144)
    at
io.netty.channel.ChannelInboundHandlerAdapter.channelUnregistered(ChannelInboundHandlerAdapter.java:5
3)
    at
io.netty.channel.AbstractChannelHandlerContext.invokeChannelUnregistered(AbstractChannelHandlerContext
.java:158)
    at
io.netty.channel.AbstractChannelHandlerContext.fireChannelUnregistered(AbstractChannelHandlerContext.ja
va:144)
    at
io.netty.channel.ChannelInboundHandlerAdapter.channelUnregistered(ChannelInboundHandlerAdapter.java:5
3)
    at
io.netty.channel.AbstractChannelHandlerContext.invokeChannelUnregistered(AbstractChannelHandlerContext
.java:158)
    at
io.netty.channel.AbstractChannelHandlerContext.fireChannelUnregistered(AbstractChannelHandlerContext.ja
va:144)
    at
io.netty.channel.ChannelInboundHandlerAdapter.channelUnregistered(ChannelInboundHandlerAdapter.java:5
3)
    at
io.netty.channel.AbstractChannelHandlerContext.invokeChannelUnregistered(AbstractChannelHandlerContext
.java:158)
    at
io.netty.channel.AbstractChannelHandlerContext.fireChannelUnregistered(AbstractChannelHandlerContext.ja
va:144)
    at
io.netty.channel.ChannelInboundHandlerAdapter.channelUnregistered(ChannelInboundHandlerAdapter.java:5
3)
    at
io.netty.channel.AbstractChannelHandlerContext.invokeChannelUnregistered(AbstractChannelHandlerContext
.java:158)
    at
io.netty.channel.AbstractChannelHandlerContext.fireChannelUnregistered(AbstractChannelHandlerContext.ja
va:144)
    at
io.netty.channel.DefaultChannelPipeline.fireChannelUnregistered(DefaultChannelPipeline.java:739)
    at io.netty.channel.AbstractChannel$AbstractUnsafe$8.run(AbstractChannel.java:659)
    at io.netty.util.concurrent.SingleThreadEventExecutor.runAllTasks(SingleThreadEventExecutor.java:357)
    at io.netty.channel.nio.NioEventLoop.run(NioEventLoop.java:357)
    at io.netty.util.concurrent.SingleThreadEventExecutor$2.run(SingleThreadEventExecutor.java:111)
    at java.lang.Thread.run(Thread.java:745)
2016-04-20 10:42:00,573 | INFO | [dispatcher-event-loop-14] | Starting task 177.0 in stage 1492.0 (TID
1996351, linux-254, PROCESS_LOCAL, 2106 bytes) | org.apache.spark.Logging$class.logInfo(Logging.scala:
59)
2016-04-20 10:42:00,574 | INFO | [task-result-getter-0] | Finished task 85.0 in stage 1492.0 (TID 1996259)
in 191336 ms on linux-254 (106/3000) | org.apache.spark.Logging$class.logInfo(Logging.scala:59)
2016-04-20 10:42:00,811 | ERROR | [Yarn application state monitor] | Yarn application has already exited
with state FINISHED! | org.apache.spark.Logging$class.logError(Logging.scala:75)
```

回答

当配置channel过期时间（spark.rpc.io.connectionTimeout）< RPC响应超时时间（spark.rpc.askTimeout），在特殊条件下（Full GC，网络延时等）消息响应时间较长，消息还没有反馈，channel又达到了过期时间，该channel就被终止了，AM端感知到channel被终止后认为driver失联，然后整个应用停止。

解决办法：在Spark客户端的“spark-defaults.conf”文件中或通过set命令进行设置。参数配置时要保证channel过期时间（spark.rpc.io.connectionTimeout）大于或等于RPC响应超时时间（spark.rpc.askTimeout）。

表 26-76 参数说明

参数	描述	默认值
spark.rpc.askTimeout	RPC响应超时时间，不配置的话默认使用spark.network.timeout的值。	120s

26.12.1.8 NodeManager 关闭导致 Executor(s)未移除

问题

在Executor动态分配打开的情况下，如果在任务执行过程中，执行NodeManager关闭动作，NodeManager关闭节点上的Executor(s)在空闲超时之后，在driver页面上未被移除。

回答

因为ResourceManager感知到NodeManager关闭时，Executor(s)已经因空闲超时而被driver请求结束。

但因为NodeManager已经关闭，这些Executor(s)实际上并不能被结束，因此driver不能感知到这些Executor(s)的LOST事件，所以并未从自身的Executor list中移除。

从而导致在driver页面上还能看到这些Executor(s)，这是YARN NodeManager关闭之后的正常现象，NodeManager再次启动后，这些Executor(s)会被移除。

26.12.1.9 Password cannot be null if SASL is enabled 异常

问题

运行Spark的应用启用了ExternalShuffle，应用出现了Task任务丢失，原因是由于java.lang.NullPointerException: Password cannot be null if SASL is enabled异常，部分关键日志如下图所示：

```
2016-05-13 12:05:27,093 | WARN | [task-result-getter-2] | Lost task 98.0 in stage 22.1 (TID 193603, linux-173. 2): FetchFailed(BlockManagerId(13, 172.168.100.13, 27337), org.apache.spark.shuffle.FetchFailedException: java.lang.NullPointerException: Password cannot be null if SASL is enabled
    at org.spark-project.guava.base.Preconditions.checkNotNull(Preconditions.java:208)
    at org.apache.spark.network.sasl.SparkSaslServer.encodePassword(SparkSaslServer.java:196)
    at org.apache.spark.network.sasl.SparkSaslServerDigestCallbackHandler.handle(SparkSaslServer.java:166)
    at com.sun.security.sasl.digest.DigestMD5Server.validateClientResponse(DigestMD5Server.java:589)
    at com.sun.security.sasl.digest.DigestMD5Server.evaluateResponse(DigestMD5Server.java:244)
    at org.apache.spark.network.sasl.SparkSaslServer.response(SparkSaslServer.java:119)
    at org.apache.spark.network.sasl.SaslRpcHandler.receive(SaslRpcHandler.java:100)
    at org.apache.spark.network.server.TransportRequestHandler.processRpcRequest(TransportRequestHandler.java:128)
    at org.apache.spark.network.server.TransportRequestHandler.handle(TransportRequestHandler.java:99)
    at org.apache.spark.network.server.TransportChannelHandler.channelRead0(TransportChannelHandler.java:104)
```

回答

造成该现象的原因是NodeManager重启。使用ExternalShuffle的时候，Spark将借用NodeManager传输Shuffle数据，因此NodeManager的内存将成为瓶颈。

在当前版本的FusionInsight中，NodeManager的默认内存只有1G，在数据量比较大（1T以上）的Spark任务下，内存严重不足，消息响应缓慢，导致FusionInsight健康检查认为NodeManager进程退出，强制重启NodeManager，导致上述问题产生。

解决方式：

调整NodeManager的内存，数据量比较大（1T以上）的情况下，NodeManager的内存至少在4G以上。

26.12.1.10 向动态分区表中插入数据时，在重试的 task 中出现"Failed to CREATE_FILE"异常

问题

向动态分区表中插入数据时，shuffle过程中大面积shuffle文件损坏（磁盘掉线、节点故障等）后，为什么会在重试的task中出现"Failed to CREATE_FILE"异常？

```
2016-06-25 15:11:31,323 | ERROR | [Executor task launch worker-0] | Exception in task 15.0 in stage 10.1 (TID 1258) | org.apache.spark.Logging$class.logError(Logging.scala:96)
org.apache.hadoop.hive.ql.metadata.HiveException:
org.apache.hadoop.ipc.RemoteException(org.apache.hadoop.hdfs.protocol.AlreadyBeingCreatedException):
Failed to CREATE_FILE /user/hive/warehouse/testdb.db/web_sales/hive-staging_hive_2016-06-25_15-09-16_999_8137121701603617850-1/-ext-10000/_temporary/0/_temporary/attempt_201606251509_0010_m_000015_0/ws_sold_date=1999-12-17/part-00015 for
DFSClient_attempt_2016
06251509_0010_m_000015_0_353134803_151 on 10.1.1.5 because this file lease is currently owned by
DFSClient_attempt_201606251509_0010_m_000015_0_-848353830_156 on 10.1.1.6
```

回答

动态分区表插入数据的最后一步是读取shuffle文件的数据，再写入到表对应的分区文件中。

当大面积shuffle文件损坏后，会引起大批量task失败，然后进行job重试。重试前Spark会将写表分区文件的句柄关闭，大批量task关闭句柄时HDFS无法及时处理。在task进行下一次重试时，句柄在NameNode端未被及时释放，即会发生"Failed to CREATE_FILE"异常。

这种现象仅会在大面积shuffle文件损坏时发生，出现异常后task会重试，重试耗时在毫秒级，影响较小，可以忽略不计。

26.12.1.11 使用 Hash shuffle 出现任务失败

问题

使用Hash shuffle运行1000000（map个数）*100000（reduce个数）的任务，运行日志中出现大量的消息发送失败和Executor心跳超时，从而导致任务失败。

回答

对于Hash shuffle，在shuffle的过程中写数据时不做排序操作，只是将数据根据Hash的结果，将各个reduce分区的数据写到各自的磁盘文件中。

这样带来的问题是如果reduce分区的数量比较大的话，将会产生大量的磁盘文件（比如：该问题中将产生1000000 * 100000 = 10¹¹个shuffle文件）。如果磁盘文件数量特别巨大，对文件读写的性能会带来比较大的影响，此外由于同时打开的文件句柄数量多，序列化以及压缩等操作需要占用非常大的临时内存空间，对内存的使用和GC带来很大的压力，从而容易造成Executor无法响应Driver。

因此，建议使用Sort shuffle，而不使用Hash shuffle。

26.12.1.12 访问 Spark 应用的聚合日志页面报“DNS 查找失败”错误

问题

采用http(s)://<spark ip>:<spark port>的方式直接访问Spark JobHistory页面时，如果当前跳转的Spark JobHistory页面不是FusionInsight代理的页面（FusionInsight代理的URL地址类似于：https://<oms ip>:20026/Spark2x/JobHistory2x/xx/），单击某个应用，再单击“AggregatedLogs”，然后单击需要查看的其中一个Executor的“logs”，此时会报如图26-9所示的错误。

图 26-9 聚合日志失败页面



回答

原因：弹出的URL地址（如https://<hostname>:20026/Spark2x/JobHistory2x/xx/history/application_xxx/jobs/），其中的<hostname>没有在Windows系统的hosts文件中添加域名信息，导致DNS查找失败无法显示此网页。

解决措施：

- 建议用户使用FusionInsight代理去访问Spark JobHistory页面，即单击如图26-10中蓝框所示的Spark WebUI的链接。

图 26-10 FusionInsight Manager 的 Spark2x 页面



- 如果用户需要不通过FusionInsight Manager访问Spark JobHistory页面，则需要将URL地址中的<hostname>更改为IP地址进行访问，或者在Windows系统的hosts文件中添加该域名信息。

26.12.1.13 由于 Timeout waiting for task 异常导致 Shuffle FetchFailed

问题

使用JDBCServer模式执行100T的TPCDS测试套，出现Timeout waiting for task异常导致Shuffle FetchFailed，Stage一直重试，任务无法正常完成。

回答

JDBCServer方式使用了ShuffleService功能，Reduce阶段所有的Executor会从NodeManager中获取数据，当数据量达到一个级别（10T级别），会出现NodeManager单点瓶颈（ShuffleService服务在NodeManager进程中），就会出现某些Task获取数据超时，从而出现该问题。

因此，当数据量达到10T级别以上的Spark任务，建议用户关闭ShuffleService功能，即在“Spark-defaults.conf”配置文件中将配置项“spark.shuffle.service.enabled”配置为“false”。

26.12.1.14 Executor 进程 Crash 导致 Stage 重试

问题

在执行大数据量的Spark任务（如100T的TPCDS测试套）过程中，有时会出现Executor丢失从而导致Stage重试的现象。查看Executor的日志，出现“Executor 532 is lost rpc with driver,but is still alive, going to kill it”所示信息，表明Executor丢失是由于JVM Crash导致的。

JVM的关键Crash错误日志，如下：

```
#  
# A fatal error has been detected by the Java Runtime Environment:  
#  
# Internal Error (sharedRuntime.cpp:834), pid=241075, tid=140476258551552  
# fatal error: exception happened outside interpreter, nmethods and vtable stubs at pc  
0x00007fcda9eb8eb1
```

回答

上述问题在Oracle官网上有类似的情况，该问题现象是Oracle JVM的缺陷，并不是平台代码引入的问题，且Spark中有对Executor的容错机制，Executor Crash之后，Stage会进入重试，可以保证任务最终可以执行完成，不会对业务产生影响。

26.12.1.15 执行大数据量的 shuffle 过程时 Executor 注册 shuffle service 失败

问题

执行超过50T数据的shuffle过程时，出现部分Executor注册shuffle service超时然后丢失从而导致任务失败的问题。错误日志如下所示：

```
2016-10-19 01:33:34,030 | WARN | ContainersLauncher #14 | Exception from container-launch with  
container ID: container_e1452_1476801295027_2003_01_004512 and exit code: 1 |  
LinuxContainerExecutor.java:397  
ExitCodeException exitCode=1:  
at org.apache.hadoop.util.Shell.runCommand(Shell.java:561)  
at org.apache.hadoop.util.Shell.run(Shell.java:472)  
at org.apache.hadoop.util.Shell$ShellCommandExecutor.execute(Shell.java:738)  
at  
org.apache.hadoop.yarn.server.nodemanager.LinuxContainerExecutor.launchContainer(LinuxContainerExecuto  
r.java:381)  
at  
org.apache.hadoop.yarn.server.nodemanager.containermanager.launcher.ContainerLaunch.call(ContainerLaun  
ch.java:312)  
at  
org.apache.hadoop.yarn.server.nodemanager.containermanager.launcher.ContainerLaunch.call(ContainerLaun  
ch.java:88)  
at java.util.concurrent.FutureTask.run(FutureTask.java:266)  
at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1142)  
at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:617)  
at java.lang.Thread.run(Thread.java:745)  
2016-10-19 01:33:34,031 | INFO | ContainersLauncher #14 | Exception from container-launch. |  
ContainerExecutor.java:300  
2016-10-19 01:33:34,031 | INFO | ContainersLauncher #14 | Container id:  
container_e1452_1476801295027_2003_01_004512 | ContainerExecutor.java:300  
2016-10-19 01:33:34,031 | INFO | ContainersLauncher #14 | Exit code: 1 | ContainerExecutor.java:300  
2016-10-19 01:33:34,031 | INFO | ContainersLauncher #14 | Stack trace: ExitCodeException exitCode=1: |  
ContainerExecutor.java:300
```

回答

由于当前数据量较大，有50T数据导入，超过了shuffle的规格，shuffle负载过高，shuffle service服务处于过载状态，可能无法及时响应Executor的注册请求，从而出现上面的问题。

Executor注册shuffle service的超时时间是5秒，最多重试3次，该参数目前不可配。

建议适当调大task retry次数和Executor失败次数。

在客户端的“spark-defaults.conf”配置文件中配置如下参数。
“spark.yarn.max.executor.failures”若不存在，则手动添加该参数项。

表 26-77 参数说明

参数	描述	默认值
spark.task.maxFailures	task retry次数。	4
spark.yarn.max.executor.failures	Executor失败次数。 关闭Executor个数动态分配功能的场景即 “spark.dynamicAllocation.enabled”参数设为“false”时。	numExecutors * 2, with minimum of 3
	Executor失败次数。 开启Executor个数动态分配功能的场景即 “spark.dynamicAllocation.enabled”参数设为“true”时。	3

26.12.1.16 在 Spark 应用执行过程中 NodeManager 出现 OOM 异常

问题

当开启Yarn External Shuffle服务时，在Spark应用执行过程中，如果当前shuffle连接过多，Yarn External Shuffle会出现“java.lang.OutOfMemoryError: Direct buffer Memory”的异常，该异常说明内存不足。错误日志如下：

```
2016-12-06 02:01:00,768 | WARN | shuffle-server-38 | Exception in connection from /192.168.101.95:53680 |
TransportChannelHandler.java:79
io.netty.handler.codec.DecoderException: java.lang.OutOfMemoryError: Direct buffer memory
    at io.netty.handler.codec.ByteToMessageDecoder.channelRead(ByteToMessageDecoder.java:153)
    at
io.netty.channel.AbstractChannelHandlerContext.invokeChannelRead(AbstractChannelHandlerContext.java:333)
    at
io.netty.channel.AbstractChannelHandlerContext.fireChannelRead(AbstractChannelHandlerContext.java:319)
    at io.netty.channel.DefaultChannelPipeline.fireChannelRead(DefaultChannelPipeline.java:787)
    at io.netty.channel.nio.AbstractNioByteChannel$NioByteUnsafe.read(AbstractNioByteChannel.java:130)
    at io.netty.channel.nio.NioEventLoop.processSelectedKey(NioEventLoop.java:511)
    at io.netty.channel.nio.NioEventLoop.processSelectedKeysOptimized(NioEventLoop.java:468)
    at io.netty.channel.nio.NioEventLoop.processSelectedKeys(NioEventLoop.java:382)
    at io.netty.channel.nio.NioEventLoop.run(NioEventLoop.java:354)
    at io.netty.util.concurrent.SingleThreadEventExecutor$2.run(SingleThreadEventExecutor.java:116)
    at java.lang.Thread.run(Thread.java:745)
Caused by: java.lang.OutOfMemoryError: Direct buffer memory
    at java.nio.Bits.reserveMemory(Bits.java:693)
    at java.nio.DirectByteBuffer.<init>(DirectByteBuffer.java:123)
    at java.nio.ByteBuffer.allocateDirect(ByteBuffer.java:311)
    at io.netty.buffer.PoolArena$DirectArena.newChunk(PoolArena.java:434)
    at io.netty.buffer.PoolArena.allocateNormal(PoolArena.java:179)
    at io.netty.buffer.PoolArena.allocate(PoolArena.java:168)
    at io.netty.buffer.PoolArena.reallocate(PoolArena.java:277)
    at io.netty.buffer.PooledByteBuf.capacity(PooledByteBuf.java:108)
    at io.netty.buffer.AbstractByteBuf.ensureWritable(AbstractByteBuf.java:251)
    at io.netty.buffer.AbstractByteBuf.writeBytes(AbstractByteBuf.java:849)
    at io.netty.buffer.AbstractByteBuf.writeBytes(AbstractByteBuf.java:841)
    at io.netty.buffer.AbstractByteBuf.writeBytes(AbstractByteBuf.java:831)
    at io.netty.handler.codec.ByteToMessageDecoder.channelRead(ByteToMessageDecoder.java:146)
    ... 10 more
```

回答

对于Yarn的Shuffle Service，其启动的线程数为机器可用CPU核数的两倍，而默认配置的Direct buffer Memory为128M，因此当有较多shuffle同时连接时，平均分配到各线程所能使用的Direct buffer Memory将较低（例如，当机器的CPU为40核，Yarn的Shuffle Service启动的线程数为80，80个线程共享进程里的Direct buffer Memory，这种场景下每个线程分配到的内存将不足2MB）。

因此建议根据集群中的NodeManager节点的CPU核数适当调整Direct buffer Memory，例如在CPU核数为40时，将Direct buffer Memory配置为512M。即配置集群NodeManager的“GC_OPTS”参数，如：

```
-XX:MaxDirectMemorySize=512M
```

📖 说明

GC_OPTS参数中-XX:MaxDirectMemorySize默认没有配置，如需配置，用户可在GC_OPTS参数中自定义添加。

具体的配置方法如下：

用户可登录FusionInsight Manager，单击“集群 > 待操作集群的名称 > 服务 > Yarn > 配置”，单击“全部配置”，单击“NodeManager > 系统”，在“GC_OPTS”参数中修改配置。

表 26-78 参数说明

参数	描述	默认值
GC_OPTS	Yarn NodeManager的GC参数。	128M

26.12.1.17 安全集群使用 HiBench 工具运行 sparkbench 获取不到 realm

问题

运行HiBench6的sparkbench任务，如Wordcount，任务执行失败。

“bench.log”中显示Yarn任务执行失败。

登录Yarn WebUI，查看对应application的失败信息，显示如下：

```
Exception in thread "main" org.apache.spark.SparkException: Unable to load YARN support
    at org.apache.spark.deploy.SparkHadoopUtil$.liftedTree$1$1(SparkHadoopUtil.scala:390)
    at org.apache.spark.deploy.SparkHadoopUtil$.yarn$lzycompute(SparkHadoopUtil.scala:385)
    at org.apache.spark.deploy.SparkHadoopUtil$.yarn(SparkHadoopUtil.scala:385)
    at org.apache.spark.deploy.SparkHadoopUtil$.get(SparkHadoopUtil.scala:410)
    at org.apache.spark.deploy.yarn.ApplicationMaster$.main(ApplicationMaster.scala:796)
    at org.apache.spark.deploy.yarn.ExecutorLauncher$.main(ApplicationMaster.scala:821)
    at org.apache.spark.deploy.yarn.ExecutorLauncher.main(ApplicationMaster.scala)
Caused by: java.lang.IllegalArgumentException: Can't get Kerberos realm
    at org.apache.hadoop.security.HadoopKerberosName.setConfiguration(HadoopKerberosName.java:65)
    at org.apache.hadoop.security.UserGroupInformation.initialize(UserGroupInformation.java:288)
    at org.apache.hadoop.security.UserGroupInformation.setConfiguration(UserGroupInformation.java:336)
    at org.apache.spark.deploy.SparkHadoopUtil.<init>(SparkHadoopUtil.scala:51)
    at org.apache.spark.deploy.yarn.YarnSparkHadoopUtil.<init>(YarnSparkHadoopUtil.scala:49)
    at sun.reflect.NativeConstructorAccessorImpl.newInstance0(Native Method)
    at sun.reflect.NativeConstructorAccessorImpl.newInstance(NativeConstructorAccessorImpl.java:62)
    at sun.reflect.DelegatingConstructorAccessorImpl.newInstance(DelegatingConstructorAccessorImpl.java:45)
    at java.lang.reflect.Constructor.newInstance(Constructor.java:423)
```

```
at java.lang.Class.newInstance(Class.java:442)
at org.apache.spark.deploy.SparkHadoopUtil$.liftedTree$1$1(SparkHadoopUtil.scala:387)
... 6 more
Caused by: java.lang.reflect.InvocationTargetException
at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
at java.lang.reflect.Method.invoke(Method.java:498)
at org.apache.hadoop.security.authentication.util.KerberosUtil.getDefaultRealm(KerberosUtil.java:88)
at org.apache.hadoop.security.HadoopKerberosName.setConfiguration(HadoopKerberosName.java:63)
... 16 more
Caused by: KrbException: Cannot locate default realm
at sun.security.krb5.Config.getDefaultRealm(Config.java:1029)
... 22 more
```

回答

失败原因是C80SPC200版本开始，创建集群不再替换/etc/krb5.conf文件，改为通过配置参数指定到客户端内krb5路径，而HiBench并不引用客户端配置文件。

解决方案：

将客户端/opt/client/KrbClient/kerberos/var/krb5kdc/krb5.conf，copy覆盖集群内所有节点的/etc/krb5.conf，注意替换前需要备份。

26.12.2 SQL 和 DataFrame

26.12.2.1 Spark SQL ROLLUP 和 CUBE 使用的注意事项

问题

假设有表src(d1, d2, m)，其数据如下：

```
1 a 1
1 b 1
2 b 2
```

对于语句select d1, sum(d1) from src group by d1, d2 with rollup其结果如下：

```
NULL 0
1 2
2 2
1 1
1 1
2 2
```

对于以上结果的第一条为什么是(NULL,0)而不是(NULL,4)。

回答

在进行rollup和cube操作时，用户通常是基于维度进行分析，需要的是度量的结果，因此不会对维度进行聚合操作。

例如当前有表src(d1, d2, m)，那么语句1 “select d1, sum(m) from src group by d1, d2 with rollup”就是对维度d1和d2进行上卷操作计算度量m的结果，因此有实际业务意义，而其结果也跟预期是一致的。但语句2 “select d1, sum(d1) from src group by d1, d2 with rollup”则从业务上无法解释。当前对于语句2所有聚合（sum/avg/max/min）结果均为0。

说明

只有在rollup和cube操作中对出现在group by中的字段进行聚合结果才是0，非rollup和cube操作其结果跟预期一致。

26.12.2.2 Spark SQL 在不同 DB 都可以显示临时表

问题

切换数据库之后，为什么还能看到之前数据库的临时表？

1. 创建一个DataSource的临时表，例如以下建表语句。

```
create temporary table ds_parquet
using org.apache.spark.sql.parquet
options(path '/tmp/users.parquet');
```

2. 切换到另外一个数据库，执行**show tables**，依然可以看到上个步骤创建的临时表。

```
0: jdbc:hive2://192.168.169.84:22550/default> show tables;
+-----+-----+
| tableName | isTemporary |
+-----+-----+
| ds_parquet | true      |
| cmb_tbl_carbon | false    |
+-----+-----+
2 rows selected (0.109 seconds)
0: jdbc:hive2://192.168.169.84:22550/default>
```

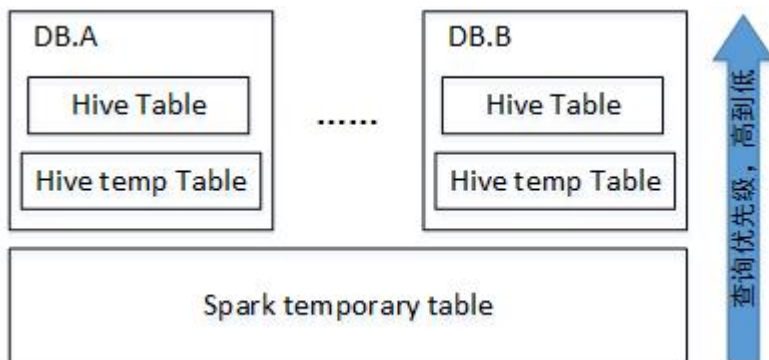
回答

Spark的表管理层次如图26-11所示，最底层是Spark的临时表，存储着使用DataSource方式的临时表，在这一个层面中没有数据库的概念，因此对于这种类型表，表名在各个数据库中都是可见的。

上层为Hive的MetaStore，该层有了各个DB之分。在每个DB中，又有Hive的临时表与Hive的持久化表，因此在Spark中允许三个层次的同名数据表。

查询的时候，Spark SQL优先查看是否有Spark的临时表，再查找当前DB的Hive临时表，最后查找当前DB的Hive持久化表。

图 26-11 Spark 表管理层次



当Session退出时，用户操作相关的临时表将自动删除。建议用户不要手动删除临时表。

删除临时表时，其优先级与查询相同，从高到低为Spark临时表、Hive临时表、Hive持久化表。如果想直接删除Hive表，不删除Spark临时表，您可以直接使用 ***drop table DbName.TableName***命令。

26.12.2.3 如何在 Spark 命令中指定参数值

问题

如果用户不希望在界面上或配置文件设置参数值，如何在Spark命令中指定参数值？

回答

Spark的配置项，不仅可以在配置文件中设置，也可以在命令中指定参数值。

在Spark客户端，应用执行命令添加如下内容设置参数值，命令执行完成后立即生效。在--conf后添加参数名称及其参数值，例如：

```
--conf spark.eventQueue.size=50000
```

26.12.2.4 SparkSQL 建表时的目录权限

问题

新建的用户，使用SparkSQL建表时出现类似如下错误：

```
0: jdbc:hive2://192.168.169.84:22550/default> create table testACL(c string);
Error: org.apache.spark.sql.execution.QueryExecutionException: FAILED: Execution Error, return code 1 from
org.apache.hadoop.hive.ql.exec.DDLTask. MetaException(message:Got exception:
org.apache.hadoop.security.AccessControlException
Permission denied: user=testACL, access=EXECUTE, inode="/user/hive/warehouse/
testacl":spark:hadoop:drwxrwx---
    at
org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkAccessAcl(FSPermissionChecker.java:403
)
    at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:306)
    at
org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkTraverse(FSPermissionChecker.java:259)
    at
org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:20
5)
    at
org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:19
0)
    at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1710)
    at
org.apache.hadoop.hdfs.server.namenode.FSDirStatAndListingOp.getFileInfo(FSDirStatAndListingOp.java:109)
    at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.getFileInfo(FSNamesystem.java:3762)
    at
org.apache.hadoop.hdfs.server.namenode.NameNodeRpcServer.getFileInfo(NameNodeRpcServer.java:1014)
    at
org.apache.hadoop.hdfs.protocolPB.ClientNamenodeProtocolServerSideTranslatorPB.getFileInfo(ClientNamen
odeProtocolServerSideTranslatorPB.java:853)
    at org.apache.hadoop.hdfs.protocol.proto.ClientNamenodeProtocolProtos$ClientNamenodeProtocol
$2.callBlockingMethod(ClientNamenodeProtocolProtos.java)
    at org.apache.hadoop.ipc.ProtobufRpcEngine$Server$ProtoBufRpcInvoker.call(ProtobufRpcEngine.java:616)
    at org.apache.hadoop.ipc.RPC$Server.call(RPC.java:973)
    at org.apache.hadoop.ipc.Server$Handler$1.run(Server.java:2089)
    at org.apache.hadoop.ipc.Server$Handler$1.run(Server.java:2085)
    at java.security.AccessController.doPrivileged(Native Method)
    at javax.security.auth.Subject.doAs(Subject.java:422)
    at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1675)
    at org.apache.hadoop.ipc.Server$Handler.run(Server.java:2083)
) (state=,code=0)
```


回答

Spark SQL建表底层调用的是Hive的接口，其建表时会在“/user/hive/warehouse”目录下新建一个以表名命名的目录，因此要求用户具备“/user/hive/warehouse”目录的读写、执行权限或具有Hive的group权限。

“/user/hive/warehouse”目录可通过hive.metastore.warehouse.dir参数指定。

26.12.2.5 为什么不同服务之间互相删除 UDF 失败

问题

不同服务之间互相删除UDF失败，例如，Spark SQL无法删除Hive创建的UDF。

回答

当前可以通过以下3种方式创建UDF：

1. 在Hive端创建UDF。
2. 通过JDBCServer接口创建UDF。用户可以通过Spark Beeline或者JDBC客户端代码来连接JDBCServer，从而执行SQL命令，创建UDF。
3. 通过spark-sql创建UDF。

删除UDF失败，存在以下两种场景：

- 在Spark Beeline中，对于其他方式创建的UDF，需要重新启动Spark服务端的JDBCServer后，才能将此类UDF删除成功，否则删除失败。在spark-sql中，对于其他方式创建的UDF，需要重新启动spark-sql后，才能将此类UDF删除成功，否则删除失败。
原因：创建UDF后，Spark服务端的JDBCServer未重启或者spark-sql未重新启动的场景，Spark所在线程的FunctionRegistry对象未保存新创建的UDF，那么删除UDF时就会出现错误。
解决方法：重启Spark服务端的JDBCServer和spark-sql，再删除此类UDF。
- 在Hive端创建UDF时未在创建语句中指定jar包路径，而是通过**add jar**命令添加UDF的jar包如**add jar /opt/test/two_udfs.jar**，这种场景下，在其他服务中删除UDF时就会出现ClassNotFound的错误，从而导致删除失败。
原因：在删除UDF时，会先获取该UDF，此时会去加载该UDF对应的类，由于创建UDF时是通过**add jar**命令指定jar包路径的，其他服务进程的classpath不存在这些jar包，因此会出现ClassNotFound的错误从而导致删除失败。
解决方法：该方式创建的UDF不支持通过其他方式删除，只能通过与创建时一致的方式删除。

26.12.2.6 Spark SQL 无法查询到 Parquet 类型的 Hive 表的新插入数据

问题

为什么通过Spark SQL无法查询到存储类型为Parquet的Hive表的新插入数据？主要有以下两种场景存在这个问题：

1. 对于分区表和非分区表，在Hive客户端中执行插入数据的操作后，会出现Spark SQL无法查询到最新插入的数据的问题。

2. 对于分区表，在Spark SQL中执行插入数据的操作后，如果分区信息未改变，会出现Spark SQL无法查询到最新插入的数据的问题。

回答

由于Spark存在一个机制，为了提高性能会缓存Parquet的元数据信息。当通过Hive或其他方式更新了Parquet表时，缓存的元数据信息未更新，导致Spark SQL查询不到新插入的数据。

对于存储类型为Parquet的Hive分区表，在执行插入数据操作后，如果分区信息未改变，则缓存的元数据信息未更新，导致Spark SQL查询不到新插入的数据。

解决措施：在使用Spark SQL查询之前，需执行Refresh操作更新元数据信息。

REFRESH TABLE table_name;

table_name为刷新的表名，该表必须存在，否则会出错。

执行查询语句时，即可获取到最新插入的数据。

26.12.2.7 cache table 使用指导

问题

cache table的作用是什么？cache table时需要注意哪些方面？

回答

Spark SQL可以将表cache到内存中，并且使用压缩存储来尽量减少内存压力。通过将表cache，查询可以直接从内存中读取数据，从而减少读取磁盘带来的内存开销。

但需要注意的是，被cache的表会占用executor的内存。尽管在Spark SQL采用压缩存储的方式来尽量减少内存开销、缓解GC压力，但当缓存的表较大或者缓存表数量较多时，将不可避免地影响executor的稳定性。

此时的最佳实践是，当不需要将表cache来实现查询加速时，应及时将表进行uncache以释放内存。可以执行命令**uncache table table_name**来uncache表。

说明

被cache的表也可以在Spark Driver UI的Storage标签里查看。

26.12.2.8 Repartition 时有部分 Partition 没数据

问题

在repartition操作时，分块数“spark.sql.shuffle.partitions”设置为4500，repartition用到的key列中有超过4000个的不同key值。期望不同key对应的数据能分到不同的partition，实际上却只有2000个partition里有数据，不同key对应的数据也被分到相同的partition里。

回答

这是正常现象。

数据分到哪个partition是通过对key的hashcode取模得到的，不同的hashcode取模后的结果有可能是一样的，那样数据就会被分到相同的partition里面，因此出现有些partition没有数据而有些partition里面有多个key对应的数据。

通过调整“spark.sql.shuffle.partitions”参数值可以调整取模时的基数，改善数据分块不均匀的情况，多次验证发现配置为质数或者奇数效果比较好。

在Driver端的“spark-defaults.conf”配置文件中调整如下参数。

表 26-79 参数说明

参数	描述	默认值
spark.sql.shuffle.partitions	shuffle操作时，shuffle数据的分块数。	200

26.12.2.9 16T 的文本数据转成 4T Parquet 数据失败

问题

使用默认配置时，16T的文本数据转成4T Parquet数据失败，报如下错误信息。

```
Job aborted due to stage failure: Task 2866 in stage 11.0 failed 4 times, most recent failure: Lost task 2866.6 in stage 11.0 (TID 54863, linux-161, 2): java.io.IOException: Failed to connect to /10.16.1.11:23124 at org.apache.spark.network.client.TransportClientFactory.createClient(TransportClientFactory.java:214) at org.apache.spark.network.client.TransportClientFactory.createClient(TransportClientFactory.java:167) at org.apache.spark.network.netty.NettyBlockTransferService$anon$1.createAndStart(NettyBlockTransferService.scala:92)
```

使用的默认配置如表26-80所示。

表 26-80 参数说明

参数	描述	默认值
spark.sql.shuffle.partitions	shuffle操作时，shuffle数据的分块数。	200
spark.shuffle.sasl.timeout	shuffle操作时SASL认证的超时时间。单位：秒。	120s
spark.shuffle.io.connectionTimeout	shuffle操作时连接远程节点的超时时间。单位：秒。	120s
spark.network.timeout	所有涉及网络连接操作的超时时间。单位：秒。	360s

回答

由于当前数据量较大，有16T，而分区数只有200，造成每个task任务过重，才会出现上面的问题。

为了解决上面问题，需要对参数进行调整。

- 增大partition数，把任务切分的更小。
- 增大任务执行过程中的超时时间。

在客户端的“spark-defaults.conf”配置文件中配置如下参数。

表 26-81 参数说明

参数	描述	建议值
spark.sql.shuffle.partitions	shuffle操作时，shuffle数据的分块数。	4501
spark.shuffle.sasl.timeout	shuffle操作时SASL认证的超时时间。单位：秒。	2000s
spark.shuffle.io.connectionTimeout	shuffle操作时连接远程节点的超时时间。单位：秒。	3000s
spark.network.timeout	所有涉及网络连接操作的超时时间。单位：秒。	360s

26.12.2.10 当表名为 table 时，执行相关操作时出现异常

问题

当创建了表名为table的表后，执行**drop table table**上报以下错误。

或者执行其他操作也会出现类似错误。

```
16/07/12 18:56:29 ERROR SparkSQLDriver: Failed in [drop table table]
java.lang.RuntimeException: [1.1] failure: identifier expected
table
^
at scala.sys.package$.error(package.scala:27)
at org.apache.spark.sql.catalyst.SqlParserTrait$class.parseTableIdentifier(SqlParser.scala:56)
at org.apache.spark.sql.catalyst.SqlParser$.parseTableIdentifier(SqlParser.scala:485)
```

回答

这是因为table为Spark SQL的关键词，不能用作表名使用。

建议用户不要使用table用作表的名字。

26.12.2.11 执行 analyze table 语句，因资源不足出现任务卡住

问题

使用spark-sql执行**analyze table**语句，任务一直卡住，打印的信息如下：

```
spark-sql> analyze table hivetable2 compute statistics;
Query ID = root_20160716174218_90f55869-000a-40b4-a908-533f63866fed
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
16/07/20 17:40:56 WARN JobResourceUploader: Hadoop command-line option parsing not performed.
```

```
Implement the Tool interface and execute your application with ToolRunner to remedy this.  
Starting Job = job_1468982600676_0002, Tracking URL = http://10-120-175-107:8088/proxy/  
application_1468982600676_0002/  
Kill Command = /opt/client/HDFS/hadoop/bin/hadoop job -kill job_1468982600676_0002
```

回答

执行 ***analyze table hivetable2 compute statistics*** 语句时，由于该sql语句会启动 MapReduce 任务。从 YARN 的 Resource Manager Web UI 页面看到，该任务由于资源不足导致任务没有被执行，表现出任务卡住的现象。

图 26-12 ResourceManager Web UI 页面

application_	Type	default	Wed Jul 20 17:40:56 +0800 2016	N/A	ACCEPTED	UNDEFINED	0	0	0	ApplicationMaster	0
analyze table hivetable2 compute statistics(Stage-0)	MAPREDUCE										
SparkSQL::192.168.169.84	SPARK		Wed Jul 20 17:40:56	N/A	RUNNING	UNDEFINED	3	3	4096	ApplicationMaster	0

建议用户执行 ***analyze table*** 语句时加上 ***noscan***，其功能与 ***analyze table hivetable2 compute statistics*** 语句相同，具体命令如下：

```
spark-sql> analyze table hivetable2 compute statistics noscan
```

该命令不用启动 MapReduce 任务，不会占用 YARN 资源，从而任务可以被执行。

26.12.2.12 为什么有时访问没有权限的 parquet 表时，在上报 “Missing Privileges” 错误提示之前，会运行一个 Job？

问题

为什么有时访问没有权限的 parquet 表时，在上报 “Missing Privileges” 错误提示之前，会运行一个 Job？

回答

Spark SQL 对用户 SQL 语句的执行逻辑是：首先解析出语句中包含的表，再获取表的元数据信息，然后对权限进行检查。

当表是 parquet 表时，元数据信息包括文件的 Split 信息。Split 信息需要调用 HDFS 的接口去读取，当表包含的文件数量很多时，串行读取 Split 信息变得缓慢，影响性能。故对此做了优化，当表包含的文件大于一定阈值（即 `spark.sql.sources.parallelSplitDiscovery.threshold` 参数值）时，会生成一个 Job，利用 Executor 的并行能力去读取，从而提升执行效率。

由于权限检查在获取表元数据之后，因此当读取的 parquet 表包含的文件数量很多时，会在报 “Missing Privileges” 之前，运行一个 Job 来并行读取元数据信息。

26.12.2.13 spark-sql 退出时打印 RejectedExecutionException 异常栈

问题

执行大数据量的 Spark 任务（如 2T 的 TPCDS 测试套），任务运行成功后，在 spark-sql 退出时概率性出现 RejectedExecutionException 的异常栈信息，相关日志如下所示：

```
16/07/16 10:19:56 ERROR TransportResponseHandler: Still have 2 requests outstanding when connection from linux-192/10.1.1.5:59250 is closed
```

```
java.util.concurrent.RejectedExecutionException: Task scala.concurrent.impl.CallbackRunnable@5fc1ab
rejected from java.util.concurrent.ThreadPoolExecutor@52fa7e19[Terminated, pool size = 0, active threads =
0, queued tasks = 0, completed tasks = 3025]
```

回答

出现上述问题的原因是：当spark-sql退出时，应用退出关闭消息通道，如果当前还有消息未处理，需要做连接关闭异常的处理，此时，如果scala内部的线程池已经关闭，就会打印RejectEdExecutionException的异常栈，如果scala内部的线程池尚未关闭就不会打印该异常栈。

因为该问题出现在应用退出时，此时任务已经运行成功，所以不会对业务产生影响。

26.12.2.14 健康检查时，误将 JDBCServer Kill

问题

健康检查方案中，在并发执行的语句达到线程池上限后依然会导致健康检查命令无法执行，从而导致健康检查程序超时，然后把Spark JDBCServer进程Kill。

回答

当前JDBCServer中存在两个线程池HiveServer2-Handler-Pool和HiveServer2-Background-Pool，其中HiveServer2-Handler-Pool用于处理session连接，HiveServer2-Background-Pool用于处理SQL语句的执行。

当前的健康检查机制是通过新增一个session连接，并在该session所在的线程中执行健康检查命令 **HEALTHCHECK**来判断SparkJDBCServer的健康状况，因此HiveServer2-Handler-Pool必须保留一个线程，用于处理健康检查的session连接和健康检查命令执行，否则将导致无法建立健康检查的session连接或健康检查命令无法执行，从而认为Spark JDBCServer不健康而被Kill。即如果当前HiveServer2-Handler-Pool的线程池数为100，那么最多支持连接99个session。

26.12.2.15 日期类型的字段作为过滤条件时匹配'2016-6-30'时没有查询结果

问题

为什么日期类型的字段作为过滤条件时匹配'2016-6-30'时没有查询结果，匹配'2016-06-30'时有查询结果。

如下图所示：“select count(*)from trxfintrx2012 a where trx_dte_par='2016-6-30'”，其中trx_dte_par为日期类型的字段，当过滤条件为“where trx_dte_par='2016-6-30'”时没有查询结果，当过滤条件为“where trx_dte_par='2016-06-30'”时有查询结果。

图 26-13 示例

```
0: jdbc:hive2://ha-cluster/default> select count(*)
0: jdbc:hive2://ha-cluster/default>   from TRXFINTRX2012 a
0: jdbc:hive2://ha-cluster/default>   where trx_dte_par = '2016-6-30';
+-----+
| _c0 |
+-----+
| 0 |
+-----+
1 row selected (0.498 seconds)
0: jdbc:hive2://ha-cluster/default> select count(*)
0: jdbc:hive2://ha-cluster/default>   from TRXFINTRX2012 a
0: jdbc:hive2://ha-cluster/default>   where trx_dte_par = '2016-06-30';
+-----+
| _c0 |
+-----+
| 8520808 |
+-----+
1 row selected (15.788 seconds)
```

回答

在Spark SQL查询语句中，当查询条件中含有日期格式的字符串时，Spark SQL不会对它做日期格式的检查，就是把它当做普通的字符串进行匹配。以上面的例子为例，如果数据格式为“yyyy-mm-dd”，那么字符串'2016-6-30'就是不正确的数据格式。

26.12.2.16 为什么在启动 spark-beeline 的命令中指定 “--hivevar” 选项无效

问题

为什么在启动spark-beeline的命令中指定 “--hivevar” 选项无效？

在MRS集群启动spark-beeline的命令中如果使用了 “--hivevar <VAR_NAME>=<var_value>” 选项自定义一个变量，在启动spark-beeline时不会报错，但在SQL语句中用到变量<VAR_NAME>时会报无法解析<VAR_NAME>的错误。

举例说明，场景如下：

1. 执行以下命令启动spark-beeline：
spark-beeline --hivevar <VAR_NAME>=<var_value>
2. 启动成功后，在spark-beeline中执行SQL语句，如 “DROP TABLE \$ {VAR_NAME}”，报无法解析VAR_NAME的错误。

回答

MRS集群因新增多session管理功能，Hive的特性 “--hivevar <VAR_NAME>=<var_value>” 在Spark中已不再支持，因此在spark-beeline的启动命令中使用 “--hivevar” 选项无效。

26.12.2.17 执行复杂 SQL 语句时报 “Code of method ... grows beyond 64 KB” 的错误

问题

当执行一个很复杂的SQL语句时，例如有多层语句嵌套，且单层语句中对字段有大量的逻辑处理（如多层嵌套的case when语句），此时执行该语句会报如下所示的错误日志，该错误表明某个方法的代码超出了64KB。

```
java.util.concurrent.ExecutionException: java.lang.Exception: failed to compile:
org.codehaus.janino.JaninoRuntimeException: Code of method "(Lorg/apache/spark/sql/catalyst/expressions/
GeneratedClass$SpecificUnsafeProjection;Lorg/apache/spark/sql/catalyst/InternalRow;)V" of class
"org.apache.spark.sql.catalyst.expressions.GeneratedClass$SpecificUnsafeProjection" grows beyond 64 KB
```

回答

在开启钨丝计划（即tungsten功能）后，Spark对于部分执行计划会使用codegen的方式来生成Java代码，但JDK编译时要求Java代码中的每个函数的长度不能超过64KB。当执行一个很复杂的SQL语句时，例如有多层语句嵌套，且单层语句中对字段有大量的逻辑处理（如多层嵌套的case when语句），这种情况下，通过codegen生成的Java代码中函数的大小就可能会超过64KB，从而导致编译失败。

规避措施：

当出现上述问题时，用户可以通过关闭钨丝计划，关闭使用codegen的方式来生成Java代码的功能，从而确保语句的正常执行。即在客户端的“spark-defaults.conf”配置文件中将“spark.sql.codegen.wholeStage”配置为“false”。

26.12.2.18 在 Beeline/JDBCServer 模式下连续运行 10T 的 TPCDS 测试套会出现内存不足的现象

问题

在Driver内存配置为10G时，Beeline/JDBCServer模式下连续运行10T的TPCDS测试套，会出现因为Driver内存不足导致SQL语句执行失败的现象。

回答

当前在默认配置下，在内存中保留的Job和Stage的UI数据个数为1000个。

当前大集群优化已增加将UI数据溢出到磁盘的优化，其溢出条件是每个Stage中的UI数据大小达到最小阈值5MB。如果每个Stage的task数较小，那么其UI数据大小可能达不到该阈值，从而导致该Stage的UI数据一直缓存在内存中，直到UI数据个数到达保留的上限值（当前默认值为1000个），旧的UI数据才会在内存中被清除。

因此，在将旧的UI数据从内存中清除之前，UI数据会占用大量内存，从而导致执行10T的TPCDS测试套时出现Driver内存不足的现象。

规避措施：

- 根据业务需要，配置合适的需要保留的Job和Stage的UI数据个数，即配置“spark.ui.retainedJobs”和“spark.ui.retainedStages”参数。详细信息请参考[常用参数](#)中的表26-48。
- 如果需要保留的Job和Stage的UI数据个数较多，可通过配置“spark.driver.memory”参数，适当增大Driver的内存。详细信息请参考[常用参数](#)中的表26-45。

26.12.2.19 连上不同的 JDBCServer，function 不能正常使用

问题

场景一：

通过add jar的方式建立永久函数，当Beeline连上不同的JDBCServer或者JDBCServer重启后都需要重新add jar。

图 26-14 场景一异常信息

```

0: jdbc:hive2://192.168.91.247:23040/default> create function al as '
-----+-----+
| result |
-----+-----+
NO rows selected (0.222 seconds)
0: jdbc:hive2://192.168.91.247:23040/default> SELECT test.al(array(1, 2, 3), array(2));
-----+-----+
| _col |
-----+-----+
| true |
-----+-----+
1 row selected (8.282 seconds)
0: jdbc:hive2://192.168.91.247:23040/default> closing: 0: jdbc:hive2://192.168.91.247:24002,192.168.154.81:24002,192.168.8.27:24002;serviceDiscoveryMode=zooKeeper;auth-conf;auth-kereberos;principal=spark/hadoop.hadoop.com@HADOOP.COM;
100-106-121-140:/opt/hadoopclient # ./spark-beeline
it's running the fl spark-beeline, it calls /opt/hadoopclient/spark/spark/bin/beeline
and helps to connect to the jdbcserver automatically
connecting to jdbc:hive2://192.168.91.247:24002,192.168.154.81:24002,192.168.8.27:24002;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=sparkthriftserver;sas
doop.hadoop.conf.HADOOP.COM;
2017-06-15 08:17:55,495 | WARN | Thread-2 | TGT refresh thread time adjusted from : Thu Jun 15 05:59:42 GMT+08:00 2017 to : Thu Jun 15 08:18:55 GMT+08:00 2017
fresh interval (60 seconds) from now. | org.apache.zookeeper.Login$.run(Login.java:177)
2017-06-15 08:17:56,743 | WARN | main | unable to load native-hadoop library for your platform... using builtin-java classes where applicable | org.apache.hadoop
java:62)
2017-06-15 08:17:56,773 | WARN | TGT Renewer for sparkuser@HADOOP.COM | Exception encountered while running the renewal command. Aborting renew thread. ExitCo
| requested option while renewing credentials
| org.apache.hadoop.security.UserGroupInformation$.run(UserGroupInformation.java:946)
Connected to: Spark SQL (version)
Driver: Hive JDBC (version 1.2.1.spark)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 1.2.1.spark by Apache Hive
[INFO] unable to bind key for unsupported operation: backward-delete-word
[INFO] unable to bind key for unsupported operation: backward-delete-word
[INFO] unable to bind key for unsupported operation: down-history
[INFO] unable to bind key for unsupported operation: up-history
[INFO] unable to bind key for unsupported operation: up-history
[INFO] unable to bind key for unsupported operation: down-history
[INFO] unable to bind key for unsupported operation: up-history
[INFO] unable to bind key for unsupported operation: up-history
[INFO] unable to bind key for unsupported operation: down-history
[INFO] unable to bind key for unsupported operation: up-history
[INFO] unable to bind key for unsupported operation: up-history
0: jdbc:hive2://192.168.8.27:23040/default> SELECT test.al(array(1, 2, 3), array(2));
Error: org.apache.spark.sql.AnalysisException: unable to load udf class (state=,code=0)
0: jdbc:hive2://192.168.8.27:23040/default> set role admin;
-----+-----+
| key | value |
-----+-----+
| role admin |
-----+-----+
1 row selected (0.465 seconds)
0: jdbc:hive2://192.168.8.27:23040/default> add jar /home/smartcare-udf-0.0.1-SNAPSHOT.jar;
-----+-----+
| result |
-----+-----+
| 0 |
-----+-----+

```

场景二:

show functions能够查到相应的函数，但是无法使用，这是由于连接上的JDBC节点上没有相应路径的jar包，添加上相应的jar包能够查询成功。

图 26-15 场景二异常信息

```

-----+-----+
| function |
-----+-----+
| stddev_pop |
| stddev_samp |
| str_to_map |
| string |
| struct |
| substr |
| substr_index |
| substr_index |
| sum |
| tan |
| tanh |
| test.al |
| timestamp |
| tinyint |
| to_date |
| to_unix_timestamp |
| to_utc_timestamp |
| translate |
| trim |
| trunc |
| ucase |
| unbase64 |
| unhex |
| unix_timestamp |
| upper |
| var_pop |
| var_samp |
| variance |
| weekofyear |
| when |
| window |
| xpath |
-----+-----+
0: jdbc:hive2://192.168.8.27:22550/default> use test;
-----+-----+
| Result |
-----+-----+
NO rows selected (0.038 seconds)
0: jdbc:hive2://192.168.8.27:22550/default> SELECT test.al(array(1, 2, 3), array(2));
Error: org.apache.spark.sql.AnalysisException: undefined function: 'test.al'. This function is neither a registered temporary function nor a permanen
7 (state=,code=0)
0: jdbc:hive2://192.168.8.27:22550/default> show functions;
-----+-----+
| function |
-----+-----+

```

回答

场景一:

add jar语句只会将jar加载到当前连接的JDBCServer的jarClassLoader，不同JDBCServer不会共用。JDBCServer重启后会创建新的jarClassLoader，所以需要重新add jar。

添加jar包有两种方式：可以在启动spark-sql的时候添加jar包，如`spark-sql --jars /opt/test/two_udfs.jar`；也可在spark-sql启动后再添加jar包，如`add jar /opt/test/two_udfs.jar`。add jar所指定的路径可以是本地路径也可以是HDFS上的路径。

场景二：

show functions会从外部的Catalog获取当前database中所有的function。SQL中使用function时，JDBCServer会加载该function对应的jar。

若jar不存在，则该function无法使用，需要重新执行`add jar`命令。

26.12.2.20 用 add jar 方式创建 function，执行 drop function 时出现问题

问题

- 问题一：
用户没有drop function的权限，能够drop成功。具体场景如下：
 - a. 在FusionInsight Manager页面上添加user1用户，给予用户admin权限，执行下列操作：

```
set role admin;add jar /home/smartcare-udf-0.0.1-SNAPSHOT.jar;create database db4;use db4;create function f11 as 'com.huawei.smartcare.dac.hive.udf.UDFArrayGreaterEqual';create function f12 as 'com.huawei.smartcare.dac.hive.udf.UDFArrayGreaterEqual';
```
 - b. 修改user1用户，取消admin权限，执行下列操作：

```
drop function db4.f11;
```

结果显示drop成功，如图26-16所示。

图 26-16 用户没有权限却 drop 成功结果

```
source /opt/${clientPath}/bigdata_env;/opt/${clientPath}/Spark2x/spark/bin/beeline -u
'jdbc:hive2://10.90.46.60:24002,10.90.46.61:24002,10.90.46.62:24002;serviceDiscoveryMode=zooKeeper;zooKe
eperNamespace=sparkthriftserver2x;saslQop=auth-conf;auth=KERBEROS;principal=spark2x/hadoop.hadoo
p.com@HADOOP.COM;' -e "drop function db4.f11;";
```

- 问题二：
用户drop function成功，show function的时候，function仍然存在。具体场景如下：
 - a. 在FusionInsight Manager页面上添加user1用户，给予用户admin权限，进入spark-beeline执行下列操作：

```
set role admin;create database db2;use db2;add jar /home/smartcare-udf-0.0.1-SNAPSHOT.jar;create function f11 as
'com.huawei.smartcare.dac.hive.udf.UDFArrayGreaterEqual';create function f12 as
'com.huawei.smartcare.dac.hive.udf.UDFArrayGreaterEqual';
```
 - b. 退出后再进入spark-beeline执行下列操作：

```
set role admin;use db2;drop function db2.f11;
```
 - c. 退出后再进入spark-beeline执行下列操作：

```
use db2;show functions;
```结果显示，被drop的function仍然存在，如图26-17所示。

图 26-17 执行 show functions 操作后的结果

```

| datediff
| day
| dayofmonth
| dayofyear
| db2.f11
| db2.f12
| decimal
| decode

```

回答

- 问题根因：**
 上述两个问题是由于多主实例模式或者多租户模式下，使用spark-beeline通过add jar的方式创建function，此function在各个JDBCServer实例之间是不可见的。执行drop function时，如果该session连接的JDBCServer实例不是创建function的JDBCServer实例，则在该session中找不到该function，而且hive默认将“hive.exec.drop.ignorenonexistent”设置为“true”，即当function不存在时，删除function操作不会报错，这样就表现出了用户没有drop function的权限，执行drop时却没有报错，让用户误以为drop成功；但重新起session时又连到创建function的JDBCServer上，因此执行show function，function仍然存在。该行为是hive的社区行为。
- 修改方案：**
 在执行drop function命令之前先执行add jar命令，则该function在有权限的情况下才能drop成功，且drop成功之后不会出现show function仍然存在的现象。

26.12.2.21 Spark2x 无法访问 Spark1.5 创建的 DataSource 表

问题

在Spark2x中访问Spark1.5创建的DataSource表时，报无法获取schema信息，导致无法访问表。

回答

- 原因分析：**
 这是由于Spark2x与Spark1.5存储DataSource表信息的格式不一致导致的。Spark1.5会将schema信息分成多个part，使用path.park.0作为key进行存储，读取时再将各个part都读取出来，重新拼成完整的信息。而Spark2x直接使用相应的key获取对应的信息。这样在Spark2x中去读取Spark1.5创建的DataSource表时，就无法成功读取到key对应的信息，导致解析DataSource表信息失败。
 而在处理Hive格式的表时，Spark2x与Spark1.5的存储方式一致，所以Spark2x可以直接读取Spark1.5创建的表，不存在上述问题。
- 规避措施：**
 Spark2x可以通过创建外表的方式来创建一张指向Spark1.5表实际数据的表，这样可以在Spark2x中读取Spark1.5创建的DataSource表。同时，Spark1.5更新过数据后，Spark2x中访问也能感知到变化，反过来一样。这样即可实现Spark2x对Spark1.5创建的DataSource表的访问。

26.12.2.22 Spark SQL 无法查询到 ORC 类型的 Hive 表的新插入数据

问题

为什么通过Spark SQL无法查询到存储类型为ORC的Hive表的新插入数据？主要有以下两种场景存在这个问题：

- 对于分区表和非分区表，在Hive客户端中执行插入数据的操作后，会出现Spark SQL无法查询到最新插入的数据的问题。
- 对于分区表，在Spark SQL中执行插入数据的操作后，如果分区信息未改变，会出现Spark SQL无法查询到最新插入的数据的问题。

回答

由于Spark存在一个机制，为了提高性能会缓存ORC的元数据信息。当通过Hive或其他方式更新了ORC表时，缓存的元数据信息未更新，导致Spark SQL查询不到新插入的数据。

对于存储类型为ORC的Hive分区表，在执行插入数据操作后，如果分区信息未改变，则缓存的元数据信息未更新，导致Spark SQL查询不到新插入的数据。

解决措施：

1. 在使用Spark SQL查询之前，需执行Refresh操作更新元数据信息：

```
REFRESH TABLE table_name;
```

*table_name*为刷新的表名，该表必须存在，否则会出错。

执行查询语句时，即可获取到最新插入的数据。

2. 使用sqark时，执行以下命令禁用Spark优化：

```
set spark.sql.hive.convertMetastoreOrc=false;
```

26.12.3 Spark Streaming

26.12.3.1 Streaming 任务打印两次相同 DAG 日志

问题

在使用Spark Streaming时，使用以下命令运行程序：

```
spark-submit -master yarn-client --conf spark.logLineage=true --jars $SPARK_HOME/jars/streamingClient/kafka-clients-0.8.2.1.jar,$SPARK_HOME/jars/streamingClient/kafka_2.11-0.8.2.1.jar,$SPARK_HOME/jars/streamingClient/spark-streaming-kafka-0-8_2.11-2.1.0.jar --class com.huawei.bigdata.spark.examples.FemaleInfoCollectionPrint /opt/female/SparkStreamingJavaExample-1.0.jar <checkpoint> <batchTime> <windowTime> <topics> <brokers>
```

在没有Kafka数据输入的情况下，日志中显示的RDD的DAG结构会在一个Batch中打印两次，相关日志如下所示：

```
-----  
Time: 1491447950000 ms  
-----
```

```
17/04/06 11:06:00 INFO SparkContext: RDD's recursive dependencies:  
(2) MapPartitionsRDD[49] at filter at FemaleInfoCollectionPrint.java:111 []  
| MapPartitionsRDD[48] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 []  
| CoGroupedRDD[47] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 []
```

```
| MapPartitionsRDD[38] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 []
|   CachedPartitions: 2; MemorySize: 8.0 B; ExternalBlockStoreSize: 0.0 B; DiskSize: 0.0 B
| ReliableCheckpointRDD[40] at print at FemaleInfoCollectionPrint.java:123 []
| ShuffledRDD[36] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 []
|   CachedPartitions: 2; MemorySize: 8.0 B; ExternalBlockStoreSize: 0.0 B; DiskSize: 0.0 B
+-(5) MapPartitionsRDD[35] at map at FemaleInfoCollectionPrint.java:81 []
|   MapPartitionsRDD[34] at filter at FemaleInfoCollectionPrint.java:81 []
|   MapPartitionsRDD[33] at map at FemaleInfoCollectionPrint.java:72 []
|   MapPartitionsRDD[32] at map at FemaleInfoCollectionPrint.java:63 []
|   KafkaRDD[31] at createDirectStream at FemaleInfoCollectionPrint.java:63 []
| ShuffledRDD[46] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 []
+-(5) MapPartitionsRDD[45] at map at FemaleInfoCollectionPrint.java:81 []
|   MapPartitionsRDD[44] at filter at FemaleInfoCollectionPrint.java:81 []
|   MapPartitionsRDD[43] at map at FemaleInfoCollectionPrint.java:72 []
|   MapPartitionsRDD[42] at map at FemaleInfoCollectionPrint.java:63 []
|   KafkaRDD[41] at createDirectStream at FemaleInfoCollectionPrint.java:63 []
17/04/06 11:06:00 INFO SparkContext: RDD's recursive dependencies: (2) MapPartitionsRDD[48] at
reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 [Memory Serialized 1x Replicated]
|   CachedPartitions: 1; MemorySize: 4.0 B; ExternalBlockStoreSize: 0.0 B; DiskSize: 0.0 B
|   CoGroupedRDD[47] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 [Memory
Serialized 1x Replicated]
|   MapPartitionsRDD[38] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 [Memory
Serialized 1x Replicated]
|   CachedPartitions: 2; MemorySize: 8.0 B; ExternalBlockStoreSize: 0.0 B; DiskSize: 0.0 B
|   ReliableCheckpointRDD[40] at print at FemaleInfoCollectionPrint.java:123 [Memory Serialized 1x
Replicated]
|   ShuffledRDD[36] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 [Memory Serialized
1x Replicated]
|   CachedPartitions: 2; MemorySize: 8.0 B; ExternalBlockStoreSize: 0.0 B; DiskSize: 0.0 B
+-(5) MapPartitionsRDD[35] at map at FemaleInfoCollectionPrint.java:81 [Memory Serialized 1x
Replicated]
|   MapPartitionsRDD[34] at filter at FemaleInfoCollectionPrint.java:81 [Memory Serialized 1x Replicated]
|   MapPartitionsRDD[33] at map at FemaleInfoCollectionPrint.java:72 [Memory Serialized 1x Replicated]
|   MapPartitionsRDD[32] at map at FemaleInfoCollectionPrint.java:63 [Memory Serialized 1x Replicated]
|   KafkaRDD[31] at createDirectStream at FemaleInfoCollectionPrint.java:63 [Memory Serialized 1x
Replicated]
|   ShuffledRDD[46] at reduceByKeyAndWindow at FemaleInfoCollectionPrint.java:98 [Memory Serialized
1x Replicated]
|   CachedPartitions: 1; MemorySize: 4.0 B; ExternalBlockStoreSize: 0.0 B; DiskSize: 0.0 B
+-(5) MapPartitionsRDD[45] at map at FemaleInfoCollectionPrint.java:81 [Memory Serialized 1x
Replicated]
|   MapPartitionsRDD[44] at filter at FemaleInfoCollectionPrint.java:81 [Memory Serialized 1x Replicated]
|   MapPartitionsRDD[43] at map at FemaleInfoCollectionPrint.java:72 [Memory Serialized 1x Replicated]
|   MapPartitionsRDD[42] at map at FemaleInfoCollectionPrint.java:63 [Memory Serialized 1x Replicated]
|   KafkaRDD[41] at createDirectStream at FemaleInfoCollectionPrint.java:63 [Memory Serialized 1x
Replicated]
-----
Time: 1491447960000 ms
-----
```

解答

该应用程序中使用了DStream中的print算子来显示结果，该算子会调用RDD中的take算子来实现底层的计算。

Take算子会以Partition为单位多次触发计算。

在该问题中，由于Shuffle操作，导致take算子默认有两个Partition，Spark首先计算第一个Partition，但由于没有数据输入，导致获取结果不足10个，从而触发第二次计算，因此会出现RDD的DAG结构打印两次的现象。

在代码中将print算子修改为foreach(collect)，该问题则不会出现。

26.12.3.2 Spark Streaming 任务一直阻塞

问题

运行一个Spark Streaming任务，确认有数据输入后，发现没有任何处理的结果。打开Web界面查看Spark Job执行情况，发现如下图所示：有两个Job一直在等待运行，但一直无法成功运行。

图 26-18 Active Jobs

Active Jobs (2)

| Job Id | Description | Submitted | Duration | Stages: Succeeded/Total |
|--------|--|---------------------|----------|-------------------------|
| 3 | print at test2StreamFromKafka.scala:31 | 2015/05/25 18:28:55 | 63.7 h | 0/3 |
| 2 | start at test2StreamFromKafka.scala:34 | 2015/05/25 18:28:55 | 63.7 h | 0/1 |

继续查看已经完成的Job，发现也只有两个，说明Spark Streaming都没有触发数据计算的任务（Spark Streaming默认有两个尝试运行的Job，就是图中两个）

图 26-19 Completed Jobs

Completed Jobs (2)

| Job Id | Description | Submitted | Duration | Stages: Succeeded/Total |
|--------|--|---------------------|----------|-------------------------|
| 1 | print at test2StreamFromKafka.scala:31 | 2015/05/25 18:28:55 | 0.7 s | 2/2 (1 skipped) |
| 0 | start at test2StreamFromKafka.scala:34 | 2015/05/25 18:28:54 | 1 s | 2/2 |

回答

经过定位发现，导致这个问题的原因是：Spark Streaming的计算核数少于Receiver的个数，导致部分Receiver启动以后，系统已经没有资源去运行计算任务，导致第一个任务一直在等待，后续任务一直在排队。从现象上看，就是如问题中的图26-18中所示，会有两个任务一直在等待。

因此，当Web出现两个任务一直在等待的情况，首先检查Spark的核数是否大于Receiver的个数。

说明

Receiver在Spark Streaming中是一个常驻的Spark Job，Receiver对于Spark是一个普通的任务，但它的生命周期和Spark Streaming任务相同，并且占用一个核的计算资源。

在调试和测试等经常使用默认配置的场景下，要时刻注意核数与Receiver个数的关系。

26.12.3.3 运行 Spark Streaming 任务参数调优的注意事项

问题

运行Spark Streaming任务时，随着executor个数的增长，数据处理性能没有明显提升，对于参数调优有哪些注意事项？

回答

在executor核数等于1的情况下，遵循以下规则对调优Spark Streaming运行参数有所帮助。

- Spark任务处理速度和Kafka上partition个数有关，当partition个数小于给定executor个数时，实际使用的executor个数和partition个数相同，其余的将会被空闲。所以应该使得executor个数小于或者等于partition个数。
- 当Kafka上不同partition数据有倾斜时，数据较多的partition对应的executor将成为数据处理的瓶颈，所以在执行Producer程序时，数据平均发送到每个partition可以提升处理的速度。
- 在partition数据均匀分布的情况下，同时提高partition和executor个数，将会提升Spark处理速度（当partition个数和executor个数保持一致时，处理速度是最快的）。
- 在partition数据均匀分布的情况下，尽量保持partition个数是executor个数的整数倍，这样将会使资源得到合理利用。

26.12.3.4 为什么提交 Spark Streaming 应用超过 token 有效期，应用失败

问题

修改kerberos的票据和HDFS token过期时间为5分钟，设置“dfs.namenode.delegation.token.renew-interval”小于60秒，提交Spark Streaming应用，超过token有效期，提示以下错误，应用失败。

```
token (HDFS_DELEGATION_TOKEN token 17410 for spark2x) is expired
```

回答

- 问题原因：
ApplicationMaster进程中有1个Credential Refresh Thread会根据 $token\ renew/周期 * 0.75$ 的时间比例上传更新后的Credential文件到HDFS上。
Executor进程中有1个Credential Refresh Thread会根据 $token\ renew/周期 * 0.8$ 的时间比例去HDFS上获取更新后的Credential文件，用来刷新UserGroupInformation中的token，避免token失效。
当Executor进程的Credential Refresh Thread发现当前时间已经超过Credential文件更新时间（即 $token\ renew/周期 * 0.8$ ）时，会等待1分钟再去HDFS上面获取最新的Credential文件，以确保AM端已经将更新后的Credential文件放到HDFS上。
当“dfs.namenode.delegation.token.renew-interval”配置值小于60秒，Executor进程起来时发现当前时间已经超过Credential文件更新时间，等待1分钟再去HDFS上面获取最新的Credential文件，而此时token已经失效，task运行失败，然后在其他Executor上重试，由于重试时间都是在1分钟内完成，所以task在其他Executor上也运行失败，导致运行失败的Executor加入到黑名单，没有可用的Executor，应用退出。
- 修改方案：
在Spark使用场景下，需设置“dfs.namenode.delegation.token.renew-interval”大于80秒。“dfs.namenode.delegation.token.renew-interval”参数描述请参[表 26-82](#)考。

表 26-82 参数说明

| 参数 | 描述 | 默认值 |
|--|--------------------------------------|----------|
| dfs.namenode.delegation.token.renew-interval | 该参数为服务器端参数，设置token renew的时间间隔，单位为毫秒。 | 86400000 |

26.12.3.5 为什么 Spark Streaming 应用创建输入流，但该输入流无输出逻辑时，应用从 checkpoint 恢复启动失败

问题

Spark Streaming应用创建1个输入流，但该输入流无输出逻辑。应用从checkpoint恢复启动失败，报错如下：

```
17/04/24 10:13:57 ERROR Utils: Exception encountered
java.lang.NullPointerException
at org.apache.spark.streaming.dstream.DStreamCheckpointData$$anonfun$writeObject$1.apply$mcV$sp(DStreamCheckpointData.scala:125)
at org.apache.spark.streaming.dstream.DStreamCheckpointData$$anonfun$writeObject$1.apply(DStreamCheckpointData.scala:123)
at org.apache.spark.streaming.dstream.DStreamCheckpointData$$anonfun$writeObject$1.apply(DStreamCheckpointData.scala:123)
at org.apache.spark.util.Utils$.tryOrIOException(Utils.scala:1195)
at
org.apache.spark.streaming.dstream.DStreamCheckpointData.writeObject(DStreamCheckpointData.scala:123)
at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
at java.lang.reflect.Method.invoke(Method.java:498)
at java.io.ObjectStreamClass.invokeWriteObject(ObjectStreamClass.java:1028)
at java.io.ObjectOutputStream.writeSerialData(ObjectOutputStream.java:1496)
at java.io.ObjectOutputStream.writeOrdinaryObject(ObjectOutputStream.java:1432)
at java.io.ObjectOutputStream.writeObject0(ObjectOutputStream.java:1178)
at java.io.ObjectOutputStream.defaultWriteFields(ObjectOutputStream.java:1548)
at java.io.ObjectOutputStream.defaultWriteObject(ObjectOutputStream.java:441)
at org.apache.spark.streaming.dstream.DStream$$anonfun$writeObject$1.apply$mcV$sp(DStream.scala:515)
at org.apache.spark.streaming.dstream.DStream$$anonfun$writeObject$1.apply(DStream.scala:510)
at org.apache.spark.streaming.dstream.DStream$$anonfun$writeObject$1.apply(DStream.scala:510)
at org.apache.spark.util.Utils$.tryOrIOException(Utils.scala:1195)
at org.apache.spark.streaming.dstream.DStream.writeObject(DStream.scala:510)
at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
at java.lang.reflect.Method.invoke(Method.java:498)
at java.io.ObjectStreamClass.invokeWriteObject(ObjectStreamClass.java:1028)
at java.io.ObjectOutputStream.writeSerialData(ObjectOutputStream.java:1496)
at java.io.ObjectOutputStream.writeOrdinaryObject(ObjectOutputStream.java:1432)
at java.io.ObjectOutputStream.writeObject0(ObjectOutputStream.java:1178)
at java.io.ObjectOutputStream.writeArray(ObjectOutputStream.java:1378)
at java.io.ObjectOutputStream.writeObject0(ObjectOutputStream.java:1174)
at java.io.ObjectOutputStream.defaultWriteFields(ObjectOutputStream.java:1548)
at java.io.ObjectOutputStream.writeSerialData(ObjectOutputStream.java:1509)
at java.io.ObjectOutputStream.writeOrdinaryObject(ObjectOutputStream.java:1432)
at java.io.ObjectOutputStream.writeObject0(ObjectOutputStream.java:1178)
at java.io.ObjectOutputStream.defaultWriteFields(ObjectOutputStream.java:1548)
at java.io.ObjectOutputStream.defaultWriteObject(ObjectOutputStream.java:441)
at org.apache.spark.streaming.DStreamGraph$$anonfun$writeObject$1.apply$mcV$sp(DStreamGraph.scala:191)
at org.apache.spark.streaming.DStreamGraph$$anonfun$writeObject$1.apply(DStreamGraph.scala:186)
```



```
at org.apache.spark.streaming.DStreamGraph$$anonfun$writeObject$1.apply(DStreamGraph.scala:186)
at org.apache.spark.util.Utils$.tryOrIOException(Utils.scala:1195)
at org.apache.spark.streaming.DStreamGraph.writeObject(DStreamGraph.scala:186)
at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
at java.lang.reflect.Method.invoke(Method.java:498)
at java.io.ObjectStreamClass.invokeWriteObject(ObjectStreamClass.java:1028)
at java.io.ObjectOutputStream.writeSerialData(ObjectOutputStream.java:1496)
at java.io.ObjectOutputStream.writeOrdinaryObject(ObjectOutputStream.java:1432)
at java.io.ObjectOutputStream.writeObject0(ObjectOutputStream.java:1178)
at java.io.ObjectOutputStream.defaultWriteFields(ObjectOutputStream.java:1548)
at java.io.ObjectOutputStream.writeSerialData(ObjectOutputStream.java:1509)
at java.io.ObjectOutputStream.writeOrdinaryObject(ObjectOutputStream.java:1432)
at java.io.ObjectOutputStream.writeObject0(ObjectOutputStream.java:1178)
at java.io.ObjectOutputStream.writeObject(ObjectOutputStream.java:348)
at org.apache.spark.streaming.Checkpoint$$anonfun$serialize$1.apply$mcV$sp(Checkpoint.scala:142)
at org.apache.spark.streaming.Checkpoint$$anonfun$serialize$1.apply(Checkpoint.scala:142)
at org.apache.spark.streaming.Checkpoint$$anonfun$serialize$1.apply(Checkpoint.scala:142)
at org.apache.spark.util.Utils$.tryWithSafeFinally(Utils.scala:1230)
at org.apache.spark.streaming.Checkpoint$.serialize(Checkpoint.scala:143)
at org.apache.spark.streaming.StreamingContext.validate(StreamingContext.scala:566)
at org.apache.spark.streaming.StreamingContext.liftedTree1$1(StreamingContext.scala:612)
at org.apache.spark.streaming.StreamingContext.start(StreamingContext.scala:611)
at com.spark.test.kafka08LifoTwoInkfk$.main(kafka08LifoTwoInkfk.scala:21)
at com.spark.test.kafka08LifoTwoInkfk.main(kafka08LifoTwoInkfk.scala)
at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
at java.lang.reflect.Method.invoke(Method.java:498)
at org.apache.spark.deploy.SparkSubmit$.org$apache$spark$deploy$SparkSubmit$
$runMain(SparkSubmit.scala:772)
at org.apache.spark.deploy.SparkSubmit$.doRunMain$1(SparkSubmit.scala:183)
at org.apache.spark.deploy.SparkSubmit$.submit(SparkSubmit.scala:208)
at org.apache.spark.deploy.SparkSubmit$.main(SparkSubmit.scala:123)
at org.apache.spark.deploy.SparkSubmit.main(SparkSubmit.scala)
```

回答

Streaming Context启动时，若应用设置了checkpoint，则需要对应用中的DStream checkpoint对象进行序列化，序列化时会用到dstream.context。

dstream.context是Streaming Context启动时从output Streams反向查找所依赖的DStream，逐个设置context。若Spark Streaming应用创建1个输入流，但该输入流无输出逻辑时，则不会给它设置context。所以在序列化时报“NullPointerException”。

解决办法：应用中如果有无输出逻辑的输入流，则在代码中删除该输入流，或添加该输入流的相关输出逻辑。

26.12.3.6 Spark Streaming 应用运行过程中重启 Kafka，Web UI 界面部分 batch time 对应 Input Size 为 0 records

问题

在Spark Streaming应用执行过程中重启Kafka时，应用无法从Kafka获取topic offset，从而导致生成Job失败。如图26-20所示，其中2017/05/11 10:57:00~2017/05/11 10:58:00为Kafka重启时间段。2017/05/11 10:58:00重启成功后对应的“Input Size”的值显示为“0 records”。

图 26-20 Web UI 界面部分 batch time 对应 Input Size 为 0 records

| Completed Batches (last 9 out of 9) | | | | | |
|-------------------------------------|------------|----------------------|---------------------|-----------------|-----------------------------|
| Batch Time | Input Size | Scheduling Delay (?) | Processing Time (?) | Total Delay (?) | Output Ops: Succeeded/Total |
| 2017/05/11 10:58:50 | 18 records | 0 ms | 0.4 s | 0.4 s | 1/1 |
| 2017/05/11 10:58:40 | 20 records | 4 s | 0.3 s | 4 s | 1/1 |
| 2017/05/11 10:58:30 | 20 records | 14 s | 0.5 s | 14 s | 1/1 |
| 2017/05/11 10:58:20 | 20 records | 23 s | 0.4 s | 24 s | 1/1 |
| 2017/05/11 10:58:10 | 20 records | 33 s | 0.5 s | 33 s | 1/1 |
| 2017/05/11 10:58:00 | 0 records | 6 ms | 43 s | 43 s | 1/1 |
| 2017/05/11 10:57:00 | 19 records | 1 ms | 0.9 s | 0.9 s | 1/1 |
| 2017/05/11 10:56:50 | 20 records | 1 ms | 0.6 s | 0.6 s | 1/1 |
| 2017/05/11 10:56:40 | 28 records | 13 ms | 5 s | 5 s | 1/1 |

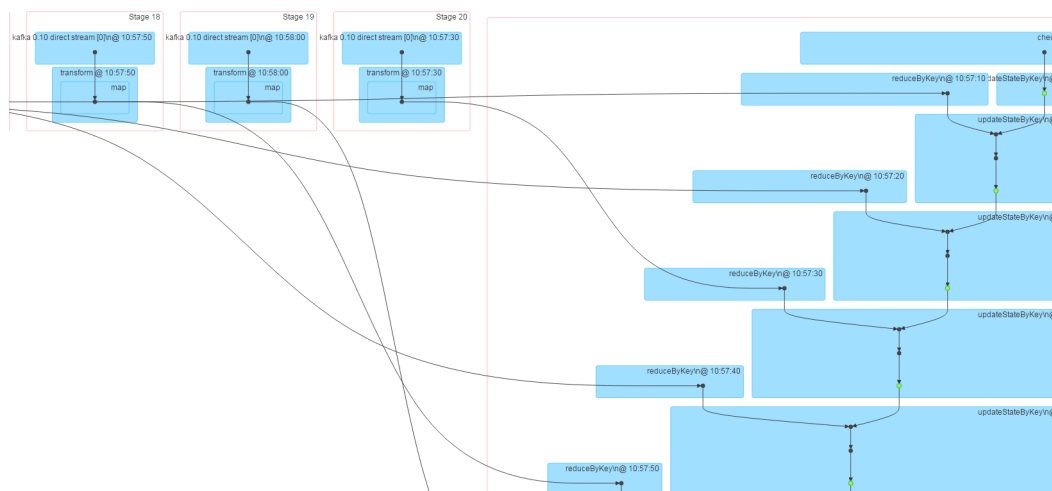
回答

Kafka重启成功后应用会按照batch时间把2017/05/11 10:57:00~2017/05/11 10:58:00 缺失的RDD补上（如图26-21所示），尽管UI界面上显示读取的数据个数为“0”，但实际上这部分数据在补的RDD中进行了处理，因此，不存在数据丢失。

Kafka重启时间段的数据处理机制如下。

Spark Streaming应用使用了state函数（例如：updateStateByKey），在Kafka重启成功后，Spark Streaming应用生成2017/05/11 10:58:00 batch任务时，会按照batch时间把2017/05/11 10:57:00~2017/05/11 10:58:00缺失的RDD补上（Kafka重启前Kafka上未读取完的数据，属于2017/05/11 10:57:00之前的batch），如图26-21所示。

图 26-21 重启时间段缺失数据处理机制



26.12.4 访问 Spark 应用获取的 restful 接口信息有误

问题

当Spark应用结束后，访问该应用的restful接口获取job信息，发现job信息中“numActiveTasks”的值是负数，如图26-22所示。

图 26-22 job 信息

```
[ {  
  "jobId" : 0,  
  "name" : "reduce at SparkPi.scala:36",  
  "submissionTime" : "2016-05-28T09:35:34.415GMT",  
  "completionTime" : "2016-05-28T09:35:35.686GMT",  
  "stageIds" : [ 0 ],  
  "status" : "SUCCEEDED",  
  "numTasks" : 2,  
  "numActiveTasks" : -1,  
  "numCompletedTasks" : 2,  
  "numSkippedTasks" : 2,  
  "numFailedTasks" : 0,  
  "numActiveStages" : 0,  
  "numCompletedStages" : 1,  
  "numSkippedStages" : 0,  
  "numFailedStages" : 0  
} ]
```

说明

numActiveTasks是指当前正在运行task的个数。

回答

通过下面两种途径获取上面的job信息：

- 配置spark.history.briefInfo.gather=true，查看JobHistory的brief信息。
- 使用Spark JobHistory2x页面访问：<https://IP:port/api/v1/<appid>/jobs/>。

job信息中“numActiveTasks”的值是根据eventlog文件中SparkListenerTaskStart和SparkListenerTaskEnd事件的个数的差值计算得到的。如果eventLog文件中有事件丢失，就可能出现上面的现象。

26.12.5 为什么从 Yarn Web UI 页面无法跳转到 Spark Web UI 界面

问题

FusionInsight版本中，在客户端采用yarn-client模式运行Spark应用，然后从Yarn的页面打开该应用的Web UI界面，出现下面的错误：

Error Occurred.

Problem accessing /proxy/application_ /

Powered by Jetty://

从YARN ResourceManager的日志看到：

```
2016-07-21 16:35:27,099 | INFO | Socket Reader #1 for port 8032 | Auth successful for mapred/hadoop.<系统域名>@<系统域名> (auth:KERBEROS) | Server.java:1388  
2016-07-21 16:35:27,105 | INFO | 1526016381@qtp-1178290888-1015 | admin is accessing unchecked  
http://10.120.169.53:23011 which is the app master GUI of  
application_1468986660719_0045 owned by spark | WebAppProxyServlet.java:393
```

```
2016-07-21 16:36:02,843 | INFO | Socket Reader #1 for port 8032 | Auth successful for hive/hadoop.<系统域名>@<系统域名> (auth:KERBEROS) | Server.java:1388
2016-07-21 16:36:02,851 | INFO | Socket Reader #1 for port 8032 | Auth successful for hive/hadoop.<系统域名>@<系统域名> (auth:KERBEROS) | Server.java:1388
2016-07-21 16:36:12,163 | WARN | 1526016381@qtp-1178290888-1015 | /proxy/application_1468986660719_0045/: java.net.ConnectException: Connection timed out | Slf4jLog.java:76
2016-07-21 16:37:03,918 | INFO | Socket Reader #1 for port 8032 | Auth successful for hive/hadoop.<系统域名>@<系统域名> (auth:KERBEROS) | Server.java:1388
2016-07-21 16:37:03,926 | INFO | Socket Reader #1 for port 8032 | Auth successful for hive/hadoop.<系统域名>@<系统域名> (auth:KERBEROS) | Server.java:1388
2016-07-21 16:37:11,956 | INFO | AsyncDispatcher event handler | Updating application attempt appattempt_1468986660719_0045_000001 with final state: FINISHING, and exit status: -1000 | RMAAppAttemptImpl.java:1253
```

回答

打开FusionInsight Manager页面，看到Yarn服务的业务IP地址为192网段。

从Yarn的日志看到，Yarn读取的Spark Web UI地址为http://10.120.169.53:23011，是10网段的IP地址。由于192网段的IP和10网段的IP不能互通，所以导致访问Spark Web UI界面失败。

修改方案：

登录10.120.169.53客户端机器，修改/etc/hosts文件，将10.120.169.53更改为相对应的192网段的IP地址。再重新运行Spark应用，这时就可以打开Spark Web UI界面。

26.12.6 HistoryServer 缓存的应用被回收，导致此类应用页面访问时出错

问题

在History Server页面中访问某个Spark应用的页面时，发现访问时出错。

查看相应的HistoryServer日志后，发现有“FileNotFound”异常，相关日志如下所示：

```
2016-11-22 23:58:03,694 | WARN | [qtp55429210-232] | /history/application_1479662594976_0001/stages/stage/ | org.sparkproject.jetty.servlet.ServletHandler.doHandle(ServletHandler.java:628)
java.io.FileNotFoundException: ${BIGDATA_HOME}/tmp/spark/jobHistoryTemp/blockmgr-5f1f6aca-2303-4290-9845-88fa94d78480/09/temp_shuffle_11f82aaf-e226-46dc-b1f0-002751557694 (No such file or directory)
```

回答

在History Server页面加载Task个数较多的Spark应用时，由于无法把全部的数据放入内存中，导致数据溢出到磁盘时，会产生前缀为“temp_shuffle”的文件。

HistoryServer默认会缓存50个Spark应用（由配置项“spark.history.retainedApplications”决定），当内存中的Spark应用个数超过这个数值时，HistoryServer会回收最先缓存的Spark应用，同时会清理掉相应的“temp_shuffle”文件。

当用户正在查看即将被回收的Spark应用时，可能会出现找不到“temp_shuffle”文件的错误，从而导致当前页面无法访问。

如果遇到上述问题，可参考以下两种方法解决。

- 重新访问这个Spark应用的HistoryServer页面，即可查看到正确的页面信息。

- 如果用户场景需要同时访问50个以上的Spark应用时，需要调大“spark.history.retainedApplications”参数的值。

请登录FusionInsight Manager管理界面，单击“集群 > 待操作集群的名称 > 服务 > Spark2x > 配置”，单击“全部配置”，在左侧的导航列表中，单击“JobHistory2x > 界面”，配置如下参数。

表 26-83 参数说明

| 参数 | 描述 | 默认值 |
|------------------------------------|--|-----|
| spark.history.retainedApplications | HistoryServer缓存的Spark应用数，当需要缓存的应用个数超过此参数值时，HistoryServer会回收最先缓存的Spark应用。 | 50 |

26.12.7 加载空的 part 文件时，app 无法显示在 JobHistory 的页面上

问题

在分组模式下执行应用，当HDFS上的part文件为空时，发现JobHistory首页面上不显示该part对应的app。

回答

JobHistory服务更新页面上的app时，会根据HDFS上的part文件大小变更与否判断是否刷新首页面的app显示信息。若文件为第一次查看，则将当前文件大小与0作比较，如果大于0则读取该文件。

分组的情况下，如果执行的app没有job处于执行状态，则part文件为空，即JobHistory服务不会读取该文件，此app也不会显示在JobHistory页面上。但若part文件大小之后有更新，JobHistory又会显示该app。

26.12.8 Spark2x 导出带有相同字段名的表，结果导出失败

问题

在Spark2x的spark-shell上执行如下语句失败：

```
val acctId = List(("49562", "Amal", "Derry"), ("00000", "Fred", "Xanadu"))
val rddLeft = sc.makeRDD(acctId)
val dfLeft = rddLeft.toDF("Id", "Name", "City")
//dfLeft.show

val acctCustId = List(("Amal", "49562", "CO"), ("Dave", "99999", "ZZ"))
val rddRight = sc.makeRDD(acctCustId)
val dfRight = rddRight.toDF("Name", "CustId", "State")
//dfRight.show
```

```
val dfJoin = dfLeft.join(dfRight, dfLeft("Id") === dfRight("CustId"), "outer")  
dfJoin.show  
  
dfJoin.repartition(1).write.format("com.databricks.spark.csv").option("delimiter", "\t").option("header", "true").option("treatEmptyValuesAsNulls", "true").option("nullValue", "").save("/tmp/outputDir")
```

回答

Spark2x中对join语句重名字段做了判断，需要修改代码保证保存的数据中无重复字段。

26.12.9 为什么多次运行 Spark 应用程序会引发致命 JRE 错误

问题

为什么多次运行Spark应用程序会引发致命JRE错误？

回答

多次运行Spark应用程序会引发致命的JRE错误，这个错误由Linux内核导致。

升级内核版本到4.13.9-2.ge7d7106-default来解决这个问题。

26.12.10 IE 浏览器访问 Spark2x 原生 UI 界面失败，无法显示此页或者页面显示错误

问题

通过IE 9、IE 10和IE 11浏览器访问Spark2x的原生UI界面，出现访问失败情况或者页面显示错误问题。

现象

访问页面失败，浏览器无法显示此页，如下图所示：



无法显示此页

在高级设置中启用 SSL 3.0、TLS 1.0、TLS 1.1 和 TLS 1.2，然后尝试再次连接

原因

IE 9、IE 10、IE 11浏览器的某些版本在处理SSL握手有问题导致访问失败。

解决方法

推荐使用Google Chrome浏览器71及其以上版本。

26.12.11 Spark2x 如何访问外部集群组件

问题

存在两个集群：cluster1 和cluster2，如何使用cluster1中的Spark2x访问cluster2中的HDFS、Hive、HBase和Kafka组件。

回答

1. 可以有条件的实现两个集群间组件互相访问，但是存在以下限制：
 - 仅允许访问一个Hive MetaStore，不支持同时访问cluster1的Hive MetaStore和cluster2的Hive MetaStore。
 - 不同集群的用户系统没有同步，因此访问跨集群组件时，用户的权限管理由对端集群的用户配置决定。比如cluster1的userA没有访问本集群HBase meta表权限，但是cluster2的userA有访问该集群HBase meta表权限，则cluster1的userA可以访问cluster2的HBase meta表。
 - 跨Manager之间的安全集群间组件互相访问，需要先配置系统互信。
2. 以下分别阐述cluster1上使用userA访问cluster2的Hive、HBase、Kafka组件。

说明

以下操作皆以用户使用FusionInsight客户端提交Spark2x应用为基础，若用户使用了自己的配置文件目录，则需要修改本应用配置目录中的对应文件，并注意需要将配置文件上传到executor端。

由于hdfs和hbase客户端访问服务端时，使用hostname配置服务端地址，因此，客户端的/etc/hosts需要保存有所有需要访问节点的hosts配置。用户可预先将对端集群节点的host添加到客户端节点的/etc/hosts文件中。

- 访问Hive MetaStore：使用cluster2中的Spark2x客户端下“conf”目录下的hive-site.xml文件，替换到cluster1中的Spark2x客户端下“conf”目录下的hive-site.xml文件。

如上操作后可以用sparksql访问hive MetaStore，如需访问hive表数据，需要按照[同时访问两个集群的HDFS](#)的操作步骤配置且指定对端集群nameservice为LOCATION后才能访问表数据。

- 访问对端集群的HBase：
 - i. 先将cluster2集群的所有Zookeeper节点和HBase节点的IP和主机名配置到cluster1集群的客户端节点的/etc/hosts文件中。
 - ii. 使用cluster2中的Spark2x客户端下“conf”目录的hbase-site.xml文件，替换到cluster1中的Spark2x客户端下“conf”目录hbase-site.xml文件。
- 访问Kafka，仅需将应用访问的Kafka Broker地址设置为cluster2中的Kafka Broker地址即可。
- 同时访问两个集群的HDFS：
 - 无法同时获取两个相同nameservice的token，因此两个HDFS的nameservice必须不同，例如：一个为hacluster，一个为test
 - 1) 从cluster2的hdfs-site.xml中获取以下配置，添加到cluster1的spark2x客户端conf目录的hdfs-site.xml中
dfs.nameservices.mappings、dfs.nameservices、
dfs.namenode.rpc-address.test.*、dfs.ha.namenodes.test、
dfs.client.failover.proxy.provider.test

参考样例如下：

```
<property>
<name>dfs.nameservices.mappings</name>
<value>[{"name":"hacluster","roleInstances":["14","15"]},
{"name":"test","roleInstances":["16","17"]}]</value>
</property>
<property>
<name>dfs.nameservices</name>
<value>hacluster,test</value>
</property>
<property>
<name>dfs.namenode.rpc-address.test.16</name>
<value>192.168.0.1:8020</value>
</property>
<property>
<name>dfs.namenode.rpc-address.test.17</name>
<value>192.168.0.2:8020</value>
</property>
<property>
<name>dfs.ha.namenodes.test</name>
<value>16,17</value>
</property>
<property>
<name>dfs.client.failover.proxy.provider.test</name>
<value>org.apache.hadoop.hdfs.server.namenode.ha.ConfiguredFailoverProxyProvider
</value>
</property>
```

- 2) 修改cluster1的spark客户端conf目录下的spark-defaults.conf配置文件中，修改spark.yarn.extra.hadoopFileSystems = hdfs://test，spark.hadoop.hdfs.externalToken.enable = true，如下所示：

```
spark.yarn.extra.hadoopFileSystems = hdfs://test
spark.hadoop.hdfs.externalToken.enable = true
```
 - 3) 应用提交命令中，需要添加--keytab 和 --principal参数，参数配置为cluster1中提交任务的用戶。
 - 4) 使用cluster1的spark客户端提交应用，即可同时访问两个hdfs服务
- 同时访问两个集群的HBase：
- i. 修改cluster1的spark客户端conf目录下的spark-defaults.conf配置文件中，修改spark.hadoop.hbase.externalToken.enable = true，如下所示：

```
spark.hadoop.hbase.externalToken.enable = true
```
 - ii. 用戶访问HBase时，需要使用对应集群的配置文件创建Configuration对象，用于创建Connection对象。
 - iii. MRS集群中支持同时获取多个HBase服务的token，以解决Executor中无法访问HBase的问题，使用方式如下：
假设需要访问本集群的HBase和cluster2的HBase，将cluster2的hbase-site.xml文件放到一个压缩包内，压缩包命名为external_hbase_conf***，提交命令时，使用--archives指定这些压缩包。

26.12.12 对同一目录创建多个外表，可能导致外表查询失败

问题

假设存在数据文件路径“/test_data_path”，用戶userA对该目录创建外表tableA，用戶userB对该目录创建外表tableB，当userB对tableB执行insert操作后，userA将查询tableA失败，出现Permission denied异常。

回答

当userB对tableB执行insert操作后，会在外表数据路径下生成新的数据文件，且文件属组是userB，当userA查询tableA时，会读取外表数据目录下的所有的文件，此时会因没有userB生成的文件的读取权限而查询失败。

实际上，不只是查询场景，还有其他场景也会出现问题。例如：inset overwrite操作将会把此目录下的其他表文件也一起复写。

由于Spark SQL当前的实现机制，如果对此种场景添加检查限制，会存在一致性问题 and 性能问题，因此未对此种场景添加限制，但是用户应避免此种用法，以避免此场景带来的各种问题。

26.12.13 访问 Spark2x JobHistory 中某个应用的原生页面时页面显示错误

问题

提交一个Spark应用，包含单个Job 百万个task。应用结束后，在JobHistory中访问该应用的原生页面，浏览器会等待较长时间才跳转到应用原生页面，若10分钟内无法跳转，则页面会显示Proxy Error信息。

图 26-23 错误信息样例

Proxy Error

```
The proxy server received an invalid response from an upstream server.  
The proxy server could not handle the request GET /Spark2x/JobHistory2x/77/history/application [redacted] /iobs/  
Reason: Error reading from remote server
```

回答

在JobHistory界面中跳转到某个应用的原生页面时，JobHistory需要回放该应用的Event log，若应用包含的事件日志较大，则回放时间较长，浏览器需要较长时间的等待。

当前浏览器访问JobHistory原生页面需经过httpd代理，代理的超时时间是10分钟，因此，如果JobHistory在10分钟内无法完成Event log的解析并返回，httpd会主动向浏览器返回Proxy Error信息。

解决方法

由于当前JobHistory开启了本地磁盘缓存功能，访问应用时，会将应用的Event log的解析结果缓存到本地磁盘中，第二次访问时，能大大加快响应速度。因此，出现此种情况时，仅需稍作等待，重新访问原来的链接即可，此时不会再出现需要长时间等待的现象。

26.12.14 对接 OBS 场景中，spark-beeline 登录后指定 loaction 到 OBS 建表失败

问题

对接OBS ECS/BMS集群，spark-beeline登录后，指定location到OBS建表报错失败。

图 26-24 错误信息

```
de-master2qCKJ:22550/> create database sparkdb location 'obs://800mrs/sparktest/sparkdb';

0.626 seconds)
de-master2qCKJ:22550/> use sparkdb;

0.072 seconds)
de-master2qCKJ:22550/> create table orc (id int,name string) using orc;
Exception: Configuration problem with provider path. (state=,code=0)
```

回答

HDFS上ssl.jceks文件权限不足，导致建表失败。

```
Caused by: org.apache.hadoop.security.AccessControlException: Permission denied: user=root, access=READ, inode="/user/spark2x/jars/8.0.2/ssl.jceks":spark2x:hadoop-rw-----
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:410)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:264)
at com.sun.security.sasl.digest.DigestTransporter.checkPermission(DigestTransporter.java:154)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:194)
at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1957)
at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1941)
at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPathAccess(FSDirectory.java:1694)
at org.apache.hadoop.hdfs.server.namenode.FSBlockLocations.getLocations(FSBlockLocations.java:175)
at org.apache.hadoop.hdfs.server.namenode.FSNameSystem.getBlockLocations(FSNameSystem.java:1990)
at org.apache.hadoop.hdfs.server.namenode.NameNodeRpcServer.getBlockLocations(NameNodeRpcServer.java:742)
at org.apache.hadoop.hdfs.protocol.proto.ClientNameNodeProtocol$ClientNameNodeProtocolTranslatorPB.getBlockLocations(ClientNameNodeProtocol$ClientNameNodeProtocolTranslatorPB.java:445)
at org.apache.hadoop.hdfs.protocol.proto.ClientNameNodeProtocol$ClientNameNodeProtocol$2.callBlockingMethod(ClientNameNodeProtocol$ClientNameNodeProtocol$2.java:101)
at org.apache.hadoop.ipc.ProtobufRpcEngine$Server$ProtobufRpcInvoker.call(ProtobufRpcEngine.java:520)
at org.apache.hadoop.ipc.Server$RpcCall.run(Server.java:985)
at org.apache.hadoop.ipc.Server$RpcCall.run(Server.java:913)
at java.security.AccessController.doPrivileged(Native Method)
at javax.security.auth.Subject.doAs(Subject.java:422)
at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1737)
at org.apache.hadoop.ipc.Server$Handler.run(Server.java:1284)
```

解决方法

1. 使用omm用户登录Spark2x所在节点，执行如下命令：
vi \${BIGDATA_HOME}/FusionInsight_Spark2x_8.1.0.1/install/FusionInsight-Spark2x-*/spark/sbin/fake_prestart.sh
2. 将“eval "\${hdfsCmd}" -chmod 600 "\${InnerHdfsDir}"/ssl.jceks >> "\${PRESTART_LOG}" 2>&1”修改成“eval "\${hdfsCmd}" -chmod 644 "\${InnerHdfsDir}"/ssl.jceks >> "\${PRESTART_LOG}" 2>&1”。
3. 重启SparkResource实例。

26.12.15 Spark shuffle 异常处理

问题

在部分场景Spark shuffle阶段会有如下异常

```
2021-08-18 02:53:08.364 INFO [shuffle-server-0-1] | DIGEST1:Unmatched MACs | javax.security.sasl.unwrap(DigestMD5Base.java:148)
2021-08-18 02:53:08.368 WARN [shuffle-server-0-1] | Exception in connection from /0000000000000000 | org.apache.spark.network.server.TransportChannelHandler.exceptionCaught(TransportChannelHandler.java:87)
io.netty.handler.codec.DecoderException: javax.security.sasl.SaslException: DIGEST-MD5: Out of order sequencing of messages from server. Got: 16 Expected: 14
at io.netty.handler.codec.MessageToMessageDecoder.channelRead(MessageToMessageDecoder.java:98)
at io.netty.channel.AbstractChannelHandlerContext.invokeChannelRead(AbstractChannelHandlerContext.java:379)
at io.netty.channel.AbstractChannelHandlerContext.read(AbstractChannelHandlerContext.java:365)
at io.netty.channel.AbstractChannelHandlerContext.fireChannelRead(AbstractChannelHandlerContext.java:357)
at org.apache.spark.network.util.TransportFrameDecoder.channelRead(TransportFrameDecoder.java:102)
at io.netty.channel.AbstractChannelHandlerContext.invokeChannelRead(AbstractChannelHandlerContext.java:379)
at io.netty.channel.AbstractChannelHandlerContext.read(AbstractChannelHandlerContext.java:365)
at io.netty.channel.DefaultChannelPipeline.readChannel(DefaultChannelPipeline.java:1410)
at io.netty.channel.AbstractChannelHandlerContext.invokeChannelRead(AbstractChannelHandlerContext.java:379)
at io.netty.channel.AbstractChannelHandlerContext.read(AbstractChannelHandlerContext.java:365)
at io.netty.channel.DefaultChannelPipeline.readChannel(DefaultChannelPipeline.java:918)
at io.netty.channel.nio.AbstractNioByteChannel$NioByteUnsafe.read(AbstractNioByteChannel.java:163)
at io.netty.channel.nio.NioEventLoop.processSelectedKey(NioEventLoop.java:714)
at io.netty.channel.nio.NioEventLoop.processSelectedKeysOptimized(NioEventLoop.java:650)
at io.netty.channel.nio.NioEventLoop.run(NioEventLoop.java:576)
at io.netty.util.concurrent.SingleThreadEventExecutor$4.run(SingleThreadEventExecutor.java:989)
at io.netty.util.internal.ThreadExecutorMap$2.run(ThreadExecutorMap.java:74)
at io.netty.util.concurrent.FastThreadLocalRunnable.run(FastThreadLocalRunnable.java:30)
at java.lang.Thread.run(Thread.java:748)
Caused by: javax.security.sasl.SaslException: DIGEST-MD5: Out of order sequencing of messages from server. Got: 16 Expected: 14
at com.sun.security.sasl.digest.DigestMD5Base$DigestTransporter.unwrap(DigestMD5Base.java:148)
at com.sun.security.sasl.digest.DigestMD5Base.unwrap(DigestMD5Base.java:131)
at org.apache.spark.network.sasl.SparkSaslServer.unwrap(SparkSaslServer.java:140)
at org.apache.spark.network.sasl.SaslEncryptionDecryptionHandler.decode(SaslEncryption.java:126)
at org.apache.spark.network.sasl.SaslEncryptionDecryptionHandler.decode(SaslEncryption.java:101)
at io.netty.handler.codec.MessageToMessageDecoder.channelRead(MessageToMessageDecoder.java:80)
at io.netty.handler.codec.MessageToMessageDecoder.channelRead(MessageToMessageDecoder.java:80)
```

解决方法

JDBC应该:

登录FusionInsight Manager管理界面，修改JDBCServer的参数
“spark.authenticate.enableSaslEncryption”值为“false”，并重启对应的实例。

客户端作业：

客户端应用在提交应用的时候，修改spark-defaults.conf配置文件的
“spark.authenticate.enableSaslEncryption”值为“false”。

27 使用 Sqoop

27.1 Sqoop 客户端使用实践

Sqoop是一款开源的工具，主要用于在Hadoop(Hive)与传统的数据库(MySQL、PostgreSQL...)间进行数据的传递，可以将一个关系型数据库（例如：MySQL、Oracle、PostgreSQL等）中的数据导进到Hadoop的HDFS中，也可以将HDFS的数据导进到关系型数据库中。

前提条件

- MRS 3.1.0及之后版本在创建集群时已勾选Sqoop组件。
- 安装客户端，具体请参考[安装客户端（3.x及之后版本）](#)。例如安装目录为“/opt/client”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- 客户端目录/Sqoop/sqoop/lib下已有对应驱动包（例mysql驱动包mysql-connector-java-5.1.47.jar），如果没有请参考[Sqoop1.4.7适配MRS 3.x集群](#)章节中的步骤3下载对应mysql包。

sqoop export（HDFS 到 MySQL）

步骤1 登录客户端所在节点。

步骤2 执行如下命令初始化环境变量。

```
source /opt/client/bigdata_env
```

步骤3 使用sqoop命令操作sqoop客户端。

```
sqoop export --connect jdbc:mysql://10.100.xxx.xxx:3306/test --username root  
--password xxx --table component13 --export-dir hdfs://hacluster/user/hive/  
warehouse/component_test3 --fields-terminated-by ',' -m 1
```

更多参数介绍请参见[Sqoop常用命令及参数介绍](#)。

表 27-1 参数说明

参数	说明
--connect	指定JDBC连接的URL，格式为： jdbc:mysql://MySQL数据库IP地址:MySQL的端口/数据库名称 。
--username	连接MySQL数据库的用户名。
-password	连接MySQL数据库的用户密码。命令中如果携带认证密码信息可能存在安全风险，在执行命令前建议关闭系统的history命令记录功能，避免信息泄露。
-table <table-name>	MySQL中用于存放导出数据的表名称。
-export-dir <dir>	需要导出的Sqoop表所在的HDFS路径。
--fields-terminated-by	指定导出数据的分隔符，与需要导出的HDFS中的数据表中的分隔符保持一致。
-m或-num-mappers <n>	启动n个map来并行导入数据，默认是4个，该值请勿高于集群的最大Map数。
-direct	快速模式，利用了数据库的导入工具，如MySQL的mysqlimport，可以比jdbc连接的方式更为高效的将数据导入到关系数据库中。
-update-key <col-name>	后面接条件列名，通过该参数可以将关系数据库中已经存在的数据进行更新操作，类似于关系数据库中的update操作。
-update-mode <mode>	更新模式，有两个值updateonly和默认的allowinsert，该参数只能在关系数据表里不存在要导入的记录时才能使用，比如要导入的hdfs中有一条id=1的记录，如果在表里已经有一条记录id=2，那么更新会失败。
-input-null-string <null-string>	可选参数，如果没有指定，则字符串null将被使用。
-input-null-non-string <null-string>	可选参数，如果没有指定，则字符串null将被使用。
-staging-table <staging-table-name>	<p>创建一个与导入目标表同样数据结构的表，将所有数据先存放在该表中，然后由该表通过一次事务将结果写入到目标表中。</p> <p>该参数是用来保证在数据导入关系数据库表的过程中的事务安全性，因为在导入的过程中可能会有多个事务，那么一个事务失败会影响到其它事务，比如导入的数据会出现错误或出现重复的记录等情况，那么通过该参数可以避免这种情况。</p>
-clear-staging-table	如果该staging-table非空，则通过该参数可以在运行导入前清除staging-table里的数据。

----结束

sqoop import (MySQL 到 Hive 表)

步骤1 登录客户端所在节点。

步骤2 执行如下命令初始化环境变量。

```
source /opt/client/bigdata_env
```

步骤3 使用sqoop命令操作sqoop客户端。

```
sqoop import --connect jdbc:mysql://10.100.xxx.xxx:3306/test --username root
--password xxx --table component --hive-import --hive-table component_test2
--delete-target-dir --fields-terminated-by "," -m 1 --as-textfile
```

表 27-2 参数说明

参数	说明
--hive-import	表示从关系型数据库中导入数据到MRS Hive中。
--delete-target-dir	若Hive中已存在目标文件，则先删除该文件再导入。
-append	将数据追加到hdfs中已经存在的dataset中。使用该参数，sqoop将把数据先导入到一个临时目录中，然后重新给文件命名到一个正式的目录中，以避免和该目录中已存在的文件重名。
-as-avrodatafile	将数据导入到一个Avro数据文件中。
-as-sequencefile	将数据导入到一个sequence文件中。
-as-textfile	将数据导入到一个普通文本文件中，生成该文本文件后，可以在hive中通过sql语句查询出结果。
-boundary-query <statement>	边界查询，在导入前先通过SQL查询得到一个结果集，然后导入的数据就是该结果集内的数据，格式如： - boundary-query 'select id,creationdate from person where id = 3' ，表示导入的数据为id=3的记录，或者 select min(<split-by>), max(<split-by>) from <table name> 。 注意：查询的字段中不能有数据类型为字符串的字段，否则会报错：java.sql.SQLException: Invalid value for getLong()。
-columns<col,col,col...>	指定要导入的字段值，格式如：-columns id,username
-direct	快速模式，利用了数据库的导入工具，如MySQL的mysqlimport，可以比jdbc连接的方式更为高效的将数据导入到关系数据库中。
-direct-split-size	在使用上面direct直接导入的基础上，对导入的流按字节数分块，特别是使用直连模式从PostgreSQL导入数据时，可以将一个到达设定大小的文件分为几个独立的文件。
-inline-lob-limit	设定大对象数据类型的最大值。

参数	说明
-m或-num-mappers	启动n个map来并行导入数据，默认是4个，该值请勿高于集群的最大Map数。
-query, -e<statement>	从查询结果中导入数据，该参数使用时必须指定-target-dir、-hive-table，在查询语句中一定要有where条件且在where条件中需要包含\$CONDITIONS。 示例：-query 'select * from person where \$CONDITIONS' -target-dir /user/hive/warehouse/person -hive-table person
-split-by<column-name>	表的列名，用来切分工作单元，一般后面跟主键ID。
-table <table-name>	关系数据库表名，数据从该表中获取。
-target-dir <dir>	指定hdfs路径。
-warehouse-dir <dir>	与-target-dir不能同时使用，指定数据导入的存放目录，适用于导入hdfs，不适合导入hive目录。
-where	从关系数据库导入数据时的查询条件，示例：-where 'id = 2'
-z,-compress	压缩参数，默认数据不压缩，通过该参数可以使用gzip压缩算法对数据进行压缩，适用于SequenceFile，text文本文件，和Avro文件。
-compression-codec	Hadoop压缩编码，默认为gzip。
-null-string <null-string>	替换null字符串，如果没有指定，则字符串null将被使用。
-null-non-string<null-string>	替换非String的null字符串，如果没有指定，则字符串null将被使用。
-check-column (col)	增量导入参数，用来作为判断的列名，如id。
-incremental (mode) append 或lastmodified	增量导入参数。 append：追加，比如对大于last-value指定的值之后的记录进行追加导入。 lastmodified：最后的修改时间，追加last-value指定的日期之后的记录。
-last-value (value)	增量导入参数，指定自从上次导入后列的最大值（大于该指定的值），也可以自己设定某一值。

----结束

Sqoop 使用样例

- sqoop import (MySQL到HDFS)
sqoop import --connect jdbc:mysql://10.100.231.134:3306/test --username root --password xxx --query 'SELECT * FROM component where

```
$CONDITIONS and component_id="MRS 1.0_002" --target-dir /tmp/  
component_test --delete-target-dir --fields-terminated-by "," -m 1 --as-  
textfile
```

- sqoop export (obs到MySQL)

```
sqoop export --connect jdbc:mysql://10.100.231.134:3306/test --username  
root --password xxx --table component14 -export-dir obs://obs-file-  
bucket/xx/part-m-00000 --fields-terminated-by ',' -m 1
```
- sqoop import (MySQL到obs)

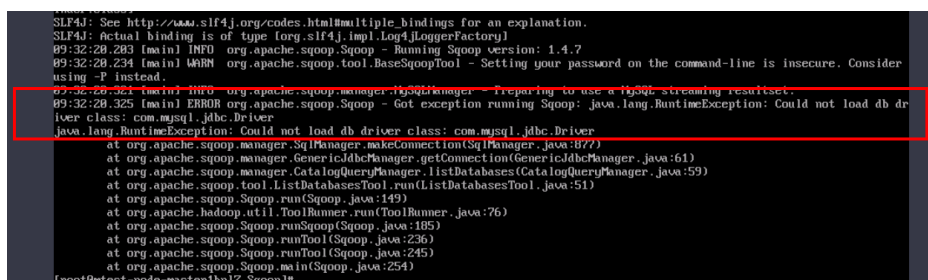
```
sqoop import --connect jdbc:mysql://10.100.231.134:3306/test --username  
root --password xxx --table component --target-dir obs://obs-file-  
bucket/xx --delete-target-dir --fields-terminated-by "," -m 1 --as-textfile
```
- sqoop import (MySQL到Hive外obs表)

```
sqoop import --connect jdbc:mysql://10.100.231.134:3306/test --username  
root --password xxx --table component --hive-import --hive-table  
component_test01 --fields-terminated-by "," -m 1 --as-textfile
```

导入或导出数据时缺少 MySQL 驱动包

若执行sqoop import或sqoop export命令报错“Could not load db driver class: com.mysql.jdbc.Driver”，如图27-1所示，则表示缺少MySQL驱动包，需在MySQL官网下载对应MySQL驱动包，解压并上传至“客户端安装目录/Sqoop/sqoop/lib”目录下，再执行Sqoop导入或导出数据命令即可

图 27-1 缺少 MySQL 驱动包报错



```
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.  
SLF4J: Actual binding is of type org.slf4j.impl.Log4jLoggerFactory  
09:32:28.283 [main] INFO org.apache.sqoop.Sqoop - Running Sqoop version: 1.4.7  
09:32:28.224 [main] WARN org.apache.sqoop.tool.BaseSqoopTool - Setting your password on the command-line is insecure. Consider  
using -P instead.  
09:32:28.321 [main] INFO org.apache.sqoop.manager.HiveManager - Preparing to use a MySQL streaming resultset.  
09:32:28.325 [main] ERROR org.apache.sqoop.Sqoop - Got exception running Sqoop: java.lang.RuntimeException: Could not load db dr  
iver class: com.mysql.jdbc.Driver  
java.lang.RuntimeException: Could not load db driver class: com.mysql.jdbc.Driver  
    at org.apache.sqoop.manager.SqlManager.makeConnection(SqlManager.java:377)  
    at org.apache.sqoop.manager.GenericJdbcManager.getConnection(GenericJdbcManager.java:61)  
    at org.apache.sqoop.manager.CatalogQueryManager.listDatabases(CatalogQueryManager.java:59)  
    at org.apache.sqoop.tool.ListDatabasesTool.run(ListDatabasesTool.java:51)  
    at org.apache.sqoop.Sqoop.run(Sqoop.java:149)  
    at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:76)  
    at org.apache.sqoop.Sqoop.runSqoop(Sqoop.java:185)  
    at org.apache.sqoop.Sqoop.runTool(Sqoop.java:236)  
    at org.apache.sqoop.Sqoop.runTool(Sqoop.java:245)  
    at org.apache.sqoop.Sqoop.main(Sqoop.java:254)  
[root@node-master1b1z2:~]#
```

27.2 Sqoop1.4.7 适配 MRS 3.x 集群

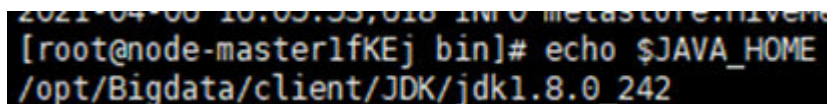
Sqoop是专为Apache Hadoop和结构化数据库（如关系型数据库）设计的高效传输大量数据的工具。客户需要在MRS中使用sqoop进行数据迁移，MRS旧版本中未自带Sqoop，客户可参考此文档自行安装使用。MRS 3.1.0及之后版本已支持创建集群时勾选Sqoop组件，请创建集群时勾选即可。

说明

本章节仅适用于MRS 3.0.2及MRS 3.0.5版本。

前提条件

已安装MRS客户端的节点，且已安装jdk环境。



```
2021-04-08 10:05:55,018 INFO metastore.HiveMetaStore  
[root@node-master1fKEj bin]# echo $JAVA_HOME  
/opt/Bigdata/client/JDK/jdk1.8.0_242
```


Sqoop1.4.7 适配步骤

- 步骤1** 下载开源sqoop-1.4.7.bin__hadoop-2.6.0.tar.gz包（下载地址<http://archive.apache.org/dist/sqoop/1.4.7/>）。
- 步骤2** 将下载的sqoop-1.4.7.bin__hadoop-2.6.0.tar.gz包放入已安装MRS客户端的节点的“/opt/Bigdata/client”目录并解压。
- tar zxvf sqoop-1.4.7.bin__hadoop-2.6.0.tar.gz**
- 步骤3** 从MySQL官网下载MySQL jdbc驱动程序“mysql-connector-java-xxx.jar”，具体MySQL jdbc驱动程序选择参见下表。

表 27-3 版本信息

jdbc驱动程序版本	MySQL版本
Connector/J 5.1	MySQL 4.1、MySQL 5.0、MySQL 5.1、MySQL 6.0 alpha
Connector/J 5.0	MySQL 4.1、MySQL 5.0 servers、distributed transaction (XA)
Connector/J 3.1	MySQL 4.1、MySQL 5.0 servers、MySQL 5.0 except distributed transaction (XA)
Connector/J 3.0	MySQL 3.x、MySQL 4.1

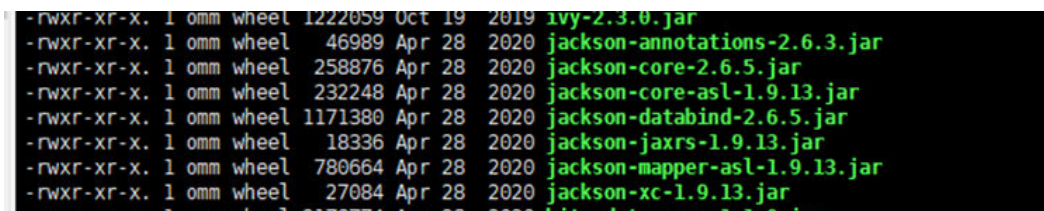
- 步骤4** 将MySQL 驱动包放入Sqoop的lib目录下（/opt/Bigdata/client/sqoop-1.4.7.bin__hadoop-2.6.0/lib）并修改jar包的属组和权限，参考图27-2的omm:wheel 和755的属组和权限。

图 27-2 MySQL 驱动包的属组和权限



- 步骤5** 使用MRS客户端中Hive的lib目录下（/opt/Bigdata/client/Hive/Beeline/lib）的jackson开头的jar包替换Sqoop的lib下的相应jar包。

图 27-3 jackson 开头的 jar



- 步骤6** 将MRS Hive客户端中（/opt/Bigdata/client/Hive/Beeline/lib）的jline的包，拷贝到Sqoop的lib下。
- 步骤7** 执行vim \$JAVA_HOME/jre/lib/security/java.policy增加如下配置：
permission javax.management.MBeanTrustPermission "register";

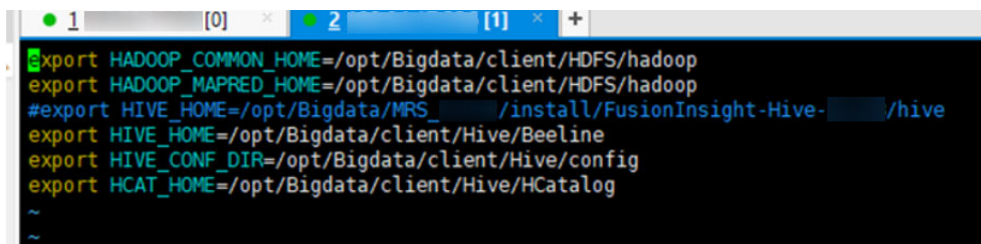
步骤8 执行如下命令，进入Sqoop的conf目录并增加配置：

```
cd /opt/Bigdata/client/sqoop-1.4.7.bin__hadoop-2.6.0/conf
cp sqoop-env-template.sh sqoop-env.sh
```

步骤9 执行vim sqoop-env.sh 设置Sqoop的环境变量，Hadoop、Hive的目录根据实际目录修改。

```
export HADOOP_COMMON_HOME=/opt/Bigdata/client/HDFS/hadoop
export HADOOP_MAPRED_HOME=/opt/Bigdata/client/HDFS/hadoop
export HIVE_HOME=/opt/Bigdata/MRS_1.9.X/install/FusionInsight-Hive-3.1.0/hive(请按照实际路径填写)
export HIVE_CONF_DIR=/opt/Bigdata/client/Hive/config
export HCAT_HOME=/opt/Bigdata/client/Hive/HCatalog
```

图 27-4 设置 Sqoop 的环境变量



步骤10 编写Sqoop脚本 例如：

```
/opt/Bigdata/FusionInsight_Current/1_19_SqoopClient/install/FusionInsight-Sqoop-1.4.7/bin/sqoop import
--connect jdbc:mysql://192.168.0.183:3306/test
--driver com.mysql.jdbc.Driver
--username 'root'
--password 'xxx'
--query "SELECT id, name FROM tbtest WHERE \$CONDITIONS"
--hcatalog-database default
--hcatalog-table test
--num-mappers 1
```

----结束

27.3 Sqoop 常用命令及参数介绍

Sqoop 常用命令介绍

表 27-4 Sqoop 常用命令介绍

命令	说明
import	数据导入到集群
export	集群数据导出
codegen	获取数据库中某张表数据生成Java并打包jar
create-hive-table	创建Hive表
eval	执行sql并查看结果
import-all-tables	导入某个数据库下的所有表到HDFS中
job	生成一个sqoop任务

命令	说明
list-databases	列举数据库名
list-tables	列举表名
merge	将HDFS不同目录下的数据合在一起并存放到指定目录
metastore	启动元数据库，记录sqoop job的元数据
help	打印帮助信息
version	打印版本信息

公用参数介绍

表 27-5 公用参数介绍

分类	参数	说明
连接数据库	--connect	连接关系型数据库的url
	--connection-manager	指定连接管理类
	--driver jdbc	连接驱动包
	--help	帮助信息
	--password	连接数据库密码
	--username	连接数据库的用户名
	--verbose	在控制台打印详细信息
import参数	--fields-terminated-by	设定字段分隔符，和Hive表或hdfs文件保持一致
	--lines-terminated-by	设定行分隔符，和hive表或hdfs文件保持一致
	--mysql-delimiters	MySQL默认分隔符设置
export参数	--input-fields-terminated-by	字段分隔符
	--input-lines-terminated-by	行分隔符
hive参数	--hive-delims-replacement	用自定义的字符替换数据中的\r\n等字符
	--hive-drop-import-delims	在导入数据到hive时，去掉\r\n等字符

分类	参数	说明
	--map-column-hive	生成hive表时可以更改字段的数据类型
	--hive-partition-key	创建分区
	--hive-partition-value	导入数据库指定分区
	--hive-home	指定hive安装目录
	--hive-import	表示操作是从关系型数据库导入到hive中
	--hive-overwrite	覆盖hive已有数据
	--create-hive-table	创建Hive表，默认false，如果目标表不存在，则会创建目标表
	--hive-table	指定hive表
	--table	关系型数据库表名
	--columns	指定需要导入的关系型数据表字段
	--query	指定查询语句，将查询结果导入
hcatalog参数	--hcatalog-database	指定hive库，使用hcatalog方式导入hive库
	--hcatalog-table	指定hive表，使用hcatalog方式导入hive表
其他参数	-m或--num-mappers	后跟数字，表示sqoop任务的分片数
	--split-by	按照某一字段进行分片，配合-m
	--target-dir	指定hdfs临时目录
	--null-string string	类型为null时替换字符串
	--null-non-string	非string类型为null时替换字符串
	--check-column	增量判断的字段
	--incremental append或lastmodified	增量导入参数 append: 追加，比如对大于last-value指定的值之后的记录进行追加导入。 lastmodified: 最后的修改时间，追加last-value指定的日期之后的记录。
	--last-value	指定一个值，用于标记增量导入
	--input-null-string	替换null字符串，如果没有指定，则字符串null将被使用。

分类	参数	说明
	--input-null-non-string	替换非String的null字符串，如果没有指定，则字符串null将被使用。

27.4 Sqoop 常见问题

27.4.1 报错找不到 QueryProvider 类

问题

报错找不到QueryProvider类。

```

2021-04-06 15:57:10,756 INFO manager.SqlManager: Using default fetchSize of 1000
2021-04-06 15:57:10,756 INFO tool.CodeGenTool: Beginning code generation
Apr 06, 2021 3:57:10 PM java.util.logging.LogManager$RootLogger log
SEVERE: Error loading factory org.apache.calcite.jdbc.CalciteJdbc41Factory
java.lang.NoClassDefFoundError: org/apache/calcite/linq4j/QueryProvider
    at java.lang.ClassLoader.defineClass1(Native Method)
    at java.lang.ClassLoader.defineClass(ClassLoader.java:757)
    at java.security.SecureClassLoader.defineClass(SecureClassLoader.java:142)
    at java.net.URLClassLoader.defineClass(URLClassLoader.java:468)
    at java.net.URLClassLoader.access$100(URLClassLoader.java:74)
    at java.net.URLClassLoader$1.run(URLClassLoader.java:369)
    at java.net.URLClassLoader$1.run(URLClassLoader.java:363)
    at java.security.AccessController.doPrivileged(Native Method)
    at java.net.URLClassLoader.findClass(URLClassLoader.java:362)
    at java.lang.ClassLoader.loadClass(ClassLoader.java:419)
    at sun.misc.Launcher$AppClassLoader.loadClass(Launcher.java:352)
    at java.lang.ClassLoader.loadClass(ClassLoader.java:352)
    at java.lang.ClassLoader.defineClass1(Native Method)
    at java.lang.ClassLoader.defineClass(ClassLoader.java:757)
    at java.security.SecureClassLoader.defineClass(SecureClassLoader.java:142)
    at java.net.URLClassLoader.defineClass(URLClassLoader.java:468)
    at java.net.URLClassLoader.access$100(URLClassLoader.java:74)
    at java.net.URLClassLoader$1.run(URLClassLoader.java:369)
    at java.net.URLClassLoader$1.run(URLClassLoader.java:363)
    at java.security.AccessController.doPrivileged(Native Method)
    at java.net.URLClassLoader.findClass(URLClassLoader.java:362)
    at java.lang.ClassLoader.loadClass(ClassLoader.java:419)
    at sun.misc.Launcher$AppClassLoader.loadClass(Launcher.java:352)
    at java.lang.ClassLoader.loadClass(ClassLoader.java:352)
    at java.lang.ClassLoader.defineClass1(Native Method)
    at java.lang.ClassLoader.defineClass(ClassLoader.java:757)

```

回答

搜索mrs客户端目录，将以下两个jar包放入sqoop的lib目录下。

```

-rwxr-xr-x. 1 omm wheel 4813045 Apr  6 15:56 calcite-core-1.19.0.jar
-rwxr-xr-x. 1 omm wheel  459944 Apr  6 16:01 calcite-linq4j-1.19.0.jar

```

27.4.2 使用 hcatalog 方式同步数据，报错 getHiveClient 方法不存在

问题

使用hcatalog方式同步数据，报错getHiveClient方法不存在。

- The authentication type 12 is not supported. Check that you have configured the pg_hba.conf file to include the client's IP address or subnet, and that it
- The authentication type 5 is not supported. Check that you have configured the pg_hba.conf file to include the client's IP address or subnet, and that it
- 问题根因：
 - 报错中type为5时：在执行sqoop import命令时，会启动MapReduce任务，由于MRS Hadoop安装目录（/opt/Bigdata/FusionInsight_HD_*/1_*_NodeManager/install/hadoop/share/hadoop/common/lib）下自带了postgre驱动包gsjdbc4-*.jar，与开源postgre服务不兼容导致报错。
 - 报错中type为12时：数据库的pg_hba.conf文件配置有误。
- 解决方案：
 - 报错中type为5时：在每台MRS NodeManager实例所在节点上移动驱动包gsjdbc4-*.jar到tmp目录下。
mv /opt/Bigdata/FusionInsight_HD_*/1_*_NodeManager/install/hadoop/share/hadoop/common/lib/gsjdbc4-*.jar /tmp
 - 报错中type为12时：调整数据库的pg_hba.conf文件，将address改成sqoop所在节点的ip。

```
# TYPE DATABASE USER ADDRESS METHOD
# "local" is for Unix domain socket connections only
local all all trust
# IPv4 local connections:
host all all 127.0.0.1/32 trust
host all all 0.0.0.0/0 md5
# IPv6 local connections:
host all all ::1/128 trust
#host all all 0.0.0.0/0 password
# Allow replication connections from localhost, by a user with the
# replication privilege.
local replication postgres trust
host replication postgres 127.0.0.1/32 trust
host replication postgres ::1/128 trust
```

- 场景二：（**export**场景）使用**sqoop export**命令抽取开源postgre到MRS hdfs或hive等。
 - 问题现象：

使用sqoop命令查询postgre表可以，但是执行sqoop export命令倒数时报错：The authentication type 5 is not supported. Check that you have configured the pg_hba.conf file to include the client's IP address or subnet, and that it
 - 问题根因：

在执行sqoop export命令时，会启动MapReduce任务，由于MRS Hadoop安装目录（/opt/Bigdata/FusionInsight_HD_*/1_*_NodeManager/install/

hadoop/share/hadoop/common/lib) 下自带了postgre驱动包gsjdbc4-*.jar，与开源postgre服务不兼容导致报错。

- 解决方案：
 - i. 在每台MRS NodeManager实例所在节点上移动驱动包gsjdbc4-*.jar到tmp目录下。
mv /opt/Bigdata/FusionInsight_HD_*/1_*_NodeManager/install/hadoop/share/hadoop/common/lib/gsjdbc4-*.jar /tmp
 - ii. 将/opt/Bigdata/client/Hive/Beeline/lib/gsjdbc4-*.jar删除。

27.4.4 使用 hive-table 方式同步数据到 obs 上的 hive 表报错

问题

使用hive-table方式同步数据到obs上的hive表报错。

```
2021-09-03 16:28:11,611 ERROR tools.DistCp: XAttrs not supported on at least one file system:
org.apache.hadoop.tools.CopyListing$XAttrsNotSupportedException: XAttrs not supported for file system:
obs://fdd-fs
    at org.apache.hadoop.tools.util.DistCpUtils.checkFileSystemXAttrSupport(DistCpUtils.java:555)
    at org.apache.hadoop.tools.DistCp.configureOutputFormat(DistCp.java:341)
    at org.apache.hadoop.tools.DistCp.createJob(DistCp.java:308)
    at org.apache.hadoop.tools.DistCp.createAndSubmitJob(DistCp.java:218)
    at org.apache.hadoop.tools.DistCp.execute(DistCp.java:197)
    at org.apache.hadoop.tools.DistCp.run(DistCp.java:155)
```

回答

修改数据同步方式，将-hive-table改成-hcatalog-table。

27.4.5 使用 hive-table 方式同步数据到 orc 表或者 parquet 表失败

问题

使用hive-table方式同步数据到orc表或者parquet表失败。

报错信息中有kite-sdk的包名。

回答

修改数据同步方式，将-hive-table改成-hcatalog-table。

27.4.6 使用 hive-table 方式同步数据报错

问题

使用hive-table方式同步数据报错：


```
at org.apache.hadoop.hive.ql.metadata.Hive.registerAllFunctionsSource(Hive.java:420) [hive-exec-0.12.3-UBUNTU12.04.1-DRACONUI.jar:0.12.3-UBUNTU12.04.1-DRACONUI]
... 41 more
14:41:42.891 [h1e1438c-07bb-43fd-9149-910a348f6e91 main] ERROR org.apache.hadoop.hive.metastore.ObjectStore - Version information not found in metastore. The process will exit.
14:41:42.892 [h1e1438c-07bb-43fd-9149-910a348f6e91 main] ERROR org.apache.hadoop.hive.metastore.RetryingHMSHandler - ExitSecurityException
at org.apache.sqoop.util.SubprocessSecurityManager.checkExit(SubprocessSecurityManager.java:83)
at java.lang.Runtime.exit(Runtime.java:107)
at java.lang.System.exit(System.java:973)
at org.apache.hadoop.hive.metastore.ObjectStore.checkSchema(ObjectStore.java:9555)
at org.apache.hadoop.hive.metastore.ObjectStore.verifySchema(ObjectStore.java:9531)
at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
at java.lang.reflect.Method.invoke(Method.java:496)
at org.apache.hadoop.hive.metastore.RawStoreProxy.invoke(RawStoreProxy.java:97)
at com.sun.proxy.$Proxy37.verifySchema(Unknown Source)
at org.apache.hadoop.hive.metastore.HiveMetaStore$HMSHandler.getMSForConf(HiveMetaStore.java:903)
at org.apache.hadoop.hive.metastore.HiveMetaStore$HMSHandler.getMS(HiveMetaStore.java:895)
at org.apache.hadoop.hive.metastore.HiveMetaStore$HMSHandler.createDefaultDB(HiveMetaStore.java:978)
at org.apache.hadoop.hive.metastore.HiveMetaStore$HMSHandler.init(HiveMetaStore.java:505)
at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
at java.lang.reflect.Method.invoke(Method.java:496)
at org.apache.hadoop.hive.metastore.RetryingHMSHandler.invokeInternal(RetryingHMSHandler.java:148)
at org.apache.hadoop.hive.metastore.RetryingHMSHandler.invoke(RetryingHMSHandler.java:109)
at org.apache.hadoop.hive.metastore.RetryingHMSHandler.<init>(RetryingHMSHandler.java:81)
at org.apache.hadoop.hive.metastore.RetryingHMSHandler.getProxy(RetryingHMSHandler.java:94)
at org.apache.hadoop.hive.metastore.HiveMetaStore.newRetryingHMSHandler(HiveMetaStore.java:9683)
at org.apache.hadoop.hive.metastore.HiveMetaStoreClient.<init>(HiveMetaStoreClient.java:185)
at org.apache.hadoop.hive.ql.metadata.SessionHiveMetaStoreClient.<init>(SessionHiveMetaStoreClient.java:96)
at sun.reflect.NativeConstructorAccessorImpl.newInstance0(Native Method)
at sun.reflect.NativeConstructorAccessorImpl.newInstance(NativeConstructorAccessorImpl.java:62)
at sun.reflect.DelegatingConstructorAccessorImpl.newInstance(DelegatingConstructorAccessorImpl.java:45)
at java.lang.reflect.Constructor.newInstance(Constructor.java:423)
at org.apache.hadoop.hive.metastore.util.JdbcUtils.newInstance(JdbcUtils.java:84)
at org.apache.hadoop.hive.metastore.RetryingMetaStoreClient.<init>(RetryingMetaStoreClient.java:97)
...

```

回答

修改hive-site.xml，加入如下值。

```
<property>
<name>hive.metastore.schema.verification</name>
<value>false</value>
</property>
```

27.4.7 使用 hcatalog 方式同步 hive parquet 表报错

问题

同步hive parquet表，其分区字段为非string类型，无法正常使用hive import导入，只能考虑使用hcatalog方式，但是hcatalog方式报错如下：

```
2021-09-28 12:12:17.629 INFO common.HCatUtil: mapreduce.lib.hcatoutput.hive.conf is set. Applying configuration differences.
2021-09-28 12:12:17.629 INFO common.HiveClientCache: Initializing cache: eviction-timeout=120 initial-capacity=50 maximum-capacity=50
2021-09-28 12:12:17.648 INFO metastore.HiveMetaStoreClient: Trying to connect to metastore with URI thrift://node-master4yP0W.a0dbfe45-7b6c-4386-83
68f7765cdd.com:9083
2021-09-28 12:12:17.649 INFO metastore.HiveMetaStoreClient: Opened a connection to metastore, current connections: 2
2021-09-28 12:12:17.651 INFO metastore.HiveMetaStoreClient: Connected to metastore.
2021-09-28 12:12:17.651 INFO metastore.RetryingMetaStoreClient: RetryingMetaStoreClient proxy=class org.apache.hive.hcatalog.common.HiveClientCache
eableHiveMetaStoreClient ugi=poseidon (auth=SIMPLE) retries=1 delay=1 lifetime=0
2021-09-28 12:12:17.875 WARN conf.HiveConf: HiveConf of name hive.http.filter.initializers does not exist
2021-09-28 12:12:17.876 WARN conf.HiveConf: HiveConf of name hive.server2.authentication.ldap.url.port does not exist
2021-09-28 12:12:17.877 INFO conf.HiveConf: current conf hive.parquet.time.zone.isLocal=true
2021-09-28 12:12:18.056 INFO hcat.SqoopHCatUtilities: Setting HCatInputFormat filter to days='20210928'
2021-09-28 12:12:18.072 WARN conf.HiveConf: HiveConf of name hive.http.filter.initializers does not exist
2021-09-28 12:12:18.072 WARN conf.HiveConf: HiveConf of name hive.server2.authentication.ldap.url.port does not exist
2021-09-28 12:12:18.073 INFO conf.HiveConf: current conf hive.parquet.time.zone.isLocal=true
2021-09-28 12:12:18.073 INFO common.HCatUtil: mapreduce.lib.hcatoutput.hive.conf is set. Applying configuration differences.
2021-09-28 12:12:18.108 ERROR tool.ImportTool: Import failed: java.io.IOException: MetaException(message:Filtering is supported only on partition k
f type string)
at org.apache.hive.hcatalog.mapreduce.HCatInputFormat.setFilter(HCatInputFormat.java:120)
at org.apache.sqoop.mapreduce.hcat.SqoopHCatUtilities.configureHCat(SqoopHCatUtilities.java:381)
at org.apache.sqoop.mapreduce.hcat.SqoopHCatUtilities.configureImportOutputFormat(SqoopHCatUtilities.java:850)
at org.apache.sqoop.mapreduce.ImportJobBase.configureOutputFormat(ImportJobBase.java:102)
at org.apache.sqoop.mapreduce.ImportJobBase.runImport(ImportJobBase.java:263)
at org.apache.sqoop.manager.SqlManager.importQuery(SqlManager.java:748)
at org.apache.sqoop.tool.ImportTool.importTable(ImportTool.java:522)
at org.apache.sqoop.tool.ImportTool.run(ImportTool.java:628)
at org.apache.sqoop.Sqoop.run(Sqoop.java:147)
at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:76)
at org.apache.sqoop.Sqoop.runSqoop(Sqoop.java:183)
at org.apache.sqoop.Sqoop.runTool(Sqoop.java:234)
```

回答

1. 修改sqoop源码SqoopHCatUtilities中的代码，将限制代码去掉。
2. 修改hive客户端中的hive-site.xml文件，修改hive.metastore.integral.jdo.pushdown参数为true。

27.4.8 使用 Hcatalog 方式同步 Hive 和 MySQL 之间的数据，timestamp 和 data 类型字段会报错

问题

使用Hcatalog方式同步Hive和MySQL之间的数据，timestamp和data类型字段会报错：

```
2021-10-20 21:16:34,034 | INFO | main | current conf hive.parquet.time.zone.isLocal=true | HiveConf.java:5506
2021-10-20 21:16:34,034 | INFO | Thread-19 | Auto-progress thread is finished. keepGoing=false | ProgressThread.java:158
2021-10-20 21:16:34,034 | WARN | main | Exception running child: java.lang.ClassCastException: org.apache.hadoop.hive.common.type.Timestamp cannot be cast to java.sql.Timestamp
    at org.apache.sqoop.mapreduce.hcat.SqoopHcatExportHelper.convertToSqoop(SqoopHcatExportHelper.java:203)
    at org.apache.sqoop.mapreduce.hcat.SqoopHcatExportHelper.convertToSqoopRecord(SqoopHcatExportHelper.java:138)
    at org.apache.sqoop.mapreduce.hcat.SqoopHcatExportMapper.map(SqoopHcatExportMapper.java:56)
    at org.apache.sqoop.mapreduce.hcat.SqoopHcatExportMapper.map(SqoopHcatExportMapper.java:35)
    at org.apache.hadoop.mapreduce.Mapper.run(Mapper.java:146)
    at org.apache.sqoop.mapreduce.AutoProgressMapper.run(AutoProgressMapper.java:64)
    at org.apache.hadoop.mapred.MapTask.runRedMapper(MapTask.java:799)
    at org.apache.hadoop.mapred.MapTask.run(MapTask.java:347)
    at org.apache.hadoop.mapred.YarnChild$1.run(YarnChild.java:183)
    at java.security.AccessController.doPrivileged(Native Method)
    at javax.security.auth.Subject.doAs(Subject.java:422)
    at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1761)
    at org.apache.hadoop.mapred.YarnChild.main(YarnChild.java:177)
| YarnChild.java:199
```

回答

- 调整Sqoop源码包中的代码，将timestamp强制转换类型和Hive保持一致。
- 将Hive中的字段类型修改为String。

28 使用 Storm

28.1 从零开始使用 Storm

用户可以在MRS集群的客户端中提交和删除Storm拓扑等基本功能。

前提条件

已安装MRS集群客户端，例如安装目录为“/opt/hadoopclient”。以下操作的客户端目录只是举例，请根据实际安装目录修改。

操作步骤

步骤1 根据业务情况，准备好客户端，登录安装客户端的节点。

请根据客户端所在位置，参考[安装客户端](#)章节，登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端目录，例如“/opt/hadoopclient”。

```
cd /opt/hadoopclient
```

步骤3 执行以下命令，配置环境变量。

```
source bigdata_env
```

步骤4 启用Kerberos认证的集群，执行以下命令认证用户身份。未启用Kerberos认证的集群无需执行。

```
kinit Storm用户
```

步骤5 执行以下命令，提交Storm拓扑：

```
storm jar 拓扑包路径 拓扑Main方法的类名称 拓扑名称
```

界面提示以下信息表示提交成功：

```
Finished submitting topology: topo1
```

步骤6 执行以下命令，查看Storm中的拓扑。启用Kerberos认证的集群，只有属于“stormadmin”或“storm”的用户可以查看所有拓扑。

```
storm list
```

步骤7 执行以下命令，删除Storm中的拓扑。

```
storm kill 拓扑名称  
----结束
```

28.2 使用 Storm 客户端

操作场景

该任务指导用户在运维场景或业务场景中使用Storm客户端。

前提条件

- 已安装客户端。例如安装目录为“/opt/hadoopclient”。
- 各组件业务用户由MRS集群管理员根据业务需要创建。安全模式下，“机机”用户需要下载keytab文件。“人机”用户第一次登录时需修改密码。（普通模式不涉及）

操作步骤

步骤1 根据业务情况，准备好客户端，登录安装客户端的节点。

请根据客户端所在位置，参考[安装客户端](#)章节，登录安装客户端的节点。

步骤2 执行以下命令，切换到客户端安装目录。

```
cd /opt/hadoopclient
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 执行以下命令，进行用户认证。（普通模式跳过此步骤）

```
kinit 组件业务用户
```

步骤5 执行命令进行客户端操作。

例如执行以下命令：

- cql
- storm

说明

同一个storm客户端不能同时连接安全和非安全的ZooKeeper。

```
----结束
```

28.3 使用客户端提交 Storm 拓扑

操作场景

用户可以根据业务需要，在集群的客户端中提交Storm拓扑，持续处理用户的流数据。启用Kerberos认证的集群，需要提交拓扑的用户属于“stormadmin”或“storm”组。

前提条件

已刷新客户端。

操作步骤

步骤1 根据业务情况，准备好客户端，登录安装客户端的节点。

请根据客户端所在位置，参考[安装客户端](#)章节，登录安装客户端的节点。

步骤2 执行以下命令，设置拓扑的jar包权限。

例如修改“/opt/storm/topology.jar”的权限：

```
chmod 600 /opt/storm/topology.jar
```

步骤3 执行以下命令，切换到客户端目录，例如“/opt/client”。

```
cd /opt/client
```

步骤4 执行以下命令，配置环境变量。

```
source bigdata_env
```

步骤5 若安装了Storm多实例，在使用Storm命令提交拓扑时，请执行以下命令加载具体实例的环境变量，否则请跳过此步骤。例如，Storm-2实例：

```
source Storm-2/component_env
```

步骤6 启用Kerberos认证的集群，执行以下命令认证用户身份。未启用Kerberos认证的集群无需执行。

```
kinit Storm用户
```

步骤7 MRS 3.x之前版本：执行以下命令，提交Storm拓扑。

```
storm jar 拓扑包路径 拓扑Main方法的类名称 拓扑名称
```

界面提示以下信息表示提交成功：

```
Finished submitting topology: topo1
```

说明

- 如果需要拓扑支持采样消息，则还需要增加参数“topology.debug”和“topology.eventlogger.executors”。
- 拓扑如何处理数据是拓扑自身行为。样例拓扑随机生成字符并分隔字符串，需要查看处理情况时，请启用采样功能并参见[查看Storm拓扑日志](#)。

步骤8 MRS 3.x及后续版本：执行以下命令，提交拓扑任务。

```
storm jar topology-jar-path class 入参列表
```

- topology-jar-path：表示拓扑的jar包所在路径。
- class：表示拓扑使用的main方法所在类名称。
- 入参列表：表示拓扑使用的main方法入参。

例如，提交WordCount计算的拓扑“/opt/storm/topology.jar”并以拓扑命名作为入参，执行：

```
storm jar /opt/storm/topology.jar  
com.huawei.storm.example.WordCountTopology topology1
```

显示以下信息表示拓扑提交成功：

```
Finished submitting topology: topology1
```

📖 说明

- 登录认证用户必须与所加载环境变量（component_env）一一对应，否则使用storm命令提交拓扑任务出错。
- 加载客户端环境变量且对应用户登录后，该用户可以在任意storm客户端下执行storm命令来提交拓扑任务，但提交拓扑命令执行完成后，提交成功的拓扑仍然在用户所对应的Storm集群中，不会出现在其他Storm集群中。
- 如果修改了集群域名，需要在提交拓扑前重新设置域名信息，进入cql语句执行命令，例如：
`set "kerberos.domain.name" = "hadoop.huawei.com"`。

步骤9 执行以下命令，查看Storm中的拓扑。启用Kerberos认证的集群，只有属于“stormadmin”或“storm”的用户可以查看所有拓扑。

```
storm list
```

```
----结束
```

28.4 访问 Storm 的 WebUI

操作场景

用户可以通过Storm的WebUI，在图形化界面使用Storm。

Storm的WebUI支持查看以下信息：

- Storm集群汇总信息
- Nimbus汇总信息
- 拓扑汇总信息
- Supervisor汇总信息
- Nimbus配置信息

前提条件

- 获取用户“admin”账号密码。“admin”密码在创建集群时由用户指定。
- 使用其他用户访问Storm WebUI，需要用户属于“storm”或“stormadmin”用户组。

操作步骤

步骤1 进入组件管理页面：

- MRS 3.x之前版本，单击集群名称，登录集群详情页面，选择“组件管理”。

📖 说明

若集群详情页面没有“组件管理”页签，请先完成IAM用户同步（在集群详情页的“概览”页签，单击“IAM用户同步”右侧的“同步”进行IAM用户同步）。

- MRS 3.x及后续版本，登录FusionInsight Manager，具体请参见[访问 FusionInsight Manager（MRS 3.x及之后版本）](#)。然后选择“集群 > 待操作的集群名称 > 服务”。

步骤2 登录Storm WebUI：

- MRS 3.x之前版本：选择“Storm”，在“Storm 概述”的“Storm Web UI”，单击任意一个UI链接，打开Storm的WebUI。

📖 说明

第一次访问Storm WebUI，需要在浏览器中添加站点信任以继续打开页面。

- MRS 3.x及后续版本：选择“Storm > 概览”，在“基本信息”的“Storm Web UI”，单击任意一个UI链接，打开Storm的WebUI。

----结束

相关任务

- 单击拓扑名称，可查看指定拓扑的详细信息、拓扑状态、Spouts信息、Bolts信息和拓扑配置。
- 在“Topology actions”区域，用户可以对拓扑执行激活、去激活、重部署、删除操作、调试、停止调试和修改日志级别，即“Activate”、“Deactivate”、“Rebalance”、“Kill”、“Debug”、“Stop Debug”、“Change Log Level”。重部署和删除操作需要设置操作执行的等待时间，单位为秒。
- 在“Topology Visualization”区域，用户可以执行拓扑可视化操作，即单击“Show Visualization”。拓扑可视化后，WebUI将显示拓扑结构图。

28.5 管理 Storm 拓扑

操作场景

用户可以使用Storm的WebUI管理拓扑。“storm”用户组的用户只能管理由自己提交的拓扑任务，“stormadmin”用户组的用户可以管理所有拓扑任务。

操作步骤

步骤1 访问Storm的WebUI，请参考[访问Storm的WebUI](#)。

步骤2 在“Topology summary”区域，单击指定的拓扑名称。

步骤3 通过“Topology actions”管理Storm拓扑。

- 激活拓扑
单击“Activate”，转化当前拓扑为激活状态。

- 去激活拓扑
单击“Deactivate”，转化当前拓扑为去激活状态。
- 重部署拓扑
单击“Rebalance”，将当前拓扑重新部署执行，需要输入执行重部署的等待时间，单位为秒。一般在集群中节点数发生变化时进行，以更好利用集群资源。
- 删除拓扑
单击“Kill”，将当前拓扑删除，需要输入执行操作的等待时间，单位为秒。
- 采样、停止采样拓扑消息
单击“Debug”，在弹出窗口输入流数据采样消息的数值，单位为百分比，表示从开始采样到停止采样这段时间内所有数据的采集比例。例如输入“10”，则采集比例为10%。
如果需要停止采样，则单击“Stop Debug”。

📖 说明

只有在提交拓扑时启用采样功能，才支持此功能。查看采样处理数据，请参见[查看Storm拓扑日志](#)。

- 修改拓扑日志级别
单击“Change Log Level”，可以为Storm日志指定新的日志信息级别。

步骤4 显示拓扑结构图。

在“Topology Visualization”区域单击“Show Visualization”，执行拓扑可视化操作。

----结束

28.6 查看 Storm 拓扑日志

操作场景

用户需要查看Storm拓扑在worker进程中的执行情况时，需要查看worker中关于拓扑的日志。如果需要查询拓扑在运行时数据处理的日志，提交拓扑并启用“Debug”功能后可以查看日志。仅启用Kerberos认证的流集群支持该场景，且用户需要是拓扑的提交者，或者加入“stormadmin”。

前提条件

- 在工作环境完成网络配置。
- 需要查看处理数据的拓扑，提交时已启用采样功能。

查看 worker 进程日志

步骤1 访问Storm的WebUI，请参考[访问Storm的WebUI](#)。

步骤2 在“Topology Summary”区域单击指定的拓扑名称，打开拓扑的详细信息。

步骤3 单击要查看日志的“Spouts”或“Bolts”任务，在“Executors (All time)”区域单击“Port”列的端口值，查看详细日志内容。

----结束

查看拓扑处理数据日志

- 步骤1** 访问Storm的WebUI，请参考[访问Storm的WebUI](#)。
- 步骤2** 在“Topology Summary”区域单击指定的拓扑名称，打开拓扑的详细信息。
- 步骤3** 单击“Debug”，输入采样数据的百分比数值，并单击“OK”开始采样。
- 步骤4** 单击拓扑的“Spouts”或“Bolts”任务，在“Component summary”单击“events”打开处理数据日志。
- 结束

28.7 Storm 常用参数

本章节内容适用于MRS 3.x及后续版本。

参数入口

参数入口，请参考[修改集群服务配置参数](#)。

参数说明

表 28-1 参数说明

配置参数	说明	默认值
supervisor.slots.ports	supervisor上能够运行workers的端口列表。每个worker占用一个端口，且每个端口只运行一个worker。通过这项配置可以设置每台机器上运行的worker数量。端口的取值范围是1024到65535，不同端口使用逗号分隔。	6700,6701,6702,6703
WORKER_GC_OPTS	supervisor启动worker时使用的jvm选项。需要根据业务中对内存等的使用来进行设置，例如是简单业务处理，建议1G，即“-Xmx1G”；如果有窗口缓存，根据窗口大小计算：每条记录大小*周期*2。	-Xms1G -Xmx1G -XX:+UseG1GC -XX:+PrintGCDetails -Xloggc:artifacts/gc.log -XX:+PrintGCDateStamps -XX:+PrintGCTimeStamps -XX:+UseGCLogFileRotation -XX:NumberOfGCLogFiles=10 -XX:GCLogFileSize=1M -XX:+HeapDumpOnOutOfMemoryError -XX:HeapDumpPath=artifacts/heapdump

配置参数	说明	默认值
default.scheduler.mode	默认调度器的调度模式。目前支持两个值，具体值与含义如下： <ul style="list-style-type: none">“AVERAGE”：使用按空闲Slot数目为优先级的调度机制“RATE”：使用按空闲Slot比率为优先级的调度机制	AVERAGE
nimbus.thrift.threads	设置主用Nimbus对外提供服务时的最大连接线程数。当Storm集群规模较大，Supervisor实例数量较多时，需要增加线程数。	512

28.8 配置 Storm 业务用户密码策略

操作场景

本章节内容适用于MRS 3.x及后续版本。

使用Storm业务用户提交一个拓扑以后，该任务需要使用提交拓扑的用户身份持续运行。在拓扑运行的过程中，worker进程可能需要正常重启以保持拓扑工作。若业务用户的密码被修改，或密码使用天数超过了默认密码策略指定的最大有效期，则会影响拓扑正常运行。MRS集群管理员需要根据企业安全要求，为Storm业务用户配置独立的密码策略。

说明

如果不为Storm业务用户配置独立的密码策略，在修改业务用户密码以后，可以删除旧的拓扑并重新提交，使拓扑继续运行。

对系统的影响

- 为Storm业务用户配置独立的密码策略后，此用户将不受Manager界面上的“密码策略”配置影响。
- 为Storm业务用户配置独立的密码策略后，如果配置了跨集群互信，请根据此密码策略，在Manager为Storm业务用户重置密码。

前提条件

MRS集群管理员已明确业务需求，并创建好“人机”用户，例如“testpol”。

操作步骤

步骤1 以“omm”用户登录集群内任意节点。

步骤2 执行以下命令，防止超时退出。

```
TMOUT=0
```

📖 说明

执行完本章节操作后，请及时恢复超时退出时间，执行命令 `TMOUT=超时退出时间`。例如：
`TMOUT=600`，表示用户无操作600秒后超时退出。

步骤3 执行以下命令，导出环境变量。

```
EXECUTABLE_HOME="${CONTROLLER_HOME}/kerberos_user_specific_binay/  
kerberos"
```

```
LD_LIBRARY_PATH=${EXECUTABLE_HOME}/lib:$LD_LIBRARY_PATH
```

```
PATH=${EXECUTABLE_HOME}/bin:$PATH
```

步骤4 执行以下命令，并输入Kerberos管理员密码，进入Kerberos管理控制台。

```
kadmin -p kadmin/admin
```

📖 说明

第一次使用“kadmin/admin”用户需要修改“kadmin/admin”密码。

界面显示如下信息，则表示已成功进入Kerberos管理控制台。

```
kadmin:
```

步骤5 执行以下命令，查看创建好的“Human-Machine”用户的具体信息。

```
getprinc 用户名
```

例如，查看“testpol”用户的详细信息：

```
getprinc testpol
```

界面显示如下信息，说明指定用户使用了默认的密码策略：

```
Principal: testpol@<系统域名>  
.....  
Policy: default
```

步骤6 执行以下命令，为Storm业务用户创建独立的密码策略，例如“streampol”：

```
addpol -maxlife 0day -minlife 0sec -history 1 -maxfailure 5 -  
failurecountinterval 5min -lockoutduration 5min -minlength 8 -minclasses 4  
streampol
```

其中“-maxlife”表示密码最大有效期，“0day”表示永不过期。

步骤7 执行以下命令，查看新创建的策略“streampol”。

```
getpol streampol
```

界面显示如下信息，说明新策略已指定密码不过期：

```
Policy: streampol  
Maximum password life: 0 days 00:00:00  
.....
```

步骤8 执行以下命令，将新的策略“streampol”应用到Storm用户“testpol”。

```
modprinc -policy streampol testpol
```

其中“streampol”是策略名称，“testpol”是用户名。

界面显示如下信息，说明指定用户的属性已修改：

```
Principal "testpol@<系统域名>" modified.
```

步骤9 执行以下命令，查看Storm用户“testpol”用户的当前信息。

```
getprinc testpol
```

界面显示如下信息，说明指定用户使用了新的密码策略：

```
Principal: testpol@<系统域名>  
.....  
Policy: streampol
```

----结束

28.9 迁移 Storm 业务至 Flink

28.9.1 概述

本章节内容适用于MRS 3.x及后续版本。

Flink从0.10.0版本开始提供了一套API可以将使用Storm API编写的业务平滑迁移到Flink平台上，只需要极少的改动即可完成。通过这项转换可以覆盖大部分的业务场景。

Flink支持两种方式的业务迁移：

1. 完整迁移Storm业务：转换并运行完整的由Storm API开发的Storm拓扑。
2. 嵌入式迁移Storm业务：在Flink的DataStream中嵌入Storm的代码，如使用Storm API编写的Spout/Bolt。

Flink提供了flink-storm包用来完成上述转换。

28.9.2 完整迁移 Storm 业务

操作场景

该任务指导用户通过Storm业务完整迁移的方式转换并运行完整的由Storm API开发的Storm拓扑。

操作步骤

步骤1 打开Storm业务工程，修改工程的pom文件，增加“flink-storm”、“flink-core”和“flink-streaming-java_2.11”的引用。如下：

```
<dependency>  
  <groupId>org.apache.flink</groupId>  
  <artifactId>flink-storm_2.11</artifactId>  
  <version>1.4.0</version>  
  <exclusions>  
    <exclusion>  
      <groupId>*</groupId>  
      <artifactId>*</artifactId>  
    </exclusion>  
  </exclusions>  
</dependency>  
<dependency>  
  <groupId>org.apache.flink</groupId>  
  <artifactId>flink-core</artifactId>  
  <version>1.4.0</version>
```

```
<exclusions>
  <exclusion>
    <groupId>*</groupId>
    <artifactId>*</artifactId>
  </exclusion>
</exclusions>
</dependency>

<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-streaming-java_2.11</artifactId>
  <version>1.4.0</version>
  <exclusions>
    <exclusion>
      <groupId>*</groupId>
      <artifactId>*</artifactId>
    </exclusion>
  </exclusions>
</dependency>
```

📖 说明

如果是非maven工程，则手动收集如上jar包，添加到工程的classpath中。

步骤2 修改拓扑提交部分代码，下面以WordCount为例：

1. Storm拓扑的构造部分保持不变，无需修改，包括使用Storm API开发的Spout和Bolt都无需修改。

```
TopologyBuilder builder = new TopologyBuilder();
builder.setSpout("spout", new RandomSentenceSpout(), 5);
builder.setBolt("split", new SplitSentenceBolt(), 8).shuffleGrouping("spout");
builder.setBolt("count", new WordCountBolt(), 12).fieldsGrouping("split", new Fields("word"));
```

2. 拓扑的提交部分需要修改，Storm的提交示例如下：

```
Config conf = new Config();
conf.setNumWorkers(3);

StormSubmitter.submitTopology("word-count", conf, builder.createTopology());
```

需要进行如下修改：

```
Config conf = new Config();
conf.setNumWorkers(3);

//将Storm的Config转化为Flink的StormConfig
StormConfig stormConfig = new StormConfig(conf);

//使用Storm的TopologBuilder构造FlinkTopology
FlinkTopology topology = FlinkTopology.createTopology(builder);

//获取StreamExecutionEnvironment
StreamExecutionEnvironment env = topology.getExecutionEnvironment();

//将StormConfig设置到Job的环境中，用于构造Bolt和Spout
//如果Bolt和Spout初始化时不需要config，则不用设置
env.getConfig().setGlobalJobParameters(stormConfig);
//执行拓扑提交
topology.execute();
```

3. 重新打包之后使用flink命令行进行提交：

```
flink run -class {MainClass} WordCount.jar
```

----结束

28.9.3 嵌入式迁移 Storm 业务

操作场景

该任务指导用户通过嵌入式迁移的方式在Flink的DataStream中嵌入Storm的代码，如使用Storm API编写的Spout/Bolt。

操作步骤

- 步骤1** 在Flink中，对Storm拓扑中的Spout和Bolt进行嵌入式转换，将之转换为Flink的Operator，代码示例如下：

```
//set up the execution environment
final StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();

//get input data
final DataStream<String> text = getTextDataStream(env);

final DataStream<Tuple2<String, Integer>> counts = text

//split up the lines in pairs (2-tuples) containing: (word,1)
//this is done by a bolt that is wrapped accordingly
.transform("CountBolt",
    TypeExtractor.getForObject(new Tuple2<String, Integer>("", 0)),
    new BoltWrapper<String, Tuple2<String, Integer>>(new CountBolt()))
//group by the tuple field "0" and sum up tuple field "1"
.keyBy(0).sum(1);
// execute program
env.execute("Streaming WordCount with bolt tokenizer");
```

- 步骤2** 修改完成后使用Flink命令进行提交。

```
flink run -class {MainClass} WordCount.jar
```

----结束

28.9.4 迁移 Storm 对接的外部安全组件业务

迁移 Storm 对接 HDFS 和 HBase 组件的业务

如果Storm的业务使用的storm-hdfs或者storm-hbase插件包进行的对接，那么在按照[完整迁移Storm业务](#)进行迁移时，需要指定特定安全参数，如下：

```
//初始化Storm的Config
Config conf = new Config();

//初始化安全插件列表
List<String> auto_tgts = new ArrayList<String>();
//添加AutoTGT插件
auto_tgts.add("org.apache.storm.security.auth.kerberos.AutoTGT");
//添加AutoHDFS插件
//如果对接HBase，则如下更改为： auto_tgts.add("org.apache.storm.hbase.security.AutoHBase");
auto_tgts.add("org.apache.storm.hdfs.common.security.AutoHDFS");

//设置安全参数
conf.put(Config.TOPOLOGY_AUTO_CREDENTIALS, auto_tgts);
//设置worker个数
conf.setNumWorkers(3);

//将Storm的Config转化为Flink的StormConfig
StormConfig stormConfig = new StormConfig(conf);

//使用Storm的TopologBuilder构造FlinkTopology
```

```
FlinkTopology topology = FlinkTopology.createTopology(builder);

//获取StreamExecutionEnvironment
StreamExecutionEnvironment env = topology.getExecutionEnvironment();

//将StormConfig设置到Job的环境变量中，用于构造Bolt和Spout
//如果Bolt和Spout初始化时不需要config，则不用设置
env.getConfig().setGlobalJobParameters(stormConfig);

//执行拓扑提交
topology.execute();
```

增加如上的安全插件配置后，可以避免HDFS Bolt和HBase Bolt在初始化过程中的无谓登录，因为Flink已经实现准备好了安全上下文，无需再登录。

迁移 Storm 对接其他安全组件的业务

如果Storm的业务使用的storm-kafka-client等插件包进行的对接时，需要注意，之前所配置的安全插件需要去掉，如下：

```
List<String> auto_tgts = new ArrayList<String>();
//keytab方式
auto_tgts.add("org.apache.storm.security.auth.kerberos.AutoTGTFromKeytab");

//将客户端配置的plugin列表写入config指定项中
//安全模式必配
//普通模式不用配置，请注释掉该行
conf.put(Config.TOPOLOGY_AUTO_CREDENTIALS, auto_tgts);
```

如上所配置的AutoTGTFromKeytab插件在进行业务迁移时，必须删除，否则会引起相应Bolt或Spout初始化时登录异常。

28.10 Storm 日志介绍

本章节内容适用于MRS 3.x及后续版本。

日志描述

日志路径：Storm相关日志的默认存储路径为“/var/log/Bigdata/storm/角色名”（运行日志），“/var/log/Bigdata/audit/storm/角色名”（审计日志）。

- Nimbus：“/var/log/Bigdata/storm/nimbus”（运行日志），“/var/log/Bigdata/audit/storm/nimbus”（审计日志）
- Supervisor：“/var/log/Bigdata/storm/supervisor”（运行日志），“/var/log/Bigdata/audit/storm/supervisor”（审计日志）
- UI：“/var/log/Bigdata/storm/ui”（运行日志），“/var/log/Bigdata/audit/storm/ui”（审计日志）
- Logviewer：“/var/log/Bigdata/storm/logviewer”（运行日志），“/var/log/Bigdata/audit/storm/logviewer”（审计日志）

日志归档规则：Storm的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过10MB的时候会自动压缩，压缩后的日志文件名规则为：“<原有日志名>.log.[编号].gz”。默认最多保留最近的20个压缩文件，压缩文件保留个数和压缩文件阈值可以配置。

审计日志压缩后的日志文件名规则为：“audit.log.[yyyy-MM-dd].[编号].zip”。该文件永远都不会删除。

表 28-2 Storm 日志列表

日志类型	日志文件名	描述
运行日志	nimbus/access.log	Nimbus用户访问日志。
	nimbus/nimbus-<PID>-gc.log	Nimbus进程的GC日志。
	nimbus/checkavailable.log	Nimbus可用性检查日志。
	nimbus/checkService.log	Nimbus可服务性检查日志。
	nimbus/metrics.log	Nimbus监控统计的日志。
	nimbus/nimbus.log	Nimbus进程运行日志。
	nimbus/postinstall.log	Nimbus安装后的工作日志。
	nimbus/prestart.log	Nimbus启动前的工作日志。
	nimbus/start.log	Nimbus启动的工作日志。
	nimbus/stop.log	Nimbus停止的工作日志。
	supervisor/access.log	Supervisor用户访问日志。
	supervisor/metrics.log	Supervisor监控统计的日志。
	supervisor/postinstall.log	Supervisor安装后的工作日志。
	supervisor/prestart.log	Supervisor启动前的工作日志。
	supervisor/start.log	Supervisor启动的工作日志。
	supervisor/stop.log	Supervisor停止的工作日志。
	supervisor/supervisor.log	Supervisor进程运行日志。
	supervisor/supervisor-<PID>-gc.log	Supervisor进程的GC日志。
	ui/access.log	UI用户访问日志。
	ui/metric.log	UI监控统计的日志。
ui/ui-<PID>-gc.log	UI进程的GC日志。	
ui/postinstall.log	UI安装后的工作日志。	

日志类型	日志文件名	描述
	ui/prestart.log	UI启动前的工作日志。
	ui/start.log	UI启动的工作日志。
	ui/stop.log	UI停止的工作日志。
	ui/ui.log	UI进程运行日志。
	logviewer/access.log	Logviewer用户访问日志。
	logviewer/metric.log	Logviewer监控统计的日志。
	logviewer/logviewer- <PID>-gc.log	Logviewer进程的GC日志。
	logviewer/logviewer.log	logviewer运行日志。
	logviewer/postinstall.log	logviewer安装后的工作日志。
	logviewer/prestart.log	logviewer启动前的工作日志。
	logviewer/start.log	logviewer启动的工作日志。
	logviewer/stop.log	logviewer停止的工作日志。
	supervisor/[topologyId]- worker-[端口号].log	Worker进程运行日志，一个端口占用一个日志文件，系统默认包含29100,29101,29102,29103,29304五个端口。
	supervisor/metadata/ [topologyid]-worker-[端口号].yaml	worker日志元数据文件，logviewer在清理日志的时候会以该文件来作为清理依据。该文件会被logviewer日志清理线程根据一定条件自动删除。
	nimbus/cleanup.log	Nimbus卸载的清理日志。
	logviewer/cleanup.log	logviewer卸载的清理日志。
	ui/cleanup.log	UI卸载的清理日志。
	supervisor/cleanup.log	Supervisor卸载的清理日志。
	leader_switch.log	Storm主备倒换运行日志。

日志类型	日志文件名	描述
审计日志	nimbus/audit.log	Nimbus审计日志。
	ui/audit.log	UI审计日志。
	supervisor/audit.log	Supervisor审计日志。
	logviewer/audit	Logviewer审计日志。

日志级别

Storm提供了如表28-3所示的日志级别。

运行日志和审计日志的级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 28-3 日志级别

级别	描述
ERROR	ERROR表示系统运行的错误信息。
WARN	WARN表示当前事件处理存在异常信息。
INFO	INFO表示记录系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 请参考[修改集群服务配置参数](#)，进入Storm的“全部配置”页面。
- 步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤3** 选择所需修改的日志级别。
- 步骤4** 保存配置，在弹出窗口中单击“确定”使配置生效。

----结束

日志格式

Storm的日志格式如下所示：

表 28-4 日志格式

日志类型	格式	示例
运行日志	<code>%d{yyyy-MM-dd HH:mm:ss,SSS} %-5p [%t] %m %logger (%F:%L) %n</code>	2015-03-11 23:04:00,241 INFO [RMI TCP Connection(2646)-10.0.0.2] The baseSleepTimeMs [1000] the maxSleepTimeMs [1000] the maxRetries [1] backtype.storm.utils.StormBoundedExponentialBackoffRetry (StormBoundedExponentialBackoffRetry.java:46)
	<code><yyyy-MM-dd HH:mm:ss,SSS><HostName><RoleName><LogLevel><Message></code>	2017-03-28 02:57:52 493 10-5-146-1 storm- INFO Nimbus start normally
审计日志	<code><用户名><用户IP><时间><操作><操作对象><操作结果></code>	UserName=storm/hadoop, UserIP=10.10.0.2, Time=Tue Mar 10 01:15:35 CST 2015, Operation=Kill, Resource=test, Result=Success

28.11 性能调优

28.11.1 Storm 性能调优

操作场景

通过调整Storm参数设置，可以提升特定业务场景下Storm的性能。

本章节适用于MRS 3.x及后续版本。

修改服务配置参数，请参考[修改集群服务配置参数](#)。

拓扑调优

当需要提升Storm数据量处理性能时，可以通过拓扑调优的操作提高效率。建议在可靠性要求不高的场景下进行优化。

表 28-5 调优参数

配置参数	默认值	调优场景
topology.acker.executors	null	Acker的执行器数量。当业务应用对可靠性要求较低，允许不处理部分数据，可设置参数值为“null”或“0”，以关闭Acker的执行器，减少流控制，不统计消息时延，提高性能。
topology.max.spout.pending	null	Spout消息缓存数，仅在Acker不为0或者不为null的情况下生效。Spout将发送到下游Bolt的每条消息加入到pending队列，待下游Bolt处理完成并确认后，再从pending队列移除，当pending队列占满时Spout暂停消息发送。增加pending值可提高Spout的每秒消息吞吐量，提高性能，但延时同步增加。
topology.transfer.buffer.size	32	每个worker进程Distruptor消息队列大小，建议在4到32之间，增大消息队列可以提升吞吐量，但延时可能会增加。
RES_CPUSET_PERCENTAGE	80	设置各个节点上的Supervisor角色实例（包含其启动并管理的Worker进程）所使用的物理CPU百分比。根据Supervisor所在节点业务量需求，适当调整参数值，优化CPU使用率。

JVM 调优

当应用程序需要处理大量数据从而占用更多的内存时，存在worker内存大于2GB的情况，推荐使用G1垃圾回收算法。

表 28-6 调优参数

配置参数	缺省值	调优场景
WORKER_GC_OPTS	-Xms1G - Xmx1G - XX:+UseG1GC - C - XX:+PrintGCDe- tails - Xloggc:artifacts/gc.log - XX:+PrintGCDateStamps - XX:+PrintGCTimeStamps - XX:+UseGCLogFileRotation - XX:NumberOfGCLogFiles=10 - XX:GCLogFile- Size=1M - XX:+HeapDumpOnOutOfMemoryError - XX:HeapDumpPath=artifacts/heapdump	应用程序内存中需要保存大量数据，worker进程使用的内存大于2G，那么建议使用G1垃圾回收算法，可修改参数值为“-Xms2G -Xmx5G -XX:+UseG1GC”。

29 使用 Tez

29.1 访问 Tez WebUI 查看任务执行结果

Tez提供Tez任务执行过程图形化展示功能，使用户可以通过界面的方式查看Tez任务执行细节。

本章节适用于MRS 3.x及后续版本。

前提条件

已安装Yarn服务的TimelineServer实例。

使用介绍

登录Manager系统，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)，在Manager界面选择“集群 > 服务 > Tez”，在“基本信息”中单击“Tez WebUI”右侧的链接，打开Tez WebUI。可查看执行的Tez任务执行细节。

29.2 Tez 常用配置参数

参数入口

在Manager系统中，选择“集群 > 服务 > Tez > 配置”，选择“全部配置”。在搜索框中输入参数名称。

本章节适用于MRS 3.x及后续版本。

参数说明

表 29-1 参数说明

配置参数	说明	缺省值
property.tez.log.dir	Tez日志目录。	/var/log/Bigdata/tez/tezui

配置参数	说明	缺省值
property.tez.log.level	Tez的日志级别。	INFO

29.3 Tez 日志介绍

本章节适用于MRS 3.x及后续版本。

日志描述

日志路径： Tez相关日志的默认存储路径为“/var/log/Bigdata/tez/角色名”。

TezUI：“/var/log/Bigdata/tez/tezui”（运行日志），“/var/log/Bigdata/audit/tez/tezui”（审计日志）。

日志归档规则： Tez的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过20MB的时候（此日志文件大小可进行配置），会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd_hh-mm-ss>.[编号].log.zip”。最多保留最近的20个压缩文件，压缩文件保留个数和压缩文件阈值可以配置。

表 29-2 Tez 日志列表

日志类型	日志文件名	描述
运行日志	tezui.out	TezUI运行环境信息日志
	tezui.log	TezUI进程的运行日志
	tezui-omm-<日期>-gc.log.<编号>	TezUI进程的GC日志
	prestartDetail.log	TezUI启动前的工作日志
	check-serviceDetail.log	TezUI服务启动是否成功的检查日志
	postinstallDetail.log	TezUI安装后的工作日志
	startDetail.log	TezUI进程启动日志
	stopDetail.log	TezUI进程停止日志
审计日志	tezui-audit.log	TezUI审计日志

日志级别

TezUI提供了如表29-3所示的日志级别。

运行日志的级别优先级从高到低分别是ERROR、WARN、INFO、DEBUG，程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 29-3 日志级别

级别	描述
ERROR	ERROR表示系统运行的错误信息。
WARN	WARN表示当前事件处理存在异常信息。
INFO	INFO表示记录系统及各事件正常运行状态信息。
DEBUG	DEBUG表示记录系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 登录Manager。
- 步骤2** 选择“集群 > 服务 > Tez > 配置”。
- 步骤3** 选择“全部配置”。
- 步骤4** 左边菜单栏中选择“TezUI > 日志”。
- 步骤5** 选择所需修改的日志级别。
- 步骤6** 单击“保存”，在弹出窗口中单击“确定”保存配置。
- 步骤7** 单击“实例”，勾选“TezUI”角色，选择“更多 > 重启实例”，输入用户密码后，在弹出窗口单击“确定”。
- 步骤8** 等待实例重启完成，配置生效。

----结束

日志格式

Tez的日志格式如下所示：

表 29-4 日志格式

日志类型	格式	示例
运行日志	<yyyy-MM-dd HH:mm:ss,SSS> <LogLevel> <产生该日志的 线程名字> <log中的 message> <日志事件的发生 位置>	2020-07-31 11:44:21,378 INFO TezUI-health-check Start health check com.XXX.tez.HealthCheck.run(HealthCheck.java:30)

日志类型	格式	示例
审计日志	<yyyy-MM-dd HH:mm:ss,SSS> <LogLevel> <产生该日志的 线程名字> <User Name><User IP><Time><Operation><Re source><Result><Detail > < 日志事件的发生位置>	2018-12-24 12:16:25,319 INFO HiveServer2-Handler- Pool: Thread-185 UserName=hive UserIP=10.153.2.204 Time=2018/12/24 12:16:25 Operation=CloseSession Result=SUCCESS Detail= org.apache.hive.service.cli.thrift. ThriftCLIService.logAuditEven t(ThriftCLIService.java:434)

29.4 Tez 常见问题

29.4.1 TezUI 无法展示 Tez 任务执行细节

问题

登录Manager界面，跳转Tez WebUI界面，已经提交的Tez任务未展示，如何解决。

回答

Tez WebUI展示的Tez任务数据，需要Yarn的TimelineServer支持，确认提交任务之前TimelineServer已经开启且正常运行。

在设置Hive执行引擎为Tez的同时，需要设置参数“yarn.timeline-service.enabled”为“true”，详情请参考[切换Hive执行引擎为Tez](#)。

29.4.2 进入 Tez WebUI 界面显示异常

问题

登录Manager界面，跳转Tez WebUI界面，显示404异常或503异常。

HTTP ERROR 404

Problem accessing /null/applicationhistory. Reason:

Not Found

Powered by Jetty:// 9.3.20.v20170531

Adapter operation failed Å» 503: Error accessing https://20026/Yarn/TimelineServer/57/ws/v1/timeline/TEZ_DAG_ID

回答

Tez WebUI依赖Yarn的TimelineServer实例，需要预先安装TimelineServer，且处于良好状态。

29.4.3 TezUI 界面无法查看 Yarn 日志

问题

登录Tez WebUI界面，单击Logs跳转yarn日志界面失败，无法加载数据。



无法访问此网站

找不到 **10-244-224-45** 的服务器 IP 地址。

请试试以下办法：

- 检查网络连接
- 检查代理服务器、防火墙和 DNS 配置
- 运行 Windows 网络诊断

ERR_NAME_NOT_RESOLVED

重新加载

回答

Tez WebUI跳转Yarn Logs界面时，目前是通过hostname进行访问，需要在windows机器，配置hostname到ip的映射。具体方法为：

修改windows机器C:\Windows\System32\drivers\etc\hosts文件，增加一行hostname到ip的映射，例：10.244.224.45 10-044-224-45，保存后重新访问正常。

29.4.4 TezUI HiveQueries 界面表格数据为空

问题

登录Manager界面，跳转Tez WebUI界面，已经提交的任务，Hive Queries界面未展示数据，如何解决。

回答

Tez WebUI展示的Hive Queries任务数据，需要设置以下3个参数：

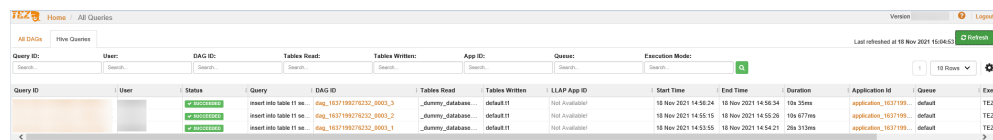
在FusionInsight Manager页面，选择“集群 > 服务 > Hive > 配置 > 全部配置 > HiveServer > 自定义”，在hive-site.xml中增加以下配置：

属性名	属性值
hive.exec.pre.hooks	org.apache.hadoop.hive.ql.hooks.ATSHook
hive.exec.post.hooks	org.apache.hadoop.hive.ql.hooks.ATSHook
hive.exec.failure.hooks	org.apache.hadoop.hive.ql.hooks.ATSHook

说明

TezUI数据展示依赖于Yarn组件的TimelineServer实例，如果TimelineServer实例故障或未启动，需设置hive自定义参数yarn-site.xml中**yarn.timeline-service.enabled=false**，否则hive任务会执行失败。

参数设置完成后，Hive Queries界面即可展示数据，但无法展示历史数据，展示效果如下：



30 使用 Yarn

30.1 Yarn 用户权限管理

30.1.1 创建 Yarn 角色

操作场景

该任务指导MRS集群管理员创建并设置Yarn的角色。Yarn角色可设置Yarn管理员权限以及Yarn队列资源管理。

📖 说明

如果当前组件使用了Ranger进行权限控制，须基于Ranger配置相关策略进行权限管理。对于MRS 3.x及后续版本集群，具体操作可参考[添加Yarn的Ranger访问权限策略](#)。

操作步骤

MRS 3.x以前版本集群执行以下操作：

步骤1 登录Manager，选择“系统设置 > 角色管理 > 添加角色”。

步骤2 在“角色名称”和“描述”输入角色名字与描述。

步骤3 设置角色“权限”请参见[表30-1](#)。

Yarn权限：

- “Cluster Admin Operations”：Yarn管理员权限。
- “Scheduler Queue”：队列资源管理。

表 30-1 设置角色

任务场景	角色授权操作
设置Yarn管理员权限	在“权限”的表格中选择“Yarn”，勾选“Cluster Admin Operations”。 说明 设置Yarn管理员权限需要重启Yarn服务，才能使保存的角色配置生效。
设置用户在指定Yarn队列提交任务的权限	1. 在“权限”的表格中选择“Yarn > Scheduler Queue”。 2. 在指定队列的“权限”列，勾选“Submit”。
设置用户在指定Yarn队列管理任务的权限	1. 在“权限”的表格中选择“Yarn > Scheduler Queue”。 2. 在指定队列的“权限”列，勾选“Admin”。

如果Yarn角色包含了某个父队列的“提交”或“管理”权限，则角色默认子队列也继承此权限，将自动添加子队列的“提交”或“管理”权限。子队列继承的权限不在“配置资源权限”表格显示被选中。

如果设置Yarn角色时仅勾选到某个父队列的“提交”权限，使用拥有该角色权限的用户提交任务时，注意需要手动指定队列名称，否则当父队列下有多个子队列时，系统并不会自动判断，从而将任务提交到了“default”队列。

步骤4 单击“确定”完成。

----结束

MRS 3.x及以后版本集群执行以下操作：

步骤1 登录Manager，选择“系统 > 权限 > 角色”。

步骤2 单击“添加角色”，然后“角色名称”和“描述”输入角色名字与描述。

步骤3 设置角色“配置资源权限”请参见[表30-2](#)。

Yarn权限：

- “集群管理操作权限”：Yarn管理员权限。
- “调度队列”：队列资源管理。

表 30-2 设置角色

任务场景	角色授权操作
设置Yarn管理员权限	在“配置资源权限”的表格中选择“待操作集群的名称 > Yarn”，勾选“集群管理操作权限”。 说明 设置Yarn管理员权限需要重启Yarn服务，才能使保存的角色配置生效。

任务场景	角色授权操作
设置用户在指定 Yarn 队列提交任务的权限	<ol style="list-style-type: none">1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Yarn > 调度队列 > root”。2. 在指定队列的“权限”列，勾选“提交”。
设置用户在指定 Yarn 队列管理任务的权限	<ol style="list-style-type: none">1. 在“配置资源权限”的表格中选择“待操作集群的名称 > Yarn > 调度队列 > root”。2. 在指定队列的“权限”列，勾选“管理”。

如果 Yarn 角色包含了某个父队列的“提交”或“管理”权限，则角色默认子队列也继承此权限，将自动添加子队列的“提交”或“管理”权限。子队列继承的权限不在“配置资源权限”表格显示被选中。

如果设置 Yarn 角色时仅勾选到某个父队列的“提交”权限，使用拥有该角色权限的用户提交任务时，注意需要手动指定队列名称，否则当父队列下有多个子队列时，系统并不会自动判断，从而将任务提交到了“default”队列。

步骤4 单击“确定”完成。

----结束

30.2 使用 Yarn 客户端提交任务

操作场景

该任务指导用户在运维场景或业务场景中使用 Yarn 客户端。

前提条件

- 已安装客户端。
例如安装目录为“/opt/client”，以下操作的客户端目录只是举例，请根据实际安装目录修改。
- 各组件业务用户由 MRS 集群管理员根据业务需要创建。安全模式下，“机机”用户需要下载 keytab 文件。“人机”用户第一次登录时需修改密码。普通模式不需要下载 keytab 文件及修改密码操作。

使用 Yarn 客户端

步骤1 安装客户端。

- MRS 3.x 之前版本请参考[安装客户端](#)章节。
- MRS 3.x 及之后版本请参考[安装客户端](#)章节。

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，切换到客户端安装目录。

```
cd /opt/client
```

步骤4 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤5 如果集群为安全模式，执行以下命令进行用户认证。普通模式集群无需执行用户认证。

```
kinit 组件业务用户
```

步骤6 直接执行Yarn命令。例如：

```
yarn application -list
```

----结束

客户端常见使用问题

- 问题一：当执行Yarn客户端命令时，客户端程序异常退出，报“java.lang.OutOfMemoryError”的错误。

这个问题是由于Yarn客户端运行时所需的内存超过了Yarn客户端设置的内存上限（默认为128MB）。

- a. 可以通过修改“<客户端安装路径>/HDFS/component_env”中的参数来修改Yarn客户端的内存上限。

- 对于MRS 3.x及之后版本集群，修改“CLIENT_GC_OPTS”参数。例如，需要设置该内存上限为1GB，则设置：

```
export CLIENT_GC_OPTS="-Xmx1G"
```

- 对于MRS 3.x之前版本集群，修改“GC_OPTS_YARN”参数。例如，需要设置该内存上限为1GB，则设置：

```
export GC_OPTS_YARN="-Xmx1G"
```

- b. 在修改完后，使用如下命令刷新客户端配置，使之生效：

```
source <客户端安装路径>/bigdata_env
```

- 问题二：如何设置Yarn客户端运行时的日志级别？

Yarn客户端运行时的日志默认输出到Console控制台，其级别默认为INFO级别。有时为了定位问题，需要开启DEBUG级别日志，可以通过导出环境变量来设置，命令如下：

```
export YARN_ROOT_LOGGER=DEBUG,console
```

在执行完以上命令后，再执行Yarn Shell命令时，即可打印出DEBUG级别日志。

如果想恢复INFO级别日志，可执行如下命令：

```
export YARN_ROOT_LOGGER=INFO,console
```

30.3 配置 Container 日志聚合功能

配置场景

Yarn提供了Container日志聚合功能，可以将各节点Container产生的日志收集到HDFS，释放本地磁盘空间。日志收集的方式有两种：

- 应用完成后将Container日志一次性收集到HDFS。
- 应用运行过程中周期性收集Container输出的日志片段到HDFS。

配置描述

参数入口：

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入[表 30-3](#)中参数名称，修改并保存配置。然后在Yarn服务“概览”页面选择“更多 > 同步配置”。同步完成后重启Yarn服务。

其中“yarn.nodemanager.remote-app-log-dir-suffix”参数还需要在Yarn的客户端进行配置，且在ResourceManager、NodeManager和JobHistory节点的配置与在Yarn的客户端的配置必须一致。

周期性收集日志功能目前仅支持MapReduce应用，且MapReduce应用必须进行相应的日志文件滚动输出配置，需要在MapReduce客户端节点的“客户端安装路径/Yarn/config/mapred-site.xml”配置文件中进行如[表30-5](#)所示的配置。

表 30-3 参数说明

参数	描述	默认值
yarn.log-aggregation-enable	<p>设置是否打开Container日志聚合功能。</p> <ul style="list-style-type: none"> 设置为“true”，表示打开该功能，日志会被收集到HDFS目录中。 设置为“false”，表示关闭该功能，表示日志不会收集到HDFS中。 <p>修改参数值后，需重启Yarn服务使其生效。</p> <p>说明</p> <ul style="list-style-type: none"> 在修改值为“false”并生效后，生效前的日志无法在WebUI中获取。 如果需要在WebUI界面上查看之前产生的日志，建议将此参数设置为“true”。 	true
yarn.nodemanager.log-aggregation.roll-monitoring-interval-seconds	<p>NodeManager周期性日志收集的时间间隔。</p> <ul style="list-style-type: none"> 设置为-1或0时，表示周期性收集日志功能关闭，日志在应用运行完成后一次性收集。 收集周期最小可设定为3600秒。当设置为大于0秒且小于3600秒时，收集周期将使用3600秒。 <p>定义NodeManager唤醒并上传日志的间隔周期。设置为-1或0表示禁用滚动监控，应用任务结束后日志汇聚。取值范围大于等于-1。</p>	-1

参数	描述	默认值
yarn.nodemanager.disk-health-checker.log-dirs.max-disk-utilization-per-disk-percentage	<p>配置Container日志目录可以占用每块磁盘上Yarn的磁盘配额的最大百分比。当日志目录占用空间超过此设定值时，将触发周期性日志收集服务启动一次周期外的日志收集活动，以释放本地磁盘空间。每个磁盘上可提供给Container logs的最大可使用率。当Container logs使用超过这个限制，会触发滚动汇聚。</p> <ul style="list-style-type: none"> 对于MRS 3.x之前的版本集群，磁盘配额最大百分比的有效取值范围为0~100，如果配置小于等于0，会被强制重置为25；如果配置大于100，则被强制重置为25。 对于MRS 3.x及后续版本集群，磁盘配额最大百分比的有效取值范围为-1~100，如果配置小于-1，会被强制重置为25；如果配置大于100，则被强制重置为25。而配置为-1时则关闭Container日志目录的磁盘容量检测功能。 <p>说明</p> <ul style="list-style-type: none"> Container日志目录实际可用磁盘百分比=YARN磁盘可用百分比（“yarn.nodemanager.disk-health-checker.max-disk-utilization-per-disk-percentage”）* 日志目录可用百分比（“yarn.nodemanager.disk-health-checker.log-dirs.max-disk-utilization-per-disk-percentage”）。 只有启用了周期性收集日志功能的应用才会在日志目录磁盘配额超过设定阈值时被触发启动日志收集。 	25
yarn.nodemanager.remote-app-log-dir-suffix	<p>设置HDFS用于存放Container日志的文件夹名称。该配置加上“yarn.nodemanager.remote-app-log-dir”，构成了Container日志的完整存放目录。目录为： “{yarn.nodemanager.remote-app-log-dir}/{user}/{yarn.nodemanager.remote-app-log-dir-suffix}”。</p> <p>说明 <i>{user}</i>为运行任务时的用户名。</p>	logs
yarn.nodemanager.log-aggregator.on-fail.remain-log-in-sec	<p>设置Container日志归集失败后日志在本地保留的时间。单位：秒。</p> <ul style="list-style-type: none"> 设置为0时，本地日志将马上删除。 设置为正数时，表示本地日志将保留这段时间。 	604800

参考[修改集群服务配置参数](#)进入Mapreduce服务参数“全部配置”界面，在搜索框中输入[表30-4](#)中参数名称。

表 30-4 参数说明

参数	描述	默认值
yarn.log-aggregation.retain-seconds	<p>汇聚日志的保存时间。单位：秒。</p> <ul style="list-style-type: none"> • 设置为-1时，表示HDFS上面的Container聚合日志将永久保留。 • 设置为0或正数时，表示HDFS上面的Container聚合日志将保留这段时间，超时将被删除。 <p>说明 当时间设置太短时，有可能会增加NameNode的负担，建议根据实际情况设置一个合理的时间值。</p>	1296000
yarn.log-aggregation.retain-check-interval-seconds	<p>设置扫描HDFS保存的Container聚合日志的间隔时间。单位：秒。</p> <ul style="list-style-type: none"> • 设置为-1或0时，间隔时间将为“yarn.log-aggregation.retain-seconds”该配置时间的十分之一。 <p>说明 当该配置设置为-1或0时，“yarn.log-aggregation.retain-seconds”不能设置为0。 • 设置为正数时，将周期性的间隔这段时间以后对HDFS上的container聚合日志进行扫描。 <p>说明 当时间设置太短时，有可能会增加NameNode的负担，建议根据实际情况设置一个合理的时间。</p> </p>	86400

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入表 30-5中参数名称。

表 30-5 MapReduce 应用日志文件滚动输出配置

参数	描述	默认值
mapreduce.task.userlog.limit.kb	MR应用程序单个task日志文件大小限制。当日志文件达到该限制时，会新建一个日志文件进行输出。设置为“0”表示不限制日志文件大小。	51200

参数	描述	默认值
yarn.app.mapreduce.task.container.log.backups	MR应用程序task日志保留的最大个数。 设置为“0”表示不滚动输出。 使用CRLA（ContainerRollingLogAppender）时任务日志备份文件的数量。默认使用CLA（ContainerLogAppender）且container日志不回滚。 当mapreduce.task.userlog.limit.kb和yarn.app.mapreduce.task.container.log.backups都大于0时，任务启用CRLA。取值范围0~999。	10
yarn.app.mapreduce.am.container.log.limit.kb	MR应用程序单个AM日志文件大小限制。单位：KB，当日志文件达到该限制时，会新建一个日志文件进行输出。设置为“0”表示不限制单个AM日志文件大小。	51200
yarn.app.mapreduce.am.container.log.backups	MR应用程序AM日志保留的最大个数。设置为“0”表示不滚动输出。使用CRLA（ContainerRollingLogAppender）时ApplicationMaster日志备份文件的数量。默认使用CLA（ContainerLogAppender）且容器日志不回滚。 当yarn.app.mapreduce.am.container.log.limit.kb和yarn.app.mapreduce.am.container.log.backups都大于0时，ApplicationMaster启用CRLA。取值范围0~999。	20
yarn.app.mapreduce.shuffle.log.backups	MR应用程序shuffle日志保留的最大个数。设置为“0”表示不滚动输出。 当yarn.app.mapreduce.shuffle.log.limit.kb和yarn.app.mapreduce.shuffle.log.backups都大于0时，syslog.shuffle将采用CRLA。取值范围0~999。	10
yarn.app.mapreduce.shuffle.log.limit.kb	MR应用程序单个shuffle日志文件大小限制，单位KB。当日志文件达到该限制时，会新建一个日志文件进行输出。设置为“0”不限制单个shuffle日志文件大小。取值范围大于等于0。	51200

30.4 启用 Yarn CGroups 功能限制 Container CPU 使用率

配置场景

CGroups是一个Linux内核特性。它可以将任务集及其子集聚合或分离成具备特定行为的分层组。在Yarn中，CGroups特性对容器（Container）使用的资源（例如CPU使用率）进行限制。本特性大大降低了限制容器CPU使用的难度。

说明

当前CGroups仅用于限制CPU使用率。

本章节适用于MRS 3.x及后续版本集群。

配置描述

有关如何配置CPU隔离与安全的CGroups功能的详细信息，请参见Hadoop官网：

MRS 3.2.0之前版本：<http://hadoop.apache.org/docs/r3.1.1/hadoop-yarn/hadoop-yarn-site/NodeManagerCgroups.html>

MRS 3.2.0及之后版本：<https://hadoop.apache.org/docs/r3.3.1/hadoop-yarn/hadoop-yarn-site/NodeManagerCgroups.html>

由于CGroups为Linux内核特性，是通过LinuxContainerExecutor进行开放。请参考官网资料对LinuxContainerExecutor进行安全配置。您可通过官网资料了解系统用户和用户组配置对应的文件系统权限。详情请参见：

MRS 3.2.0之前版本：<http://hadoop.apache.org/docs/r3.1.1/hadoop-project-dist/hadoop-common/SecureMode.html#LinuxContainerExecutor>

MRS 3.2.0及之后版本：<https://hadoop.apache.org/docs/r3.3.1/hadoop-project-dist/hadoop-common/SecureMode.html#LinuxContainerExecutor>

说明

- 请勿修改对应文件系统中各路径所属的用户、用户组及对应的权限，否则可能导致本功能异常。
- 当参数“yarn.nodemanager.resource.percentage-physical-cpu-limit”配置过小，导致可使用的核不足1个时，例如4核节点，将此参数设置为20%，不足1个核，那么将会使用系统全部的核。Linux的一些版本不支持Quota模式，例如Cent OS。在这种情况下，可以使用CPUset模式。

配置cpuset模式，即Yarn只能使用配置的CPU，需要在Manager界面添加以下配置。

表 30-6 cpuset 配置

参数	描述	默认值
yarn.nodemanager.linux-container-executor.cgroups.cpu-set-usage	设置为“true”时，应用以cpuset模式运行。	false

配置strictcpuset模式，即Container只能使用配置的CPU，需要在Manager界面添加以下配置。

表 30-7 CPU 硬隔离参数配置

参数	描述	默认值
yarn.nodemanager.linux-container-executor.cgroups.cpu-set-usage	设置为“true”时，应用以cpuset模式运行。	false
yarn.nodemanager.linux-container-executor.cgroups.cpuset.strict.enabled	设置为true时，Container只能使用配置的CPU。	false

要从cpuset模式切换到Quota模式，必须遵循以下条件：

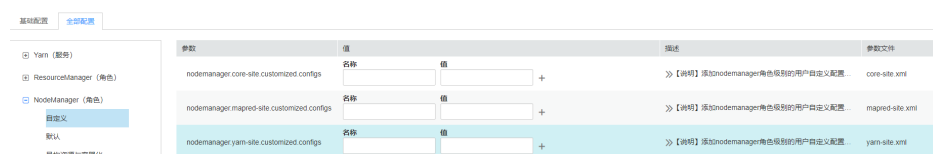
- 配置“yarn.nodemanager.linux-container-executor.cgroups.cpu-set-usage” = “false”。
- 删除“/sys/fs/cgroup/cpuset/hadoop-yarn/”路径下container文件夹（如果存在）。
- 删除“/sys/fs/cgroup/cpuset/hadoop-yarn/”路径下cpuset.cpus文件中设置的所有CPU。

操作步骤

步骤1 登录Manager系统。选择“集群 > 服务 > Yarn > 配置”，选择“全部配置”。

步骤2 在左侧导航栏选择“NodeManager > 自定义”，找到yarn-site.xml文件。

步骤3 添加表30-6和表30-7中的参数为自定义参数。



根据配置文件与参数作用，在“yarn-site.xml”所在行“名称”列输入参数名，在“值”列输入此参数的参数值。

单击“+”增加自定义参数。

步骤4 单击“保存”，在弹出的“保存配置”窗口中确认修改参数，单击“确定”。界面提示“操作成功”，单击“完成”，配置保存成功。

保存完成后请重新启动配置过期的Yarn服务以使配置生效。

----结束

30.5 Yarn 企业级能力增强

30.5.1 配置 Yarn 权限控制开关

配置场景

在安全模式的多租户场景下，一个集群可以支持多个用户使用以及支持多个用户任务提交、运行，用户之间不可见，需要有一个权限控制机制，使用户的任务信息不被其他用户获取。

例如，用户A提交的应用正在运行，此时用户B登录系统并查看应用列表，用户B不应该访问到A用户的应用信息。

配置描述

- 查看Yarn服务配置参数
参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入[表30-8](#)中参数名称。

表 30-8 参数描述

参数	描述	默认值
yarn.acl.enable	Yarn权限控制启用开关。	true
yarn.webapp.filter-entity-list-by-user	严格视图启用开关，开启后，登录用户只能查看该用户有权限查看的内容。当要开启该功能时，同时需要设置参数“yarn.acl.enable”为true。 说明 此参数适用于MRS 3.x及后续版本集群。	true

- 查看MapReduce服务配置参数
参考[修改集群服务配置参数](#)进入MapReduce服务参数“全部配置”界面，在搜索框中输入[表30-9](#)中参数名称。

表 30-9 参数描述

参数	描述	默认值
mapreduce.cluster.acls.enabled	MR JobHistoryServer权限控制启用开关。该参数为客户端参数，当JobHistoryServer服务端开启权限控制之后该参数生效。	true

参数	描述	默认值
yarn.webapp.filter-entity-list-by-user	MR JobHistoryServer严格视图启用开关，开启后，登录用户只能查看该用户有权限查看的内容。该参数为JobHistoryServer的服务端参数，表示JHS开启了权限控制，但是否要对某一个特定的Application进行控制，是由客户端参数：“mapreduce.cluster.acls.enabled”决定。 说明 此参数适用于MRS 3.x及后续版本集群。	true

须知

以上配置会影响restful API和shell命令结果，即以上配置开启后，restful API调用和shell命令运行所返回的内容只包含调用用户有权查看的信息。

当“yarn.acl.enable”或“mapreduce.cluster.acls.enabled”设置为“false”时，即关闭Yarn或MapReduce的权限校验功能。此时任何用户都可以在Yarn或MapReduce上提交任务和查看任务信息，存在安全风险，请谨慎使用。

30.5.2 手动指定运行 Yarn 任务的用户

本章节适用于MRS 3.x及后续版本集群。

配置场景

目前Yarn支持启动NodeManager的用户运行所有用户提交的任务，也支持以提交任务的用户运行任务。

配置描述

在Manager系统中，选择“集群 > 服务 > Yarn > 配置”，选择“全部配置”。在搜索框中输入参数名称。

表 30-10 参数描述

参数	描述	默认值
yarn.nodemanager.linux-container-executor.user	运行任务的用户。	默认为空。 说明 默认为空，实际以提交任务的用户来运行任务。
yarn.nodemanager.container-executor.class	启动任务的executor。	org.apache.hadoop.yarn.server.nodemanager.EnhancedLinuxContainerExecutor

说明

- “yarn.nodemanager.linux-container-executor.user”配置运行Container的用户。默认空表示运行Container的用户就是提交任务的用户。该参数仅在“yarn.nodemanager.container-executor.class”配置为“org.apache.hadoop.yarn.server.nodemanager.EnhancedLinuxContainerExecutor”时有效。
- 非安全模式下，当“yarn.nodemanager.linux-container-executor.user”设置为omm时，也需设置“yarn.nodemanager.linux-container-executor.nonsecure-mode.local-user”为omm。
- 建议“yarn.nodemanager.linux-container-executor.user”和“yarn.nodemanager.container-executor.class”这两个参数都采用默认值，这样安全性更高。

30.5.3 配置 AM 失败重试次数

配置场景

在资源不足导致ApplicationMaster启动失败的情况下，调整如下参数值，提高容错性，保证客户端应用的正常运行。

配置描述

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入[表 30-11](#)中参数名称。

表 30-11 参数说明

参数	描述	默认值
yarn.resource manager.am.max-attempts	ApplicationMaster重试次数，增加重试次数，可以防止资源不足导致的ApplicationMaster启动失败问题。适用于所有ApplicationMaster的全局设置。每个ApplicationMaster都可以使用API设置一个单独的最大尝试次数，但这个次数不能大于全局的最大次数。如果大于，ResourceManager将会覆写这个单独的最大尝试次数。以允许至少一次重试。取值范围大于等于1。	5

30.5.4 配置 AM 自动调整分配内存

本章节适用于MRS 3.x及后续版本集群。

配置场景

启动该配置的过程中，ApplicationMaster在创建Container时，分配的内存会根据任务总数的浮动自动调整，资源利用更加灵活，提高了客户端应用运行的容错性。

配置描述

参数入口：

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称“mapreduce.job.am.memory.policy”。

配置说明：

配置项的默认值为空，此时不会启动自动调整的策略，ApplicationMaster的内存仍受“yarn.app.mapreduce.am.resource.mb”配置项的影响。

配置参数的值由5个数值组成，中间使用“:”与“,”分隔，格式为：

baseTaskCount:taskStep:memoryStep,minMemory:maxMemory，在键入时会严格校验格式。

表 30-12 配置数值说明

数值名称	描述	设定要求
baseTaskCount	任务总量基数，只有当应用的task总数（map端与reduce端之和）不小于该值时配置才会起作用。	不能为空且大于零。
taskStep	任务增量步进，与memoryStep共同决定内存调整量。	不能为空且大于零。
memoryStep	内存增量步进，在“yarn.app.mapreduce.am.resource.mb”配置的基础上对内存向上调整。	不能为空且大于零，单位：MB。
minMemory	内存自动调整下限，若调整后的内存不大于该值，仍保持“yarn.app.mapreduce.am.resource.mb”的配置。	不能为空且大于零，且不大于maxMemory的设定值。 单位：MB
maxMemory	内存自动调整上限，若调整后的内存超过该值，则使用该值作为最终调整值。	不能为空且大于零，且不小于minMemory的设定值。 单位：MB

配置示例

配置情况：

- yarn.app.mapreduce.am.resource.mb=1536
- mapreduce.job.am.memory.policy=100:10:50,1200:2000
- 某应用task总数=120

计算过程：

调整后内存=1536+[(120-100) /10]*50=1636，满足1200<1636且2000>1636，最终ApplicationMaster内存会设定为1636MB。

若memStep修改为250，调整后内存=1536+[(120-100) /10]*250=2136，超过maxMemory=2000的限制，最终ApplicationMaster内存会设定为2000MB。

说明

对于计算后的调整值低于设定的“minMemory”值的情形，虽然此时配置不会生效但后台仍然会打印出这个调整值，用于为用户提供“minMemory”参数调整的依据，保证配置可以生效。

30.5.5 配置 AM 作业自动保留

本章节适用于MRS 3.x及后续版本集群。

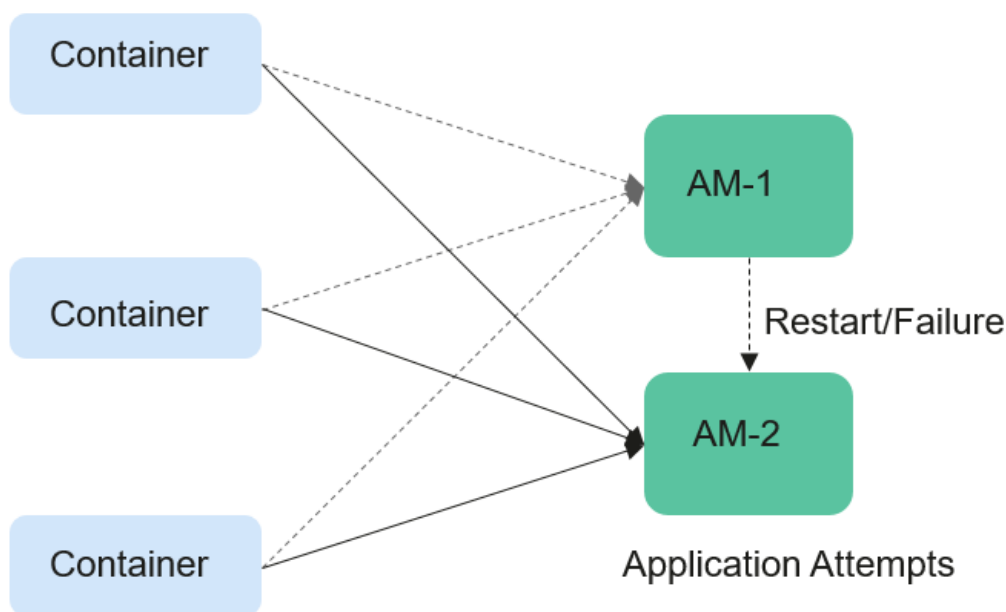
配置场景

在YARN中，ApplicationMaster(AM)与Container类似，都运行在NodeManager(NM)上（本文中忽略未管理的AM）。AM可能由于多种原因崩溃、退出或关闭。如果AM停止运行，ResourceManager(RM)会关闭ApplicationAttempt中管理的所有Container，其中包括当前在NM上运行的所有Container。RM会在另一计算节点上启动新的ApplicationAttempt。

对于不同类型的應用，希望以不同方式处理AM重启的事件。MapReduce类应用的目标是不丢失任务，但允许丢失当前运行的Container。但是对于长周期的YARN服务而言，用户可能并不希望由于AM的故障而导致整个服务停止运行。

YARN支持在新的ApplicationAttempt启动时，保留之前Container的状态，因此运行中的作业可以继续无故障的运行。

图 30-1 AM 作业保留



配置描述

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。

根据[表30-13](#)，对如下参数进行设置。

表 30-13 AM 作业保留相关参数

参数	说明	默认值
yarn.app.mapreduce.am.work-preserve	是否开启AM作业保留特性。	false
yarn.app.mapreduce.am.umbilical.max.retries	AM作业保留特性中，运行的容器尝试恢复的最大次数。	5
yarn.app.mapreduce.am.umbilical.retry.interval	AM作业保留特性中，运行的容器尝试恢复的时间间隔。单位：毫秒。	10000
yarn.resourcemanager.am.max-attempts	ApplicationMaster的重试次数。增加重试次数可以避免当资源不足时造成AM启动失败。 适用于所有ApplicationMaster的全局设置。每个ApplicationMaster都可以使用API设置一个单独的最大尝试次数，但这个次数不能大于全局的最大次数。如果大于了，那ResourceManager将会覆写这个单独的最大尝试次数。取值范围大于等于1。	2

30.5.6 配置 Yarn 数据访问通道协议

配置场景

服务端配置了web访问为https通道，如果客户端没有配置，默认使用http访问，客户端和服务端的配置不同，就会导致访问结果显示乱码。在客户端和服务端配置相同的“yarn.http.policy”参数，可以防止客户端访问结果显示乱码。

操作步骤

步骤1 在Manager系统中，选择“集群 > 服务 > Yarn > 配置”，选择“全部配置”，在搜索框中输入参数名称“yarn.http.policy”。

- 安全模式下配置为“HTTPS_ONLY”。
- 普通模式下配置为“HTTP_ONLY”。

步骤2 以客户端安装用户，登录安装客户端的节点。

步骤3 执行以下命令，进入客户端安装路径。

```
cd /opt/client
```

步骤4 执行以下命令编辑“yarn-site.xml”文件。

```
vi Yarn/config/yarn-site.xml
```

修改“yarn.http.policy”的参数值。

安全模式下，“yarn.http.policy”配置成“HTTPS_ONLY”。

普通模式下，“yarn.http.policy”配置成“HTTP_ONLY”。

步骤5 执行:wq命令保存。

步骤6 重启客户端使配置生效。

----结束

30.5.7 配置自定义调度器的 WebUI

配置场景

如果用户在ResourceManager中配置了自定义的调度器，可以通过以下配置项为其配置相应的Web展示页面及其他Web应用。

配置描述

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。

表 30-14 配置自定义调度器的 WebUI

参数	描述	默认值
hadoop.http.rmwebapp.scheduler.page.classes	在RM WebUI中为自定义调度器加载相应的web页面。仅当“yarn.resourcemanager.scheduler.class”配置为自定义调度器时此配置项生效。	-
yarn.http.rmwebapp.external.classes	在RM的Web服务中加载用户自定义的web应用。	-

30.5.8 配置 NodeManager 角色实例使用的资源

操作场景

如果部署NodeManager的各个节点硬件资源（如CPU核数、内存总量）不一样，而NodeManager可用硬件资源设置为相同的值，可能造成性能浪费或状态异常，需要修改各个NodeManager角色实例的配置，使硬件资源得到充分利用。

对系统的影响

保存新的配置需要重启NodeManager角色实例，此时对应的角色实例不可用。

操作步骤

MRS 3.x之前的版本集群执行以下操作：

步骤1 登录MRS控制台，选择“现有集群”，单击集群名称。选择“组件管理 > Yarn > 实例”。

步骤2 单击“角色”列“NodeManager”角色实例名称，并切换到“实例配置”。单击“基础配置”下拉菜单，选择“全部配置”，在搜索框中输入以下参数。

- 步骤3** “yarn.nodemanager.resource.cpu-vcores”设置当前节点上NodeManager可使用的虚拟CPU核数，建议按节点实际逻辑核数的1.5到2倍配置。
“yarn.nodemanager.resource.memory-mb”设置当前节点上NodeManager可使用的物理内存大小，建议按节点实际物理内存大小的75%~90%配置。

 **说明**

“yarn.scheduler.maximum-allocation-vcores”可配置单个Container最多CPU可用核数，
“yarn.scheduler.maximum-allocation-mb”可配置单个Container最大内存可用值。不支持实例级别的修改，需要在Yarn服务的配置中修改参数值，并重启Yarn服务。

- 步骤4** 单击“保存配置”，勾选“重新启动受影响的服务或实例”，单击“确定”。重启NodeManager角色实例。

界面提示“操作成功。”，单击“完成”，NodeManager角色实例成功启动。

----**结束**

MRS 3.x及后续版本集群也可执行以下操作：

- 步骤1** 登录FusionInsight Manager，选择“集群 > 待操作集群的名称 > 服务 > Yarn > 实例”。

- 步骤2** 单击部署NodeManager节点对应角色实例名称，并切换到“实例配置”，选择“全部配置”。

- 步骤3** “yarn.nodemanager.resource.cpu-vcores”设置当前节点上NodeManager可使用的虚拟CPU核数，建议按节点实际逻辑核数的1.5到2倍配置。
“yarn.nodemanager.resource.memory-mb”设置当前节点上NodeManager可使用的物理内存大小，建议按节点实际物理内存大小的75%配置。

 **说明**

“yarn.scheduler.maximum-allocation-vcores”可配置单个Container最多CPU可用核数，
“yarn.scheduler.maximum-allocation-mb”可配置单个Container最大内存可用值。不支持实例级别的修改，需要在Yarn服务的配置中修改参数值，并重启Yarn服务。

- 步骤4** 单击“保存”，单击“确定”。重启NodeManager角色实例。

界面提示“操作成功”，单击“完成”，NodeManager角色实例成功启动。

----**结束**

30.5.9 配置 ResourceManager 重启后自动加载 Container 信息

配置场景

YARN Restart特性包含两部分内容：ResourceManager Restart和NodeManager Restart。

- 当启用ResourceManager Restart时，升主后的ResourceManager就可以通过加载之前的主ResourceManager的状态信息，并通过接收所有NodeManager上container的状态信息，重构运行状态继续执行。这样应用程序通过定期执行检查点操作保存当前状态信息，就可以避免工作内容的丢失。
- 当启用NodeManager Restart时，NodeManager在本地保存当前节点上运行的container信息，重启NodeManager服务后通过恢复此前保存的状态信息，就不会丢失在此节点上运行的container进度。

配置描述

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。

ResourceManager Restart特性配置如下。

表 30-15 ResourceManager Restart 参数配置

参数	描述	默认值
yarn.resourcemanager.recovery.enabled	设置是否让ResourceManager在启动后恢复状态。如果设置为true，那yarn.resourcemanager.store.class也必须设置。	true
yarn.resourcemanager.store.class	指定用于保存应用程序和任务状态以及证书内容的state-store类。	MRS 3.x之前的版本集群： org.apache.hadoop.yarn.server.resourcemanager.recovery.ZKRMStateStore MRS 3.x及后续版本集群： org.apache.hadoop.yarn.server.resourcemanager.recovery.AsyncZKRMStateStore
yarn.resourcemanager.zk-state-store.parent-path	ZKRMStateStore在ZooKeeper上的保存目录。	/rmstore
yarn.resourcemanager.work-preserving-recovery.enabled	启用ResourceManager Work preserving功能。该配置仅用于YARN特性验证。	true
yarn.resourcemanager.state-store.async.load	对已完成的application采用ResourceManager异步恢复方式。	MRS 3.x之前的版本集群：false MRS 3.x及后续版本集群：true
yarn.resourcemanager.zk-state-store.num-fetch-threads	启用异步恢复功能，增加工作线程的数量可以加快恢复ZK中保存的任务信息的速度，取值范围大于0。	MRS 3.x之前的版本集群：1 MRS 3.x及后续版本集群：20

NodeManager Restart特性配置如下。

表 30-16 NodeManager Restart 参数配置

参数	描述	默认值
yarn.nodemanager.recovery.enabled	当Nodemanager重启时是否启用日志失败收集功能，是否恢复未完成的Application。	true
yarn.nodemanager.recovery.dir	NodeManager用于保存container状态的本地目录。适用于MRS 3.x及后续版本集群。	\${SRV_HOME}/tmp/yarn-nm-recovery
yarn.nodemanager.recovery.supervised	NodeManager是否在监控下运行。开启此特性后NodeManager在退出后不会清理containers，NodeManager会假设自己会立即重启和恢复containers。	true

30.6 Yarn 性能调优

30.6.1 调整 Yarn 任务抢占机制

操作场景

抢占任务可精简队列中的job运行并提高资源利用率，由ResourceManager的capacity scheduler实现，其简易流程如下：

1. 假设存在两个队列A和B。其中队列A的capacity为25%，队列B的capacity为75%。
2. 初始状态下，任务1发送给队列A，此任务需要75%的集群资源。之后任务2发送到了队列B，此任务需要50%的集群资源。
3. 任务1将会使用队列A提供的25%的集群资源，并从队列B获取的50%的集群资源。队列B保留25%的集群资源。
4. 启用抢占任务特性，则任务1使用的资源将会被抢占。队列B会从队列A中获取25%的集群资源以满足任务2的执行。
5. 当任务2完成后，集群中存在足够的资源时，任务1将重新开始执行。

操作步骤

参数入口：

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。

表 30-17 Preemption 配置

参数	描述	默认值
yarn.resourcemanager.scheduler.monitor.enable	根据“yarn.resourcemanager.scheduler.monitor.policies”中的策略，启用新的scheduler监控。设置为“true”表示启用监控，并根据scheduler的信息，启动抢占的功能。设置为“false”表示不启用。	false
yarn.resourcemanager.scheduler.monitor.policies	设置与scheduler配合的“SchedulingEditPolicy”的类的清单。	org.apache.hadoop.yarn.server.resourcemanager.monitor.capacity.ProportionalCapacityPreemptionPolicy
yarn.resourcemanager.monitor.capacity.preemption.observe_only	<ul style="list-style-type: none"> 设置为“true”，则执行策略，但是不对集群资源进程抢占操作。 设置为“false”，则执行策略，且根据策略启用集群资源抢占的功能。 	false
yarn.resourcemanager.monitor.capacity.preemption.monitoring_interval	根据策略监控的时间间隔，单位为毫秒。如果将该参数设置为更大的值，容量检测将不那么频繁地运行。	3000
yarn.resourcemanager.monitor.capacity.preemption.max_wait_before_kill	应用发送抢占需求到停止container（释放资源）的时间间隔，单位为毫秒。取值范围大于等于0。 默认情况下，若ApplicationMaster15秒内没有终止container，ResourceManager等待15秒后会强制终止。	15000
yarn.resourcemanager.monitor.capacity.preemption.total_preemption_per_round	在一个周期内能够抢占资源的最大的比例。可使用这个值来限制从集群回收容器的速度。计算出了期望的总抢占值之后，策略会伸缩回这个限制。	0.1
yarn.resourcemanager.monitor.capacity.preemption.max_ignored_over_capacity	集群中资源总量乘以此配置项的值加上某个队列（例如队列A）原有的资源量为资源抢占盲区。当队列A中的任务实际使用的资源超过该抢占盲区时，超过部分的资源将会被抢占。取值范围：0~1。 说明 设置的值越小越有利于资源抢占。	0

参数	描述	默认值
yarn.resourcemanager.monitor.capacity.preemption.natural_termination_factor	<p>设置抢占目标，Container只会抢占所配置比例的资源。</p> <p>示例，如果设置为0.5，则在5*“yarn.resourcemanager.monitor.capacity.preemption.max_wait_before_kill”的时间内，任务会回收所抢占资源的近95%。即接连抢占5次，每次抢占待抢占资源的0.5，呈几何收敛，每次的时间间隔为“yarn.resourcemanager.monitor.capacity.preemption.max_wait_before_kill”。</p> <p>取值范围：0~1。</p>	1

30.6.2 手动配置 Yarn 任务优先级

操作场景

集群的资源竞争场景如下：

1. 提交两个低优先级的应用Job 1和Job 2。
2. 正在运行中的Job 1和Job 2有部分task处于running状态，但由于集群或队列资源容量有限，仍有部分task未得到资源而处于pending状态。
3. 提交一个较高优先级的应用Job 3，此时会出现如下资源分配情况：当Job 1和Job 2中running状态的task运行结束并释放资源后，Job 3中处于pending状态的task将优先得到这部分新释放的资源。
4. Job 3完成后，资源释放给Job 1、Job 2继续执行。

用户可以在YARN中配置任务的优先级。任务优先级是通过ResourceManager的调度器实现的。

操作步骤

设置参数“mapreduce.job.priority”，使用命令行接口或API接口设置任务优先级。

- 命令行接口。
提交任务时，添加“-Dmapreduce.job.priority=<priority>”参数。
<priority>可以设置为：
 - VERY_HIGH
 - HIGH
 - NORMAL
 - LOW
 - VERY_LOW
- API接口。
用户也可以使用API配置对象的优先级。

设置优先级，可通过`Configuration.set("mapreduce.job.priority", <priority>)`或`Job.setPriority(JobPriority priority)`设置。

30.6.3 Yarn 节点配置调优

操作场景

合理配置大数据集群的调度器后，还可通过调节每个节点的可用内存、CPU资源及本地磁盘的配置进行性能调优。

具体包括以下配置项：

- 可用内存
- CPU虚拟核数
- 物理CPU使用百分比
- 内存和CPU资源的协调
- 本地磁盘

操作步骤

若您需要对参数配置进行调整，具体操作请参考[修改集群服务配置参数](#)。

- **可用内存**

除了分配给操作系统、其他服务的内存外，剩余的资源应尽量分配给YARN。通过如下配置参数进行调整。

例如，如果一个container默认使用512M，则内存使用的计算公式为：
 $512M * \text{container数}$ 。

默认情况下，Map或Reduce container会使用1个虚拟CPU内核和1024MB内存，ApplicationMaster使用1536MB内存。

参数	描述	默认值
yarn.nodemanager.resourcement.memory-mb	设置可分配给容器的物理内存数量。单位：MB，取值范围大于0。 建议配置成节点物理内存总量的75%~90%。若该节点有其他业务的常驻进程，请降低此参数值给该进程预留足够运行资源。如果节点的总物理内存空间较大，且无其他业务的常驻进程时，该参数可配置为：总物理内存 - NodeManager的常驻进程所占内存。	MRS 3.x及之后：16384 MRS 3.x之前：8192

- **CPU虚拟核数**

建议将此配置设定在逻辑核数的1.5~2倍之间。如果上层计算应用对CPU的计算能力要求不高，可以配置为2倍的逻辑CPU。

参数	描述	默认值
yarn.nodemanager.resource.cpu-vcores	<p>表示该节点上YARN可使用的虚拟CPU个数，默认是8。</p> <p>目前推荐将该值设置为逻辑CPU核数的1.5~2倍之间。</p> <ul style="list-style-type: none"> • 若任务为计算密集型，该参数可设置为与逻辑CPU核数一致。 • 若任务为非计算密集型资源，该参数可设置为逻辑CPU核数的1.5~2倍之间。 • 若任务所使用的CPU核数与内存资源差异较大时，CPU资源可参考实际的内存资源进行配置。例如大部分任务使用1核3G，如果“yarn.nodemanager.resource.memory-mb”设置380G，那么该参数设置为128。 	8

• **物理CPU使用百分比**

建议预留适量的CPU给操作系统和其他进程（数据库、HBase等）外，剩余的CPU核都分配给YARN。可以通过如下配置参数进行调整。

参数	描述	默认值
yarn.nodemanager.resource.percentage-physical-cpu-limit	<p>表示该节点上YARN可使用的物理CPU百分比。默认是90，即不进行CPU控制，YARN可以使用节点全部CPU。该参数只支持查看，可通过调整YARN的RES_CPUSSET_PERCENTAGE参数来修改本参数值。注意，目前推荐将该值设为可供YARN集群使用的CPU百分数。</p> <p>例如：当前节点除了YARN服务外的其他服务（如HBase、HDFS、Hive等）及系统进程使用CPU为20%左右，则可以供YARN调度的CPU为1-20%=80%，即配置此参数为80。</p>	90

• **本地磁盘**

由于本地磁盘会提供给MapReduce写job执行的中间结果，数据量大。因此配置的原则是磁盘尽量多，且磁盘空间尽量大，单个达到百GB以上规模更好。简单的做法是配置和data node相同的磁盘，只在最下一级目录上不同即可。

 **说明**

多个磁盘之间使用逗号隔开。

参数	描述	默认值
yarn.nodemanager.log-dirs	<p>日志存放地址（可配置多个目录）。</p> <p>容器日志的存储位置。默认值为%{@auto.detect.datapart.nm.logs}。如果有数据分区，基于该数据分区生成一个类似/srv/BigData/hadoop/data1/nm/containerlogs,/srv/BigData/hadoop/data2/nm/containerlogs的路径清单。如果没有数据分区，生成默认路径/srv/BigData/yarn/data1/nm/containerlogs。除了使用表达式以外，还可以输入完整的路径清单，比如/srv/BigData/yarn/data1/nm/containerlogs或/srv/BigData/yarn/data1/nm/containerlogs,/srv/BigData/yarn/data2/nm/containerlogs。这样数据就会存储在所有设置的目录中，一般会是在不同的设备中。为保证磁盘IO负载均衡，需要提供几个路径且每个路径都对应一个单独的磁盘。应用程序的本地化后的日志目录存在于相对路径/application_%{appid}中。单独容器的日志目录，即container_{\$contid}，是该路径下的子目录。每个容器目录都含容器生成的stderr、stdin及syslog文件。要新增目录，比如新增/srv/BigData/yarn/data2/nm/containerlogs目录，应首先删除/srv/BigData/yarn/data2/nm/containerlogs下的文件。之后，为/srv/BigData/yarn/data2/nm/containerlogs赋予跟/srv/BigData/yarn/data1/nm/containerlogs一样的读写权限，再将/srv/BigData/yarn/data1/nm/containerlogs修改为/srv/BigData/yarn/data1/nm/containerlogs,/srv/BigData/yarn/data2/nm/containerlogs。可以新增目录，但不要修改或删除现有目录。否则，NodeManager的数据将丢失，且服务将不可用。</p> <p>【默认值】%{@auto.detect.datapart.nm.logs}</p> <p>【注意】请谨慎修改该项。如果配置不当，将造成服务不可用。当角色级别的该配置项修改后，所有实例级别的该配置项都将被修改。如</p>	%{@auto.detect.datapart.nm.logs}

参数	描述	默认值
	果实例级别的配置项修改后，其他实例的该配置项的值保持不变。	

参数	描述	默认值
yarn.nodemanager.local-dirs	<p>本地化后的文件的存储位置。默认值为%</p> <p>{@auto.detect.datapart.nm.localdir}。如果有数据分区，基于该数据分区生成一个类似/srv/BigData/hadoop/data1/nm/localdir,/srv/BigData/hadoop/data2/nm/localdir的路径清单。如果没有数据分区，生成默认路径/srv/BigData/yarn/data1/nm/localdir。除了使用表达式以外，还可以输入完整的路径清单，比如/srv/BigData/yarn/data1/nm/localdir或/srv/BigData/yarn/data1/nm/localdir,/srv/BigData/yarn/data2/nm/localdir。这样数据就会存储在所有设置的目录中，一般会是在不同的设备中。为保证磁盘IO负载均衡，需要提供几个路径且每个路径都对应一个单独的磁盘。应用程序的本地化后的文件目录存在于相对路径/usercache/{user}/appcache/application_{appid}中。单独容器的工作目录，即container_{contid}，是该路径下的子目录。要新增目录，比如新增/srv/BigData/yarn/data2/nm/localdir目录，应首先删除/srv/BigData/yarn/data2/nm/localdir下的文件。之后，为/srv/BigData/hadoop/data2/nm/localdir赋予跟/srv/BigData/hadoop/data1/nm/localdir一样的读写权限，再将/srv/BigData/yarn/data1/nm/localdir修改为/srv/BigData/yarn/data1/nm/localdir,/srv/BigData/yarn/data2/nm/localdir。可以新增目录，但不要修改或删除现有目录。否则，NodeManager的数据将丢失，且服务将不可用。</p> <p>【默认值】% {@auto.detect.datapart.nm.localdir}</p> <p>【注意】请谨慎修改该项。如果配置不当，将造成服务不可用。当角色级别的该配置项修改后，所有实例级别的该配置项都将被修改。如果实例级别的配置项修改后，其他实例的该配置项的值保持不变。</p>	% {@auto.detect.datapart.nm.localdir}

30.7 Yarn 运维管理

30.7.1 Yarn 常用配置参数

队列资源分配

Yarn服务提供队列给用户使用，用户分配对应的系统资源给各队列使用。完成配置后，您可以单击“刷新队列”按钮或者重启Yarn服务使配置生效。

参数入口：

MRS 3.x之前的版本集群执行以下操作：

用户在MRS控制台上，选择“租户管理 > 资源分布策略”。

参数说明以default为例，其他队列的配置类似，单击“修改”编辑。

表 30-18 参数说明

配置参数	说明	默认值
资源容量	队列的资源容量（百分比）。当系统非常繁忙时，应保证每个队列的容量得到满足，而如果每个队列应用程序较少，可将剩余资源共享给其他队列。注意，所有队列的容量之和应小于100。	20
最大资源容量	队列的资源使用上限（百分比）。由于存在资源共享，因此一个队列使用的资源量可能超过其容量，而最多使用资源量可通过该参数限制。	100

MRS 3.x及后续版本集群执行以下操作：

用户可在Manager系统中，选择“租户资源 > 动态资源计划 > 队列配置”。



参数说明以修改Superior调度器的default租户为例，其他队列的配置类似，单击“修改”编辑。

表 30-19 队列配置参数

参数名	描述
AM最多占有资源（%）	表示当前队列内所有Application Master所占的最大资源百分比。

参数名	描述
每个YARN容器最多分配核数	表示当前队列内单个YARN容器可分配的最多核数，默认为-1，表示取值范围内不限制。
每个YARN容器最大分配内存（MB）	表示当前队列内单个YARN容器可分配的最大内存，默认为-1，表示取值范围内不限制。
最多运行任务数	表示当前队列最多同时可执行任务的数目，默认为-1，表示取值范围内不限制（为空意义相同），为0表示不可执行任务。取值范围为-1 ~ 2147483647。
每个用户最多运行任务数	表示每个用户在当前队列中最多同时可执行任务的数目，默认为-1，表示取值范围内不限制（为空意义相同），为0表示不可执行任务。取值范围为-1 ~ 2147483647。
最多挂起任务数	表示当前队列最多同时可挂起任务的数目，默认为-1，表示取值范围内不限制（为空意义相同），为0表示不可挂起任务。取值范围为-1 ~ 2147483647。
资源分配规则	表示单个用户任务间的资源分配规则，包括FIFO和FAIR。一个用户若在当前队列上提交了多个任务，FIFO规则代表一个任务完成后再执行其他任务，按顺序执行。FAIR规则代表各个任务同时获取到资源并平均分配资源。
默认资源标签	表示在指定资源标签（Label）的节点上执行任务。
Active状态	<ul style="list-style-type: none"> ACTIVE表示当前队列可接受并执行任务。 INACTIVE表示当前队列可接受但不执行任务，若提交任务，任务将处于挂起状态。
Open状态	<ul style="list-style-type: none"> OPEN表示当前队列处于打开状态。 CLOSED表示当前队列处于关闭状态，若提交任务，任务直接会被拒绝。

在 UI 显示 container 日志

默认情况下，系统会将container日志收集到HDFS中。如果您不需要将container日志收集到HDFS中，可以配置参数见[表30-20](#)。具体配置操作请参考[修改集群服务配置参数](#)。

表 30-20 参数说明

配置参数	说明	默认值
yarn.log-aggregation-enable	<p>设置是否将container日志收集到HDFS中。</p> <ul style="list-style-type: none"> • 设置为true，表示日志会被收集到HDFS目录中。默认目录为“{yarn.nodemanager.remote-app-log-dir}/{user}/{thisParam}”，该路径可通过界面上的“yarn.nodemanager.remote-app-log-dir-suffix”参数进行配置。 • 设置为false，表示日志不会收集到HDFS中。 <p>修改参数值后，需重启Yarn服务使其生效。</p> <p>说明 在修改值为false并生效后，生效前的日志无法在UI中获取。您可以在“yarn.nodemanager.remote-app-log-dir-suffix”参数指定的路径中获取到生效前的日志。 如果需要在UI上查看之前产生的日志，建议将此参数设置为true。</p>	true

在 WebUI 显示更多历史作业

默认情况下，Yarn WebUI界面支持任务列表分页功能，每个分页最多显示5000条历史作业，总共最多保留10000条历史作业。如果您需要在WebUI上查看更多的作业，可以配置参数如[表30-21](#)。具体配置操作请参考[修改集群服务配置参数](#)。

表 30-21 参数说明

配置参数	说明	默认值
yarn.resourcemanager.max-completed-applications	设置在WebUI总共显示的历史作业数量。	10000
yarn.resourcemanager.webapp.pagination.enable	是否开启Yarn WebUI的任务列表后台分页功能。	true
yarn.resourcemanager.webapp.pagination.threshold	开启Yarn WebUI的任务列表后台分页功能后，每个分页显示的最大作业数量。	5000

说明

- 显示更多的历史作业，会影响性能，增加打开Yarn WebUI的时间，建议开启后台分页功能，并根据实际硬件性能修改“yarn.resourcemanager.max-completed-applications”参数。
- 修改参数值后，需重启Yarn服务使其生效。

30.7.2 Yarn 日志介绍

日志描述

Yarn相关日志的默认存储路径如下：

- ResourceManager: “/var/log/Bigdata/yarn/rm”（运行日志），“/var/log/Bigdata/audit/yarn/rm”（审计日志）
- NodeManager: “/var/log/Bigdata/yarn/nm”（运行日志），“/var/log/Bigdata/audit/yarn/nm”（审计日志）

日志归档规则：Yarn的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过50MB的时候，会自动压缩，压缩后的日志文件名规则为：“<原有日志名>-<yyyy-mm-dd_hh-mm-ss>.[编号].log.zip”。最多保留最近的100个压缩文件，压缩文件保留个数可以在Manager界面中配置。

日志归档规则：

表 30-22 Yarn 日志列表

日志类型	日志文件名	描述
运行日志	hadoop-<SSH_USER>-<process_name>-<hostname>.log	Yarn组件日志，记录Yarn组件运行时候所产生的大部分日志。
	hadoop-<SSH_USER>-<process_name>-<hostname>.out	Yarn运行环境信息日志。
	<process_name>-<SSH_USER>-<DATE>-<PID>-gc.log	垃圾回收日志。
	yarn-haCheck.log	ResourceManager主备状态检测日志。
	yarn-service-check.log	Yarn服务健康状态检查日志。
	yarn-start-stop.log	Yarn服务启停操作日志。
	yarn-prestart.log	Yarn服务启动前集群操作的记录日志。
	yarn-postinstall.log	Yarn服务安装后启动前的工作日志。
	hadoop-commission.log	Yarn入服日志。
	yarn-cleanup.log	Yarn服务卸载时候的清理日志。
	yarn-refreshqueue.log	Yarn刷新队列日志。
	upgradeDetail.log	升级日志记录。
	stderr/stdin/syslog	Yarn服务上运行的应用所对应的container日志。

日志类型	日志文件名	描述
	yarn-application-check.log	Yarn服务上运行的应用检查日志。
	yarn-appsummary.log	Yarn服务上运行的应用的运行结果日志。
	yarn-switch-resourcemanager.log	Yarn主备倒换运行日志。
	ranger-yarn-plugin-enable.log	Yarn启用Ranger鉴权的日志
	yarn-nodemanager-period-check.log	Yarn nodemanager的周期检查日志
	yarn-resourcemanager-period-check.log	Yarn resourcemanager的周期检查日志
	hadoop.log	Hadoop的客户端日志
	env.log	实例启停前的环境信息日志。
审计日志	yarn-audit-<process_name>.log	Yarn操作审计日志。
	ranger-plugin-audit.log	
	SecurityAuth.audit	Yarn安全审计日志。

日志级别

Yarn中提供了如表30-23所示的日志级别。其中日志级别优先级从高到低分别是OFF、FATAL、ERROR、WARN、INFO、DEBUG。程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 30-23 日志级别

级别	描述
FATAL	FATAL表示当前事件处理存在严重错误信息。
ERROR	ERROR表示当前事件处理存在错误信息。
WARN	WARN表示当前事件处理存在异常告警信息。
INFO	INFO表示记录系统及各事件正常运行状态信息
DEBUG	DEBUG表示记录系统及系统的调试信息

如果您需要修改日志级别，请执行如下操作：

步骤1 参考[修改集群服务配置参数](#)，进入Yarn服务“全部配置”页面。

步骤2 在左边菜单栏中选择所需修改的角色所对应的日志菜单。

步骤3 选择所需修改的日志级别。

步骤4 单击“保存配置”，在弹出窗口中单击“确定”使配置生效。

说明

配置完成后立即生效，不需要重启服务。

---结束

日志格式

Yarn的日志格式如下所示：

表 30-24 日志格式

日志类型	格式	示例
运行日志	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> < 产生该日志的线程名字> <log 中的message> <日志事件的发 生位置>	2014-09-26 14:18:59,109 INFO main Client environment:java.compiler= <NA> org.apache.zookeeper.Enviro nment.logEnv(Environment. java:100)
审计日志	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> < 产生该日志的线程名字> <log 中的message> <日志事件的发 生位置>	2014-09-26 14:24:43,605 INFO main-EventThread USER=omm OPERATION=refreshAdmin Acls TARGET=AdminService RESULT=SUCCESS org.apache.hadoop.yarn.ser ver.resourcemanager.RMAu ditLogger\$LogLevel \$6.printLog(RMAuditLogger. java:91)

30.7.3 配置 Yarn 本地化日志级别

本章节适用于MRS 3.x及后续版本集群。

配置场景

container本地化默认的日志级别是INFO。用户可以通过配置“yarn.nodemanager.container-localizer.java.opts”来改变日志级别。

配置描述

在Manager系统中，选择“集群 > 待操作集群的名称 > 服务 > Yarn > 配置”，选择“全部配置”，在NodeManager的配置文件中“yarn-site.xml”中配置下面的参数来更改日志级别。

表 30-25 参数描述

参数	描述	默认值
yarn.nodemanager.container-localizer.java.opts	附加的jvm参数是提供给本地化container进程使用的。	-Xmx256m -Djava.security.krb5.conf=\${KRB5_CONFIG}

默认值-Xmx256m -Djava.security.krb5.conf=\${KRB5_CONFIG}和默认日志级别是INFO。为了更改container本地化的日志级别，添加下面的内容。

```
-Dhadoop.root.logger=<LOG_LEVEL>,localizationCLA
```

示例:

为了更改本地化日志级别为DEBUG，参数值应该为

```
-Xmx256m -Dhadoop.root.logger=DEBUG,localizationCLA
```

 **说明**

允许的日志级别是：FATAL，ERROR，WARN，INFO，DEBUG，TRACE和ALL。

30.7.4 检测 Yarn 内存使用情况

配置场景

针对所提交应用的内存使用无法预估的情况，可以通过修改服务端的配置项控制是否对内存使用进行检测。

若不检测内存使用，Container会占用内存直到内存溢出；若检测内存使用，当内存使用超过配置的内存大小时，相应的Container会被kill掉。

配置描述

参考[修改集群服务配置参数](#)进入Yarn服务参数“全部配置”界面，在搜索框中输入参数名称。

表 30-26 参数说明

参数	描述	默认值
yarn.nodemanager.vmem-check-enabled	是否进行虚拟内存检测的开关。如果任务使用的内存量超出分配值，则直接将任务强制终止。 <ul style="list-style-type: none"> • 设置为true时，进行虚拟内存检测； • 设置为false时，不进行虚拟内存检测。 	MRS 3.x之前的版本集群:false MRS 3.x及后续版本集群:true

参数	描述	默认值
yarn.nodemanager.pmem-check-enabled	是否进行物理内存检测的开关。如果任务使用的内存量超出分配值，则直接将任务强制终止。 <ul style="list-style-type: none">• 设置为true时，进行物理内存检测；• 设置为false时，不进行物理内存检测。	true

30.7.5 更改 NodeManager 的存储目录

操作场景

Yarn NodeManager定义的存储目录不正确或Yarn的存储规划变化时，MRS集群管理员需要在Manager中修改NodeManager的存储目录，以保证Yarn正常工作。NodeManager的存储目录包含本地存放目录“yarn.nodemanager.local-dirs”和日志目录“yarn.nodemanager.log-dirs”。适用于以下场景：

- 更改NodeManager角色的存储目录，所有NodeManager实例的存储目录将同步修改。
- 更改NodeManager单个实例的存储目录，只对单个实例生效，其他节点NodeManager实例存储目录不变。

对系统的影响

- 更改NodeManager角色的存储目录需要停止并重新启动集群，集群未启动前无法提供服务。
- 更改NodeManager单个实例的存储目录需要停止并重新启动实例，该节点NodeManager实例未启动前无法提供服务。
- 服务参数配置如果使用旧的存储目录，需要更新为新目录。
- 更改NodeManager的存储目录以后，需要重新下载并安装客户端。

前提条件

- 在各个数据节点准备并安装好新磁盘，并格式化磁盘。
- 规划好新的目录路径，用于保存旧目录中的数据。
- 准备好MRS集群管理员用户admin。

操作步骤

MRS 3.x之前的版本集群执行以下操作：

步骤1 检查环境。

1. 登录MRS控制台，在左侧导航栏选择“现有集群”，单击集群名称。选择“组件管理”，查看Yarn的“健康状态”是否为“良好”。
 - 是，执行[步骤1.3](#)。
 - 否，Yarn状态不健康，执行[步骤1.2](#)。
2. 请先修复Yarn异常，任务结束。

3. 确定修改NodeManager的存储目录场景。
 - 更改NodeManager角色的存储目录，执行**步骤2**。
 - 更改NodeManager单个实例的存储目录，执行**步骤3**。

步骤2 更改NodeManager角色的存储目录。

1. 选择“现有集群”，单击集群名称。选择“组件管理 > Yarn > 停止”，停止Yarn服务。
2. 登录弹性云服务器，以root用户登录到安装Yarn服务的各个节点中，执行如下操作。
 - a. 创建目标目录。
例如目标目录为“`${BIGDATA_DATA_HOME}/data2`”：
执行**`mkdir ${BIGDATA_DATA_HOME}/data2`**
 - b. 挂载目标目录到新磁盘。
例如挂载“`${BIGDATA_DATA_HOME}/data2`”到新磁盘。
 - c. 修改新目录的权限。
例如新目录路径为“`${BIGDATA_DATA_HOME}/data2`”：
执行**`chmod 750 ${BIGDATA_DATA_HOME}/data2 -R`**和**`chown omm:wheel ${BIGDATA_DATA_HOME}/data2 -R`**
3. 在MRS控制台界面，选择“现有集群”，单击集群名称。选择“组件管理 > Yarn > 实例”，选择对应主机的NodeManager实例，单击“实例配置”，“选择”“全部配置”。
将配置项“yarn.nodemanager.local-dirs”或“yarn.nodemanager.log-dirs”修改为新的目标目录。
例如：如果修改“yarn.nodemanager.local-dirs”参数，则将其值修改为“`/srv/BigData/data2/nm/localdir`”。如果修改“yarn.nodemanager.log-dirs”参数，则将其值修改为“`/srv/BigData/data2/nm/containerlogs`”。
4. 单击“保存配置”，勾选“重新启动受影响的服务或实例”，单击“确定”。重启Yarn服务。
界面提示“操作成功”，单击“完成”，Yarn成功启动，任务结束。

步骤3 更改NodeManager单个实例的存储目录。

1. 选择“现有集群”，单击集群名称。选择“组件管理 > Yarn > 实例”，勾选需要修改存储目录的NodeManager单个实例，选择“更多 > 停止实例”。
2. 登录弹性云服务器，以root用户登录到这个NodeManager节点，执行如下操作。
 - a. 创建目标目录。
例如目标目录为“`${BIGDATA_DATA_HOME}/data2`”：
执行**`mkdir ${BIGDATA_DATA_HOME}/data2`**。
 - b. 挂载目标目录到新磁盘。
例如挂载“`${BIGDATA_DATA_HOME}/data2`”到新磁盘。
 - c. 修改新目录的权限。
例如新目录路径为“`${BIGDATA_DATA_HOME}/data2`”：
执行**`chmod 750 ${BIGDATA_DATA_HOME}/data2 -R`**和**`chown omm:wheel ${BIGDATA_DATA_HOME}/data2 -R`**。
3. 在MRS控制台，单击指定的NodeManager实例并切换到“实例配置”。

将配置项“yarn.nodemanager.local-dirs”或“yarn.nodemanager.log-dirs”修改为新的目标目录。

例如：如果修改“yarn.nodemanager.local-dirs”参数，则将其值修改为“/srv/BigData/data2/nm/localdir”。如果修改“yarn.nodemanager.log-dirs”参数，则将其值修改为“/srv/BigData/data2/nm/containerlogs”。

4. 单击“保存配置”，勾选“重新启动受影响的服务或实例”。单击“确定”。重启NodeManager实例。

界面提示“操作成功”，单击“完成”，NodeManager实例启动成功。

----结束

MRS 3.x及后续版本集群也可执行以下操作：

步骤1 检查环境。

1. 登录Manager，选择“集群 > 待操作集群的名称 > 服务”查看Yarn的状态“运行状态”是否为“良好”。
 - 是，执行**1.c**。
 - 否，Yarn状态不健康，执行**1.b**。
2. 修复Yarn异常，任务结束。
3. 确定修改NodeManager的存储目录场景。
 - 更改NodeManager角色的存储目录，执行**2**。
 - 更改NodeManager单个实例的存储目录，执行**3**。

步骤2 更改NodeManager角色的存储目录。

1. 选择“集群 > 待操作集群的名称 > 服务 > Yarn > 停止服务”，停止Yarn服务。
2. 以root用户登录到安装Yarn服务的各个节点中，执行如下操作。

- a. 创建目标目录。

例如目标目录为“\${BIGDATA_DATA_HOME}/data2”：

执行**mkdir \${BIGDATA_DATA_HOME}/data2**

- b. 挂载目标目录到新磁盘。

例如挂载“\${BIGDATA_DATA_HOME}/data2”到新磁盘。

- c. 修改新目录的权限。

例如新目录路径为“\${BIGDATA_DATA_HOME}/data2”：

执行**chmod 750 \${BIGDATA_DATA_HOME}/data2 -R**和**chown omm:wheel \${BIGDATA_DATA_HOME}/data2 -R**

3. 在Manager管理界面，选择“集群 > 待操作集群的名称 > 服务 > Yarn > 实例”，选择对应主机的NodeManager实例，单击“实例配置”，选择“全部配置”。

将配置项“yarn.nodemanager.local-dirs”或“yarn.nodemanager.log-dirs”修改为新的目标目录。

例如：如果修改“yarn.nodemanager.local-dirs”参数，则将其值修改为“/srv/BigData/data2/nm/localdir”。如果修改“yarn.nodemanager.log-dirs”参数，则将其值修改为“/srv/BigData/data2/nm/containerlogs”。

4. 单击“保存”，单击“确定”。重启Yarn服务。

界面提示“操作成功”，单击“完成”，Yarn成功启动，任务结束。

步骤3 更改NodeManager单个实例的存储目录。

1. 选择“集群 > 待操作集群的名称 > 服务 > Yarn > 实例”，勾选需要修改存储目录的NodeManager单个实例，选择“更多 > 停止实例”。
2. 以root用户登录到这个NodeManager节点，执行如下操作。
 - a. 创建目标目录。
例如目标目录为“`${BIGDATA_DATA_HOME}/data2`”：
执行`mkdir ${BIGDATA_DATA_HOME}/data2`。
 - b. 挂载目标目录到新磁盘。
例如挂载“`${BIGDATA_DATA_HOME}/data2`”到新磁盘。
 - c. 修改新目录的权限。
例如新目录路径为“`${BIGDATA_DATA_HOME}/data2`”：
执行`chmod 750 ${BIGDATA_DATA_HOME}/data2 -R`和`chown omm:wheel ${BIGDATA_DATA_HOME}/data2 -R`。
3. 在Manager管理界面，单击指定的NodeManager实例并切换到“实例配置”。
将配置项“`yarn.nodemanager.local-dirs`”或“`yarn.nodemanager.log-dirs`”修改为新的目标目录。
例如：如果修改“`yarn.nodemanager.local-dirs`”参数，则将其值修改为“`/srv/BigData/data2/nm/localdir`”。如果修改“`yarn.nodemanager.log-dirs`”参数，则将其值修改为“`/srv/BigData/data2/nm/containerlogs`”。
4. 单击“保存”，单击“确定”。重启NodeManager实例。
界面提示“操作成功”，单击“完成”，NodeManager实例启动成功。

----结束

30.8 Yarn 常见问题

30.8.1 任务完成后 Container 挂载的文件目录未清除

问题

使用了CGroups功能的场景下，任务完成后Container挂载的文件目录未清除。

回答

即使任务失败，Container挂载的目录也应该被清除。

上述问题是由于删除动作超时导致的。完成某些任务所使用的时间已远超过删除时间。

为避免出现这种场景，您可以参考[修改集群服务配置参数](#)，进入Yarn“全部配置”页面。在搜索框搜索“`yarn.nodemanager.linux-container-executor.cgroups.delete-timeout-ms`”配置项来修改删除时间的时长。参数值的单位为毫秒。

30.8.2 作业执行失败时会发生 HDFS_DELEGATION_TOKEN 到期的异常

问题

安全模式下，为什么作业执行失败时会发生HDFS_DELEGATION_TOKEN到期的异常？

回答

HDFS_DELEGATION_TOKEN到期的异常是由于token没有更新或者超出了最大生命周期。

在token的最大生命周期内确保下面的参数值大于作业的运行时间。

“dfs.namenode.delegation.token.max-lifetime” = “604800000”（默认是一星期）

参考[修改集群服务配置参数](#)，进入HDFS“全部配置”页面，在搜索框搜索该参数。

说明

建议在token的最大生命周期内参数值为多倍小时数。

30.8.3 重启 YARN，本地日志不被删除

问题

在以下两种情况下重启YARN，本地日志不会被定时删除，将被永久保留。

- 在任务运行过程中，重启YARN，本地日志不被删除。
- 在任务完成，日志归集失败后定时清除日志前，重启YARN，本地日志不被删除。

回答

NodeManager有重启恢复机制，详情请参见：

https://hadoop.apache.org/docs/r3.1.1/hadoop-yarn/hadoop-yarn-site/NodeManager.html#NodeManager_Restart

可以参考[修改集群服务配置参数](#)，进入Yarn“全部配置”页面。需将NodeManager的“yarn.nodemanager.recovery.enabled”配置项为“true”后才生效，默认为“true”，这样在YARN重启的异常场景时会定时删除多余的本地日志，避免问题的出现。

30.8.4 执行任务时 AppAttempts 重试次数超过 2 次还没有运行失败

问题

系统默认的AppAttempts运行失败的次数为2。

为什么在执行任务时，AppAttempts重试次数超过2次还没有运行失败？

回答

在执行任务过程中，若ContainerExitStatus的返回值为ABORTED、PREEMPTED、DISKS_FAILED、KILLED_BY_RESOURCEMANAGER这四种状态之一时，系统不会将其计入failed attempts中。

因此出现上述的问题后，只有当真正失败尝试2次之后才会运行失败。

30.8.5 在 ResourceManager 重启后，应用程序会移回原来的队列

问题

将应用程序从一个队列移到另一个队列时，为什么在RM（ResourceManager）重启后，应用程序会被移回原来的队列？

回答

这是RM的使用限制，应用程序运行过程中移动到别的队列，此时RM重启，RM并不会在状态存储中存储新队列的信息。

假设用户提交一个MR任务到叶子队列test11上。当任务运行时，删除叶子队列test11，这时提交队列自动变为lost_and_found队列（找不到队列的任务会被放入lost_and_found队列中），任务暂停运行。要启动该任务，用户将任务移动到叶子队列test21上。在将任务移动到叶子队列test21后，任务继续运行，此时RM重启，重启后显示提交队列为lost_and_found队列，而不是test21队列。

发生上述情况的原因是，任务未完成时，RM状态存储中存储的还是应用程序移动前的队列状态。唯一的解决办法就是等RM重启后，再次移动应用程序，将新的队列状态信息写入状态存储中。

30.8.6 YARN 资源池的所有节点都被加入黑名单，任务一直处于运行状态

问题

为什么YARN资源池的所有节点都被加入黑名单，而YARN却没有释放黑名单，导致任务一直处于运行状态？

回答

在YARN中，当一个APP的节点被AM（ApplicationMaster）加入黑名单的数量达到一定比例（默认值为节点总数的33%）时，该AM会自动释放黑名单，从而不会出现由于所有可用节点都被加入黑名单而任务无法获取节点资源的现象。

在资源池场景下，假设该集群上有8个节点，通过NodeLabel特性将集群划分为两个资源池，pool A和pool B，其中pool B包含两个节点。用户提交了一个任务App1到pool B，由于HDFS空间不足，App1运行失败，导致pool B的两个节点都被App1的AM加入了黑名单，根据上述原则，2个节点小于8个节点的33%，所以YARN不会释放黑名单，使得App1一直无法得到资源而保持运行状态，后续即使被加入黑名单的节点恢复，App1也无法得到资源。

由于上述原则不适用于资源池场景，所以目前可通过调整客户端参数（路径为“客户端安装路径/Yarn/config/yarn-site.xml”）“yarn.resourcemanager.am-

scheduling.node-blacklisting-disable-threshold”为： $(\text{nodes number of pool} / \text{total nodes}) * 33\%$ 解决该问题。

30.8.7 ResourceManager 持续主备倒换

问题

RM（ResourceManager）在多个任务（比如2000个任务）正常并发运行时出现持续的主备倒换，导致YARN服务不可用。

回答

产生上述问题的原因是，full GC（GarbageCollection）时间过长，超出了RM与ZK（ZooKeeper）之间定期交互时长的阈值，导致RM与ZK失联，从而造成RM主备倒换。

在多任务情况下，RM需要保存多个任务的鉴权信息，并通过心跳传递给各个NM（NodeManager），即心跳Response。心跳Response的生命周期短，默认值为1s，一般可以在JVM minor GC时被回收，但在多任务的情况下，集群规模较大，比如5000节点，多个节点的心跳Response会占用大量内存，导致JVM在minor GC时无法完全回收，无法回收的内存持续累积，最终触发JVM的full GC。JVM的GC都是阻塞式的，即在GC过程中不执行任何作业，所以若full GC的时间过长，超出了RM与ZK之间定期交互时长的阈值，就会出现主备倒换。

登录FusionInsight Manager，选择“集群 > 服务 > Yarn > 配置 > 全部配置”，在左侧选择“Yarn > 自定义”，在“yarn.yarn-site.customized.configs”中添加“yarn.resourcemanager.zk-timeout-ms”参数来增大RM与ZK之间定期交互时长的阈值（参数值的范围为小于等于90000毫秒），可以解决RM持续主备倒换的问题。

30.8.8 当一个 NodeManager 处于 unhealthy 的状态 10 分钟时，新应用程序失败

问题

当一个NM（NodeManager）处于unhealthy的状态10分钟时，新应用程序失败。

回答

当nodeSelectPolicy为SEQUENCE，且第一个连接到RM的NM不可用时，RM会在“yarn.nm.liveness-monitor.expiry-interval-ms”属性中指定的周期内，一直尝试为同一个NM分配任务。

可以通过两种方式来避免上述问题：

- 使用其他的nodeSelectPolicy，如RANDOM。
- 参考[修改集群服务配置参数](#)，进入Yarn“全部配置”页面。在搜索框搜索以下参数，通过“yarn-site.xml”文件更改以下属性：
“yarn.resourcemanager.am-scheduling.node-blacklisting-enabled” = “true”；
“yarn.resourcemanager.am-scheduling.node-blacklisting-disable-threshold” = “0.5”。

30.8.9 Superior 通过 REST 接口查看已结束或不存在的 applicationID，页面提示 Error Occurred

问题

Superior通过REST接口查看已结束或不存在的applicationID，返回的页面提示Error Occurred。

回答

用户提交查看applicationID的请求，访问REST接口“https://<SS_REST_SERVER>/ws/v1/sscheduler/applications/{application_id}”。

由于Superior Scheduler只存储正在运行的applicationID，所以当查看的是已结束或不存在的applicationID，服务器会响应给浏览器“404”的状态码。但是由于chrome浏览器访问该REST接口时，优先以“application/xml”的格式响应，该行为会导致服务器端处理出现异常，所以返回的页面会提示“Error Occurred”。而IE浏览器访问该REST接口时，优先以“application/json”的格式响应，服务器会正确响应给浏览器“404”的状态码。

30.8.10 Superior 调度模式下，单个 NodeManager 故障可能导致 MapReduce 任务失败

问题

在Superior调度模式下，如果出现单个NodeManager故障，可能会导致Mapreduce任务失败。

回答

正常情况下，当一个application的单个task的attempt连续在一个节点上失败3次，那么该application的AppMaster就会将该节点加入黑名单，之后AppMaster就会通知调度器不要继续调度task到该节点，从而避免任务失败。

但是默认情况下，当集群中有33%的节点都被加入黑名单时，调度器会忽略黑名单节点。因此，该黑名单特性在小集群场景下容易失效。比如，集群只有3个节点，当1个节点出现故障，黑名单机制失效，不管task的attempt在同一个节点失败多少次，调度器仍然会将task继续调度到该节点，从而导致application因为task失败达到最大attempt次数（MapReduce默认4次）而失败。

规避手段：

在“客户端安装路径/Yarn/config/yarn-site.xml”文件中修改“yarn.resourceamanager.am-scheduling.node-blacklisting-disable-threshold”参数以百分比的形式配置忽略黑名单节点的阈值。建议根据集群规模，适当增大该参数的值，如3个节点的集群，建议增大到50%。

说明

Superior调度器的框架设计是基于时间的异步调度，当NodeManager故障后，ResourceManager无法快速的感知到NodeManager已经出了问题(默认10mins)，因此在此期间，Superior调度器仍然会向该节点调度task，从而导致任务失败。

30.8.11 当应用程序从 `lost_and_found` 队列移动到其它队列时，应用程序不能继续执行

问题

当删除一个有部分应用程序正在运行的队列，这些应用程序会被移动到“`lost_and_found`”队列上。当这些应用程序移回运行正常的队列时，某些任务会被挂起，不能正常运行。

回答

如果应用程序没有设置标签表达式，那么该应用程序上新增的container/resource将使用其所在队列默认的标签表达式。如果队列没有默认的标签表达式，则将其标签表达式设置为“`default label`”。

当应用程序（`app1`）提交到队列（`Q1`）上时，应用程序上新增的container/resource使用队列默认的标签表达式（“`label1`”）。若`app1`正在运行时`Q1`被删除，则`app1`被移动到“`lost_and_found`”队列上。由于“`lost_and_found`”队列没有标签表达式，其标签表达式设置为“`default label`”，此时`app1`上新增的container/resource也将其标签表达式设置为“`default label`”。当`app1`被移回正常运行的队列（例如，`Q2`）时，如果`Q2`支持调用`app1`中的所有标签表达式（包含“`label1`”和“`default label`”），则`app1`能正常运行直到结束；如果`Q2`仅支持调用`app1`中的部分标签表达式（例如，仅支持调用“`default label`”），那么`app1`在运行时，拥有“`label1`”标签表达式的部分任务的资源请求将无法获得资源，从而被挂起，不能正常运行。

因此当把应用程序从“`lost_and_found`”队列移动到其它运行正常的队列上时，需要保证目标队列能够调用该应用程序的所有标签表达式。

建议不要删除正在运行应用程序的队列。

30.8.12 如何限制存储在 ZKstore 中的应用程序诊断消息的大小

问题

如何限制存储在ZKstore中的应用程序诊断消息的大小？

回答

在某些情况下，已经观察到诊断消息可能无限增长。由于诊断消息存储在状态存储中，不建议允许诊断消息无限增长。因此，需要有一个属性参数用于设置诊断消息的最大大小。

若您需要设置“`yarn.app.attempt.diagnostics.limit.kc`”参数值，具体操作参考[修改集群服务配置参数](#)，进入Yarn“全部配置”页面，在搜索框搜索以下参数。

表 30-27 参数描述

参数	描述	默认值
yarn.app.attempt.diagnostics.limit.kc	定义每次应用连接的诊断消息的数据大小，以千字节为单位（字符数*1024）。当使用ZooKeeper来存储应用程序的行为状态时，需要限制诊断消息的大小，以防止YARN拖垮ZooKeeper。如果将“yarn.resourcemanager.state-store.max-completed-applications”设置为一个较大的数值，则需要减小该属性参数的值以限制存储的总数据大小。	64

30.8.13 为什么将非 ViewFS 文件系统配置为 ViewFS 时 MapReduce 作业运行失败

问题

为什么将非ViewFS文件系统配置为ViewFS时MR作业运行失败？

回答

通过集群将非ViewFS文件系统配置为ViewFS时，ViewFS中的文件夹的用户权限与默认NameService中的非ViewFS不同。因为目录权限不匹配，所以已提交的MR作业运行失败。

在集群中配置ViewFS的用户，需要检查并校验目录权限。在提交作业之前，应按照默认NameService文件夹权限更改ViewFS文件夹权限。

下表列出了ViewFS中配置的目录的默认权限结构。如果配置的目录权限与下表不匹配，则必须相应地更改目录权限。

表 30-28 ViewFS 中配置的目录的默认权限结构

参数	描述	默认值	默认值及其父目录的默认权限
yarn.nodemanager.remote-app-log-dir	在默认文件系统上（通常是HDFS），指定NM应将日志聚合到哪个目录。	logs	777
yarn.nodemanager.remote-app-log-archive-dir	将日志归档的目录。	-	777
yarn.app.mapreduce.am.staging-dir	提交作业时使用的 staging 目录。	/tmp/hadoop-yarn/staging	777

参数	描述	默认值	默认值及其父目录的默认权限
mapreduce.jobhistory.intermediate-done-dir	MapReduce作业记录历史文件的目录。	<code>\${yarn.app.mapreduce.am.staging-dir}/history/done_intermediate</code>	777
mapreduce.jobhistory.done-dir	由MR JobHistory Server管理的历史文件的目录。	<code>\${yarn.app.mapreduce.am.staging-dir}/history/done</code>	777

30.8.14 开启 Native Task 特性后，Reduce 任务在部分操作系统运行失败

问题

开启Native Task特性后，Reduce任务在部分操作系统运行失败。

回答

运行包含Reduce的Mapreduce任务时，通过-Dmapreduce.job.map.output.collector.class=org.apache.hadoop.mapred.nativetask.NativeMapOutputCollectorDelegator命令开启Native Task特性，任务在部分操作系统运行失败，日志中提示错误“version 'GLIBCXX_3.4.20' not found”。该问题原因是操作系统的GLIBCXX版本较低，导致该特性依赖的libnativetask.so.1.0.0库无法加载，进而导致任务失败。

规避手段：

设置配置项mapreduce.job.map.output.collector.class的值为org.apache.hadoop.mapred.MapTask\$MapOutputBuffer。

31 使用 ZooKeeper

31.1 使用 ZooKeeper 客户端

Zookeeper是一个开源的，高可靠的，分布式一致性协调服务。Zookeeper设计目标是用来解决那些复杂，易出错的分布式系统难以保证数据一致性的。不必开发专门的协同应用，十分适合高可用服务保持数据一致性。

背景信息

在使用客户端前，除主管理节点以外的客户端，需要下载并更新客户端配置文件。

操作步骤

MRS 2.x及以前版本集群执行以下操作：

步骤1 下载客户端配置文件。

1. 登录MRS控制台，在左侧导航栏选择“现有集群”，单击待操作集群的名称。
2. 选择“组件管理”。

说明

针对MRS 1.8.10及之前版本，登录MRS Manager页面，具体请参见[访问MRS Manager](#)，然后选择“服务管理”。

3. 单击“下载客户端”。
“客户端类型”选择“仅配置文件”，“下载路径”选择“服务器端”，单击“确定”开始生成客户端配置文件，文件生成后默认保存在主管理节点“/tmp/MRS-client”。

图 31-1 仅下载客户端的配置文件



步骤2 登录MRS Manager的主管理节点。

1. 在MRS控制台，选择“现有集群”，单击集群名称，在“节点管理”页签中查看节点名称，名称中包含“master1”的节点为Master1节点，名称中包含“master2”的节点为Master2节点。

MRS Manager的主备管理节点默认安装在集群Master节点上。在主备模式下，由于Master1和Master2之间会切换，Master1节点不一定是MRS Manager的主管理节点，需要在Master1节点中执行命令，确认MRS Manager的主管理节点。命令请参考[步骤2.4](#)。

2. 以root用户使用密码方式登录Master1节点。操作方法，请参见[登录集群节点](#)章节。
3. 切换至omm用户。

```
sudo su - root
su - omm
```

4. 执行以下命令确认MRS Manager的主管理节点。

```
sh ${BIGDATA_HOME}/om-0.0.1/sbin/status-oms.sh
```

回显信息中“HAActive”参数值为“active”的节点为主管理节点（如下例中“mgtomsdat-sh-3-01-1”为主管理节点），参数值为“standby”的节点为备管理节点（如下例中“mgtomsdat-sh-3-01-2”为备管理节点）。

NodeName	HostName	HAVersion	StartTime	HAActive
192-168-0-30	mgtomsdat-sh-3-01-1	V100R001C01	2014-11-18 23:43:02	active
192-168-0-24	mgtomsdat-sh-3-01-2	V100R001C01	2014-11-21 07:14:02	standby

5. 使用root用户登录MRS Manager的主管理节点，例如“192-168-0-30”节点，并执行以下命令切换到omm用户。

```
sudo su - omm
```

步骤3 执行以下命令切换到客户端安装目录。例如“/opt/client”。

```
cd /opt/client
```

步骤4 执行以下命令，更新主管理节点的客户端配置。

```
sh refreshConfig.sh /opt/client 客户端配置文件压缩包完整路径
```

例如，执行命令：

```
sh refreshConfig.sh /opt/client/tmp/MRS-client/MRS_Services_Client.tar
```

界面显示以下信息表示配置刷新更新成功：

```
ReFresh components client config is complete.  
Succeed to refresh components client config.
```

说明

对于MRS 1.8.5及之后版本集群，步骤**步骤1~步骤4**的操作也可以参考[更新客户端（3.x之前版本）](#)页面的方法二操作。

步骤5 在Master节点使用客户端。

1. 在已更新客户端的主管理节点，例如“192-168-0-30”节点，执行以下命令切换到客户端目录。

```
cd /opt/client
```

2. 执行以下命令配置环境变量。

```
source bigdata_env
```

3. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，具体请参见[准备开发用户](#)创建对应用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit MRS集群用户
```

例如，kinit zookeeperuser。

4. 直接执行Zookeeper组件的客户端命令。

```
zkCli.sh -server <zookeeper安装节点ip>:<port>
```

例如：**zkCli.sh -server node-master1DGhZ:2181**

步骤6 运行Zookeeper客户端命令。

1. 创建ZNode。
create /test
2. 查看ZNode信息。
ls /
3. 向ZNode中写入数据。
set /test "zookeeper test"
4. 查看写入ZNode中的数据。
get /test
5. 删除创建的ZNode。
delete /test

----**结束**

MRS 3.x及以后版本集群执行以下操作：

步骤1 下载客户端配置文件。

1. 登录FusionInsight Manager页面，具体请参见[访问FusionInsight Manager（MRS 3.x及之后版本）](#)。
2. 选择“集群 > 待操作集群的名称 > 概览 > 更多 > 下载客户端”。
3. 下载集群客户端。
“选择客户端类型”选择“仅配置文件”，选择平台类型，单击“确定”开始生成客户端配置文件，文件生成后默认保存在主管理节点“/tmp/FusionInsight-Client/”。

步骤2 登录Manager的主管理节点。

1. 以root用户登录任意部署Manager的节点。
2. 执行以下命令确认主备管理节点。

```
sh ${BIGDATA_HOME}/om-server/om/sbin/status-oms.sh
```

界面打印信息中“HAActive”参数值为“active”的节点为主管理节点（如下例中“node-master1”为主管理节点），参数值为“standby”的节点为备管理节点（如下例中“node-master2”为备管理节点）。

```
HAMode
double
NodeName      HostName      HAVersion      StartTime      HAActive
HAAllResOK    HARunPhase
192-168-0-30  node-master1  V100R001C01    2020-05-01 23:43:02  active
normal        Activated
192-168-0-24  node-master2  V100R001C01    2020-05-01 07:14:02  standby
normal        Deactivated
```

3. 以root用户登录主管理节点，并执行以下命令切换到omm用户。

```
sudo su - omm
```

步骤3 执行以下命令切换到客户端安装目录。例如“/opt/client”。

```
cd /opt/client
```

步骤4 执行以下命令，更新主管理节点的客户端配置。

```
sh refreshConfig.sh /opt/client 客户端配置文件压缩包完整路径
```

例如，执行命令：

```
sh refreshConfig.sh /opt/client /tmp/FusionInsight-Client/
FusionInsight_Cluster_1_Services_Client.tar
```

界面显示以下信息表示配置刷新更新成功：

```
ReFresh components client config is complete.
Succeed to refresh components client config.
```

步骤5 在Master节点使用客户端。

1. 在已更新客户端的主管理节点，例如“192-168-0-30”节点，执行以下命令切换到客户端目录。

```
cd /opt/client
```

2. 执行以下命令配置环境变量。

```
source bigdata_env
```

3. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户，具体请参见[角色管理](#)配置拥有对应权限的角色，参考[创建用户](#)为用户绑定对应角色。如果当前集群未启用Kerberos认证，则无需执行此命令。

kinit MRS 集群用户

例如，`kinit zookeeperuser`。

4. 直接执行Zookeeper组件的客户端命令。

`zkCli.sh -server <zookeeper安装节点ip>:<port>`

例如：`zkCli.sh -server node-master1DGhZ:2181`

步骤6 运行Zookeeper客户端命令。

1. 创建ZNode。
`create /test`
2. 查看ZNode信息。
`ls /`
3. 向ZNode中写入数据。
`set /test "zookeeper test"`
4. 查看写入ZNode中的数据。
`get /test`
5. 删除创建的ZNode。
`delete /test`

----结束

31.2 配置 ZooKeeper ZNode ACL

操作场景

该操作指导用户对ZooKeeper的znode设置权限。

ZooKeeper通过访问控制列表（ACL）来对znode进行访问控制。ZooKeeper客户端为znode指定ACL，ZooKeeper服务器根据ACL列表判定某个请求znode的客户端是否有对应操作的权限。ACL设置涉及如下四个方面。

- 查看ZooKeeper中znode的ACL。
- 增加ZooKeeper中znode的ACL。
- 修改ZooKeeper中znode的ACL。
- 删除ZooKeeper中znode的ACL。

ZooKeeper的ACL权限说明：

ZooKeeper目前支持create, delete, read, write, admin五种权限，且ZooKeeper对权限的控制是znode级别的，而且不继承，即对父znode设置权限，其子znode不继承父znode的权限。ZooKeeper中znode的默认权限为**world:anyone:cdrwa**，即任何用户都有所有权限。

说明

ACL有三部分：

第一部分是认证类型，如world指所有认证类型，sasl是kerberos认证类型；

第二部分是账号，如anyone指的是任何人；

第三部分是权限，如cdrwa指的是拥有所有权限。

特别的，由于普通模式启动客户端不需要认证，sasl认证类型的ACL在普通模式下将不能使用。本文所有涉及sasl方式的鉴权操作均是在安全集群中进行。

表 31-1 Zookeeper 的五种 ACL

权限说明	权限简称	权限详情
创建权限	create(c)	可以在当前znode下创建子znode
删除权限	delete(d)	删除当前的znode
读权限	read(r)	获取当前znode的数据，可以列出当前znode所有的子znodes
写权限	write(w)	向当前znode写数据，写入子znode
管理权限	admin(a)	设置当前znode的权限

对系统的影响

须知

修改ZooKeeper的ACL是高危操作。修改ZooKeeper中znode的权限，可能会导致其他用户无权限访问该znode，导致系统功能异常。另外在3.5.6及以后版本，用户对于getAcl操作需要有读权限。

前提条件

- 已安装ZooKeeper客户端。例如安装目录为“/opt/client”。
- 已获取MRS集群管理员用户和密码。

操作步骤

启动ZooKeeper客户端

步骤1 以root用户登录安装了ZooKeeper客户端的服务器。

步骤2 进入客户端安装目录。

```
cd /opt/client
```

步骤3 执行以下命令配置环境变量。

```
source bigdata_env
```

步骤4 执行以下命令认证用户身份，并输入用户密码（任意有权限的用户，这里以userA为例，普通模式不涉及）。

kinit userA

步骤5 在ZooKeeper客户端执行以下命令，进入ZooKeeper命令行。

sh zkCli.sh -server ZooKeeper任意实例所在节点业务平面IP.clientPort

默认的“clientPort”为“2181”

例如：**sh zkCli.sh -server 192.168.0.151:2181**

步骤6 登录ZooKeeper客户端后，使用ls命令，可以查看ZooKeeper中的znode列表。例如，可以查看根目录znode列表。

ls /

```
[zk: 192.168.0.151:2181(CONNECTED) 1] ls /  
[hadoop-flag, hadoop-ha, test, test2, test3, test4, test5, test6, zookeeper]
```

查看ZooKeeper znode ACL信息

步骤7 启动ZooKeeper客户端。

步骤8 使用getAcl命令，可以查看znode。如下命令，可以查看到之前创建的名为test的znode的ACL权限。

getAcl /znode名称

```
[zk: 192.168.0.151:2181(CONNECTED) 2] getAcl /test  
'world,'anyone  
: cdrwa
```

增加ZooKeeper znode ACL信息

步骤9 启动ZooKeeper客户端。

步骤10 查看旧的ACL信息，查看当前账号是否有权限修改该znode的ACL信息的权限（a权限），如果没有权限，需要kinit登录有权限的用户，并重新启动ZooKeeper客户端。

getAcl /znode名称

```
[zk: 192.168.0.151:2181(CONNECTED) 3] getAcl /test  
'world,'anyone  
: cdrwa
```

步骤11 使用setAcl命令增加权限。设置新权限命令如下：

setAcl /test world:anyone:cdrwa,sasl:用户名@<系统域名>:权限值

例如对test的znode，需要增加userA用户的权限：

setAcl /test world:anyone:cdrwa,sasl:userA@HADOOP.COM:cdrwa

📖 说明

增加权限时，需要保留已有权限。新增加权限和旧的权限用英文逗号隔开，新增加权限有三个部分：

- 第一部分是认证类型，如sasl指使用kerberos认证；
- 第二部分是账号，如userA@HADOOP.COM指的是userA用户；
- 第三部分是权限，如cdrwa指的是拥有所有权限。

步骤12 setAcl后，可以使用getAcl命令查看增加权限是否成功：

getAcl /znode名称

```
[zk: 192.168.0.151:2181(CONNECTED) 4] getAcl /test
'world,'anyone
: cdrwa
'sasl,'userA@<系统域名>
: cdrwa
```

修改ZooKeeper znode ACL信息

步骤13 启动ZooKeeper客户端。

步骤14 查看旧的ACL信息，查看当前账号是否有权限修改该znode的ACL信息的权限（a权限），如果没有权限，需要kinit登录有权限的用户，并重新启动ZooKeeper客户端。

getAcl /znode名称

```
[zk: 192.168.0.151:2181(CONNECTED) 5] getAcl /test
'world,'anyone
: cdrwa
'sasl,'userA@<系统域名>
: cdrwa
```

步骤15 使用setAcl命令修改权限。设置新权限命令如下：

setAcl /test sasl:用户名@<系统域名>:权限值

例如仅保留userA用户的所有权限，删除anyone用户的rw权限。

setAcl /test sasl:userA@HADOOP.COM:cdrwa

步骤16 setAcl后，可以使用getAcl命令查看修改权限是否成功：

getAcl /znode名称

```
[zk: 192.168.0.151:2181(CONNECTED) 6] getAcl /test
'sasl,'userA@<系统域名>
: cdrwa
```

删除ZooKeeper znode ACL信息

步骤17 启动ZooKeeper客户端。

步骤18 查看旧的ACL信息，查看当前账号是否有权限修改该znode的ACL信息的权限（a权限），如果没有权限，需要kinit登录有权限的用户，并重新启动ZooKeeper客户端。

getAcl /znode名称

```
[zk: 192.168.0.151:2181(CONNECTED) 5] getAcl /test
'world,'anyone
: rw
'sasl,'userA@<系统域名>
: cdrwa
```

步骤19 使用setAcl命令增加权限。设置新权限命令如下：

setAcl /test sasl:用户名@<系统域名>:权限值

例如，仅保留userA用户是所有权限，取消anyone用户的rw权限。

setAcl /test sasl:userA@HADOOP.COM:cdrwa

步骤20 setAcl后，可以使用getAcl命令查看修改权限是否成功

getAcl /znode名称


```
[zk: 192.168.0.151:2181(CONNECTED) 6] getAcl /test
'sasl,'userA@<系统域名>
: cdrwa
```

----结束

31.3 ZooKeeper 常用配置参数

参数入口：

请参考[修改集群服务配置参数](#)，进入ZooKeeper“全部配置”页面。在搜索框中输入参数名称。

表 31-2 参数说明

配置参数	说明	默认值
skipACL	是否跳过ZooKeeper节点的权限检查。	no
maxClientCnxns	ZooKeeper的最大连接数，在连接数多的情况下，建议增加。	2000
LOG_LEVEL	日志级别，在调试的时候，可以改为DEBUG。	INFO
acl.compare.shortName	当Znode的ACL权限认证类型为SASL时，是否仅使用principal的用户名部分进行ACL权限认证。	true
synclimit	Follower与leader进行同步的时间间隔（单位为tick）。如果在指定的时间内leader没响应，连接将不能被建立。	15
tickTime	一次tick的时间（毫秒），它是ZooKeeper使用的基本时间单位，心跳、超时的时间都由它来规定。	4000

说明

ZooKeeper内部时间由参数ticktime和参数synclimit控制，如需调大ZooKeeper内部超时时间，需要调大客户端连接ZooKeeper的超时时间。

31.4 ZooKeeper 日志介绍

日志描述

日志存储路径：“/var/log/Bigdata/zookeeper/quorumpeer”（运行日志），
“/var/log/Bigdata/audit/zookeeper/quorumpeer”（审计日志）

日志归档规则：ZooKeeper的日志启动了自动压缩归档功能，缺省情况下，当日志大小超过30MB的时候，会自动压缩。最多保留20个压缩文件，压缩文件保留个数可以在Manager界面中配置。

表 31-3 ZooKeeper 日志列表

日志类型	日志文件名	描述
运行日志	zookeeper-<SSH_USER>-<process_name>-<hostname>.log	ZooKeeper系统日志，记录 ZooKeeper系统运行时候所产生的大部分日志。
	check-serviceDetail.log	ZooKeeper服务启动是否成功的检查日志。
	zookeeper-<SSH_USER>-<DATA>-<PID>-gc.log	ZooKeeper垃圾回收日志。
	instanceHealthDetail.log	ZooKeeper实例健康状态检查日志
	zookeeper-omm-server-<hostname>.out	ZooKeeper运行异常退出日志。
	zk-err-<zkpid>.log	ZooKeeper致命错误日志。
	java_pid<zkpid>.hprof	ZooKeeper内存溢出日志。
	funcDetail.log	ZooKeeper实例启动日志。
	zookeeper-period-check.log	ZooKeeper实例健康检查日志。
	zookeeper-period-check-java.log	ZooKeeper配额监控周期检查日志。
审计日志	zk-audit-quorumpeer.log	ZooKeeper操作审计日志。

日志级别

ZooKeeper中提供了如表31-4所示的日志级别。日志级别优先级从高到低分别是 FATAL、ERROR、WARN、INFO、DEBUG。程序会打印高于或等于所设置级别的日志，设置的日志等级越高，打印出来的日志就越少。

表 31-4 日志级别

级别	描述
FATAL	FATAL表示当前事件处理出现严重错误信息，可能导致系统崩溃。
ERROR	ERROR表示当前事件处理出现错误信息，系统运行出错。
WARN	WARN表示当前事件处理存在异常信息，但认为是正常范围，不会导致系统出错。
INFO	INFO表示系统及各事件正常运行状态信息。
DEBUG	DEBUG表示系统及系统的调试信息。

如果您需要修改日志级别，请执行如下操作：

- 步骤1** 参考[修改集群服务配置参数](#)章节，进入ZooKeeper服务“全部配置”页面。
- 步骤2** 左边菜单栏中选择所需修改的角色所对应的日志菜单。
- 步骤3** 选择所需修改的日志级别。
- 步骤4** 单击“保存”，在弹出窗口中单击“确定”使配置生效。

 **说明**

配置完成后立即生效，不需要重启服务。

----**结束**

日志格式

ZooKeeper的日志格式如下所示：

表 31-5 日志格式

日志类型	组件	格式	示例
运行日志	zookeeper quorumpeer	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生 该日志的线程名字 > <log中的 message> <日志事 件的发生位置>	2020-01-20 16:33:43,816 INFO main Defaulting to majority quorums org.apache.zookee per.server.quorum. QuorumPeerConfi g.parseProperties(QuorumPeerConfi g.java:335)

日志类型	组件	格式	示例
审计日志	zookeeper quorumpeer	<yyyy-MM-dd HH:mm:ss,SSS> <Log Level> <产生 该日志的线程名字 > <log中的 message> <日志事 件的发生位置>	2020-01-20 16:33:54,313 INFO CommitProcessor: 13 session=0xd4b067 9daea0000 ip=10.177.112.145 operation=create znode target=ZooKeeper Server znode=/zk- write-test-2 result=success org.apache.zookee per.ZKAuditLogger \$LogLevel \$5.printLog(ZKAu ditLogger.java:70)

31.5 ZooKeeper 常见问题

31.5.1 创建大量 ZNode 后 ZooKeeper Server 启动失败

问题

创建大量ZNode后，ZooKeeper集群处于故障状态不能自动恢复，尝试重启失败，ZooKeeper Server日志显示如下内容：

Follower:

```
2016-06-23 08:00:18,763 | WARN | QuorumPeer[myid=26](plain=/10.16.9.138:2181)(secure=disabled) |
Exception when following the leader |
org.apache.zookeeper.server.quorum.Follower.followLeader(Follower.java:93)
java.net.SocketTimeoutException: Read timed out
    at java.net.SocketInputStream.socketRead0(Native Method)
    at java.net.SocketInputStream.socketRead(SocketInputStream.java:116)
    at java.net.SocketInputStream.read(SocketInputStream.java:170)
    at java.net.SocketInputStream.read(SocketInputStream.java:141)
    at java.io.BufferedInputStream.fill(BufferedInputStream.java:246)
    at java.io.BufferedInputStream.read(BufferedInputStream.java:265)
    at java.io.DataInputStream.readInt(DataInputStream.java:387)
    at org.apache.jute.BinaryInputArchive.readInt(BinaryInputArchive.java:63)
    at org.apache.zookeeper.server.quorum.QuorumPacket.deserialize(QuorumPacket.java:83)
    at org.apache.jute.BinaryInputArchive.readRecord(BinaryInputArchive.java:99)
    at org.apache.zookeeper.server.quorum.Learner.readPacket(Learner.java:156)
    at org.apache.zookeeper.server.quorum.Learner.registerWithLeader(Learner.java:276)
    at org.apache.zookeeper.server.quorum.Follower.followLeader(Follower.java:75)
    at org.apache.zookeeper.server.quorum.QuorumPeer.run(QuorumPeer.java:1094)
2016-06-23 08:00:18,764 | INFO | QuorumPeer[myid=26](plain=/10.16.9.138:2181)(secure=disabled) |
shutdown called | org.apache.zookeeper.server.quorum.Follower.shutdown(Follower.java:198)
java.lang.Exception: shutdown Follower
    at org.apache.zookeeper.server.quorum.Follower.shutdown(Follower.java:198)
```

```
at org.apache.zookeeper.server.quorum.QuorumPeer.stopFollower(QuorumPeer.java:1141)
at org.apache.zookeeper.server.quorum.QuorumPeer.run(QuorumPeer.java:1098)
```

Leader:

```
2016-06-23 07:30:57,481 | WARN | QuorumPeer[myid=25](plain=/10.16.9.136:2181)(secure=disabled) |
Unexpected exception | org.apache.zookeeper.server.quorum.QuorumPeer.run(QuorumPeer.java:1108)
java.lang.InterruptedExcepion: Timeout while waiting for epoch to be acked by quorum
at org.apache.zookeeper.server.quorum.Leader.waitForEpochAck(Leader.java:1221)
at org.apache.zookeeper.server.quorum.Leader.lead(Leader.java:487)
at org.apache.zookeeper.server.quorum.QuorumPeer.run(QuorumPeer.java:1105)
2016-06-23 07:30:57,482 | INFO | QuorumPeer[myid=25](plain=/10.16.9.136:2181)(secure=disabled) |
Shutdown called | org.apache.zookeeper.server.quorum.Leader.shutdown(Leader.java:623)
java.lang.Exception: shutdown Leader! reason: Forcing shutdown
at org.apache.zookeeper.server.quorum.Leader.shutdown(Leader.java:623)
at org.apache.zookeeper.server.quorum.QuorumPeer.stopLeader(QuorumPeer.java:1149)
at org.apache.zookeeper.server.quorum.QuorumPeer.run(QuorumPeer.java:1110)
```

回答

创建大量节点后，follower与leader同步时数据量大，在集群数据同步限定时间内不能完成同步过程，导致超时，各个ZooKeeper Server启动失败。

参考[修改集群服务配置参数](#)章节，进入ZooKeeper服务“全部配置”页面。不断尝试调大ZooKeeper配置文件“zoo.cfg”中的“syncLimit”和“initLimit”两参数值，直到ZooKeeperServer正常。

表 31-6 参数说明

参数	描述	默认值
syncLimit	follower与leader进行同步的时间间隔（时长为ticket时长的倍数）。如果在该时间范围内leader没响应，连接将不能被建立。	15
initLimit	follower连接到leader并与leader同步的时间（时长为ticket时长的倍数）。	15

如果将参数“initLimit”和“syncLimit”的参数值均配置为“300”之后，ZooKeeper Server仍然无法恢复，则需确认没有其他应用程序正在kill ZooKeeper。例如，参数值为“300”，ticket时长为2000毫秒，即同步限定时间为300*2000ms=600s。

可能存在以下场景，在ZooKeeper中创建的数据过大，需要大量时间与leader同步，并保存到硬盘。在这个过程中，如果ZooKeeper需要运行很长时间，则需确保没有其他监控应用程序kill ZooKeeper而判断其服务停止。

31.5.2 为什么 ZooKeeper Server 出现 java.io.IOException: Len 的错误日志

问题

在父目录中创建大量的znode之后，当ZooKeeper客户端尝试在单个请求中获取该父目录中的所有子节点时，将返回失败。

客户端日志，如下所示：

```
2017-07-11 13:17:19,610 [myid:] - WARN [New I/O worker #3:ClientCnxnSocketNetty
$ZKClientHandler@468] - Exception caught: [id: 0xb66cbb85, /10.18.97.97:49192 ->
```

```
10.18.97.97/10.18.97.97:2181] EXCEPTION: java.nio.channels.ClosedChannelException
java.nio.channels.ClosedChannelException
at org.jboss.netty.handler.ssl.SslHandler$6.run(SslHandler.java:1580)
at org.jboss.netty.channel.socket.ChannelRunnableWrapper.run(ChannelRunnableWrapper.java:40)
at org.jboss.netty.channel.socket.nio.AbstractNioWorker.executeInIoThread(AbstractNioWorker.java:71)
at org.jboss.netty.channel.socket.nio.NioWorker.executeInIoThread(NioWorker.java:36)
at org.jboss.netty.channel.socket.nio.AbstractNioWorker.executeInIoThread(AbstractNioWorker.java:57)
at org.jboss.netty.channel.socket.nio.NioWorker.executeInIoThread(NioWorker.java:36)
at org.jboss.netty.channel.socket.nio.AbstractNioChannelSink.execute(AbstractNioChannelSink.java:34)
at org.jboss.netty.handler.ssl.SslHandler.channelClosed(SslHandler.java:1566)
at org.jboss.netty.channel.Channels.fireChannelClosed(Channels.java:468)
at org.jboss.netty.channel.socket.nio.AbstractNioWorker.close(AbstractNioWorker.java:376)
at org.jboss.netty.channel.socket.nio.NioWorker.read(NioWorker.java:93)
at org.jboss.netty.channel.socket.nio.AbstractNioWorker.process(AbstractNioWorker.java:109)
at org.jboss.netty.channel.socket.nio.AbstractNioSelector.run(AbstractNioSelector.java:312)
at org.jboss.netty.channel.socket.nio.AbstractNioWorker.run(AbstractNioWorker.java:90)
at org.jboss.netty.channel.socket.nio.NioWorker.run(NioWorker.java:178)
at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1142)
at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:617)
at java.lang.Thread.run(Thread.java:745)
```

Leader节点的日志，如下所示：

```
2017-07-11 13:17:33,043 [myid:1] - WARN [New I/O worker #7:NettyServerCnxn@445] - Closing
connection to /10.18.101.110:39856
java.io.IOException: Len error 45
at org.apache.zookeeper.server.NettyServerCnxn.receiveMessage(NettyServerCnxn.java:438)
at org.apache.zookeeper.server.NettyServerCnxnFactory
$CnxnChannelHandler.processMessage(NettyServerCnxnFactory.java:267)
at org.apache.zookeeper.server.NettyServerCnxnFactory
$CnxnChannelHandler.messageReceived(NettyServerCnxnFactory.java:187)
at org.jboss.netty.channel.SimpleChannelHandler.handleUpstream(SimpleChannelHandler.java:88)
at org.jboss.netty.channel.DefaultChannelPipeline.sendUpstream(DefaultChannelPipeline.java:564)
at org.jboss.netty.channel.DefaultChannelPipeline.sendUpstream(DefaultChannelPipeline.java:559)
at org.jboss.netty.channel.Channels.fireMessageReceived(Channels.java:268)
at org.jboss.netty.channel.Channels.fireMessageReceived(Channels.java:255)
at org.jboss.netty.channel.socket.nio.NioWorker.read(NioWorker.java:88)
at org.jboss.netty.channel.socket.nio.AbstractNioWorker.process(AbstractNioWorker.java:109)
at org.jboss.netty.channel.socket.nio.AbstractNioSelector.run(AbstractNioSelector.java:312)
at org.jboss.netty.channel.socket.nio.AbstractNioWorker.run(AbstractNioWorker.java:90)
at org.jboss.netty.channel.socket.nio.NioWorker.run(NioWorker.java:178)
at org.jboss.netty.util.ThreadRenamingRunnable.run(ThreadRenamingRunnable.java:108)
at org.jboss.netty.util.internal.DeadLockProofWorker$1.run(DeadLockProofWorker.java:42)
at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1142)
at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:617)
at java.lang.Thread.run(Thread.java:745)
```

回答

在单个父目录中创建大量的znode后，当客户端尝试在单个请求中获取所有子节点时，服务端将无法返回，因为结果将超出可存储在znode上的数据的最大长度。

为了避免这个问题，应该根据客户端应用的实际情况将“jute.maxbuffer”参数配置为一个更高的值。

“jute.maxbuffer”只能设置为Java系统属性，且没有zookeeper前缀。如果要将“jute.maxbuffer”的值设为X，在ZooKeeper客户端或服务端启动时传入以下系统属性：-Djute.maxbuffer=X。

例如，将参数值设置为4MB：-Djute.maxbuffer=0x400000。

表 31-7 配置参数

参数	描述	默认值
jute.maxbuffer	指定可以存储在znode中的数据的最大长度。单位是Byte。默认值为0xfffff，即低于1MB。 说明 如果更改此选项，则必须在所有服务器和客户端上设置该系统属性，否则将出现问题。	0xfffff

31.5.3 为什么 ZooKeeper 节点上 netcat 命令无法正常运行

问题

为什么在Zookeeper服务器上启用安全的netty配置时，四个字母的命令不能与linux的 *netcat*命令一起使用？

例如：

```
echo stat /netcat host port
```

回答

Linux的 *netcat*命令没有与Zookeeper服务器安全通信的选项，所以当启用安全的netty配置时，它不能支持Zookeeper四个字母的命令。

为了避免这个问题，用户可以使用下面的Java API来执行四个字母的命令。

```
org.apache.zookeeper.client.FourLetterWordMain
```

例如：

```
String[] args = new String[]{host, port, "stat"};  
org.apache.zookeeper.client.FourLetterWordMain.main(args);
```

说明

*netcat*命令只能用于非安全的netty配置。

31.5.4 如何查看哪个 ZooKeeper 实例是 Leader

问题

如何查看ZooKeeper实例的角色是Leader还是Follower？

回答

1. 登录集群Manager管理界面，选择“集群 > 服务 > ZooKeeper > 实例”。
2. 单击相应的quorumpeer实例名称，进入对应实例的详情页面。
3. 查看该实例的“服务器状态”。

31.5.5 使用 IBM JDK 时客户端无法连接 ZooKeeper

问题

使用IBM的JDK的情况下客户端连接ZooKeeper失败。

回答

可能因为IBM的JDK和普通JDK的jaas.conf文件格式不一样。

在使用IBM JDK时，建议使用如下jaas.conf文件模板，其中“useKeytab”中的文件路径必须以“file://”开头，后面为绝对路径。

```
Client {  
  com.ibm.security.auth.module.Krb5LoginModule required  
  useKeytab="file:///D:/install/HbaseClientSample/conf/user.keytab"  
  principal="hbaseuser1"  
  credsType="both";  
};
```

31.5.6 ZooKeeper 客户端刷新 TGT 失败

问题

ZooKeeper客户端刷新TGT失败，无法连接ZooKeeper。报错内容如下：

```
Login: Could not renew TGT due to problem running shell command: '*/kinit -R'; exception  
was:org.apache.zookeeper.Shell$ExitCodeException: kinit: Ticket expired while renewing credentials
```

回答

ZooKeeper使用系统命令**kinit -R**对票据进行刷新，当前MRS版本已经取消了该命令的功能，如需运行长任务，建议使用keytab方式完成鉴权功能。

在“客户端安装路径/ZooKeeper/zookeeper/conf/jaas.conf”配置文件中设置属性“useTicketCache=false”，设置“useKeyTab=true”，并指明keytab路径。

31.5.7 使用 deleteall 命令删除大量 znode 时偶现报错 “Node does not exist”

问题

客户端连接非Leader实例，使用deleteall命令删除大量znode时，报错Node does not exist，但是stat命令能够获取到node状态。

回答

由于网络问题或者数据量大导致leader和follower数据不同步。

解决方法是客户端连接到Leader实例进行删除操作。

具体过程是首先根据[如何查看哪个ZooKeeper实例是Leader](#)查看Leader所在节点IP，使用连接客户端命令**zkCli.sh -server Leader节点IP:2181**成功后进行deleteall命令删除操作。具体操作请参见[使用ZooKeeper客户端](#)。

32 常见操作

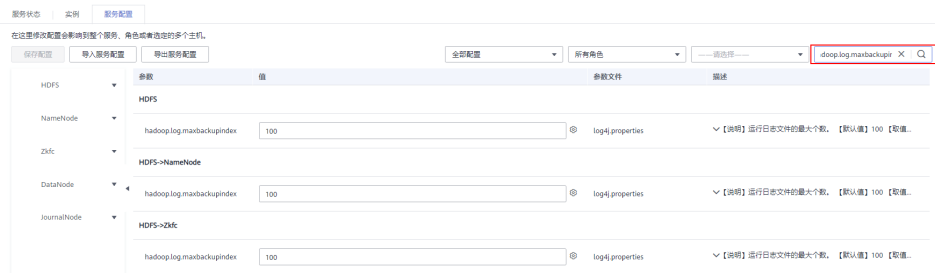
32.1 修改集群服务配置参数

- MRS 3.x之前版本，用户可直接通过MRS管理控制台的集群管理页面修改各服务配置参数：
 - a. 登录MRS控制台，在左侧导航栏选择“现有集群”，单击集群名称。
 - b. 选择“组件管理 > 服务名称 > 服务配置”。

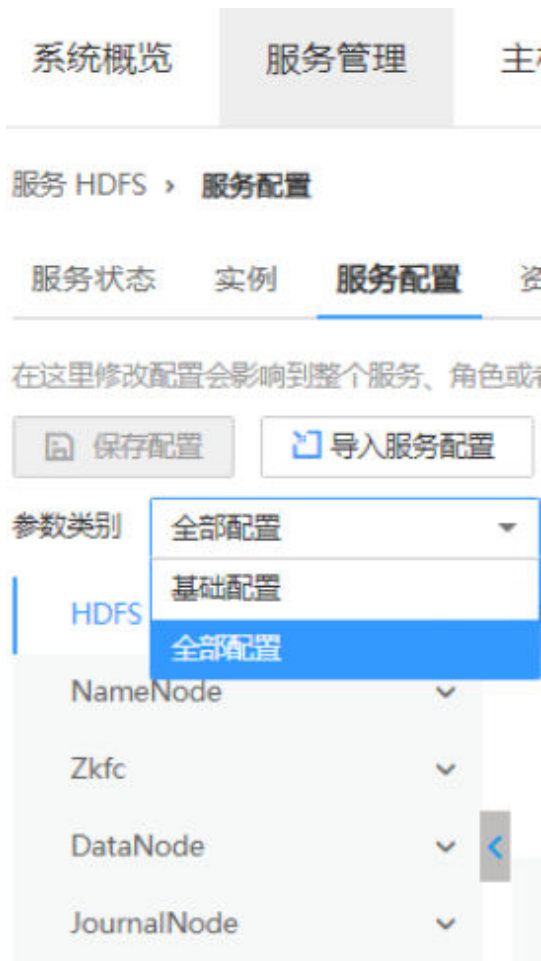
默认显示“基础配置”，如果需要修改更多参数，请选择“全部配置”，界面上将显示该服务的全部配置参数导航树，导航树从上到下的一级节点分别为服务名称和角色名称。展开一级节点后显示参数分类。（下图以HDFS组件为例）



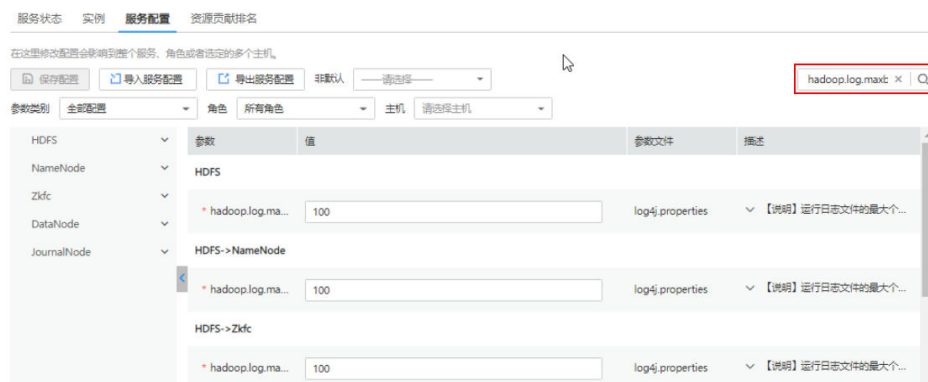
- c. 在导航树选择指定的参数分类，并在右侧修改参数值。
不确定参数的具体位置时，支持在右上角输入参数名，系统将实时进行搜索并显示结果。（下图以HDFS组件为例）



- d. 单击“保存配置”，并在确认对话框中单击“是”。
 - e. 等待界面提示“操作成功”，单击“完成”，配置已修改。
查看集群是否存在配置过期的服务，如果存在，需重启对应服务或角色实例使配置生效。也可在保存配置时直接勾选“重新启动受影响的服务或实例。”。
- MRS 3.x之前的版本，服务配置参数均支持登录MRS Manager进行修改：
 - a. 登录MRS Manager。
 - b. 单击“服务管理”。
 - c. 单击服务视图中指定的服务名称。
 - d. 单击“服务配置”。默认显示“基础配置”，如果需要修改更多参数，请选择“全部配置”，界面上将显示该服务的全部配置参数导航树，导航树从上到下的一级节点分别为服务名称和角色名称。展开一级节点后显示参数分类。（下图以HDFS组件为例）



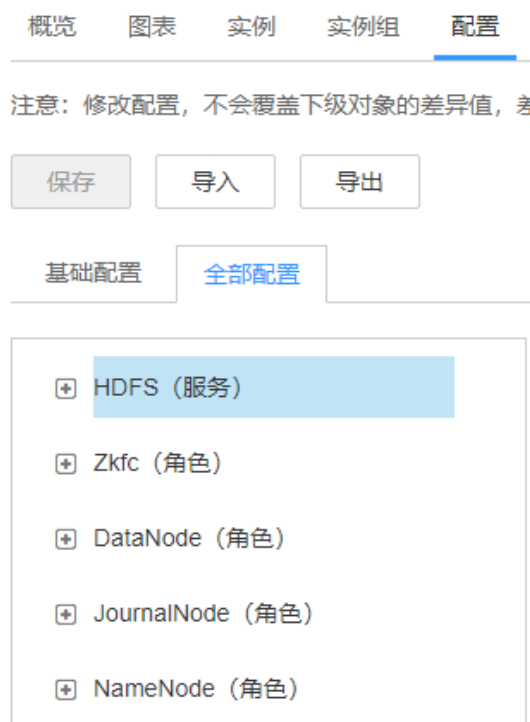
- e. 在导航树选择指定的参数分类，并在右侧修改参数值。
 不确定参数的具体位置时，支持在右上角输入参数名，Manager将实时进行搜索并显示结果。（下图以HDFS组件为例）



- f. 单击“保存配置”，并在确认对话框中单击“是”。
- g. 等待界面提示“操作成功”，单击“完成”，配置已修改。
 查看集群是否存在配置过期的服务，如果存在，需重启对应服务或角色实例使配置生效。也可在保存配置时直接勾选“重新启动受影响的服务或实例。”
- MRS 3.x及后续版本，服务配置参数均支持登录FusionInsight Manager进行修改：

- a. 登录FusionInsight Manager。
- b. 选择“集群 > 服务”。
- c. 单击服务视图中指定的服务名称。
- d. 单击“配置”。

默认显示“基础配置”，如果需要修改更多参数，请选择“全部配置”，界面上将显示该服务的全部配置参数导航树，导航树从上到下的一级节点分别为服务名称和角色名称。展开一级节点后显示参数分类。（下图以HDFS组件为例）



- e. 在导航树选择指定的参数分类，并在右侧修改参数值。

不确定参数的具体位置时，支持在右上角输入参数名，Manager将实时进行搜索并显示结果。（下图以HDFS组件为例）



- f. 单击“保存”，并在确认对话框中单击“确定”。
- g. 等待界面提示“操作成功”，单击“完成”，配置已修改。

查看集群是否存在配置过期的服务，如果存在，需重启对应服务或角色实例使配置生效。

32.2 访问集群 Manager

32.2.1 访问 MRS Manager（MRS 3.x 之前版本）

操作场景

MRS 3.x之前版本集群使用MRS Manager对集群进行监控、配置和管理，用户可以在MRS控制台页面打开Manager管理页面。

访问 MRS Manager

步骤1 登录MRS管理控制台页面。

步骤2 单击“现有集群”，在集群列表中单击指定的集群名称，进入集群信息页面。

步骤3 单击“前往 Manager”，打开“访问MRS Manager页面”。

- 若用户创建集群时已经绑定弹性公网IP，如[图32-1](#)所示：
 - a. 选择待添加的安全组规则所在安全组，该安全组在创建群时配置。
 - b. 添加安全组规则，默认填充的是用户访问公网IP地址9022端口的规则，如需开放多个IP段为可信范围用于访问MRS Manager页面，请参考[步骤6~步骤9](#)。如需对安全组规则进行查看，修改和删除操作，请单击“管理安全组规则”。

说明

- 自动获取的访问公网IP与用户本机IP不一致，属于正常现象，无需处理。
- 9022端口为knox的端口，需要开启访问knox的9022端口权限，才能访问MRS Manager服务。
- c. 勾选“我确认xx.xx.xx.xx为可信任的公网访问IP，并允许从该IP访问MRS Manager页面。”

图 32-1 添加访问 MRS Manager 的安全组规则

访问MRS Manager页面

访问MRS Manager页面需绑定弹性公网IP以及添加安全组规则。 [了解更多](#)

弹性公网IP   [管理弹性公网IP](#) 

安全组 

添加安全组规则   [管理安全组规则](#)

我确认  为可信任的公网访问IP，并允许从该IP访问MRS Manager页面。

确定

取消

- 若用户创建集群时暂未绑定弹性公网IP，如[图32-2](#)所示：
 - 在弹性公网IP下拉框中选择可用的弹性公网IP或单击“管理弹性公网IP”购买弹性公网IP。
 - 选择待添加的安全组规则所在安全组，该安全组在创建群时配置。
 - 添加安全组规则，默认填充的是用户访问公网IP地址9022端口的规则，如需开放多个IP段为可信范围用于访问MRS Manager页面，请参考[步骤6](#)~[步骤9](#)。如需对安全组规则进行查看，修改和删除操作，请单击“管理安全组规则”。

说明

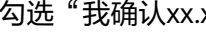
- 自动获取的访问公网IP与用户本机IP不一致，属于正常现象，无需处理。
 - 9022端口为knox的端口，需要开启访问knox的9022端口权限，才能访问MRS Manager服务。
- 勾选“我确认为可信任的公网访问IP，并允许从该IP访问MRS Manager页面。”

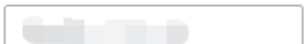
图 32-2 绑定弹性公网 IP

访问MRS Manager页面

访问MRS Manager页面需绑定弹性公网IP以及添加安全组规则。 [了解更多](#)

弹性公网IP  [管理弹性公网IP](#) 

安全组 

添加安全组规则  [管理安全组规则](#)

我确认  为可信任的公网访问IP，并允许从该IP访问MRS Manager页面。

确定

取消

步骤4 单击“确定”，进入MRS Manager登录页面。

步骤5 输入默认用户名“admin”及创建集群时设置的密码，单击“登录”进入MRS Manager页面。

步骤6 在MRS管理控制台，在“现有集群”列表，单击指定的集群名称，进入集群信息页面。

说明

如需给其他用户开通访问MRS Manager的权限，请执行**步骤6-步骤9**，添加对应用户访问公网的IP地址为可信范围。

步骤7 单击弹性公网IP后边的“添加安全组规则”，如**图32-3**所示。

图 32-3 集群详情



步骤8 进入“添加安全组规则”页面，添加需要开放权限用户访问公网的IP地址段并勾选“我确认这里设置的授权对象是可信任的公网访问IP范围，禁止使用0.0.0.0/0,否则会有安全风险。”如图32-4所示。

图 32-4 添加安全组规则



默认填充的是用户访问公网的IP地址，用户可根据需要修改IP地址段，如需开放多个IP段为可信范围，请重复执行**步骤6-步骤9**。如需对安全组规则进行查看，修改和删除操作，请单击“管理安全组规则”。

步骤9 单击“确定”完成安全组规则添加。

----结束

为其他用户开通访问 MRS Manager 的权限

步骤1 在MRS管理控制台，在“现有集群”列表，单击指定的集群名称，进入集群信息页面。

步骤2 单击弹性公网IP后边的“添加安全组规则”，如**图32-5**所示。

图 32-5 集群详情



步骤3 进入“添加安全组规则”页面，添加需要开放权限用户访问公网的IP地址段并勾选“我确认这里设置的授权对象是可信任的公网访问IP范围，禁止使用0.0.0.0/0,否则会有安全风险。”如图32-6所示。

图 32-6 添加安全组规则



默认填充的是用户访问公网的 IP 地址，用户可根据需要修改 IP 地址段，如需开放多个 IP 段为可信范围，请重复执行 [步骤1-步骤4](#)。如需对安全组规则进行查看，修改和删除操作，请单击“管理安全组规则”。

步骤4 单击“确定”完成安全组规则添加。

----结束

32.2.2 访问 FusionInsight Manager（MRS 3.x 及之后版本）

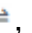
操作场景

MRS 3.x 及之后版本的集群使用 FusionInsight Manager 对集群进行监控、配置和管理。用户在集群安装后可使用账号登录 FusionInsight Manager。

当前支持以下几种方式访问 FusionInsight Manager，请根据实际情况选择。

- [通过弹性 IP 访问 FusionInsight Manager。](#)
- [通过云专线访问 FusionInsight Manager。](#)
- [通过 ECS 访问 FusionInsight Manager。](#)

其中弹性 IP 访问和专线访问可以在 MRS 集群管理控制台上进行切换，具体切换操作步骤如下：

登录 MRS 管理控制台，单击待操作的 MRS 集群，在集群“概览”页面单击“集群管理页面”后的 ，在打开的页面中进行切换。

说明

如果不能正常登录组件的 WebUI 页面，请参考 [通过 ECS 访问 FusionInsight Manager](#) 方式访问 FusionInsight Manager。

集群处于以下状态时无法访问 FusionInsight Manager：

启动中、停止中、停止、删除中、已删除、冻结。

通过弹性 IP 访问 FusionInsight Manager

步骤1 登录MRS管理控制台页面。

步骤2 单击“现有集群”，在集群列表中单击指定的集群名称，进入集群信息页面。

步骤3 单击“集群管理页面”后的“前往 Manager”，在弹出的窗口中配置弹性IP信息。

1. 若创建MRS集群时暂未绑定弹性公网IP，在“弹性公网IP”下拉框中选择可用的弹性公网IP。若用户创建集群时已经绑定弹性公网IP，直接执行**步骤3.2**。

📖 说明

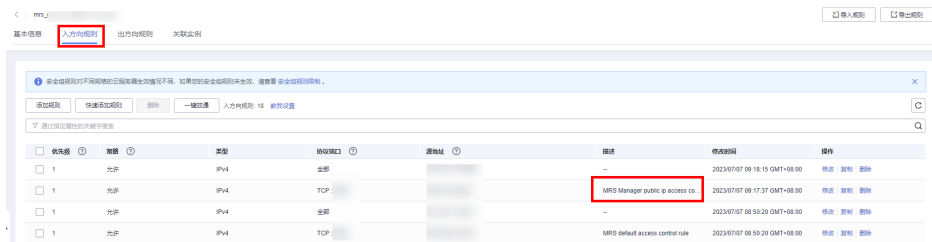
- 如果没有弹性公网IP，可先单击“管理弹性公网IP”购买弹性公网IP后，然后在弹性公网IP下拉框中选择购买的弹性公网IP。
 - 如果在使用完后需要解绑或释放弹性公网IP，请登录“弹性公网IP”界面，在待操作的弹性公网IP后，单击“操作”列的“解绑”或“更多 > 释放”。
 - 如果已创建弹性公网IP，但在绑定时无法找到，可能是由于该弹性公网IP被其他集群绑定，请先在弹性公网IP界面解绑，然后再为当前集群绑定。
2. 在“安全组”中选择当前集群所在的安全组，该安全组在创建集群时配置或集群自动创建。

📖 说明

- 创建自定义集群时，安全组可配置提前创建的安全组或保持默认“自动创建”；快速创建集群时，安全组由集群自动创建。
 - 安全组名称可在集群的“概览”界面的“安全组”查看。
3. 添加安全组规则，默认填充的是用户访问弹性IP地址的规则，如需开放多个IP段为可信范围用于访问Manager页面，请参考**步骤6~步骤9**。如需对安全组规则进行查看，修改和删除操作，请单击“管理安全组规则”。

📖 说明

添加安全组规则会在“安全组 > 入方向规则”（页面入口：单击“管理安全组规则”）页签列表中描述列自动增加“MRS Manager public ip access control rule”，便于用户识别。



4. 勾选确认信息后，单击“确定”。

📖 说明

单击“前往 Manager”右侧的  按钮，可以切换访问FusionInsight Manager的方式，云专线访问请参考**通过云专线访问FusionInsight Manager**。

步骤4 单击“确定”，进入Manager登录页面。

步骤5 输入默认用户名“admin”及创建集群时设置的密码，单击“登录”进入Manager页面。

步骤6 在MRS管理控制台，在“现有集群”列表，单击指定的集群名称，进入集群信息页面。

说明

如需给其他用户开通访问Manager的权限，请执行**步骤6**~**步骤9**，添加对应用户访问公网的IP地址为可信范围。

步骤7 单击弹性公网IP后边的“添加安全组规则”如**图32-7**所示。

图 32-7 集群详情页面



步骤8 进入“添加安全组规则”页面，添加需要开放权限用户访问公网的IP地址段并勾选“我确认这里设置的公网IP/端口号是可信任的公网访问IP范围，我了解使用0.0.0.0/0会带来安全风险”如**图32-8**所示。

图 32-8 添加安全组规则



默认填充的是用户访问公网的IP地址，用户可根据需要修改IP地址段，如需开放多个IP段为可信范围，请重复执行**步骤6**-**步骤9**。如需对安全组规则进行查看，修改和删除操作，请单击“管理安全组规则”。

步骤9 单击“确定”完成安全组规则添加。

----结束

通过云专线访问 FusionInsight Manager

操作前请确保云专线服务可用，并已打通本地数据中心到线上VPC的连接通道。云专线详情请参考[什么是云专线](#)。

步骤1 登录MRS管理控制台。

步骤2 单击集群名称进入集群详情页。

步骤3 在集群详情页的“概览”页签，单击“集群管理页面”右侧的“前往 Manager”。

步骤4 “访问方式”选择“专线访问”，并勾选“我确认已打通本地与浮动IP的网络，可使用专线直接访问MRS Manager。”。

浮动IP为MRS为您访问MRS Manager页面自动分配的IP地址，使用专线访问MRS Manager之前您确保云专线服务已打通本地数据中心到线上VPC的连接通道。

访问MRS Manager页面

您可以选择使用弹性公网IP或者专线访问MRS Manager页面。 [了解更多](#)

访问方式

EIP访问

专线访问

浮动IP 

192 . 168 . 2 . 186

我确认已打通本地与浮动IP的网络，可使用专线直接访问MRS Manager。

确定

取消

步骤5 单击“确定”，进入MRS Manager登录页面，用户名使用“admin”，密码为创建集群时设置的admin密码。

----结束

通过 ECS 访问 FusionInsight Manager

步骤1 进入MRS管理控制台。

步骤2 在“现有集群”列表中，单击指定的集群名称。

记录集群的“可用区”、“虚拟私有云”、“集群管理页面”、“安全组”。

步骤3 在管理控制台首页服务列表中选择“弹性云服务器”，进入ECS管理控制台，创建一个新的弹性云服务器。

- 弹性云服务器的“可用区”、“虚拟私有云”、“安全组”，需要和待访问集群的配置相同。
- 选择一个Windows系统的公共镜像。例如，选择一个标准镜像“Windows Server 2012 R2 Standard 64bit(40GB)”。
- 其他配置参数详细信息，请参见[购买弹性云服务器](#)。

📖 说明

如果ECS的安全组和Master节点的“默认安全组”不同，用户可以选择以下任一种方法修改配置：

- 将ECS的安全组修改为Master节点的默认安全组，请参见[更改安全组](#)。
- 在集群Master节点和Core节点的安全组添加两条安全组规则使ECS可以访问集群，“协议”需选择为“TCP”，“端口”需分别选择“28443”和“20009”。请参见[创建安全组](#)。

如果界面提示“添加安全组规则失败”，请检查安全组配额是否不足，请增加配额或删除不再使用的安全组规则。

步骤4 在EIP管理控制台，申请一个弹性IP地址，并与ECS绑定。

步骤5 登录弹性云服务器。

登录ECS需要Windows系统的账号、密码，弹性IP地址以及配置安全组规则。具体请参见[Windows云服务器登录方式](#)。

步骤6 在Windows的远程桌面中，打开浏览器访问Manager。

Manager访问地址为“集群管理页面”地址。访问时需要输入集群的用户名和密码，例如“admin”用户。

📖 说明

- 如果使用其他集群用户访问Manager，第一次访问时需要修改密码。新密码需要满足集群当前的用户密码复杂度策略。请咨询管理员。
- 默认情况下，在登录时输入5次错误密码将锁定用户，需等待5分钟自动解锁。

步骤7 注销用户退出Manager时移动鼠标到右上角 ，然后单击“注销”。

----结束

32.3 使用 MRS 客户端

32.3.1 安装客户端（3.x 及之后版本）

操作场景

该操作指导用户在MRS集群创建成功后安装MRS集群所有服务（不包含Flume）的客户端。Flume客户端安装请参见[安装Flume客户端](#)。

客户端可以安装在集群内的节点上，也可以安装在集群外的节点上。

修改集群内组件的服务端配置后，建议重新安装客户端，否则客户端与服务端版本将不一致，可能影响正常使用。

操作视频

本视频将以MRS 3.1.0版本集群为例为您介绍手动安装及使用客户端的操作方法。

📖 说明

因不同版本操作界面可能存在差异，相关视频供参考，具体以实际环境为准。

前提条件

- 待安装客户端节点为集群外节点时，该节点必须能够与集群内节点网络互通，否则安装会失败。
- 待安装客户端节点必须启用NTP服务，并保持与服务端的时间一致，否则安装会失败。
- 在节点上安装客户端可以使用root或任意操作系统用户进行操作，要求该用户对客户端文件存放目录和安装目录具有操作权限，两个目录的权限为“755”。
本章节以使用操作系统用户“user_client”安装客户端进行举例，安装目录为“/opt/hadoopclient”。
- 使用omm和root以外的用户安装客户端时，若“/var/tmp/patch”目录已存在，需将此目录权限修改为“777”，将此目录内的日志权限修改为“666”。

集群内节点安装客户端

步骤1 获取客户端软件包。

参考[访问FusionInsight Manager（MRS 3.x及之后版本）](#)访问Manager，单击“集群”在“概览”选择“更多 > 下载客户端”，弹出“下载集群客户端”提示框，选择相关下载参数后单击“确定”。

图 32-9 下载客户端



说明

- 在Manager主页下载的客户端软件包，包含了集群内所有服务（除Flume之外）的客户端。如果只需要安装单个服务的客户端，请选择“集群 > 服务 > 服务名称 > 更多 > 下载客户端”，弹出“下载客户端”信息提示框。
- MRS 3.3.0及之后版本，可以在主页中直接单击“下载客户端”。

表 32-1 客户端下载参数说明

参数	描述	示例
选择客户端类型	<ul style="list-style-type: none">完整客户端：包含完整客户端软件包及配置文件，适用于非开发任务场景。仅配置文件：仅下载客户端配置文件，适用于应用开发任务中，完整客户端已下载并安装后，管理员通过Manager界面修改了服务端配置，开发人员需要更新客户端配置文件的场景。	完整客户端
选择平台类型	<p>客户端类型必须与待安装客户端的节点架构匹配，否则客户端会安装失败。LTS版本集群仅支持下载与Manager平台类型一致的客户端软件包。</p> <ul style="list-style-type: none">x86_64：可以部署在X86平台的客户端软件包。aarch64：可以部署在TaiShan服务器的客户端软件包。	x86_64
仅保存到如下路径	<p>指定客户端软件包在主OMS节点的存放路径。</p> <ul style="list-style-type: none">勾选“仅保存到如下路径”：自定义客户端软件包在主OMS节点的存放路径，omm用户需拥有该目录的读、写与执行权限。如未修改保存路径，文件生成后将默认保存在集群主OMS节点的“/tmp/FusionInsight-Client”。不勾选“仅保存到如下路径”：文件生成后将自动下载并保存至本地，安装客户端时需将其上传至待安装客户端节点的指定目录。	勾选“仅保存到如下路径”

步骤2 复制客户端软件包到待安装客户端节点的指定目录。

客户端软件包生成后默认保存在集群主OMS节点，若需要在集群内其他节点上安装客户端，需以omm用户登录主OMS节点，执行以下命令复制软件包到指定节点，否则跳过本步骤。

例如复制到“/tmp/clienttemp”目录：

```
scp -p /tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_Client.tar 待安装客户端节点的IP地址:/tmp/clienttemp
```

步骤3 以待安装客户端的用户（如user_client）登录将要安装客户端的节点。

📖 说明

在节点上安装客户端可以使用root或其他任意操作系统用户进行操作，要求该用户对客户端文件存放目录和安装目录具有操作权限，两个目录的权限为“755”。

步骤4 解压客户端软件包。

1. 进入客户端软件包所在的目录，例如“/tmp/clienttemp”。

```
cd /tmp/clienttemp
```

2. 执行如下命令解压安装包获取“FusionInsight_Cluster_1_Services_ClientConfig.tar”：
tar -xvf FusionInsight_Cluster_1_Services_Client.tar
3. 执行**sha256sum**命令校验解压得到的文件，检查回显信息与sha256文件里面的内容是否一致，例如：
sha256sum -c FusionInsight_Cluster_1_Services_ClientConfig.tar.sha256
FusionInsight_Cluster_1_Services_ClientConfig.tar: OK
4. 执行以下命令解压“FusionInsight_Cluster_1_Services_ClientConfig.tar”：
tar -xvf FusionInsight_Cluster_1_Services_ClientConfig.tar

步骤5 进入客户端软件包目录，执行如下命令安装客户端到指定目录。

```
cd FusionInsight_Cluster_1_Services_ClientConfig
```

```
./install.sh 客户端安装目录
```

例如执行**./install.sh /opt/hadoopclient**命令安装客户端，等待客户端安装完成。

```
...
The component client is installed successfully
```

📖 说明

- 客户端安装目录可以不存在，会自动创建。但如果存在，则必须为空，目录路径不能包含空格。且客户端安装目录路径只能包含大写字母、小写字母、数字以及下字符。
- 卸载客户端请手动删除客户端安装目录。
- 如果要求安装后的客户端仅能被该安装用户使用，请在安装时加“-o”参数，例如执行**./install.sh /opt/hadoopclient -o**命令安装客户端。

步骤6 客户端安装完成后，可参考“各组件客户端使用实践”使用客户端。

----结束

集群外节点安装客户端

步骤1 创建一个满足要求的弹性云服务器，要求如下：

- 已准备一个Linux弹性云服务器，主机操作系统及版本建议参见[表32-2](#)。

表 32-2 参考列表

CPU架构	操作系统	支持的版本号
x86计算	Euler	EulerOS 2.5
	SUSE	SUSE Linux Enterprise Server 12 SP4 (SUSE 12.4)
	Red Hat	Red Hat-7.5-x86_64 (Red Hat 7.5)
	CentOS	CentOS-7.6版本 (CentOS 7.6)
鲲鹏计算 (ARM)	Euler	EulerOS 2.8
	CentOS	CentOS-7.6版本 (CentOS 7.6)

同时为弹性云服务分配足够的磁盘空间，例如“40GB”。

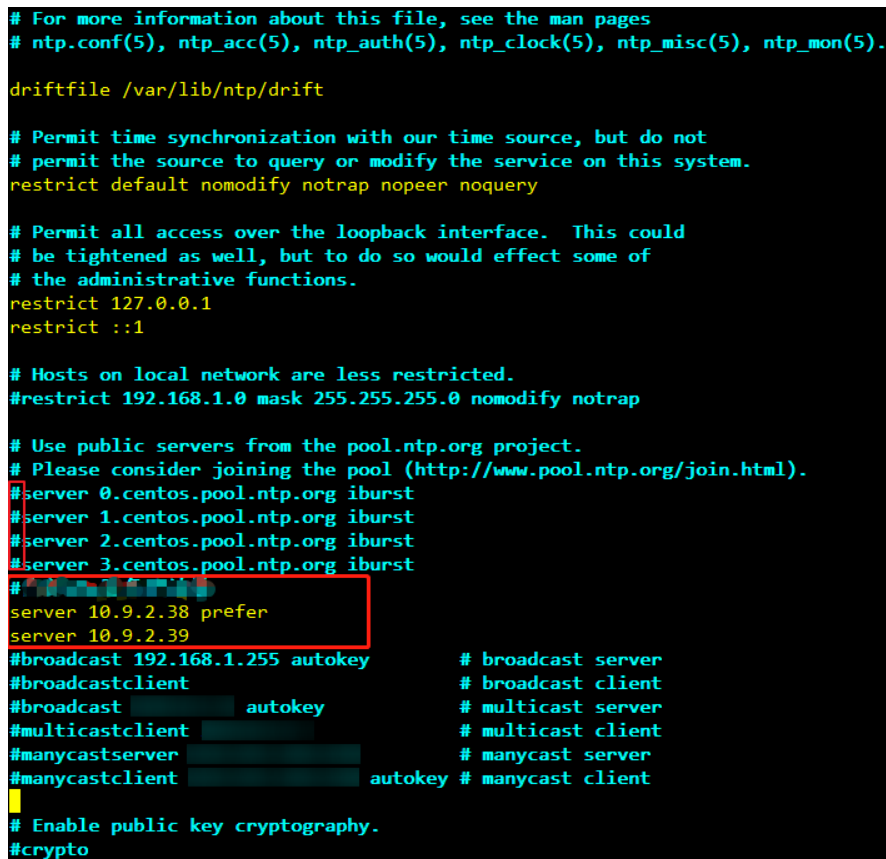
- 弹性云服务器的VPC需要与MRS集群在同一个VPC中。
- 弹性云服务器的安全组需要和MRS集群Master节点的安全组相同。
- 弹性云服务器操作系统已安装NTP服务，且NTP服务运行正常。
若未安装，在配置了yum源的情况下，可执行**yum install ntp -y**命令自行安装。
- 需要允许用户使用密码方式登录Linux弹性云服务器（SSH方式）。
- MRS集群安全组入方向将所有端口对客户节点放开，具体操作请参考[添加安全组规则](#)。

步骤2 执行NTP时间同步，使集群外节点的时间与MRS集群时间同步。

1. 执行**vi /etc/ntp.conf**命令编辑NTP客户端配置文件，并增加MRS集群中Master节点的IP并注释掉其他server的地址。

```
server master1_ip prefer
server master2_ip
```

图 32-10 增加 Master 节点的 IP



```
# For more information about this file, see the man pages
# ntp.conf(5), ntp_acc(5), ntp_auth(5), ntp_clock(5), ntp_misc(5), ntp_mon(5).

driftfile /var/lib/ntp/drift

# Permit time synchronization with our time source, but do not
# permit the source to query or modify the service on this system.
restrict default nomodify notrap nopeer noquery

# Permit all access over the loopback interface. This could
# be tightened as well, but to do so would effect some of
# the administrative functions.
restrict 127.0.0.1
restrict ::1

# Hosts on local network are less restricted.
#restrict 192.168.1.0 mask 255.255.255.0 nomodify notrap

# Use public servers from the pool.ntp.org project.
# Please consider joining the pool (http://www.pool.ntp.org/join.html).
#server 0.centos.pool.ntp.org iburst
#server 1.centos.pool.ntp.org iburst
#server 2.centos.pool.ntp.org iburst
#server 3.centos.pool.ntp.org iburst
#server 4.centos.pool.ntp.org iburst
server 10.9.2.38 prefer
server 10.9.2.39
#broadcast 192.168.1.255 autokey # broadcast server
#broadcastclient # broadcast client
#broadcast # autokey # multicast server
#multicastclient # multicast client
#manycastserver # manycast server
#manycastclient # autokey # manycast client

# Enable public key cryptography.
#crypto
```

2. 执行**service ntpd stop**命令关闭NTP服务。
3. 执行如下命令，手动同步一次时间：
/usr/sbin/ntpdate 192.168.10.8

📖 说明

192.168.10.8为主Master节点的IP地址。

4. 执行**service ntpd start**或**systemctl restart ntpd**命令启动NTP服务。

5. 执行`ntpstat`命令查看时间同步结果。

步骤3 获取客户端软件包。

参考[访问FusionInsight Manager（MRS 3.x及之后版本）](#)访问Manager，单击“集群”在“概览”选择“更多 > 下载客户端”，弹出“下载集群客户端”提示框，选择相关下载参数后单击“确定”。

图 32-11 下载客户端



说明

- 在Manager主页下载的客户端软件包，包含了集群内所有服务（除Flume之外）的客户端。如果只需要安装单个服务的客户端，请选择“集群 > 服务 > 服务名称 > 更多 > 下载客户端”，弹出“下载客户端”信息提示框。
- MRS 3.3.0及之后版本，可以在主页中直接单击“下载客户端”。

表 32-3 客户端下载参数说明

参数	描述	示例
选择客户端类型	<ul style="list-style-type: none"> 完整客户端：包含完整客户端软件包及配置文件，适用于非开发任务场景。 仅配置文件：仅下载客户端配置文件，适用于应用开发任务中，完整客户端已下载并安装后，管理员通过Manager界面修改了服务端配置，开发人员需要更新客户端配置文件的场景。 	完整客户端

参数	描述	示例
选择平台类型	<p>客户端类型必须与待安装客户端的节点架构匹配，否则客户端会安装失败。LTS版本集群仅支持下载与Manager平台类型一致的客户端软件包。</p> <ul style="list-style-type: none"> x86_64：可以部署在X86平台的客户端软件包。 aarch64：可以部署在TaiShan服务器的客户端软件包。 	x86_64
仅保存到如下路径	<p>指定客户端软件包在主OMS节点的存放路径。</p> <ul style="list-style-type: none"> 勾选“仅保存到如下路径”：自定义客户端软件包在主OMS节点的存放路径，omm用户需拥有该目录的读、写与执行权限。如未修改保存路径，文件生成后将默认保存在集群主OMS节点的“/tmp/FusionInsight-Client”。 不勾选“仅保存到如下路径”：文件生成后将自动下载并保存至本地，安装客户端时需将其上传至待安装客户端节点的指定目录。 	勾选“仅保存到如下路径”

步骤4 复制客户端软件包到待安装客户端节点的指定目录。

客户端软件包生成后默认保存在集群主OMS节点，需以omm用户登录主OMS节点，执行以下命令复制软件包到指定弹性云服务器节点。

例如复制到“/tmp/clienttemp”目录：

```
scp -p /tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_Client.tar 待安装客户端节点的IP地址:/tmp/clienttemp
```

步骤5 以待安装客户端的用户（如user_client）登录将要安装客户端的节点。

说明

在节点上安装客户端可以使用root或其他任意操作系统用户进行操作，要求该用户对客户端文件存放目录和安装目录具有操作权限，两个目录的权限为“755”。

步骤6 解压客户端软件包。

1. 进入客户端软件包所在的目录，例如“/tmp/clienttemp”。

```
cd /tmp/clienttemp
```

2. 执行如下命令解压安装包获取“FusionInsight_Cluster_1_Services_ClientConfig.tar”：

```
tar -xvf FusionInsight_Cluster_1_Services_Client.tar
```

3. 执行sha256sum命令校验解压得到的文件，检查回显信息与sha256文件里面的内容是否一致，例如：

```
sha256sum -c FusionInsight_Cluster_1_Services_ClientConfig.tar.sha256
FusionInsight_Cluster_1_Services_ClientConfig.tar: OK
```

4. 执行以下命令解压“FusionInsight_Cluster_1_Services_ClientConfig.tar”：

```
tar -xvf FusionInsight_Cluster_1_Services_ClientConfig.tar
```

步骤7 检查客户端网络连接。

1. 确保客户端所在主机能与解压目录下“hosts”文件（例如“/tmp/FusionInsight_Cluster_1_Services_ClientConfig/hosts”）中所列出的各主机在网络上互通。
2. 当客户端所在主机不是集群中的节点时，需要在客户端所在节点的“/etc/hosts”文件（更改此文件需要root用户权限）中设置集群所有节点主机名和业务平面IP地址映射，主机名和IP地址请保持一一对应，可执行以下步骤在hosts文件中导入集群的域名映射关系。
 - a. 切换至root用户或者其他具有修改hosts文件权限的用户。
`su - root`
 - b. 进入客户端解压目录。
`cd /tmp/clienttemp/FusionInsight_Cluster_1_Services_ClientConfig`
 - c. 执行`cat realm.ini >> /etc/hosts`，将域名映射关系导入到hosts文件中。

说明

- 当客户端所在主机不是集群中的节点时，配置客户端网络连接，可避免执行客户端命令时出现错误。
- 如果采用yarn-client模式运行Spark任务，请在“客户端安装目录/Spark/spark/conf/spark-defaults.conf”文件中添加参数“spark.driver.host”，并将参数值设置为客户端的IP地址。
- 当采用yarn-client模式时，为了Spark WebUI能够正常显示，需要在Yarn的主备节点（即集群中的ResourceManager节点）的hosts文件中，配置客户端的IP地址及主机名对应关系。

步骤8 以待安装客户端的用户（`user_client`）登录将要安装客户端的节点，进入客户端软件包目录，执行如下命令安装客户端到指定目录。

```
cd /tmp/clienttemp/FusionInsight_Cluster_1_Services_ClientConfig
./install.sh 客户端安装目录
```

例如执行`./install.sh /opt/hadoopclient`命令安装客户端，等待客户端安装完成。

```
...
The component client is installed successfully
```

说明

- 客户端安装目录可以不存在，会自动创建。但如果存在，则必须为空，目录路径不能包含空格。且客户端安装目录路径只能包含大写字母、小写字母、数字以及下划线。
- 卸载客户端请手动删除客户端安装目录。
- 如果要求安装后的客户端仅能被该安装用户使用，请在安装时加“-o”参数，例如执行`./install.sh /opt/hadoopclient -o`命令安装客户端。

步骤9 客户端安装完成后，可参考“各组件客户端使用实践”使用客户端。

----结束

32.3.2 安装客户端（3.x 之前版本）

操作场景

用户需要使用MRS客户端。MRS集群客户端可以安装在集群内的Master节点或者Core节点，也可以安装在集群外节点上。

MRS 3.x之前版本集群在集群创建后，在主Master节点默认安装有客户端，可以直接使用，安装目录为“/opt/client”。

MRS 3.x及之后版本客户端的安装请参考[安装客户端（3.x及之后版本）](#)。

说明

如果集群外的节点已安装客户端且只需要更新客户端，请使用安装客户端的用户例如“root”。

在 Core 节点安装客户端

1. 登录MRS Manager页面，选择“服务管理 > 下载客户端”下载客户端安装包至主管理节点。

说明

如仅需更新客户端配置文件，请参考[更新客户端（3.x之前版本）](#)页面的方法二操作。

2. 使用IP地址搜索主管理节点并使用VNC登录主管理节点。
3. 在主管理节点，执行以下命令切换用户。

```
sudo su - omm
```

4. 在MRS管理控制台，查看指定集群“节点管理”页面的“IP”地址。
记录需使用客户端的Core节点IP地址。
5. 在主管理节点，执行以下命令，将客户端安装包从主管理节点文件拷贝到当前Core节点：

```
scp -p /tmp/MRS-client/MRS_Services_Client.tar Core节点的IP地址:/opt/client
```

6. 使用“root”登录Core节点。
Master节点支持Cloud-Init特性，Cloud-init预配置的用户名“root”，密码为创建集群时设置的密码。
7. 执行以下命令，安装客户端：

```
cd /opt/client
tar -xvf MRS_Services_Client.tar
tar -xvf MRS_Services_ClientConfig.tar
cd /opt/client/MRS_Services_ClientConfig
./install.sh 客户端安装目录
```

例如，执行命令：

```
./install.sh /opt/client
```

8. 客户端的使用请参见[使用MRS客户端](#)。

使用 MRS 客户端

1. 在已安装客户端的节点，执行**sudo su - omm**命令切换用户。执行以下命令切换到客户端目录：

```
cd /opt/client
```
2. 执行以下命令配置环境变量：

```
source bigdata_env
```
3. 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

kinit MRS集群用户

例如，**kinit admin**。

说明

启用Kerberos认证的MRS集群默认创建“admin”用户账号，用于集群管理员维护集群。

4. 直接执行组件的客户端命令。

例如：使用HDFS客户端命令查看HDFS根目录文件，执行**hdfs dfs -ls /**。

在集群外节点上安装客户端

步骤1 创建一个满足要求的弹性云服务器，要求如下：

- 针对MRS 3.x之前版本的集群，需要先确认当前MRS集群节点的CPU架构。针对MRS 3.x之前版本的集群，该弹性云服务器的CPU架构请和MRS集群节点保持一致，MRS 3.x及之后版本MRS客户端兼容两种CPU架构。
- 已准备一个弹性云服务器，主机操作系统及版本请参见**表32-4**。

表 32-4 参考列表

CPU架构	操作系统	支持的版本号
x86计算	Euler	- 可用：Euler OS 2.2 - 可用：Euler OS 2.3 - 可用：Euler OS 2.5
鲲鹏计算 (ARM)	Euler	可用：Euler OS 2.8

例如，用户可以选择操作系统为**Euler**的弹性云服务器准备操作。

同时为弹性云服务分配足够的磁盘空间，例如“40GB”。

- 弹性云服务器的VPC需要与MRS集群在同一个VPC中。
- 弹性云服务器的安全组需要和MRS集群Master节点的安全组相同。
如果不同，请修改弹性云服务器安全组或配置弹性云服务器安全组的出入规则允许MRS集群所有安全组的访问。
- 需要允许用户使用密码方式登录Linux弹性云服务器（SSH方式），请参见**SSH密码方式登录**。
- MRS集群安全组入方向将所有端口对客户端节点放开，具体操作请参考**添加安全组规则**。

步骤2 登录MRS Manager页面，具体请参见**访问MRS Manager（MRS 3.x之前版本）**，然后选择“服务管理”。

步骤3 单击“下载客户端”。

步骤4 在“客户端类型”选择“完整客户端”。

步骤5 在“下载路径”选择“远端主机”。

步骤6 将“主机IP”设置为ECS的IP地址，设置“主机端口”为“22”，并将“存放路径”设置为“/tmp”。

- 如果使用SSH登录ECS的默认端口“22”被修改，请将“主机端口”设置为新端口。
- “存放路径”最多可以包含256个字符。

步骤7 “登录用户”设置为“root”。

如果使用其他用户，请确保该用户对保存目录拥有读取、写入和执行权限。

步骤8 在“登录方式”选择“密码”或“SSH私钥”。

- 密码：输入创建集群时设置的root用户密码。
- SSH私钥：选择并上传创建集群时使用的密钥文件。

步骤9 单击“确定”开始生成客户端文件。

若界面显示以下提示信息表示客户端包已经成功保存。单击“关闭”。客户端文件请到下载客户端时设置的远端主机的“存放路径”中获取。

下载客户端文件到远端主机成功。

若界面显示以下提示信息，请检查用户名密码及远端主机的安全组配置，确保用户名密码正确，及远端主机的安全组已增加SSH(22)端口的入方向规则。然后从**步骤2**执行重新下载客户端。

连接到服务器失败，请检查网络连接或参数设置。

说明

生成客户端会占用大量的磁盘IO，不建议在集群处于安装中、启动中、打补丁中等非稳态场景下载客户端。

步骤10 使用VNC方式，登录弹性云服务器。参见[远程登录（VNC方式）](#)。

所有镜像均支持Cloud-init特性。Cloud-init预配置的用户名“root”，密码为创建集群时设置的密码。首次登录建议修改。

步骤11 执行ntp时间同步，使集群外节点的时间与MRS集群时间同步。

1. 检查安装NTP服务有没有安装，未安装请执行`yum install ntp -y`命令自行安装。
2. 执行`vim /etc/ntp.conf`命令编辑NTP客户端配置文件，并增加MRS集群中Master节点的IP并注释掉其他server的地址。

```
server master1_ip prefer  
server master2_ip
```

图 32-12 增加 Master 节点的 IP

```
# For more information about this file, see the man pages
# ntp.conf(5), ntp_acc(5), ntp_auth(5), ntp_clock(5), ntp_misc(5), ntp_mon(5).

driftfile /var/lib/ntp/drift

# Permit time synchronization with our time source, but do not
# permit the source to query or modify the service on this system.
restrict default nomodify notrap nopeer noquery

# Permit all access over the loopback interface. This could
# be tightened as well, but to do so would effect some of
# the administrative functions.
restrict 127.0.0.1
restrict ::1

# Hosts on local network are less restricted.
#restrict 192.168.1.0 mask 255.255.255.0 nomodify notrap

# Use public servers from the pool.ntp.org project.
# Please consider joining the pool (http://www.pool.ntp.org/join.html).
#server 0.centos.pool.ntp.org iburst
#server 1.centos.pool.ntp.org iburst
#server 2.centos.pool.ntp.org iburst
#server 3.centos.pool.ntp.org iburst
#server 10.9.2.38 prefer
#server 10.9.2.39
#broadcast 192.168.1.255 autokey # broadcast server
#broadcastclient # broadcast client
#broadcast # autokey # multicast server
#multicastclient # multicast client
#manycastserver # manycast server
#manycastclient # autokey # manycast client

# Enable public key cryptography.
#crypto
```

3. 执行 `service ntpd stop` 命令关闭 NTP 服务。
4. 执行如下命令，手动同步一次时间：
`/usr/sbin/ntpdate 192.168.10.8`

📖 说明

192.168.10.8 为主 Master 节点的 IP 地址。

5. 执行 `service ntpd start` 或 `systemctl restart ntpd` 命令启动 NTP 服务。
6. 执行 `ntpstat` 命令查看时间同步结果。

步骤 12 在弹性云服务器，切换到 `root` 用户，并将 **步骤 6** 中“存放路径”中的安装包复制到目录“`/opt`”，例如“存放路径”设置为“`/tmp`”时命令如下。

```
sudo su - root
cp /tmp/MRS_Services_Client.tar /opt
```

步骤 13 在“`/opt`”目录执行以下命令，解压压缩包获取校验文件与客户端配置包。

```
tar -xvf MRS_Services_Client.tar
```

步骤 14 执行以下命令，校验文件包。

```
sha256sum -c MRS_Services_ClientConfig.tar.sha256
```

界面显示如下：

```
MRS_Services_ClientConfig.tar: OK
```

步骤15 执行以下命令，解压“MRS_Services_ClientConfig.tar”。

```
tar -xvf MRS_Services_ClientConfig.tar
```

步骤16 执行以下命令，安装客户端到新的目录，例如“/opt/Bigdata/client”。安装时自动生成目录。

```
sh /opt/MRS_Services_ClientConfig/install.sh /opt/Bigdata/client
```

查看安装输出信息，如有以下结果表示客户端安装成功：

```
Components client installation is complete.
```

步骤17 验证弹性云服务器节点是否与集群Master节点的IP是否连通？

例如，执行以下命令：**ping Master节点IP地址**

- 是，执行**步骤18**。
- 否，检查VPC、安全组是否正确，是否与MRS集群在相同VPC和安全组，然后执行**步骤18**。

步骤18 执行以下命令配置环境变量：

```
source /opt/Bigdata/client/bigdata_env
```

步骤19 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit MRS集群用户
```

例如，**kinit admin**

步骤20 执行组件的客户端命令。

例如，执行以下命令查看HDFS目录：

```
hdfs dfs -ls /
```

----结束

32.3.3 更新客户端（3.x 及之后版本）

集群提供了客户端，可以在连接服务端、查看任务结果或管理数据的场景中使用。用户如果在Manager修改了服务配置参数并重启了服务，已安装的客户端需要重新下载并安装，或者使用配置文件更新客户端。

更新客户端配置

方法一：

步骤1 访问[FusionInsight Manager（MRS 3.x及之后版本）](#)，在“集群”下拉列表中单击需要操作的集群名称。

步骤2 选择“更多 > 下载客户端 > 仅配置文件”。

此时生成的压缩文件包含所有服务的配置文件。



步骤3 是否在集群的节点中生成配置文件？

- 是，勾选“仅保存到如下路径”，单击“确定”开始生成客户端文件，文件生成后默认保存在主管理节点“/tmp/FusionInsight-Client”。支持自定义其他目录且 **omm** 用户拥有目录的读、写与执行权限。然后执行 **步骤4**。
- 否，单击“确定”指定本地的保存位置，开始下载完整客户端，等待下载完成，然后执行 **步骤4**。

步骤4 使用WinSCP工具，以客户端安装用户将压缩文件保存到客户端安装的目录，例如“/opt/hadoopclient”。

步骤5 解压软件包。

例如下载的客户端文件为“FusionInsight_Cluster_1_Services_Client.tar”执行如下命令进入客户端所在目录，解压文件到本地目录。

```
cd /opt/hadoopclient
```

```
tar -xvf FusionInsight_Cluster_1_Services_Client.tar
```

步骤6 校验软件包。

执行sha256sum命令校验解压得到的文件，检查回显信息与sha256文件里面的内容是否一致，例如：

```
sha256sum -c  
FusionInsight_Cluster_1_Services_ClientConfig_ConfigFiles.tar.sha256
```

```
FusionInsight_Cluster_1_Services_ClientConfig_ConfigFiles.tar: OK
```

步骤7 解压获取配置文件。

```
tar -xvf FusionInsight_Cluster_1_Services_ClientConfig_ConfigFiles.tar
```

步骤8 在客户端安装目录下执行如下命令，使用配置文件更新客户端。

```
sh refreshConfig.sh 客户端安装目录 配置文件所在目录
```

例如，执行以下命令：

```
sh refreshConfig.sh /opt/hadoopclient /opt/hadoopclient/  
FusionInsight_Cluster_1_Services_ClientConfig_ConfigFiles
```

界面显示以下信息表示配置刷新更新成功：

```
Succeed to refresh components client config.
```

----结束

方法二：

步骤1 以root用户登录客户端安装节点。

步骤2 进入客户端安装的目录，例如“/opt/hadoopclient”，执行以下命令更新配置文件：

```
cd /opt/hadoopclient
```

```
sh autoRefreshConfig.sh
```

步骤3 按照提示输入FusionInsight Manager管理员用户名，密码以及OMS浮动IP地址。

说明

OMS浮动IP地址获取方式：远程登录Master2节点，执行“ifconfig”命令，系统回显中“eth0:wsom”表示OMS浮动IP地址，请记录“inet”的实际参数值。如果在Master2节点无法查询到OMS浮动IP地址，请切换到Master1节点查询并记录。如果只有一个Master节点时，直接在该Master节点查询并记录。

步骤4 输入需要更新配置的组件名，组件名之间使用“,”分隔。如需更新所有组件配置，可直接单击回车键。

界面显示以下信息表示配置刷新更新成功：

```
Succeed to refresh components client config.
```

----结束

32.3.4 更新客户端（3.x 之前版本）

说明

本章节适用于MRS 3.x之前版本的集群。MRS 3.x及之后版本，请参考[更新客户端（3.x及之后版本）](#)。

更新客户端配置文件

操作场景

MRS集群提供了客户端，可以在连接服务端、查看任务结果或管理数据的场景中使用。用户使用MRS的客户端时，如果在MRS Manager修改了服务配置参数并重启了服务，需要先下载并更新客户端配置文件。

用户创建集群时，默认在集群所有节点的“/opt/client”目录安装保存了原始客户端。集群创建完成后，仅Master节点的客户端可以直接使用，Core节点客户端在使用前需要更新客户端配置文件。

操作步骤

方法一：所有版本集群均支持使用。

步骤1 登录MRS Manager页面，具体请参见[访问MRS Manager（MRS 3.x之前版本）](#)，然后选择“服务管理”。

步骤2 单击“下载客户端”。

“客户端类型”选择“仅配置文件”，“下载路径”选择“服务器端”，单击“确定”开始生成客户端配置文件，文件生成后默认保存在主管理节点“/tmp/MRS-client”。文件保存路径支持自定义。

图 32-13 下载客户端配置文件

下载客户端

警告：生成客户端会占用大量的磁盘IO，不建议在集群处于安装中、启动中、打补丁中等非稳态场景进行“下载客户端”操作。

* 客户端类型 完整客户端 仅配置文件

* 下载路径 服务器端 远端主机

仅保存到服务器如下路径，如果存在客户端文件，会覆盖路径下已有的客户端文件。

* 客户端路径

确定

取消

步骤3 查询并登录主Master节点。

步骤4 若在集群内使用客户端，执行以下命令切换到omm用户，若在集群外使用客户端，请切换到root用户：

```
sudo su - omm
```

步骤5 执行以下命令切换客户端目录：

```
cd {客户端安装目录}
```

步骤6 执行以下命令，更新客户端配置：

```
sh refreshConfig.sh 客户端安装目录客户端配置文件压缩包完整路径
```

例如，执行命令：

```
sh refreshConfig.sh /opt/Bigdata/client /tmp/MRS-client/  
MRS_Services_Client.tar
```

界面显示以下信息表示配置刷新更新成功：

```
ReFresh components client config is complete.  
Succeed to refresh components client config.
```

----结束

方法二：

步骤1 集群安装完成之后，执行以下命令切换到omm用户，若在集群外使用客户端，请切换到root用户。

```
sudo su - omm
```

步骤2 执行以下命令切换客户端目录。

```
cd {客户端安装目录}
```

步骤3 执行以下命令并按照提示输入MRS Manager有下载权限的用户名和密码（例如，用户名为admin，密码为创建集群时设置的密码），更新客户端配置。

```
sh autoRefreshConfig.sh
```

步骤4 命令执行后显示如下信息，其中XXX表示集群安装的组件名称，如需更新全部组件配置，单击“Enter”键，如需更新部分组件配置，请输入需要更新的组件名称，多个组件名称以逗号分隔。

```
Components "xxx" have been installed in the cluster. Please input the comma-separated names of the components for which you want to update client configurations. If you press Enter without inputting any component name, the client configurations of all components will be updated:
```

界面显示以下信息表示配置更新成功：

```
Succeed to refresh components client config.
```

界面显示以下信息表示用户名或者密码错误：

```
login manager failed,Incorrect username or password.
```

📖 说明

- 该脚本会自动连接到集群并调用refreshConfig.sh脚本下载并刷新客户端配置文件。
- 客户端默认使用安装目录下文件Version中的“wsom=xxx”所配置的浮动IP刷新客户端配置，如需刷新为其他集群的配置文件，请执行本步骤前修改Version文件中“wsom=xxx”的值为对应集群的浮动IP地址。

----结束

全量更新主 Master 节点的原始客户端

场景描述

用户创建集群时，默认在集群所有节点的“/opt/client”目录安装保存了原始客户端。以下操作以“/opt/Bigdata/client”为例进行说明。

- MRS普通集群，在console页面提交作业时，会使用master节点上预置安装的客户端进行作业提交。
- 用户也可使用master节点上预置安装的客户端来连接服务端、查看任务结果或管理数据等

对集群安装补丁后，用户需要重新更新master节点上的客户端，才能保证继续使用内置客户端功能。

操作步骤

步骤1 登录MRS Manager页面，具体请参见[访问MRS Manager（MRS 3.x之前版本）](#)，然后选择“服务管理”。

步骤2 单击“下载客户端”。

“客户端类型”选择“完整客户端”，“下载路径”选择“服务器端”，单击“确定”开始生成客户端配置文件，文件生成后默认保存在主管理节点“/tmp/MRS-client”。文件保存路径支持自定义。

步骤3 查询并登录主Master节点。

步骤4 在弹性云服务器，切换到root用户，并将安装包复制到目录“/opt”。

```
sudo su - root
```

```
cp /tmp/MRS-client/MRS_Services_Client.tar /opt
```

步骤5 在“/opt”目录执行以下命令，解压压缩包获取校验文件与客户端配置包。

```
tar -xvf MRS_Services_Client.tar
```

步骤6 执行以下命令，校验文件包。

```
sha256sum -c MRS_Services_ClientConfig.tar.sha256
```

界面显示如下：

```
MRS_Services_ClientConfig.tar: OK
```

步骤7 执行以下命令，解压“MRS_Services_ClientConfig.tar”。

```
tar -xvf MRS_Services_ClientConfig.tar
```

步骤8 执行以下命令，移动旧客户端到“/opt/Bigdata/client_bak”目录下。

```
mv /opt/Bigdata/client /opt/Bigdata/client_bak
```

步骤9 执行以下命令，安装客户端到新的目录，客户端路径必须为“/opt/Bigdata/client”。

```
sh /opt/MRS_Services_ClientConfig/install.sh /opt/Bigdata/client
```

查看安装输出信息，如有以下结果表示客户端安装成功：

```
Components client installation is complete.
```

步骤10 执行以下命令，修改“/opt/Bigdata/client”目录的所属用户和用户组。

```
chown omm:wheel /opt/Bigdata/client -R
```

步骤11 执行以下命令配置环境变量：

```
source /opt/Bigdata/client/bigdata_env
```

步骤12 如果当前集群已启用Kerberos认证，执行以下命令认证当前用户。如果当前集群未启用Kerberos认证，则无需执行此命令。

```
kinit MRS集群用户
```

例如, **kinit admin**

步骤13 执行组件的客户端命令。

例如，执行以下命令查看HDFS目录：

```
hdfs dfs -ls /
```

----结束

全量更新备 Master 节点的原始客户端

步骤1 参见**步骤1~步骤3**登录备Master节点，执行如下命令切换到omm用户。

```
sudo su - omm
```

步骤2 在备master节点上执行如下命令，从主master节点拷贝下载的客户端包。


```
scp omm@master1节点IP地址:/tmp/MRS-client/MRS_Services_Client.tar /tmp/  
MRS-client/
```

 说明

- 该命令以master1节点为主master节点为例。
- 目的路径以备master节点的/tmp/MRS-client/目录为例，请根据实际路径修改。

步骤3 参见[步骤4](#)~[步骤13](#)，更新备Master节点的客户端。

----结束