

云容器引擎 Autopilot 快速入门

文档版本 01
发布日期 2024-10-12



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

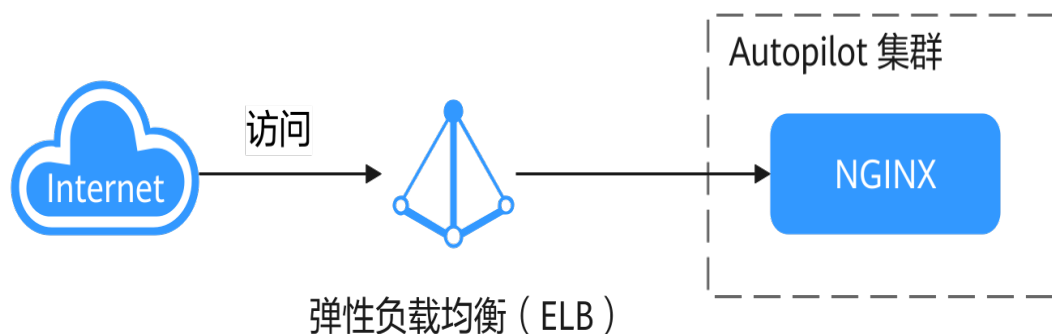
网址：<https://www.huaweicloud.com/>

目录

1 在 CCE Autopilot 集群中部署 Nginx 工作负载.....	1
2 使用 HPA 策略实现工作负载弹性伸缩.....	13

1 在 CCE Autopilot 集群中部署 Nginx 工作负载

CCE Autopilot集群是云容器引擎服务推出的Serverless版集群，提供免运维的容器服务，您无需购买和管理节点即可部署应用，降低运维成本，同时提升应用的可靠性和扩展性。Nginx是一款高性能的开源HTTP服务器和反向代理服务器，广泛用于处理高并发、负载均衡和静态资源服务。本示例以Nginx为例，帮助您了解如何创建CCE Autopilot集群以及在集群中部署工作负载。



操作流程

操作步骤	步骤说明	费用说明
准备工作	您需要注册华为账号，并为账户充值。	不涉及费用。
步骤一：首次开通CCE并进行授权	当您的账号在当前区域中首次使用CCE时，您需要为CCE进行授权。	不涉及费用。
步骤二：创建CCE Autopilot 集群	在CCE服务中创建CCE Autopilot集群，从而简化Kubernetes集群的管理和操作。	涉及集群管理和终端节点等费用，具体请参见 计费说明 。

操作步骤	步骤说明	费用说明
步骤三：创建并访问Nginx工作负载	在集群中创建工作负载以运行您的容器，并为其创建一个服务，然后您就可以从公网访问您的应用。	涉及Pod费用，具体请参见 计费说明 。
后续操作：释放资源	如果您在完成实践后不需要继续使用集群，请及时清理资源以免产生额外扣费。	不涉及费用。


准备工作

- 您需要注册华为账号并完成实名认证，详情请参见[注册华为账号并开通华为云和个人实名认证](#)。
- 您需要确保账户有足够的资金，以免创建集群失败，具体操作请参见[账户充值](#)。

步骤一：首次开通 CCE 并进行授权

由于CCE在运行中对计算、存储、网络以及监控等各类云服务资源都存在依赖关系，因此当您首次登录CCE控制台时，CCE将自动请求获取当前区域下的云资源权限，从而更好地为您提供服务。如果您在当前区域已完成授权，可忽略本步骤。

步骤1 使用华为账号登录[CCE控制台](#)。

步骤2 单击管理控制台左上角的，选择区域。

步骤3 在首次登录某个区域的CCE控制台时将跳出“授权说明”，请您在仔细阅读后单击“确定”。

当您同意授权后，CCE将在IAM中创建名为“cce_admin_trust”委托，统一对您的其他云服务资源进行操作，并且授予其Tenant Administrator权限。Tenant Administrator拥有除IAM管理外的全部云服务管理员权限，用于对CCE所依赖的其他云服务资源进行调用，且该授权仅在当前区域生效。您可前往“统一身份认证服务 IAM > 委托”页签，单击“cce_admin_trust”查看各区域的授权记录。关于资源委托详情，您可参考[委托](#)进行了解。

说明

由于CCE对其他云服务有许多依赖，如果没有Tenant Administrator权限，可能会因为某个服务权限不足而影响CCE功能的正常使用。因此在使用CCE服务期间，请不要自行删除或者修改“cce_admin_trust”委托。

----结束

步骤二：创建 CCE Autopilot 集群

在CCE服务中创建CCE Autopilot集群，从而简化Kubernetes集群的管理和操作。

步骤1 返回[CCE控制台](#)。

- 如果您的帐号在当前区域未创建过集群，请在当前页面单击“购买集群”或“购买CCE Autopilot集群”，进入购买页。
- 如果您的帐号在当前区域已创建过集群，请在左侧菜单栏选择集群管理，单击右上角“购买集群”，进入购买页。

步骤2 配置集群基础信息。

本示例中仅解释必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[购买Autopilot集群](#)。



参数	示例	参数说明
集群类型	CCE Autopilot 集群	<p>CCE支持创建多种类型集群，满足各种业务需求，提供高可靠、安全的商业级容器集群服务。</p> <ul style="list-style-type: none"> ● CCE Standard集群：标准版本集群，提供高可靠、安全的商业级容器集群服务。 ● CCE Turbo集群：拥有更高性能的云原生网络，提供云原生混部调度能力，可实现更高的资源利用率和更广的全场景覆盖。 ● CCE Autopilot集群：Serverless版集群，提供免运维的容器服务，可以大幅降低运维成本，提高应用程序的可靠性和可扩展性。 <p>了解集群类型的更多内容，请参见集群对比。</p>
集群名称	autopilot-example	新建集群的名称。
企业项目	default	<p>该参数仅对开通企业项目的企业客户账号显示，不显示时请忽略。</p> <p>企业项目是一种资源管理单位，可跨区域归类资源，方便企业按部门或项目组集中管理。了解企业项目的更多内容，请参见项目管理。</p> <p>请根据需要选择适合的企业项目，如果没有特殊要求，可以选择default。</p>
集群版本	v1.28	集群安装的Kubernetes软件版本，建议选择最新版本。

步骤3 配置集群网络信息。

本示例中仅解释必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[购买Autopilot集群](#)。



参数	示例	参数说明
虚拟私有云	vpc-autopilot	选择集群所在的虚拟私有云VPC。如果没有可选项，单击右侧“新建虚拟私有云”创建，具体请参见 创建虚拟私有云和子网 。集群创建后不可修改。
容器子网	subnet-502f	选择容器所在子网。如果没有可选项，单击右侧“新建子网”创建，具体请参见 创建虚拟私有云和子网 。容器子网决定了集群下容器的数量上限，集群创建后支持新增子网。
服务网段	10.247.0.0/16	同一集群下容器互相访问时使用的Service资源网段，决定Service资源的上限。集群创建后不可修改。
镜像访问	-	为确保您的集群节点可以从容器镜像服务中拉取镜像，默认使用所选的VPC中已有的终端节点，否则系统将为您新建SWR和OBS的终端节点。 终端节点将产生一定费用，详情请参见 价格计算器（VPC终端节点） 。
配置SNAT	开启	默认开启，开启后您的集群可以通过NAT网关访问公网。默认使用所选的VPC中已有的NAT网关，否则系统将会为您自动创建一个默认规格的NAT网关并绑定弹性公网IP，并自动配置SNAT规则。 使用NAT网关将产生一定费用，详情请参见 价格计算器（NAT网关） 。

步骤4 单击“下一步：插件选择”，选择创建集群时需要安装的插件。

本示例中，仅选择默认安装的必选插件。

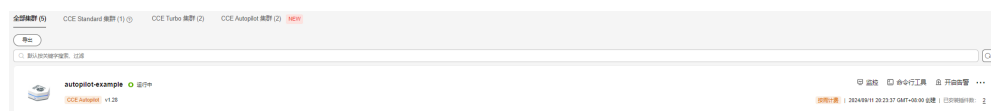
步骤5 单击“下一步：插件配置”，对选择的插件进行配置，其中默认插件无法进行配置。

本示例中，默认插件无需配置。

步骤6 单击“下一步：确认配置”，显示集群资源清单，确认无误后，单击“提交”。

创建集群预计需要5-10分钟左右，请耐心等待。

创建成功后，集群管理中对对应集群的状态为运行中。



----结束

步骤三：创建并访问 Nginx 工作负载

在集群中创建Nginx工作负载，将应用程序或服务部署到容器环境中，实现资源的高效利用和自动化管理。同时，为该工作负载创建负载均衡类型的服务，使您能够从公网访问应用。本节介绍两种方式创建和访问Nginx工作负载，即控制台方式和kubect命令行方式。

使用控制台方式

步骤1 单击新建的集群名称，进入集群控制台。

步骤2 左侧菜单栏选择“工作负载”，单击右上角“创建工作负载”，进入创建页。

步骤3 配置工作负载基本信息。

本示例中仅解释必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[创建工作负载](#)，您可以根据工作负载类型选择适合的参考文档。



参数	示例	参数说明
负载类型	无状态负载 Deployment	<p>工作负载是一种对Pod的抽象管理方式，用于定义和控制Pod的创建、运行和生命周期。通过工作负载，您可以批量管理和自动化控制多个Pod的行为，如伸缩、更新和恢复。</p> <ul style="list-style-type: none"> 无状态负载（Deployment）：管理无状态应用，支持上线部署、滚动升级、创建副本和恢复上线。 有状态负载（StatefulSet）：管理有状态应用，确保每个Pod能够拥有独立的状态。 普通任务（Job）：一次性任务，完成后Pod自动删除。 定时任务（CronJob）：基于时间的Job，指定时间周期内运行指定的Job。 <p>了解工作负载的更多内容，请参见工作负载概述。</p> <p>Nginx主要用于处理请求转发、负载均衡和静态内容分发，不需要在本地保存任何持久性数据，因此本示例将Nginx部署为无状态负载。</p>
负载名称	nginx	请填写工作负载的名称。

参数	示例	参数说明
命名空间	default	命名空间是Kubernetes集群中的抽象概念，可以将集群中的资源或对象划分为一个组，且不同命名空间中的数据彼此隔离，您可以根据需要创建并使用命名空间。 集群创建后会默认生成default命名空间，如果没有特殊要求，可以直接选择default命名空间。
实例数量	1	工作负载中Pod的数量。

步骤4 配置工作负载容器信息。

本示例中仅解释必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[创建工作负载](#)，您可以根据工作负载类型选择适合的参考文档。



参数	示例	参数说明
镜像名称	nginx	单击“选择镜像”，在弹出的窗口中切换至“镜像中心”，选择公共镜像。
镜像版本	latest	选择需要部署的镜像版本。
CPU配额	0.25cores	CPU资源限制值，即CPU资源限制值，即允许容器使用的CPU最大值，防止占用过多资源，默认0.25cores。
内存配额	512MiB	内存资源限制值，即允许容器使用的内存最大值。如果超过，容器会被终止，默认512MiB。

步骤5 配置工作负载服务信息。

单击“服务配置”下的加号，进入创建服务页面。

本示例中需要外部访问Nginx，所以访问类型设置为负载均衡。

本示例仅解释负载均衡的必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[服务 \(Service\)](#)，您可以根据服务类型选择适合的参考文档。

✕

创建服务

Service名称

访问类型

集群内访问

通过集群的内部IP暴露服务，只能在集群内部访问

负载均衡

通过ELB负载均衡对外提供服务，高可用，超高性能，稳定安全

① 集群外访问推荐选择负载均衡访问类型

服务亲和 集群级别 ⓘ

负载均衡器

独享型
网络型 (TCP/UDP) & 应...
选择已有
elb-nginx
🔍 创建负载均衡器

仅支持集群所在 VPC vpc-autopilot 下、实例规格支持网络型 & 应用型的独享型负载均衡实例 [查看约束与限制](#)

负载均衡配置：分配策略：加权轮询算法；会话保持类型：不启用； [编辑](#)

我已阅读 [《负载均衡使用须知》](#)

健康检查

不启用
全局检查
自定义检查

协议：TCP | 检查周期（秒）：5 | 超时时间（秒）：10 | 最大重试次数：3 [🔗](#)

端口配置

协议	容器端口	服务端口	监听器前端协议	操作
TCP	80	8080	TCP	删除
+				

监听器配置

访问控制 未配置

高级配置 [+ 添加高级配置](#)

注解

键 = 值 确认添加 [使用指南](#)

取消
确定

参数	示例	参数说明
Service名称	nginx	请填写服务的名称。
访问类型	负载均衡	<p>选择服务类型，即服务访问的方式。</p> <ul style="list-style-type: none"> 集群内访问：通过集群的内部IP暴露服务，只能在集群内部访问。 负载均衡：通过弹性负载均衡（ELB）对外提供服务，即能够从公网访问到工作负载。 <p>了解服务类型的更多内容，请参见服务（Service）。</p>
负载均衡器	<p>独享型</p> <p>网络型(TCP/UDP)&应用型(HTTP/HTTPS)</p> <p>选择已有</p> <p>elb-nginx</p>	<ul style="list-style-type: none"> 如果已有负载均衡（ELB）实例，可以选择已有ELB。 <p>说明</p> <p>使用已有的ELB时，仅支持集群所在VPC下、实例规格支持网络型的独享型ELB实例。</p> <ul style="list-style-type: none"> 如果没有请选择“自动创建”，创建一个负载均衡器，并绑定弹性公网IP，具体操作请参见创建负载均衡类型的服务。

参数	示例	参数说明
端口配置	协议: TCP	负载均衡监听器端口协议。
	容器端口: 80	容器中应用启动监听的端口, 该容器端口需和应用对外提供的监听端口一致。 使用nginx镜像请设置为80。
	服务端口: 8080	ELB将会使用该端口创建监听器, 提供外部流量访问入口, 可自定义。

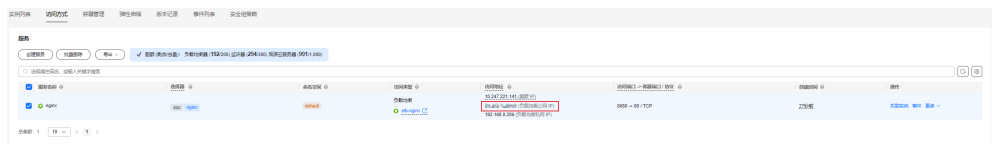
步骤6 单击右下角“创建工作负载”。

创建成功后, 无状态工作负载列表中对工作负载的状态为运行中。

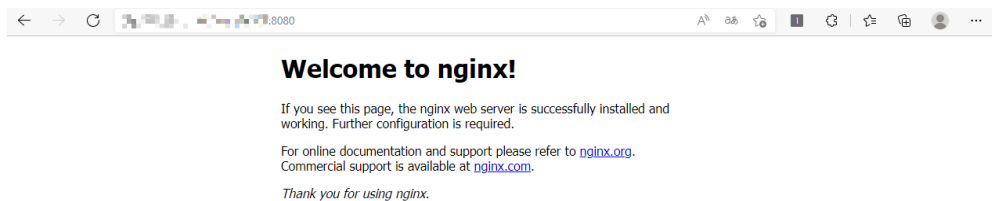


步骤7 获取Nginx的外部访问地址。

单击Nginx工作负载名称, 进入工作负载详情页。“访问方式”页签中, “负载均衡公网IP: 访问端口”即为外部访问地址。



步骤8 在浏览器中输入“负载均衡公网IP: 访问端口”, 即可成功访问应用。



----结束

使用 kubectl 命令行方式

该步骤涉及命令行操作, 您可以使用以下两种方式进行相关操作:

- 通过集群内命令行工具进行操作, 该命令行工具已经配置kubectl命令, 并已连接集群, 更多信息请参见[通过CloudShell连接集群](#)。
- 通过ECS虚拟机进行操作, 该ECS需与集群处于同一VPC, 并绑定弹性公网IP, 具体操作请参见[快速购买和使用Linux ECS](#)。此外, 您还需要安装kubectl命令, 并[通过kubectl连接集群](#)。

以第一种方法为例, 介绍如何使用kubectl命令行方式创建Nginx工作负载。

步骤1 单击新建的集群名称，进入集群控制台。

步骤2 单击右上角“命令行工具”，进入CloudShell页面。

📖 说明

目前，只有部分局点支持CloudShell连接集群，具体情况请以控制台为准。



步骤3 创建YAML文件nginx-deployment.yaml，用于配置nginx工作负载，文件名称可自定义。

```
vim nginx-deployment.yaml
```

文件内容如下：

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx # 工作负载名称
spec:
  replicas: 1 # 实例数量
  selector:
    matchLabels: # 选择器，用于选择带有特定标签的资源
      app: nginx
  template:
    metadata:
      labels: # 标签
        app: nginx
    spec:
      containers:
        - image: nginx:latest # 镜像名称：镜像版本
          name: nginx
          imagePullSecrets:
            - name: default-secret
```

输入完成后，**Esc**键退出编辑，输入:**wq**保存。

步骤4 执行以下命令，创建工作负载。

```
kubectl create -f nginx-deployment.yaml
```

回显如下，表示已经开始创建工作负载。

```
deployment.apps/nginx created
```

步骤5 执行以下命令，查看工作负载状态。

```
kubectl get deployment
```

回显如下，如果工作负载创建的Pod皆为可用状态，则表示创建成功。

```
NAME      READY   UP-TO-DATE   AVAILABLE   AGE
nginx     1/1     1             1           4m59s
```

回显内容的参数解析如下：

参数	示例	参数说明
NAME	ngi nx	工作负载的名称。

参数	示例	参数说明
READY	1/1	表示工作负载的可用状态，显示为“可用Pod个数/期望Pod个数”。
UP-TO-DATE	1	指当前工作负载已经完成更新的Pod数。
AVAILABLE	1	工作负载可用的Pod个数。
AGE	4m59s	工作负载已经运行的时间。

步骤6 创建负载均衡类型的服务，并关联已创建的工作负载nginx。

本示例基于已有的弹性负载均衡（ELB）实例创建服务，如果您需要自动创建ELB请参考[通过kubectl命令行创建-自动创建ELB](#)。

创建YAML文件nginx-elb-svc.yaml，用于配置负载均衡服务，文件名称可自定义。

```
vim nginx-elb-svc.yaml
```

文件内容如下：

```
apiVersion: v1
kind: Service
metadata:
  name: nginx # 服务的名称
  annotations:
    kubernetes.io/elb.id: <your_elb_id> # ELB ID, 替换为实际值
    kubernetes.io/elb.class: performance # 负载均衡器类型
spec:
  selector:
    app: nginx
  ports:
  - name: service0
    port: 8080
    protocol: TCP
    targetPort: 80
  type: LoadBalancer
```

输入完成后，**Esc**键退出编辑，输入:**wq**保存。

参数	示例	参数说明
kubernetes.io/elb.id	405ef586-0397-45c3-bfc4-xxx	已有的ELB的ID。 说明 使用已有的ELB时，仅支持集群所在VPC下、实例规格支持网络型的独享型ELB实例。 获取方式： 进入网络控制台首页，选择“弹性负载均衡 > 我的ELB”，找到对应的ELB名称，名称下方即为对应ID。
kubernetes.io/elb.class	performance	负载均衡器类型，仅支持performance类型，即独享型负载均衡。

参数	示例	参数说明
selector	app: nginx	选择器，服务将流量发送给对应标签的Pod。 须知 selector取值需与工作负载YAML文件中matchLabels参数取值一致，本示例中为app: nginx。
ports.port	8080	弹性负载均衡（ELB）将会使用该端口创建监听器，提供外部流量访问入口，可自定义。
ports.protocol	TCP	负载均衡监听器端口协议。
ports.targetPort	80	Service访问目标容器的端口，此端口与容器中运行的应用强相关。 使用nginx镜像请设置为80。

步骤7 执行以下命令创建服务。

```
kubectl create -f nginx-elb-svc.yaml
```

回显如下，表示服务已创建。

```
service/nginx created
```

步骤8 执行以下命令查看服务。

```
kubectl get svc
```

回显如下，表示工作负载访问方式已设置成功。

NAME	TYPE	CLUSTER-IP	EXTERNAL-IP	PORT(S)	AGE
kubernetes	ClusterIP	10.247.0.1	<none>	443/TCP	18h
nginx	LoadBalancer	10.247.56.18	xx.xx.xx.xx,xx.xx.xx.xx	8080:30581/TCP	5m8s

步骤9 在浏览器中输入“外部访问地址：服务端口”，即可成功访问应用。其中外部访问地址为EXTERNAL-IP对应的第一个IP地址的，服务端口为8080。



Welcome to nginx!

If you see this page, the nginx web server is successfully installed and working. Further configuration is required.

For online documentation and support please refer to nginx.org.
Commercial support is available at nginx.com.

Thank you for using nginx.

----结束

后续操作：释放资源

如果您无需继续使用集群，请及时释放资源，避免产生额外的费用。

步骤1 返回**CCE控制台**，在左侧导航栏中选择“集群管理”。

步骤2 找到需要删除的集群，单击集群卡片右上角的“...”，并单击“删除集群”。

步骤3 在弹出的“删除集群”窗口中，勾选需要删除的所有资源，并根据页面提示删除NAT网关、SNAT出网EIP和终端节点等资源。

删除集群

×

您确定要删除以下 1 台集群吗？删除操作无法恢复，请谨慎操作

集群名称	集群版本	创建时间
autopilot-example	v1.28	2024/09/11 20:23:37 GMT+08:00

- 1. 删除集群会删除集群下工作负载和服务，相关业务将无法恢复。
- 2. 集群关联创建的VPC级别的资源(如终端节点、NAT 网关、SNAT 出网 EIP)，删除集群时默认保留，请确认其他集群或者地方没有重用该资源，再进行删除操作。

请选择您要同步删除的资源：

资源类型	操作
网络资源 (ELB等)	<input type="checkbox"/> 删除 (仅删除自动创建的ELB资源)
NAT网关	VPC级别的资源，请前往网络控制台删除。 前往删除
SNAT出网EIP	VPC级别的资源，请前往网络控制台删除。 前往删除
终端节点	VPC级别的资源，请前往网络控制台删除。 前往删除

如果您确定要删除，请输入 DELETE

DELETE

*注：删除中时请勿关闭当前弹窗或刷新页面，删除完成后弹窗会自动关闭，否则可能导致部分资源残留

否

是

步骤4 在确认框中输入“DELETE”，单击“是”，开始执行删除集群操作。

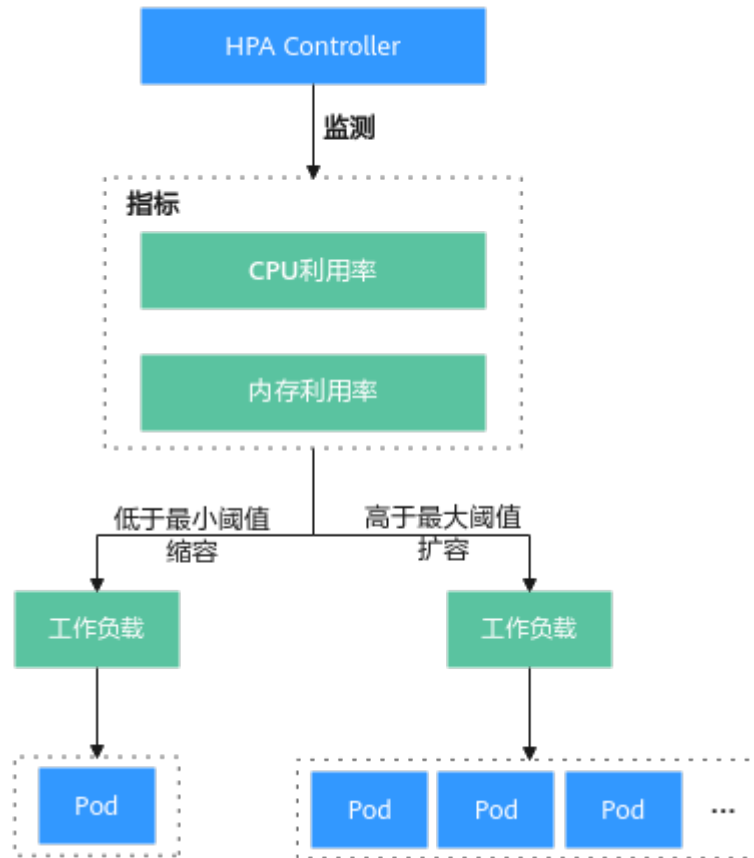
删除集群需要花费1~3分钟，请耐心等待。

---结束

2 使用 HPA 策略实现工作负载弹性伸缩

企业应用的流量大小不是每时每刻都一样，有高峰，有低谷，如果每时每刻都要保持能够扛住高峰流量的机器数目，那么成本会很高。一般解决该问题的办法是根据流量大小或资源占用率自动调节工作负载的数量，也就是弹性伸缩。在CCE Standard/Turbo集群中，配置工作负载弹性伸缩需要制定节点和负载的弹性联合策略，即CA（Cluster AutoScaling）策略和HPA（Horizontal Pod Autoscaling），其中CA负责节点弹性伸缩（避免资源不足导致工作负载创建失败），HPA负责工作负载弹性伸缩。而在CCE Autopilot集群中，您无需对节点的部署、管理和安全性进行维护，只需设置HPA策略，即可按需自动调整Pod实例数量，简化了资源管理和运维流程。此外，CCE Autopilot集群具有高性能的Serverless容器资源底座和多级资源池预热技术，可以实现秒级弹性扩缩，帮助您快速上线新应用，灵活应对市场变化。

HPA策略主要用来控制工作负载创建的Pod实例数量，通过周期性检查Pod的CPU利用率和内存利用率等数据，计算实现HPA资源目标值所需的Pod实例数量，进而调整工作负载的replicas字段。本示例主要介绍CCE Autopilot集群如何通过HPA策略实现工作负载弹性伸缩。



操作流程

操作步骤	步骤说明	费用说明
准备工作	您需要注册华为账号，并为账户充值。	不涉及费用。
步骤一：首次开通 CCE 并进行授权	当您的账号在当前区域中首次使用 CCE 时，您需要为 CCE 进行授权。	不涉及费用。
步骤二：创建 CCE Autopilot 集群	在 CCE 服务中创建 CCE Autopilot 集群，从而简化 Kubernetes 集群的管理和操作。	涉及集群管理和终端节点等费用，具体请参见 计费说明 。
步骤三：创建算力密集型应用并上传 SWR	创建算力密集型应用（用于压力测试），创建完成的应用需要上传至容器镜像服务（SWR），以便在集群中进行部署和管理。	涉及创建 ECS 实例，会产生云服务器和弹性公网 IP 等费用，具体请参见 计费概述 。
步骤四：使用 hpa-example 镜像创建工作负载	使用构建的 hpa-example 镜像创建工作负载，并为该工作负载创建负载均衡类型服务，使您能够从公网访问应用。	涉及 Pod 费用，具体请参见 计费说明 。
步骤五：创建 HPA 策略	创建 HPA 策略并与工作负载关联，从而控制工作负载创建的 Pod 副本数量。	不涉及费用。

操作步骤	步骤说明	费用说明
步骤六：检查工作负载能否自动进行弹性伸缩	检查配置的HPA策略是否有效，即工作负载hpa-example能否自动进行弹性伸缩。	工作负载弹性伸缩时，会涉及Pod实例数量的增加与减少，会涉及Pod费用的变更，具体请参见 计费说明 。
后续操作：释放资源	如果您在完成实践后不需要继续使用集群，请及时清理资源以免产生额外扣费。	不涉及费用。


准备工作

- 您需要注册华为账号并完成实名认证，详情请参见[注册华为账号并开通华为云和个人实名认证](#)。
- 您需要确保账户有足够的资金，以免创建集群失败，具体操作请参见[账户充值](#)。

步骤一：首次开通 CCE 并进行授权

由于CCE在运行中对计算、存储、网络以及监控等各类云服务资源都存在依赖关系，因此当您首次登录CCE控制台时，CCE将自动请求获取当前区域下的云资源权限，从而更好地为您提供服务。如果您在当前区域已完成授权，可忽略本步骤。

步骤1 使用华为账号登录[CCE控制台](#)。

步骤2 单击管理控制台左上角的，选择区域。

步骤3 在首次登录某个区域的CCE控制台时将跳出“授权说明”，请您在仔细阅读后单击“确定”。

当您同意授权后，CCE将在IAM中创建名为“cce_admin_trust”委托，统一对您的其他云服务资源进行操作，并且授予其Tenant Administrator权限。Tenant Administrator拥有除IAM管理外的全部云服务管理员权限，用于对CCE所依赖的其他云服务资源进行调用，且该授权仅在当前区域生效。您可前往“统一身份认证服务 IAM > 委托”页签，单击“cce_admin_trust”查看各区域的授权记录。关于资源委托详情，您可参考[委托](#)进行了解。

说明

由于CCE对其他云服务有许多依赖，如果没有Tenant Administrator权限，可能会因为某个服务权限不足而影响CCE功能的正常使用。因此在使用CCE服务期间，请不要自行删除或者修改“cce_admin_trust”委托。

----结束

步骤二：创建 CCE Autopilot 集群

在CCE服务中创建CCE Autopilot集群，从而简化Kubernetes集群的管理和操作。

步骤1 返回[CCE控制台](#)。

- 如果您的帐号在当前区域未创建过集群，请在当前页面单击“购买集群”或“购买CCE Autopilot集群”，进入购买页。

- 如果您的帐号在当前区域已创建过集群，请在左侧菜单栏选择集群管理，单击右上角“购买集群”，进入购买页。

步骤2 配置集群基础信息。

本示例中仅解释必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[购买Autopilot集群](#)。



参数	示例	参数说明
集群类型	CCE Autopilot 集群	<p>CCE支持创建多种类型集群，满足各种业务需求，提供高可靠、安全的商业级容器集群服务。</p> <ul style="list-style-type: none"> CCE Standard集群：标准版本集群，提供高可靠、安全的商业级容器集群服务。 CCE Turbo集群：拥有更高性能的云原生网络，提供云原生混部调度能力，可实现更高的资源利用率和更广的全场景覆盖。 CCE Autopilot集群：Serverless版集群，提供免运维的容器服务，可以大幅降低运维成本，提高应用程序的可靠性和可扩展性。 <p>了解集群类型的更多内容，请参见集群对比。</p>
集群名称	autopilot-example	新建集群的名称。
企业项目	default	<p>该参数仅对开通企业项目的企业客户账号显示，不显示时请忽略。</p> <p>企业项目是一种资源管理单位，可跨区域归类资源，方便企业按部门或项目组集中管理。了解企业项目的更多内容，请参见项目管理。</p> <p>请根据需要选择适合的企业项目，如果没有特殊要求，可以选择default。</p>
集群版本	v1.28	集群安装的Kubernetes软件版本，建议选择最新版本。

步骤3 配置集群网络信息。

本示例中仅解释必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[购买Autopilot集群](#)。



参数	示例	参数说明
虚拟私有云	vpc-autopilot	选择集群所在的虚拟私有云VPC。如果没有可选项，单击右侧“新建虚拟私有云”创建，具体请参见 创建虚拟私有云和子网 。集群创建后不可修改。
容器子网	subnet-502f	选择容器所在子网。如果没有可选项，单击右侧“新建子网”创建，具体请参见 创建虚拟私有云和子网 。容器子网决定了集群下容器的数量上限，集群创建后支持新增子网。
服务网段	10.247.0.0/16	同一集群下容器互相访问时使用的Service资源网段，决定Service资源的上限。集群创建后不可修改。
镜像访问	-	为确保您的集群节点可以从容器镜像服务中拉取镜像，默认使用所选的VPC中已有的终端节点，否则系统将会为您新建SWR和OBS的终端节点。 终端节点将产生一定费用，详情请参见 价格计算器（VPC终端节点） 。
配置SNAT	开启	默认开启，开启后您的集群可以通过NAT网关访问公网。默认使用所选的VPC中已有的NAT网关，否则系统将会为您自动创建一个默认规格的NAT网关并绑定弹性公网IP，自动配置SNAT规则。 使用NAT网关将产生一定费用，详情请参见 价格计算器（NAT网关） 。

步骤4 单击“下一步：插件选择”，选择创建集群时需要安装的插件。

本示例中，选择默认安装的必选插件。

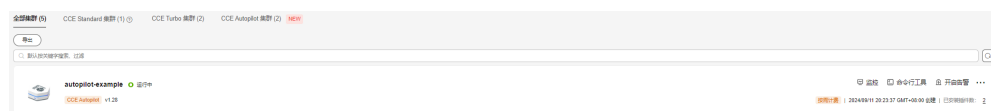
步骤5 单击“下一步：插件配置”，对选择的插件进行配置，其中默认插件无法进行配置。

本示例中，默认插件无需配置。

步骤6 单击“下一步：确认配置”，显示集群资源清单，确认无误后，单击“提交”。

创建集群预计需要5-10分钟左右，请耐心等待。

创建成功后，集群管理中对对应集群的状态为运行中。



----结束

步骤三：创建算力密集型应用并上传 SWR

创建算力密集型应用主要用于压力测试，以验证部署的HPA策略是否能够根据CPU和内存资源的使用情况自动调整Pod实例的数量。创建完成的应用需要上传至容器镜像服务（SWR），以便在集群中进行部署和管理。该步骤涉及命令行操作，您需要准备一台Linux系统的ECS虚拟机，该ECS与集群处于同一VPC，并绑定弹性公网IP，具体操作请参见[快速购买和使用Linux ECS](#)。本示例以“CentOS 7.9 64bit(40GiB)”操作系统为例，为您演示如何基于PHP创建算力密集型应用并上传SWR的具体操作。

步骤1 登录ECS虚拟机，具体操作请参见[通过CloudShell登录Linux ECS](#)。

步骤2 安装Docker。

1. 检查ECS是否安装Docker，如果已安装可以跳过该步骤。

```
docker --version
```

回显如下，则说明未安装Docker。

```
-bash: docker: command not found
```

2. 安装并运行Docker。

```
yum install docker
```

设置Docker在系统启动时自动启动，并立即开始运行。

```
systemctl enable docker
```

```
systemctl start docker
```

3. 检查安装结果。

```
docker --version
```

回显如下，则说明Docker安装成功。

```
Docker version 1.13.1, build 7d71120/1.13.1
```

步骤3 创建算力密集型应用，并制作镜像。

1. 创建一个名为index.php的PHP文件，文件名称可自定义。该文件在接收到用户请求时，先进行1000000次平方根循环计算，即“0.0001+0.01+0.1+...”，然后返回“OK!”。

```
vim index.php
```

文件内容如下：

```
<?php
$x = 0.0001;
for ($i = 0; $i <= 1000000; $i++)
    $x += sqrt($x);
}
echo "OK!";
?>
```

输入完成后，**Esc**键退出编辑，输入:wq保存。

2. 编写Dockerfile制作镜像。

```
vim Dockerfile
```

文件内容如下：

```
# 使用PHP的Docker官方镜像
FROM php:5-apache
# 将本地的index.php文件复制到容器指定目录，这个文件将被用作默认的网页
COPY index.php /var/www/html/index.php
# 修改index.php文件的权限，使其对所有用户可读和可执行，确保文件在Web服务器上能够被正确访问
RUN chmod a+rx index.php
```

3. 执行如下命令构建镜像，镜像名称为hpa-example，版本为latest，名称和版本可自定义。

```
docker build -t hpa-example:latest .
```

步骤4 将创建的hpa-example镜像上传至SWR。

1. 创建镜像组织。若已有组织，该步骤可以跳过。

登录[SWR控制台](#)，右侧单击“组织管理”，单击页面右上角“创建组织”。在创建组织界面输入“组织名称”，单击“确定”。

创建组织

① 规则

1. 组织名称，全局唯一。
2. 当前租户最多可创建5个组织。
3. 推荐您创建的每一个组织对应一个公司、下属部门或者个人用户。

示例

1. 以公司、部门作为组织：cloud-hangzhou、cloud-develop。
2. 以个人作为组织：john。

组织名称

2. ECS登录SWR。

在左侧导航栏选择“我的镜像”，单击右侧“客户端上传”。在弹出的页面中单击“生成登录指令”，单击“复制图标”，复制“临时登录指令”。

在ECS执行复制的登录指令，如下所示。

```
docker login -u cn-east-3xxx swr.cn-east-3.myhuaweicloud.com
```

回显如下，说明登录成功。

```
Login Succeeded
```

3. 为hpa-example镜像添加标签，代码结构如下：

```
docker tag {镜像名称1:版本名称1} {镜像仓库地址}/{组织名称}/{镜像名称2:版本名称2}
```

具体示例：

```
docker tag hpa-example:latest swr.cn-east-3.myhuaweicloud.com/test/hpa-example:latest
```

参数	示例	参数说明
{镜像名称1:版本名称1}	hpa-example:latest	请替换为本地所要上传的实际镜像的名称和版本名称。
{镜像仓库地址}	swr.cn-east-3.myhuaweicloud.com	请替换为 登录指令 中的末尾域名，该域名即为镜像仓库地址。
{组织名称}	test	请替换为 已创建的镜像组织 。
{镜像名称2:版本名称2}	hpa-example:latest	请替换为SWR镜像仓库中需要显示的镜像名称和镜像版本，您可以自定义。

4. 上传镜像至镜像仓库，代码结构如下：

```
docker push {镜像仓库地址}/{组织名称}/{镜像名称2:版本名称2}
```

具体示例：

```
docker push swr.cn-east-3.myhuaweicloud.com/test/hpa-example:latest
```

回显结果如下，则说明上传镜像成功。

```
6d6b9812c8ae: Pushed
...
fe4c16cbf7a4: Pushed
latest: digest: sha256:eb7e3bbdxxx size: xxx
```

返回[SWR控制台](#)，在“我的镜像”页面，执行刷新操作后可查看到对应的镜像信息。

----结束

步骤四：使用 hpa-example 镜像创建工作负载

本示例将使用构建的hpa-example镜像创建工作负载，并为该工作负载创建负载均衡类型服务，使您能够从公网访问应用。本节介绍两种方式创建该工作负载，即控制台方式和kubectl命令行方式。

使用控制台方式

步骤1 返回[CCE控制台](#)。单击新建的集群名称，进入集群控制台。

步骤2 左侧菜单栏选择“工作负载”，单击右上角“创建工作负载”，进入创建页。

步骤3 配置工作负载基本信息。

本示例中仅解释必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[创建工作负载](#)，您可以根据工作负载类型选择适合的参考文档。



参数	示例	参数说明
负载类型	无状态负载 Deployment	工作负载是一种对Pod的抽象管理方式，用于定义和控制Pod的创建、运行和生命周期。通过工作负载，您可以批量管理和自动化控制多个Pod的行为，如伸缩、更新和恢复。 <ul style="list-style-type: none">无状态负载（Deployment）：管理无状态应用，支持上线部署、滚动升级、创建副本和恢复上线。有状态负载（StatefulSet）：管理有状态应用，确保每个Pod能够拥有独立的状态。普通任务（Job）：一次性任务，完成后Pod自动删除。定时任务（CronJob）：基于时间的Job，指定时间周期内运行指定的Job。 了解工作负载的更多内容，请参见 工作负载概述 。
负载名称	hpa-example	请填写工作负载的名称。

参数	示例	参数说明
命名空间	default	命名空间是Kubernetes集群中的抽象概念，可以将集群中的资源或对象划分为一个组，且不同命名空间中的数据彼此隔离，您可以根据需要创建并使用命名空间。 集群创建后会默认生成default命名空间，如果没有特殊要求，可以直接选择default命名空间。
实例数量	1	工作负载中Pod的数量。

步骤4 配置工作负载容器信息。

本示例中仅解释必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[创建工作负载](#)，您可以根据工作负载类型选择适合的参考文档。



参数	示例	参数说明
镜像名称	hpa-example	单击“选择镜像”，在弹出的窗口中切换至“我的镜像”，选择上传的hpa-example镜像。
镜像版本	latest	选择需要部署的镜像版本。
CPU配额	0.25cores	CPU资源限制值，即CPU资源限制值，即允许容器使用的CPU最大值，防止占用过多资源，默认0.25cores。
内存配额	512MiB	内存资源限制值，即允许容器使用的内存最大值。如果超过，容器会被终止，默认512MiB。

步骤5 配置工作负载服务信息。

单击“服务配置”下的加号，进入创建服务页面。

本示例中需要外部访问hpa-example，所以访问类型设置为负载均衡。

本示例仅解释负载均衡的必要参数，其他参数保留默认值。关于配置参数的详细说明请参见[服务 \(Service\)](#)，您可以根据服务类型选择适合的参考文档。

创建服务
✕

Service名称

访问类型

集群内访问

通过集群的内部IP暴露服务，只能够在集群内部访问

负载均衡

通过ELB负载均衡对外部提供服务，高可用，超高性能，稳定安全

1 集群外访问推荐选择负载均衡访问类型

服务亲和 集群级别 ?

负载均衡器

独享型
网络型 (TCP/UDP) & ...
选择...
elb-hpa-example
创建负载均衡器

仅支持集群所在 VPC vpc-autopilot 下、实例规格支持网络型 & 应用型的独享型负载均衡实例。 [查看约束与限制](#)

负载均衡配置：分配策略：加权轮询算法；会话保持：不启用； [编辑](#)

☑ 我已阅读 [《负载均衡使用须知》](#)

健康检查

不启用
全局检查
自定义检查

协议：TCP | 检查周期（秒）：5 | 超时时间（秒）：10 | 最大重试次数：3 [🔗](#)

端口配置

协议	容器端口	服务端口	监听器前端协议	操作
TCP	80	80	TCP	删除
+				

监听器配置

访问控制 未配置

高级配置 [+ 添加高级配置](#)

注解

=
确认添加
使用指南

参数	示例	参数说明
Service名称	hpa-example	请填写服务的名称。
访问类型	负载均衡	<p>选择服务类型，即服务访问的方式。</p> <ul style="list-style-type: none"> 集群内访问：通过集群的内部IP暴露服务，只能够在集群内部访问。 负载均衡：通过弹性负载均衡（ELB）对外部提供服务，即能够从公网访问到工作负载。 <p>了解服务类型的更多内容，请参见服务（Service）。</p>
负载均衡器	独享型 网络型(TCP/UDP)&应用型(HTTP/HTTPS) 选择已有 elb-hpa-example	<ul style="list-style-type: none"> 如果已有负载均衡（ELB）实例，可以选择已有ELB。 <p>说明</p> <p>使用已有的ELB时，仅支持集群所在VPC下、实例规格支持网络型的独享型ELB实例。</p> <ul style="list-style-type: none"> 如果没有请选择“自动创建”，创建一个负载均衡器，并绑定弹性公网IP，具体操作请参见创建负载均衡类型的服务。

参数	示例	参数说明
端口配置	协议: TCP	负载均衡监听器端口协议。
	容器端口: 80	容器中应用启动监听的端口, 该容器端口需和应用对外提供的监听端口一致。
	服务端点: 80	ELB将会使用该端口创建监听器, 提供外部流量访问入口, 可自定义。

步骤6 单击右下角“创建工作负载”。

创建成功后, 无状态工作负载列表中对工作负载的状态为运行中。



步骤7 获取hpa-example的外部访问地址。

单击hpa-example工作负载名称, 进入工作负载详情页。在“访问方式”页签下可以看到hpa-example的IP地址, 其中“负载均衡公网IP: 访问端口”就是外部访问地址。



---结束

使用 kubectl 命令行方式

本节继续使用**步骤三: 创建算力密集型应用并上传SWR**中的ECS。此外, 您还需要安装kubectl命令, 并**通过kubectl连接集群**。

步骤1 安装kubectl命令行工具。

您可以尝试执行**kubectl version**命令判断是否已安装kubectl, 如果已经安装kubectl, 则可跳过此步骤。

本文以Linux环境为例安装和配置kubectl, 更多安装方式请参考**安装kubectl**。

1. 登录ECS虚拟机, 下载kubectl。

```
cd /home
curl -LO https://dl.k8s.io/release/{v1.28.0}/bin/linux/amd64/kubectl
```

其中{v1.28.0}为指定的版本号, 请根据集群版本进行替换。

2. 安装kubectl。

```
chmod +x kubectl
mv -f kubectl /usr/local/bin
```

步骤2 为kubectl命令行工具配置访问Kubernetes集群的凭证。

1. 返回**CCE控制台**, 并单击集群名称进入集群总览页。

2. 在集群总览页中找到“连接信息”版块。单击kubectl后的“配置”，查看kubectl的连接信息。
3. 在弹出页面中选择“内网访问”，然后下载对应的配置文件。
4. 登录已安装kubectl客户端的虚拟机，将上一步中下载的配置文​​件（以 kubeconfig.yaml为例）复制到/home目录下。

5. 将kubectl认证文件保持至\$HOME/.kube目录下的config文件中。

```
cd /home
mkdir -p $HOME/.kube
mv -f kubeconfig.yaml $HOME/.kube/config
```

6. 执行kubectl命令验证集群的连通性。
以查看集群信息为例，执行以下命令。

```
kubectl cluster-info
```

回显内容如下，则说明集群连接成功。

```
Kubernetes control plane is running at https://xx.xx.xx.xx:5443
CoreDNS is running at https://xx.xx.xx.xx:5443/api/v1/namespaces/kube-system/services/coredns.dns/proxy
To further debug and diagnose cluster problems, use 'kubectl cluster-info dump'.
```

步骤3 创建YAML文件hpa-example.yaml，用于配置hpa-example工作负载，文件名称可自定义。

```
vim hpa-example.yaml
```

文件内容如下：

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: hpa-example # 工作负载名称
spec:
  replicas: 1 # 实例数量
  selector:
    matchLabels: # 选择器，用于选择带有特定标签的资源
      app: hpa-example
  template:
    metadata:
      labels: # 标签
        app: hpa-example
    spec:
      containers:
        - name: container-1
          image: 'swr.cn-east-3.myhuaweicloud.com/test/hpa-example:latest' # 替换为您上传到SWR的镜像地址
      resources:
        limits: # limits与requests建议取值保持一致，避免扩缩容过程中出现震荡
          cpu: 250m
          memory: 512Mi
        requests:
          cpu: 250m
          memory: 512Mi
      imagePullSecrets:
        - name: default-secret
```

📖 说明

“swr.cn-east-3.myhuaweicloud.com/test/hpa-example:latest”需替换为您上传到SWR的镜像地址，镜像地址即为步骤[上传镜像至镜像仓库](#)中docker push后内容。

输入完成后，**Esc**键退出编辑，输入:wq保存。

步骤4 执行以下命令，创建工作负载。

```
kubectl create -f hpa-example.yaml
```

回显如下，表示已经开始创建工作负载。

```
deployment.apps/hpa-example created
```

步骤5 执行以下命令，查看工作负载状态。

```
kubectl get deployment
```

回显如下，如果READY值为1/1，则说明工作负载创建的Pod皆为可用状态，即创建成功。

```
NAME          READY  UP-TO-DATE  AVAILABLE  AGE
hpa-example   1/1    1            1           4m59s
```

回显内容的参数解析如下：

参数	示例	参数说明
NAME	hpa-example	工作负载的名称。
READY	1/1	表示工作负载的可用状态，显示为“可用Pod个数/期望Pod个数”。
UP-TO-DATE	1	指当前工作负载已经完成更新的Pod数。
AVAILABLE	1	工作负载可用的Pod个数。
AGE	4m59s	工作负载已经运行的时间。

步骤6 创建负载均衡类型的服务，并关联已创建的工作负载hpa-example。

本示例基于已有的弹性负载均衡（ELB）实例创建服务，如果您需要自动创建ELB请参考[通过kubectl命令行创建-自动创建ELB](#)。

创建YAML文件nginx-elb-svc.yaml，用于配置负载均衡服务，文件名称可自定义。

```
vim hpa-example-elb-svc.yaml
```

文件内容如下：

```
apiVersion: v1
kind: Service
metadata:
  name: hpa-example # 服务的名称
  annotations:
    kubernetes.io/elb.id: <your_elb_id> # ELB ID, 替换为实际值
    kubernetes.io/elb.class: performance # 负载均衡器类型
spec:
  selector:
    app: hpa-example
  ports:
  - name: service0
    port: 80
    protocol: TCP
    targetPort: 80
  type: LoadBalancer
```

输入完成后，**Esc**键退出编辑，输入:**wq**保存。

参数	示例	参数说明
kubernetes.io/elb.id	405ef586-0397-45c3-bfc4-xxx	已有的ELB的ID。 获取方式： 进入 网络控制台 ，选择“弹性负载均衡 > 我的ELB”，找到对应的ELB名称，名称下方即为对应ID。
kubernetes.io/elb.class	performance	负载均衡器类型，仅支持performance类型，即独享型负载均衡。 说明 负载均衡类型的服务对接已有的独享型ELB时，该独享型ELB必须支持网络型（TCP/UDP）规格。
selector	app: hpa-example	选择器，服务将流量发送给对应标签的Pod。 须知 selector取值需与工作负载YAML文件中matchLabels参数取值一致，本示例中为app: hpa-example。
ports.port	8080	弹性负载均衡（ELB）将会使用该端口创建监听器，提供外部流量访问入口，可自定义。
ports.protocol	TCP	负载均衡监听器端口协议。
ports.targetPort	80	Service访问目标容器的端口，此端口与容器中运行的应用强相关。

步骤7 执行以下命令创建服务。

```
kubectl create -f hpa-example-elb-svc.yaml
```

回显如下，表示服务已创建。

```
service/hpa-example created
```

步骤8 执行以下命令查看服务。

```
kubectl get svc
```

回显内容如下，表示工作负载访问方式已设置成功。您可以通过“外部访问地址：服务端口”访问该工作负载，其中外部访问地址为EXTERNAL-IP对应的第一个IP地址的，服务端口为8080。

NAME	TYPE	CLUSTER-IP	EXTERNAL-IP	PORT(S)	AGE
kubernetes	ClusterIP	10.247.0.1	<none>	443/TCP	18h
hpa-example	LoadBalancer	10.247.56.18	xx.xx.xx.xx,xx.xx.xx.xx	8080:30581/TCP	5m8s

----结束

步骤五：创建 HPA 策略

HPA策略（Horizontal Pod Autoscaler）主要用来控制Pod的水平伸缩，通过周期性检查Pod的度量数据，计算实现HPA策略中资源目标值所需的副本数量，进而调整目标资源的replicas字段（工作负载创建的Pod实例数量）。了解HPA策略的更多信息，请参见[HPA工作原理](#)。本节介绍两种方式创建HPA策略，即控制台方式和kubectl命令行方式。

使用控制台方式

本示例中仅解释必要参数，其他参数保留默认值。关于默认参数的详细说明，请参见 [Pod水平自动扩缩](#)。

步骤1 返回集群总览页，左侧菜单栏选择“策略”，单击右上角“创建HPA策略”，进入创建页。

步骤2 填写策略基础信息。

策略名称	<input type="text" value="hpa-example"/>
集群名称	autopilot-example
命名空间	<input type="text" value="default"/> 创建命名空间
关联工作负载	<input type="text" value="hpa-example"/> 创建工作负载

参数	示例	参数说明
策略名称	hpa-example	请填写策略的名称。
命名空间	default	命名空间是Kubernetes集群中的抽象概念，可以将集群中的资源或对象划分为一个组，且不同命名空间中的数据彼此隔离，您可以根据需要创建并使用命名空间。 集群创建后会默认生成default命名空间，如果没有特殊要求，可以直接选择default命名空间。
关联工作负载	hpa-example	请选择 步骤四：使用hpa-example镜像创建工作负载 中创建的工作负载。 说明 弹性伸缩策略会根据名称匹配关联工作负载，若相同命名空间下的多个工作负载中app的label相同，会导致弹性伸缩策略不符合预期。

步骤3 进行“策略配置”。

策略配置

实例范围: - 策略生效时，工作负载实例将在此范围内伸缩

策略配置: [策略默认](#) [自定义](#)

选择基础策略以采用 K8s 社区策略的默认行为进行负载伸缩。选择自定义策略用户可以自定义策略窗口、中点、收缩策略等实现更灵活的策略。未配置的策略将采用社区策略的默认值。社区默认行为说明

系统策略:

指标	策略值	策略范围	操作
CPU使用率	<input type="text" value="50"/> %	<input type="text" value="30"/> % - <input type="text" value="70"/> %	删除
+			

自定义策略:

自定义策略名称	策略策略	策略值	策略范围	操作
+				

参数	示例	参数说明
实例范围	1-10	表示工作负载创建的Pod实例的伸缩范围。
伸缩配置	系统默认	用于选择伸缩配置策略，有以下两种配置策略： <ul style="list-style-type: none">系统默认：稳定窗口和步长等策略直接采用K8S社区推荐的默认行为，更多信息请参见Pod水平自动扩缩。自定义：用户可以自定义稳定窗口和步长等策略实现更灵活的配置，未配置参数将采用社区推荐的默认值，更多信息请参见Pod水平自动扩缩。
系统策略	指标：cpu利用率	HPA可以通过跟踪相关指标的资源用量，触发Pod实例的扩缩操作。 <ul style="list-style-type: none">CPU利用率：工作负载容器组（Pod）的实际使用量/申请量。内存利用率：工作负载容器组（Pod）的实际使用量/申请量。
	期望值：50%	表示所选指标的期望值，可以用来计算需要伸缩的实例数，即“需要伸缩的实例数=（当前指标值/期望值）×当前实例数”。
	容忍范围：30%-70%	表示伸缩的触发范围，当指标处于范围内时不会触发伸缩，期望值必须在容忍范围内。

步骤4 单击“创建”。您可以在“HPA策略”列表中看到已创建的策略，最新状态为“已启动”。



----结束

使用 kubectl 命令行方式

步骤1 创建YAML文件hpa-policy.yaml，用于配置HPA策略，文件名称可自定义。

```
vim hpa-policy.yaml
```

文件内容如下：

```
apiVersion: autoscaling/v2
kind: HorizontalPodAutoscaler
metadata:
  name: hpa-example # HPA策略的名称
  namespace: default # 命名空间
annotations:
```

```
# 定义伸缩的触发条件：CPU利用率在30%到70%之间时，不会进行伸缩操作，反之，则会进行伸缩操作。
extendedhpa.metrics: '[{"type":"Resource","name":"cpu","targetType":"Utilization","targetRange":
{"low":"30","high":"70"}}]'
spec:
  scaleTargetRef: # 关联的工作负载信息
    kind: Deployment
    name: hpa-example
    apiVersion: apps/v1
  minReplicas: 1
  maxReplicas: 10
  metrics:
  - type: Resource # 设置监控的资源
    resource:
      name: cpu # 设置资源为cpu，还可以设置为memory
      target:
        type: Utilization # 设置指标为利用率
        averageUtilization: 50 # HPA控制器会维持伸缩目标中的Pod的平均资源利用率在50%
```

输入完成后，**Esc**键退出编辑，输入:wq保存。

步骤2 执行以下命令创建HPA策略。

```
kubectl create -f hpa-policy.yaml
```

回显如下，表示策略已创建。

```
horizontalpodautoscaler.autoscaling/hpa-example created
```

步骤3 执行以下命令查看创建的HPA策略。

```
kubectl get hpa
```

回显结果如下：

NAME	REFERENCE	TARGETS	MINPODS	MAXPODS	REPLICAS	AGE
hpa-example	Deployment/hpa-example	<unknown>/50%	1	10	0	10s

----结束

步骤六：检查工作负载能否自动进行弹性伸缩

本节介绍两种方式检查工作负载hpa-example能否自动进行弹性伸缩，即控制台方式和kubectl命令行方式。

使用控制台方式

步骤1 登录**步骤三：创建算力密集型应用并上传SWR**中使用的ECS虚拟机。


步骤2 执行以下命令循环访问工作负载，其中“ip:port”为工作负载的访问地址，与**步骤7**中的“负载均衡公网IP：访问端口”一致。

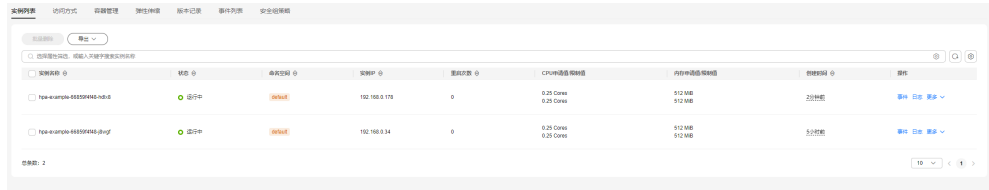
```
while true;do wget -q -O- http://ip:port; done
```

回显结果如下：

```
OK!
OK!
...
```

步骤3 观察工作负载自动扩容过程。

返回**CCE控制台**。单击新建的集群名称，进入集群控制台。左侧导航栏单击“工作负载”，单击“工作负载名称”，单击搜索栏右侧，可以发现工作负载hpa-example创建的Pod实例个数由1个扩容至2个。



步骤4 观察工作负载自动缩容过程。

返回ECS，利用Ctrl+c停止访问工作负载。

返回控制台界面，单击搜索栏右侧 。值得注意的是，工作负载缩容的稳定窗口默认为300秒，即缩容策略成功触发后，不会立即对目标进行缩容，需要先冷却300秒，以防止短期波动造成影响。



----结束

使用 kubectl 命令行方式

步骤1 执行以下命令循环访问工作负载，其中“ip:port”为工作负载的访问地址，与**步骤8**中的“外部访问地址：服务端口”相对应。

```
while true;do wget -q -O- http://ip:port; done
```

回显结果如下：

```
OK!
OK!
...
```

步骤2 新建一个终端，执行以下命令观察工作负载自动扩容过程。

```
kubectl get hpa hpa-policy --watch
```

回显结果如下：

NAME	REFERENCE	TARGETS	MINPODS	MAXPODS	REPLICAS	AGE
hpa-example	Deployment/hpa-example	24%/50%	1	10	1	17m
hpa-example	Deployment/hpa-example	100%/50%	1	10	1	17m
hpa-example	Deployment/hpa-example	100%/50%	1	10	2	18m
hpa-example	Deployment/hpa-example	57%/50%	1	10	2	19m
hpa-example	Deployment/hpa-example	57%/50%	1	10	2	20m

参数	示例	说明
NAME	hpa-example	HPA策略的名称。 本示例中HPA策略名称为hpa-example。
REFERENCE	Deployment/hpa-example	HPA策略管理的对象。 本示例中为无状态工作负载hpa-example。

参数	示例	说明
TARGETS	24%/50%	HPA策略监控指标的当前值和期望值。 本示例中监控指标为CPU利用率，24%表示CPU利用率的当前值，50%表示CPU利用率的期望值。
MINPODS	1	HPA策略允许的最小Pod实例数量。 本示例中允许的最小Pod实例数量为1，当Pod实例数量为1时，即使当前CPU利用率低于容忍范围，也不会进行缩容。
MAXPODS	10	HPA策略允许的最大Pod实例数量。 本示例中允许的最大Pod实例数量为10，当Pod实例数量为10时，即使当前CPU利用率高于容忍范围，也不会进行扩容。
REPLICAS	1	当前实际运行的Pod实例数量。
AGE	17m	HPA策略已存在的时间。

第2条记录中，工作负载的CPU利用率为100%，超过70%，触发弹性伸缩。第3条记录中，工作负载将Pod实例数量由1扩容至2。第5条记录中，工作负载的CPU利用率降至57%，属于30%~70%，并未触发弹性伸缩。

步骤3 观察负载自动缩容过程。

在**步骤1**所在的终端中利用Ctrl+c停止访问工作负载。

在**步骤2**所在的终端中，继续观察回显结果。值得注意的是，工作负载缩容的稳定窗口默认为300秒，即缩容策略成功触发后，不会立即对目标进行缩容，需要先冷却300秒，以防止短期波动造成影响。

```
NAME          REFERENCE          TARGETS  MINPODS  MAXPODS  REPLICAS  AGE
...
hpa-example   Deployment/hpa-example  19%/50%  1        10       2         30m
hpa-example   Deployment/hpa-example  0%/50%   1        10       2         31m
hpa-example   Deployment/hpa-example  0%/50%   1        10       2         35m
hpa-example   Deployment/hpa-example  0%/50%   1        10       1         36m
...
```

第1条记录中，工作负载的CPU利用率为19%，低于30%，触发弹性伸缩。但此时处于冷却时间，因此工作负载并未立即缩容。第4条记录中，工作负载将Pod实例数量由2缩容至1，达到最低实例数量，不再触发弹性伸缩。

----结束

后续操作：释放资源

如果您无需继续使用集群和ECS，请及时释放资源，避免产生额外的费用。

步骤1 返回[CCE控制台](#)，在左侧导航栏中选择“集群管理”。

步骤2 找到需要删除的集群，单击集群卡片右上角的“...”，并单击“删除集群”。

步骤3 在弹出的“删除集群”窗口中，勾选需要删除的所有资源，并根据页面提示删除NAT网关、SNAT出网EIP和终端节点等资源。



在确认框中输入“DELETE”，单击“是”，开始执行删除集群操作。

删除集群需要花费1~3分钟，请耐心等待。

步骤4 进入云服务器控制台，左侧导航栏单击“弹性云服务器”，找到对应的ECS，右侧单击“更多”，单击“删除”。

在删除界面中勾选“删除云服务器绑定的弹性公网IP地址”和“删除云服务器挂载的数据盘”，单击“下一步”。



在确认框中输入“DELETE”，单击“确定”，开始执行删除ECS操作。
删除ECS需要花费0.5~1分钟，请耐心等待。



---结束