

弹性伸缩服务

产品介绍

文档版本 06
发布日期 2022-11-15



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。

目录

1 图解弹性伸缩	1
2 什么是弹性伸缩?	3
3 弹性伸缩的优势	5
4 生命周期	9
5 使用限制	13
6 区域和可用区	15
7 计费标准	17
8 安全	18
8.1 责任共担.....	18
8.2 数据保护技术.....	19
8.3 审计与日志.....	19
8.4 监控安全风险.....	19
8.5 认证证书.....	20
9 与其他服务的关系	22
10 权限管理	25
11 基本概念	28
12 修订记录	30

1 图解弹性伸缩



初识弹性伸缩

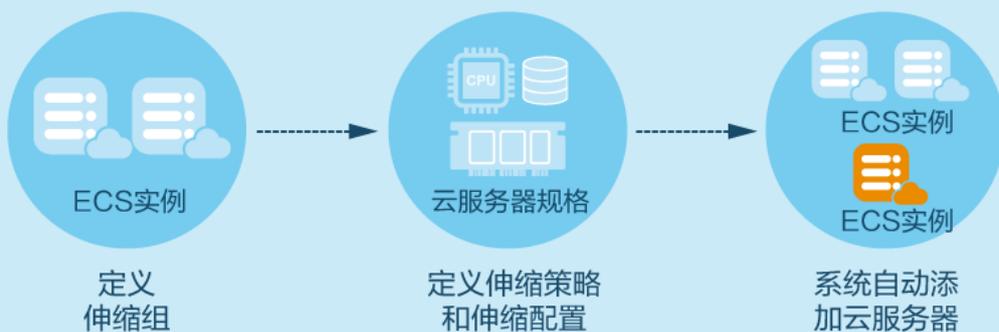


01

什么是弹性伸缩服务

弹性伸缩 (Auto Scaling) 是根据用户的业务需求, 通过策略自动调整其业务资源的服务, 简称AS。

您可以根据业务需求自行定义伸缩配置和伸缩策略, 降低应对业务变化和高峰压力时人为反复调整资源的工作量, 帮助您节约资源和人力成本。



每年年初、年终购买火车票时, 应用系统业务需求无法满足, 怎么办?

向应用系统中添加更多的服务器-----造成资源浪费

平均需求所需容量-----无法满足高峰期的业务需求

02

产品优势有哪些

当您向应用程序中添加弹性伸缩后可解决上述问题, 同时还可获得其他优势。

2 什么是弹性伸缩?

弹性伸缩简介

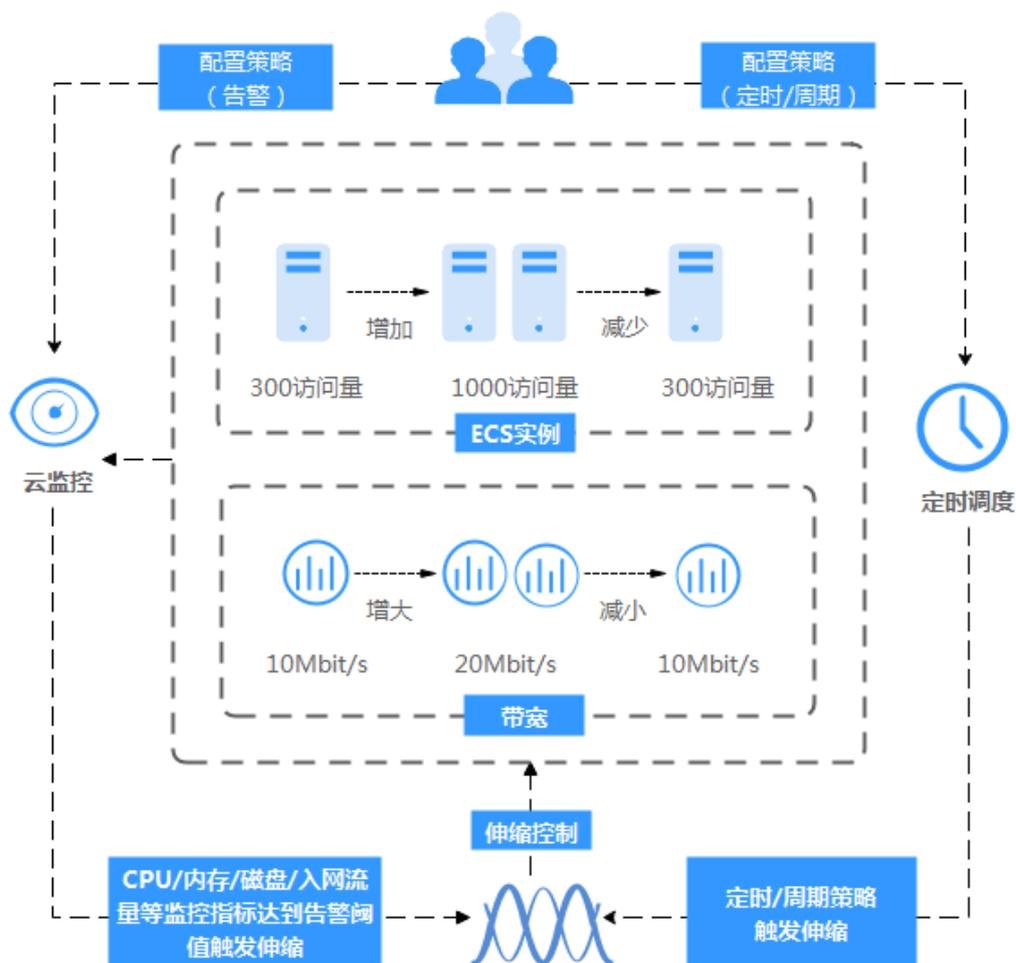
弹性伸缩（Auto Scaling，以下简称AS）是根据用户的业务需求，通过设置伸缩规则来自动增加/缩减业务资源。当业务需求增长时，AS自动为您增加弹性云服务器（ECS）实例或带宽资源，以保证业务能力；当业务需求下降时，AS自动为您缩减弹性云服务器（ECS）实例或带宽资源，以节约成本。AS支持自动调整弹性云服务器和带宽资源。

产品架构

通过伸缩控制可以实现弹性云服务器（ECS）实例伸缩和带宽伸缩：

- 伸缩控制：配置策略设置指标阈值/伸缩活动执行的时间，通过云监控监控指标是否达到阈值，通过定时调度，实现伸缩控制。
- 配置策略：可以根据业务需求，配置告警策略/定时策略/周期策略。
- 配置告警策略：可配置CPU、内存、磁盘、入网流量等监控指标。
- 配置定时策略：通过配置触发时间可以配置定时策略。
- 配置周期策略：通过配置重复周期、触发时间、生效时间可以配置周期策略。
- 云监控监控到所配置的告警策略中的某些指标达到告警阈值，从而触发伸缩活动，实现ECS实例的增加/减少或带宽的增大/减小。
- 到达所配置的触发时间时，触发伸缩活动，实现ECS实例的增加/减少或带宽的增大/减小。

图 2-1 弹性伸缩产品架构图



访问方式

公有云提供了Web化的服务管理平台，即管理控制台和基于HTTPS请求的API（Application programming interface）管理方式。

- API方式
如果用户需要将公有云平台上的弹性伸缩服务集成到第三方系统，用于二次开发，请使用API方式访问弹性伸缩服务，具体操作请参见《[弹性伸缩API参考](#)》。
- 控制台方式
其他相关操作，请使用管理控制台方式访问弹性伸缩服务。
如果用户已注册公有云，可直接登录管理控制台，从主页选择“弹性伸缩”。

3 弹性伸缩的优势

弹性伸缩服务可根据用户的业务需求，通过策略自动调整其业务的资源。具有自动调整资源、节约成本开支、提高可用性和容错能力的优势。适用以下场景：

- 访问流量较大的论坛网站，业务负载变化难以预测，需要根据实时监控到的云服务器CPU使用率、内存使用率等指标对云服务器数量进行动态调整。
- 电商网站，在进行大型促销活动时，需要定时增加云服务器数量和带宽大小，以保证促销活动顺利进行。
- 视频直播网站，每天14:00~16:00播出热门节目，每天都需要在该时段增加云服务器数量，增大带宽大小，保证业务的平稳运行。

自动调整资源

弹性伸缩能够实现应用系统自动按需调整资源，即在业务增长时能够实现自动增加实例数量和带宽大小，以满足业务需求，业务下降时能够实现应用系统自动扩容，保障业务平稳运行。

- 按需调整云服务器资源

向应用系统中添加弹性伸缩，能够实现按需调整资源，即能够在业务增长时增加实例，业务下降时减少实例，这样加强了应用系统的成本管理。调整资源主要包括以下几种方式：

- 动态调整资源

动态调整资源是通过告警策略的触发来调整资源。

- 计划调整资源

计划调整资源是通过定时策略或周期策略的触发来调整资源。

- 手工调整资源

通过修改期望实例数或手动移入、移出实例来调整资源。

例如，运行在公有云上的基本Web应用程序。此应用程序允许乘客购买火车票。在每年中期时段，人员流动性较低，此应用程序的使用率较低。每年年底和年初，人员流动性较高，因此对此应用程序的需求会显著提高。一般系统会采用添加足够多的服务器，如[图3-1](#)所示，或添加处理应用程序平均需求所需的容量，如[图3-2](#)所示，来满足业务需求。但这两种方案会造成资源浪费或无法满足高峰期的需求。当您给应用程序中添加弹性伸缩后，弹性伸缩会自动根据需求调整服务器的数量，如[图3-3](#)所示，为您节约成本并且满足高峰期的需求。

图 3-1 服务器资源冗余

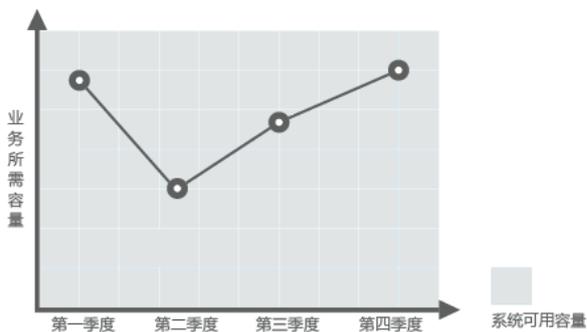


图 3-2 服务器资源不足

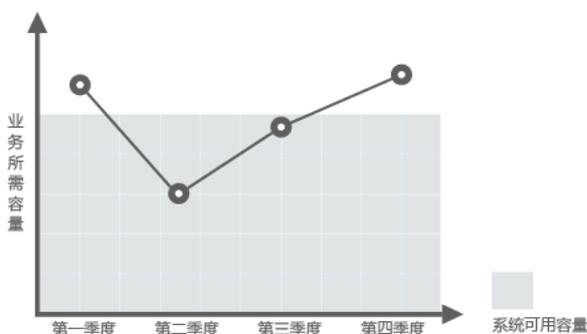
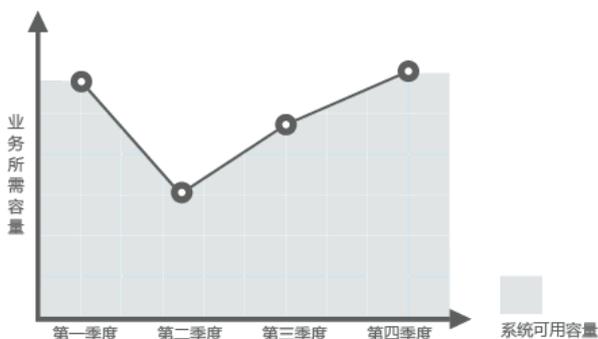


图 3-3 向应用程序中添加弹性伸缩



- 按需调整带宽资源

弹性伸缩能够实现按需调整带宽，即能够在业务增长时扩大带宽，业务下降时减小带宽，加强了应用系统的成本管理。

您可以根据实际情况选择如下伸缩带宽策略来实现按需调整IP带宽：

- 告警策略

可设置出网流量、出网带宽等告警触发条件，系统检测到触发条件满足时，会自动调整带宽的大小。

- 定时策略

系统可根据定时策略在固定的时间自动将带宽增大、减小或者调整到固定的值。

- 周期策略

系统可根据周期策略周期性的调整带宽大小，减少了人工重复设置带宽的工作量。

以告警策略的使用为例说明如下：

某视频直播网站，在不同时间段业务负载变化难以预测，需要根据出网流量、入网流量等指标在10Mbit/s到30Mbit/s之间动态调整带宽资源。弹性伸缩可以实现自动按需调整带宽，很好的解决这个问题。您只需选择需要调整的弹性公网IP，同时创建两个告警策略，一个策略设置在出网流量大于XXXbyte时，增加2Mbit/s，限制值为30Mbit/s；另一个策略在出网流量小于XXXbyte时，减少2Mbit/s，限制值为10Mbit/s。

● 按可用区均匀分配实例

按可用区均匀分配实例是指尽可能地将实例均匀的分布在不同的可用区中，来降低电力、网络等可能出现的故障对整个系统稳定性的影响。

区域指弹性云服务器云主机所在的物理位置。每个区域包含许多不同的称为“可用区”的位置，即在同一区域下，电力、网络隔离的物理区域，可用区之间内网互通，不同可用区之间物理隔离。每个可用区都被设计成不受其他可用区故障影响的模式，并提供低价、低延迟的网络连接，以连接到同一地区其他可用区。

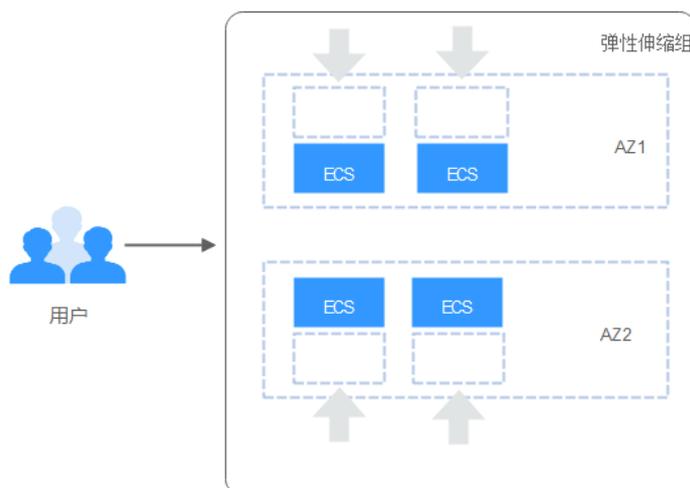
伸缩组可以包含来自同一区域的一个或多个可用区的实例。在资源调整时，弹性伸缩会通过实例分配和再均衡两种方法尽可能的将实例均匀分配到可用区中。

实例分配

弹性伸缩尝试在为伸缩组使用的可用区之间均匀分配实例。弹性伸缩通过尝试向实例最少的可用区中移入新实例来实现此目标。

例如，伸缩组目前有四个实例均匀分布在两个可用区内，若该伸缩组下一个伸缩活动增加四个实例时，会在两个可用区内分别增加两个实例，以实现可用区之间均匀分配实例。

图 3-4 均匀实例分配

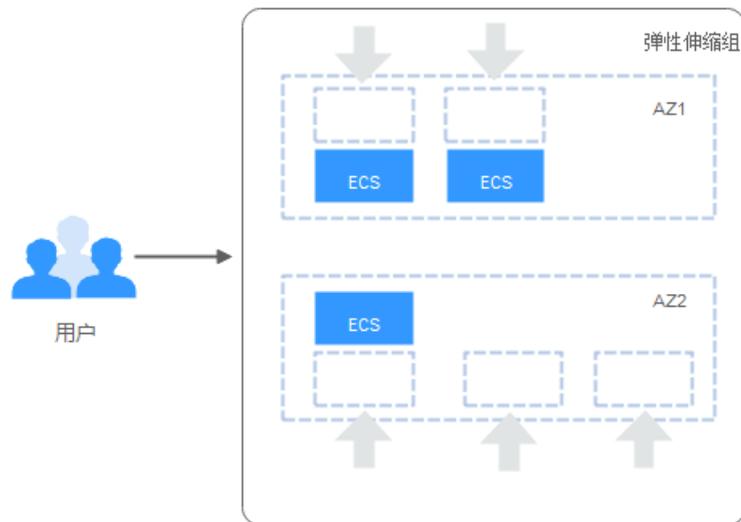


再均衡

手工加入或移出实例后，伸缩组中的实例没有均匀分配在可用区时，后续进行的伸缩活动会优先在可用区内均匀分配实例。

例如，伸缩组中目前有三个实例分布在两个可用区内，若该伸缩组下一个伸缩活动增加五个实例时，会在有两个实例的可用区内增加两个实例，在有一个实例的可用区增加三个实例，以实现可用区之间均匀分配实例。

图 3-5 再均衡



加强成本管理

弹性伸缩能够实现按需使用实例和带宽，并自动调整系统中的资源，节省了资源和人为调整资源带来的损耗，为您极大程度节约了成本。

提高可用性

弹性伸缩可确保应用系统始终拥有合适的容量以满足当前流量需求。

弹性伸缩和负载均衡结合使用

当您在使用弹性伸缩时，业务增长时应用系统自动扩容，业务下降时应用系统自动缩容，在伸缩组添加和删除实例时，须确保所有实例均可分配到应用程序的流量。弹性伸缩和负载均衡结合使用可以解决这个问题。

使用负载均衡后，伸缩组会自动地将加入伸缩组的实例绑定负载均衡监听器。访问流量将通过负载均衡监听器自动分发到伸缩组内的所有实例，提高了应用系统的可用性。若伸缩组中的实例上部署了多个业务，还可以添加多个负载均衡监听器到伸缩组，同时监听多个业务，从而提高业务的可扩展性。

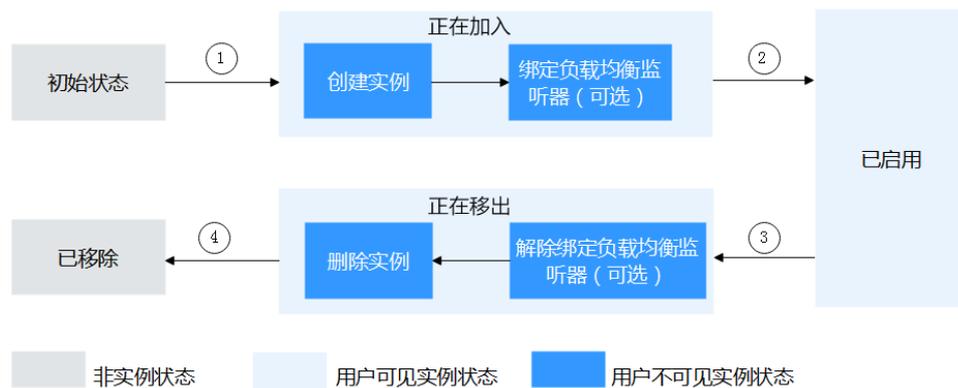
提高容错能力

弹性伸缩可以检测到应用系统中实例的运行状况，并启用新实例以替换运行状况不佳的实例。

4 生命周期

伸缩组中的实例生命周期，从创建实例开始，到该实例从伸缩组中移除结束。
伸缩组中未添加生命周期挂钩时，实例生命周期内状态之间的过渡如图4-1所示。

图 4-1 实例生命周期内状态之间的过渡



触发条件②和④表示系统自发触发实例状态的改变。

表 4-1 实例的状态

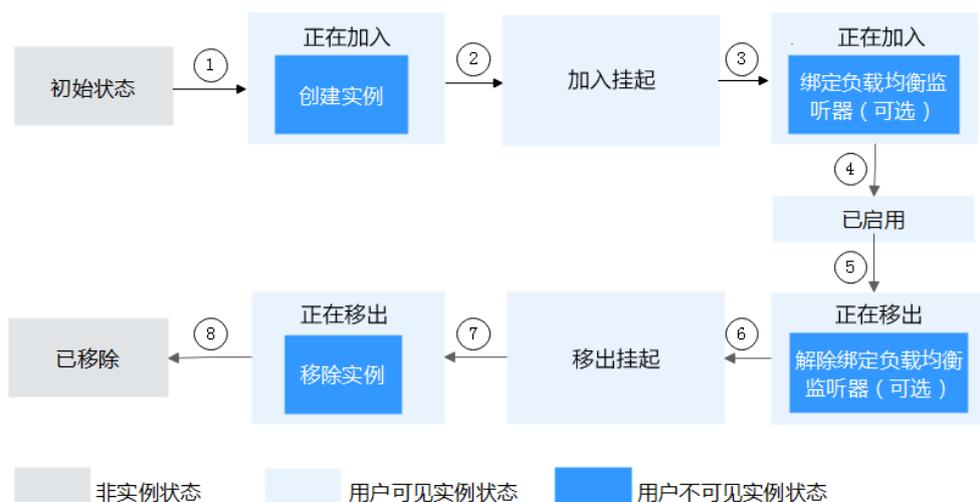
实例所处状态	子状态	实例状态含义	触发条件
初始状态	-	即实例还没状态。	触发条件①包括有两种情况，只要其中一种情况就能够触发实例进入“正在加入”状态。 <ul style="list-style-type: none"> • 手动修改期望实例数或满足伸缩策略的条件时，触发伸缩活动，进行扩容。 • 手动添加已有实例至伸缩组。
正在加入	创建实例 绑定负载均衡监听器 (可选)	在触发条件①的作用下，伸缩组开始扩容，创建实例。 在触发条件①的作用下，创建实例完成后，实例绑定负载均衡监听器。	

实例所处状态	子状态	实例状态含义	触发条件
已启用	-	实例进入伸缩组，开始接受处理业务流量。	触发条件③包括有三种情况，只要其中一种情况就能够触发实例从“已启用”状态到“正在移出”状态： <ul style="list-style-type: none"> • 手动修改期望实例数或满足伸缩策略的条件时，触发伸缩活动，进行扩容。 • 实例进入启用状态后，开始健康检查，健康检查失败后，移出实例。 • 手动将实例移出伸缩组。
正在移出	解除绑定负载均衡监听器（可选）	在触发条件③的作用下，伸缩组开始缩容，实例解除绑定负载均衡监听器。	
	删除实例	实例解除绑定负载均衡监听器后，从伸缩组中移出。	
已移除	-	实例在伸缩组中的生命周期已结束，即实例没有状态。	-

通过手动添加实例和伸缩活动向伸缩组添加实例，实例经过正在加入、已启用和正在移出状态后，实例将移出伸缩组。

伸缩组中已添加生命周期挂钩后，实例生命周期内状态之间的过渡如图4-2所示。当伸缩组进行伸缩活动时，实例将被生命周期挂钩挂起并置于等待状态，实例将保持此状态直至超时时间结束或者用户手动回调。用户能够在实例保持等待状态的时间段内执行自定义操作，例如，用户可以在新移入的实例上安装或配置软件，也可以实例终止前从实例中下载日志文件。

图 4-2 实例生命周期内状态之间的过渡



触发条件②、④、⑥、⑧表示系统自发触发实例状态的改变。

表 4-2 实例状态

实例所处状态	子状态	实例状态含义	触发条件含义
初始状态	-	即实例还没状态。	触发条件①包括有两种情况，只要其中一种情况就能够触发实例进入“正在加入”状态。
正在加入	创建实例	在触发条件①的作用下，伸缩组开始扩容，创建实例。	<ul style="list-style-type: none"> • 手动修改期望实例数或满足伸缩策略的条件时，触发伸缩活动，进行扩容。 • 手动添加已有实例至伸缩组。
加入挂起	-	正在加入伸缩组的实例被生命周期挂钩挂起，将实例至于等待的状态。	触发条件③包括有两种情况，只要其中一种情况就能够触发实例从“加入挂起”到“正在加入”状态。
正在加入	绑定负载均衡监听器（可选）	在触发条件③的作用下，实例将继续正在加入伸缩组，绑定负载均衡监听器。	<ul style="list-style-type: none"> • 默认回调操作 • 手动回调操作
已启用	-	实例进入伸缩组，开始接受处理业务流量。	触发条件⑤包括有三种情况，只要其中一种情况就能够触发实例从“已启用”状态到“正在移出”状态：
正在移出	解除绑定负载均衡监听器（可选）	在触发条件⑤的作用下，伸缩组开始缩容，实例解除绑定负载均衡监听器。	<ul style="list-style-type: none"> • 手动修改期望实例数或满足伸缩策略的条件时，触发伸缩活动，进行缩容。 • 实例进入启用状态后，开始健康检查，健康检查失败后，移出实例。 • 手动将实例移出伸缩组。
移出挂起	-	正在移出伸缩组的实例被生命周期挂钩挂起，将实例至于等待的状态。	触发条件⑦包括有两种情况，只要其中一种情况就能够触发实例从“移出挂起”到“正在移出”状态。
正在移出	删除实例	在触发条件⑦的作用下，实例将继续正在移出伸缩组，删除实例。	<ul style="list-style-type: none"> • 默认回调操作 • 手动回调操作
已移除	-	实例在伸缩组中的生命周期已结束，即实例没有状态。	-

通过手动添加实例和伸缩活动向伸缩组添加实例，实例经过正在加入、加入挂起、正在加入、已启用、正在移出、移出挂起和正在移出状态后，实例将移出伸缩组。

5 使用限制

功能限制

在应用系统中添加弹性伸缩后，使用时有以下功能限制：

- 弹性伸缩的云服务器中运行的应用需要是无状态、可横向扩展的。

📖 说明

- 无状态：**关于应用的既往事务，没有任何记录和参考，每项事务处理均是从头开始。无状态应用运行的实例不会在本地存储需要持久化的数据。
例如：可以将无状态事务看作一台自动售货机：一个请求对应一个响应。
- 有状态：**是可以周而复始、反复发生的应用和流程，操作是在之前的事务背景下执行的，当前事务可能会受到之前事务的影响。
有状态应用运行的实例会在本地存储需要持久化的数据。
例如：可以将有状态事务看作网上银行或电子邮件，有上下文记录。
- 弹性伸缩会自动释放云服务器，所以弹性伸缩组内的云服务器不可以保存应用的状态信息（例如session）和相关数据（如数据库、日志等）。如果应用中需要云服务器保存状态或日志信息，可以考虑把相关信息保存到独立的服务器中。
- 弹性伸缩无法纵向扩展，即弹性伸缩无法自动升降ECS实例的vCPU和内存等配置。

配额限制

弹性伸缩对用户的资源数量或容量做的配额限制如[表5-1](#)所示。

表 5-1 配额一览表

类别	描述	默认值
弹性伸缩组	用户可以创建的最多伸缩组个数。	50
弹性伸缩配置	用户可以创建的最多伸缩配置个数。	200
弹性伸缩策略	某个弹性伸缩组下可以创建的最多伸缩策略个数。	10

类别	描述	默认值
弹性伸缩实例	某个弹性伸缩组下可以创建的最多实例个数。	300
伸缩带宽策略	用户最多可以创建的伸缩带宽策略个数。	10
生命周期挂钩	某个弹性伸缩组内最多可添加的生命周期挂钩个数。	5
通知	某个弹性伸缩组最多可以配置的通知个数。	5
标签	某个弹性伸缩组最多可以添加的标签个数。	10

6 区域和可用区

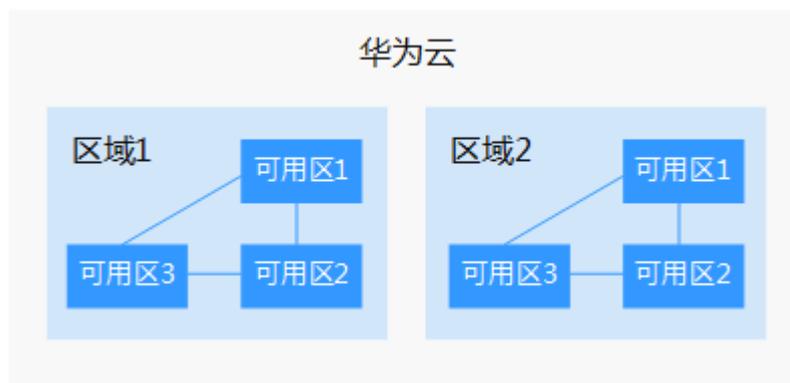
什么是区域、可用区？

区域和可用区用来描述数据中心的位置，您可以在特定的区域、可用区创建资源。

- 区域（Region）：从地理位置和网络时延维度划分，同一个Region内共享弹性计算、块存储、对象存储、VPC网络、弹性公网IP、镜像等公共服务。Region分为通用Region和专属Region，通用Region指面向公共租户提供通用云服务的Region；专属Region指只承载同一类业务或只面向特定租户提供业务服务的专用Region。
- 可用区（AZ，Availability Zone）：一个AZ是一个或多个物理数据中心的集合，有独立的风火水电，AZ内逻辑上再将计算、网络、存储等资源划分成多个集群。一个Region中的多个AZ间通过高速光纤相连，以满足用户跨AZ构建高可用性系统的需求。

图6-1阐明了区域和可用区之间的关系。

图 6-1 区域和可用区



目前，华为云已在全球多个地域开放云服务，您可以根据需求选择适合自己的区域和可用区。更多信息请参见华为云全球站点。

如何选择区域？

选择区域时，您需要考虑以下几个因素：

- 地理位置

一般情况下，建议就近选择靠近您或者您的目标用户的区域，这样可以减少网络时延，提高访问速度。

- 在除中国大陆以外的亚太地区有业务的用户，可以选择“中国-香港”、“亚太-曼谷”或“亚太-新加坡”区域。
- 在非洲地区有业务的用户，可以选择“非洲-约翰内斯堡”区域。
- 在拉丁美洲地区有业务的用户，可以选择“拉美-圣地亚哥”区域。

 说明

“拉美-圣地亚哥”区域位于智利。

- 资源的价格

不同区域的资源价格可能有差异，请参见华为云服务价格详情。

如何选择可用区？

是否将资源放在同一可用区内，主要取决于您对容灾能力和网络时延的要求。

- 如果您的应用需要较高的容灾能力，建议您将资源部署在同一区域的不同可用区内。
- 如果您的应用要求实例之间的网络延时较低，则建议您将资源创建在同一可用区内。

区域和终端节点

当您通过API使用资源时，您必须指定其区域终端节点。有关华为云的区域和终端节点的更多信息，请参阅[地区和终端节点](#)。

7 计费标准

弹性伸缩服务本身不收取费用，但伸缩组自动创建的按需付费实例需要支付相应的费用，实例的计费标准请参见[计费说明](#)。实例使用的弹性公网IP也需支付相应的费用，弹性公网IP的计费标准请参见[计费说明](#)。伸缩组进行减容时，自动创建的实例会被移出伸缩组并删除，删除后将不再收取费用。而之前通过手动移入的实例只会被移出伸缩组，系统仍会收取该实例的使用费用。若您不再需要使用该实例，请自行在ECS页面进行退订。

例如，弹性伸缩进行扩容活动创建了两台实例，使用一个小时后，进行了缩容活动，这两台实例被移出伸缩组并删除了，则系统只收取这两台实例使用一小时产生的费用。

8 安全

8.1 责任共担

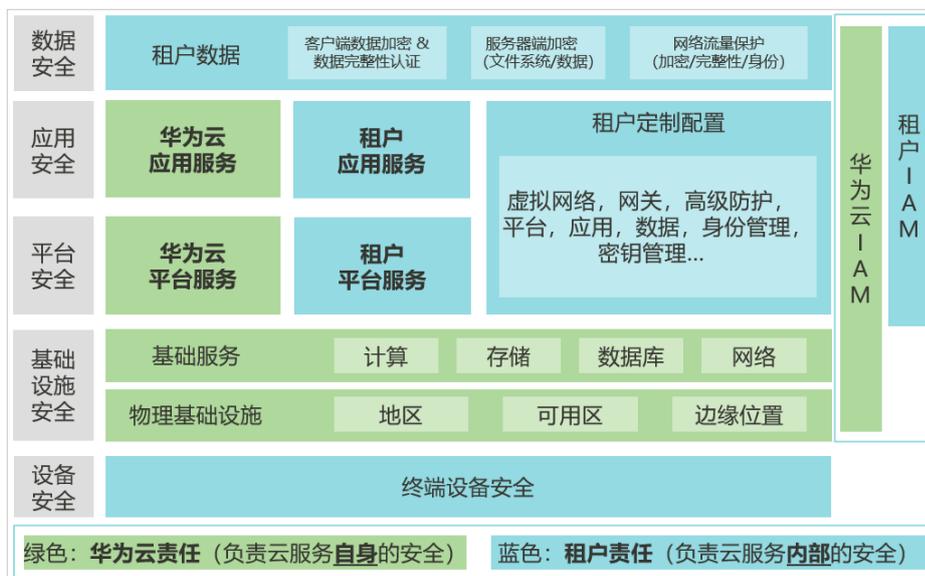
华为云秉承“将对网络和业务安全性保障的责任置于公司的商业利益之上”。针对层出不穷的云安全挑战和无孔不入的云安全威胁与攻击，华为云在遵从法律法规业界标准的基础上，以安全生态圈为护城河，依托华为独有的软硬件优势，构建面向不同区域和行业的完善云服务安全保障体系。

安全性是华为云与您的共同责任，如图8-1所示。

- **华为云**：负责云服务自身的安全，提供安全的云。华为云的安全责任在于保障其所提供的 IaaS、PaaS 和 SaaS 类云服务自身的安全，涵盖华为云数据中心的物理环境设施和运行其上的基础服务、平台服务、应用服务等。这不仅包括华为云基础设施和各项云服务技术的安全功能和性能本身，也包括运维运营安全，以及更广义的安全合规遵从。
- **租户**：负责云服务内部的安全，安全地使用云。华为云租户的安全责任在于对使用的 IaaS、PaaS 和 SaaS 类云服务内部的安全以及对租户定制配置进行安全有效的管理，包括但不限于虚拟网络、虚拟主机和访客虚拟机的操作系统，虚拟防火墙、API 网关和高级安全服务，各项云服务，租户数据，以及身份账号和密钥管理等方面的安全配置。

《[华为云安全白皮书](#)》详细介绍华为云安全性的构建思路与措施，包括云安全战略、责任共担模型、合规与隐私、安全组织与人员、基础设施安全、租户服务与租户安全、工程安全、运维运营安全、生态安全。

图 8-1 华为云安全责任共担模型



8.2 数据保护技术

用户加密，是指用户通过提供的加密特性，对弹性云服务器资源进行加密，从而提升数据的安全性。如果使用加密的弹性云服务器创建弹性伸缩配置，那么创建出来的伸缩配置，加密方式与原云服务器保持一致。请参见[用户加密](#)。

8.3 审计与日志

云审计服务（Cloud Trace Service, CTS），是华为云安全解决方案中专业的日志审计服务，提供对各种云资源操作记录的收集、存储和查询功能，可用于支撑安全分析、合规审计、资源跟踪和问题定位等常见应用场景。

弹性伸缩支持使用云审计记录服务资源操作。云审计记录的操作类型有三种，通过云平台帐户登录管理控制台执行的操作，通过云服务支持的API执行的操作，以及系统内部触发的操作。

在您的应用系统中启用云审计服务后，将在日志文件记录对弹性伸缩执行的API调用的操作。您可以在云审计服务管理控制台查询近7天内的操作记录。如果需要保存7天之前的操作记录，您可以通过对象存储服务（Object Storage Service，以下简称 OBS），将操作记录实时同步保存至OBS。

- CTS的详细介绍和开通配置方法请参见[CTS快速入门](#)。
- 云审计服务支持的AS操作列表请参见[记录弹性伸缩](#)。
- 查看审计日志请参见[查看审计日志](#)。

8.4 监控安全风险

监控指标

弹性伸缩支持云监控的监控指标，用户可以通过云监控检索弹性伸缩服务产生的监控指标和告警信息。查看弹性伸缩支持的监控指标请参见[监控指标说明](#)。

健康度检查

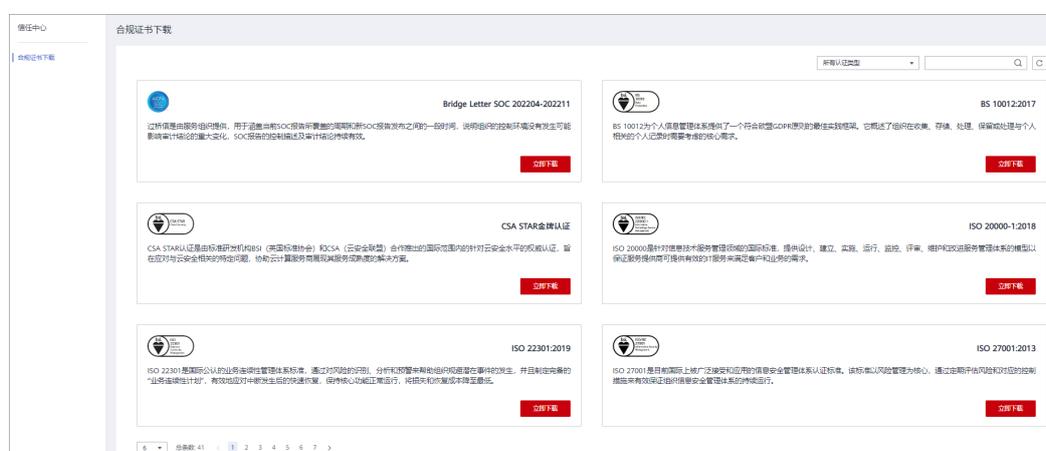
为使用户更好地掌握自己的弹性云服务器运行状态，云平台提供了云监控。您可以查看伸缩组的监控指标详情，更好地了解弹性云服务器的各项性能指标。详情请参考[查看监控指标数据](#)。

8.5 认证证书

合规证书

华为云服务及平台通过了多项国内外权威机构（ISO/SOC/PCI等）的安全合规认证，用户可自行[申请下载](#)合规资质证书。

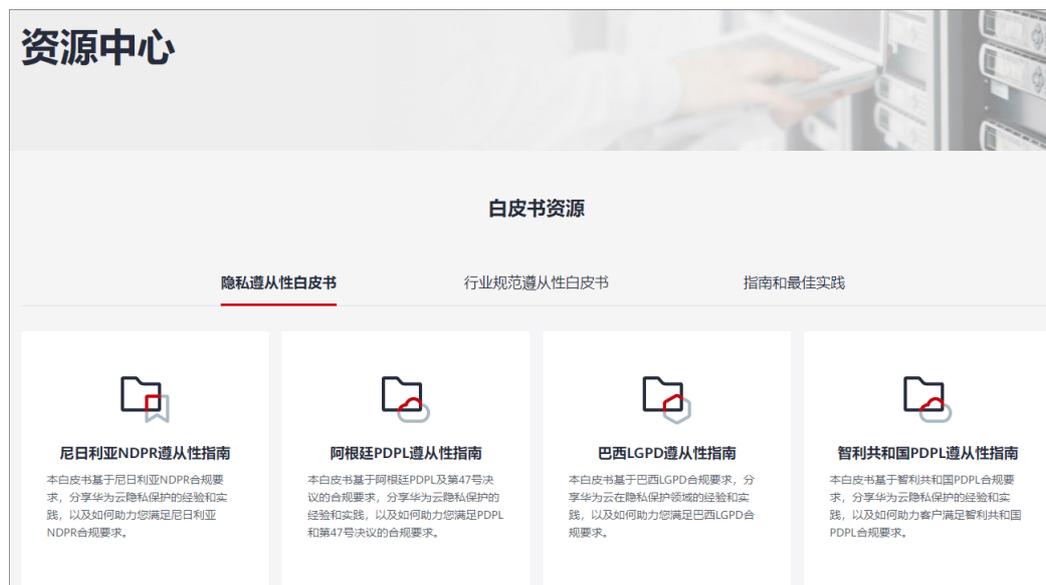
图 8-2 合规证书下载



资源中心

华为云还提供以下资源来帮助用户满足合规性要求，具体请查看[资源中心](#)。

图 8-3 资源中心



销售许可证&软件著作权证书

另外，华为云还提供了以下销售许可证及软件著作权证书，供用户下载和参考。具体请查看[合规资质证书](#)。

图 8-4 销售许可证&软件著作权证书



9 与其他服务的关系

除直接使用弹性伸缩提供的对资源进行调整的功能外，若您同时购买了云服务中的其他产品，可以结合其他产品一起使用，能满足您多种场景下对云产品的需求。

弹性伸缩服务与周边服务的依赖关系如图9-1所示。

图 9-1 弹性伸缩服务与其他服务的关系示意图

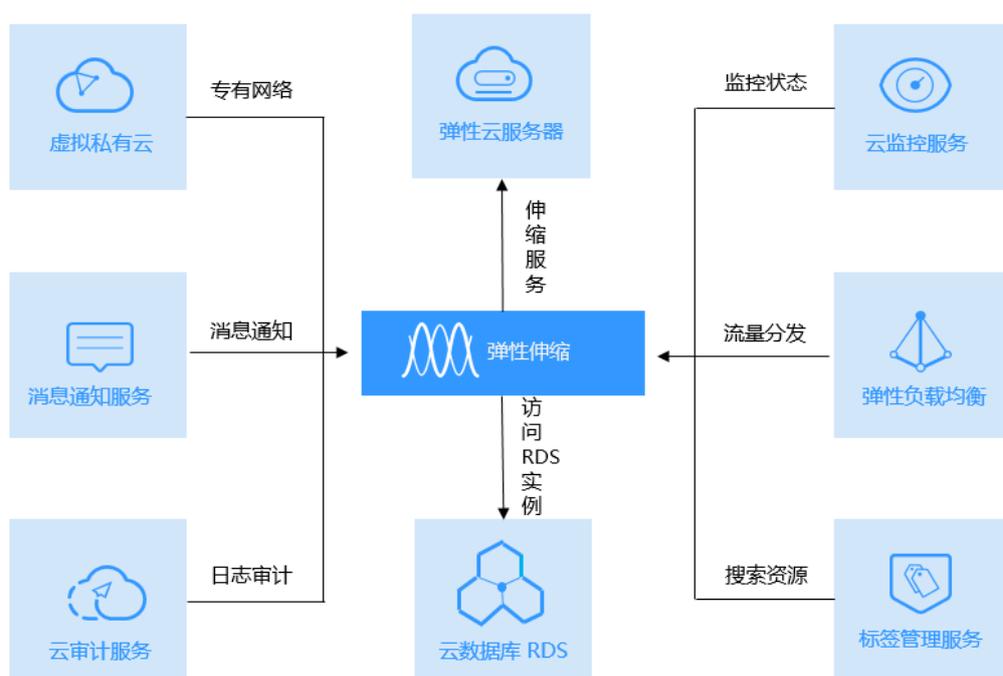


表 9-1 弹性伸缩与其他服务的关系

服务名称	说明	交互功能	相关内容
弹性负载均衡 (Elastic Load Balance)	当配置了负载均衡服务后, 弹性伸缩组在添加或移除云服务器时, 自动会为云服务器绑定或解绑负载均衡监听器。 AS支持ELB的前提是: 弹性伸缩组和负载均衡器必须处于同一VPC内。	使伸缩组中每一个实例均可分配到应用程序流量	添加负载均衡器到伸缩组
云监控服务 (Cloud Eye)	弹性伸缩配置了告警触发策略时, 会根据云监控的告警条件触发弹性伸缩活动。	通过监控伸缩组内实例的状态指标调节资源。	弹性伸缩支持的监控指标
弹性云服务器 (Elastic Cloud Server)	弹性伸缩活动中添加的云服务器可以通过弹性云服务器进行管理和维护。	自动调整弹性云服务器数量	动态扩展资源
虚拟私有云 (Virtual Private Cloud)	弹性伸缩支持自动调整虚拟私有云中创建的弹性公网IP带宽或共享带宽大小。	自动调整带宽大小	创建伸缩带宽策略
消息通知服务 (Simple Message Notification)	用户使用消息通知功能后, 系统会将伸缩组的多种情况及时推送给用户, 便于用户及时了解伸缩组的状态。	消息通知	为伸缩组配置通知
云审计服务 (Cloud Trace Service)	开通云审计服务后, 可以记录弹性伸缩相关的操作事件, 便于日后的查询、审计和回溯。	日志审计	记录弹性伸缩

服务名称	说明	交互功能	相关内容
标签管理服务 (Tag Management Service)	当您具有许多相同类型的弹性伸缩资源时，标签可以为您提供灵活的资源管理能力。	标签	标记伸缩组和实例
云数据库服务 (Relational Database Service)	伸缩出来的实例，可以直接访问RDS实例的前提条件是： <ul style="list-style-type: none">● 该实例与目标RDS实例必须处于同一VPC内；● 该实例必须处于目标RDS实例所属安全组允许访问的范围内；	伸缩出来的实例可以访问云数据库实例	通过内网连接RDS for MySQL实例

10 权限管理

如果您需要对华为云上购买的弹性伸缩（Auto Scaling，简称AS）资源，为企业中的员工设置不同的访问权限，以达到不同员工之间的权限隔离，您可以使用统一身份认证服务（Identity and Access Management，简称IAM）进行精细的权限管理。该服务提供用户身份认证、权限分配、访问控制等功能，可以帮助您安全的控制华为云资源的访问。

通过IAM，您可以在帐号中给员工创建IAM用户，并授权控制他们对华为云资源的访问范围。例如您的员工中有负责软件开发的人员，您希望他们拥有AS的使用权限，但是不希望他们拥有删除伸缩组等高危操作的权限，那么您可以使用IAM为开发人员创建用户，通过授予仅能使用伸缩组，但是不允许删除伸缩组的权限策略，控制他们对AS资源的使用范围。

如果华为云帐号已经能满足您的要求，不需要创建独立的IAM用户进行权限管理，您可以跳过本章节，不影响您使用AS服务的其它功能。

IAM是华为云提供权限管理的基础服务，无需付费即可使用，您只需要为您帐号中的资源进行付费。关于IAM的详细介绍，请参见[IAM产品介绍](#)。

AS 权限

默认情况下，管理员创建的IAM用户没有任何权限，需要将其加入用户组，并给用户组授予策略或角色，才能使得用户组中的用户获得对应的权限，这一过程称为授权。授权后，用户就可以基于被授予的权限对云服务进行操作。

AS部署时通过物理区域划分，为项目级服务，授权时，“作用范围”需要选择“区域级项目”，然后在指定区域（如华北-北京1）对应的项目（cn-north-1）中设置相关权限，并且该权限仅对此项目生效；如果在“所有项目”中设置权限，则该权限在所有区域项目中都生效。访问AS时，需要先切换至授权区域。

权限根据授权精细程度分为角色和策略。

- **角色：** IAM最初提供的一种根据用户的工作职能定义权限的粗粒度授权机制。该机制以服务为粒度，提供有限的服务相关角色用于授权。由于华为云各服务之间存在业务依赖关系，因此给用户授予角色时，可能需要一并授予依赖的其他角色，才能正确完成业务。角色并不能满足用户对精细化授权的要求，无法完全达到企业对权限最小化的安全管控要求。
- **策略：** IAM最新提供的一种细粒度授权的能力，可以精确到具体服务的操作、资源以及请求条件等。基于策略的授权是一种更加灵活的授权方式，能够满足企业对权限最小化的安全管控要求。例如：针对AS服务，管理员能够控制IAM用户仅

能对某一类云服务器资源进行指定的管理操作。多数细粒度策略以API接口为粒度进行权限拆分，AS支持的API授权项请参见[策略及授权项说明](#)。

如表10-1所示，包括了AS的所有系统权限。

表 10-1 AS 系统权限

策略名称	描述	类别	依赖关系
AutoScaling FullAccess	对弹性伸缩全部资源的所有执行权限。	系统策略	无
AutoScaling ReadOnlyAccess	弹性伸缩只读权限，对弹性伸缩全部资源的只读权限。	系统策略	无
AutoScaling Administrator	对弹性伸缩全部资源的所有执行权限。	系统角色	依赖ELB Administrator、CES Administrator、Server Administrator和Tenant Administrator角色，在同项目中勾选依赖的角色。

表10-2列出了AS常用操作与系统权限的授权关系，您可以参照该表选择合适的系统权限。

表 10-2 常用操作与系统权限的关系

操作	AutoScaling FullAccess	AutoScaling ReadOnlyAccess	AutoScaling Administrator
创建伸缩组	√	x	√
修改伸缩组	√	x	√
查询伸缩组详情	√	√	√
删除伸缩组	√	x	√
创建伸缩配置	√	x	√
创建伸缩策略	√	x	√
创建伸缩带宽策略	√	x	√

相关链接

- [IAM产品介绍](#)
- [创建用户并授予使用AS权限](#)
- [策略及授权项说明](#)

11 基本概念

伸缩组

伸缩组是具有相同应用场景的实例的集合，是启停伸缩策略和进行伸缩活动的基本单位。

伸缩配置

伸缩配置是伸缩组内实例（弹性云服务器）的模板，定义了伸缩组内待添加的实例的规格数据。包括云服务器类型、vCPU、内存、镜像、磁盘、登录方式等。

伸缩策略

伸缩策略可以触发伸缩活动，是对伸缩组中实例数量进行调整的一种方式。伸缩策略规定了伸缩活动触发需要满足的条件及需要执行的操作，当满足伸缩条件时，系统会自动触发一次伸缩活动。

伸缩活动

伸缩组中增加或减少实例的过程称为伸缩活动。伸缩活动的目的是使应用系统中当前实例数和期望实例数保持一致，或达到已设置的伸缩策略触发条件时，执行增加或减少实例数量的操作，保证业务正常运行。

冷却时间

为了避免告警策略频繁触发，必须设置冷却时间。冷却时间是指冷却伸缩活动的时间，单位为秒。在每次伸缩活动完成之后，系统开始计算冷却时间。伸缩组在冷却时间内，会拒绝告警策略的触发，其他类型的伸缩策略（如定时策略和周期策略）及手动触发不受限制。

例如：冷却时间设置为300秒，定时策略设置了10:32进行伸缩活动，10:30告警触发的伸缩活动结束，则在10:30-10:35时间内，伸缩组会拒绝新告警触发的伸缩活动，但不会拒绝在10:32时定时策略触发的伸缩活动；若10:36定时策略触发的伸缩活动结束，则冷却时间为10:36-10:41。

伸缩带宽

伸缩带宽可以根据用户配置的伸缩带宽策略自动调整带宽资源。弹性伸缩仅支持对按需购买的弹性公网IP带宽和共享带宽进行调整，不支持对包年包月的带宽进行调整。

12 修订记录

版本日期	变更说明
2022-11-15	第六次正式发布。 本次变更如下： 增加“安全”章节。
2021-08-27	第五次正式发布。 本次变更如下： 增加AS支持伸缩带宽的描述。
2020-10-19	第四次正式发布。 本次变更如下： 增加“访问方式”章节。
2020-02-10	第三次正式发布。 本次变更如下： AS系统权限名称变更。
2019-09-30	第二次正式发布。
2018-04-30	第一次正式发布。