

Auto Scaling

User Guide

Issue 13
Date 2021-10-30



Copyright © Huawei Technologies Co., Ltd. 2024. All rights reserved.

No part of this document may be reproduced or transmitted in any form or by any means without prior written consent of Huawei Technologies Co., Ltd.

Trademarks and Permissions



HUAWEI and other Huawei trademarks are trademarks of Huawei Technologies Co., Ltd.

All other trademarks and trade names mentioned in this document are the property of their respective holders.

Notice

The purchased products, services and features are stipulated by the contract made between Huawei and the customer. All or part of the products, services and features described in this document may not be within the purchase scope or the usage scope. Unless otherwise specified in the contract, all statements, information, and recommendations in this document are provided "AS IS" without warranties, guarantees or representations of any kind, either express or implied.

The information in this document is subject to change without notice. Every effort has been made in the preparation of this document to ensure accuracy of the contents, but all statements, information, and recommendations in this document do not constitute a warranty of any kind, express or implied.

Security Declaration

Vulnerability

Huawei's regulations on product vulnerability management are subject to the *Vul. Response Process*. For details about this process, visit the following web page:

<https://www.huawei.com/en/psirt/vul-response-process>

For vulnerability information, enterprise customers can visit the following web page:

<https://securitybulletin.huawei.com/enterprise/en/security-advisory>

Contents

1 AS Group	1
1.1 Creating an AS Group	1
1.2 (Optional) Adding a Load Balancer to an AS Group	8
1.3 Changing the AS Configuration for an AS Group	8
1.4 Enabling an AS Group	9
1.5 Disabling an AS Group	10
1.6 Modifying an AS Group	10
1.7 Deleting an AS Group	11
2 AS Configuration	13
2.1 Creating an AS Configuration	13
2.2 Creating an AS Configuration from an Existing ECS Instance	13
2.3 Creating an AS Configuration from a New Specifications Template	16
2.4 Copying an AS Configuration	21
2.5 Deleting an AS Configuration	21
3 AS Policy	23
3.1 Overview	23
3.2 Creating an AS Policy	24
3.3 Managing AS Policies	33
4 Scaling Action	36
4.1 Dynamic Scaling	36
4.2 Scheduled Scaling	38
4.3 Manual Scaling	38
4.4 Configuring an Instance Removal Policy	40
4.5 Viewing a Scaling Action	41
4.6 Managing Lifecycle Hooks	41
4.7 Configuring Instance Protection	48
4.8 Putting an Instance Into the Standby State	49
5 Bandwidth Scaling	51
5.1 Creating a Bandwidth Scaling Policy	51
5.2 Viewing Details About a Bandwidth Scaling Policy	58
5.3 Managing a Bandwidth Scaling Policy	58

6 AS Group and Instance Monitoring.....	61
6.1 Health Check.....	61
6.2 Configuring Notifications for an AS Group.....	62
6.3 Recording AS Resource Operations.....	63
6.4 Querying Real-Time Traces.....	66
6.5 Adding Tags to AS Groups and Instances.....	69
6.6 Monitoring Metrics.....	70
6.7 Viewing Monitoring Metrics.....	75
6.8 Setting Monitoring Alarm Rules.....	76
7 Permissions Management.....	77
7.1 Creating a User and Granting AS Permissions.....	77
7.2 AS Custom Policies.....	78
A Change History.....	80

1 AS Group

1.1 Creating an AS Group

Scenarios

An AS group consists of a collection of instances and AS policies that have similar attributes and apply to the same application scenario. An AS group is the basis for enabling or disabling AS policies and performing scaling actions. The pre-configured AS policy automatically adds or deletes instances to or from an AS group, or maintains a fixed number of instances in an AS group.

When creating an AS group, specify an AS configuration for it. Additionally, add one or more AS policies for the AS group.

Creating an AS group involves the configuration of the maximum, minimum, and expected numbers of instances and the associated load balancer.

Notes

ECS types available in different AZs may vary. When creating an AS group, choose an AS configuration that uses an ECS type available in the AZs used by the AS group.

- If the ECS type specified in the AS configuration is not available in any of the AZs used by the AS group, the following situations will occur:
 - If the AS group is disabled, it cannot be enabled again later.
 - If the AS group is enabled, its status will become abnormal when instances are added to it.
- If the ECS type specified in the AS configuration is only available in certain AZs used by the AS group, the ECS instances added by a scaling action are only deployed in the AZs where that ECS type is available. As a result, the instances in the AS group may not be evenly distributed.

Procedure

1. Log in to the management console.

2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
3. Click **Create AS Group**.
4. Set parameters, such as **Name**, **Max. Instances**, **Min. Instances**, and **Expected Instances**. [Table 1-1](#) describes the key parameters to be configured.

Table 1-1 AS group parameters

Parameter	Description	Example Value
Region	A region is where the AS group is deployed. Resources in different regions cannot communicate with each other over internal networks. For lower network latency and faster access to your resources, select the region nearest to your target users.	N/A
AZ	An AZ is a physical location where resources use independent power supply and networks. AZs are physically isolated but interconnected through an internal network. <ul style="list-style-type: none"> • If you require high availability, buy servers in different AZs. • If you require low network latency, buy servers in the same AZ. 	N/A
Multi-AZ Scaling Policy	This parameter can be set to Balanced or Sequenced . <ul style="list-style-type: none"> • Balanced: When scaling out an AS group, the system preferentially distributes ECS instances evenly across AZs used by the AS group. If it fails in the target AZ, it automatically selects another AZ based on the sequenced policy. • Sequenced: When scaling out an AS group, the system distributes ECS instances to the AZ selected according to the order in which AZs are specified. <p>NOTE This parameter needs to be configured when two or more AZs are selected.</p>	Balanced
Name	Specifies the name of the AS group to be created. The name contains 1 to 64 characters and consists of only letters, digits, underscores (_), and hyphens (-).	N/A
Max. Instances	Specifies the maximum number of ECS instances in an AS group.	1

Parameter	Description	Example Value
Expected Instance s	Specifies the expected number of ECS instances in an AS group. After an AS group is created, you can change this value, which will trigger a scaling action. The number of expected instances cannot be smaller than the minimum number of instances or greater than the maximum number of instances.	0
Min. Instance s	Specifies the minimum number of ECS instances in an AS group.	0
AS configur ation	Specifies the required AS configuration for the AS group. An AS configuration defines the specifications of the ECS instances to be added to an AS group. The specifications include the ECS image and system disk size. You need to create the required AS configuration before creating an AS group.	N/A
VPC	Provides a network for your ECS instances. All ECS instances in the AS group are deployed in this VPC.	N/A
Subnet	You can select up to five subnets. The AS group automatically binds all NICs to the created ECS instances. The first subnet is used by the primary NIC of an ECS instance by default, and other subnets are used by extension NICs of the instance.	N/A

Parameter	Description	Example Value
Load Balancing	<p>This parameter is optional. A load balancer automatically distributes traffic across all instances in an AS group to balance their service load. It improves the fault tolerance of your applications and expands application service capabilities.</p> <p>NOTE</p> <ul style="list-style-type: none">• Up to six load balancers can be added to an AS group.• After multiple load balancers are added to an AS group, multiple services can be concurrently listened to, thereby improving service scalability. If ELB health check is selected for Health Check Method, when any one of the listeners detects that an instance becomes unhealthy, AS will replace it with a new one. <p>If you select Elastic load balancer, configure the following parameters:</p> <ul style="list-style-type: none">• Load Balancer• Backend ECS Group• Backend Port: specifies the port on which a backend ECS listens for traffic.• Weight: determines the portion of requests a backend ECS processes compared to other backend ECSs added to the same listener. For more information about load balancing, see Elastic Load Balance User Guide.	N/A

Parameter	Description	Example Value
Instance Removal Policy	<p>Controls which instances are first to be removed during scale in. If specified conditions are met, scaling actions are triggered to remove instances. You can choose from any of the following instance removal policies:</p> <ul style="list-style-type: none"> • Oldest instance created from oldest AS configuration: The oldest instance created from the oldest configuration is removed from the AS group first. • Newest instance created from oldest AS configuration: The newest instance created from the oldest configuration is removed from the AS group first. • Oldest instance: The oldest instance is removed from the AS group first. • Newest instance: The latest instance is removed from the AS group first. <p>NOTE Manually added ECS instances are the last to be removed. If AS does remove a manually added instance, it only removes the instance from the AS group. It does not delete instance. If multiple manually added instances must be removed, AS preferentially removes the earliest-added instance first.</p>	Oldest instance created from oldest AS configuration
EIP	<p>If EIP has been selected in an AS configuration for an AS group, an EIP is automatically bound to the ECS instance added by a scaling action to the AS group. If you select Release, the EIP bound to an instance is released when the instance is removed from the AS group. Otherwise, the system unbinds the EIP from the instance, but does not release it when the instance is removed from the AS group.</p>	N/A
Data Disk	<p>If Data Disk is configured in the AS configuration used by the AS group, a data disk will be automatically created and attached to the ECS instances added during a scaling action to the AS group.</p> <p>If you select Release, the data disks attached to an instance will be deleted when the instance is removed from the AS group. Otherwise, the system detaches the data disks from the instance, but does not release them when the instance is removed from the AS group.</p>	N/A

Parameter	Description	Example Value
Health Check Method	<p>When a health check detects an unhealthy ECS instance, AS replaces it with a new one. You can choose from either of the following health check methods:</p> <ul style="list-style-type: none"> • ECS health check: checks ECS instance health status. If an instance is stopped or deleted, it is considered to be unhealthy. This method is selected by default. Using this method, the AS group periodically evaluates the running status of each instance based on the health check results. If the health check results show that an instance is unhealthy, AS removes the instance from the AS group. • ELB health check: determines ECS running status using a load balancing listener. When a load balancing listener detects that an instance is unhealthy, AS removes the instance from the AS group. 	N/A
Health Check Interval	Specifies the length of time between health checks. You can set a health check interval, such as 5 minutes, 15 minutes, 1 hour, or 3 hours, based on the service requirements.	5 minutes
Enterprise Project	<p>Specifies the enterprise project to which the AS group belongs. If an enterprise project is configured for an AS group, ECSs created in this AS group also belong to this enterprise project. If you do not specify an enterprise project, the default enterprise project will be used.</p> <p>NOTE</p> <ul style="list-style-type: none"> • Value default indicates the default enterprise project. Resources that are not allocated to any enterprise projects under your account are displayed in the default enterprise project. • Enterprise project is an upgraded version of IAM. It allocates and manages resources of different projects. 	N/A

Parameter	Description	Example Value
Tag	<p>If you have many resources of the same type, you can use tags to manage your resources. You can identify specified resources quickly using the tags allocated to them.</p> <p>If you have configured tag policies for AS, add tags to AS groups based on the tag policies. If you add a tag that does not comply with the tag policies, AS groups may fail to be created. Contact your administrator to learn more about tag policies.</p> <p>Each tag contains a key and a value. You can specify the key and value for each tag.</p> <ul style="list-style-type: none">• Key<ul style="list-style-type: none">- The key must be specified.- The key must be unique to the AS group.- The key can include up to 36 characters. It cannot contain non-printable ASCII characters (0-31) or the following characters: =*<>\\, /• Value<ul style="list-style-type: none">- The value is optional.- A key can have only one value.- The value can include up to 43 characters. It cannot contain non-printable ASCII characters (0-31) or the following characters: =*<>\\, /	N/A

5. Click **Next**.
6. On the displayed page, you can use an existing AS configuration or create an AS configuration.
7. Click **Next**.
8. (Optional) Add an AS policy to an AS group.

On the displayed page, click **Add AS Policy**.

Configure the required parameters, such as the **Policy Type**, **Scaling Action**, and **Cooldown Period**.

 **NOTE**

- If a scaling action is triggered by an AS policy, the cooldown period is whatever configured for that AS policy.
- If a scaling action is triggered by manually changing the expected number of instances or by other actions, the cooldown period is whatever configured for the AS group. The default cooldown period is 300 seconds.

9. Click **Create Now**.

10. Check the AS group, AS configuration, and AS policy information. Click **Submit**.
11. Confirm the creation result and go back to the **AS Groups** page as prompted. After the AS group is created, its status changes to **Enabled**.

1.2 (Optional) Adding a Load Balancer to an AS Group

Elastic Load Balance (ELB) automatically distributes incoming traffic across multiple backend servers based on configured forwarding policies. ELB expands the service capabilities of applications and improves their availability by eliminating single points of failure (SPOFs).

If ELB functions are required, perform the operations provided in this section to add a load balancer to your AS group. The load balancer added to an AS group distributes application traffic to all instances in the AS group when an instance is added to or deleted from the AS group.

Only a created load balancer can be bound to an AS group, and the AS group and load balancer must be in the same VPC. For details about how to create a load balancer, see *Elastic Load Balance User Guide*. To add a load balancer for an AS group, perform the following operations:

- When creating an AS group, configure parameter **Load Balancing** to add a load balancer. For details, see [Creating an AS Group](#).
- If an AS group has no scaling action ongoing, modify parameter **Load Balancing** to add a load balancer. For details, see [Modifying an AS Group](#).

1.3 Changing the AS Configuration for an AS Group

Scenarios

If you need to change the specifications of ECS instances in an AS group, changing the AS configuration used by the AS group is an easy way to help you get there.

Effective Time of New AS Configuration

After you change the AS configuration for an AS group, the new AS configuration will not be used until any ongoing scaling actions are complete.

For example, if there is a scaling action ongoing for an AS group, and you change the AS configuration of the AS group from **as-config-A** to **as-config-B**, **as-config-A** is still used for the instances that are being added in the ongoing scaling action.

as-config-B will take effect in the next scaling action.

Figure 1-1 Changing the AS configuration

Name	Status	Specifications	Image	System Disk	Data Disks	Login Mode	Created	Billing Mode	Operation
as-config-B	Unbound	kc1.large2 2 vCPUs 4 GB	CentOS 8.0 64bit with ARM	High I/O 40 GB	0	Password	Oct 09, 2020 15:17:52 ...	Pay-per-use	Copy Delete
as-config-A	Unbound	s2.small1 1 vCPUs 1 GB	CentOS 7.6 64bit	High I/O 40 GB	0	Password	Oct 09, 2020 14:42:33 ...	Pay-per-use	Copy Delete

Procedure

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. Click the name of the AS group for which you want to change the AS configuration. On the **Basic Information** page, click **Change Configuration** to the right of **Configuration Name**.
You can also locate the row containing the target AS group and choose **More > Change Configuration** in the **Operation** column.
4. In the displayed **Change AS Configuration** dialog box, select another AS configuration to be used by the AS group.
5. Click **OK**.

1.4 Enabling an AS Group

Scenarios

You can enable an AS group to automatically scale capacity in or out.

After an AS group is enabled, its status changes to **Enabled**. AS monitors the AS policy and triggers a scaling action for AS groups only in **Enabled** state. After an AS group is enabled, AS triggers a scaling action to automatically add or remove instances if the number of instances in the AS group is different from the expected number of instances.

- Only AS groups in the **Abnormal** state can be forcibly enabled. You can choose **More > Forcibly Enable** to enable an abnormal AS group. Forcibly enabling an AS group does not have adverse consequences.
- After you create an AS group and add an AS configuration to an AS group, the AS group is automatically enabled.

Enabling an AS Group

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. In the AS group list, locate the row containing the target AS group and click **Enable** in the **Operation** column. You can also click the AS group name and then **Enable** in the upper right corner of the page to enable the AS group.
4. In the **Enable AS Group** dialog box, click **Yes**.

Forcibly Enabling an AS Group

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. In the AS group list, locate the row containing the target AS group and select **Forcibly Enable** from the **More** drop-down list in the **Operation** column. You

can also click the AS group name and then **Forcibly Enable** in the upper right corner of the page to enable the AS group.

4. In the **Forcibly Enable AS Group** dialog box, click **Yes**.

1.5 Disabling an AS Group

Scenarios

If you need to stop an instance in an AS group for configuration or upgrade, disable the AS group before performing the operation. This prevents the instance from being deleted in a health check. When the instance status is restored, you can enable the AS group again.

If a scaling action keeps failing and being retried (the failure cause can be viewed on the **Elastic Cloud Server** page) for an AS group, use either of the following methods to stop the action from being repeated:

- Disable the AS group. Then, after the scaling action fails, it will not be retried. Enable the AS group again when the environment recovers or after replacing the AS configuration.
- Disable the AS group and change the expected number of instances to the number of existing instances. Then after the scaling action fails, the scaling action will not be retried.

After an AS group is disabled, its status changes to **Disabled**. AS does not automatically trigger any scaling actions for a **Disabled** AS group. When an AS group has an in-progress scaling action, the scaling action does not stop immediately after the AS group is disabled.

You can disable an AS group when its status is **Enabled** or **Abnormal**.

Procedure

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. In the AS group list, locate the row containing the target AS group and click **Disable** in the **Operation** column. You can also click the AS group name and then **Disable** in the upper right corner of the page to disable the AS group.
4. In the **Disable AS Group** dialog box, click **Yes**.

1.6 Modifying an AS Group

Scenarios

You can modify an AS group if needed. The settings of the following parameters can be changed: **Name**, **Max. Instances**, **Min. Instances**, **Expected Instances**, **Health Check Method**, **Health Check Interval**, **Instance Removal Policy**, **Cooldown Period**, and **Multi-AZ Scaling Policy**.

 NOTE

Changing the value of **Expected Instances** will trigger a scaling action. AS will automatically increase or decrease the number of instances to the value of **Expected Instances**.

If the AS group is not enabled, contains no instance, and has no scaling action ongoing, you can modify **Subnet** configurations. If an AS group has no scaling action ongoing, you can modify its **AZ** and **Load Balancing** configurations.

Procedure

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. In the AS group list, locate the AS group you want to modify and choose **More > Modify** in the **Operation** column.
You can also click the AS group name to switch to the **Overview** page, and click **Modify** in the upper right corner.
4. In the **Modify AS Group** dialog box, modify related data, for example, the expected number of instances.
5. Click **OK**.

1.7 Deleting an AS Group

Scenarios

You can delete an AS group when it is no longer required.

- If an AS group is not required during a specified period, you are advised to disable it but not delete it.
- For an AS group that has instances or ongoing scaling actions, if you attempt to forcibly delete the AS group and remove and delete the instances in the AS group, the AS group enters the deleting state, rejects new scaling requests, waits until the ongoing scaling action completes, and removes all instances from the AS group. Then, the AS group is automatically deleted. Instances automatically created are removed and deleted, but instances manually added are only removed out of the AS group. During this process, other operations cannot be performed in the AS group.

 CAUTION

Forcibly deleting an AS group may not delete ECS instances in the group.

- When an AS group is deleted, its AS policies and the alarm rules generated based on those AS policies will be automatically deleted.

Procedure

1. Log in to the management console.

2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. In the AS group list, locate the row containing the target AS group and choose **More > Delete** in the **Operation** column.
4. In the displayed **Delete AS Group** dialog box, click **Yes**.

2 AS Configuration

2.1 Creating an AS Configuration

An AS configuration defines the specifications of the ECS instances to be added to an AS group. The specifications include the ECS image and system disk size.

Scenarios

- When you create an AS group, create a new AS configuration or use an existing AS configuration.
- Create the required AS configuration on the **Instance Scaling** page.
- Change the AS configuration on the AS group details page.

Methods

- Create an AS configuration from an existing ECS instance.
If you create an AS configuration from an existing ECS instance, the vCPU, memory, image, disk, and ECS type are the same as those of the selected instance by default. For details, see [Creating an AS Configuration from an Existing ECS Instance](#).
- Create an AS configuration from a new specifications template.
If you have special requirements on the ECSs for resource expansion, use a new specifications template to create the AS configuration. For details, see [Creating an AS Configuration from a New Specifications Template](#).

2.2 Creating an AS Configuration from an Existing ECS Instance

Scenarios

You can use an existing ECS instance to rapidly create an AS configuration. In such a case, the parameter settings, such as the ECS type, vCPUs, memory, image, and disk settings (including size, type, encryption, and key) in the AS configuration are the same as those of the selected instance by default.

Procedure

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
3. Click **Create AS Configuration**.
4. Set the parameters for the AS configuration. [Table 2-1](#) lists the AS configuration parameters.

Table 2-1 AS configuration parameters

Parameter	Description	Example Value
Region	A region is where an AS configuration resides.	N/A
Name	Specifies the name of an AS configuration.	N/A
Configuration Template	Choose Use existing ECS > Select ECS . In such a case, the parameter settings, such as the ECS type, vCPUs, memory, image, and disk settings (including size, type, encryption, and key) in the AS configuration are the same as those of the selected instance by default.	Use existing ECS
EIP	An EIP is a static public IP address bound to an ECS in a VPC. Using the EIP, the ECS provides services externally. The following options are provided: <ul style="list-style-type: none"> • Do not use An ECS without an EIP cannot access the Internet. However, it can still be used as a service ECS or deployed in a cluster on a private network. • Automatically assign An EIP with a dedicated bandwidth is automatically assigned to each ECS. The bandwidth size is configurable. <p>NOTE If you select Automatically assign, specify Type, Billed By, and Bandwidth.</p>	Automatically assign

Parameter	Description	Example Value
Login Mode	<p>An ECS can be authorized using a key pair or a password.</p> <ul style="list-style-type: none">• Key pair Keys are used for authenticating the users who attempt to log in to ECS instances. If you select this mode, create or import a key pair on the Key Pair page.<p>NOTE If you use an existing key, make sure that you have saved the key file locally. Without the key, you will not be able to log in to your instance.</p>• Password The initial password of user root (for Linux) or user Administrator (for Windows) is used for authentication. You can log in to an ECS instance using a username and password combination.	Admin@123
Advanced Settings	<p>This allows you to configure User Data and ECS Group.</p> <p>You can select Do not configure or Configure now.</p>	N/A
User Data	<p>Enables an ECS to automatically inject user data when the ECS starts for the first time. This configuration is optional. If this function is enabled, the ECS automatically injects user data during its first startup.</p> <p>For details, see Elastic Cloud Server User Guide.</p> <p>The following two methods are available:</p> <ul style="list-style-type: none">• As text: allows you to enter the user data in the text box below.• As file: allows you to inject a script or other files when you create an ECS instance. <p>NOTE</p> <ul style="list-style-type: none">• For Linux, if you use password authentication, this function is not supported.• If the selected image does not support user data injection, this function is not supported.	N/A
ECS Group	<p>An ECS group allows you to create ECSs on different hosts to improve service reliability.</p> <p>For details, see Managing ECS Groups.</p>	N/A

5. Click **Create Now**.
6. If you want to use the newly created AS configuration, add it to the AS group. For details, see [Changing the AS Configuration for an AS Group](#).

2.3 Creating an AS Configuration from a New Specifications Template

Scenarios

If you have special requirements on the ECS instances for resource expansion, use a new specifications template to create the AS configuration. In such a case, ECS instances that have the specifications specified in the template will be added to the AS group in scaling actions.

Procedure

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
3. Click **Create AS Configuration**.
4. Set the parameters for the AS configuration. [Table 2-2](#) lists the AS configuration parameters.

Table 2-2 AS configuration parameters

Parameter	Description	Example Value
Region	A region is where an AS configuration resides.	N/A
Name	Specifies the name of the AS configuration to be created.	N/A
Configuration Template	Select Create new template . If this option is selected, configure parameters, such as the vCPUs, memory, image, disk, and ECS type, to create a new AS configuration.	Create new template

Parameter	Description	Example Value
CPU Architecture	<p>Both x86 and Kunpeng CPU architectures are available:</p> <ul style="list-style-type: none"> • x86: The x86-based CPU architecture uses Complex Instruction Set Computing (CISC). • Kunpeng: The Kunpeng-based CPU architecture uses Reduced Instruction Set Computing (RISC). <p>NOTE This parameter is displayed only when both x86-based and Kunpeng-based ECSs are available in the current region.</p>	x86
Specifications	<p>The public cloud provides various ECS types for different application scenarios. For more information, see Elastic Cloud Server User Guide.</p> <p>Configure the ECS specifications, including vCPUs, memory, image type, and disk, according to the ECS type.</p>	Memory-optimized ECS
Image	<ul style="list-style-type: none"> • Public image A public image is a standard, widely used image. It contains an OS and preinstalled public applications and is available to all users. You can configure the applications or software in the public image as needed. • Private image A private image is an image available only to the user who created it. It contains an OS, preinstalled public applications, and the user's private applications. Using a private image to create ECSs frees you from configuring multiple ECSs repeatedly. • Shared image A shared image is a private image shared by another public cloud user. 	Public image

Parameter	Description	Example Value
Disk	<p>Includes system and data disks.</p> <ul style="list-style-type: none">• System Disk Common I/O: uses Serial Advanced Technology Attachment (SATA) drives to store data. High I/O: uses serial attached SCSI (SAS) drives to store data. General Purpose SSD: uses solid state disk (SSD) drives to store data. Extreme SSD: uses enhanced solid state disk (ESSD) drives to store data. Ultra-high I/O: uses solid state disk (SSD) drives to store data. <p>If a full-ECS image is used, the system disk is restored using the disk backup. On the console, you can only change the volume type and size. In addition, the volume cannot be smaller than the disk backup.</p> <p>NOTE Different ECS flavors support different disk types. The supported disk types will be displayed on the management console.</p> <ul style="list-style-type: none">• Data Disk You can create multiple data disks for an ECS instance. In addition, you can specify a data disk image for exporting data. <p>If the image you selected is a full-ECS image, you can change the volume type and size and encryption attributes of the data disk restored using the disk backup. Ensure that the disk is at least as big as the disk backup. The encryption attributes can only be modified if the disk backup is in the same region as the disk.</p>	Common I/O for System Disk
Security Group	Controls ECS access within or between security groups by defining access rules. ECSs added to a security group are protected by the access rules you define.	N/A

Parameter	Description	Example Value
EIP	<p>An EIP is a static public IP address bound to an ECS in a VPC. Using the EIP, the ECS provides services externally.</p> <p>The following options are provided:</p> <ul style="list-style-type: none">• Do not use: An ECS without an EIP cannot access the Internet. However, it can still be used as a service ECS or deployed in a cluster on a private network.• Automatically assign: An EIP with a dedicated bandwidth is automatically assigned to each ECS. You can set the bandwidth size. If you select Automatically assign, you need to specify EIP Type, Billed By, and Bandwidth. <p>NOTE If you select Automatically assign, you need to specify EIP Type, Billed By, and Bandwidth.</p>	Automatically assign
Bandwidth	<p>Specifies the bandwidth size in Mbit/s.</p> <p>NOTE</p> <ul style="list-style-type: none">• This parameter is available only when EIP is set to Automatically assign.• If you select Dedicated, you can select Bandwidth or Traffic for Billed By.• The shared bandwidth can be billed only by bandwidth. You can select a shared bandwidth to which the EIP is to be added.	100
Login Mode	<p>An ECS can be authorized using a key pair or a password.</p> <ul style="list-style-type: none">• Key pair Keys are used for authenticating the users who attempt to log in to ECS instances. If you select this mode, create or import a key pair on the Key Pair page. <p>NOTE If you use an existing key, make sure that you have saved the key file locally.</p> <ul style="list-style-type: none">• Password The initial password of user root (for Linux) or user Administrator (for Windows) is used for authentication. You can log in to an ECS instance using a username and password combination.	Admin@123

Parameter	Description	Example Value
Advanced Settings	This parameter allows you to configure ECS Group , and User Data . You can select Do not configure or Configure now .	N/A
User Data	Enables an ECS to automatically inject user data when the ECS starts for the first time. This configuration is optional. After this function is enabled, the ECS automatically injects user data during its first startup. For details, see Elastic Cloud Server User Guide . The following methods are available: <ul style="list-style-type: none">• As text: allows you to enter the user data in the text box below.• As file: allows you to inject script files or other files when you create an ECS. NOTE <ul style="list-style-type: none">• For Linux, if you use password authentication, this function is not supported.• If the selected image does not support user data injection, this function is not supported.	N/A
ECS Group	An ECS group allows you to create ECSs on different hosts to improve service reliability.	N/A

5. Click **Create Now**. The system displays a message indicating that the AS configuration is successfully created.
6. If you want to use the newly created AS configuration, add it to the AS group. For details, see [Creating an AS Group](#) or [Changing the AS Configuration for an AS Group](#).

Notes on Multiple Flavors in an AS Configuration

AS configuration supports multiple flavors to minimize the probability of capacity expansion failures due to insufficient or unavailable flavors and ensure that capacity expansion succeeds during peak hours.

A maximum of 10 flavors can be selected for an AS configuration.

Applicable Scenario

- No special requirement for the instance flavors created in the AS group
- Requiring higher success ratio and low latency of creating instances in an AS group

- Requiring instances with high specifications
- Services that are stateless and can be horizontally scaled

The AS group sorts multiple flavors in either of the following ways:

- **Sequenced:** During AS group expansion, flavors are used based on the sequence they are selected. When the first flavor is insufficient or the instance fails to be created due to other reasons, the system attempts to create an instance of the second flavor, and so on.
- **Cost-centered:** During AS group expansion, the flavor with the minimum cost comes first. When creating an instance in an AS group, the system selects the flavor with the minimum cost. If the instance cannot be created, the system selects one with the minimum cost from the remained flavors, and so on.

2.4 Copying an AS Configuration

Scenarios

You can copy an existing AS configuration.

When copying an AS configuration, you can modify parameter settings, such as the configuration name, ECS specifications, and image of the existing AS configuration to rapidly add a new AS configuration.

Procedure

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
3. Click the **AS Configurations** tab, locate the row containing the target AS configuration, and click **Copy** in the **Operation** column.
4. On the **Copy AS Configuration** page, modify parameter settings, such as **Name**, **Specifications**, and **Image**, and configure the ECS login mode based on service requirements.
5. Click **OK**.

2.5 Deleting an AS Configuration

Scenarios

When you no longer need an AS configuration, you can delete it as long as the AS configuration is not used by an AS group. You can delete a single AS configuration or delete them in batches.

Procedure

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.

3. Click the **AS Configurations** tab page, locate the row containing the target AS configuration, and click **Delete** in the **Operation** column to delete this AS configuration. You can also select multiple AS configurations to be deleted and click **Delete** in the upper part of the AS configuration list to delete them all at once.

3 AS Policy

3.1 Overview

AS policies can trigger scaling actions to adjust bandwidth or the number of instances in an AS group. An AS policy defines the conditions for triggering a scaling action and the operation that will be performed. When the triggering condition is met, a scaling action is triggered automatically.

NOTE

If multiple AS policies are applied to an AS group, a scaling action is triggered as long as any of the AS policies is invoked, provided that the AS policies do not conflict with each other.

The number of instances in the AS group will never exceed the specified maximum and minimum numbers of instances.

Restrictions

A maximum of 10 AS policies can be created for an AS group.

AS supports the following policies:

- Alarm policy: AS automatically adjusts the number of instances in an AS group or sets the number of instances to the configured value when an alarm is generated for a configured metric, such as CPU Usage.
- Scheduled policy: AS automatically increases or decreases the number of instances in an AS group or sets the number of instances to the configured value at a specified time.
- Periodic policy: AS automatically increases or decreases the number of instances in an AS group or sets the number of instances to the configured value at a configured interval, such as daily, weekly, and monthly.

Resource Adjustment Modes

- Dynamic
AS adjusts the number of instances or bandwidth when an alarm policy is triggered.

This mode is suitable for scenarios where workloads are unpredictable. Alarm policies are used to trigger scaling actions based on real-time monitoring data (such as CPU usage) to dynamically adjust the number of instances in the AS group.

- Planned

AS adjusts the number of instances or bandwidth when a periodic or scheduled policy is triggered.

This mode is suitable for scenarios where workloads are periodic.

- Manual

AS allows you to adjust resources by manually adding instances to an AS group, removing instances from an AS group, or changing the expected number of instances.

3.2 Creating an AS Policy

Scenarios

You can create different types of AS policies. In an AS policy, you can define the conditions for triggering a scaling action and what operation to be performed. When the conditions are met, AS automatically triggers a scaling action to adjust the number of instances in the AS group.

This section describes how to create alarm-based, scheduled, or periodic AS policy for an AS group.

Creating an Alarm Policy

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
3. Locate the row containing the target AS group and click **View AS Policy** in the **Operation** column.
4. On the **AS Policies** page, click **Add AS Policy**.
5. Set the parameters listed in [Table 3-1](#).

Table 3-1 AS policy parameters

Parameter	Description	Example Value
Policy Name	Specifies the name of the AS policy to be created.	as-policy-p6g5
Policy Type	Select Alarm .	Alarm

Parameter	Description	Example Value
Alarm Rule	<p>Specifies whether a new alarm rule is to be created (Create) or an existing alarm rule will be used (Use existing).</p> <p>For details about how to use an existing alarm rule, see Setting Monitoring Alarm Rules.</p> <p>If you choose to create an alarm, system monitoring and custom monitoring are supported.</p> <ul style="list-style-type: none"> • System monitoring requires the parameters in Table 3-2. • Custom monitoring requires the parameters in Table 3-3. <p>NOTE If you select Use existing, do not choose alarm rules that are:</p> <ul style="list-style-type: none"> • Used to monitor all resources. • Created by associating templates. 	N/A

Parameter	Description	Example Value
Scaling Action	<p>Specifies an action and the number or percentage of instances.</p> <p>The following scaling action options are available:</p> <ul style="list-style-type: none">• Add Adds instances to an AS group when the scaling action is performed.• Reduce Removes instances from an AS group when the scaling action is performed.• Set to Maintains a fixed number of instances in an AS group.	<ul style="list-style-type: none">• Add 1 instance• Add 10% instances The number of instances to be added is 10% of the current number of instances in the AS group. If the product of the current number of instances and the percentage is not an integer, AS automatically rounds the value up or down:<ul style="list-style-type: none">– If the value is greater than 1, it is rounded down. For example, value 12.7 is rounded down to 12.– If the value is greater than 0 but less than 1, it is rounded up to 1. For example, value 0.67 is rounded up to 1. <p>For example, there are 10 instances in an AS group, and the scaling action is Add 15% instances. When the AS policy is triggered, AS calculates the number of instances to be added is 1.5 and rounds 1.5 down to 1. After the scaling action is complete, there are 11 instances in the AS group.</p>

Parameter	Description	Example Value
Cooldown Period	<p>To prevent an alarm-based policy from being repeatedly triggered by the same event, you can set a cooldown period.</p> <p>A cooldown period (in seconds) is the period of time between two scaling actions. AS recounts the cooldown period after a scaling action is complete.</p> <p>During the cooldown period, AS denies all scaling requests triggered by alarm-based policies. Scaling requests triggered manually or by scheduled or periodic policies are not affected.</p> <p>For example, suppose that the cooldown period is set to 300 seconds (5 minutes), and a scheduled policy is specified to trigger a scaling action at 10:32, and a previous scaling action triggered by an alarm policy ends at 10:30. Any alarm-triggered scaling action will then be denied during the cooldown period from 10:30 to 10:35, but the scaling action scheduled for 10:32 will still take place. If the scheduled scaling action ends at 10:36, a new cooldown period starts at 10:36 and ends at 10:41.</p>	300

Parameter	Description	Example Value
	<p>NOTE</p> <ul style="list-style-type: none"> • If a scaling action is triggered by an AS policy, the cooldown period is whatever configured for that AS policy. • If a scaling action is triggered by manually changing the expected number of instances or by other actions, the cooldown period is whatever configured for the AS group. The default cooldown period is 300 seconds. • When an AS group scales out, scale-in requests triggered manually or by scheduled or periodic policies will be executed immediately after the scale-out is complete, without being affected by the cooldown period. • When an AS group scales in, scale-out requests triggered manually or by scheduled or periodic policies will be executed immediately after the scale-in is complete, without being affected by the cooldown period. 	

Table 3-2 System monitoring parameters

Parameter	Description	Example Value
Rule Name	Specifies the name of the alarm rule.	as-alarm-7o1u
Monitoring Type	Specifies the type of monitoring metrics, which can be System monitoring or Custom monitoring . Select System monitoring .	System monitoring
Trigger Condition	Select monitoring metrics supported by AS and set alarm conditions for the metrics.	CPU Usage Max. >70%
Monitoring Interval	Specifies the interval at which the alarm status is updated based on the alarm rule.	5 minutes

Parameter	Description	Example Value
Consecutive Occurrences	Specifies the number of sampling points when an alarm is triggered. If Consecutive Occurrences is set to n , the sampling points of the alarm rule are the sampling points in n consecutive sampling periods. The alarm rule status does not change to Alarm unless all sampling points breach the threshold configured by the alarm rule.	3

Table 3-3 Custom monitoring parameters

Parameter	Description	Example Value
Rule Name	Specifies the name of the alarm rule.	as-alarm-7o1u
Monitoring Type	Select Custom monitoring . Custom monitoring meets monitoring requirements in various scenarios.	Custom monitoring
Resource Type	Specifies the name of the service for which the alarm rule is configured.	AGT.ECS
Dimension	Specifies the metric dimension of the alarm rule.	instance_id
Monitored Object	Specifies the resources to which the alarm rule applies.	N/A
Trigger Condition	Select monitoring metrics supported by AS and set alarm conditions for the metrics.	CPU Usage Max. >70%
Monitoring Interval	Specifies the interval at which the alarm status is updated based on the alarm rule.	5 minutes
Consecutive Occurrences	Specifies the number of sampling points when an alarm is triggered. If Consecutive Occurrences is set to n , the sampling points of the alarm rule are the sampling points in n consecutive sampling periods. The alarm rule status does not change to Alarm unless all sampling points breach the threshold configured by the alarm rule.	3

6. Click **OK**.

The newly added AS policy is displayed on the **AS Policy** tab. In addition, the AS policy is in **Enabled** state by default.

Creating a Scheduled or Periodic Policy

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
3. Locate the row containing the target AS group and click **View AS Policy** in the **Operation** column.
4. On the **AS Policies** page, click **Add AS Policy**.
5. Configure the parameters listed in [Table 3-4](#).

Table 3-4 Parameter description

Parameter	Description	Example Value
Policy Name	Specifies the name of the AS policy to be created.	as-policy-p6g5
Policy Type	Select Scheduled or Periodic for expanding resources at a specified time. If you select Periodic , you are required to configure two more parameters: <ul style="list-style-type: none"> • Period <ul style="list-style-type: none"> - Day - Week - Month • Time Range Specifies the time range during which the AS policy can be triggered. 	Day 2023/03/01 00:00:00 - 2023/03/31 23:59:59 In this example, the AS policy will trigger a scaling action every day in March, and will become invalid from April 1, 2023 00:00:00.
Time Zone	The default value is GMT +08:00 . GMT+08:00 is 8:00 hours ahead of Greenwich Mean Time.	GMT+08:00
Triggered At	Specifies the time at which the AS policy is triggered. NOTE The selected triggering time must fall inside the effective time range of the policy.	N/A

Parameter	Description	Example Value
Scaling Action	<p>Specifies an action and the number of instances.</p> <p>The following scaling action options are available:</p> <ul style="list-style-type: none"> • Add Adds instances to an AS group when the scaling action is performed. • Reduce Removes instances from an AS group when the scaling action is performed. • Set to Sets the expected number of instances in an AS group to a specified value. 	<ul style="list-style-type: none"> • Add 1 instance • Add 10% instances The number of instances to be added is 10% of the current number of instances in the AS group. If the product of the current number of instances and the percentage is not an integer, AS automatically rounds the value up or down: <ul style="list-style-type: none"> - If the value is greater than 1, it is rounded down. For example, value 12.7 is rounded down to 12. - If the value is greater than 0 but less than 1, it is rounded up to 1. For example, value 0.67 is rounded up to 1. <p>For example, suppose that there are 10 instances in an AS group, and the scaling action is Add 15% instances. When the AS policy is triggered, AS calculates the number of instances to be added is 1.5 and rounds 1.5 down to 1. After the scaling action is complete, there are 11 instances in the AS group.</p>

Parameter	Description	Example Value
Cooldown Period	<p>To prevent an alarm-based policy from being repeatedly triggered by the same event, you can set a cooldown period.</p> <p>A cooldown period (in seconds) is the period of time between two scaling actions. AS recounts the cooldown period after a scaling action is complete.</p> <p>During the cooldown period, AS denies all scaling requests triggered by alarm-based policies. Scaling requests triggered manually or by scheduled or periodic policies are not affected.</p> <p>For example, suppose that the cooldown period is set to 300 seconds (5 minutes), and a scheduled policy is specified to trigger a scaling action at 10:32, and a previous scaling action triggered by an alarm policy ends at 10:30. Any alarm-triggered scaling action will then be denied during the cooldown period from 10:30 to 10:35, but the scaling action scheduled for 10:32 will still take place. If the scheduled scaling action ends at 10:36, a new cooldown period starts at 10:36 and ends at 10:41.</p>	300

Parameter	Description	Example Value
	<p>NOTE</p> <ul style="list-style-type: none"> • If a scaling action is triggered by an AS policy, the cooldown period is whatever configured for that AS policy. • If a scaling action is triggered by manually changing the expected number of instances or by other actions, the cooldown period is whatever configured for the AS group. The default cooldown period is 300 seconds. • When an AS group scales out, scale-in requests triggered manually or by scheduled or periodic policies will be executed immediately after the scale-out is complete, without being affected by the cooldown period. • When an AS group scales in, scale-out requests triggered manually or by scheduled or periodic policies will be executed immediately after the scale-in is complete, without being affected by the cooldown period. 	

6. Click **OK**.

The newly added AS policy is displayed on the **AS Policy** tab. In addition, the AS policy is in **Enabled** state by default.

 **NOTE**

If you have created scheduled or periodic AS policies that are invoked at the same time, AS will execute the one created later. This constraint does not apply to alarm-triggered AS policies.

3.3 Managing AS Policies

Scenarios

An AS policy specifies the conditions for triggering a scaling action as well as the operation that will be performed. If the conditions are met, a scaling action is triggered automatically.

This section describes how to manage an AS policy, including modifying, enabling, disabling, executing, and deleting an AS policy.

Modifying an AS Policy

If a particular AS policy cannot meet service requirements, you can modify the parameter settings of the policy.

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. Locate the row containing the target AS group and click **View AS Policy** in the **Operation** column. On the displayed page, locate the row containing the target AS policy and choose **More > Modify** in the **Operation** column.
4. In the displayed **Modify AS Policy** dialog box, modify the parameters and click **OK**.

Enabling an AS Policy

An AS policy can trigger scaling actions only when it and the AS group are both enabled. You can enable one or more AS policies for an AS group as required.

- Before enabling multiple AS policies, ensure that the AS policies do not conflict with one another.
- An AS policy can be enabled only when its status is **Disabled**.

Locate the row containing the target AS group and click **View AS Policy** in the **Operation** column. On the displayed page, locate the row containing the target AS policy and click **Enable** in the **Operation** column. To concurrently enable multiple AS policies, select these AS policies and click **Enable** in the upper part of the AS policy list.

Disabling an AS Policy

If you do not want a particular AS policy to trigger any scaling actions within a specified period of time, you can disable it.

- If all of the AS policies configured for an AS group are disabled, no scaling action will be triggered for this AS group. However, if you manually change the value of **Expected Instances**, a scaling action will still be triggered.
- You can disable an AS policy only when its status is **Enabled**.

Locate the row containing the target AS group and click **View AS Policy** in the **Operation** column. On the displayed page, locate the row containing the target AS policy and click **Disable** in the **Operation** column. To concurrently disable multiple AS policies, select these AS policies and click **Disable** in the upper part of the AS policy list.

Manually Executing an AS Policy

You can make the number of instances in an AS group reach the expected number of instances immediately by manually executing an AS policy.

- You can manually execute an AS policy if the scaling conditions configured in the AS policy are not met.
- You can manually execute an AS policy only when the AS group and AS policy are both in **Enabled** state.

Locate the row containing the target AS group and click **View AS Policy** in the **Operation** column. On the displayed page, locate the row containing the target AS policy and click **Execute Now** in the **Operation** column.

 **NOTE**

If **Policy Type** is set to **Alarm** and **Alarm Policy Type** to **Refined scaling**, the scaling policy cannot be executed immediately.

Deleting an AS Policy

You can delete an AS policy that will not be used for triggering scaling actions.

An AS policy can be deleted even when the scaling action triggered by the policy is in progress. Deleting the AS policy does not affect a scaling action that has already started.

Locate the row containing the target AS group and click **View AS Policy** in the **Operation** column. On the displayed page, locate the row containing the target AS policy and choose **More > Delete** in the **Operation** column.

To concurrently delete multiple AS policies, select these AS policies and click **Delete** in the upper part of the AS policy list.

4 Scaling Action

4.1 Dynamic Scaling

Before using AS to perform scaling actions, you must specify how to perform the scaling actions to dynamically expand resources.

If the demands change frequently, you can configure alarm-based policies to scale resources. When the conditions for invoking an AS policy are met, AS automatically changes the expected number of instances to trigger a scaling action to scale up or down resources. For details about how to create an alarm policy, see [Creating an AS Policy](#).

Consider a train ticket booking application. If the CPU usage of the instances that run the application goes up to 90%, an instance needs to be added to ensure that services run properly. If the CPU usage drops down to 30%, an instance needs to be deleted to prevent resource waste. To meet the requirements, you can configure two alarm policies. One policy is used to add one instance if the maximum CPU usage exceeds 90%. For details, see [Figure 4-1](#). The other policy is used to remove an instance if the minimum CPU usage drops below 30%. For details, see [Figure 4-2](#).

Figure 4-1 Alarm policy 01

Add AS Policy

Policy Name:

Policy Type: Alarm Scheduled Periodic
Policies of this type are applied only when their associated alarm rules are enabled. [View alarm rules](#)

Alarm Rule: Create Use existing

Rule Name:

Monitoring Type: System monitoring Custom monitoring

Trigger Condition: %
The metrics that can be monitored vary somewhat by OS. [Learn more](#)
To select a metric whose name starts with (Agent), make sure the Agent has been installed on all instances in the AS group. [Learn more](#)

Monitoring Interval:

Consecutive Occurrences: ?

Alarm Policy Type: Simplified scaling Refined scaling

Scaling Action:

Cooldown Period (s): ?

Figure 4-2 Alarm policy 02

Add AS Policy

Policy Name:

Policy Type: Alarm Scheduled Periodic
Policies of this type are applied only when their associated alarm rules are enabled. [View alarm rules](#)

Alarm Rule: Create Use existing

Rule Name:

Monitoring Type: System monitoring Custom monitoring

Trigger Condition: %
The metrics that can be monitored vary somewhat by OS. [Learn more](#)
To select a metric whose name starts with (Agent), make sure the Agent has been installed on all instances in the AS group. [Learn more](#)

Monitoring Interval:

Consecutive Occurrences: ?

Alarm Policy Type Simplified scaling Refined scaling

Scaling Action Reduce 1 instances

Cooldown Period (s) 900 ?

OK Cancel

4.2 Scheduled Scaling

To satisfy demands that change regularly, you can configure a scheduled or periodic policy to scale resources at specified time or periodically. For details about how to create a scheduled or periodic policy, see [Creating an AS Policy](#).

Take an online course selection web application as an example. This application is frequently used when a semester starts and seldom used during other parts of the year. You can configure two scheduled policies to scale resources at the beginning of each semester. The first policy is used to add an instance when the course selection starts, and the second policy is used to remove an instance when the course selection ends.

4.3 Manual Scaling

Scenarios

You can change the size of an AS group manually. You can either add or remove instances to or from the AS group, or modify the expected number of instances of the AS group.

Procedure

Adding an instance to an AS group

Before you add an instance to an AS group, ensure that the conditions below are met.

Table 4-1 Conditions for manually adding an instance to an AS group

Item	Condition
AS group	<ul style="list-style-type: none"> The AS group is in the Enabled status. The AS group does not have ongoing scaling actions. The number of instances to be added plus the expected number of instances cannot exceed the maximum number of instances of the AS group.
Instance	<ul style="list-style-type: none"> The instance to be added is not a member of another AS group. The instance is in the same VPC as the AS group.

 **NOTE**

- A maximum of 10 instances can be added to an AS group at a time.
- If the AS group has an attached load balancer, the instances will be associated with the load balancer.

To add instances to an AS group, perform the following steps:

1. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
2. Click the **AS Groups** tab and then the name of the target AS group.
3. On the AS group details page, click the **Instances** tab and then **Add**.
4. Select the instances to be added and click **OK**.

Removing an instance from an AS group

You can remove an instance from an AS group, update the instance or fix an instance fault, and add the instance back to the AS group. After the instance is removed from the AS group, it no longer processes any application traffic.

For example, you can change AS configuration for an AS group at any time. New instances will be created using the new configuration, but existing instances in the AS group are not affected. To update the existing instances, you can stop them so that they can be replaced automatically. You can also remove the instances from the AS group, update them, and then add them back to the AS group.

When you remove instances from an AS group, consider the restrictions below.

Table 4-2 Constraints on manually removing an instance from an AS group

Item	Constraint
AS group	<ul style="list-style-type: none"> • The AS group is in the Enabled status. • The AS group does not have ongoing scaling actions.
Instance	<ul style="list-style-type: none"> • The instances are in the Enabled lifecycle status. • The instances are not used by SDRS.

 **NOTE**

- A maximum of 50 instances can be removed from to an AS group at a time.
- If the number of instances you are removing decreases the number of instances in the AS group below the minimum number of instances allowed, AS launches new instances to maintain the expected capacity.
- If you remove instances from an AS group that has an associated load balancer, the instances will be dissociated from the load balancer.

To remove an instance from an AS group, perform the following steps:

1. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
2. Click the **AS Groups** tab and then the name of the target AS group.

3. Click the **Instances** tab, locate the row containing the desired instance, and click **Remove** or **Remove and Delete** in the **Operation** column.

To remove multiple instances from the AS group, select the check boxes in front of them and click **Remove** or **Remove and Delete**.

To remove all instances from the AS group, select the check box on the left of **Name** and click **Remove** or **Remove and Delete**.

NOTE

- If the instances you want to remove were automatically added to the AS group, they are billed on a pay-per-use basis by default. You can:
 - Remove the instances from the AS group by choosing **Remove**.
 - Remove the instances from the AS group and delete them by choosing **Remove and Delete**.
- If you have changed the billing mode of these instances from pay-per-use to yearly/monthly, you can only remove them from the AS group even when you choose **Remove and Delete**.
- If the instances were manually added to the AS group, they can only be removed. They cannot be removed and deleted.

Changing the expected number of instances

Manually change the expected number of instances to add or reduce the number of instances in an AS group for expanding resources.

For details, see [Modifying an AS Group](#).

4.4 Configuring an Instance Removal Policy

When instances are automatically removed from your AS group, the instances that are not in the currently used AZs will be removed first. Then the instance removal policy you select will be applied.

AS supports the following instance removal policies:

- **Oldest instance:** The oldest instance is removed from the AS group first. Use this policy if you want to upgrade instances in an AS group to a new ECS type. You can gradually replace instances of the old type with instances of the new type.
- **Newest instance:** The newest instance is removed from the AS group first. Use this policy if you want to test a new AS configuration but do not want to keep it in production.
- **Oldest instance created from oldest AS configuration:** The oldest instance created from the oldest configuration is removed from the AS group first. Use this policy if you want to update an AS group and phase out the instances created from a previous AS configuration.
- **Newest instance created from oldest AS configuration:** The newest instance created from the oldest configuration is removed from the AS group first.

 NOTE

Manually added instances are the last to be removed, and if AS does remove a manually added instance, it only removes the instance. It does not delete the instance. If multiple manually added instances must be removed, AS preferentially removes the earliest-added instance first.

4.5 Viewing a Scaling Action

Scenarios

This section describes how to check whether a scaling action has been performed and how to view scaling action details.

Viewing Monitoring Data


The following steps illustrate how to view scaling actions of an AS group.

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
3. Click the **AS Groups** tab and then the name of the target AS group.
4. Click the **Monitoring** tab and view scaling actions. On the **Monitoring** page, you can view changes in the number of instances and metrics such as CPU Usage.

Viewing Historical Scaling Actions

The following steps illustrate how to view the historical records of scaling actions of an AS group.

1. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**.
2. Click the **AS Groups** tab and then the name of the target AS group.
3. Click the **Scaling Actions** tab. This page displays historical scaling actions of an AS group, including instance scaling and load balancer migration.

Scaling Action ID, **Status**, **Scaling Action Type**, **Description**, **Start Time**, and **End Time** of scaling actions are displayed. Click  before the scaling action ID to view the resource name, status, and failure cause. You can also use the filtering function in the upper right corner to view scaling actions in a specified period.

4.6 Managing Lifecycle Hooks

Lifecycle hooks enable you to flexibly control addition and removal of ECS instances in AS groups and manage the lifecycle of ECS instances in AS groups. [Figure 4-3](#) shows the instance lifecycle when no lifecycle hook is added to an AS group.

Figure 4-3 Instance lifecycle when no lifecycle hook is added to an AS group

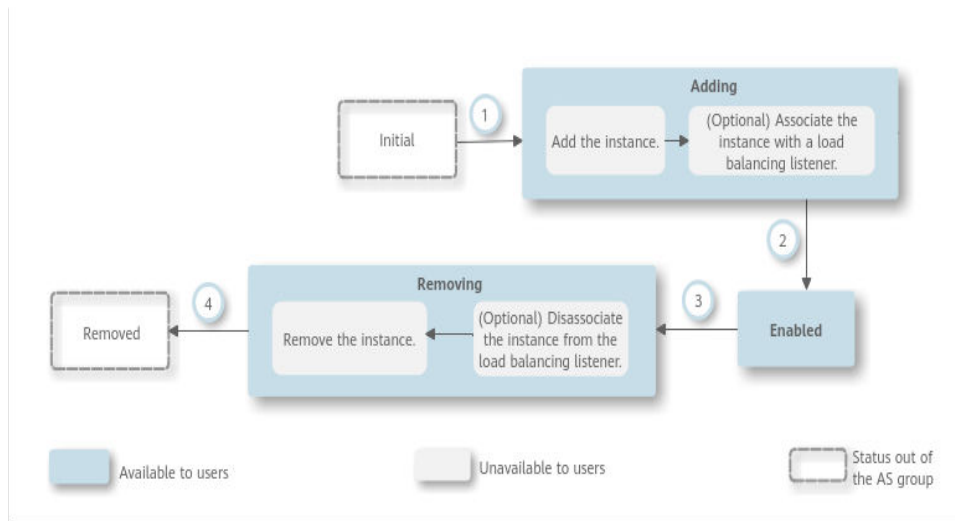
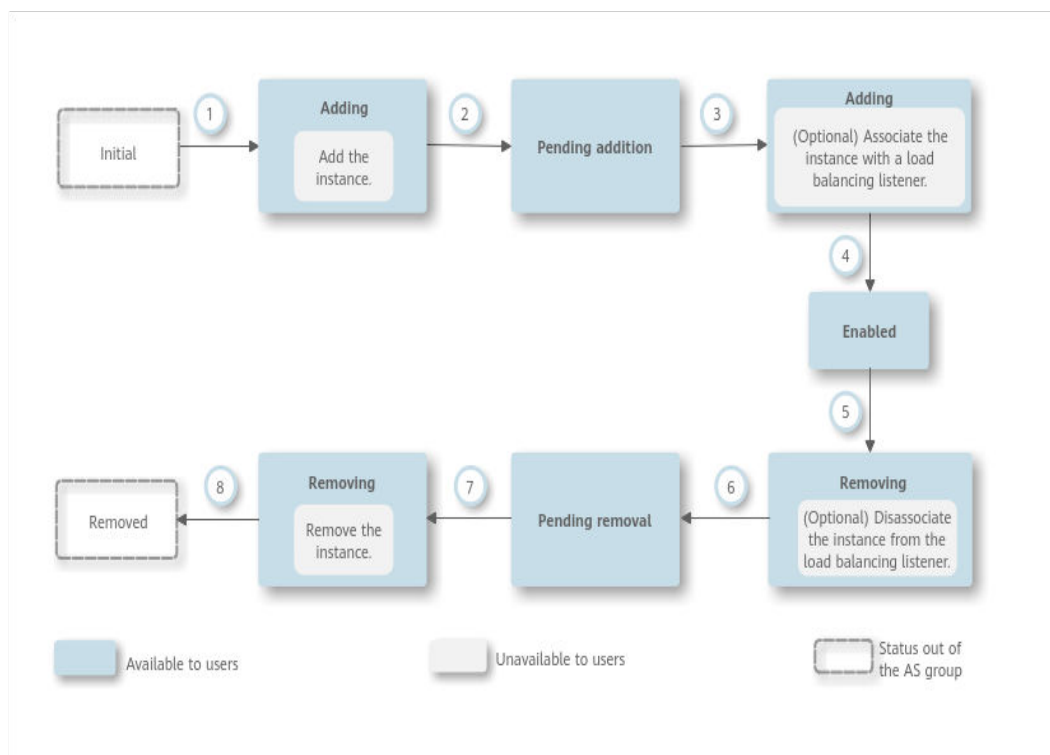


Figure 4-4 shows the instance lifecycle when a lifecycle hook is added to an AS group.

Figure 4-4 Instance lifecycle when a lifecycle hook is added to an AS group



When the AS group scales in or out, the added lifecycle hooks are triggered, the scaling action is suspended, and the instance being added or removed is put into a wait state, as shown in 2 and 6 in **Figure 4-4**. During this period of time, you can perform some custom operations on the instance. For example, you can install or configure software on an instance being added to the AS group. A suspended scaling action will be resumed if either of the following occurs:

- The timeout duration ends.
Assume that you have set the timeout period to 3,600s by referring to section [Table 4-3](#). The suspended scaling action will be automatically resumed if the timeout duration (3,600s) ends.
- A callback action is performed to move the instance out of the wait state. For details, see [Performing a Callback Action](#).

Application Scenarios

- Instances newly added to an AS group need to be initialized before they are bound to a load balancer listener. Initialization means the software is installed and configured and the instance is fully ready to accept traffic.
- To remove an instance from an AS group, it needs to be first unbound from the load balancer listener, stops accepting new requests, and finishes processing any accepted requests.
- Before instances are removed from an AS group, you may need to back up data or download logs.
- Other scenarios where custom operations need to be performed

How Lifecycle Hooks Work

After you add lifecycle hooks to an AS group, they work as follows:

- Adding an ECS instance to an AS group
When an instance is initialized and added to an AS group, a lifecycle hook of the **Instance adding** type is automatically triggered. The instance enters the **Pending addition** state, that is, the instance is suspended by the lifecycle hook. If you have configured a notification object, the system sends a message to the object. After receiving the message, you can perform custom operations, for example, installing software on the instance. The instance remains in a wait state until you complete the custom operations and perform a callback action (see [Performing a Callback Action](#)) or the timeout duration ends. After the instance moves out of a wait state, the specified default callback action will take place.
 - **Continue:** The instance will be added to the AS group.
 - **Abandon:** The instance will be deleted and a new instance will be created.

If you have configured multiple **Instance adding** lifecycle hooks, all of them will be triggered when an instance is added to the AS group. If the default callback action of any lifecycle hook is **Abandon**, the instance will be deleted and a new instance will be created. If the default callback action of all lifecycle hooks is **Continue**, the instance is added to the AS group after suspension by the last lifecycle hook is complete.

- Removing an instance from an AS group
When an instance is removed from an AS group, the instance enters the **Removing** state. After a lifecycle hook is triggered, the instance enters the **Pending removal** state. The system sends messages to the configured notification object. After receiving the message, you can perform custom operations, such as uninstalling software and backing up data. The instance remains in the wait state until you finish the custom operations and perform

the default callback operation or the timeout duration ends. After the instance moves out of a wait state, the specified default callback action will take place.

- **Continue:** The instance is removed from the AS group.
- **Abandon:** The instance is removed from the AS group.

If you have configured multiple lifecycle hooks, and the default callback action of all lifecycle hooks is **Continue**, the instance will be removed from the AS group until suspension by the remaining lifecycle hooks time out. If the default callback action of any lifecycle hook is **Abandon**, the instance will be directly removed from the AS group.

Constraints

- You can add, modify, or delete a lifecycle hook when the AS group does not perform a scaling action.
- Up to five lifecycle hooks can be added to one AS group.

Adding a Lifecycle Hook

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. Click the name of the AS group to which the lifecycle hook is to be added. On the AS group details page, click the **Lifecycle Hooks** tab and then **Add Lifecycle Hook**.
4. In the displayed **Add Lifecycle Hook** dialog box, set the parameters listed in [Table 4-3](#).

Table 4-3 Parameter description

Parameter	Description	Example Value
Hook Name	Specifies the lifecycle hook name. The name can contain letters, digits, underscores (_), and hyphens (-), and cannot exceed 32 characters.	we12_w
Hook Type	Specifies the lifecycle hook type. The value can be Instance adding or Instance removal . Instance adding puts an instance that is being added to an AS group to Pending addition state. Instance removal puts an instance that is being removed from an AS group to Pending removal state.	Instance adding

Parameter	Description	Example Value
Default Callback Action	<p>Specifies the action that the system takes when an instance moves out of a wait state.</p> <p>The default callback action for an Instance adding lifecycle hook can be Continue or Abandon:</p> <ul style="list-style-type: none"> • Continue: If multiple lifecycle hooks are configured for the AS group, and the default callback action of all the hooks is Continue, the system will continue to add the instance to the AS group until the all lifecycle hooks time out. • Abandon: If multiple lifecycle hooks are configured for the AS group, and the default callback action of one lifecycle hook is Abandon, the system will delete the instance and create another one without waiting for the remaining lifecycle hooks to time out. <p>The default callback action for an Instance removal lifecycle hook can be Continue or Abandon:</p> <ul style="list-style-type: none"> • Continue: If only one lifecycle hook is configured for the AS group, the system will remove the instance from the AS group. If multiple lifecycle hooks are configured for the AS group, and the default callback actions of all the hooks are Continue, the system will continue to remove the instance from the AS group until all lifecycle hooks time out. • Abandon: If multiple lifecycle hooks are configured for the AS group, and the default callback action of one lifecycle hook is Abandon, the system will continue to remove the instance from the AS group without waiting for the remaining lifecycle hooks to time out. 	Continue
Timeout Duration (s)	<p>Specifies the amount of time for the instances to remain in a wait state. The value ranges from 60s to 86,400s.</p> <p>You can extend the timeout duration or perform a Continue or Abandon action before the timeout duration ends. For more information about callback actions, see Performing a Callback Action.</p>	3600

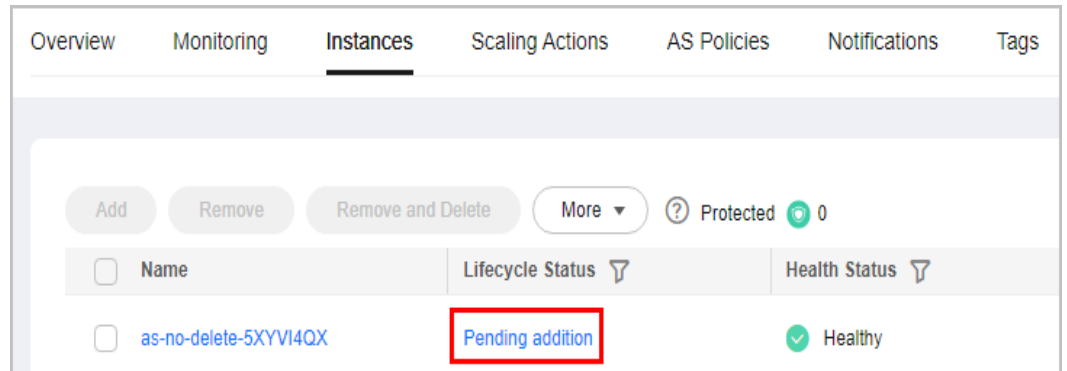
Parameter	Description	Example Value
Notification Topic	<p>Specifies a notification object for a lifecycle hook. For details, see "Creating a Topic" in Simple Message Notification User Guide. When an instance is suspended by the lifecycle hook, the system sends a notification to the object. This notification contains the basic instance information, your custom notification content, and the token for controlling lifecycle actions. An example notification is as follows:</p> <pre>{ "service": "AutoScaling", "tenant_id": "93075aa73f6a4fc0a3209490cc57181a", "lifecycle_hook_type": "INSTANCE_LAUNCHING", "lifecycle_hook_name": "test02", "lifecycle_action_key": "4c76c562-9688-45c6-b685-7fd732df310a", "notification_metadata": "xxxxxxxxxxxx", "scaling_instance": { "instance_id": "89b421e4-5fa6-4733-bf40-6b07a8657256", "instance_name": "as-config-kxeg_RM6OCREY", "instance_ip": "192.168.0.202" }, "scaling_group": { "scaling_group_id": "fe376277-50a6-4e36-bdb0-685da85f1a82", "scaling_group_name": "as-group-wyz01", "scaling_config_id": "16ca8027-b6cc-45fc-af2d-5a79996f685d", "scaling_config_name": "as-config-kxeg" } }</pre>	N/A
Notification Message	After a notification object is configured, the system sends your custom notification to the object.	N/A

5. Click **OK**.
The added lifecycle hook is displayed on the **Lifecycle Hooks** page.

Performing a Callback Action

1. On the **AS Groups** page, click the name of the target AS group.
2. On the displayed page, click the **Instances** tab.
3. Locate the instance that has been suspended by the lifecycle hook and click Pending addition or Pending removal in the **Lifecycle Status** column, as shown in [Figure 4-5](#).

Figure 4-5 Performing a callback action

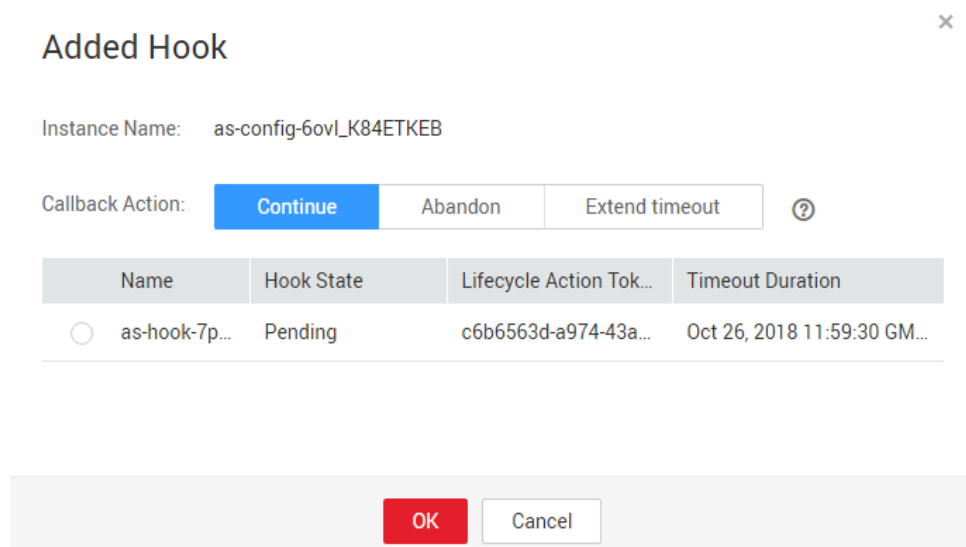


NOTE

Callback actions can only be performed on instances that have been suspended by a lifecycle hook.

4. In the displayed **Added Hook** dialog box, view the suspended instance and all the lifecycle hooks, and perform callback actions on lifecycle hooks.

Figure 4-6 Added Hook dialog box



Callback actions include:

- **Continue**
- **Abandon**
- **Extend timeout**

If you have performed custom operations before the timeout duration ends, select **Continue** or **Abandon** to complete the lifecycle actions. For details about **Continue** and **Abandon**, see [Table 4-3](#). If you need more time for custom operations, select **Extend timeout** to extend the timeout duration. Then, the timeout duration will be extended by 3600 seconds each time.

Modifying a Lifecycle Hook

On the **Lifecycle Hooks** page, locate the target lifecycle hook and click **Modify** in the **Operation** column, see [Table 4-3](#) for parameters. You can modify the parameter except **Hook Name**, such as **Hook Type**, **Default Callback Action**, and **Timeout Duration**.

Deleting a Lifecycle Hook

On the **Lifecycle Hooks** page, locate the target lifecycle hook and click **Delete** in the **Operation** column.

4.7 Configuring Instance Protection

Scenarios

To control whether an instance can be removed automatically from an AS group, use instance protection. Once configured, when AS automatically scales in the AS group, the instance that is protected will not be removed.

Prerequisites

Instance protection does not protect instances from the following:

- Health check replacement if the instance fails health checks
- Manual removal

NOTE

- Instance protection does not protect unhealthy instances because such instances cannot provide services.
- By default, instance protection does not take effect on the ECSs that are newly created in or added to an AS group.
- If an instance is removed from an AS group, its instance protection setting is lost.

Enabling Instance Protection

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. Click the name of the target AS group.
4. Click the **Instances** tab. Select one or more instances and choose **Enable Instance Protection** from the **More** drop-down list. In the displayed **Enable Instance Protection** dialog box, click **Yes**.

Disabling Instance Protection

1. Log in to the management console.
1. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
2. Click the name of the target AS group.

3. Click the **Instances** tab. Select one or more instances and choose **Disable Instance Protection** from the **More** drop-down list. In the displayed **Disable Instance Protection** dialog box, click **Yes**.

4.8 Putting an Instance Into the Standby State

If you want to stop distributing traffic to some instances in your AS group but do not want to remove them from the AS group, you can put the instances on standby. You can put one or more instances in your AS group on standby, and then stop or restart these instances without worrying about they are removed from the AS group.

Application Scenarios

You cannot control the lifecycle of ECS instances in an AS group. The AS group removes unhealthy instances and does not allow you to stop or restart these instances. As a result, some ECS functions are unavailable. For example, you cannot reset the password, or reinstall or change the OS of these instances.

By putting ECS instances into standby state, you can control their lifecycle and perform operations on them as needed, such as stopping the instances. This facilitates management of instances in your AS group and is helpful in a number of different scenarios.

- If you want to change the OS of an ECS added by a scaling action or stop the ECS, you can set the ECS to standby mode. Then you can perform all operations supported by the ECS service. After completing the operations, cancel standby mode for the ECS.

For example, you can change the AS configuration for your AS group at any time. This configuration will be used by any instance that is created in the AS group. However, the AS group does not update instances that are running. You can stop these instances, and the AS group will replace them. Alternatively, you can set the instances to standby mode, update software on them, and then cancel standby mode for them.

- If an instance in your AS group associated with a load balancer becomes faulty, you can set the instance to standby mode, after which the load balancer will no longer distribute access traffic to the instance. Then you can log in to the instance, locate and rectify the fault, and restart the instance. After the instance recovers, cancel standby mode for the instance to receive traffic again.

Working Rules

- Set instances to standby mode.

After you set an instance to standby mode, the instance will be automatically unbound from the load balancer associated with the AS group. The instance is still in the AS group, but no health check will be performed on the instance. In this case, load on other instances will increase. To reduce load on other instances and ensure proper service running, you can select **Add the same number of new instances to the AS group** when setting the instance to standby mode.

 **NOTE**

- An instance can be set to standby mode only when the instance is enabled and the AS group has no ongoing scaling action.
- Scaling actions will not remove standby instances from the AS group.
- You can manually remove standby instances from the AS group.
- Cancel standby mode for instances.

After you cancel standby mode for an instance, it will be in running state and receive traffic again. If a load balancer is associated with the AS group, the instance will be bound to the load balancer. After the instance starts running properly, health check will be performed on it again.

 **NOTE**

Standby mode can be canceled for an instance only when the instance is in standby mode and the AS group has no ongoing scaling action.

Setting Instances to Standby Mode

1. Log in to the management console.
2. Click **Service List**.
3. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
4. Click the name of the target AS group.
5. Click the **Instances** tab. Select one or more instances, click **More**, and select **Set to Standby** in the drop-down list. In the displayed dialog box, select **Add the same number of new instances to the AS group** as you need and click **Yes**.

Canceling Standby Mode for Instances

1. Log in to the management console.
1. Click **Service List**.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. Click the name of the target AS group.
4. Click the **Instances** tab. Select one or more instances, click **More**, and select **Cancel Standby** in the drop-down list. In the displayed dialog box, click **Yes**.

5 Bandwidth Scaling

5.1 Creating a Bandwidth Scaling Policy

Scenarios

You can automatically adjust your purchased EIP bandwidth and shared bandwidth using a bandwidth scaling policy. This section describes how to create a bandwidth scaling policy.

When creating a bandwidth scaling policy, you need to configure basic information. The system supports three types of bandwidth scaling policies: alarm-based, scheduled, and periodic.

The basic information for creating a bandwidth scaling policy includes the policy name, resource type, policy type, and trigger condition.

Creating an Alarm-based Bandwidth Scaling Policy

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Bandwidth Scaling**.
3. Click **Create Bandwidth Scaling Policy**.
4. Set parameters, such as the policy name, policy type, and trigger condition. For details, see [Table 5-1](#).

Table 5-1 Alarm policy parameters

Parameter	Description	Example Value
Region	Specifies the region where the AS group resides.	N/A
Policy Name	Specifies the name of the bandwidth scaling policy. The name consists of only letters, digits, underscores (_), and hyphens (-).	N/A

Parameter	Description	Example Value
Resource Type	Specifies the type of the bandwidth to be adjusted. You can select EIP or Shared bandwidth .	EIP
EIP	Specifies the public network IP address whose bandwidth needs to be scaled. NOTE Only bandwidths of pay-per-use EIPs can be scaled.	N/A
Policy Type	Select Alarm .	Alarm
Alarm Rule	You can use an existing alarm rule or create a new one. Alternatively, click Create Alarm Rule on the right side of the Rule Name parameter and create an alarm rule on the Alarm Rules page. For details, see Creating an Alarm Rule . To create an alarm rule, configure the following parameters: <ul style="list-style-type: none">● Rule Name Specifies the name of the new alarm rule, for example, as-alarm-7o1u.● Trigger Condition Select a monitoring metric and trigger condition based on the metric. Table 5-2 lists the supported monitoring metrics. An example value is Outbound Traffic Avg. > 100 bit/s.● Monitoring Interval Specifies the period for the metric, for example, 5 minutes.● Consecutive Occurrences Specifies the number of consecutive periods in which the triggering condition is met for triggering a scaling action.	N/A

Parameter	Description	Example Value
Scaling Action	<p>Specifies the execution action in the AS policy. The following scaling action options are available:</p> <ul style="list-style-type: none"> • Add When a scaling action is triggered, the bandwidth is increased. • Reduce When a scaling action is triggered, the bandwidth is decreased. • Set to The bandwidth is set to a fixed value. <p>NOTE The step (minimum unit for bandwidth adjustment) varies depending on the bandwidth value range. The bandwidth will be automatically adjusted to the nearest value according to the actual step.</p> <ul style="list-style-type: none"> • If the bandwidth is less than or equal to 300 Mbit/s, the default step is 1 Mbit/s. • If the bandwidth ranges from 300 Mbit/s to 1000 Mbit/s, the default step is 50 Mbit/s. • If the bandwidth is greater than 1000 Mbit/s, the default step is 500 Mbit/s. 	N/A
Cooldown Period	<p>A cooldown period (in seconds) is the period of time between two scaling actions. AS recounts the cooldown period after a scaling action is complete. During the cooldown period, AS denies all scaling requests triggered by alarm-based policies. Scaling requests triggered manually or by scheduled or periodic policies are not affected.</p>	300s

Table 5-2 Monitoring metrics supported by the alarm policy

Metric	Description
Inbound Bandwidth	Indicates the network rate of inbound traffic.
Inbound Traffic	Indicates the network traffic going out of the cloud platform.
Outbound Bandwidth	Indicates the network rate of outbound traffic.
Outbound Traffic	Indicates the network traffic going out of the cloud platform.
Outbound Bandwidth Usage	Indicates the usage of network rate of outbound traffic in the unit of percentage.

- After setting the parameters, click **Create Now**.
The newly created bandwidth scaling policy is displayed on the **Bandwidth Scaling** page and is in **Enabled** state by default.

Creating an Alarm Rule

When creating an alarm-based bandwidth scaling policy, you can click **Create Alarm Rule** to the right of **Rule Name** to create an alarm rule. To do so, perform the following operations:

- Click **Create Alarm Rule** to the right of **Rule Name** to switch to the **Alarm Rules** page of Cloud Eye.
- On the **Alarm Rules** page, click **Create Alarm Rule** in the upper right corner.
- Set parameters based on [Figure 5-1](#) and [Table 5-3](#). For more information about how to set alarm rules, see [Cloud Eye User Guide](#).

Figure 5-1 Creating an alarm rule


The screenshot displays the 'Create Alarm Rule' configuration interface in the Cloud Eye console. Key elements include:

- Resource Type:** Elastic IP and Bandwidth...
- Dimension:** Bandwidths
- Monitoring Scope:** Specific resources
- Method:** Create manually
- Alarm Policy:** Outbound Bandwidth, Max., 5 minutes, 3 consecutive, ≥, 500 bit/s
- Alarm Severity:** Major
- Alarm Notification:** (Input field)

A line graph shows the bandwidth usage for 'ecs-transitvp-band...' over time, with a peak around 48,454 bit/s. The graph includes a red vertical line at 500 bit/s, indicating the alarm threshold.

Table 5-3 Key parameters for creating an alarm rule

Parameter	Description	Example Value
Name	Specifies the name of the alarm rule.	alarm-bandwidth
Description	(Optional) Provides supplementary information about the alarm rule.	N/A
Enterprise Project	Specifies the enterprise project the alarm rule belongs to. Only users with the enterprise project permissions can view and manage the alarm rule.	default
Resource Type	Specifies the name of the service to which the alarm rule applies. Set this parameter to Elastic IP and Bandwidth .	Elastic IP and Bandwidth
Dimension	Specifies the item of the monitored service. Bandwidth scaling adjusts the bandwidth. Therefore, set this parameter to Bandwidths .	Bandwidths
Monitoring Scope	Specifies the resources to which the alarm rule applies. Set this parameter to Specific resources . Search for resources by bandwidth name or ID, which can be obtained on the page providing details about the target EIP.	Specific resources
Method	There are three options: Associate template , Use existing template , and Configure manually . NOTE After an associated template is modified, the policies contained in this alarm rule to be created will be modified accordingly.	Configure manually
Alarm Policy	Specifies the alarm policy for triggering the alarm rule. Set this parameter as required. For details about the monitoring metrics, see Table 5-2 .	N/A

4. After setting the parameters, click **Create**.
5. On the **Create Bandwidth Scaling Policy** page, click  to the right of **Rule Name**, and select the created alarm rule.

Alternatively, create your desired alarm rule on the **Cloud Eye** page before creating a bandwidth scaling policy. Ensure that the specific resources selected during alarm rule creation are the bandwidth of the EIP selected for the bandwidth scaling policy to be created. After the alarm rule is created, you can select the rule when creating a bandwidth scaling policy.

Creating a Scheduled or Periodic Bandwidth Scaling Policy

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Bandwidth Scaling**.
3. Click **Create Bandwidth Scaling Policy**.
4. Set parameters, such as the policy name, resource type, policy type, and trigger condition. For details, see [Table 5-4](#).

Table 5-4 Scheduled or periodic policy parameters

Parameter	Description	Example Value
Region	Specifies the region where the AS group resides.	N/A
Policy Name	Specifies the name of the bandwidth scaling policy. The name consists of only letters, digits, underscores (_), and hyphens (-).	as-policy-p6g5
Resource Type	Specifies the type of the bandwidth to be managed. You can select EIP or Shared bandwidth .	EIP
EIP	Specifies the public network IP address whose bandwidth needs to be scaled. This parameter is mandatory when Resource Type is set to EIP . NOTE Only bandwidths of pay-per-use EIPs can be scaled.	N/A
Shared Bandwidth	Specifies the shared bandwidth to be scaled. This parameter is mandatory when Resource Type is set to Shared bandwidth .	N/A

Parameter	Description	Example Value
Policy Type	<p>Specifies the policy type. You can select a scheduled or periodic policy.</p> <p>If you select Periodic, you are required to configure two more parameters:</p> <ul style="list-style-type: none"> • Time Range Specifies the time range during which the AS policy can be triggered. • Period <ul style="list-style-type: none"> - Day - Week - Month 	<p>Day 2023/03/01 00:00:00 - 2023/03/31 23:59:59</p> <p>In this example, the AS policy will trigger a scaling action every day in March, and will become invalid from April 1, 2023 00:00:00.</p>
Triggered At	<p>Specifies the time at which the AS policy is triggered.</p> <p>NOTE The selected triggering time must fall inside the effective time range of the policy.</p>	N/A
Scaling Action	<p>Specifies the action to be performed.</p> <p>The following scaling action options are available:</p> <ul style="list-style-type: none"> • Add When a scaling action is triggered, the bandwidth is increased. • Reduce When a scaling action is triggered, the bandwidth is decreased. • Set to The bandwidth is set to a fixed value. <p>NOTE The step (minimum unit for bandwidth adjustment) varies depending on the bandwidth value range. The bandwidth will be automatically adjusted to the nearest value according to the actual step.</p> <ul style="list-style-type: none"> • If the bandwidth is less than or equal to 300 Mbit/s, the default step is 1 Mbit/s. • If the bandwidth ranges from 300 Mbit/s to 1000 Mbit/s, the default step is 50 Mbit/s. • If the bandwidth is greater than 1000 Mbit/s, the default step is 500 Mbit/s. 	N/A

Parameter	Description	Example Value
Cooldown Period	A cooldown period (in seconds) is the period of time between two scaling actions. AS recounts the cooldown period after a scaling action is complete. During the cooldown period, AS denies all scaling requests triggered by alarm-based policies. Scaling requests triggered manually or by scheduled or periodic policies are not affected.	300s

5. After setting the parameters, click **Create Now**.

5.2 Viewing Details About a Bandwidth Scaling Policy

Scenarios

You can view details about a bandwidth scaling policy, including its basic information and execution logs. Policy execution logs record details about policy execution.

Procedure

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Bandwidth Scaling**.
3. On the **Bandwidth Scaling** page, click the name of a bandwidth scaling policy to go to the page showing its basic information and view its details. You can view basic information about the scaling policy, including **Policy Type**, **Trigger Condition**, and **Scaling Action**.

Viewing Execution Logs of a Bandwidth Scaling Policy

In the **Policy Execution Logs** area on the bandwidth scaling policy details page, you can view the policy execution logs. Policy execution logs record the execution status, execution time, original value, and target value of a bandwidth scaling policy.

5.3 Managing a Bandwidth Scaling Policy

Scenarios

You can adjust the bandwidth through a bandwidth scaling policy.

This section describes how to manage bandwidth scaling policies, including enabling, disabling, modifying, deleting, and immediately executing a bandwidth scaling policy.

 NOTE

The bandwidth scaling policy configured for a released EIP still occupies the policy quota. Only the account and its IAM users with the global permission can manage the bandwidth scaling policy.

Enabling a Bandwidth Scaling Policy

A bandwidth scaling policy can be enabled only when its status is **Disabled**.

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Bandwidth Scaling**.
3. In the bandwidth scaling policy list, locate the row containing the target policy and click **Enable** in the **Operation** column.
4. In the displayed **Enable Bandwidth Scaling Policy** dialog box, click **Yes**.

Disabling a Bandwidth Scaling Policy

A bandwidth scaling policy can be disabled only when its status is **Enabled**.

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Bandwidth Scaling**.
3. In the bandwidth scaling policy list, locate the row containing the target policy and click **Disable** in the **Operation** column.
4. In the displayed **Disable Bandwidth Scaling Policy** dialog box, click **Yes**.

 NOTE

After a bandwidth scaling policy is disabled, its status changes to **Disabled**. AS does not automatically trigger any scaling action based on a **Disabled** bandwidth scaling policy.

Modifying a Bandwidth Scaling Policy

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Bandwidth Scaling**.
3. In the bandwidth scaling policy list, locate the row containing the target policy and click the policy name to switch to its details page.
Click **Modify** in the upper right corner of the page.
You can also locate the row containing the target policy, click **More** in the **Operation** column, and select **Modify**.
4. Modify parameters. You can modify the following parameters of a bandwidth scaling policy: **Policy Name**, **EIP**, **Policy Type**, **Scaling Action**, and **Cooldown Period**.
5. Click **OK**.

 NOTE

A bandwidth scaling policy which is being executed cannot be modified.

Deleting a Bandwidth Scaling Policy

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Bandwidth Scaling**.
3. In the bandwidth scaling policy list, locate the row containing the target policy, click **More** in the **Operation** column, and select **Delete**.
4. In the displayed **Delete Bandwidth Scaling Policy** dialog box, click **Yes**.
You can also select one or more scaling policies and click **Delete** above the list to delete one or more scaling policies.

NOTE

- You can delete a bandwidth scaling policy when you no longer need it. If you do not need it only during a specified period of time, you are advised to disable rather than delete it.
- A bandwidth scaling policy can be deleted only when it is not being executed.

Executing a Bandwidth Scaling Policy

By executing a bandwidth scaling policy, you can immediately adjust the bandwidth to that configured in the bandwidth scaling policy, instead of having to wait until the trigger condition is met.

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Bandwidth Scaling**.
3. In the bandwidth scaling policy list, locate the row that contains the target policy and click **Execute Now** in the **Operation** column.
4. In the displayed **Execute Bandwidth Scaling Policy** dialog box, click **Yes**.

You can also go to the bandwidth scaling policy details page and click **Execute Now** in the upper right corner.

NOTE

- A bandwidth scaling policy can be executed only when the policy is enabled and no other bandwidth scaling policy is being executed.
- Executing a bandwidth scaling policy does not affect automatic adjustment of the bandwidth when the trigger condition of the policy is met.
- If **Policy Type** is set to **Alarm** and **Alarm Policy Type** to **Refined scaling**, the bandwidth scaling policy cannot be executed immediately.

6 AS Group and Instance Monitoring

6.1 Health Check

Health Check Methods

A health check removes unhealthy instances from an AS group. Then, AS adds new instances to the AS group so that the number of instances is the same as the expected number. There are two types of AS group health checks.

- **ECS health check:** checks ECS instance running status. If an instance is stopped or deleted, it is considered abnormal. **ECS health check** is the default health check mode for an AS group. The AS group periodically uses the check result to determine the running status of every instance in the AS group. If the health check results show that an instance is unhealthy, AS removes the instance from the AS group.
- **ELB health check:** determines ECS instance running status using a load balancing listener. If the AS group uses load balancers, the health check method can also be **ELB health check**.

If you add multiple load balancers to an AS group, an ECS instance is considered to be healthy only when all load balancers detect that the instance status is healthy. If any load balancer detects that an instance is unhealthy, the instance will be removed from the AS group.

In both **ECS health check** and **ELB health check** modes, AS removes unhealthy instances from the AS group. Whether a removed instance will be deleted depends on how the instance is added to the AS group.

Table 6-1 Instance removal and deletion rules

Instance Type	Description	Billing Mode	Removed If Unhealthy	Deleted When Removed
Automatically added instances	Instances automatically created and added to an AS group in a scaling action	Pay-per-use By default, this type of instances is billed on a pay-per-use basis.	Yes	Yes
		Yearly/Monthly The billing mode of an automatically added instance can be manually changed from pay-per-use to yearly/monthly.	Yes	No
Manually added instances	Instances manually created and added to an AS group	Pay-per-use	Yes	No
		Yearly/Monthly	Yes	No

NOTE

If you need to perform custom operations on an instance that is being deleted or created, you can use [lifecycle hooks](#).

Constraints

- Even when an AS group is disabled, AS still checks the health of instances in the AS group, but does not remove unhealthy instances.
- AS does not check the health of instances in standby state.

6.2 Configuring Notifications for an AS Group**Scenarios**

After the SMN service is provisioned, you can promptly send AS group information, such as successful instance increasing, failed instance increasing, successful instance decreasing, failed instance decreasing, or AS group exception to the user using the SMN service. This helps the user learn the AS group status.

To configure notifications for an AS group, you need to specify a notification event and topic. You need to create a notification topic on the SMN console. When the

notification scenario matched with the notification topic appears, the AS group sends a notification to the subscribers.

A maximum of five notifications can be configured for an AS group.

Procedure

1. Log in to the management console.
1. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
2. Click the name of the target AS group. On the AS group details page, click the **Notifications** tab and then click **Add Notification**.
3. Set the parameters listed in [Table 6-2](#).

Table 6-2 Parameter description

Parameter	Description	Example Value
Event	When at least one of the following conditions is met, SMN sends a notification to the user: <ul style="list-style-type: none">• Instance creation succeeds• Instance removal succeeds• Errors occur in an AS group• Instance creation fails• Instance removal fails	N/A
Topic	Select a created topic. For details about how to create a topic, see Creating a Topic .	N/A

4. Click **OK**.

6.3 Recording AS Resource Operations

Scenarios

AS can use the Cloud Trace Service (CTS) to record resource operations. CTS can record operations performed on the management console, operations performed by calling APIs, and operations triggered within the cloud system.

If you have enabled CTS, when a call is made to the AS API, the operation will be reported to CTS which will then deliver the operation record to a specified OBS bucket for storage. With CTS, you can record operations associated with AS for later query, audit, and backtrack operations.

Obtaining AS Information in CTS

After you enable CTS in the application system, the system logs the API calling operations performed on AS resources. On the **Cloud Trace Service** console, you

can view operation records for the last 7 days. For details, see [Querying Real-Time Traces](#). To obtain more operation records, you can enable the Object Storage Service (OBS) and synchronize operation records to the OBS in real time.

Table 6-3 lists the AS operations that can be recorded by CTS.

Table 6-3 AS operations that can be recorded by CTS

Operation	Resource Type	Trace Name
Creating an AS group	scaling_group	createScalingGroup
Modifying an AS group	scaling_group	modifyScalingGroup
Deleting an AS group	scaling_group	deleteScalingGroup
Enabling an AS group	scaling_group	enableScalingGroup
Disabling an AS group	scaling_group	disableScalingGroup
Performing operations on an AS group	scaling_group	operateScalingGroup
Creating an AS configuration	scaling_configuration	createScalingConfiguration
Deleting an AS configuration	scaling_configuration	deleteScalingConfiguration
Deleting AS configurations in a batch	scaling_configuration	batchDeleteScalingConfiguration
Creating an AS policy	scaling_policy	createScalingPolicy
Modifying an AS policy	scaling_policy	modifyScalingPolicy
Deleting an AS policy	scaling_policy	deleteScalingPolicy
Enabling an AS policy	scaling_policy	enableScalingPolicy
Disabling an AS policy	scaling_policy	disableScalingPolicy

Operation	Resource Type	Trace Name
Executing an AS policy	scaling_policy	executeScalingPolicy
Performing operations on an AS policy	scaling_policy	operateScalingPolicy
Enabling AS policies in a batch	scaling_policy	batchEnableScalingPolicies
Disabling AS policies in a batch	scaling_policy	batchDisableScalingPolicies
Removing an instance	scaling_instance	removeInstance
Removing instances in batches	scaling_instance	batchRemoveInstances
Adding instances in batches	scaling_instance	batchAddInstances
Performing operations on instances in batches	scaling_instance	batchOperateInstance
Enabling instance protection in a batch	scaling_instance	batchProtectInstances
Disabling instance protection in a batch	scaling_instance	batchUnprotectInstances
Putting instances into standby in a batch	scaling_instance	batchEnterStandbyInstances
Configuring a notification	scaling_notification	putScalingNotification
Deleting a notification	scaling_notification	deleteScalingNotification
Creating a lifecycle hook	scaling_lifecycle_hook	createLifecycleHook

Operation	Resource Type	Trace Name
Modifying a lifecycle hook	scaling_lifecycle_hook	modifyLifecycleHook
Deleting a lifecycle hook	scaling_lifecycle_hook	deleteLifecycleHook

6.4 Querying Real-Time Traces

Scenarios

After you enable CTS and the management tracker is created, CTS starts recording operations on cloud resources. After a data tracker is created, the system starts recording operations on data in OBS buckets. CTS stores operation records generated in the last seven days.


This section describes how to query and export operation records of the last seven days on the CTS console.





- [Viewing Real-Time Traces in the Trace List of the New Edition](#)
- [Viewing Real-Time Traces in the Trace List of the Old Edition](#)

Constraints


- Traces of a single account can be viewed on the CTS console. Multi-account traces can be viewed only on the **Trace List** page of each account, or in the OBS bucket or the **CTS/system** log stream configured for the management tracker with the organization function enabled.
- You can only query operation records of the last seven days on the CTS console. To store operation records for more than seven days, you must configure an OBS bucket to transfer records to it. Otherwise, you cannot query the operation records generated seven days ago.
- After performing operations on the cloud, you can query management traces on the CTS console 1 minute later and query data traces on the CTS console 5 minutes later.



Viewing Real-Time Traces in the Trace List of the New Edition

1. Log in to the management console.
2. Click  in the upper left corner and choose **Management & Governance > Cloud Trace Service**. The CTS console is displayed.
3. Choose **Trace List** in the navigation pane on the left.
4. On the **Trace List** page, use advanced search to query traces. You can combine one or more filters.
 - **Trace Name:** Enter a trace name.
 - **Trace ID:** Enter a trace ID.
 - **Resource Name:** Enter a resource name. If the cloud resource involved in the trace does not have a resource name or the corresponding API

- operation does not involve the resource name parameter, leave this field empty.
- **Resource ID:** Enter a resource ID. Leave this field empty if the resource has no resource ID or if resource creation failed.
 - **Trace Source:** Select a cloud service name from the drop-down list.
 - **Resource Type:** Select a resource type from the drop-down list.
 - **Operator:** Select one or more operators from the drop-down list.
 - **Trace Status:** Select **normal**, **warning**, or **incident**.
 - **normal:** The operation succeeded.
 - **warning:** The operation failed.
 - **incident:** The operation caused a fault that is more serious than the operation failure, for example, causing other faults.
 - Time range: Select **Last 1 hour**, **Last 1 day**, or **Last 1 week**, or specify a custom time range.
5. On the **Trace List** page, you can also export and refresh the trace list, and customize the list display settings.
- Enter any keyword in the search box and click  to filter desired traces.
 - Click **Export** to export all traces in the query result as an .xlsx file. The file can contain up to 5000 records.
 - Click  to view the latest information about traces.
 - Click  to customize the information to be displayed in the trace list. If **Auto wrapping** is enabled () , excess text will move down to the next line; otherwise, the text will be truncated. By default, this function is disabled.
6. For details about key fields in the trace structure, see [Trace Structure](#) and [Example Traces](#).
7. (Optional) On the **Trace List** page of the new edition, click **Go to Old Edition** in the upper right corner to switch to the **Trace List** page of the old edition.

Viewing Real-Time Traces in the Trace List of the Old Edition

1. Log in to the management console.
2. Click  in the upper left corner and choose **Management & Governance > Cloud Trace Service**. The CTS console is displayed.
3. Choose **Trace List** in the navigation pane on the left.
4. Each time you log in to the CTS console, the new edition is displayed by default. Click **Go to Old Edition** in the upper right corner to switch to the trace list of the old edition.
5. Set filters to search for your desired traces. The following filters are available:
 - **Trace Type, Trace Source, Resource Type, and Search By:** Select a filter from the drop-down list.

- If you select **Resource ID** for **Search By**, specify a resource ID.
 - If you select **Trace name** for **Search By**, specify a trace name.
 - If you select **Resource name** for **Search By**, specify a resource name.
 - **Operator**: Select a user.
 - **Trace Status**: Select **All trace statuses**, **Normal**, **Warning**, or **Incident**.
 - Time range: You can query traces generated during any time range in the last seven days.
 - Click **Export** to export all traces in the query result as a CSV file. The file can contain up to 5000 records.
6. Click **Query**.
 7. On the **Trace List** page, you can also export and refresh the trace list.
 - Click **Export** to export all traces in the query result as a CSV file. The file can contain up to 5000 records.
 - Click  to view the latest information about traces.
 8. Click  on the left of a trace to expand its details.

Trace Name	Resource Type	Trace Source	Resource ID	Resource Name	Trace Status	Operator	Operation Time	Operation
createDockerConfig	dockerlogincmd	SWR	-	dockerlogincmd	normal		Nov 16, 2023 10:54:04 GMT+08:00	View Trace

```

request
trace_id
code 200
trace_name createDockerConfig
resource_type dockerlogincmd
trace_rating normal
api_version
message createDockerConfig, Method: POST Url=/v2/manage/utlils/secret, Reason:
source_ip
domain_id
trace_type ApiCall
            
```

9. Click **View Trace** in the **Operation** column. The trace details are displayed.

View Trace ✕

```

{
  "request": "",
  "trace_id": "",
  "code": "200",
  "trace_name": "createDockerConfig",
  "resource_type": "dockerlogincmd",
  "trace_rating": "normal",
  "api_version": "",
  "message": "createDockerConfig, Method: POST Url=/v2/manage/utlils/secret, Reason:",
  "source_ip": "",
  "domain_id": "",
  "trace_type": "ApiCall",
  "service_type": "SWR",
  "event_type": "system",
  "project_id": "",
  "response": "",
  "resource_id": "",
  "tracker_name": "system",
  "time": "Nov 16, 2023 10:54:04 GMT+08:00",
  "resource_name": "dockerlogincmd",
  "user": {
    "domain": {
      "name": "",
      "id": ""
    }
  }
}
            
```

10. For details about key fields in the trace structure, see [Trace Structure](#) and [Example Traces](#).
11. (Optional) On the **Trace List** page of the old edition, click **New Edition** in the upper right corner to switch to the **Trace List** page of the new edition.

6.5 Adding Tags to AS Groups and Instances

Scenarios

If you have many resources of the same type, you can use a tag to manage resources flexibly. You can identify specified resources quickly using the tags allocated to them.

Using a tag, you can assign custom data to each AS group. You can organize and manage AS groups, for example, classify AS group resources by usage, owner, or environment.

Each tag contains a key and a value. You can specify the key and value for each tag. A key can be a category associated with certain values, such as usage, owner, and environment.

For example, if you want to distinguish between the test environment and production environment, you can allocate a tag with the key **environment** to each AS group. For the test environment, the key value is **test** and for the production environment, the key value is **production**. You are advised to use one or more groups of consistent tags to manage your AS group resources.

After you allocate a tag to an AS group, the system will automatically add the tag to the instances automatically created in the AS group. If you add a tag to an AS group or modify the tag, the new tag will be added to the ECSs automatically created in the AS group. Creating, deleting, or modifying the tag of an AS group will have no impact on the ECSs in the AS group.

Restrictions of Using Tags

You must observe the following rules when using tags:

- Each AS group can have a maximum of 10 tags added to it.
- Each tag contains a key and a value.
- You can set the tag value to an empty character string.
- If you delete an AS group, all tags of it will also be deleted.
- If you have configured tag policies for AS, add tags to AS groups based on the tag policies. If you add a tag that does not comply with the tag policies, AS groups may fail to be created or the tags may fail to be added or modified. Contact your administrator to learn more about tag policies.

Adding a Tag to an AS Group

1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. Click the AS group name. On the AS group details page, click the **Tags** tab and then click **Add Tag**.
4. Set the parameters listed in [Table 6-4](#).

Table 6-4 Tag naming rules

Parameter	Requirement	Example Value
Tag Key	<ul style="list-style-type: none"> The value cannot be empty. An AS group has a unique key. A key can contain a maximum of 36 characters, including digits, letters, underscores (_), hyphens (-), and Unicode characters from \u4e00 to \u9fff. 	Organization
Tag Value	<ul style="list-style-type: none"> The value can be an empty character string. A key can have only one value. A tag value can contain a maximum of 43 characters, including digits, letters, underscores (_), periods (.), hyphens (-), and Unicode characters from \u4e00 to \u9fff. 	Apache

5. Click **OK**.

Modifying or Deleting Tags of an AS Group

1. Log in to the management console.
1. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
2. Click the AS group name. On the **Overview** page, click the **Tags** tab.
3. Locate the row that contains the tag and click **Edit** or **Delete** in the **Operation** column.

After clicking **Edit**, configure required parameters. For details, see [Table 6-4](#).

After you click **Delete**, the added tag will be deleted.

6.6 Monitoring Metrics

Description

This section describes the monitoring metrics reported by AS to Cloud Eye and defines the namespace for the metrics. You can use Cloud Eye to query monitoring metrics and alarms generated by AS.

Namespace

SYS.AS

Monitoring metrics

[Table 6-5](#) lists the AS metrics supported by Cloud Eye.

Table 6-5 AS metrics

Metric ID	Metric	Description	Value Range	Monitored Object	Monitoring Interval (Raw Data)
cpu_util	CPU Usage	CPU usage of an AS group Formula: Total CPU usage of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Percent	≥0%	AS group	5 minutes
mem_util	Memory Usage	Memory usage of an AS group Formula: Total memory usage of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Percent NOTE This metric is unavailable if the image has no VM Tools installed.	≥0%	AS group	5 minutes
instance_num	Number of Instances	Number of available ECS instances in an AS group Formula: Total number of ECS instances in Enabled state in the AS group	≥0	AS group	5 minutes
network_incoming_bytes_rate_inband	Inband Incoming Rate	Number of incoming bytes per second on an ECS in an AS group Formula: Total inband incoming rates of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Byte/s	≥0 Byte/s	AS group	5 minutes

Metric ID	Metric	Description	Value Range	Monitored Object	Monitoring Interval (Raw Data)
network_outgoing_bytes_rate_inband	Inband Outgoing Rate	Number of outgoing bytes per second on an ECS in an AS group Formula: Total inband outgoing rates of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Byte/s	≥0 Byte/s	AS group	5 minutes
disk_read_bytes_rate	Disks Read Rate	Number of bytes read from an AS group per second Formula: Total disks read rates of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Byte/s	≥0 Byte/s	AS group	5 minutes
disk_write_bytes_rate	Disks Write Rate	Number of bytes written to an AS group per second Formula: Total disks write rates of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Byte/s	≥0 Byte/s	AS group	5 minutes
disk_read_requests_rate	Disks Read Requests	Number of read requests per second sent to an ECS disk in an AS group Formula: Total disks read rates of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Request/s	≥0 request/s	AS group	5 minutes

Metric ID	Metric	Description	Value Range	Monitored Object	Monitoring Interval (Raw Data)
disk_write_requests_rate	Disks Write Requests	Number of write requests per second sent to an ECS disk in an AS group Formula: Total disks write rates of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Request/s	≥0 request/s	AS group	5 minutes
cpu_usage	(Agent) CPU Usage	Agent CPU usage of an AS group Formula: Total Agent CPU usage of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Percent	0-100 %	AS group	1 minute
mem_usedPercent	(Agent) Memory Usage	Agent memory usage of an AS group Formula: Total Agent memory usage of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Percent	0-100 %	AS group	1 minute
load_average1	(Agent) 1-Minute Load Average	Average CPU load of all ECSs in an AS group in the last 1 minute Unit: none	≥0	AS group	1 minute
load_average5	(Agent) 5-Minute Load Average	Average CPU load of all ECSs in an AS group in the last 5 minutes Unit: none	≥0	AS group	1 minute

Metric ID	Metric	Description	Value Range	Monitored Object	Monitoring Interval (Raw Data)
load_ave rage15	(Agent) 15- Minute Load Average	Average CPU load of all ECSs in an AS group in the last 15 minutes Unit: none	≥0	AS group	1 minute
gpu_usa ge_gpu	(Agent) GPU Usage	Agent GPU usage of an AS group Formula: Total Agent GPU usage of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Percent	0-100 %	AS group	1 minute
gpu_usa ge_mem	(Agent) Video Memory Usage	Agent GPU memory usage of an AS group Formula: Total Agent GPU memory usage of all ECS instances in an AS group/Number of ECS instances in the AS group Unit: Percent	0-100 %	AS group	1 minute

 **NOTE**

Monitoring metrics are classified into metrics with Agent and without Agent. For some OSs, you need to install the Agent to obtain the corresponding monitoring metrics. In this case, select the monitoring metrics with Agent, for example, (Agent) Memory Usage.

 **NOTE**

OSs determine whether the **Memory Usage**, **Inband Outgoing Rate**, and **Inband Incoming Rate** metrics are supported. For details, see [Elastic Cloud Server User Guide](#).
Before using Agent monitoring metrics, make sure that the Agent has been installed on the instances in the AS group. For details, see [How Do I Install the Agent Plug-in on the Instances in an AS Group to Use Agent Monitoring Metrics?](#)

Dimension

Key	Value
AutoScalingGroup	AS group ID

6.7 Viewing Monitoring Metrics

Scenarios

The cloud platform provides Cloud Eye to help you obtain the running status of your ECS instances. This section describes how to view details of AS group metrics to obtain information about the status of the ECS instances in the AS group.

Prerequisites

The ECS instance is running properly.


NOTE

- Monitoring metrics such as **CPU Usage** and **Disks Read Rate** are available only when there is at least one instance in an AS group. If not, only the **Number of Instances** metric is available.
- Monitoring data is not displayed for a stopped, faulty, or deleted instance. After such an instance restarts or recovers, the monitoring data is available.


Viewing Monitoring Metrics on Auto Scaling


1. Log in to the management console.
2. Under **Compute**, click **Auto Scaling**. In the navigation pane on the left, choose **Instance Scaling**. Then click the **AS Groups** tab.
3. On the **AS Groups** page, find the AS group to view monitoring data and click its name.
4. Click the **Monitoring** tab to view the monitoring data.

You can view data of the last one, three, 12, or 24 hours, or last 7 days. If you want to view data for a longer time range, click **View details** to go to the

Cloud Eye page, hover your mouse over a graph, and click .

Viewing Monitoring Metrics on Cloud Eye

1. Log in to the management console.
2. Click  in the upper left corner to select a region and a project.
3. Under **Management & Governance**, select **Cloud Eye**.
4. In the navigation pane on the left, choose **Cloud Service Monitoring > Auto Scaling**.
5. Locate the row that contains the target AS group and click **View Metric** in the **Operation** column to view monitoring data.

You can view data of the last one, three, 12, or 24 hours, or last 7 days. Hover your mouse over a graph and click  to view data for a longer time range.

NOTE


It can take a period of time to obtain and transfer the monitoring data. Therefore, wait for a while and then check the data.

6.8 Setting Monitoring Alarm Rules

Scenarios

Setting alarm rules allows you to customize the monitored objects and notification policies and determine the running status of your ECS instances at any time.

Procedure

1. Log in to the management console.
2. Click  in the upper left corner and select the desired region and project.
3. Under **Management & Governance**, select **Cloud Eye**.
4. In the navigation pane, choose **Alarm Management > Alarm Rules**.
5. On the **Alarm Rules** page, click **Create Alarm Rule** to create an alarm rule for the AS service or modify an existing alarm rule of the AS service.
6. After setting the parameters, click **Create**.

NOTE

- For more information about how to set alarm rules, see *Cloud Eye User Guide*.
- You can create alarm rules on the Cloud Eye console to dynamically expand resources.

7 Permissions Management

7.1 Creating a User and Granting AS Permissions

Scenarios

IAM can help you implement fine-grained permissions control over your AS resources. With IAM, you can:

- Create IAM users for employees based on your enterprise's organizational structure. Each IAM user will have their own security credentials for accessing AS resources.
- Grant only the permissions required for users to perform a specific task.
- Use **IAM** to entrust a Huawei Cloud account or cloud service to perform efficient O&M on your AS resources.

If your Huawei Cloud account does not require individual IAM users, skip this section.

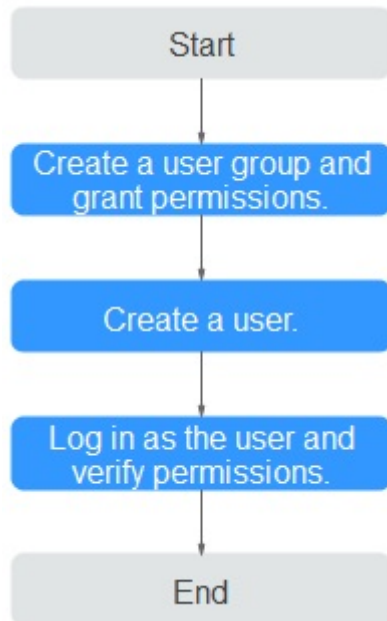
This section describes the procedure for granting permissions (see [Figure 7-1](#)).

Prerequisites

Learn about the permissions (see [Permissions Management](#)) supported by AS and choose policies or roles according to your requirements. For the permissions of other services, see [System Permissions](#).

Process Flow

Figure 7-1 Process for granting AS permissions



1. **Create a user group and assign permissions to it.**
Create a user group on the IAM console, and attach the **AutoScalingReadOnlyAccess** policy to the group.
2. **Create an IAM user and add it to the user group.**
Create a user on the IAM console and add the user to the group created in 1.
3. **Log in** and verify permissions.
Log in to the AS console as the created user, and verify the user's permissions for AS.
 - Choose **Service List > Auto Scaling**. Then, click **Create AS Group** on the AS console. If a message appears indicating that you have insufficient permissions to perform the operation, the **AutoScalingReadOnlyAccess** policy has already taken effect.
 - Choose any other service in the **Service List**. If a message appears indicating that you have insufficient permissions to access the service, the **AutoScalingReadOnlyAccess** policy has already taken effect.

7.2 AS Custom Policies

Scenarios

Custom policies can be created to supplement the system-defined policies of AS. For the actions that can be added to custom policies, see [Permissions Policies and Supported Actions](#).

You can create custom policies in either of the following ways:

- Visual editor: Select cloud services, actions, resources, and request conditions. This does not require knowledge of policy syntax.
- JSON: Edit JSON policies from scratch or based on an existing policy.

For operation details, see [Creating a Custom Policy](#). The following section contains examples of common AS custom policies.

Example Custom Policies

- Example 1: Allowing users to remove instances from an AS group and create AS configurations

```
{
  "Version": "1.1",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "as:instances:delete",
        "as:configs:create"
      ]
    }
  ]
}
```

- Example 2: Denying AS group deletion

A policy with only "Deny" permissions must be used in conjunction with other policies to take effect. If the permissions assigned to a user contain both "Allow" and "Deny", the "Deny" permissions take precedence over the "Allow" permissions.

The following method can be used if you need to assign permissions of the **AutoScaling FullAccess** policy to a user but you want to prevent the user from deleting AS groups. Create a custom policy for denying AS group deletion, and attach both policies to the group to which the user belongs. Then, the user can perform all operations on AS except deleting AS groups. The following is an example of a deny policy:

```
{
  "Version": "1.1",
  "Statement": [
    {
      "Action": [
        "as:groups:delete"
      ],
      "Effect": "Deny"
    }
  ]
}
```

A Change History

Released On	Description
2021-10-30	This issue is the thirteenth official release. Modified the following content: Added section "Permissions Management."
2020-10-19	This issue is the twelfth official release. Modified the following content: Added section "Access Methods."
2019-08-30	This issue the eleventh official release. Modified the following content: Added Agent monitoring metrics in Monitoring Metrics .
2019-01-30	This issue is the tenth official release. Modified the following content: <ul style="list-style-type: none">• Optimized content about the AS policy and bandwidth scaling policy.• Added Putting an Instance Into the Standby State.

Released On	Description
2018-11-30	<p>This issue is the ninth official release.</p> <p>Modified the following content:</p> <ul style="list-style-type: none">• Added the operations to view historical scaling actions to Viewing a Scaling Action.• Organized FAQs by category.• Added FAQs "What Is the Expected Number of Instances?", "What Operation Will Be Suspended After an AS Group Is Disabled?", "Why Do Instances in an AS Group Frequently Fail Health Checks and Are Deleted and Then Created Repeatedly?", "What Are the Conditions to Trigger an Alarm in the AS Policy?", and "Do I Need to Configure an EIP in an AS Configuration When a Load Balancer Has Been Enabled in an AS Group?"• Optimized FAQs "What Can I Do to Enable My Application to Be Automatically Deployed on an Instance?" and "What Is a Cooldown Period? Why Is It Required?"
2018-09-30	<p>This issue is the eighth official release.</p> <p>Modified the following content:</p> <ul style="list-style-type: none">• Added parameter ECS Group to Creating an AS Configuration from an Existing ECS Instance and Creating an AS Configuration from a New Specifications Template.• Added parameter Bandwidth Type to Creating an AS Configuration from an Existing ECS Instance and Creating an AS Configuration from a New Specifications Template.
2018-08-30	<p>This issue is the seventh official release.</p> <p>Modified the following content:</p> <ul style="list-style-type: none">• The cooldown period starts after a scaling action is complete and the default cooldown period is 300 seconds.• Added FAQ "What Is a Cooldown Period? How Is It Calculated?"• Added the application scenarios of the scaling bandwidth in Bandwidth Scaling.

Released On	Description
2018-07-30	<p>This issue is the sixth official release.</p> <p>Modified the following content:</p> <ul style="list-style-type: none"> • Optimized "Creating an AS Group Quickly", "Creating an AS Group", and "Dynamically Expanding Resources". • Added FAQ "Will the Data on an Instance Be Retained After the Instance Is Removed from an AS Group and Deleted?" • Added FAQ "Can AS Scale Capacity Based on Custom Monitoring of Cloud Eye?" • Added FAQ "Can AS Automatically Scale Up and Down vCPUs, Memory, and Bandwidth of ECSs?" • Added option custom monitoring when Alarm is set to Create Alarm Rule to Dynamic Scaling. • Added 10s and 1 min as new options for the health check interval in Creating an AS Group.
2018-06-30	<p>This issue is the fifth official release.</p> <p>Modified the following content:</p> <ul style="list-style-type: none"> • Optimized "Creating an AS Group" and "Creating an AS Configuration". • Added FAQ "What Can I Do If an AS Policy Fails to Be Enabled?" • Added FAQ "Can I Add Yearly/Monthly ECSs?" • Added FAQ "How Do I Prevent ECSs in an AS Group from Being Removed Automatically?"
2018-05-30	<p>This issue is the fourth official release.</p> <p>Modified the following content:</p> <ul style="list-style-type: none"> • Added the bandwidth scaling feature. • Added parameter "Security Group" to Creating an AS Configuration from a New Specifications Template. • Added "View Audit Logs" to Recording AS Resource Operations. • Added Monitoring Metrics. • Added Viewing Monitoring Metrics. • Added Setting Monitoring Alarm Rules. • Added FAQ "What Can I Do If the AS Group Fails to Be Enabled?" • Added FAQ "How Should I Handle Unhealthy Instances in an AS Group?"

Released On	Description
2018-04-30	This issue is the third official release. Modified the following content: <ul style="list-style-type: none">• Added FAQ "How Do I Delete an ECS Created in a Scaling Action?"
2018-03-30	This issue is the second official release. Modified the following content: <ul style="list-style-type: none">• Optimized description of user data injection.• Optimized description of the instance health check.
2017-12-31	This issue is the first official release.