Alibaba Cloud

Elastic Compute Service

Best Practices

Document Version: 20220711

(-) Alibaba Cloud

Legal disclaimer

Alibaba Cloud reminds you to carefully read and fully understand the terms and conditions of this legal disclaimer before you read or use this document. If you have read or used this document, it shall be deemed as your total acceptance of this legal disclaimer.

- You shall download and obtain this document from the Alibaba Cloud website or other Alibaba Cloudauthorized channels, and use this document for your own legal business activities only. The content of this document is considered confidential information of Alibaba Cloud. You shall strictly abide by the confidentiality obligations. No part of this document shall be disclosed or provided to any third party for use without the prior written consent of Alibaba Cloud.
- 2. No part of this document shall be excerpted, translated, reproduced, transmitted, or disseminated by any organization, company or individual in any form or by any means without the prior written consent of Alibaba Cloud.
- 3. The content of this document may be changed because of product version upgrade, adjustment, or other reasons. Alibaba Cloud reserves the right to modify the content of this document without notice and an updated version of this document will be released through Alibaba Cloud-authorized channels from time to time. You should pay attention to the version changes of this document as they occur and download and obtain the most up-to-date version of this document from Alibaba Cloud-authorized channels.
- 4. This document serves only as a reference guide for your use of Alibaba Cloud products and services. Alibaba Cloud provides this document based on the "status quo", "being defective", and "existing functions" of its products and services. Alibaba Cloud makes every effort to provide relevant operational guidance based on existing technologies. However, Alibaba Cloud hereby makes a clear statement that it in no way guarantees the accuracy, integrity, applicability, and reliability of the content of this document, either explicitly or implicitly. Alibaba Cloud shall not take legal responsibility for any errors or lost profits incurred by any organization, company, or individual arising from download, use, or trust in this document. Alibaba Cloud shall not, under any circumstances, take responsibility for any indirect, consequential, punitive, contingent, special, or punitive damages, including lost profits arising from the use or trust in this document (even if Alibaba Cloud has been notified of the possibility of such a loss).
- 5. By law, all the contents in Alibaba Cloud documents, including but not limited to pictures, architecture design, page layout, and text description, are intellectual property of Alibaba Cloud and/or its affiliates. This intellectual property includes, but is not limited to, trademark rights, patent rights, copyrights, and trade secrets. No part of this document shall be used, modified, reproduced, publicly transmitted, changed, disseminated, distributed, or published without the prior written consent of Alibaba Cloud and/or its affiliates. The names owned by Alibaba Cloud shall not be used, published, or reproduced for marketing, advertising, promotion, or other purposes without the prior written consent of Alibaba Cloud. The names owned by Alibaba Cloud include, but are not limited to, "Alibaba Cloud", "Aliyun", "HiChina", and other brands of Alibaba Cloud and/or its affiliates, which appear separately or in combination, as well as the auxiliary signs and patterns of the preceding brands, or anything similar to the company names, trade names, trademarks, product or service names, domain names, patterns, logos, marks, signs, or special descriptions that third parties identify as Alibaba Cloud and/or its affiliates.
- 6. Please directly contact Alibaba Cloud for any errors of this document.

Document conventions

Style Description		Example
<u> Danger</u>	A danger notice indicates a situation that will cause major system changes, faults, physical injuries, and other adverse results.	Danger: Resetting will result in the loss of user configuration data.
<u> </u>	A warning notice indicates a situation that may cause major system changes, faults, physical injuries, and other adverse results.	Warning: Restarting will cause business interruption. About 10 minutes are required to restart an instance.
Notice	A caution notice indicates warning information, supplementary instructions, and other content that the user must understand.	Notice: If the weight is set to 0, the server no longer receives new requests.
? Note	A note indicates supplemental instructions, best practices, tips, and other content.	Note: You can use Ctrl + A to select all files.
>	Closing angle brackets are used to indicate a multi-level menu cascade.	Click Settings> Network> Set network type.
Bold	Bold formatting is used for buttons , menus, page names, and other UI elements.	Click OK.
Courier font	Courier font is used for commands	Run the cd /d C:/window command to enter the Windows system folder.
Italic	Italic formatting is used for parameters and variables.	bae log listinstanceid Instance_ID
[] or [a b]	This format is used for an optional value, where only one item can be selected.	ipconfig [-all -t]
{} or {a b}	This format is used for a required value, where only one item can be selected.	switch {active stand}

Table of Contents

1.	Quick reference	07
2.	Best practices for instance type selection	12
3.	Provisioning methods of ECS instances	21
4.	Batch configure sequential names or hostnames for multiple in	25
5.	Use auto provisioning group-related API operations to create	30
6.	Best practices for heterogeneous computing services	39
7.	Best practices for preemptible instances	40
	7.1. Query the price of a preemptible instance	40
	7.2. Select a bidding mode for a preemptible instance	42
	7.3. Simulate a preemptible instance interruption event	46
	7.4. Receive a preemptible instance interruption event	51
8.	.Security	57
	8.1. Best practices of the security group (part 1)	57
	8.2. Best practices for ECS security groups (part 2)	59
	8.3. Best practices for ECS security groups (part 3)	64
	8.4. Best practices for ECS data security	67
	8.5. Make ECS instances more secure	68
	8.6. Configure interconnection of instances in the classic netw	71
	8.7. Modify the default port used by an instance to accept con	74
	8.8. Use logs in Windows instances	78
	8.9. Network isolation within a basic security group	79
	8.10. Security group quintuple rules	81
	8.11. Use Cloud Firewall to control access between ECS instanc	83
	8.12. Enable or disable SELinux	86
	8.13. Revoke the authorization for internal network communica	88
	8.14. Authorize internal network communication between ECS i	89

9.Data recovery	92
9.1. Handle low disk space on Windows instances	92
9.2. Restore data in Linux instances	- 94
9.3. Restore data in Windows instances	100
10.Configuration preference	103
10.1. Transfer ECS instance data	103
10.2. Increase data throughput through read/write splitting	108
10.3. Change the preferred language of a Windows instance	115
10.4. Boot a Linux ECS instance into single user mode	119
11.Block Storage	125
11.1. Resize partitions and file systems of Linux system disks	125
11.2. Resize partitions and file systems of Linux data disks	130
11.3. Use LVs for Linux	148
11.3.1. Use LVM to create a logical volume	148
11.3.2. Resize an LV by using LVM	153
11.4. Create a RAID array for a Linux instance	155
11.5. Modify the UUID of a disk	159
11.6. Configure UUIDs in the fstab file to automatically attach	161
11.7. Shrink a disk	165
11.8. Encrypt data stored on ECS resources	166
12.Best practices for tag design	171
13.Set the boot mode of custom images to the UEFI mode by c	174
14.Best practices for using custom images	176
14.1. Overview	176
14.2. Use OOS to update custom images	177
14.3. Packer: machine images as code	180
14.3.1. Benefits of using Packer to create custom images	180
14.3.2. Alicloud Image Builder parameters used to implement	186

14.4. Create and import a custom image	190
15.Monitor	196
15.1. Use CloudMonitor to monitor websites deployed on ECS i	196
16.Use RAM roles to access other Alibaba Cloud services	202
17.Networks	207
17.1. Best practices for configuring public bandwidth	207
17.2. Best practices for testing network performance	211
18.Automate O&M based on status change events of ECS instanc	222
19.Use the snapshot and image features to migrate instance dat	230
20.Disaster recovery solutions	234
21.Deploy a highly available architecture	237
21.1. Deploy a highly available architecture	237
21.2. Replicate ECS instances	238
21.3. Configure an SLB instance	239
21.4. Migrate self-managed databases to ApsaraDB RDS instan	242
22.Use Analytics Zoo and bfloat16 to accelerate AI applications	244

1.Quick reference

Common operations on ECS instancesECS learning path

When you use Elastic Compute Service (ECS), you may need to perform various operations on different resources, such as connecting to instances, replacing operating systems, resizing disks, upgrading or downgrading instance configurations, and using snapshots or images. This topic describes common operations on ECS resources.

Limits

- For information about the usage notes of ECS instances, see Usage notes.
- For information about the limits of ECS resources, see 使用限制 and View and increase instance quotas.

•

Create and manage ECS instances

- You can perform the following steps to manage the lifecycle of an ECS instance:
 - i. Create an instance by using the wizard
 - ii. Connect to an ECS instance
 - iii. Stop an instance
 - iv. Release an instance
- If the instance type or network configuration of your instance cannot meet your business requirements, you can change the instance type, IP address, and maximum public bandwidth.
 - Subscription instances:
 - Upgrade the instance types of subscription instances
 - Downgrade the configurations of an instance during renewal
 - o Pay-as-you-go instances:
 - Change the instance type of a pay-as-you-go instance
 - Modify the bandwidth configurations of pay-as-you-go instances
 - o IP addresses of ECS instances:
 - Change the public IP address of an instance
 - Convert the public IP address of a VPC-type instance to an EIP
- If the operating system of your instance cannot meet your requirements, you can replace the operating system. For more information, see Change the operating system.
- You can use the following features to manage ECS instances in a fine-grained manner:
 - User data
 - Instance metadata
 - Instance identity
 - Instance RAM roles

Manage the billing of instances

• Subscription instances:

You can use one of the following methods to renew subscription instances:

- Manually renew an instance
- Enable auto-renewal for an instance
- Downgrade the configurations of an instance during renewal
- Pay-as-you-go instances:

You can enable the economical mode for pay-as-you-go instances. For more information, see Economical mode.

- Change the billing methods of instances:
 - Change the billing method of an ECS instance from pay-as-you-go to subscription
 - o Change the billing method of an instance from subscription to pay-as-you-go

Improve cost-effectiveness

- You can purchase preemptible instances to reduce costs and use preemptible instances in conjunction with Auto Provisioning for automated provisioning of instances. For more information, see Create an auto provisioning group and Create a preemptible instance.
- You can purchase reserved instances to improve the flexibility of paying for instances and reduce costs. For more information, see Purchase reserved instances.

Create and manage disks

If you want to use a disk as a data disk, you can perform the following steps:

- 1. Create a disk.
- 2. Attach a data disk.
- 3. Partition and format a data disk on a Linux instance or Partition and format a data disk on a Windows instance.
- 4. Create a snapshot to back up data. For more information, see Create a snapshot of a disk.
- 5. If the storage capacity of an existing disk cannot meet your requirements, resize the disk. For more information, see the following topics:
 - Resize disks online for Linux instances
 - Resize disks offline for Linux instances
 - Resize disks online for Windows instances
 - Resize disks offline for Windows instances
- 6. If a data error occurs on a disk, use a snapshot created at a specific point in time to roll back the disk. For more information, see Roll back a disk by using a snapshot.
- 7. If you want to restore a disk to its initial state, re-initialize the disk. For more information, see Re-initialize a data disk.
- 8. Detach a data disk.
- 9. Release a disk.

Create and manage snapshots

You can perform the following steps to use a snapshot:

1. Create a snapshot. You can use one of the following methods to manually or automatically create

a snapshot:

- o Create a snapshot of a disk.
- Use an automatic snapshot policy to automatically create snapshots on a regular basis. For more information, see Apply or disable an automatic snapshot policy.
- 2. View the snapshot size.
- 3. Delete snapshots that are no longer needed to save storage space. For more information, see Reduce snapshot fees.

The following section describes the use scenarios of snapshots:

- To copy or restore data, you can use a snapshot to create or roll back a disk. For more information, see Create a disk from a snapshot and Roll back a disk by using a snapshot.
- To deploy an environment, you can use a system disk snapshot to create a custom image and then use the custom image to create instances. For more information, see Create a custom image from a snapshot and Create an ECS instance by using a custom image.

Create and manage custom images

Only custom images can be managed in the ECS console. You can use a custom image to deploy a business environment in a quick manner. You can use one of the following methods to obtain a custom image:

- Create a custom image from a snapshot.
- Create a custom image from an instance.
- Use Packer to create a custom image.
- Copy custom images across regions. For more information, see Copy custom images.
- Share custom images across accounts. For more information, see the "Share a custom image" section in Share or unshare custom images.
- Import custom images.
- Create and import an on-premises image by using Packer.

You can export custom images to back up environments. For more information, see Export a custom image

Create and manage security groups

You can perform the following steps to create and manage a security group.



- 1. Create a security group.
- 2. Add a security group rule.
- 3. Add an ECS instance to a security group.
- 4. Delete a security group rule.
- 5. Delete a security group.

You can clone a security group across regions and network types to simplify business deployment. For more information, see Clone a security group.

If new security group rules disrupt your online business, you can restore all or some of the security group rules. For more information, see Restore security group rules.

Create and attach instance RAM roles

You can perform the following steps to create and attach an instance RAM role.

- 1. (Optional) Authorize a RAM user to manage an instance RAM role. For more information, see Authorize a RAM user to manage an instance RAM role.
- 2. Create and attach an instance RAM role. For more information, see Attach an instance RAM role to an ECS instance.
- 3. Replace the instance RAM role based on your requirements. For more information, see Replace an instance RAM role.

Create and manage SSH key pairs

You can perform the following steps to create and manage an SSH key pair:

- 1. Create an SSH key pair or Import an SSH key pair.
- 2. Bind an SSH key pair to an instance.
- 3. Connect to a Linux instance by using an SSH key pair.
- 4. Unbind an SSH key pair.
- 5. Delete an SSH key pair.

Create and manage ENIs

You can perform the following steps to create and manage an elastic network interface (ENI).



- 1. Create an ENI.
- 2. Bind an ENI to an instance or Bind an ENI when you create an instance.
- 3. (Optional) Configure a secondary ENI.
- 4. Assign secondary private IP addresses.
- 5. Unbind an ENI.
- 6. Delete an ENI.

Use tags

You can use tags to manage resources to enhance efficiency. You can perform the following steps to use a tag:

- 1. Add a tag.
- 2. Search for resources by tag.
- 3. Delete or unbind a tag.

Create and manage launch templates

Launch templates help you create ECS instances that have the identical configurations. You can perform the following steps to use a launch template:

1. Create a launch template.

- 2. Create a launch template version.
- 3. Delete a launch template and a template version.

Create and manage deployment sets

Deployment sets help you implement high availability for underlying applications. You can perform the following steps to use a deployment set:

- 1. Create a deployment set.
- 2. Create an ECS instance in a deployment set.
- 3. Change the deployment set of an instance.
- 4. Delete a deployment set.

Use Cloud Assistant

Cloud Assistant allows you to send remote commands to ECS instances without the need to use jump servers. You can perform the following steps to use Cloud Assistant:

- 1. (Optional) Manually install and configure the Cloud Assistant client on some ECS instances. For more information, see Install the Cloud Assistant client.
- 2. Create a command.
- 3. Run a command.
- 4. Query execution results and fix common problems.

2.Best practices for instance type selection

Before you select an instance type, you must understand the key features of each instance type and instance family. Based on your understanding of these features, you can select an instance type that best suits your needs and take full advantage of the elasticity and flexibility of ECS in scenarios such as insufficient resources, retirement of instance families or instance types, and the use of preemptible instances.

This topic describes how to select enterprise-level instance families, instead of entry-level instance families (also called shared instance families). For information about how to select entry-level instance families, see Shared instance families or Overview.

Learn about instance families

When you create an ECS instance, you want to select the most cost-effective and stable instance type that can suit your performance, price, and workload needs. Alibaba Cloud provides a variety of instance families that have different specifications of vCPUs, memory, network performance, and throughput and are suited to different business scenarios. An instance family is further divided into multiple instance types. Instance family names are in the ecs. <Instance family> format, and instance type names are in the ecs. <Instance family>.<Instance f

- ecs: is the product code of ECS.
- < Instance family>: consists of lowercase letters and numbers.
 - Lowercase letters are abbreviations that indicate the performance fields of instance families. Examples:
 - c: compute optimized (computational)
 - g: general purpose (general)
 - r: memory optimized (ram)
 - ne: enhanced network performance (network enhanced)
 - The generation to which each instance family of the same type belongs is identified by a number.
 A larger number indicates a newer generation. Instance families of newer generations are more cost-effective and deliver higher performance.
- n in <nx>large: indicates the number of vCPUs that the instance type has. The greater n is, the more vCPUs the instance type has.

For example, ecs.g6.2xlarge represents an instance type that belongs to the general purpose instance family g6 and has eight vCPUs. The g6 instance family is a newer-generation one than the g5 family.

Select instance types

You can use one of the following methods to view details about instance families and instance types to select an appropriate instance type:

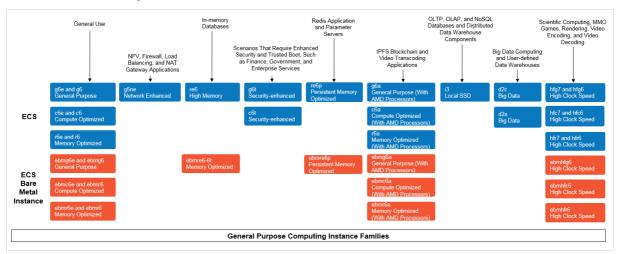
- Instance families: See this topic to learn details of instance families without logging on with an Alibaba Cloud account.
- DescribeInstanceTypes: Call this API operation to obtain the most recent performance specifications of instance types. You must log on with an Alibaba Cloud account to perform this operation.

aliyun ecs DescribeInstanceTypes --InstanceTypeFamily ecs.g6

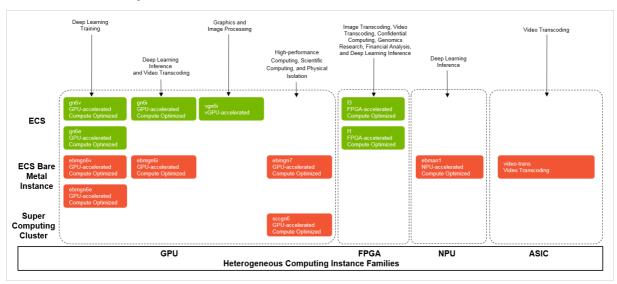
- Pricing tab of the ECS product page: Access this tab to view the pricing information and most recent special offers and estimate your costs.
- Custom Launch tab: Click this tab on the instance buy page of the ECS console to learn more about instance purchase instructions in the Instance Type section of the Basic Configurations step.

Select instance families based on scenarios

The following figure shows some ECS general purpose computing instance families and the business scenarios to which they are suited.

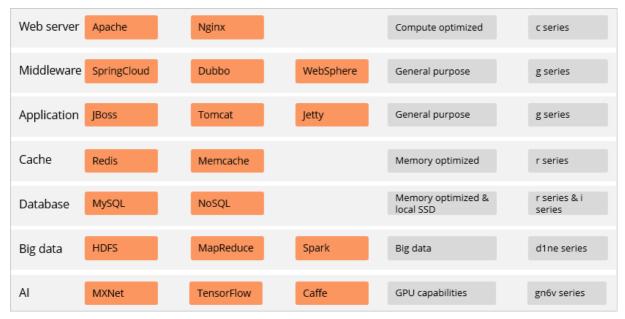


The following figure shows some ECS heterogeneous computing instance families and the business scenarios to which they are suited.



Select instance families based on applications

If you are using software or applications similar to the ones listed in the following figure, select corresponding instance families from the right side of the figure.



Select instance families based on user-defined services

If you are using user-defined services, select instance families based on your services and the selection principles.

Application type	Common application	Selection principle	Recommended instance family
Load balancing	NGINX	Connections can be frequently established. • CPU computing power: high. • Amount of memory: not large.	c6e, hfc7, and g5ne
RPC	SOFADubbo	A large amount of memory is available for network connection-intensive workloads.	g6e and g6
Cache	RedisMemcacheSolo	CPU computing power: not high.Amount of memory: large.	 r6e and re6 Elastic Block Storage: standard SSDs or enhanced SSDs (ESSDs)
Configuration center	ZooKeeper	A large number of I/O operations generated when applications initiate negotiations can be handled. • CPU computing power: not high. • Amount of memory: not large.	 c6e and c6 Elastic Block Storage: standard SSDs or ESSDs

Application type	Common application	Selection principle	Recommended instance family
Message queues	KafkaRabbitMQ	Cloud disks are preferred for message integrity. CPU computing power: not high. vCPU-to-memory ratio: 1:1. Storage: not high.	 c6e and c6 Elastic Block Storage: standard SSDs or ESSDs
Container orchestration	Kubernetes	The ECS bare metal instances and containers are combined to maximize computing power.	ebmc6e, ebmg6e, ebmc6, and ebmg6
Large table storage	HBase	 Typically, d series instance families are suitable. If your business requires ultra-high IOPS, you can select i series instance families. 	d2c and d2si3
	MySQLNoSQL	 If your business requires scalable storage, you can select ESSDs. If your business is I/O-sensitive, i series instance families are recommended. 	 c6e, g6e, and r6e Elastic Block Storage: ESSDs i3
Dat abase SQLServer		 Windows provides single-channel I/O configuration while high I/O capabilities are required. ESSDs are recommended. The logical and physical sectors of ECS instances are set to 4 KB in size. 	 c6e, g6e, and r6e Elastic Block Storage: ESSDs
Text search	Elasticsearch	 Instance types that have high vCPU-to-memory ratios are recommended. I/O capabilities can meet the requirements for exporting database data in the .es format. 	 g6e and g6 Elastic Block Storage: ESSDs d2c and d2s
Real-time computing	• Flink • Blink	You can select general purpose instance families and cloud disks or select d series instance families.	d2c and d2s
Offline computing	HadoopHDFSCDH	d series instance families are recommended.	d2c and d2s

Recommended instance families for common scenarios, game applications, and live streaming

These scenarios are CPU compute-intensive. We recommend that you select an instance type with a relatively balanced CPU-to-memory ratio such as 1:2, use an ultra disk as the system disk, and use standard SSDs or ESSDs as data disks. For scenarios that require higher network performance such as on-screen video comments, you can select an instance type with higher specifications to improve the packet forwarding rates.

Scenario category	Scenario	Recommended instance family	Performance requirement	CPU-to-memory ratio
	Balanced performance applications and backend applications	g series instance families, such as g6e	Medium clock speed, compute- intensive	1:4
Common scenarios	Applications with high packet forwarding rates	g series instance families, such as g6e	High packet forwarding rate, compute-intensive	1:4
	High-performance computing	hfc series instance families, such as hfc7	High clock speed, compute-intensive	1:2
Game applications	High-performance client games	hfc series instance families, such as hfc7	High clock speed	1:2
	Mobile or web games	g series instance families, such as g6e	Medium clock speed	1:4
Live streaming	Video forwarding	g series instance families, such as g6e	Medium clock speed, compute- intensive	1:4
	On-screen video comments	g series instance families, such as g6e	High packet forwarding rate, compute-intensive	1:4

Recommended instance families for big data scenarios such as Hadoop, Spark, and Kafka

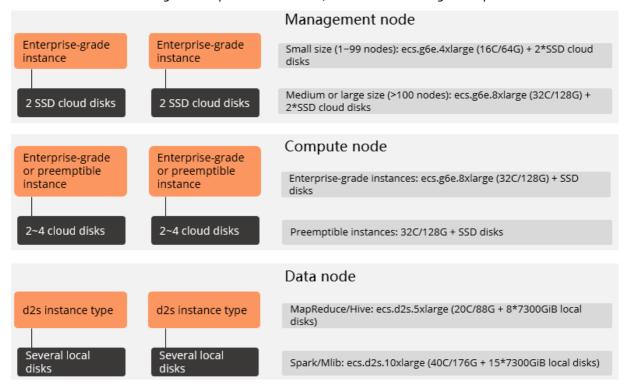
In these scenarios, the performance requirements are not the same for each node. You must balance the performance of each node, including the computing performance, storage throughput, and network performance.

- Management nodes: Select an instance family in the same manner as you would in common scenarios.
 For more information, see the Recommended instance families for common scenarios, game applications, and live streaming section.
- Compute nodes: Select an instance family in the same manner as you would in common scenarios. For more information, see the Recommended instance families for common scenarios, game applications,

and live streaming section. You must select instance types based on the cluster size. For example, you can select ecs.g6e.4xlage for a cluster that consists of less than 100 nodes and ecs.g6e.8xlage for a cluster that consists of more than 100 nodes.

Note Preemptible instances can be used as compute nodes to improve cost-effectiveness.
For more information, see Overview.

• Data nodes: require high storage throughput, high network throughput, and balanced CPU-to-memory ratios. We recommend that you use the d series big data instance families. For example, you can select ecs.d2s.5xlarge for MapReduce and Hive, and ecs.d2s.10xlarge for Spark and MLib.



Recommended instance families for databases, caches, and search scenarios

Typically, these scenarios require CPU-to-memory ratios greater than 1:4. Some software in these scenarios is sensitive to latency and storage I/O capabilities. We recommend that you select instance families whose memory is cost-effective.

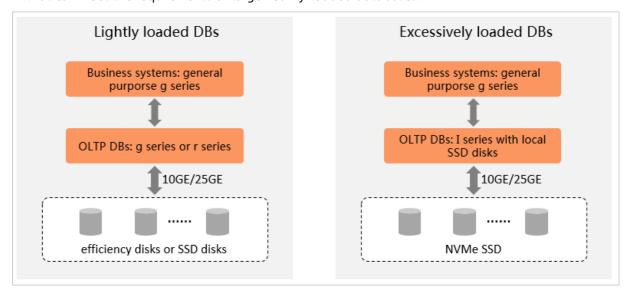
Scenario category	Scenario	Recommended instance family	CPU-to-memory ratio	Dat a disk
	High-performance and dependent on high availability in the application layer	i series instance families	1:4	Local SSDs, ultra disks, and standard SSDs

Relational detabris ecategory	Scenario	Recommended instance family	CPU-to-memory ratio	Data disk
	Small and medium-sized databases	g series instance families, or other instance families that have a CPU- to-memory ratio of 1:4	1:4	Ultra disks and standard SSDs
	High-performance databases	r series instance families	1:8	Ultra disks and standard SSDs
Distributed cache	Medium memory usage	g series instance families, or other instance families that have a CPU- to-memory ratio of 1:4	1:4	Ultra disks and standard SSDs
	High memory usage	r series instance families	1:8	Ultra disks and standard SSDs
	High performance and high availability in the application layer	i series instance families	1:4	Local SSDs, ultra disks, and standard SSDs
NoSQL database	Small and medium-sized databases	g series instance families, or other instance families that have a CPU- to-memory ratio of 1:4	1:4	Ultra disks and standard SSDs
	High-performance databases	r series instance families	1:8	Ultra disks and standard SSDs
ElasticSearch	Small clusters that rely on disks to ensure high data availability	g series instance families, or other instance families that have a CPU- to-memory ratio of 1:4	1:4	Ultra disks and standard SSDs
	Large clusters with high availability	d series instance families	1:4	Local SSDs, ultra disks, and standard SSDs

Databases are used in the following examples. Traditionally, business systems directly connect to online transaction processing (OLTP) databases and data redundancy is implemented by RAIDs. However, you can use ECS provided by Alibaba Cloud to flexibly deploy both lightly and heavily loaded databases.

• Lightly loaded databases: use enterprise-level instance families with cloud disks, which are more cost-effective.

• Heavily loaded databases: require high storage IOPS and low read/write latency. We recommend that you use i series instance families equipped with local SSDs. The local SSDs are high I/O NVMe SSDs that can meet the requirements of large heavily-loaded databases.

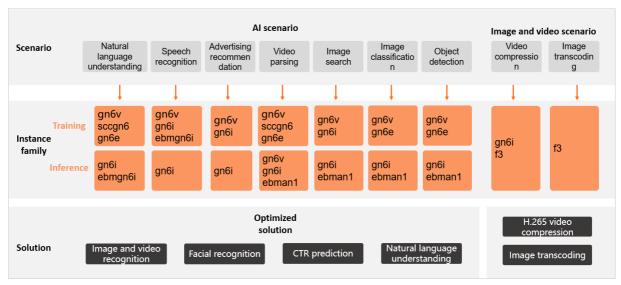


Recommended instance families for scenarios such as deep learning and image processing

In these scenarios, applications require high-performance GPU accelerators. The following GPU-to-CPU ratios are recommended for different scenarios:

- Deep learning training: A GPU-to-CPU ratio that ranges from 1:8 to 1:12 is recommended.
- General-purpose deep learning: A GPU-to-CPU ratio that ranges from 1:4 to 1:48 is recommended.
- Image recognition and inference: A GPU-to-CPU ratio that ranges from 1:4 to 1:12 is recommended.
- Speech recognition and synthesis inference: A GPU-to-CPU ratio that ranges from 1:16 to 1:48 is recommended.

The following figure shows the GPU- and FPGA-accelerated instance families recommended for common AI and image or video processing scenarios.



Check and adjust your selection

After you select an instance type to create an instance and start to use the instance, we recommend that you check whether the instance type is suitable based on the performance monitoring information.

Assume that you select the ecs.g6e.xlarge instance type. If you find that the CPU utilization is low, we recommend that you log on to the instance to check whether the memory usage is high. If the memory usage is high, you can change to an instance type that has a more suitable CPU-to-memory ratio within another instance family. For more information, see the following topics:

- ECS monitoring service
- View the monitoring data of a disk
- Overview

For scenarios such as insufficient resources, retirement of instance families, a change to a more cost-effective instance family, or an upgrade of instance configurations, you can change the instance configurations based on instance family features. For more information, see Overview of instance configuration changes and Instance families that support instance type changes.

3.Provisioning methods of ECS instances

Elastic Compute Service (ECS) instances can be provisioned by using multiple methods, such as individual provisioning, batch provisioning, high-availability deployment, and automatic cluster creation. You can use these provisioning methods by logging on to the ECS console or by calling API operations to create instances in various scenarios.

Manually create one or more instances

Scenarios: Multiple ECS instances of the same instance type and billing method need to be created in the same zone at a time.

Creation methods:

- Use the ECS console:
 - o Create an instance by using the wizard

Specify configurations in the instance creation wizard.

Create an ECS instance by using a custom image

Select a custom image within your Alibaba Cloud account and specify other configurations in the instance creation wizard.

• Purchase an ECS instance of the same configuration

Use the configurations of an existing instance and confirm the configurations in the instance creation wizard.

Create an instance by using a launch template

Select a launch template and confirm the configurations in the instance creation wizard.

- Call the Runinstances operation:
 - RunInstances
 - o Batch create ECS instances

Number of instances that can be created: If you create instances by using the ECS console, the number of instances that can be created at a time depends on your ECS usage. If you create instances by calling the Runinstances operation, you can create up to 100 instances per call.

Reboot Instance Create Instance Image Pending Starting Running Start Stop Delete Instance Instance Instance Delete Instance Deleted Stopped Stopping Stable Transitory

The following figure shows the lifecycle of an instance created by using the ECS console or by calling the Runinstances operation.

You can also call the <u>Createinstance</u> operation to create an ECS instance. The created instance enters the Stopped state, and you must manually start the instance.

Deploy instances across multiple physical servers to improve availability of the cluster (deployment set)

Scenarios: ECS instances need to be deployed across different physical servers to provide computing power for applications that require high availability and underlying disaster recovery.

Creation method: Create a deployment set and then create ECS instances in the deployment set. You can create instances by using the ECS console or by calling the RunInstances or CreateInstance operation.

Number of instances that can be created: depends on how instances are created. If you create instances by using the ECS console or by calling the Runinstances operation, you can create up to 20 instances at a time. If you create instances by calling the Createinstance operation, you can create a single instance at a time.

Limits:

- Up to 20 ECS instances can be created in each deployment set within a zone.
- Only specific ECS instance types are supported. For more information, see Overview.
- Billing methods: Subscription and pay-as-you-go are supported. Preemptible instances are not supported.

Procedure:

- Use the ECS console:
 - i. Create a deployment set
 - ii. Create an ECS instance in a deployment set
- Call API operations:
 - i. CreateDeploymentSet
 - ii. RunInstances or CreateInstance

Automatically create instance clusters at minimal costs (Auto Provisioning)

Scenarios: Instance clusters that use different billing methods need to be deployed across instance types and zones. This method is suitable for scenarios where stable computing power must be provisioned in a quick manner and preemptible instances are used to reduce costs.

Creation method: Create an auto provisioning group to automatically batch create ECS instances.

Number of instances that can be created: Up to 1,000 instances can be created by an auto provisioning group.

Limits: Pay-as-you-go instances and preemptible instances are supported, whereas subscription instances are not.

Procedure:

- Use the ECS console: Create an auto provisioning group
- Call an API operation: CreateAutoProvisioningGroup
- Best practices for calling an API operation related to Auto Provisioning to batch create ECS instances: Use auto provisioning group-related API operations to create multiple ECS instances at the same time

Automatically create and release instances (Auto Scaling)

Scenarios: Instance clusters that use different billing methods need to be maintained across instance types and zones when service loads fluctuate from time to time.

Creation method: Create a scaling group and a scaling task. The scaling group automatically batch creates or releases ECS instances.

Number of instances that can be created:

- Up to 1,000 ECS instances can be created during a single scaling activity.
- Up to 1,000 ECS instances can be created by a single scaling group.

Limits: Pay-as-you-go instances and preemptible instances are supported. You can manually add an existing subscription instance to a scaling group, but cannot create a subscription instance in the scaling group.

Procedure:

- Use the ECS console:
 - i. Scale ECS instances
 - ii. Automatically add ECS instances or Automatically remove ECS instances
- Call API operations:
 - i. CreateScalingGroup
 - ii. CreateScalingConfiguration
 - iii. CreateScalingRule
 - iv. CreateScheduledTask

Auto Scaling also provides simplified features to make provisioning more efficient and shorten the lead time from requirement to provisioning. For example, you can configure ECS instances in scaling groups to be automatically associated with Server Load Balancer (SLB) instances and ApsaraDB RDS instances. You can also configure lifecycle hooks to perform custom operations on the ECS instances in scaling groups. You can use Auto Scaling to obtain the scalability that meets your business needs. For information about the best practices for Auto Scaling, see the following topics:

- Build a scalable web application
- Save your money with Auto Scaling
- Deploy a high-availability compute cluster

4.Batch configure sequential names or hostnames for multiple instances

You can create multiple ECS instances by using the ECS console or by calling the RunInstances operation. When you create multiple ECS instances, you can customize the instance names or hostnames to facilitate subsequent instance management. This topic describes how to batch configure sequential instance names or hostnames for multiple instances.

Background information

You can batch configure sequential instance names or hostnames by specifying a sequence or by using automatic sorting.

This topic describes how to use the ECS console or call an API operation to configure sequential instance names and host names for three instances in four scenarios.

- By using the ECS console:
 - Scenario 1: Set the instance names or hostnames of three instances to be sorted in a specified sequence by using the ECS console
 - Scenario 2: Set the instance names or hostnames of three instances to be sorted in an automatic sequence by using the ECS console
- By calling an API operation:
 - Scenario 3: Set the instance names or host names of three instances to be sorted in a specified sequence by calling the Runinstances operation
 - Scenario 4: Set the instance names or hostnames of three instances to be sorted in an automatic sequence by calling the Runinstances operation

For more information about specific configuration rules, see the Naming conventions, Specify a sequence, and Automatic sorting sections.

Scenario 1: Set the instance names or hostnames of three instances to be sorted in a specified sequence by using the ECS console

In this example, instance names or hostnames are set to be sorted in the specified sequence in the ECS console. For more information about other configurations, see Create an instance by using the wizard.

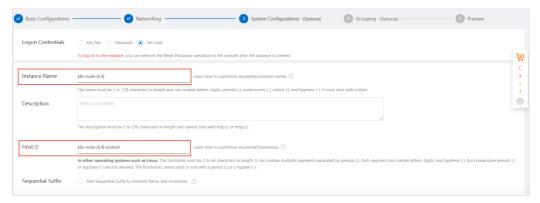
- 1. Go to the Custom Launch tab in the ECS console.
- 2. Complete the configurations in the **Basic Configurations** and **Networking** steps. In this example, the instance quantity is set to 3 on the **Basic Configurations** step.
- 3. In the ${\bf System\ Configurations}$ step, configure the parameters.

Set the values of Instance Name and Host in the following format: name_prefix[begin_number,bits] name_suffix. For more information about how to specify sorting rules, see the Specify a sequence section.

In this example, both the instance names and host names are specified to start with k8s-node- and increment starting from 0006, and the host names are specified to end with -ecshost. Set Instance Name to k8s-node-[6,4] and Host to k8s-node-[6,4]-ecshost.

? Note In this example, the sorting rules only specify instance names and hostnames. Sequential Suffix is disabled by default.

Example of specifying a sequence



4. Complete the configurations in the **Grouping** step and confirm the order.

You can view the new instances on the **Instances** page. In this example, the generated instance names are k8s-node-0006, k8s-node-0007, and k8s-node-0008, and the generated host names are k8s-node-0006-ecshost, k8s-node-0007-ecshost, and k8s-node-0008-ecshost.

Scenario 2: Set the instance names or hostnames of three instances to be sorted in an automatic sequence by using the ECS console

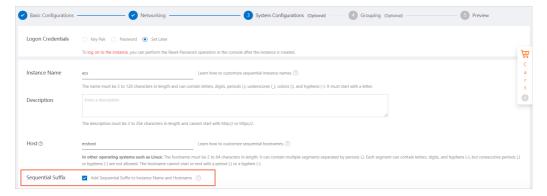
In this example, instance names or host names are set to be automatically sorted in the ECS console. For more information, see Create an instance by using the wizard.

- 1. Go to the Custom Launch tab in the ECS console.
- Complete the configurations in the Basic Configurations and Networking step.
 In this example, the instance quantity is set to 3 on the Basic Configurations tab.
- 3. In the **System Configurations** step, configure the parameters.

Select Sequential Suffix. The system adds sequential suffixes to the Instance Name and Host values and sorts the values. The added suffix starts from 001 and increments with each instance. For rules of automatic sorting, see the Automatic sorting section.

In this example, set Instance Name to ecs and Host to ecshost.

Example of automatic sorting



4. Complete the configurations in the **Grouping** step and confirm the order.

You can view the new instances on the **Instances** page. In this example, the generated instance names are ecs001, ecs002, and ecs003, and the generated hostnames are ecshost001, ecshost002, and ecshost003.

Scenario 3: Set the instance names or hostnames of three instances to be sorted in a specified sequence by calling the RunInstances operation

This section describes how to configure parameters used to specify a sequence. For more information about other parameters, see RunInstances.

Set the values of InstanceName and HostName in the following format: name_prefix[begin_number,bits]name_suffix. For more information about how to specify sorting rules, see the Specify a sequence section.

In this example, the instance names and hostnames of the three instances are specified to start with k8s-node- and increment starting from 0006, and the hostnames are specified to end with -ecshost. Configure the following parameters:

• Amount: 3

InstanceName: k8s-node-[6,4]HostName: k8s-node-[6,4]-ecshost

Note In this example, the sorting rules only specify instance names and host names. UniqueSuffix is disabled by default.

In this example, the generated instance names are k8s-node-0006, k8s-node-0007, and k8s-node-0008, and the generated hostnames are k8s-node-0006-ecshost, k8s-node-0007-ecshost, and k8s-node-0008-ecshost.

Scenario 4: Set the instance names or hostnames of three instances to be sorted in an automatic sequence by calling the RunInstances operation

This section describes how to configure parameters used to specify the automatic sorting. For more information about other parameters, see Runinstances.

Set UniqueSuffix to *true*. The system adds sequential suffixes to the InstanceName and HostName values and sorts the values. The added suffix starts from 001 and increments with each instance. For rules of automatic sorting, see the Automatic sorting section.

In this example, the instance names or hostnames of three instances are set to be automatically sorted. Configure the following parameters:

• Amount: 3

InstanceName: ecsHostName: ecshostUniqueSuffix: true

In this example, the generated instance names are ecs001, ecs002, and ecs003, and the generated host names are ecshost001, ecshost002, and ecshost003.

Naming conventions

- Instance name: The name must be 2 to 128 characters in length and must start with a letter. It can contain letters, digits, periods (.), underscores (_), colons (:), and hyphens (-).
- Hostname:
 - For Windows instances, the hostname must be 2 to 15 characters in length and can contain letters, digits, and hyphens (-). It cannot start or end with a hyphen (-), contain consecutive hyphens (-), or contain only digits.
 - For other operating systems such as Linux, the hostname must be 2 to 64 characters in length. You can use periods (.) to separate a name into multiple segments. Each segment can contain letters, digits, and hyphens (-). However, the hostname cannot contain consecutive periods (.) or hyphens (-). It cannot start or end with a period (.) or a hyphen (-).

Specify a sequence

Set the parameter value in the following format: name_prefix[begin_number,bits]name_suffix.

Fields

Field	Description	Example
name_prefix	The prefix of instance names or hostnames. Note A prefix is required when you specify sequential names. Otherwise, names are regarded as common names.	k8s-node-
[begin_number,bit s]	The ordered numeric values for instance names or hostnames. After you specify this field, the numbers in the names increment in sequence. • begin_number: the number from which the ordered numeric values of instances start. Valid values: 0 to 999999. Default value: 0. • bits: the number of digits of an ordered value. Valid values: 1 to 6. Default value: 6. Notice • The [begin_number,bits] field cannot contain spaces. • If the value of begin_number has more digits than the value of bits, the value of bits is 6. • A maximum of 999,999 ECS instances can share the same prefix and suffix. Extra instances use 999999 as the prefix.	[0,6]
name_suffix	The suffix of the instance names or hostnames.	-ecshost

Parameter examples

Example	Generated names (for three instances)
k8s-node-[]-ecshost or k8s-node-[,]-ecshost	k8s-node-000000-ecshost, k8s-node-000001- ecshost, and k8s-node-000002-ecshost
k8s-node-[99]-ecshost or k8s-node-[99,]-ecshost	k8s-node-000099-ecshost, k8s-node-000100- ecshost, and k8s-node-000101-ecshost
k8s-node-[99,1]-ecshost	k8s-node-000099-ecshost, k8s-node-000100- ecshost, and k8s-node-000101-ecshost
k8s-node-[999998]-ecshost	k8s-node-999998-ecshost, k8s-node-999999- ecshost, and k8s-node-99999-ecshost
k8s-node-[0,4]	k8s-node-0000, k8s-node-0001, and k8s-node-0002

Automatic sorting

When you create multiple instances, you can enable the automatic sorting feature to automatically add sequential suffixes to the instance names and host names. The sequential numbers in the suffix range from 001 to 999.



Note By default, the automatic sorting feature is disabled.

Parameter examples

Format (instance name or hostname)	Example	Generated names (for three instances)
Common names	ecs	ecs001, ecs002, and ecs003
Sequential names to be sorted in		k8s-node-000000-ecshost001, k8s-node-000001-ecshost002, and k8s-node-000002- ecshost003
a specified sequence: name_prefix[begin_number,bits]n ame_suffix	k8s-node-[]-ecshost or k8s-node- [,]-ecshost	Note The specified sequence and automatic sorting take effect at the same time.
		k8s-node-0000, k8s-node-0001, and k8s-node-0002
Sequential names to be sorted in a specified sequence: name_prefix[begin_number,bits]	k8s-node-[0,4]	Note <i>name_suffix</i> is not specified, and automatic sorting does not take effect.

5.Use auto provisioning grouprelated API operations to create multiple ECS instances at the same time

You can use ECS API operations to create multiple pay-as-you-go ECS instances at the same time. Compared with the RunInstances operation, the CreateAutoProvisioningGroup operation can help you create a large number of ECS instances at the same time in a simpler and more stable manner.

Context

RunInstances is commonly used when your business requires you to create multiple pay-as-you-go ECS instances. You can call RunInstances to create up to 100 ECS instances at the same time. However, if more than 100 ECS instances are required to be created at the same time, you may encounter technical bottlenecks in use of RunInstances. For more information, see Possible issues caused by using RunInstances.

Note If you are familiar with the technical bottlenecks caused by using Runinstances, you can skip the preceding section.

To meet your requirements of creating more than 100 ECS instances at the same time, Alibaba Cloud provides auto provisioning groups. You can call the CreateAutoProvisioningGroup operation to create an auto provisioning group and use the group to deploy an instance cluster across different billing methods, instance families, and zones within one click. CreateAutoProvisioningGroup is more suitable than RunInstances when you create a large number of ECS instances. For more information, see Comparison between RunInstances and CreateAutoProvisioningGroup and Benefits of auto provisioning groups.

Comparison between RunInstances and CreateAutoProvisioningGroup

The following table compares the features of the RunInstances and CreateAutoProvisioningGroup operations to help you choose an appropriate operation to create ECS instances.

ltem	RunInstances	CreateAutoProvisioningGroup
Maximum number of instances that can be created at the same time	100.	1,000 (up to 10,000 vCPUs).
Capacity delivery form	Specific number of instances.	Specific number of instances, specific number of vCPUs, and instance types of different weights
Support for multiple zones	No.	Yes.

Item	RunInstances	CreateAutoProvisioningGroup
Support for multiple instance types	No.	Yes.
Support for multiple disk categories	No.	Yes.
Provision of policies that you can use to create instances	No.	Yes. The following policies are provided: • For pay-as-you-go instances • Cost optimization policy: The auto provisioning group selects the lowest-cost instance types from the candidate instance types to create instances. • Priority-based policy: The auto provisioning group attempts to create instances based on the priorities configured for candidate instance types. • For preemptible instances • Cost optimization policy: The auto provisioning group selects the lowest-cost instance types from the candidate instance types to create instances. • Balanced distribution policy: The auto provisioning group evenly distributes instances across the specified zones. • Capacity-optimized distribution policy: The auto provisioning group selects the optimal combinations of instance types and zones to create instances based on resource availability.
Delivery stability	Affected by resource availability.	Multiple combinations of instance types and zones reduce the impacts of resource availability.
Response method	Returns a response in real time.	Returns a response in real time.

The following example scenarios show you how to use CreateAutoProvisioningGroup in place of RunInstances to create ECS instances:

• Assume that you called the Runinstances operation to batch create instances of a single instance

type in a single zone. When you call the CreateAutoProvisioningGroup operation instead to batch create instances, you need only to configure a combination of instance types and zones.

• Assume that you manually configure a business deployment plan when you called the Runinstances operation to batch create instances. You can call the CreateAutoProvisioningGroup operation instead to deploy instances that have multiple categories of disks attached across multiple instance types and zones, based on different instance provisioning policies provided by Auto Provisioning.

For example, you manually configured a business deployment plan that traversed multiple instance types and zones to create instances by calling the RunInstances operation. This improves the success rate of creating instances. If you call the CreateAutoProvisioningGroup operation instead to batch create a large number of instances, you need only to configure multiple combinations of instance types and zones, and select an appropriate policy. The created auto provisioning group automatically batch creates instances based on your configurations and policy.

Notice Limits apply to policies that are used for auto provisioning groups to create instances. A maximum of 1,000 instances can be created at the same time. If you specify the weight of each specified instance type by using the weightedCapacity parameter, the maximum weighted capacity that is created one time is 10,000.

Possible issues caused by using RunInstances

Due to the limits of the RunInstances operation, you may encounter the following issues when you call RunInstances to create instances.

Issue	Description	Solution
A limited number of instances can be created at the same time.	You can call the Runinstances operation to create a maximum of 100 ECS instances at the same time.	If you want to create more than 100 ECS instances, you must call this operation multiple times in a recurring or concurrent manner.

Issue	Description	Solution
The success rate of creating instances is not ensured.	You can call the RunInstances operation to create instances of a single instance type in a single zone. When you call this operation to create multiple ECS instances at the same time, you may encounter the following issues: Instances cannot be created due to insufficient resources of a specified instance type. Instances of a specified instance type can no longer be created due to the retirement of the instance type. Instances cannot be created because the specified instance type is unavailable in the specified zone. Instances cannot be created because the specified instance type do not support specified disk categories.	Instances cannot be created mostly due to insufficient resources. To solve this issue, we recommend that you call the DescribeAvailableResource operation to query available instance types in zones before you call the RunInstances operation to create ECS instances. You can manually configure multiple combinations of instance types and zones to have sufficient resources to create instances. The complex method of creating instances helps ensure high business delivery stability. After you configure multiple combinations of instance types and zones, you must create an appropriate policy to create instances. For example, you can create 100 ECS instances in sequence based on multiple combinations that you configured. If only 50 ECS instances can be created by using the available resources provided by the first combination, you must use the second combination to create the remaining 50 ECS instances. Instance types have limits. You can call the DescribeAvailableResource operation to query the limits. Then, you can develop fault tolerance solutions based on the limits to avoid the impacts caused by changes in limits. Note You can also learn about the limits of instance types based on the descriptions in documents about instance families. For more information, see Instance family. Example scenario: The ecs.g6e.large instance type supports only the disk category of enhanced SSDs (ESSDs) and Beijing Zone X (cn-beijing-x) does not support ESSDs.

Issue	Description	Solution
Simple policies that are used to create instances cannot meet diversified business requirements.	Runinstances can only be used to create instances of a single instance type in a single zone. If you need to deploy instances in multiple zones to implement geo-disaster recovery or create ECS instances at minimal costs, you must create a business deployment plan on your own to ensure that instances are deployed. A business deployment plan that you create on your own may have the following issues: • High development costs. You may need to deal with a series of issues when you use the business deployment plan that you configured. For example, you may encounter issues about how to create ECS instances when resources are insufficient and how to ensure computing power at the lowest cost of preemptible instances when you scale up your servers. • Low stability and less professionalism. A business deployment plan that you built on your own for the resources provided by Alibaba Cloud is less likely to be built in a professional manner. You also cannot test the plan. This brings risks to the production environment.	To resolve the issue, you may contact Alibaba Cloud for assistance.

Benefits of auto provisioning groups

Alibaba Cloud provides auto provisioning groups to address the issues that you may encounter when you call the RunInstances operation to create multiple ECS instances at the same time. Auto provisioning groups allow you to deploy an instance cluster across different billing methods, instance families, and zones within one click. Auto provisioning groups provide stable computing power, alleviate the instability caused by recycled preemptible instances, and eliminate the need to manually create instances. The following table describes the benefits of auto provisioning groups.

Benefit	Description
Allow you to create a large number of ECS instances at the same time.	You can use an auto provisioning group to create up to 1,000 ECS instances at the same time.

Benefit	Description
Allow you to create instances that have multiple categories of disks attached across instance types and zones.	Auto provisioning groups allow you to configure up to 10 combinations of instance types and zones, and allow you to select up to five disk categories. This ensures high availability when you create ECS instances at the same time. Example scenarios: When you create ECS instances based on the balanced distribution policy provided by an auto provisioning group, you can configure multiple instance types and zones. As required by the policy, instances need to be created in a relatively balanced manner across zones. However, if the required number of instances cannot be created in a zone, the auto provisioning group attempts to create the same number of instances in other zones. If you specify multiple disk categories, the auto provisioning group creates disks for instances based on the priorities of the disk categories, and changes a disk category if the disk category is unavailable. Plate If all disk categories are unavailable, the auto provisioning group no longer uses this method but selects another creation method.
Support multiple policies to create instances	 Auto provisioning groups provide the following policies to create instances: For pay-as-you-go instances Cost optimization policy: The auto provisioning group selects the lowest-cost instance types from the candidate instance types to create instances. Priority-based policy: The auto provisioning group attempts to create instances based on the priorities configured for candidate instance types. For preemptible instances Cost optimization policy: The auto provisioning group selects the lowest-cost instance types from the candidate instance types to create instances. Balanced distribution policy: The auto provisioning group evenly distributes instances across the specified zones. Capacity-optimized distribution policy: The auto provisioning group selects the optimal combinations of instance types and zones to create instances based on resource availability.

Benefit	Description
Improve availability of preemptible instances	The demand for preemptible instances is increasing due to the advantages in prices of preemptible instances. However, the prices of preemptible instances fluctuate, and the instances are recycled when bidding prices become lower than current market prices after protection periods. This makes preemptible instances difficult to manage. Auto provisioning groups improve the availability of preemptible instances while low costs are maintained. The following items describe how to improve the availability of preemptible instances:
	 You can use the default cost optimization policy. If you use this policy, the auto provisioning group attempts to create instances in ascending order of instance type prices each time the group scales out.
	 For preemptible instances, the resources based on instance types are different from those based on zones. You can configure multiple combinations of instances types and zones for preemptible instances to reduce the probability that no resources are available for all the combinations.
	 You can configure multiple disk categories when you create an auto provisioning group. This ensures that the system can select appropriate disk categories to create instances.
	• You can configure the SpotInstancePoolsToUseCount parameter to ensure that preemptible instances are created by using the lowest-cost combinations of instance types and zones. This can avoid the issue that computing power significantly reduces when instances of an instance type are recycled.

Best practices for creating auto provisioning groups by calling CreateAutoProvisioningGroup

This section provides the sample Java code that is used to create an auto provisioning group (CreateAutoProvisioningGroup). This section also describes how to call the CreateAutoProvisioningGroup operation.

- 1. Install ECS SDK for Java and the Alibaba Cloud SDK core library. For more information, see Install ECS SDK for Java.
- 2. Write the Java code that is used to call the CreateAutoProvisioningGroup operation. Sample code:

```
CreateAutoProvisioningGroupRequest request = new CreateAutoProvisioningGroupRequest();
request.setRegionId(regionId);
request.setLaunchConfigurationImageId(RequestHelper.IMAGE ID);
request.setLaunchConfigurationSecurityGroupId(securityGroupId);
request.setTotalTargetCapacity(totalTargetCapacity);
request.setPayAsYouGoTargetCapacity(payAsYouGoTargetCapacity);
request.setSpotTargetCapacity(spotTargetCapacity);
request.setLaunchConfigurationSystemDiskCategory("cloud ssd");
request.setLaunchConfigurationSystemDiskSize(40);
request.setAutoProvisioningGroupType("instant");
// Configure the policy that is used to create preemptible instances.
request.setSpotAllocationStrategy("lowest-price");
request.setSpotInstancePoolsToUseCount(spotInstancePoolsToUseCount);
// Configure the policy that is used to create pay-as-you-go instances.
request.setPayAsYouGoAllocationStrategy("prioritized");
request.setMaxSpotPrice(maxSpotPrice);
// Configure a maximum of 10 combinations of instance types and zones.
request.setLaunchTemplateConfigs(launchTemplateConfigs);
request.setClientToken(clientToken);
CreateAutoProvisioningGroupResponse response = client.getAcsResponse(request);
```

Sample response in the JSON format:

```
"autoProvisioningGroupId": "apg-***",
"launchResults":[
        "instanceIds":[
            "i_****
        ],
        "instanceType":"ecs.c5.large",
        "spotStrategy": "NoSpot",
        "zoneId": "cn-shanghai-b"
    },
        "instanceIds":[],
        "instanceType": "ecs.c5.large",
        "spotStrategy": "NoSpot",
        "zoneId": "cn-shanghai-b",
        "errorCode" : "Invalid.Parameter",
        "errorMsg" : "Specific Parameter 'imageId' is not valid"
"requestId": "20DA1E9F-BF7F-4BE7-8204-E4DE58E4FC7B"
```

When you call the CreateAutoProvisioningGroup operation to create an auto provisioning group, you need only to configure items that are used to create a large number instances at the same time, without the need to worry about the creation process. The auto provisioning group creates instances in a diligent manner.

? Note If instances are created in a diligent manner, the system attempts to create instances by using another combination of instance types and zones when instances cannot be created by using one combination. This method requires more time to create instances. In addition, the actual creation result may deviate from the creation policy that you configured.

Related information

- Overview
- Configure an auto provisioning group

6.Best practices for heterogeneous computing services

This topic describes the best practices for heterogeneous computing services. Select from the following topics based on your business scenario to learn about the associated best practices.

Elastic GPU Service

• Deploy an NGC environment on a GPU-accelerated instance

Describes how to deploy a NVIDIA GPU Cloud (NGC) environment on a GPU-accelerated instance. In the example, the TensorFlow deep learning framework is used.

• Use RAPIDS to accelerate machine learning tasks on a GPU-accelerated instance

Describes how to use the NGC-based Real-time Acceleration Platform for Integrated Data Science (RAPIDS) libraries that are installed on a GPU-accelerated instance to accelerate data science tasks and machine learning tasks and improve the efficiency of using computing resources.

FPGA as a Service (FaaS)

- Best practices for the Register Transfer Level (RTL) design on FPGA-accelerated instances
 - Use RTL Compiler on an f1 instance

Describes how to implement the RTL design on an f1 instance.

o Project modes and directories used by RTL

Describes the project modes and directories used by the RTL compiler and provides a sample framework to help you understand how to use RTL.

• Use the RTL design on an f3 instance

Describes how to implement the RTL design on an f3 instance.

- Best practices for using Open Computing Language (OpenCL) on FPGA-accelerated instances
 - Use OpenCL on an f1 instance

Describes how to use OpenCL on an f1 instance to create an image and burn the image to an FPGA.

Overview of the FaaS f3 SDAccel development environment

Describes the FaaS f3 SDAccel development environment. The FaaS f3 SDAccel development environment is based on Xilinx SDAccel dynamic 5.0. You can develop and apply the FaaS f3 SDAccel development environment based on OpenCL.

• Use OpenCL on an f3 instance

Describes how to use OpenCL on an f3 instance to create an image and burn the image to an FPGA.

7.Best practices for preemptible instances

7.1. Query the price of a preemptible instance

If you are a developer, you can refer to the sample Java code provided in this topic to query the latest price of a preemptible instance.

Prerequisites

• An Alibaba Cloud account is created and its AccessKey pair is obtained.

You must configure the AccessKey pair of your Alibaba Cloud account to use Alibaba Cloud Elastic Compute Service (ECS) SDK for Java. For more information about how to obtain an AccessKey pair, see Obtain an AccessKey pair.

• ECS SDK for Java is installed in the development environment.

You must add the following dependencies to the Maven project. For more information, see Install ECS SDK for Java.

```
<dependencies>
       <dependency>
          <groupId>com.aliyun</groupId>
           <artifactId>aliyun-java-sdk-ecs</artifactId>
          <version>4.23.10
       </dependency>
       <dependency>
          <groupId>com.aliyun</groupId>
          <artifactId>aliyun-java-sdk-core</artifactId>
          <version>4.0.8
       </dependency>
       <dependency>
          <groupId>commons-lang
          <artifactId>commons-lang</artifactId>
          <version>2.6</version>
       </dependency>
       <dependency>
          <groupId>com.alibaba
          <artifactId>fastjson</artifactId>
           <version>1.2.68
       </dependency>
</dependencies>
```

Sample code

In this topic, the QuerySpotLatestPrice Java class is used. The following code provides an example on how to query the latest price of a preemptible instance by calling the DescribePrice operation.

```
import com.aliyuncs.DefaultAcsClient;
import com.alivuncs.IAcsClient;
```

```
import com.aliyuncs.ecs.model.v20140526.DescribePriceRequest;
import com.aliyuncs.ecs.model.v20140526.DescribePriceResponse;
import com.aliyuncs.profile.DefaultProfile;
 * Call the DescribePrice operation to query the latest price of a preemptible instance.
public class QuerySpotLatestPrice {
   private static IAcsClient client;
    // Specify the region ID of the preemptible instance.
   static String regionId = "cn-hangzhou";
   // Set the network type of the preemptible instance to Virtual Private Cloud (VPC).
   static String resourceType = "instance";
    static String instanceNetworkType = "vpc";
   // Specify the instance type of the preemptible instance.
   static String instanceType = "ecs.g6.8xlarge";
   // Set the preemption policy to SpotAsPriceGo.
   static String spotStrategy = "SpotAsPriceGo";
    // Set spotDuration to the protection period within which to retain the preemptible ins
tance. If you cannot determine the protection period, set the value to 0.
   static Integer spotDuration = 1;
    // Specify the zone ID of the preemptible instance.
    static String zoneId = "cn-hangzhou-i";
    public static void main(String[] args) throws Exception {
       client = Initialization();
       describePrice(client);
    public static void describePrice(IAcsClient client) throws Exception {
        // Set the parameters of the DescribePrice operation and send the request.
        DescribePriceRequest request = new DescribePriceRequest();
       request.setRegionId(regionId);
       request.setResourceType(resourceType);
        request.setInstanceType(instanceType);
        request.setInstanceNetworkType(instanceNetworkType);
        request.putQueryParameter("spotStrategy", spotStrategy);
        request.putQueryParameter("spotDuration", spotDuration);
        request.putQueryParameter("zoneId", zoneId);
        // Obtain the response and the latest price of the preemptible instance.
        DescribePriceResponse describePriceResponse = client.getAcsResponse(request);
        System.out.println("Preemptible instance price:"+describePriceResponse.getPriceInfo
().getPrice().getTradePrice()+ "Yuan");
   private static IAcsClient Initialization() {
         * Initialize request parameters.
         * Set the <your-access-key-id> variable to the AccessKey ID of your Alibaba Cloud
account.
         * Set the <your-access-key-secret> variable to the AccessKey secret of your Alibab
a Cloud account.
        DefaultProfile profile = DefaultProfile.getProfile(regionId, "<your-access-key-id>"
, "<your-access-key-secret>");
       return new DefaultAcsClient(profile);
    }
```

The following figure shows a sample response.

Preemptible instance price:7.9298Yuan

Process finished with exit code 0

7.2. Select a bidding mode for a preemptible instance

Preemptible instances support multiple bidding modes. You can select a bidding mode based on your business requirements. This topic compares and analyzes different bidding modes of preemptible instances and provides best practices for bidding modes.

Context

Preemptible instances are a type of on-demand instances that are offered at a discounted price compared to pay-as-you-go instances. Preemptible instances are designed to minimize Elastic Compute Service (ECS) instance costs in specific scenarios. Preemptible instances support the following bidding modes:

- Set Maximum Price (SpotWithPriceLimit)
- Use Automatic Bid (SpotAsPriceGo)

This topic also provides a bidding mode that combines Use Automatic Bid (SpotAsPriceGo) with the ACS-ECS-AlarmWhenDiscountAndPriceExceedsThresholdInMultiZoneAndInstanceType public template of Operation Orchestration Service (OOS). This mode helps control costs while reducing the probability of preemptible instance interruptions.

Comparison of bidding modes

• Mode 1: Set Maximum Price (SpotWithPriceLimit)

You must specify a maximum hourly price to bid for the instance type. If the spot price (current market price per hour) of the instance type exceeds your specified maximum hourly price due to price fluctuations, a preemptible instance interruption event is triggered.

Analysis:

- Advantage: Costs are controlled to ensure that the instance type fees do not exceed the maximum hourly price that you specify.
- Disadvantage: If the instance type price drastically fluctuates, the probability of instance interruptions increases and instance stability decreases.
- Applicable scenario: business that has strict budget and price requirements for ECS instances, of which the budget cannot be exceeded.
- Mode 2: Use Automatic Bid (SpotAsPriceGo)

In this mode, the preemptible instance is billed based on the spot price.

Analysis:

• Advantage: The instance is not interrupted even if the instance type price drastically fluctuates, and instance stability is improved.

- Disadvantages: It is difficult to control costs. Increases in instance type prices cannot be detected and costs can spiral uncontrollably.
- Applicable scenario: business that has no strict cost requirements and requires to reduce costs as much as possible while improving instance stability.
- Mode 3: Combination of Use Automatic Bid (SpotAsPriceGo) and OOS

In the Use Automatic Bid (SpotAsPriceGo) mode, increases in instance type prices cannot be detected. When the ACS-ECS-AlarmWhenDiscountAndPriceExceedsThresholdInMultiZoneAndInstanceT public template of OOS is combined with Use Automatic Bid (SpotAsPriceGo), OOS sends a notification when the price of the preemptible instance exceeds the threshold of the instance type that you specify. This way, you can manage the preemptible instance.

Analysis:

- Advantage: Increases in instance type prices can be detected and the instance is stable.
- Disadvantage: The Use Automatic Bid (SpotAsPriceGo) mode must be combined with OOS, which increases costs.
- Applicable scenario: business that requires the ability to detect price increases while improving instance stability.

The following table compares the bidding modes.

Bidding mode	Instance interruption probability	Instance stability	Cost optimization	Cost controllability
Set Maximum Price (SpotWithPriceLimit)	High	Low	High	High
Use Automatic Bid (SpotAsPriceGo)	Low	High	Slightly high	Slightly low
Combination of Use Automatic Bid (SpotAsPriceGo) and OOS	Low	High	Slightly high	High

Best practices for bidding modes

Preemptible instances support the following bidding modes. You can select a bidding mode based on your business requirements:

- If you require a strict control over your budgets but do not require high instance stability, you can select Set Maximum Price (SpotWithPriceLimit).
- If you require high instance stability but do not require a strict control over your costs, you can select Use Automatic Bid (SpotAsPriceGo).

If you require both instance stability and cost control, you can refer to the following steps to use the SpotAsPriceGo+OOS mode. You can use the Use Automatic Bid (SpotAsPriceGo) mode to increase instance stability and use the ACS-ECS-

AlarmWhenDiscountAndPriceExceedsThresholdInMultiZoneAndInstanceType public template of OOS to monitor the price of the preemptible instance. When the price exceeds the threshold that you specify, the system sends a notification.

1. Make preparations.

 i. (Optional) Create a preemptible instance whose bidding mode is Use Automatic Bid (SpotAsPriceGo).

For more information, see Create a preemptible instance. If you have created a preemptible instance, skip this step.

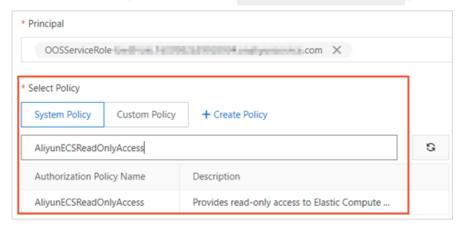
ii. Query the details of the preemptible instance.

For more information, see View instance information. In this example, two preemptible instances are created, as described in the following table. In subsequent steps, the O&M scripts for price monitoring are set based on the details of the instances.

Instance name	Region and zone	Instance type
Preemptible instance 1	Hangzhou Zone I	ecs.c5.xlarge
Preemptible instance 2	Hangzhou Zone K	ecs.r6.xlarge

iii. Create the oosserviceRole role for OOS by using Resource Access Management (RAM).

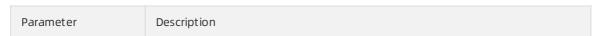
For more information, see Grant RAM permissions to OOS. In the Select Policy section of the Grant Permission page, select only the AliguneCSReadOnlyAccess policy.



iv. Add a DingTalk chat bot.

When you add a DingTalk chatbot, select **Custom Keywords** for **Security Settings** and enter Monitor in the Custom Keywords field. In subsequent steps, OOS sends notifications by using the webhook URL of the chatbot.

- 2. Log on to the OOS console.
- 3. In the left-side navigation pane, click **Scheduled O&M**.
- 4. On the **Scheduled O&M** page, click **Create**.
- 5. On the **Create Scheduled O&M** page, configure the parameters described in the following table and click **Execute Now**:



Parameter	Description		
Set Scheduled Task	 Select Executed Periodically and configure the following parameters: Execution Frequency: Use the default settings. The task is executed once per hour. Note The prices of preemptible instances change infrequently within a short period of time. We recommend that you set the value to 1 hour. Time Zone for Periodic Execution: Specify the time zone based on your region. In this example, the default setting is used. End Time for Period Execution: Specify the end time of the task based on your business requirements. In this example, the default setting is used. 		
Select Template	Enter ACS-ECS-AlarmWhenDiscountAndPriceExceedsThresholdInMultiZon eAndInstanceType in the search box. Then, select the template.		
Configure Parameters			
Advanced	Use the default settings. You can configure the parameters based on your business requirements.		

After the scheduled task is executed, you can view the state of the scheduled task on the

Scheduled O&M page.



- 6. Wait for the scheduled task to be executed or test the scheduled task.
 - After the scheduled task is executed, the system monitors the bid price of the preemptible
 instances in real time. If the price exceeds the specified threshold, a notification is sent by using
 the DingTalk chatbot.
 - In actual scenarios, the prices of preemptible instances do not fluctuate frequently and the
 execution result of the scheduled task is difficult to be verified. You can create a scheduled task
 and set the threshold to a value that triggers an alert. For example, if the price of a preemptible
 instance is USD 0.080 per hour, the threshold can be set to 0.040. Because the price of the
 preemptible instance remains higher than the threshold, an alert is triggered when the scheduled
 task is executed, and a notification is sent by using the DingTalk chatbot.

7.3. Simulate a preemptible instance interruption event

Preemptible instance interruption events are triggered events. When you develop an interruption event handler for a preemptible instance, you cannot debug the code. Therefore, Alibaba Cloud provides methods to simulate preemptible instance interruption events for easy debugging of O&M programs.

Context

Alibaba Cloud CloudMonitor is required to simulate preemptible instance interruption events. You can use the CloudMonitor console or call CloudMonitor API operations to simulate interruption events.

You can use the following methods to simulate a preemptible instance interruption event:

- Method 1: Use the CloudMonitor console to simulate an interruption event
- Method 2: Call CloudMonitor API operations to simulate an interruption event

Method 1: Use the CloudMonitor console to simulate an interruption event

This section describes how to use the CloudMonitor console to simulate an interruption event. In this example, an interruption event-triggered alert rule is configured by using CloudMonitor, and the event is delivered to a Message Service (MNS) queue.

- 1.
- 2.
- 3.
- 4.
- 5. In the Create / Modify Event Alert panel, set the following parameters and click OK:
 - Alert Rule Name: Specify a name for the alert rule. Example: preemptible instance interruption

event alert.

- Event Type: Select System Event.
- o Product Type: Select ECS.
- Event Type: Select Status Notification.
- o Event Level: Select WARN.
- Event Name: Select Instance:PreemptibleInstanceInterruption.
- Resource Range: Use the default setting.
- Alert Type: Select notification methods based on your business requirements. In this example, select MNS queue.

After the alert rule is created, you can view the rule on the **Alert Rules** tab.

- 6. Find the alert rule for preemptible instance interruption events and clicktest in the Actions column.
- 7. In the Create event test panel, modify the JSON file and then click OK.

Replace the resource information in the JSON file with the information about the preemptible instance for which you want to simulate an interruption event. The JSON file contains the following content. Take note of the following variables:

- Replace the Alibaba Cloud account ID variable with the ID of your Alibaba Cloud account.
- Replace the <resource-id> and i-abcdef variables with the ID of the preemptible instance.
- Replace the <region ID> variable with the region ID of the preemptible instance.

```
"product": "ECS",
    "resourceId": "acs:ecs:cn-shanghai: Alibaba Cloud account ID :instance/ <resource-i
d> ",
    "level": "WARN",
    "instanceName": "instanceName",
    "regionId": "<region ID> ",
    "groupId": "0",
    "name": "Instance:PreemptibleInstanceInterruption",
    "content": {
        "instanceId": "i-abcdef",
        "action": "delete"
        },
        "status": "Normal"
}
```

After debugging starts, alert notifications are sent based on the notification methods that you specified.

Method 2: Call CloudMonitor API operations to simulate an interruption event

This section describes how to call CloudMonitor API operations to simulate an interruption event. In this example, Alibaba Cloud CloudMonitor SDK for Java is used to configure an interruption event-triggered alert rule and deliver the event to an MNS queue.

1. Make preparations.

i. Obtain the AccessKey ID of your Alibaba Cloud account.

For more information, see Obtain an AccessKey pair.

ii. Install CloudMonitor SDK for Java in the development environment.

Add the following dependencies to the Maven project. For more information, see CloudMonitor SDK for Java.

```
<dependency>
  <groupId>com.aliyun</groupId>
  <artifactId>tea-openapi</artifactId>
    <version>0.0.13</version>

</dependency>
  <dependency>
    <groupId>com.aliyun</groupId>
    <artifactId>cms20190101</artifactId>
    <version>1.0.1</version>
</dependency></dependency></dependency></dependency></dependency></dependency>
```

2. Call the PutEventRule operation to create an interruption event-triggered alert rule for a preemptible instance.

Sample Java code:

```
import com.aliyun.cms20190101.models.*;
import com.aliyun.teaopenapi.models.*;
 * Call the PutEventRule operation to create or modify an event-triggered alert rule.
public class Sample {
      // Initialize request parameters.
       public static com.aliyun.cms20190101.Client createClient(String accessKeyId, String
accessKeySecret) throws Exception {
               // Specify the AccessKey ID and AccessKey secret of your Alibaba Cloud account.
               \texttt{Config config = new Config().setAccessKeyId(accessKeyId).setAccessKeySecret(accessKeyId).setAccessKeySecret(accessKeyId).setAccessKeyId(accessKeyId).setAccessKeyId(accessKeyId).setAccessKeyId(accessKeyId).setAccessKeyId(accessKeyId).setAccessKeyId(accessKeyId).setAccessKeyId(accessKeyId).setAccessKeyId(accessKeyId(accessKeyId)).setAccessKeyId(accessKeyId(accessKeyId)).setAccessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(acces)KeyId(accessKeyId(accessKeyId(accessKeyId(accessKeyId(
essKeySecret);
               // Specify the endpoint of the API operation.
               config.endpoint = "metrics.cn-hangzhou.aliyuncs.com";
               return new com.aliyun.cms20190101.Client(config);
       public static void main(String[] args ) throws Exception {
               java.util.List<String> args = java.util.Arrays.asList(args_);
                 * Set the <your-access-key-id> variable to the AccessKey ID of your Alibaba Cl
oud account.
                 * Set the <your-access-key-secret> variable to the AccessKey secret of your Al
ibaba Cloud account.
               com.aliyun.cms20190101.Client client = Sample.createClient("<your-access-key-id</pre>
>", "<your-access-key-secret>");
               {\tt PutEventRuleRequest.PutEventRuleRequestEventPattern\ eventPattern0 = new\ PutEven}
tRuleRequest.PutEventRuleRequestEventPattern()
                               // Specify the type of the event-triggered alert rule.
                               .setEventTypeList(java.util.Arrays.asList("*"))
                               // Specify the level of the event-triggered alert rule.
                                .setLevelList(java.util.Arrays.asList("*"))
                               // Specify the name of the event-triggered alert rule.
                               .setNameList(java.util.Arrays.asList("Instance:PreemptibleInstanceInter
ruption"))
                               // Specify the Alibaba Cloud service to which the event-triggered alert
rule is applied.
                               .setProduct("ECS");
               PutEventRuleRequest putEventRuleRequest = new PutEventRuleRequest()
                               // Specify a name for the alert rule.
                               .setRuleName("spot release event test")
                               .setEventPattern(java.util.Arrays.asList(eventPattern0))
                               \ensuremath{//} Specify the type of the event-triggered alert rule.
                                .setEventType("SYSTEM")
                               // Specify the state of the event-triggered alert rule.
                               .setState("ENABLED");
                // Obtain the response.
               client.putEventRule(putEventRuleRequest);
       }
}
```

3. Create an MNS queue by using Alibaba Cloud MNS.

For the sample code on how to create an MNS queue, see Step 2: Create a queue.

4. Call the PutEventRuleTargets operation to set notification methods for the alert rule and deliver messages to the MNS queue.

Sample Java code:

```
import com.aliyun.cms20190101.models.*;
import com.aliyun.teaopenapi.models.*;
^{\star} Call the PutEventRuleTargets operation to add or modify the MNS queue.
public class Sample {
   // Initialize request parameters.
    public static com.aliyun.cms20190101.Client createClient(String accessKeyId, String
accessKeySecret) throws Exception {
        # Specify the AccessKey ID and AccessKey secret of your Alibaba Cloud account.
        Config config = new Config().setAccessKeyId(accessKeyId).setAccessKeySecret(acc
essKeySecret);
        // Specify the endpoint of the API operation.
        config.endpoint = "metrics.cn-hangzhou.aliyuncs.com";
        return new com.aliyun.cms20190101.Client(config);
    public static void main(String[] args ) throws Exception {
         * Set the <your-access-key-id> variable to the AccessKey ID of your Alibaba Cl
oud account.
         * Set the <your-access-key-secret> variable to the AccessKey secret of your Al
ibaba Cloud account.
        */
        com.aliyun.cms20190101.Client client = Sample.createClient("<your-access-key-id
>", "<your-access-key-secret>");
        {\tt PutEventRuleTargetsRequest.PutEventRuleTargetsRequestMnsParameters\ mnsParameter}
\verb|s0| = \verb|new PutEventRuleTargetsRequest.PutEventRuleTargetsRequestMnsParameters()|
                // Specify the region to which the MNS queue belongs.
                .setRegion("cn-hangzhou")
                // Specify the ID of the MNS queue.
                .setId("1")
                \ensuremath{//} Specify the name of the MNS queue.
                .setQueue("mq-test");
        PutEventRuleTargetsRequest putEventRuleTargetsRequest = new PutEventRuleTargets
Request()
                // Specify the name of the alert rule.
                 .setRuleName("spot release event test")
                .setMnsParameters(java.util.Arrays.asList(
                        mnsParameters0
                ));
        // Obtain the response.
        client.putEventRuleTargets(putEventRuleTargetsRequest);
```

5. Call the SendDryRunSystemEvent operation to send the request to simulate an interruption event.

Sample Java code:

```
import com.aliyun.cms20190101.models.*;
import com.aliyun.teaopenapi.models.*;
 * Call the SendDryRunSystemEvent operation to debug a system event of an Alibaba Cloud
resource.
 */
public class Sample {
        // Initialize request parameters.
        public static com.aliyun.cms20190101.Client createClient(String accessKeyId, String
accessKeySecret) throws Exception {
                # Specify the AccessKey ID and AccessKey secret of your Alibaba Cloud account.
                Config config = new Config().setAccessKeyId(accessKeyId).setAccessKeySecret(acc
essKeySecret);
                // Specify the endpoint of the API operation.
                config.endpoint = "metrics.cn-hangzhou.aliyuncs.com";
                return new com.aliyun.cms20190101.Client(config);
        public static void main(String[] args ) throws Exception {
                  ^{\star} Set the <your-access-key-id> variable to the AccessKey ID of your Alibaba Cl
                  * Set the <your-access-key-secret> variable to the AccessKey secret of your Al
ibaba Cloud account.
                com.aliyun.cms20190101.Client client = Sample.createClient("<your-access-key-id
>", "<your-access-key-secret>");
                SendDryRunSystemEventRequest sendDryRunSystemEventRequest = new SendDryRunSyste
mEventRequest()
                                 // Specify the name of the Alibaba Cloud service.
                                 .setProduct("ecs")
                                 // Specify the name of the alert event.
                                 .setEventName("Instance:PreemptibleInstanceInterruption")
                                 // Specify the content of the event that you want to simulate.
                                  \verb|ghai:133160284996***: instance/i-abcdef\", \verb|"level\":\"WARN\", \"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"instanceName\":\"inst
eName\",\"regionId\":\"cn-beijing\",\"name\":\"Instance:PreemptibleInstanceInterruption
\",\"content\": {\"instanceId\":\"i-abcdef****\",\"action\":\"delete\"},\"status\":\"No
rmal\"}");
                // Obtain the response.
                client.sendDryRunSystemEvent(sendDryRunSystemEventRequest);
}
```

6. Receive messages of the simulated event in the MNS queue.

For the sample code on how to receive messages from a queue, see Step 4: Receive and delete the message.

7.4. Receive a preemptible instance interruption event

Alibaba Cloud Cloud Monitor provides the resource monitoring feature. You can use Cloud Monitor to monitor preemptible instance interruption events in real time. When a preemptible instance interruption event occurs, you can receive the interruption event by using the notification methods that you specified.

Context

Preemptible instances may be interrupted. 5 minutes before a preemptible instance is interrupted, an interruption event is triggered for the instance. If your business is sensitive to instance interruptions, you must receive and handle preemptible instance interruption events in a timely manner. This topic describes the best practices for receiving preemptible instance interruption events by using Message Service (MNS) or Function Compute.

- Method 1: Use MNS to receive interruption events
- Method 2: Use Function Compute to receive interruption events

In addition to the best practices provided in this topic, you can also call Elastic Compute Service (ECS) API operations to receive preemptible instance interruption events. For more information, see Query the interruption events of preemptible instances.

Method 1: Use MNS to receive interruption events

- 1. Create a queue by using MNS.
 - For more information, see Create a queue.
- 2. Create an event-triggered alert rule by using CloudMonitor.

For more information, see Create a system event-triggered alert rule. Configure the following parameters:

- Alert Rule Name: Specify a name for the alert rule. Example: preemptible instance interruption event alert.
- Event Type: Select System Event.
- Product Type: Select ECS.
- Event Type: Select **Status Notification**.
- Event Level: Select WARN.
- Event Name: Select Instance:PreemptibleInstanceInterruption.
- Resource Range: Use the default setting.
- Alert Type: Select notification methods based on your business requirements. For example, select MNS queue and configure the created queue.

After the alert rule is created, you can view the rule on the Alert Rules tab.

- 3. Find the alert rule for preemptible instance interruption events and clicktest in the Actions column.
- 4. In the Create event test panel, modify the JSON file and then click OK.

Replace the resource information in the JSON file with the information about the preemptible instance for which you want to simulate an interruption event. The JSON file contains the following content. Take note of the following variables:

- Replace the Alibaba Cloud account ID variable with the ID of your Alibaba Cloud account.
- Replace the <resource-id> and i-abcdef variables with the ID of the preemptible instance.
- Replace the <region ID> variable with the region ID of the preemptible instance.

```
"product": "ECS",
    "resourceId": "acs:ecs:cn-shanghai: Alibaba Cloud account ID :instance/ <resource-i
d> ",
    "level": "WARN",
    "instanceName": "instanceName",
    "regionId": "<region ID> ",
    "groupId": "0",
    "name": "Instance:PreemptibleInstanceInterruption",
    "content": {
        "instanceId": "i-abcdef",
        "action": "delete"
    },
    "status": "Normal"
}
```

After debugging starts, alert notifications are sent based on the notification methods that you specified.

5. Receive interruption events by using MNS.

For more information, see Receive a message.

You can also use MNS SDK for Java to associate interruption events with your business. For more information, see Release notes of the SDK for Java.

Method 2: Use Function Compute to receive interruption events

- 1. Create a Function Compute service.
- 2. Create a function.

i. ..

iii. In the Basic Settings section, set the parameters as required and click Create.

- (Optional) Name: Specify a name for the function. Example: testSpotInstance.
- Runtime Environments: Select Python 2.7.
- Function Trigger Mode: Select Event-triggered.
- Instance Category: Select Elastic Instance.
- Memory Capacity: Select 512 MB from the drop-down list.

After the function is created, you are navigated to the Code tab of the function details page.

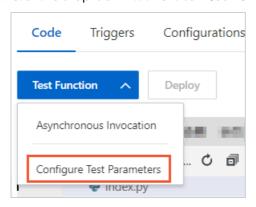
- 3. Configure the function code.
 - i. In the integrated development environment (IDE), find and click the default index.py file.

- ii. Replace the content of the index.py file with the following code, and then click Deploy.

 Take note of the following variables:
 - Replace the *Alibaba Cloud account ID* variable with the ID of your Alibaba Cloud account.
 - Replace the *<resource-id>* and *i-abcdef* variables with the ID of the preemptible instance.
 - Replace the <region ID> variable with the region ID of the preemptible instance.

```
# -*- coding: utf-8 -*-
import logging
import json, random, string, time
LOGGER = logging.getLogger()
clt = None
def handler(event, context):
    "product": "ECS",
   "resourceId": "acs:ecs:cn-shanghai: Alibaba Cloud account ID :instance/ <resour
   "level": "WARN",
    "instanceName": "instanceName",
   "regionId": "<region ID>",
   "groupId": "0",
    "name": "Instance:PreemptibleInstanceInterruption",
    "content": {
        "instanceId": "i-abcdef",
        "action": "delete"
    "status": "Normal"
  evt = json.loads(event)
  content = evt.get("content");
  regionId = evt.get("regionId");
  instanceId = content.get("instanceId");
  LOGGER.info( regionId + " " + instanceId + " termination ongoing");
```

iii. Click the drop-down list next to Test Function and select Configure Test Parameters.



- iv. In the Configure Test Parameters dialog box, set the following parameters and click OK:
 - Event Template: Use the default setting.
 - Event Name: Specify an event name. Example: preemptible instance interruption event.
 - Field: Replace the default content with the following code. The variables must be consistent with those in the index.py file that you configured.

```
"product": "ECS",
    "resourceId": "acs:ecs:cn-shanghai: Alibaba Cloud account ID :instance/ <reso
urce-id>",
    "level": "WARN",
    "instanceName": "instanceName",
    "regionId": "<region ID>",
    "groupId": "0",
    "name": "Instance:PreemptibleInstanceInterruption",
    "content": {
        "instanceId": "i-abcdef",
        "action": "delete"
    },
    "status": "Normal"
}
```

4. On the Code tab, click Test Function.

You can use the test function feature to view the sample output of the function code. The following figure shows the output returned in this example. The output contains the ID and region ID of the instance for which the interruption event is triggered.



5. If the function code runs normally, create an event-triggered alert rule by using CloudMonitor.

For more information, see Create a system event-triggered alert rule. Configure the following parameters:

- Alert Rule Name: Specify a name for the alert rule. Example: preemptible instance interruption event alert.
- Event Type: Select System Event.
- Product Type: Select ECS.
- Event Type: Select Status Notification.
- o Event Level: Select WARN.
- Event Name: Select Instance: PreemptibleInstanceInterruption.
- Resource Range: Use the default setting.
- Alert Type: Select notification methods based on your business requirements. For example, select **Function service** and configure the created Function Compute service.

After the alert rule is created, you can view the rule on the Alert Rules tab.

- 6. Find the alert rule for preemptible instance interruption events and clicktest in the Actions column.
- 7. In the Create event test panel, modify the JSON file and then click OK.

Replace the resource information in the JSON file with the information about the preemptible instance for which you want to simulate an interruption event. The JSON file contains the following content. Take note of the following variables:

- Replace the Alibaba Cloud account ID variable with the ID of your Alibaba Cloud account.
- Replace the <resource-id> and i-abcdef variables with the ID of the preemptible instance.
- Replace the <region ID> variable with the region ID of the preemptible instance.

```
"product": "ECS",
    "resourceId": "acs:ecs:cn-shanghai: Alibaba Cloud account ID :instance/ <resource-i
d> ",
    "level": "WARN",
    "instanceName": "instanceName",
    "regionId": "<region ID> ",
    "groupId": "0",
    "name": "Instance:PreemptibleInstanceInterruption",
    "content": {
        "instanceId": "i-abcdef",
        "action": "delete"
    },
    "status": "Normal"
}
```

After debugging starts, alert notifications are sent based on the notification methods that you specified.

8.Security 8.1. Best practices of the security group (part 1)

This article introduces how to configure the inbound rules of security groups.

Like a virtual firewall, a security group controls network access for one or more ECS instances. It is an important means of security isolation. When creating an ECS instance, you must select a security group. You can also add security group rules to control outbound and inbound access for all ECS instances in the same security group.

Before configuring the inbound rules for a security group, you should have learnt about the following information:

- Security group restrictions
- Default security group rules
- Set the inbound access of a security group
- Set the outbound access of a security group

General suggestions for security group practices

Before you work with security groups, read the following suggestions:

- The most important rule: A security group should be used as a whitelist.
- The "minimum authorization" principle should be observed when you configure the inbound or outbound rules for applications. For example, you can allow a specific port (such as port 80).
- It is not recommended to use one security group to manage all applications, because requirements must be different at different layers.
- For distributed applications, different security groups should be used for different application types. For example, you should use different security groups for the Web, Service, Database and Cache layers to apply different inbound/outbound rules and permissions.
- There is no need to set a separate security group for every instance, as this would unnecessarily add to management costs.
- VPC should be preferred.
- Do not assign Internet addresses to resources that require no Internet access.
- Keep the rules of each security group as concise as possible. A single instance can join up to five security groups, and a security group can contain up to 200 security group rules, so an instance may be subject to hundreds of security group rules at the same time. You can aggregate all the assigned security rules to determine whether inbound or outbound traffic is permitted or not. However, overly complicated rules for a single security group can increase management complexity. For this reason, it is recommended to keep the rules of each security group as concise as possible.
- The ECS console allows you to clone a security group and security group rules. If you want to modify an active security group and its rules, you should clone the security group and modify the cloned security group, avoiding any impacts on online applications.

? Note Adjusting inbound or outbound rules of active security groups can be risky. Therefore, do not update those rules at will unless you know what you are doing.

Set inbound access rules of security groups

The following are some suggestions about inbound rules of a security group.

Do not use the 0.0.0.0/0 inbound rule

It is a common mistake to permit all inbound access without any restrictions. Using 0.0.0.0/0 means that all ports are open to external access. This is extremely insecure. The correct practice is to deny external access to all the ports first. Whitelist items should be configured for security groups. For example, if you need to expose web services, you should only open common TCP ports such as 80, 8080 and 443 by default. All other ports should be disabled.

```
{ "IpProtocol": "tcp", "FromPort": "80", "ToPort": "80", "SourceCidrIp": "0.0.0.0/0", "Policy": "accept"},
{ "IpProtocol": "tcp", "FromPort": "8080", "ToPort": "8080", "SourceCidrIp": "0.0.0.0/0", "Policy": "accept"},
{ "IpProtocol": "tcp", "FromPort": "443", "ToPort": "443", "SourceCidrIp": "0.0.0.0/0", "Policy": "accept"},
```

Disable unneeded inbound rules

If your current inbound rules include 0.0.0.0/0, review the ports and services that must be exposed for your applications. If you do not want some ports to directly provide services for external applications, add denial rules for them. For example, if you have installed MySQL database services on the server, port 3306 should not be exposed to the Internet by default. You can add a denial rule, as shown below. Set the priority value to 100, which is the lowest priority.

```
{ "IpProtocol" : "tcp", "FromPort" : "3306", "ToPort" : "3306", "SourceCidrIp" : "0.0.0.0/0
", "Policy": "drop", Priority: 100} ,
```

This setting prevents any other ports from accessing port 3306. However, this can block normal service requests as well. For this reason, you can authorize resources of another security group for inbound access.

Authorize another security group for inbound access

Different security groups adopt inbound and outbound rules in accordance with the minimum authorization principle. Different application layers should use different security groups with corresponding inbound and outbound rules.

For example, different security groups are configured for distributed applications. However, directly authorizing IP addresses or CIDR network segments can be very difficult as different security groups cannot intercommunicate on the Internet. In this situation, you can authorize all resources of another security group to be directly accessible. For example, sg-web and sg-database security groups are created respectively for the Web and Database layers of your applications. In sg-database, you can add the following rule to authorize all resources in the sg-web security group to access port 3306.

```
{ "IpProtocol": "tcp", "FromPort": "3306", "ToPort": "3306", "SourceGroupId": "sg-web", "Policy": "accept", Priority: 2},
```

Authorize another CIDR for inbound access

In classic networks, controlling network segments is difficult and you are recommended to use security group IDs to authorize inbound rules.

In VPC networks, you can plan IP addresses on your own and use different vSwitches to set different IP domains. Therefore, in VPC networks, you can deny any access by default but authorize access for your own VPC, namely directly authorizing trusted CIDR network segments.

```
{ "IpProtocol" : "icmp", "FromPort" : "-1", "ToPort" : "-1", "SourceCidrIp" : "10.0.0.0/24" , Priority: 2} , 
{ "IpProtocol" : "tcp", "FromPort" : "0", "ToPort" : "65535", "SourceCidrIp" : "10.0.0.0/24 ", Priority: 2} , 
{ "IpProtocol" : "udp", "FromPort" : "0", "ToPort" : "65535", "SourceCidrIp" : "10.0.0.0/24 ", Priority: 2} ,
```

Steps and instructions for changing security group rules

Changing security group rules can interrupt network communication among instances. To prevent required network communication from being impacted, try to permit required instances with the method below and then execute security group policies to narrow down your changes.

Note After narrowing down the changes, check that service applications are running correctly before performing other required changes.

- Create a new security group, add instances that need mutual access to it, and then perform the changes.
- If the authorization type is **Security Group**, add the bound security group IDs of peer instances that require intercommunication into the authorization rules of the security group.
- If the authorization type is CIDR, add Intranet IP addresses of peer instances that require intercommunication into the authorization rules of the security group.

For detailed instructions, see Add a security group rule.

8.2. Best practices for ECS security groups (part 2)

This topic describes the best practices on how to add a security group rule, cancel a security group rule, add ECS instances to a security group, and remove ECS instances from a security group.

Network types

Alibaba Cloud provides two types of networks: classic network and virtual private cloud (VPC). These network types support different security group rules.

- For the classic network, you can configure inbound rules for internal network traffic, outbound rules for internal network traffic, inbound rules for Internet traffic, and outbound rules for Internet traffic.
- For VPCs, you can configure inbound and outbound rules for internal network traffic.

Different security groups may have different network types. Classic network-type ECS instances can be added only to classic network-type security groups. VPC-type ECS instances can be added only to security groups in the VPCs where the instances are located.

Basic knowledge of internal network communication between security groups

Before you begin, you must learn about the following items about internal network communication between security groups:

- The instances in different security groups of the same account are inaccessible to each other over the internal network. This applies to both the classic network and VPCs. Therefore, the classic network-type ECS instances are secure over the internal network.
- If you have two ECS instances in different security groups and the instances are not inaccessible to each other over the internal network as you expected, you must check the internal network rules of your security groups. If an internal network rule meets one of the following conditions, we recommend that you reconfigure the rule:
 - Traffic on all ports is allowed.
 - The authorization object is specified in the form of CIDR block when the SourceCidrlp parameter is set to 0.0.0.0/0 or 10.0.0.0/8. For classic network-type security groups, the preceding internal network rule allows external resources to access the security group.
- If you want to allow resources in different security groups to communicate with each other, you must configure security group rules to allow mutual access between the security groups. To allow access over the internal network, you must specify security groups as authorization objects, instead of CIDR blocks.

Attributes of security rules

Security rules describe access permissions by using the following attributes:

- Policy: the access control policy. Valid value: accept and drop.
- Priority: the priority. Security group rules are assigned priorities in descending order based on the time when the rules are created. The priority of a security rule ranges from 1 to 100. The default value is 1, which indicates the highest priority. A greater value indicates a lower priority.
- NicType: the network type. If you specify a security group as an authorization object by setting SourceGroupId instead of SourceCidrlp, you must set NicType to *intranet*.
- Description:
 - IpProtocol: the IP protocol. Valid values: *tcp, udp, icmp, gre,* and *all.* A value of *all* indicates all protocols.
 - PortRange: the port range that corresponds to the IP protocol.
 - If the value of IpProtocol is *tcp* or *udp*, the valid port numbers are 1 to 65535. You can specify a port range in the format of <start port number>/<end port number>. For example, 1/200 indicates that the port range is 1-200. If the value is 200/1, an error is reported when the API operation is called.
 - If the value of IpProtocol is icmp, gre, or all, the port range is -1/-1, which indicates all ports.
 - To control access from a security group, you must set SourceGroupId to the ID of the security group. In this case, you can choose to or not to configure SourceGroupOwnerAccount based on whether the rule is to control access between security groups across accounts.
 SourceGroupOwnerAccount specifies the account to which the source security group belongs.
 - To control access from specific IP addresses, you must set the SourceCidrlp parameter to the IP addresses in CIDR notation.

Create a rule to allow inbound requests

When you create a security group by using the console or by calling an API operation, a default rule is added to the security group to *deny all* inbound requests. If the default rule is not suitable, configure inbound rules on your own.

For example, if you want to enable port 80 which is connected to the Internet and then use the port to provide HTTP services externally, do not impose limits on CIDR blocks but set SourceCidrlp to 0.0.0.0/0 to allow inbound requests. Configure the following parameters to create an inbound rule to allow all inbound requests. The parameters and values outside the brackets are used in the ECS console, whereas, the parameters and values inside the brackets are used in ECS API operations.

- NIC Type (NicType): Public (internet). If the security group is of the VPC type, NIC Type is automatically set to Internal (intranet). The security group can be accessed by other security groups over elastic IP addresses (EIPs).
- Action (Policy): Allow (accept).
- Rule Direction: Inbound (inbound).
- Protocol Type (IpProtocol): Custom TCP (tcp).
- Port Range (PortRange): 80/80.
- Authorized Object (SourceCidrlp): 0.0.0.0/0.
- Priority (Priority): 1.

? Note The preceding values are applicable to security group rules for Internet traffic. For security group rules for internal network traffic, we recommend that you do not use CIDR blocks as authorization objects. For more information, see the **Create an internal network rule for a security group rule to allow internal network communication** section in this topic.

Create a rule to deny inbound requests

To deny inbound requests, you need only to configure a Deny policy that has a low priority. Then, you can configure another rule that has a higher priority to take precedence over this rule if necessary. For example, you can configure the following parameters to create a rule that denies access over port 6379.

- NIC Type (NicType): Internal (intranet).
- Action (Policy): Forbid (drop).
- Rule Direction: Inbound (inbound).
- Protocol Type (IpProtocol): Custom TCP (tcp).
- Port Range (PortRange): 6379/6379.
- Authorized Object (SourceCidrlp): 0.0.0.0/0.
- Priority (Priority): 100.

Create an internal network rule for a security group rule to allow internal network communication

By default, no inbound rules for internal network traffic are enabled for classic network-type ECS instances. Exercise caution for internal network authorization.

Note For security reasons, we recommend that you do not enable authorization based on CIDR blocks.

For ECS instances, internal IP addresses frequently change and the CIDR blocks to which these IP addresses belong vary. For this reason, we recommend that you allow access to classic network-type ECS instances over the internal network by creating security group rules instead of specifying CIDR blocks.

For example, if you build a Redis cluster in the sg-redis security group and want to permit only specific computers such as those in the sg-web security group to access the servers in the Redis clusters, you need only to create a security group rule that specifies the ID of the source security group, instead of CIDR blocks, as the authorization object.

- NIC Type (NicType): Internal (intranet).
- Action (Policy): Allow (accept).
- Rule Direction: Inbound (inbound).
- Protocol Type (IpProtocol): Custom TCP (tcp).
- Port Range (PortRange): 6379/6379.
- Authorized Object (SourceGroupId): sg-web.
- Priority (Priority): 1.

For VPC-type instances, if you planned an IP address range by using multiple vSwitches, you can use the CIDR blocks in inbound security group rules. However, if your VPC CIDR blocks are not clear, we recommend that you prioritize security groups over CIDR blocks as authorization objects when you configure inbound rules.

Add ECS instances that require mutual access to the same security group

You can add an ECS instance to up to five security groups. ECS instances in the same security group can communicate with each other over the internal network. If you have multiple security groups but do not want to configure multiple security group rules, you can create a new security group and add the instances that have the requirement for internal network communication to the new security group.

However, we recommend that you do not add all ECS instances to the same security group because this makes the configuration of security group rules messy. For a large or medium-sized application, each server group is assigned a different role. We recommend that you configure appropriate inbound and out bound rules for each server group.

In the ECS console, you can add an ECS instance to a security group by following the instructions in Add an ECS instance to a security group.

If you are familiar with Alibaba Cloud OpenAPI, you can use OpenAPI to add multiple instances at the same time. For more information, see Query an ECS instance. The following Python snippet shows how to add multiple ECS instances to a security group at the same time.

```
def join sg(sg id, instance id):
   request = JoinSecurityGroupRequest()
   request.set InstanceId(instance id)
   request.set SecurityGroupId(sg id)
   response = send request (request)
   return response
# send open api request
def send request (request):
    request.set accept format('json')
   try:
       response str = clt.do action(request)
       logging.info(response str)
       response detail = json.loads(response str)
       return response detail
    except Exception as e:
        logging.error(e)
```

Remove an ECS instance from a security group

If an ECS instance is added to an inappropriate security group, your services may be exposed or blocked. In this case, you can remove the ECS instance from the security group. Before you remove an ECS instance from a security group, make sure that the ECS instance also belongs to another security group.

Note After an ECS instance is removed from a security group, this instance and other instances in the security group cannot access each other. We recommend that you perform sufficient tests before you remove the instance.

The following Python snippet shows how to remove an ECS instance from a security group.

```
def leave sq(sq id, instance id):
   request = LeaveSecurityGroupRequest()
   request.set InstanceId(instance id)
   request.set SecurityGroupId(sg id)
   response = send request(request)
   return response
# send open api request
def send request (request):
   request.set accept format('json')
    try:
       response str = clt.do action(request)
       logging.info(response str)
       response detail = json.loads(response str)
       return response detail
    except Exception as e:
        logging.error(e)
```

Define appropriate names and tags for security groups

Appropriate names and descriptions can help you identify the meanings of complicated rule combinations. You can change names and descriptions of security groups.

You can configure tags for security groups. Then, you can group and manage security groups by using tags. You can configure tags by using the ECS console or by calling API operations. For more information, see Create or bind a tag.

Delete security groups that are no longer needed

Security group rules of security groups serve as whitelists and blacklists. We recommend that you delete security groups that are no longer needed to avoid issues that may occur if you add an ECS instance to an inappropriate security group.

8.3. Best practices for ECS security groups (part 3)

In practice, you may have a tendency to place all Elastic Compute Service (ECS) instances within a single security group to reduce initial configuration workloads. However, in the long run, this practice may make business system interactions more complex and less controllable. When you add rules to or remove rules from the security group, you cannot identify the impact scope of the rules.

Context

Plan and define security groups properly to make your systems easier to adjust. Sort out the services provided by applications and classify the applications into different levels. We recommend that you plan different security groups for different business and add different rules to the security groups.

Create different security groups

• Create different security groups for ECS instances that provide Internet-facing services and those that provide internal network-facing services.

The applications on an ECS instance may be accessible over the Internet regardless of whether the instance provides Internet-facing services, including proactively exposing specific ports (such as ports 80 and 443) for external access and passively providing port forwarding rules (such as NAT port forwarding rules and forwarding rules based on the system-assigned public IP address or elastic IP address of the instance).

In the two preceding scenarios, use the strictest security group rules for the instance. We recommend that you first add a rule to deny access from all protocols on all ports and then add rules to allow access only to ports required by external services, such as ports 80 and 443. The security group contains only Internet-facing ECS instances. It is easy to adjust the rules in the security group.

The security group that contains Internet-facing ECS instances must have clear and simple responsibilities to ensure that the instances provide no services other than primary services. For example, for MySQL and Redis applications, we recommend that you deploy them on ECS instances that do not provide Internet access, and then configure rules to allow access from authorized security groups to the instances.

Assume that an Internet-facing ECS instance named ECS_WEB_1 and some applications belong to the SG_CURRENT security group. You can perform the following steps to change the security group of the ECS_WEB_1 instance:

- i. Identify the protocols and ports used to provide Internet-facing services, such as ports 80 and 443.
- ii. Create a security group named SG_WEB and add a rule with the following attributes to the SG_WEB security group. For more information, see Create a security group.

■ Action: Allow

■ Protocol type: All

Port range: 80/80 and 443/443Authorization object: 0.0.0.0/0

iii. Add a rule with the following attributes to the SG_CURRENT security group to allow access from the SG_WEB security group. For more information, see Add a security group rule.

Action: Allow

- Protocol type: AllPort range: -1/-1
- Authorization object: SG_WEB
- Priority: Select a value in the range of 1 to 100 based on your requirements.
- iv. Add the ECS_WEB_1 instance to the SG_WEB security group.
 - a.
 - b.
 - c. Find the SG_WEB security group and click Manage Instances in the Actions column.
 - d. Click Add Instance.
 - e. In the Add Instance dialog box, select ECS_WEB_1 from the drop-down list and click **OK**. Confirm that data can be transferred to or from the ECS_WEB_1 instance and that the network of the instance works normally.
- v. Remove the ECS_WEB_1 instance from the SG_CURRENT security group.
 - a.
 - b.
 - c. Find the SG CURRENT security group and click Manage Instances in the Actions column.
 - d. Select the ECS_WEB_1 instance and click Remove from Security Group.
 - e. In the Remove ECS Instance from Security Group message, click $\mathbf{O}\mathbf{K}.$
 - f. Check the network connectivity of the instance to ensure that data can be transferred to or from the instance and that the network of the instance works normally.

If the ECS_WEB_1 instance fails the network connectivity check, add the instance back to the SG_CURRENT security group. Then, check whether ports in the SG_WEB security group are exposed as expected. If yes, proceed with security group changes. If no, modify the port settings.

- vi. Proceed with security group changes.
- Use different security groups for different applications

In most cases, different operating systems in a production environment do not belong to the same application group for load balancing. To provide different services, operating systems must have different ports exposed and blocked. We recommend that you assign different operating systems to different security groups.

For example, for a Linux operating system, you may need to expose TCP port 22 to allow SSH connections. For a Windows operating system, you may need to expose TCP port 3389 to allow Remote Desktop Protocol (RDP) connections.

If instances that use the same image but provide different services do not need to communicate with each other over the internal network, we recommend that you also assign these instances to different security groups. This helps decouple image types from security groups, simplify subsequent changes to security group rules, and ensure that instances have simple responsibilities.

When you plan and add new applications, you must properly plan security groups in addition to using vSwitches to define subnets. Use CIDR blocks and security groups to distinguish yourself as a service provider or a consumer.

For specific change procedures, see the preceding section.

• Use different security groups in production and testing environments

To better isolate systems, you may build multiple testing environments and a single production environment in actual development. You may need to configure different security group rules for different environments to properly isolate networks. This way, you can prevent changes made for testing purposes from being uploaded to the production environment and affecting the stability of the production environment.

You can use security groups to confine access domains of applications and avoid communication between production and testing environments. You can also assign different security groups to different testing environments to block traffic between the environments and improve development efficiency.

Only assign public IP addresses to ECS instances that require Internet access

For the classic network or virtual private clouds (VPCs), proper allocation of public IP addresses facilitates Internet access management and reduces the attack surface of systems. In VPCs, we recommend that you allocate ECS instances that require Internet access to several specified vSwitches. This makes it easy to audit and distinguish the instances and helps avoid unexpected exposure of ECS instances for Internet access.

Most distributed applications have different layers and groups. For ECS instances that do not provide Internet access, we recommend that you do not assign public IP addresses. If multiple instances that provide Internet access, we recommend that you configure Server Load Balancer (SLB) to distribute Internet traffic to improve system availability and avoid single points of failure. For more information, visit the Server Load Balancer product page.

For ECS instances that do not need to access the Internet, we recommend that you do not assign public IP addresses. In VPCs, if your ECS instances do not have public IP addresses but do require Internet access, we recommend that you use NAT gateways to provide Internet proxy services for the instances. You need only to configure SNAT entries to provide Internet access for specific CIDR blocks or subnets. This way, you can avoid exposing services to the Internet after public IP addresses are assigned, when only outbound Internet access is required. For more information, see Create and manage SNAT entries.

Principle of least privilege

Security groups are expected to work like whitelists. Therefore, open and expose as few ports as you need and allocate as few public IP addresses as you need. You can associate elastic IP addresses with online ECS instances to allow easy access to the instances for task log queries and troubleshooting. However, this operation exposes the instances to the Internet. We recommend that you use jump servers instead to manage ECS instances to improve security.

Use a jump server

A jump server has much higher permissions. You must use tools to thoroughly record and audit relevant operations. We recommend that you assign the jump server to a dedicated vSwitch in a VPC and then associate the corresponding elastic IP address or NAT port forwarding table with the jump server.

Create a dedicated security group named SG_BRIDGE and open required ports. For example, open TCP port 22 for Linux operating systems and RDP port 3389 for Windows operating systems. To restrict inbound access to the security group, you can allow only specified public IP addresses of your enterprise to access the security group. This way, you can reduce the probability that resources are scanned or accessed.

Add the ECS instance that functions as a jump server to the security group. To allow the jump server to access ECS instances in another security group, you can configure a rule that allows access to that security group for the jump server. For example, you can add a rule to a security group named SG_CURRENT to allow access from the SG_BRIDGE security group to specific protocols and ports.

When you use the jump server for SSH communication, we recommend that you use SSH key pairs to log on to ECS instances. For more information, see Overview.

In summary, proper planning of security groups makes it easy for you to scale applications and improve the security of your system.

8.4. Best practices for ECS data security

This topic describes the O&M best practices for improving the data security of ECS instances.

Intended users

This topic applies to individuals and small and medium-sized enterprises that are new to Alibaba Cloud.

Back up data regularly

Data backup is the foundation of disaster recovery, and can reduce the risk of data loss due to system failures, operational errors, and security problems. ECS provides the snapshot feature to meet the data backup requirements of most users. You can select a method of creating snapshots based on your business requirements. For more information, see Create a normal snapshot and Execute or disable an automatic snapshot policy.

We recommend that you create automatic snapshots on a daily basis and store the snapshots for at least seven days. This significantly improves disaster tolerance and minimizes potential data losses.

Design security domains

You can build private networks by using virtual private clouds (VPCs) to separately host servers of different security levels in your enterprise. This way, servers do not interfere each other over an interconnected network.

We recommend that you create a VPC. Then, specify a private IP address range in CIDR notation, and configure route tables and gateways of the VPC. You can store important data in the created VPC, which is logically isolated from the Internet. Then, you can use an elastic IP address (EIP) or a jumper server to manage data for daily O&M purposes. For more information, see Create a VPC.

Configure security group rules

Security groups are an essential tool to isolate networks. You can use security groups to configure network access control for one or more ECS instances. You can control the access to and from an ECS instance over the network by configuring security group rules. Specifically, you can restrict the inbound and outbound access on a port, and allow or deny access to and from specific IP addresses. This reduces the attack surface and makes ECS instances more secure.

Consider the following example. By default, port 22 is used as the remote access port on Linux ECS instances. Security risks arise if this port is open to external devices. You can authorize only specific local IP addresses to access the instances over port 22 when you configure security group rules. If you have higher security requirements, you can use third-party virtual private network (VPN) products to encrypt logon data.

Configure strong logon passwords

Weak passwords are one of the most common vulnerabilities and can lead to data leaks. To prevent security risks that arise from weak passwords, we recommend that you use a complex password for server access. The password must be at least eight characters in length, and contain multiple character types such as uppercase letters, lowercase letters, digits, and special characters. We also recommend that you change the password on a regular basis.

Improve the security of server ports

When ECS instances provide Internet access services, the service ports of the ECS instances are exposed. If you enable more service ports, the server is exposed to higher security risks. We recommend that you open only the required service ports, change the default service ports to high-numbered ports such as those numbers higher than 30000, and implement access control on these service ports.

For example, use internal connections for relational database service (RDS) to avoid exposing service ports to the Internet. If an external connection is necessary for RDS, change the default connection port from port 3306 to a port numbered 30001 or higher. Then, authorize client IP addresses based on your business requirements.

Install a WAF

Application vulnerabilities are security defects that can be exploited by hackers to illegally access data from web applications, cache, databases, and storage products. Common application vulnerabilities include SQL injection, XSS cross-site scripting, Webshell upload, backdoor protection, command injection, illegal HTTP requests, common vulnerabilities of web servers, unauthorized access to core files, and passthrough. Application vulnerabilities are different from system vulnerabilities, and are difficult to fix. You must consider security baseline issues before you design applications. To ensure website security and availability, we recommend that you use Web Application Firewall (WAF) to prevent various application attacks. For more information about how to deploy and use WAF, see Quick start.

8.5. Make ECS instances more secure

An ECS instance serves as a virtual machine. Typically, on-premises virtual machines are protected against attacks and intrusions. ECS instances also need security protection. You must implement effective security measures in conjunction with the inherent protection of Alibaba Cloud.

Prerequisites

Context

Lack of security protection for ECS instances may cause adverse effects. For example:

- DDoS attacks interrupt your business.
- Trojans tamper with or attack your web pages.
- Data leak caused by injection affects the normal operation of ECS.

You can use the following methods to improve the security of your ECS instances:

- Configure security groups
- Enable Anti-DDoS Basic
- Access Security Center
- Access Web Application Firewall

Configure security groups

A security group is a virtual firewall that provides stateful packet inspection (SPI). Security groups can serve the following purposes:

- Control access to one or more ECS instances. Security group rules can allow or deny inbound or outbound Internet and internal network traffic for ECS instances associated with security groups.
- If security groups are not planned properly or do not contain strict enough rules, they will be at a much greater risk of attack.

To add a rule to the security group associated with the ECS instance, perform the following operations:

- 1.
- 2.
- 3.
- 4. Find the security group to which to add a rule. Click Add Rules in the Actions column.
- 5. Click Add Security Group Rule.
- 6. In the dialog box that appears, configure the parameters.

For example, to allow a specific IP address to access the ECS instance, configure an inbound rule for Internet traffic. Assume that the instance runs the Linux operating system and you want to allow only the specific IP address to access the instance over port 22.

- i. Add the following inbound rule for Internet traffic:
 - Set Action to Allow.
 - Set Protocol Type to Custom TCP.
 - Set Port Range to 22/22.
 - Set Authorization Type to IPv4 CIDR Block.
 - Set **Authorization Object** to the CIDR block that is allowed to access to the ECS instance, in the format of x.x.x.x/xx (IP address/subnet mask). In this example, the CIDR block is 10.x.x.x /32. The priority is 1.

- ii. Add another inbound rule for Internet traffic:
 - Set Action to Forbid.
 - Set Protocol Type to Custom TCP.
 - Set Port Range to 22/22.
 - Set Authorization Type to IPv4 CIDR Block.
 - Set Authorization Object to 0.0.0.0/0. The priority is 2.

7. Click OK.

After the two rules are created, they can implement the following effects:

- The allow rule with a priority of 1 takes precedence for traffic from 10.x.x.x to port 22.
- The deny rule with a priority of 2 takes precedence for traffic from other IP addresses to port 22.

Enable Anti-DDoS Basic

Distributed denial of service (DDoS) attacks use client and server technologies to combine multiple computers into an attack platform and attack one or more targets simultaneously so that the impact of the denial-of-service (DoS) attack is multiplied.

Alibaba Cloud Security can defend against Layer 3 to Layer 7 DDoS attacks, including SYN Flood, UDP Flood, ACK Flood, ICMP Flood, DNS Flood, and HTTP Flood attacks. Anti-DDoS Basic provides up to 5 Gbit/s default DDoS protection free of charge. By default, Anti-DDoS Basic is enabled on ECS instances. Anti-DDoS eliminates the need to purchase expensive traffic scrubbing devices and allows you to maintain the access speed during DDoS attacks. Anti-DDoS help ensure your bandwidth regardless of the usage of other users. This ensures the availability and stability of your business. After an ECS instance is created, you can set the scrubbing thresholds. For more information, see Configure a traffic scrubbing threshold.

Alibaba Cloud has also launched the Security Credibility program, which provides improved DDoS protection based on a security credit score. Users that meet the criteria can obtain free protection against DDoS attacks up to 100 Gbit/s. In the Anti-DDoS Basic console, you can check your current security credibility score and details, as well as scoring criteria. For more information, see Security Credibility.

Access Security Center

Security Center is a unified security management system that recognizes, analyzes, and warns of security threats in real time. With security capabilities such as ransomware protection, anti-virus protection, web tamper protection, and compliance assessments, users can automate security operations, responses, and threat tracing to secure cloud and local servers and meet regulatory compliance requirements.

The Security Center agent is a security plug-in installed on your local servers. You must install this agent on your servers before you can enable Security Center features. For more information about how to install the Security Center agent, see Install the Security Center agent.

Note When you purchase an ECS instance, you can select the **Security Enhancement** check box to automatically install the agent and activate Security Center Basic edition.

The Basic edition of Security Center is available by default. The Basic edition only scans for the following risks: unusual logons to servers, vulnerabilities, and configuration risks in cloud services. To use advanced features such as vulnerability fixing and virus detection and removal, you must log on to the Security Center console.

Access Web Application Firewall

Web Application Firewall (WAF) is implemented based on the big data capabilities of Apsara Stack Security. This module protects web applications against common attacks reported by OWASP, such as SQL injections, XSS, vulnerability exploits in web server plugins, trojan attacks, and unauthorized access. WAF blocks malicious visits to avoid data from being compromised and ensure the security and availability of your websites.

WAF has the following benefits:

- WAF can handle various web application attacks to ensure web security and availability of a website
 without installing software or hardware or modifying website configuration and code. In addition to
 powerful web protection capabilities, WAF can customize protection for specific websites. WAF is
 used to protect web applications in fields such as finance, e-commerce, O2O, Internet Plus, gaming,
 governments, and insurance.
- Without WAF, you may be vulnerable to web intrusions such as data leaks, HTTP floods, and trojans.

For more information about how to access WAF, see Deploy WAF.

Alibaba Cloud provides multiple security services to safeguard ECS instances. You can choose appropriate methods to enhance systems and data protection, prevent intrusion into ECS instances, and ensure stability and reliability.

8.6. Configure interconnection of instances in the classic network

A security group is similar to a firewall for an instance. To ensure security, you must follow the least privilege principle when you configure security group rules for instances. This topic describes four recommended methods to configure interconnection of instances in the classic network.

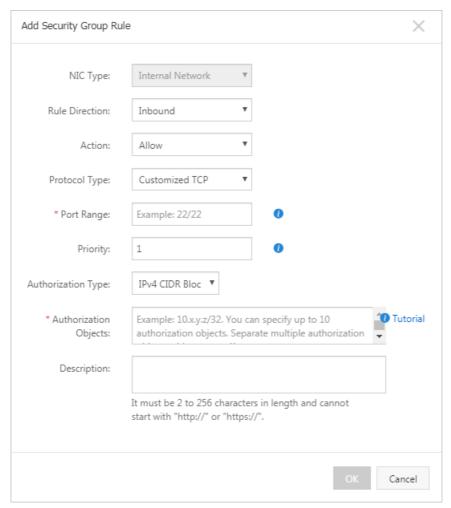
Method 1: Authorize access to a single IP address

- Scenario: This method is applicable to interconnection between a small number of instances over the internal network.
- Advantages: This method involves simple and clear security group rules to authorize access to a single IP address.
- Disadvantages: When you attempt to interconnect a large number of instances over the internal network, you are limited of the 200 security group rules and be burdened by a high maintenance workload.

To authorize access to a single IP address, perform the following steps:

- 1. Find an instance to be interconnected and click the instance ID.
- 2. On the Instance Details page, click the Security Groups tab.
- 3. Find the security group to be configured and click Add Rules in the Actions column.
- 4. Click the **Inbound** tab.
- 5. Click Add Security Group Rule.
- 6. Configure the following parameters for the security group rule:
 - o Action: Allow.
 - Priority: Configure this parameter based on your actual needs. Default value: 1.

- **Protocol Type**: Select a protocol type from the drop-down list based on your actual needs.
- o Port Range: Configure a port range based on your actual needs.
- Authorization Object: Enter the private IP address of the instance to be interconnected. The IP address must be in the *a.b.c.d/32* format. The subnet mask must be /32.



7. Click OK.

Method 2: Add instances to the same security group

- Scenario: If your application architecture is relatively simple, you can add all instances to the same basic security group. By default, instances in the same basic security group can access each other over the internal network.
- Advantages: Security group rules are simple and clear.
- Disadvantage: This method is applicable only to a simple network architecture. When the network architecture is adjusted, the authorization method must be modified accordingly.

For more information about the procedure, see Add an ECS instance to a security group.

Method 3: Add instances to a security group created solely for interconnection

• Scenario: You can add the instances to a dedicated basic security group for interconnection. This method is recommended for a network architecture that has multiple layers of applications.

- Advantages: This method is easy to implement and allows you to quickly establish interconnection between instances. The method is applicable to complicated network architectures.
- Disadvantages: The instances are added to multiple security groups and the security group rules may be complex.

To add instances to a security group created solely for interconnection, perform the following steps:

- 1. Create a basic security group that has a name configured such as security group for interconnection. No rules are required for the new security group.
- 2. Add the instances to be interconnected to the new security group. Instances in the same basic security group can communicate with each other.

Method 4: Authorize security groups

- Scenario: You can add the instances to a dedicated security group for interconnection. This method is recommended for a network architecture that has multiple layers of applications.
- Advantages: This method is easy to implement and allows you to quickly establish interconnection between instances. The method is applicable to complicated network architectures.
- Disadvantages: The instances are added to multiple security groups and the security group rules may be complex.

To authorize security groups, perform the following steps:

- 1. Find an instance to be interconnected and click the instance ID.
- 2. On the Instance Details page, click the Security Groups tab.
- 3. Find the security group to be configured and click **Add Rules** in the Actions column.
- 4. Click the **Inbound** tab.
- 5. Click Add Security Group Rule.
- 6. Configure the following parameters for the security group rule:
 - o Action: Allow.
 - Priority: Configure this parameter based on your actual needs. Default value: 1.
 - Protocol Type: Select a protocol type from the drop-down list based on your actual needs.
 - Port Range: Configure a port range based on your actual needs.
 - Authorization Object:
 - Instance in the same account as the current instance: Enter the security group ID to which the instance to be interconnected belongs.
 - Instance in a different account from the current instance: Enter the ID of the account and the ID of the security group to which the instance to be interconnected belongs. The IDs must be in the Account ID/security group ID> format.
- 7. Click OK.

Suggestions

If you determine that an inappropriate level of access has been granted by using the applied security group, we recommend that you downgrade the level of access by using the following procedure.



In the preceding figure, **Delete 0.0.0.0** means deleting the original security group rule that allows inbound traffic from 0.0.0.0/0.

If one or more security group rules are improperly configured, the interconnection between your instances may be affected. We recommend that you back up security group rules before you change them so that you can easily recover the rules.

A security group maps the role of an instance in the overall application architecture. We recommend that you plan the firewall rules based on the application architecture. For example, in a common three-tier web application architecture, you can plan three security groups and add instances deployed with applications or databases to the corresponding security groups.

- Web layer security group: allows port 80.
- Application layer security group: allows port 8080.
- Database layer security group: allows port 3306.

8.7. Modify the default port used by an instance to accept connections

This topic describes how to modify the default port used by an Elastic Compute Service (ECS) instance to accept connections.

Modify the default port used by a Windows instance to accept connections

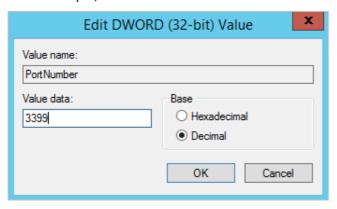
This section describes how to modify the default port used by a Windows instance to accept connections. In this example, Windows Server 2012 is used.

- Connect and log on to the Windows instance.
 For more information, see Connect to a Windows instance by using password authentication.
- 2. Modify the value of the **Port Number** registry subkey.

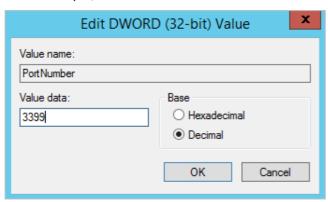
74

- i. Press Win+Rto open the Run command window.
- ii. Enter regedit.exe and press the Enter key to open the registry editor.
- iii. In the left-side navigation pane, choose HKEY_LOCAL_MACHINE > System > CurrentControlSet > Control > Terminal Server > Wds > rdpwd > Tds > tcp.
- iv. Find and right-click **Port Number** in the right list, and select **Modify...**.

v. In the Edit DWORD (32-bit) Value dialog box, enter a new port number in the **Value data** field. In this example, enter 3399. Select **Decimal** in the **Base** section and click **OK**.



- vi. In the left-side navigation pane, choose HKEY_LOCAL_MACHINE > System > CurrentControlSet > Control > Terminal Server > WinStations > RDP-Tcp.
- vii. Find and right-click **Port Number** in the right list, and select **Modify...**.
- viii. In the Edit DWORD (32-bit) Value dialog box, enter a new port number in the **Value data** field. In this example, enter 3399. Select **Decimal** in the **Base** section and click **OK**.



- 3. Restart the instance in the ECS console.
 - For more information, see Restart an instance.
- 4. Add security group rules to the security group of the instance to allow connections to the new port.
 - For more information, see Add security group rules.
- 5. Connect to the instance. In the Remote Desktop Connection dialog box, enter <IP address of the instance>:<New port number> in the Computer field and click Connect to connect to the instance.



? Note Mac Remote Desktop Connection can only be used to connect to the instance over the default port 3389.

Modify the default port used by a Linux instance to accept connections

This section describes how to modify the default port used by a Linux instance to accept connections. In the example, CentOS 6.8 and CentOS 7.7 are used.

- Connect and log on to the Linux instance.
 For more information, see Connect to a Linux instance by using a password.
- 2. Run the following command to back up the sshd configuration file:

```
cp /etc/ssh/sshd_config /etc/ssh/sshd_config_bak
```

- 3. Modify the port number of sshd.
 - i. Run the following command to edit the *sshd_config* configuration file:

```
vim /etc/ssh/sshd_config
```

ii. Press the I key to enter the edit mode.

iii. Add a new port to accept connections.

In this example, add port 1022. Enter Port 1022 under Port 22.

```
# If you want to change the port on a SELinux system, you have to tell
# SELinux about this change.
# semanage port -a -t ssh_port_t -p tcp #PORTNUMBER
#
#Port 22
#AddressFamily any
#ListenAddress 0.0.0.0
#ListenAddress ::

# If you want to change the port on a SELinux system, you have to tell
# SELinux about this change.
# semanage port -a -t ssh_port_t -p tcp #PORTNUMBER
#
#Port 22
Port 1022
#AddressFamily any
#ListenAddress 0.0.0.0
#ListenAddress ::
```

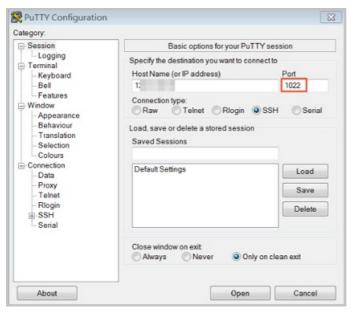
- iv. Press the Esc key, enter: wq, and then press the Enter key to save and close the file.
- 4. Run one of the following commands to restart sshd. After sshd is restarted, you can log on to the Linux instance by using SSH port 1022.
 - o If the Linux instance runs CentOS 7 or later, or Alibaba Cloud Linux 2, run the following command:

```
systemctl restart sshd
```

o If the Linux instance runs CentOS 6, run the following command:

```
/etc/init.d/sshd restart
```

- 5. Add security group rules to the security group of the instance to allow traffic over TCP port 1022. For more information, see Add a security group rule.
- 6. Use an SSH client to connect to the instance to check whether traffic over the new port is allowed. Enter the new port number in the **Port** field. In this example, enter 1022.





After the port number is modified, you cannot use the default port 22 to connect to the instance.

8.8. Use logs in Windows instances

Logs monitor events occurring in the system and record hardware, software, and system issues. When an instance is attacked or an issue occurs on an application, you can locate the critical problems based on logs. This improves work efficiency and instance security. This topic uses Windows Server 2012 R2 as an example to describe how to use and analyze four types of logs.

Prerequisites

Context

Windows logs can be further divided into system logs, application logs, security logs, and application and service logs.

View logs in Windows Event Viewer

Perform the following steps to open Event Viewer:

- 1. Click Start and open the Run dialog box.
- 2. Run the eventywr command in the dialog box to open Event Viewer.
- 3. View the following four types of logs in **Event Viewer**:

Note For IDs of all error logs found by using log-viewing methods described in this topic, you can find the corresponding solutions in the Microsoft Knowledge Base.

System logs

In the left-side navigation pane, choose **Windows Logs > System** to view system logs.

System logs contain events recorded by Windows system components. For example, the failure of a driver or other system components to load during startup is recorded in the system log.

The types of events recorded by system components are predetermined by Windows.

Application logs

In the left-side navigation pane, choose **Windows Logs > Application** to view application logs.

Application logs contain events logged by applications. For example, a database program can record a file error in the application log.

The types of the events recorded in application logs are determined by developers.

Security logs

In the left-side navigation pane, choose **Windows Logs > Security** to view security logs.

Security logs contain valid and invalid logon attempts and events related to resource use, such as creating, opening, or deleting files or other objects.

The types of the events recorded in security logs are determined by administrators. For example, if logon auditing is enabled, the security logs will record logon attempts.

o Application and service logs

An application and service log is a new type of event log. Application and service logs contain events from a single application program or component rather than events that can affect the whole system.

Modify the log path and back up logs

By default, logs are stored on the system disk. The maximum size of logs is 20 MB. If the limit is exceeded, previous events will be overwritten. You can modify the log path as needed.

Perform the following steps to modify the log path and back up logs:

- 1. In the left-side navigation pane of **Event Viewer**, click **Windows Logs**.
- 2. In the right-side list, right-click a log name and choose **Properties** from the shortcut menu.
- 3. In the **Log Properties** dialog box that appears, modify the following parameters:
 - Log path
 - Maximum log size (KB)
 - When maximum event log size is reached:

8.9. Network isolation within a basic security group

Security groups are virtual firewalls that provide stateful packet inspection (SPI) and packet filtering. By default, ECS instances that belong to the same basic security group can access each other over all protocols and ports. Alibaba Cloud provides a variety of network communication policies to allow you to isolate ECS instances within a basic security group.

Internal isolation rules of security groups

The default network communication policy of security groups contains the following items:

- Instances within a basic security group can access each other over all protocols and ports. Instances within an advanced security group are isolated from each other.
- Instances in different security groups are isolated from each other.

Note To allow ECS instances that belong to different security groups to access each other, you can use security group rules for authorization.

To isolate ECS instances within a basic security group from each other, you can modify the network communication policy for the basic security group. For more information, see Modify a network communication policy.

When you configure internal isolation rules for security groups, take note of the following items:

- Internal isolation rules are configured only for specified basic security groups. These isolation rules do not affect the effect of default network communication policies on advanced security groups and other basic security groups.
- Internal isolation rules of security groups implements isolation between network interface controllers,

instead of between ECS instances. If multiple elastic network interfaces (ENIs) are bound to an instance, you must configure internal isolation rules for the security group to which each ENI belongs.

• Internal isolation rules have the lowest priority. When you configure an internal isolation rule to isolate ECS instances within a security group, make sure that no user-defined security group rules apply in the security group to allow communication between the ECS instances.

In the following cases, ECS instances within a security group can still access each other:

- The ECS instances belong to multiple security groups, and no internal isolation rules are configured for one or more of the security groups.
- Internal isolation rules are configured for a security group, while an access control list (ACL) is configured to permit communication between ECS instances within the security group.
 - Note For more information about ACL, see Overview of network ACLs.

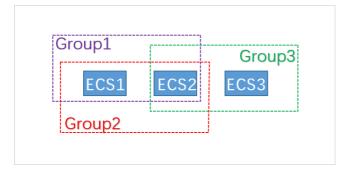
Modify a network communication policy

You can call the ModifySecurityGroupPolicy operation to modify the network communication policy for a basic security group to isolate ECS instances within the security group.

Note By default, ECS instances within an advanced security group are isolated from each other. You cannot modify network communication policies for advanced security groups.

Case analysis

In this example, Group1 and Group2 are basic security groups. ECS1, ECS2, and ECS3 are ECS instances. The following figure shows the relationships between the ECS instances and the security groups:



- Group1 contains ECS1 and ECS2 and is configured with internal isolation rules.
- Group2 contains ECS2 and ECS3 and uses the default network communication policy.

The following table describes whether the ECS instances are isolated from each other.

Instance	Isolated	Description
ECS1 and ECS2	Yes	ECS1 and ECS2 belong to Group1. Group1 is configured with internal isolation rules. Therefore, ECS1 and ECS2 are isolated from each other.
ECS2 and ECS3	No	ECS2 and ECS3 belong to Group2. Group2 uses the default network communication policy. Therefore, ECS2 and ECS3 can access each other.

Instance	Isolated	Description
ECS1 and ECS3	Yes	ECS1 and ECS3 belong to different security groups. By default, ECS instances within different security groups are isolated from each other. Therefore, ECS1 and ECS3 cannot access each other.

8.10. Security group quintuple rules

Security groups are used to configure network access control for one or more ECS instances. As an important means of security isolation, security groups logically isolate security domains on the cloud. Security group quintuple rules allow you to precisely control the following five parameters: source IP address, source port, destination IP address, destination port, and transport layer protocol.

Background information

Previously, security group rules have the following characteristics:

- The inbound rules support only the settings of the source IP address, destination port, and transport layer protocol.
- The outbound rules support only the settings of the destination IP address, destination port, and transport layer protocol.

In most scenarios, these security group rules provide simple configurations, but have the following drawbacks:

- You cannot specify a range of source ports in an inbound rule. Inbound traffic over all ports is allowed by default.
- You cannot specify the destination IP address in an inbound rule. Inbound traffic from all IP addresses within a security group is allowed by default.
- You cannot specify a range of source ports in an outbound rule. Outbound traffic over all ports is allowed by default.
- You cannot specify the source IP address in an outbound rule. Outbound traffic from all IP addresses within a security group is allowed by default.

Definition

A quintuple rule includes the following parameters: source IP address, source port, destination IP address, destination port, and transport layer protocol.

Quintuple rules are designed to provide more fine-grained control over the preceding five parameters, while completely compatible with existing security group rules.

Example quintuple out bound rule:

```
Source IP address: 172.16.1.0/32
Source port: 22
Destination IP address: 10.0.0.1/32
Destination port: no restriction
Transport layer protocol: TCP
Action: Forbid
```

The example out bound rule indicates that TCP access from 172.16.1.0/32 to 10.0.0.1/32 over port 22 is denied.

Scenarios

- Some platform products use solutions from third-party vendors to provide users with network services. To prevent unauthorized access from these products to ECS instances of the users, quintuple rules are required to control inbound and outbound traffic more precisely.
- If ECS instances in a security group are configured to be isolated from each other and you want to allow specified ECS instances to communicate with each other, you must configure quintuple rules.

Configure quintuple rules

You can call API operations to configure quintuple rules.

- To create an inbound security group rule, see AuthorizeSecurityGroup.
- To create an outbound security group rule, see AuthorizeSecurityGroupEgress.
- To delete an inbound security group rule, see RevokeSecurityGroup.
- To delete an outbound security group rule, see RevokeSecurityGroupEgress.

Parameters

The following table describes the parameters of security group rules.

Parameter	Meaning in inbound rules	Meaning in outbound rules
SecurityGroupId	The ID of the security group to which the current inbound rule belongs. This is the ID of the destination security group.	The ID of the security group to which the current outbound rule belongs. This is the ID of the source security group.
Dest Cidrlp	 Optional. The range of destination IP addresses. If DestCidrIp is specified, you can control the range of destination IP addresses in an inbound rule more precisely. If DestCidrIp is not specified, the range of IP addresses in an inbound rule includes all IP addresses in the security group with the specified SecurityGroupId. 	The range of destination IP addresses. Either DestGroupId or DestCidrIp must be specified. If both are specified, DestCidrIp takes priority.
PortRange	Required. The range of destination ports.	Required. The range of destination ports.
Dest GroupId	Manual input is not allowed. The value of DestGroupId must be the same as that of SecurityGroupId.	The ID of the destination security group. You must specify either DestGroupId or DestCidrIp. If you specify both parameters, DestCidrIp takes priority.
SourceGroupId	The ID of the source security group ID. You must specify either SourceGroupId or SourceCidrlp. If you specify both parameters, SourceCidrlp takes priority.	Manual input is not allowed. The value of SourceGroupId must be the same as that of SecurityGroupId.

Parameter	Meaning in inbound rules	Meaning in outbound rules
SourceCidrlp	The range of source IP addresses. You must specify either SourceGroupId or SourceCidrIp. If you specify both parameters, SourceCidrIp takes priority.	 Optional. The range of source IP addresses. If SourceCidrIp is specified, you can control the range of source IP addresses in an outbound rule more precisely. If SourceCidrIp is not specified, the range of IP addresses in an outbound rule includes all IP addresses in the security group with the specified SecurityGroupId.
SourcePortRange	Optional. The range of source ports. If this parameter is not specified, no source ports are restricted.	Optional. The range of source ports. If this parameter is not specified, no source ports are restricted.

8.11. Use Cloud Firewall to control access between ECS instances

Cloud Firewall allows you to centrally manage east-west traffic between ECS instances and north-south traffic between the Internet and ECS instances. This topic describes how to configure Cloud Firewall and view business relationships.

Prerequisites

•

- Cloud Firewall is granted permissions on cloud resources. For more information, see Authorize Cloud Firewall to access other cloud resources.
- Cloud Firewall is of the Enterprise Edition or Ultimate Edition. For more information, see Billing method.

Context

Cloud Firewall features quick switch for firewalls, intrusion detection, outbound connection blocking, traffic analysis, and logging. Cloud Firewall contains internal firewalls, Internet firewalls, and VPC firewalls. For more information about Cloud Firewall terms, see Cloud Firewall and Glossary.

Internal firewalls are used to control east-west traffic and use security group capabilities at the underlying layer. To control east-west traffic between ECS instances, you can create policy groups for internal firewalls in the Cloud Firewall console or add rules to security groups in the ECS console. The configurations of Cloud Firewall and ECS security groups are automatically synchronized. You can also configure application groups to view the access relationships between ECS instances. Based on the access status, you can optimize the policies.

Internet firewalls are used to control north-south traffic between the Internet and ECS instances. You can create inbound or outbound policies and implement policy hardening based on intrusion prevention. For more information, see Traffic analysis overview and Overview of access control policies.

You can use Cloud Firewall in the following scenarios:

• Domain name-based access control.

- Application-based access control.
- Automatic interception of outbound connections started by victim servers.
- Provides access logs of the last six months in compliance with Multi-Level Protection Scheme (MLPS) requirements.

Configure internal firewalls

After you publish a policy group in the Cloud Firewall console, the data of the policy group is immediately synchronized to the corresponding ECS security group. After you configure a security group in the ECS console, data in the security group is synchronized to the corresponding policy group on a daily basis. You can view the synchronization result only the next day. After you purchase Cloud Firewall Enterprise Edition or Ultimate Edition, you can centrally manage the east-west access control policies in the Cloud Firewall console.

Perform the following operations to configure an internal firewall:

- 1. Log on to the Cloud Firewall console.
- 2. In the left-side navigation pane, click Access Control.
- 3. On the page that appears, click Internal Firewall.
 - Source indicates the source of the policy group.
 - o Custom indicates the policy group is created in Cloud Firewall.
 - **Security Group Synchronization** indicates that the policy group is synchronized from an ECS security group.
 - **Application Group Synchronization** indicates that the policy group is synchronized from an application group.
- 4. Click Create Policy Group.
- 5. Configure Name, VPC, Instance ID, Description, and Template, and then click Submit.
 - **? Note** After you configure **VPC**, the policy group is assigned to the region to which the specified VPC belongs, such as in **China** (**Hangzhou**).
- 6. (Optional)Find the target policy group and click **Configure Policy** in the **Actions** column to create a policy.
- 7. Click **Publish** in the corresponding **Actions** column. After the policy is published, it is synchronized to an ECS security group. Follow these steps to view the synchronization result:
 - i. Log on to the ECS console.
 - ii. Select the region where the policy group resides, such as China (Hangzhou).
 - iii. In the left-side navigation pane, choose Network & Security > Security Groups.
 - iv. Set the filter condition to **Security Group Name**. Enter the policy group name in the search bar and click **Search**. If a security group with the same name as the policy group is displayed, the synchronization is successful.

After you configure the internal firewall, you can control the access between ECS instances. You can also configure application groups in Cloud Firewall to visualize business relationships.

View business relationships

In Cloud Firewall, a business group contains all application groups related to a specific business. For example, a web portal business group contains the web application groups and database application groups. An application group is a collection of applications that provide the same or similar services. For example, you can add all ECS instances that are deployed with MySQL to a database application group.

Perform the following operations to view the relationship between ECS instances:

- 1. Log on to the Cloud Firewall console.
- 2. In the left-side navigation pane, choose Business Visualization > Application Groups.
- 3. Create a business group.
 - i. Click the Business Groups tab.
 - ii. Click Create Business Group.
 - iii. Specify the **Name** parameter. For example, you can specify the name as database business or web business.
 - iv. Specify the **Description** parameter.
 - v. Specify the **Importance Degree** parameter. For example, you can set this parameter to **Important**.
- 4. Create an application group.
 - i. Click the **Application Groups** tab.
 - ii. Click Create Application Group.
 - iii. Specify the **Name** parameter. For example, you can specify the name as database business or web business.
 - iv. Specify the **Description** parameter.
 - v. Specify the **Importance Degree** parameter. For example, you can set this parameter to **Important**.
 - vi. Specify the Business Group parameter. For example, you can set this parameter to **Select Existing Business Group**.
 - vii. Select a business group, such as database business or web business.
- 5. Assign applications.
 - i. Select a VPC, such as China (Hangzhou) vpc-xxx.
 - ii. Assign applications based on business requirements. For example, assign all ECS instances that are deployed with MySQL to a web application group.
- 6. In the left-side navigation pane, click **Applications Group**.
- 7. Select a VPC, such as China (Hangzhou) vpc-xxx to view the access relationships between business groups. You can also go to application groups to view the access relationships.

Related information

References

- Configure the access control policy that only allow access to a certain port of outbound to inbound traffic
- Best practices to defend against worms from C&C servers
- Best practices for database security defense

8.12. Enable or disable SELinux

Security-enhanced Linux (SELinux) is a Linux kernel feature that provides a security policy-based protection mechanism for access control. This topic describes how to enable or disable SELinux and avoid system boot failures.

Prerequisites

An ECS instance is created from an Alibaba Cloud public image or a custom image.

Note If the custom image that you use was created from imported local files or migrated from the source server in Server Migration Center (SMC), ensure that SELinux is disabled on the source server before migration.

Context

Typically, enabled SELinux can enhance system security. However, it can damage files in the operating system and lead to system boot failures. If your enterprise or team has high requirements on security and SELinux must be enabled for your operating systems, you can perform operations in this topic to enable SELinux and avoid system boot failures. In this topic, the CentOS 7.2 64-bit operating system is used.

Enable SELinux

1.

2. Run the following command on an instance to modify the config file of SELinux:

```
vi /etc/selinux/config
```

3. Find SELINUX=disabled , press the I key to enter the edit mode, and then enable SELinux by modifying this parameter.

```
# This file controls the state of SELinux on the system.

# SELINUX= can take one of these three values:

# enforcing - SELinux security policy is enforced.

# permissive - SELinux prints warnings instead of enforcing.

# disabled - No SELinux policy is loaded.

SELINUX=disabled

# SELINUXTYPE= can take one of three values:

# targeted - Targeted processes are protected,

# minimum - Modification of targeted policy. Only selected processes are protected.

# mls - Multi Level Security protection.

SELINUXTYPE=targeted
```

You can modify the parameter to one of the following modes:

- SELINUX=enforcing: indicates that all security policy violations will be prohibited.
- SELINUX=permissive: indicates that security policy violations will not be prohibited but will be recorded in the operation logs.
- 4. Press the Esc key and run the :wq command to save and close the file.

- Note After you modify the config file, you must restart the instance for the modification to take effect. However, if you restart the instance directly, the system may fail to start. You need to create an autorelabel file under the root directory before you restart the instance.
- 5. Create the hidden autorelabel file under the root directory. After the instance is restarted, SELinux automatically relabels all system files.

```
touch /.autorelabel
```

6. Restart the ECS instance.

```
shutdown -r now
```

Check SELinux status

1.

2. Run the getenforce command to check the status of SELinux.

The return value can be enforcing or permissive . The return value in this topic is enforcing .

```
[root@ecs01 conf]# getenforce
Enforcing
```

3. Run the sestatus command to query more information about SELinux.

```
[root@ecs01 conf]# sestatus
SELinux status:
                                enabled
SELinuxfs mount:
                                /sys/fs/selinux
SELinux root directory:
                                /etc/selinux
Loaded policy name:
                                targeted
Current mode:
                                enforcing
Mode from config file:
                                enforcing
Policy MLS status:
                                enabled
Policy deny_unknown status:
                                allowed
Max kernel policy version:
                                28
```

If the return value of SELinux status is enabled, SELinux is enabled.

Disable SELinux

1.

2. Run the getenforce command to check the status of SELinux.

If the return value is <code>enforcing</code> , SELinux is enabled.

- 3. Disable SELinux temporarily or permanently.
 - Run the setenforce 0 command to disable SELinux temporarily.

- o Disable SELinux permanently.
 - a. Run the following command to edit the config file of SELinux:

```
vi /etc/selinux/config
```

b. Find SELINUX=enforcing , press the I key to enter the edit mode, and then modify the parameter to SELINUX=disabled .

```
# This file controls the state of SELinux on the system.

# SELINUX= can take one of these three values:

# enforcing - SELinux security policy is enforced.

# permissive - SELinux prints warnings instead of enforcing.

# disabled - No SELinux policy is loaded.

SELINUX=enforcing

# SELINUXTYPE= can take one of three two values:

# targeted - Targeted processes are protected,

# minimum - Modification of targeted policy. Only selected processes are protected.

# mls - Multi Level Security protection.

SELINUXTYPE=targeted
```

- c. Press the Esc key and run the :wq command to save and close the file.
- d. Restart the ECS instance.

```
shutdown -r now
```

e. Run the <code>getenforce</code> command to check the status of SELinux. If the return value is <code>disabled</code>.

What's next

You can create a custom image from an ECS instance that has SELinux enabled. Then, you can create more SELinux-enabled instances from this custom image.

8.13. Revoke the authorization for internal network communication between ECS instances in different accounts through the API

If you have authorized internal network communication between ECS instances across different accounts within the same region, you can revoke security group authorization by calling the API operation.

Prerequisites

- •
- •

Context

In this topic, the RevokeSecurityGroup operation is used to revoke authorized security group rules. Before you start, you must prepare the following information:

- Account name: the name of the account that you use to log on to the ECS console.
- Security group IDs of the ECS instances: the IDs of the security groups to which the instances involved belong.

You can query the security group IDs in the ECS console or by calling the DescribeSecurityGroupReferences operation.

• Region IDs of the ECS instances: See . *cn-beijing* is used in this example.

Assume that the information of the two accounts is as follows.

Account	Account name	Security group	Security group ID
Account A	a@aliyun.com	sg1	sg- bp1azkttqpldxgtedXXX
Account B	b@aliyun.com	sg2	sg- bp15ed6xe1yxeycg7XXX

Procedure

1. Run the following command for Account A:

aliyun ecs RevokeSecurityGroup --SecurityGroupId sg-bplazkttqpldxgtedXXX --RegionId cn-beijing --IpProtocol all --PortRange -1/-1 --SourceGroupId sg-bp15ed6xe1yxeycg7XXX --SourceGroupOwnerAccount b@aliyun.com --NicType intranet

2. Run the following command for Account B:

aliyun ecs RevokeSecurityGroup --SecurityGroupId sg-bp15ed6xe1yxeycg7XXX --RegionId cn-beijing --IpProtocol all --PortRange -1/-1 --SourceGroupId sg-bp1azkttqpldxgtedXXX --SourceGroupOwnerAccount a@aliyun.com --NicType intranet

Related information

- RevokeSecurityGroup
- DescribeSecurityGroupReferences

8.14. Authorize internal network communication between ECS instances in different accounts by using the API

This topic describes how to authorize internal network communication between ECS instances that are in the same region but belong to different accounts.

Prerequisites

Context

You can authorize internal network communication in one of the following modes:

- Authorize internal network communication between ECS instances. You can authorize internal communication between two ECS instances that belong to the same account.
- Authorize internal network communication between accounts. You can authorize internal network communication between ECS instances in two security groups that belong to two different accounts within the same region, including those to be purchased after the authorization is complete.

Note To enable internal network communication between different accounts, you need to authorize communication between security groups in each account. These ECS instances can then communicate over the internal network. If you modify the configurations of a security group, all instances in the security group as well as the services running on these instances are affected. Use caution when you perform this operation.

Security groups are virtual firewalls for ECS instances. Security groups do not provide communication and networking capabilities. After you authorize internal network communication between instances that belong to different security groups, ensure that the instances can establish internal network connection.

- If all instances are of the classic network type, they must be in the same region to communicate with each other.
- VPCs are isolated by default. If all instances are of the VPC type, these instances cannot
 communicate with each other. We recommend that you allow ECS instances to communicate over a
 public network or through Express Connect, VPN Gateway, or Cloud Enterprise Network (CEN). For
 more information, see Express Connect, VPN Gateway, and CEN.
- If instances are of different network types, establish a ClassicLink connection to allow communication between these instances. For more information, see Connect a classic network to a VPC.
- If instances are in different regions, we recommend that you allow ECS instances to communicate over a public network or through Express Connect, VPN Gateway, or CEN. For more information, see Express Connect, VPN Gateway, and CEN.

Authorize internal network communication between ECS instances

1. Query internal IP addresses and security group IDs of the two ECS instances.

You can use the console or call the DescribeInstances operation to obtain security group IDs of the instances. The following table lists the information of the two ECS instances.

Instance	IP address	Security group	Security group ID
Instance A	10.0.0.1	sg1	sg-bp1azkttqpldxgte****
Instance B	10.0.0.2	sg2	sg-bp15ed6xe1yxeycg****

2. Add a rule to sg1 to allow inbound traffic from 10.0.0.2.

```
aliyun ecs AuthorizeSecurityGroup --SecurityGroupId sg-bplazkttqpldxgte**** --RegionId cn-qingdao --IpProtocol all --PortRange=-1/-1 --SourceCidrIp 10.0.0.2 --NicType intran et
```

3. Add a rule to sg2 to allow inbound traffic from 10.0.0.1.

aliyun ecs AuthorizeSecurityGroup --SecurityGroupId sg-bp15ed6xelyxeycg**** --RegionId cn-qingdao --IpProtocol all --PortRange=-1/-1 --SourceCidrIp 10.0.0.1 --NicType intran et



- In the preceding commands, the region ID *cn-qingdao* is for reference only. Replace it with your actual region ID.
- In the preceding commands, the AuthorizeSecurityGroup operation is called to add inbound rules to security groups. Specify the SecurityGroupId and SourceCidrIp parameters.
- 4. After a few minutes, run the **ping** command to check whether the two ECS instances can communicate with each other over the internal network.

Authorize internal network communication between accounts

1. Query names and security group IDs of the two accounts.

You can use the console or call the DescribeInstances operation to obtain security group IDs of the ECS instances. The following table lists the information of two accounts.

Account	Account ID	Security group	Security group ID
Account A	a@aliyun.com	sg1	sg-bp1azkttqpldxgte****
Account B	b@aliyun.com	sg2	sg-bp15ed6xe1yxeycg****

2. Add a rule to sg1 to allow inbound traffic from sg2.

aliyun ecs AuthorizeSecurityGroup --SecurityGroupId sg-bplazkttqpldxgte**** --RegionId cn-qingdao --IpProtocol all --PortRange=-1/-1 --SourceGroupId sg-bp15ed6xelyxeycg7XXX --SourceGroupOwnerAccount b@aliyun.com --NicType intranet

3. Add a rule to sg2 to allow inbound traffic from sg1.

aliyun ecs AuthorizeSecurityGroup --SecurityGroupId sg-bp15ed6xe1yxeycg**** --RegionId cn-qingdao --IpProtocol all --PortRange=-1/-1 --SourceGroupId sg-bp1azkttqpldxgtedXXX --SourceGroupOwnerAccount a@aliyun.com --NicType intranet

? Note

- In the preceding commands, the region ID *cn-qingdao* is for reference only. Replace it with your actual region ID.
- In the preceding commands, the AuthorizeSecurityGroup operation is called to add inbound rules to security groups. Specify the SecurityGroupId, SourceGroupId, and SourceGroupOwnerAccount parameters.
- 4. After a few minutes, run the **ping** command to check whether the ECS instances can communicate with each other over the internal network.

9.Data recovery9.1. Handle low disk space on Windows instances

This topic describes how to handle low disk space on Windows instances and the best practices of daily disk management for Windows instances.

Context

The methods described in this topic apply to Windows Server 2003 and later. In this topic, Windows Server 2012 R2 is used.

Solutions and best practices

When you are running out of available disk space on a Windows instance, you can take one of the following measures:

- Release disk space
- Resize the disk

Develop good habits for using disk space with these best practices:

- Compress files for archiving
- Delete unneeded applications on a regular basis
- Configure disk monitoring

Release disk space

When you are running out of available disk space on a Windows instance, you can clear unnecessary files from the disk. Perform the following operations:

- 1. Find the files that are taking up the most disk space.
 - i. Connect to the Windows instance. For more information, see Connect to a Windows instance from a local client.
 - ii. Click Start and then click This PC.
 - iii. Click the disk to be cleared, and press Ctrl+F.
 - iv. In the top navigation bar, choose **Search > Size**, and filter out large files on the specified disk.
 - ? Note You can also customize a file size range for retrieval. For example:
 - Enter Size: > 500 MB to search for files that are larger than 500 MB in size.
 - Enter Size: > 100 MB < 500 MB to search for files that are between 100 MB and 500 MB in size.
- 2. Delete files that you no longer need.

We recommend that you use the Windows Disk Cleanup utility to delete unneeded files such as log files, and empty the recycle bin. The Windows Disk Cleanup utility is not installed on the instance by default and must be installed manually. Perform the following operations to install the utility and use it to delete files:

- i. In the taskbar, click the **Server Manager** icon.
- ii. In the upper-right corner, choose Manage > Add Roles and Features.
- iii. In the Add Roles and Features Wizard dialog box, keep the default settings and click Next until the Features module is displayed. Select Ink and Handwriting Services and Desktop Experience, and click Next.
- iv. Click Inst all.
- v. After the utility is installed, the system prompts you to restart the instance. You must restart the instance manually. After the instance restarts, verify that Desktop Experience has been installed.
- vi. Click **Start**. In the top search box, enter **Disk Cleanup**. In the dialog box that appears, select the disk that you want to delete and click **OK**.

Resize the disk

When you are running out of available disk space on a Windows instance, you can resize the disks. For more information, see Resize disks online for Windows instances or Resize disks offline for Windows instances.

Compress files for archiving

For files that are generated on a regular basis, you can compress them for archiving to improve disk usage. We recommend that you use WinRAR to compress files. The following example describes how to configure a backup policy for compressing files:

- 1. Download and install WinRAR. Download link: WinRAR and RAR archiver downloads. WinRAR is used in this example.
- 2. After WinRAR is installed, find the file to be compressed, right-click the file, and then select **Add** to archive...
- 3. In the Settings dialog box, click the Backup tab and select Generate archive name by mask. Do not click OK.
- 4. Click the **General** tab and then click **Browse** to specify the path in which to save the archive file. Click **Profiles...** and select **Save current settings to a new profile...**
- 5. In the Profile parameters dialog box, set the Profile name parameter and select Save archive name, Save selected file names, and Create shortcut on desktop. Click OK.
- 6. In the Archive name and parameters dialog box, click OK. A shortcut key to the archive file is generated on the desktop.
- 7. Press Win+R to open the Run dialog box. Run the control command to open the Control Panel page. On the Control Panel page, click System and Security and then click Scheduled tasks. In the Task Scheduler dialog box, click Create Basic Task...
- 8. In the Create Basic Task Wizard dialog box, enter a new task name and click Next.
- 9. Select a trigger for the task and click Next. Select $Start\ a\ program\ and\ click\ Next$.
- 10. In the Start a Program step, set the **Program/script**: parameter. To set this parameter, find the previously generated shortcut key. Right-click the shortcut key and then select **Properties**. In the dialog box that appears, copy the value of the **Target**: field.
- 11. In the Create Basic Task Wizard dialog box, paste the copied value in the **Program/script**: field of the **Start a Program** step. Click **Finish**.

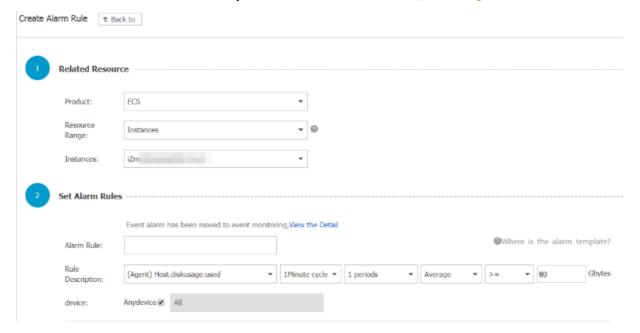
After you configure the backup policy, you can delete expired archive files on a regular basis to avoid taking up a large amount of disk space.

Delete unneeded applications on a regular basis

You can delete unneeded applications on a regular basis by clicking **Programs and Features** on the **Control Panel** page.

Configure disk monitoring

Monitoring plug-ins are pre-installed on ECS instances. You can log on to the Cloud Monitor console to create disk alert rules. In this way, you can receive alerts when disk usage exceeds a configured threshold and clean the disk in a timely manner. For more information, see Manage alert rules.



9.2. Restore data in Linux instances

When you troubleshoot disks, you may encounter the loss of data disk partitions. This topic describes data disk partition loss in Linux and the solutions. This topic also describes the common mistakes and best practices for using disks to avoid the risk of data loss.

Prerequisites

 A snapshot is created for the data disk that lost a partition. If errors occur during data restoration, you can use the snapshot to roll back the data disk to the state before the restoration. For more information, see Create a snapshot of a disk and Roll back a disk by using a snapshot.

Context

In a Linux instance, you can use one of the following tools to restore data on a data disk:

- fdisk: a tool provided by Linux for partitioning disks.
- testdisk: a tool used to restore disk partitions or data in Linux. By default, the tool is not provided in Linux. You must install it on your own. For example, you can run the **yum install -y testdisk** command to install testdisk in CentOS.
- partprobe: a tool provided by Linux. The tool is used to enable the kernel to re-read partitions without restarting the system.

> Document Version: 20220711

Methods

After you restart a Linux instance, the partition or data of data disks may be lost. This may be because you did not set the partition to be mounted automatically upon startup of the instance in the *etc/fstab* file. In this case, you can manually mount the data disk partition. If the system prompts partition table loss when you manually mount the data disk partition, you can use one of the following methods to restore the partition or data:

- Restore a partition by using fdisk
- Restore a partition by using testdisk
- Restore data by using testdisk

Restore a partition by using fdisk

Typically, default values apply to the starting and ending sectors of a partition when you partition a data disk. You can first run the **f disk** command to restore the partition. For more information, see Format a data disk for a Linux instance.

```
[root@Aliyun ~]# fdisk /dev/xvdb
Welcome to fdisk (util-linux 2.23.2).

changes will remain in memory only, until you decide to write them.
Be careful before using the write command.

Command (m for help): n
Partition type:
    p primary (0 primary, 0 extended, 4 free)
    e extended
Select (default p): p
Partition number (1-4, default 1): 1
First sector (2048-10485759, default 2048):
Using default value 2048
Last sector, +sectors or +size{K,M,G} (2048-10485759, default 10485759):
Using default value 10485759
Partition 1 of type Linux and of size 5 GiB is set

Command (m for help): w
The partition table has been altered!

calling ioctl() to re-read partition table.
Syncing disks.
[root@Aliyun ~]# mount /dev/xvd
xvda xvda1 xvdb xvdb1
[root@Aliyun ~]# mount /dev/xvdb
xvdb xvdb1
[root@Aliyun ~]# mount /dev/xvdb
xvdb xvdb1
[root@Aliyun ~]# mount /dev/xvdb1 /mnt/
[root@Aliyun ~]# mount /dev/xvdb1 /mnt/
[root@Aliyun ~]# ls /mnt/
123.sh confiqclient data diamond install_edsd.sh install.sh ip.qz
```

If the preceding operation cannot restore the partition, you can use test disk to restore the partition.

Restore a partition by using testdisk

A disk device named /dev/xvdb is used in this example. To use test disk to restore a partition, follow these steps:

1. Run the **test disk /dev/xvdb** command (you can replace the device name), select **Proceed** (default value), and then press the Enter key:

```
TestDisk 7.0, Data Recovery Utility, April 2015
Christophe GRENIER <grenier@cgsecurity.org>
http://www.cgsecurity.org

TestDisk is free software, and comes with ABSOLUTELY NO WARRANTY.

Select a media (use Arrow keys, then press Enter):
>Disk /dev/xvdb - 5368 MB / 5120 MiB

>[Proceed ] [ Quit ]

Note: Disk capacity must be correctly detected for a successful recovery. If a disk listed above has incorrect size, check HD jumper settings, BIOS detection, and install the latest OS patches and disk drivers.
```

2. Select the partition table type for scanning. Typically, the default value *Intel* is selected. Select *EFI GPT* if your data disk uses the GPT format.

```
TestDisk 7.0, Data Recovery Utility, April 2015
Christophe GRENIER <grenier@cgsecurity.org>
http://www.cgsecurity.org
Disk /dev/xvdb - 5368 MB / 5120 MiB
Please select the partition table type, press Enter when done.
Intel I Intel/PC partition
  EFI GPT
Humax ]
Mac ]
None ]
Sun ]
             EFI GPT partition map (Mac i386, some x86_64...)
             Humax partition table
  [Humax
  Гмас
             Apple partition map
             Non partitioned media
             Sun Solaris partition
XBox partition
  Sun
  XBOX
 [Return ] Return to disk selection
Note: Do NOT select 'None' for media with only a single partition. It's very
rare for a disk to be 'Non-partitioned'.
```

3. Select *Analyse* and press the Enter key.

```
Disk /dev/xvdb - 5368 MB / 5120 MiB
CHS 652 255 63 - sector size=512

Analyse | Analyse current partition structure and search for lost partitions Filesystem Utils [Geometry Change disk geometry Options Modify options | Modify options | Modify options | MBR Code | Write TestDisk MBR code to first sector | Delete | Delete all data in the partition table | Quit | Return to disk selection

Note: Correct disk geometry is required for a successful recovery. 'Analyse' process may give some warnings if it thinks the logical geometry is mismatched.
```

4. Select Quick Search and press the Enter key if the partition information is not displayed.

The partition information is displayed in the command output, as shown in the following figure.

```
Disk /dev/xvdb - 5368 MB / 5120 MiB - CHS 652 255 63
Partition Start End Size in sectors

** Linux 0 32 33 652 180 40 10483712

Structure: Ok. Use Up/Down Arrow keys to select partition.
Use Left/Right Arrow keys to CHANGE partition characteristics:
*=Primary bootable P=Primary L=Logical E=Extended D=Deleted
Keys A: add partition, L: load backup, T: change type, P: list files,
Enter: to continue
```

- 5. Select the partition and press the Enter key.
- 6. Select Write to save the partition.
 - Note Select Deeper Search to continue to search if the expected partition is not listed.

```
Disk /dev/xvdb - 5368 MB / 5120 MiB - CHS 652 255 63

Partition Start End Size in sectors

1 * Linux 0 32 33 652 180 40 10483712

[ Quit ] [Deeper Search] | Write | partition Structure to disk
```

7. Press the Ykey to save the partition.

```
TestDisk 7.0, Data Recovery Utility, April 2015
Christophe GRENIER <grenier@cgsecurity.org>
http://www.cgsecurity.org
Write partition table, confirm ? (Y/N)
```

- 8. Run the partprobe /dev/xvdb command (you can replace the device name) to manually refresh the partition table.
- 9. Mount the partition again and view the data in the data disk.

```
[root@Aliyun home]# mount /dev/xvdb1 /mnt/
[root@Aliyun home]# ls /mnt/
123.sh configclient data diamond install_edsd.sh install.sh ip.gz logs lost+found test
```

Restore data by using testdisk

In some cases, you can scan and locate a disk partition by using test disk. However, you cannot save the partition. In this case, you can directly restore the data. Follow these steps:

- 1. Scan and locate a disk partition by using testdisk. For more information, see Step 1 to Step 4 in Restore a partition by using testdisk.
- 2. Press the *P* key to list files.

 The following figure shows the command output.

* Linux Directory /		3,000	0 32 33 652 180 40 10483712
drwxr-xr-x	0	0	4096 21-Feb-2017 11:57 .
drwxr-xr-x	0	0	4096 21-Feb-2017 11:57
drwx	0	0	16384 21-Feb-2017 11:56 lost+found
-rw-rr	0	0	1701 21-Feb-2017 11:57 install_edsd.sh
-rw-rr	0	0	5848 21-Feb-2017 11:57 install.sh
-rw-rr	0	0	12136 21-Feb-2017 11:57 ip.gz
-rw-rr	0	0	0 21-Feb-2017 11:57 test
drwxr-xr-x	0 0 0	0	4096 21-Feb-2017 11:57 123.sh
drwxr-xr-x	0	0	4096 21-Feb-2017 11:57 configclient
drwxr-xr-x	0	0	4096 21-Feb-2017 11:57 data
drwxr-xr-x	0	0	4096 21-Feb-2017 11:57 diamond
drwxr-xr-x	0	0	4096 21-Feb-2017 11:57 logs
			Next
Use Right to q to quit C to copy	change t, : to the s	direct select	ory, h to hide deleted files the current file, a to select all files files. c to copy the current file

- 3. Select the file that you want to restore and press the Ckey.
- 4. Select the destination directory. In this example, the file is restored to the /home directory.

```
Please select a destination where /ip.gz will be copied.
Keys: Arrow keys to select another directory
       C when the destination is correct
       Q to quit
Directory /
                     0
                             0
                                      4096 11-Jan-2017 09:32 .
 drwxr-xr-x
                                      4096 11-Jan-2017 09:32
                     0
                             0
 drwxr-xr-x
                                      4096 25-Jul-2016 16:23 boot
2940 21-Feb-2017 12:30 dev
4096 21-Feb-2017 12:12 etc
                     0
                             0
 dr-xr-xr-x
 drwxr-xr-x
                     0
                             0
 drwxr-xr-x
                     0
                             0
                                      4096 16-Feb-2017 11:48 home
>drwxr-xr-x
                     О
                             0
                                     16384 12-May-2016 19:58 lost+found
4096 12-Aug-2015 22:22 media
4096 21-Feb-2017 11:57 mnt
4096 12-Aug-2015 22:22 opt
0 16-Feb-2017 21:35 proc
4096 21-Feb-2017 11:57 root
 drwx----
                             O
 drwxr-xr-x
                     0
                             0
 drwxr-xr-x
                     0
                             0
                     0
                             0
 drwxr-xr-x
                     0
                             0
 dr-xr-xr-x
                     0
                             0
 dr-xr-x---
                                       560 21-Feb-2017 12:12 run
                     0
                             0
 drwxr-xr-x
                                      4096 12-Aug-2015 22:22 srv
0 16-Feb-2017 21:35 sys
                     0
                             0
 drwxr-xr-x
 dr-xr-xr-x
                     0
                             0
                                      4096 21-Feb-2017 12:34 tmp
                     0
 drwxrwxrwt
                             0
                                      4096 16-Feb-2017 11:48 usr
                     0
                             0
 drwxr-xr-x
                                      4096 16-Feb-2017 21:35 var
                     0
                             0
 drwxr-xr-x
                                              3-May-2016 13:48 bin
                     0
                             0
 lrwxrwxrwx
                                              3-May-2016 13:48 lib
                     0
                             0
                                          7
 1rwxrwxrwx
                                              3-May-2016 13:48 lib64
 1rwxrwxrwx
                     0
                             0
 1rwxrwxrwx
                     0
                             0
                                              3-May-2016 13:48 sbin
                                          8
```

If Copy done! 1 ok, 0 failed is displayed, the file is copied, as shown in the following figure.

* Linux Directory /			0	32 33	652	180 40	10483712
Copy done! 1	ok, 0	failed					
drwxr-xr-x	0	0				11:57	
drwxr-xr-x	0	0				11:57	
drwx	0	0					lost+found
-rw-rr	0	0	1701	21-F6	eb-2017	11:57	install_edsd.sh
-rw-rr	0	0	5848	21-F6	eb-2017	11:57	install.sh
>-rw-rr	0	0	12136	21-F€	eb-2017	11:57	ip.gz
-rw-rr	0	0	0	21-F6	eb-2017	11:57	test
drwxr-xr-x	0	0	4096	21-F6	eb-2017	11:57	123.sh
drwxr-xr-x	0	0	4096	21-F6	eb-2017	11:57	configclient
drwxr-xr-x	0	0				11:57	
drwxr-xr-x	0	0	4096	21-F6	b-2017	11:57	diamond
drwxr-xr-x	0	0				11:57	

5. Switch to the /home directory to view details.

If the file is displayed, as shown in the following figure, the file is restored.

```
[root@Aliyun /]# ls /home/
admin<mark> ip.qz</mark>
[root@Aliyun /]#
```

Common mistakes and best practices

Data is the core asset of users. A large number of users build websites and databases such as MySQL, MongoDB, and Redis on ECS instances. Data loss may cause huge risks to businesses. This section describes the common mistakes and best practices in data security.

Common mist akes

The underlying storage of Alibaba Cloud is based on triplicate technology. Therefore, some users consider that no risk of data loss in the operating system exists. This is a misunderstanding. The three copies of data stored in the underlying layer provide physical layer protection for data disks. However, if errors occur to the cloud disk logic in the system, such as infection with viruses, accidental data deletion, and file system damage, the data may still be lost. You must use technologies such as snapshots and geo-redundancy to ensure data security. For more information about three copies, see Triplicate storage.

• Best practices

Data disk partition restoration and data restoration are the final solutions to data loss problems, but they may not restore data as expected. We recommend that you follow the best practices to create automatic or manual snapshots for data and run different backup schemes to maximize your data security.

Apply automatic snapshot policies

Automatic snapshot policies are applied to system and data disks to create automatic snapshots for the disks. Note that after the system disk is replaced, the instance expires, or the disk is manually released, the corresponding automatic snapshots may be released.

If you want automatic snapshots of a disk to be released along with the disk, you can select **Delete Automatic Snapshots While Releasing Disk** in the **Modify Disk Property** dialog box in the ECS console. If you want to retain the automatic snapshots, you can clear this option.

For more information, see Snapshot FAQ and Create an automatic snapshot policy.

Create manual snapshots

Before you perform important or high-risk operations, you must manually create snapshots for the disk. These operations include:

- Update the kernel.
- Upgrade or change applications.
- Rest ore data on disks.

Before you restore a disk, you must create a snapshot for the disk. After the snapshot is created, you can perform other operations.

OSS backup, offline backup, and geo-redundancy
 You can back up important data by using OSS backup, offline backup, or geo-redundancy.

9.3. Restore data in Windows instances

When you troubleshoot disks, you may encounter the loss of data disk partitions. This topic describes data disk partition loss in Windows and the solutions. This topic also describes the common mistakes and best practices for using disks to avoid the risk of data loss.

Prerequisites

•

• A snapshot is created for the data disk that lost a partition. If errors occur during data restoration, you can use the snapshot to roll back the data disk to the state before the restoration. For more information, see Create a snapshot of a disk and Roll back a disk by using a snapshot.

Context

In a Windows instance, you can use one of the following tools to restore data on a data disk:

- Disk Management: a tool provided by Windows for partitioning and formatting data disks.
- Data restoration software: typically, commercial software. You can download the software from the official websites. The software is used for restoring data in file systems to which exceptions occur.

Status of the disk is Foreign and no partitions are displayed

In **Disk Management** of Windows, the disk is in the **Foreign** state and no partitions are displayed.

Solution:

Right-click the space next to Foreign, select Import Foreign Disks, and click OK.

Status of the disk is Offline and no partitions are displayed

In Disk Management of Windows, the disk is in the Offline state and no partitions are displayed.

Solution:

Right-click the space next to **Offline**, select **Online**, and then click **OK**.

No driver letter assigned

In **Disk Management** of Windows, you can find the information of the data disk, but no drive letter is assigned to the data disk.

Solution:

Right-click the primary partition of the disk (such as **Disk 1**), select **Change Drive Letter and Paths**, and then follow the prompt to complete other operations.

Error occurred during storage enumeration

In **Disk Management** of Windows, you cannot view data disks. An error message similar to **An error occurred during storage enumeration** is reported in the system log.

Note The reported content may be An error occurred during enumeration of volumes based on your operating system version.

Solution:

- 1. Start Windows PowerShell.
- 2. Run the winrm quickconfig command to restore data.

When Make these changes [y/n]? is displayed, enter y to run the command.

After the restoration, you can find the data disk in Disk Management.

Data disk in the RAW format

In some cases, you may find that the data disk in Windows is in the RAW format.

If the file system of a disk is unrecognizable to Windows, the disk is displayed in the RAW format. Typically, this occurs when the partition table or boot sector that records the type or location of the file system is lost or damaged. The following common causes may lead to the loss or damage:

- Safely remove hardware is not used when the external disk is disconnected.
- Disk problems caused by power outages or unexpected shutdowns.
- Hardware faults occur.
- Underlying disk-related drivers or applications. For example, DiskProbe can be used to directly modify the disk table structure.
- Computer viruses.

For information about how to restore data disks, visit Dskprobe Overview in the Microsoft documentation.

Windows also contains a large variety of free or commercial data restoration software to restore lost data. For example, you can use Disk Genius to scan and restore expected documents.

Common mistakes and best practices

Data is the core asset of users. A large number of users build websites and databases such as MySQL, MongoDB, and Redis on ECS instances. Data loss may cause huge risks to businesses. The following section describes the common mistakes and best practices in data security.

Common mistakes

The underlying storage of Alibaba Cloud is based on triplicate technology. Therefore, some users consider that data loss will not occur in the operating system. This is a misunderstanding. The three copies of data stored in the underlying layer provide physical layer protection for data disks. However, if errors occur to the cloud disk logic in the system, such as infection with viruses, accidental data deletion, and file system damage, the data may still be lost. In this case, you must use technologies such as snapshots and geo-redundancy to ensure data security.

Best practices

Data disk partition restoration and data restoration are the final solutions to data loss problems, but they may not restore data as expected. We recommend that you follow the best practices to create automatic or manual snapshots for data and run different backup schemes to maximize your data security.

Apply automatic snapshot policies

Automatic snapshot policies are applied to system and data disks to create automatic snapshots for the disks. Note that after the system disk is replaced, the instance expires, or the disk is manually released, the corresponding automatic snapshots may be released.

If you want automatic snapshots of a disk to be released along with the disk, you can select **Delete Automatic Snapshots While Releasing Disk** in the **Modify Disk Property** dialog box in the ECS console. If you want to retain the automatic snapshots, you can clear this option.

For more information, see Snapshot FAQ and Delete automatic snapshots while releasing a disk.

Create manual snapshots

Before you perform important or high-risk operations, you must manually create snapshots for disks. These operations include:

- Update the kernel.
- Upgrade or change applications.
- Restore data on disks.

Before you restore a disk, you must create a snapshot for the disk. After the snapshot is created, you can perform other operations.

• OSS backup, offline backup, and geo-redundancy

You can back up important data by using OSS backup, offline backup, or geo-redundancy.

10.Configuration preference 10.1. Transfer ECS instance data

This topic describes the file transfer principle and typical transfer methods on Unix-like, Linux, and Windows platforms in Alibaba Cloud ECS. Additionally, this topic compares these methods to help you choose appropriate file transfer methods that meet your specific requirements.

File transfer principle

File transfer, also known as file data communication, is a form of information transfer that transmits file data between data sources and data sinks. In a file transfer process, the OS extracts file data to the memory for temporary storage, and then duplicates the data to the destination. Encryption adds a secure layer to a file, while duplication transfers the encrypted file as a whole to another location. Decryption is needed only when a compressed package is opened. A large file cannot be transferred as a whole between hosts immediately because the transfer is a continuous process. If any interruption occurs during the transfer, the file will not exist in the destination path. If multiple files are transferred, they are transferred separately and sequentially. If any interruption occurs during the transfer, the files that are being transferred or have not yet been transferred will fail, but the transferred files are transferred successfully. A compressed package is considered as one file regardless of how many files the package contains.

Multiple file transfer tools, such as Netcat, FTP, SCP, and NFS, can be used to transfer files. The following sections describe the features and usage of some typical file transfer tools.

Netcat

Net cat is a powerful and versatile networking tool with optimal file transfer capabilities.

Parameter descriptions

Parameter	Description
-g <gateway></gateway>	Specifies up to eight long-distance communication gateways for the router.
-G < number of indicators>	Specifies the number of source routing indicators. The value is a multiple of 4.
-i <delay in="" seconds=""></delay>	Specifies the time interval for sending messages and scanning the communications port.
-l	Enables the listening mode to control received data.
-o <output file=""></output>	Specifies the name of the file where the transferred data is dumped and saved in hexadecimal character codes.
-P <communication port=""></communication>	Specifies the communication port used by the local host.
-r	Specifies the communication port between the local host and the remote host.
-u	Enables the UDP transfer protocol.

Parameter	Description
-V	Shows the command running process.
-w <timeout in="" seconds=""></timeout>	Specifies the waiting time for a connection.
-z	Enables the zero input/output mode, which is used only for scanning the communications port.
-n	Uses IP addresses instead of the DNS.

Examples of usage

1. Scan ports 21-24 (for example, IP address 192.168.2.34).

```
nc -v -w 2 192.168.2.34 -z 21-24
```

Response example:

```
nc: connect to 192.168.2.34 port 21 (tcp) failed: Connection refused Connection to 192.168.2.34 22 port [tcp/ssh] succeeded!
nc: connect to 192.168.2.34 port 23 (tcp) failed: Connection refused nc: connect to 192.168.2.34 port 24 (tcp) failed: Connection refused
```

- 2. Copy files from 192.168.2.33 to 192.168.2.34.
- Run the following command at 192.168.2.34: nc-1 1234 > test.txt .
- Run the following command at 192.168.2.33: nc192.168.2.34 < test.txt .
- 3. Run the following nc commands to operate Memcached as needed:
- To store data, run the command printf "set key 0 10 6rnresultrn" |nc 192.168.2.34 11211 .
- To obtain data, run the command printf "get keyrn" |nc 192.168.2.34 11211 .
- To delete data, run the command printf "delete keyrn" |nc 192.168.2.34 11211 .
- To view the status, run the command printf "statsrn" | nc 192.168.2.34 11211 .
- To simulate the top command to view the status, run the command watch "echo stats" |nc 192.1 68.2.34 11211 .
- To clear the cache, run the following command:

```
printf "flush_allrn" |nc 192.168.2.34 11211  # This operation cannot be undone.
```

SCP

The use of Secure Copy (SCP) commands is similar to that of RCP commands. The difference is that SCP commands provide higher security protection by prompting users to enter a password for verification. Therefore, we recommend that you use SCP commands instead of RCP commands. SCP commands use SSH to transfer data and use the same authentication model as SSH to provide the same security protection. SSH is a reliable protocol that provides security for remote logon sessions and other network services. With SSH, you can effectively prevent information leakage during remote management. SCP is an SSH-based application. Therefore, it requires that the machines involved in data transfer support SSH.

Features

Similar to RCP, SCP can retain the file attributes on a specific file system and retain the copied subdirectories that need recursion.

SCP provides better file transfer confidentiality. Overall, SCP is suitable for users with high data security requirements.

Examples of usage

If you do not want to enter your username and password every time you use SCP commands to copy files between two machines, you can configure SSH.

To generate an RSA key, run the following command:

When you are prompted to enter the path and password to save the key, you can press Enter to use the default path and a null password. Then, the generated public key is saved in /.ssh/id_rsa.pub, and the private key is saved in /.ssh/id_rsa. You can copy the content of the public key from this key pair to the /.ssh/authorized_keys file in the machine that you want to access. In this way, you do not need to enter your password when you next access this machine.

Copy a file between two Linux hosts

Basic command format:

```
scp [optional parameter] file_source file_target
```

To copy a file from a local directory to a remote directory, run one of the following four commands:

```
scp local file remote username@remote ip:remote folder
scp local file remote username@remote ip:remote file
scp local file remote ip:remote folder
scp local file remote ip:remote file
```

Note In the first and second commands, user names are specified and the password must be entered after the commands are executed. In the first command, a remote directory is specified and the file name remains the same. In the second command, a file name is specified.

In the third and fourth commands, user names are not specified and the password must be entered after the commands are executed. In the third command, a remote directory is specified and the file name remains the same. In the fourth command, a file name is specified.

To copy a file from a remote directory to a local directory, run the following commands:

```
scp root@www.cumt.edu.cn:/home/root/others/music /home/space/music/i.mp3
scp -r www.cumt.edu.cn:/home/root/others/ /home/space/music/
```

Rsync

Rsync is a file synchronization and transfer tool for Linux or Unix. As an alternative to RCP, Rsync may be used through RSH or SSH, or run in daemon mode. In daemon mode, the Rsync server opens port 873 for client connections. During client connections, the Rsync server will verify the password. If the password is correct, the file can be transferred. During the first connection, the entire file is transferred. During the subsequent connections, only incremental data of the file is synchronized.

Rsync Installation methods

Note You can use the package manager of your OS to install Rsync.

```
sudo apt-get install rsync
                                # Install Rsync online in Debian and Ubuntu.
slackpkg install rsync
                                # Install Rsync online using Slackware.
yum install rsync
                                # Install Rsync in Fedora and Red Hat.
```

Install Rsync through source code compilation:

```
wget http://rsync.samba.org/ftp/rsync/src/rsync-3.0.9.tar.gz
tar xf rsync-3.0.9.tar.gz
cd rsync-3.0.9
./configure && make && make install
```

Parameter descriptions

Parameter	Description
-V	Specifies the output mode.
-a	Specifies the archive mode. It meas that files are transferred recursively and all file attributes are retained. This parameter is equivalent to the combined parameter - rlptgoD.

Parameter	Description
-r	Transfers subdirectories recursively.
-l	Retains soft links.
-р	Retains file permissions.
-t	Retains file time information.
-g	Retains file group information.
-0	Retains file owner information.
-D	Retains device file information.
-H	Retains hard links.
-S	Processes sparse files explicitly to save space for DST files.
-Z	Compresses backup files during transfer.

Six work modes of Rsync

• To copy local files from the /home/coremail directory to the /cmbak directory, run the following command:

```
rsync -avSH /home/coremail/ /cmbak/
```

• To copy files from a local machine to a remote machine, run the following command:

```
rsync -av /home/coremail/ 192.168.11.12:/home/coremail/
```

• To copy files from a remote machine to a local machine, run the following command:

```
rsync -av 192.168.11.11:/home/coremail/ /home/coremail/
```

• To copy files from a remote Rsync server (running in daemon mode) to a local machine, run the following command:

```
rsync -av root@172.16.78.192::www /databack
```

• To copy files from a local machine to a remote Rsync server (running in daemon mode), run the following command. This work mode is started when the DST path information contains the "::" delimiter.

```
rsync -av /databack root@172.16.78.192::www
```

• To show the file list of a remote machine, run the following command:

```
rsync -v rsync://192.168.11.11/data
```

Description of the Rsync configuration file

```
cat/etc/rsyncd.conf
                              # The contents are as follows:
port = 873
                              # Specify the port number.
uid = nobody
                                \# Specify the UID of the daemon process when the module tr
ansfers files.
gid = nobody
                                # Specify the GID of the daemon process when the module t
ransfers files.
                                # Use chroot to enter the directories in the file system.
use chroot = no
max connections = 10
                                # Specify the maximum concurrent connections.
strict modes = yes
                                # Specify whether to check the permissions of password-pro
tected files.
pid file = /usr/local/rsyncd/rsyncd.pid
                                           # Specify PID files.
lock file = /usr/local/rsyncd/rsyncd.lock  # Specify the lock file that supports the max
imum concurrent connections. By default, the lock file is /var/run/rsyncd.lock.
motd file = /usr/local/rsyncd/rsyncd.motd #Define server information and write the rsyn
log file = /usr/local/rsyncd/rsync.log  # Specify the log of the Rsync server.
log format = %t %a %m %f %b
syslog facility = local3
timeout = 300
                                        # custom module
[conf]
path = /usr/local/nginx/conf
                                        # Specify the directory to be backed up.
comment = Nginx conf
ignore errors
                                        # Ignore some IO errors.
                                       # To allow the client to upload files, set the val
read only = no
ue to no. Otherwise, set the value to yes.
                                       # To allow the client to download files, set the v
write only = no
alue to no. Otherwise, set the value to yes.
hosts allow = 192.168.2.0/24
                                       # Specify an allowed IP address.
hosts deny = *
                                        # Specify a denied IP address.
list = false
                                        # Use the module list upon request.
uid = root
gid = root
auth users = backup
                                        # Specify a connection user name, which is irrelev
ant to Linux user names.
secrets file = /etc/rsyncd.pass
                                        # Specify the password file.
```

10.2. Increase data throughput through read/write splitting

Typically, system performance decreases when reads and writes occur in the same database server. To improve overall system performance and optimize user experience, you can reduce the load of your primary database through read/write splitting. This topic describes how to use MySQL Proxy to split read and write operations.

Prerequisites

Context

At the application layer, read/write splitting is implemented through coding. Before you enter the service layer, Aspect-Oriented Programming (AOP) is used to determine whether to use the read database or the write database. The method names can be used to implement the target action. For example, the read database is used for method names that start with query, find, or get, and the write database is used for others.

Advantages:

- The program automatically switches among multiple data sources with ease.
- Middleware is not required.
- Theoretically, all databases are supported.

Disadvant ages:

- Manual operations are not supported.
- Dat a sources cannot be dynamically added.

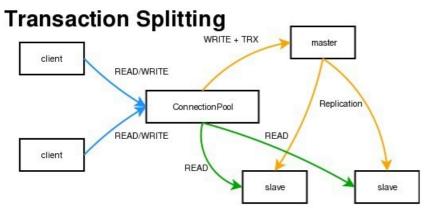
You can use one of the following methods to split read and write operations at the system layer:

- Distributed Relational Database Service (DRDS)
- MySQL Proxy

The following section describes how to use MySQL Proxy to split read and write operations.

MySQL Proxy is a simple program that is situated between your client and MySQL server and can monitor, analyze, or transform communication between the server and the client. It can serve a wide variety of purposes, such as load balancing, fault query and analysis, and query filtering and modification.

The following figure shows the principle of MySQL Proxy.



MySQL Proxy is an intermediate-layer proxy that acts as a connection pool to forward connection requests from front end applications to the backend database. MySQL Proxy can perform complex connection control and filtering to implement read/write splitting and load balancing by using the Lua script. MySQL Proxy allows applications to access the backend database smoothly. The applications only need to be connected to the listening port of MySQL Proxy. In this case, the proxy server may become a single point of failure (SPOF). You can use multiple proxy servers to implement redundancy. Therefore, you only need to configure multiple proxy connections in the connection pool of the application server.

Advantages:

- Read/write splitting can be implemented without modifying the source program.
- Data sources can be added dynamically without restarting the program.

Disadvant ages:

- The program relies on the middleware, which makes it difficult to switch databases.
- Performance decreases because the middleware serves as a forwarding proxy.

Procedure

Perform the following operations to use MySQL Proxy to split read and write operations:

- 1. Step 1. Preparations
- 2. Step 2. Configure read/write splitting
- 3. Step 3. Grant permissions
- 4. Step 4. Test read/write splitting

Step 1. Preparations

The following section describes the environment:

- Primary database IP address: 121.40.xx.xx
- Secondary database IP address: 101.37.xx.xx
- MySQL Proxy IP address: 116.62.xx.xx

Perform the following operations to prepare for the installation:

- 1. Create three ECS instances and install MySQL.
- 2. Build primary/secondary databases and ensure data consistency between them.
- 3. Modify the MySQL configuration file of the primary/secondary environment.
 - Primary environment:

```
vim /etc/my.cnf
[mysqld]
server-id=202  #Set the unique ID of the server. The default ID is 1.
log-bin=mysql-bin  # Enable binary logs.
```

Secondary environment:

```
[mysqld]
server-id=203
```

4. Restart the MySQL service on the primary/secondary servers.

```
/etc/init.d/mysqld restart
```

5. Create an account on the primary server and grant permissions to the secondary server.

```
mysql -uroot -p95c7586783
grant replication slave on *.* to 'syncms'@'Enter secondary-IP address' identified by '
123456';
flush privileges;
```

6. Check the status of the primary database.

```
mysql> show master status;
```

7. Configure the secondary database.

```
change master to master_host='Enter primary-IP address', master_user='syncms', master_password='123456', master_log_file='mysql-bin.000005', master_log_pos=602;
```

8. Start the secondary synchronization process and check the status.

```
start slave;
show slave status\G
```

9. Verify the synchronization between the primary and secondary databases.

i. Write data to the *testproxy.test1* table in the primary database.

```
mysql> create database testproxy;
mysql> create table testproxy.test1(ID int primary key,name char(10) not null);
mysql> insert into testproxy.test1 values(1,'one');
mysql> insert into testproxy.test1 values(2,'two');
mysql> select * from testproxy.test1;
```

ii. Run the following command in the secondary database to query data in the *testproxy.test1* table:

```
select * from testproxy.test1;
```

If the content in *test proxy.test1* is the same as that in the primary database, data is synchronized between the the primary and secondary databases.

Step 2. Configure read/write splitting

Perform the following operations to configure read/write splitting:

1. Inst all MySQL Proxy.

```
wget https://cdn.mysql.com/archives/mysql-proxy/mysql-proxy-0.8.5-linux-glibc2.3-x86-64
bit.tar.gz
mkdir /alidata
tar xvf mysql-proxy-0.8.5-linux-glibc2.3-x86-64bit.tar.gz
mv mysql-proxy-0.8.5-linux-glibc2.3-x86-64bit/ /alidata/mysql-proxy-0.8.5
```

2. Set environment variables.

```
vim /etc/profile  #Add the following information:
PATH=$PATH:/alidata/mysql-proxy-0.8.5/bin
export $PATH
source /etc/profile  #Validate the environment variables.
mysql-proxy -V
```

3. Set the read/write splitting parameters.

```
cd /alidata/mysql-proxy-0.8.5/share/doc/mysql-proxy/
vim rw-splitting.lua
```

MySQL Proxy will detect client connections. If the number of connections does not exceed the preset value of min_idle_connections, read/write splitting will not be performed. By default, read/write splitting will be triggered when there are at least four connections and at most eight connections. To simplify the test of read/write splitting, set the number of connections to one at least and two at most. For the production environment, you can set the number based on the actual conditions.

Before the modification:

After the modification:

4. Copy the Lua administration script *admin.lua* to the directory where the read/write splitting script *r w-splitting.lua* is located.

```
cp /alidata/mysql-proxy-0.8.5/lib/mysql-proxy/lua/admin.lua /alidata/mysql-proxy-0.8.5/
share/doc/mysql-proxy/
```

Step 3. Grant permissions

Perform the following operations to grant permissions:

1. Grant permissions in the primary database. The permissions will also be granted in the secondary database due to synchronization between the primary and secondary databases.

```
mysql -uroot -p95c7586783
grant all on *.* to 'mysql-proxy'@'Enter <MySQL Proxy IP>' identified by '123456';
flush privileges;
```

2. Start MySQL Proxy.

```
mysql-proxy --daemon --log-level=debug --log-file=/var/log/mysql-proxy.log --plugins=pr oxy -b Enter <primary-IP address>:3306 -r Enter secondary-IP:3306 --proxy-lua-script="/alidata/mysql-proxy-0.8.5/share/doc/mysql-proxy/rw-splitting.lua" --plugins=admin --adm in-username="admin" --admin-password="admin" --admin-lua-script="/alidata/mysql-proxy-0.8.5/share/doc/mysql-proxy/admin.lua"
```

3. Check the port and related processes.

```
netstat -tpln
[root@
                                 ~]# netstat -tpln
Active Internet connections (only servers)
Proto Recv-Q Send-Q Local Address
                                                 Foreign Address
                                                                           State
                                                                                        PID/Program name
                   0 0.0.0.0:22
                                                 0.0.0.0:*
                                                                           LISTEN
                                                                                        826/sshd
tcp
                    0 0.0.0.0:4040
                                                 0.0.0.0:*
                                                                           LISTEN
                                                                                        22767/mysql-proxy
tcp
            0
                    0 0.0.0.0:4041
                                                 0.0.0.0:*
                                                                           LISTEN
                                                                                        22767/mysql-proxy
 tcp
 ps -ef | grep mysql
                                ]# ps -ef | grep mysql
[root@
                                      00:00:00 /alidata/my
                   1 0 10:59 2
         22767
                                                                proxy-0.8.5/libexec/m
root.
                                                                                            l-proxy --daemon -
og-level=debug --log-file=/var/log/mysql-proxy.log --plugins=
6 --proxy-lua-script=/alidata/mysql-proxy-0.8.5/share/doc/mys
                                         l-proxy.log --plugins=proxy -b
                                                                                     :3306 -r
                                                                 1-proxy/rw-splitting.lua --plugins=admin --ad
min-username=admin --admin-password=admin --admin-lua-script=/alidata/mys
                                                                              11-proxy-0.8.5/share/doc/mys
xy/admin.lua
         22794 22602 0 11:02 pts/0 00:00:00 grep --color=auto m
```

Step 4. Test read/write splitting

Perform the following operations to test read/write splitting:

1. Disable secondary replication.

```
stop slave;
```

2. Log on to the backend of MySQL Proxy.

```
mysql -u admin -padmin -P 4041 -h MySQL-Proxy-IP select * from backends; #Check the status.
```

The first connection will be established to the primary database.

```
mysql -umysql-proxy -p123456 -h 116.62.xx.xx -P 4040
insert into testproxy.test1 values(3,'three'); #Add a data record. Secondar
y replication is disabled. Therefore, the record exists in the primary database but doe
s not exist in the secondary database.
```

Create additional test connections. If the data displayed in the *testproxy.test1* table in the primary database is the same as that in the secondary database, read/write splitting is successful.

```
mysql -umysql-proxy -p123456 -h 116.62.xx.xx -P 4040
select * from testproxy.test1;
```

10.3. Change the preferred language of a Windows instance

This topic describes how to download a language pack from Windows Update and set the language as the preferred language of a Windows instance. An ECS instance that runs Windows Server 2016 (English) is used in the examples in this topic.

Prerequisites

Context

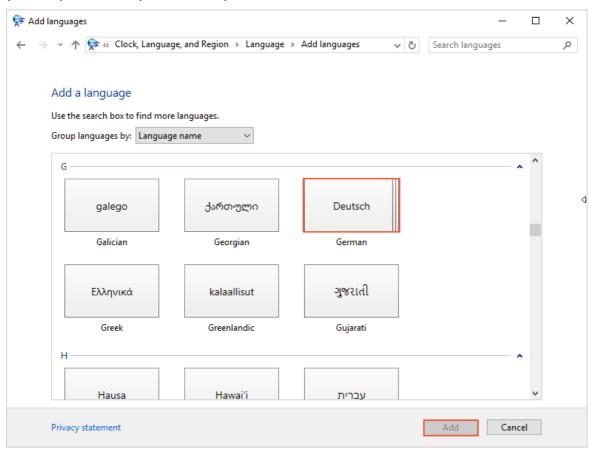
Alibaba Cloud ECS provides only Chinese and English editions of Windows Server public images. If you want to use another language on an ECS instance, such as Arabic, German, Russian, or Japanese, you can download the language pack from Windows Update and set the language as the preferred language of the instance. German is used in the examples in the following procedure. This procedure applies to instances that run Windows Server 2012 or later. After the preferred language of the instance is changed to German, you can use the instance to create a custom image. The custom image also uses German and German keyboard settings. You can then create as many instances as required from this custom image.

Procedure

- 1. Connect to the target Windows instance. For more information, see Overview.
- 2. Open the PowerShell module.
- 3. Run the following commands to temporarily disable Windows Server Update Services (WSUS):

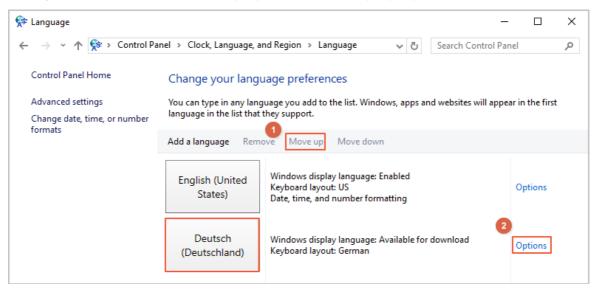
```
Set-ItemProperty -Path 'HKLM:\SOFTWARE\Policies\Microsoft\Windows\WindowsUpdate\AU' -Na me UseWUServer -Value 0
Restart-Service -Name wuauserv
```

- 4. Open the Control Panel, and choose Clock, Language, and Region > Language > Add a language.
- 5. In the Add languages dialog box, select a language and click Add. In this example, Deutsch (German) > Deutsch (Deutschland) is selected.

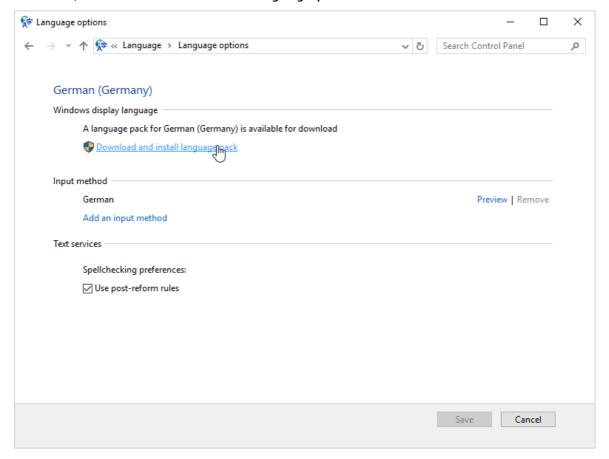


6. Select the language and click **Move up** to change the language priority. In this example, **Deutsch** (**Deutschland**) is selected.

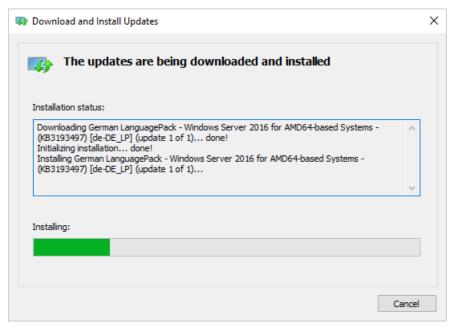
7. Click Options next to the selected language to check for language updates online.



8. Wait for three minutes while the instance checks for updates. If an update is available for download, click **Download and install language pack**.



Wait until the installation is complete.



- 9. Restart the instance by using the ECS console. For more information, see Restart an instance.
- 10. Connect to the Windows instance again. The display language is now Deutsch (German).
- 11. Open the PowerShell ISE module and run the following commands to enable WSUS:

```
Set-ItemProperty -Path 'HKLM:\SOFTWARE\Policies\Microsoft\Windows\WindowsUpdate\AU' -Na me UseWUServer -Value 1
Restart-Service -Name wuauserv
```

12. Open **Windows Update**, check for security updates, and re-install all the security updates that are already installed before the language change.

What's next

Create multiple instances with the same language settings.

- 1. Log on to the ECS console.
- 2. Create a custom image from the Windows instance with the new display language. For more information, see Create a custom image from an instance.
- 3. Create a specified number of instances from the custom image. For more information, see Create an ECS instance by using a custom image.



10.4. Boot a Linux ECS instance into single user mode

This topic describes how to boot an ECS instance that is created from a CentOS, Debian, SUSE Linux Enterprise Server (SLES), or Ubuntu image into single user mode.

Prerequisites

•

• An ECS instance is created. For more information, see <u>Creation method overview</u>. In this example, an ECS instance of the ecs.g6.large instance type is created.

Context

Single user mode is one of the modes in which Linux distributions are booted. GRUB can be used to boot a Linux distribution into single user mode. After the system enters single user mode, you will have system administrator permissions and can modify all system configurations. This mode is usually used in the following scenarios:

- Change the system password
- Troubleshoot boot failures
- Fix system exceptions
- Maint ain partitions of hard disk drives (HDD)

Notice In single user mode, you can modify critical system configurations. We recommend that you set this mode only when necessary and proceed with caution.

References

- Example 1: How a Cent OS ECS instance enters single user mode
- Example 2: How a Debian ECS instance enters single user mode
- Example 3: How a SLES instance enters single user mode
- Example 4: How an Ubuntu ECS instance enters single user mode

Example 1: How a CentOS ECS instance enters single user mode

In this example, an ECS instance running a CentOS 8.0 64-bit operating system is used.

1. Connect to the ECS instance.

For more information, see Connect to a Linux instance by using a password.

(?) Note If you connect an instance by using Workbench or SSH commands, you are not directed to the system boot page when you restart the instance by using a command. Therefore, we recommend that you do not connect to an instance by using Workbench or SSH commands.

2. Run the reboot command to restart the ECS instance. When the interface for selecting a boot system appears, press the *E* key to go to the configuration interface for boot options.

The following figure shows the configuration interface for boot options.

3. Move the pointer to the line that starts with linux by using the arrow keys on the keyboard. Replace the content from ro to the end of the line with rw init=/bin/sh creashkernel=auto

The following figure shows the information after the change is made.

4. Press the Ctrl+X composite key or the F10 key.

The system enters single user mode. The following figure shows the interface for resetting the system password.

```
[ 4.179153] systemd-journald[195]: Received SIGTERM fromsPIDr1t(systemd). kernel.
[ 4.196405] sh: 18 output lines suppressed due to ratelimiting
sh-4.4# passwd
Changing password for user root.
New password:
```

Example 2: How a Debian ECS instance enters single user mode

In this example, an ECS instance running a Debian 10.2 64-bit operating system is used.

1. Connect to the ECS instance.

For more information, see Connect to a Linux instance by using a password.

Note If you connect an instance by using Workbench or SSH commands, you are not directed to the system boot page when you restart the instance by using a command. Therefore, we recommend that you do not connect to an instance by using Workbench or SSH commands.

2. Run the reboot command to restart the ECS instance. When the configuration interface for kernel options appears, press the *E* key to go to the GRUB interface.

The following figure shows the GRUB interface.

```
GNU GRUB version 2.02+dfsg1-20
<u>s</u>etparams 'Debian GNU/Linux'
        load_video
        insmod gzio
        if [ x$grub_platform = xxen ]; then insmod xzio; insmod lzopio; fi
        insmod part_msdos
        insmod ext2
        if [ x$feature_platform_search_hint = xy ]; then
         search --no-floppy --fs-uuid --set=root fb96ed16-
         search --no-floppy --fs-uuid --set=root fb96ed16-7
                     'Loading Linux 4.19.0-6-amd64 ...'
                     /boot/vmlinuz-4.19.0-6-amd64 root=UUID= 934-20b46983
bc80 ro  vga=792 console=tty0 console=ttyS0,115200n8 net.ifnames=0 noibrs quiet
                     'Loading initial ramdisk ...'
/boot/initrd.img-
        initrd
   Minimum Emacs-like screen editing is supported. TAB lists completions. Press Ctrl-x or F10 to boot, Ctrl-c or F2 for a command-line or ESC to discard edits and return
```

3. Move the pointer to the line that starts with linux by using the arrow keys on the keyboard. Append single to the end of the line.

The following figure shows the information after the change is made.

4. Press the *Ctrl+X* composite key or the *F10* key to start the system. Enter the password of the root user.

The system enters single user mode.

```
Give root password for maintenance
(or press Control–D to continue):
root@db1:~# _
```

Example 3: How a SLES instance enters single user mode

In this example, an ECS instance running a SLES 15 SP1 64-bit operating system is used.

1. Connect to the ECS instance.

For more information, see Connect to a Linux instance by using a password.

- Note If you connect an instance by using Workbench or SSH commands, you are not directed to the system boot page when you restart the instance by using a command. Therefore, we recommend that you do not connect to an instance by using Workbench or SSH commands.
- 2. Run the reboot command to restart the ECS instance. When the configuration interface for kernel options appears, press the *E* key to go to the GRUB interface.

The following figure shows the GRUB interface.

```
GNU GRUB version 2.02
setparams 'SLES 15-SP1'
        load_video
       set gfxpayload=keep
        insmod gzio
        insmod part_msdos
insmod ext2
        set root='hd0,msdos1'
        if [ x$feature_platform_search_hint = xy ]; then
         search --no-floppy --fs-uuid --set=root --hint='hd0,msdos1'
2b38e357cd3c
                                                                        c4e5\
92f a-
       else
          search --no-floppy --fs-uuid --set=root c4e59
   Minimum Emacs-like screen editing is supported. TAB lists
   completions. Press Ctrl-x or F10 to boot, Ctrl-c or F2 for
   a command-line or ESC to discard edits and return to the GRUB menu
```

3. Move the pointer to the line that starts with $_{linux}$ by using the arrow keys on the keyboard. Append $_{single}$ to the end of the line.

The following figure shows the information after the change is made.

4. Press the *Ctrl+X* composite key or the *F10* key to start the system. Enter the password of the root user.

The system enters single user mode.

```
Booting a command list

Loading Linux 4.12.14-197.29-default ...

Loading initial ramdisk ...

Give root password for maintenance

(or press Control-D to continue):

sles01:~ # _
```

Example 4: How an Ubuntu ECS instance enters single user mode

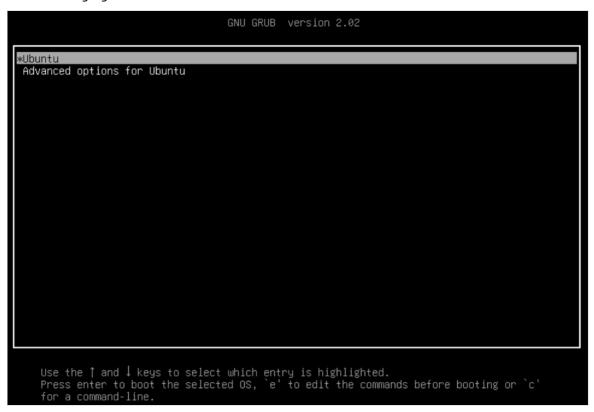
In this example, an ECS instance running an Ubuntu 18.04 64-bit operating system is used.

1. Connect to the ECS instance.

For more information, see Connect to a Linux instance by using a password.

- Note If you connect an instance by using Workbench or SSH commands, you are not directed to the system boot page when you restart the instance by using a command. Therefore, we recommend that you do not connect to an instance by using Workbench or SSH commands.
- 2. Run the reboot command to restart the ECS instance. During the restarting process, press the *Sh* if t key to go to the GRUB interface.

The following figure shows the GRUB interface.



- 3. Select the Advanced options for Ubuntu in the second line of the GRUB interface and press the *Ent er* key.
- 4. Select the recovery mode in the second line on the interface that appears and press the *E* key to edit the boot options.

```
Ubuntu, with Linux 4.15.0-88-generic

*Ubuntu, with Linux 4.15.0-88-generic (recovery mode)

Ubuntu, with Linux 4.15.0-55-generic

Ubuntu, with Linux 4.15.0-55-generic (recovery mode)
```

5. On the editing page, move the pointer to the line that starts with linux by using the arrow keys on the keyboard. Replace the content from linux to the end of the line with linux single linux in linux by using the arrow keys on the keyboard. Replace the content from linux to the end of the line with linux single linux in linux by using the arrow keys on the keyboard. Replace the content from linux to the end of the line with linux single linux in linux single linux

The following figure shows the information after the change is made.

6. Press the *Ctrl+X* composite key or the *F10* key.

The system enters single user mode. The following figure shows the interface for resetting the system password.

```
root@(none):/# passwd
Enter new UNIX password:
Retype new UNIX password:
passwd: password updated successfully
```

11.Block Storage

11.1. Resize partitions and file systems of Linux system disks

This topic describes how to use the growpart or xfsprogs tool to resize the partitions and file systems of Linux system disks.

Prerequisites

Before you resize the partitions and file systems of a system disk, make sure that the following requirements are met:

1. A snapshot is created to back up data.

To prevent data loss caused by accidental changes, we recommend that you create snapshots to back up your data. For more information, see Create a snapshot of a disk.

2. A data disk is resized in the ECS console.

If no data disks are resized, perform the operations described in Step 2: Resize the disk in the ECS console to resize a data disk.

3.

- 4. The growpart or xf sprogs tool is installed based on your operating system.
 - o Alibaba Cloud Linux 2 and Cent OS 7

```
Run the yum install <Package name> command. Example:

yum install cloud-utils-growpart xfsprogs -y
```

o Ubuntu 14, Ubuntu 16, Ubuntu 18, and Debian 9

```
Run the apt install <Package name> command. Example:

apt install cloud-guest-utils xfsprogs -y
```

Debian 8, openSUSE 42.3, openSUSE 13.1, and SUSE Linux Enterprise Server 12 SP2
 Use the upstream version of growpart or xfsprogs.

? Note If the partitions or file systems cannot be resized due to issues with growpart or xfsprogs, we recommend that you re-install the tool.

- 5. The uname -a command is run to check the kernel version of the instance operating system.
 - If the kernel version of your instance operating system is 3.6.0 or later, see the Procedure for instances with kernels 3.6.0 or later section of this topic.
 - If the kernel version of your instance operating system such as CentOS 6, Debian 7, or SUSE Linux Enterprise Server 11 SP4 is earlier than 3.6.0, you must restart the instance by using the Elastic Compute Service (ECS) console or by calling an API operation to resize the partitions or file systems. For more information, see the Procedure for instances with kernels earlier than 3.6.0 section of this topic.

Context

The procedures in this topic apply to disks that have the following partition formats and file system types:

- Partition formats: master boot record (MBR) and GUID Partition Table (GPT)
- File system types: ext, xfs, and btrfs

Procedure for instances with kernels 3.6.0 or later

In this example, Alibaba Cloud Linux 2.1903 LTS 64-bit is used to describe how to resize the partitions and file systems of the system disk.

Note The commands in the example also apply to CentOS 7.

1. Run the following command to check the size of the system disk:

```
fdisk -1
```

In the following example command output, the size of the /dev/vda disk is 100 GiB:

2. Run the following command to check the sizes of the partitions and the types of the file systems:

```
df -Th
```

In the following example command output, the size of the /dev/vda1 partition is 40 GiB, and the file system type is ext4:

3. Run the following command to resize a partition:

```
growpart <DeviceName> <PartionNumber>
```

In this command, *<DeviceName>* is the device name of the system disk, and *<PartionNumber>* is the partition number. You must use a space to separate the device name from the partition number.

In the following example command output, the first partition of the system disk is resized:

[root@ecshost ~]# growpart /dev/vda 1
CHANGED: partition=1 start=2048 old: size=83883999 end=83886047 new: size=209713119 end
=209715167

? Note

- If a single disk has successive partitions, you need only to resize the last partition to resize the disk. For example, the /dev/vda system disk has three partitions named /dev/vda1 , /dev/vda2 , and /dev/vda3 , you can run the growpart /dev/vda 3 command to resize the /dev/vda system disk.
- o If the unexpected output in sfdisk --version [sfdisk, from util-linux 2.23.2] error appears when you run the growpart /dev/vda 1 command, change the character encoding type to resolve the issue. For more information, see FAQ.
- 4. Resize the file system.

Run the df -Th command to check the file system type. Then, run one of the following commands to resize the file system based on the file system type.

• If the file system is of an ext type such as ext3 and ext4, run the following command to resize the file system:

```
resize2fs <PartitionName>
```

In the following example command output, the file system of the /dev/vda1 partition is resized:

```
[root@ecshost ~]# resize2fs /dev/vda1
resize2fs 1.42.9 (28-Dec-2013)
Filesystem at /dev/vda1 is mounted on /; on-line resizing required
old_desc_blocks = 3, new_desc_blocks = 7
The filesystem on /dev/vda1 is now 26214139 blocks long.
```

o If the file system type is xfs, run the following command to resize the file system:

```
xfs_growfs <mountpoint>
```

In the following example command output, the file system of the /dev/vda1 partition is resized. The mount point of the /dev/vda1 partition is the / root directory.

```
[root@ecshost ~]# xfs growfs /
meta-data=/dev/vda1
                            isize=512 agcount=13, agsize=1310656 blks
                            sectsz=512 attr=2, projid32bit=1
                            crc=1 finobt=1, sparse=1, rmapbt=0
                            reflink=1
                            bsize=4096 blocks=15728379, imaxpct=25
data
                           sunit=0 swidth=0 blks
naming =version 2
                           bsize=4096 ascii-ci=0, ftype=1
log =internal log
                           bsize=4096 blocks=2560, version=2
                           sectsz=512 sunit=0 blks, lazy-count=1
realtime =none
                            extsz=4096 blocks=0, rtextents=0
data blocks changed from 15728379 to 20971259
```

Note The command to run may vary based on the version of xfs_growfs. You can run the xfs_growfs --help command to check the corresponding command to run.

o If the file system type is btrfs, run the following command to resize the file system:

```
btrfs filesystem resize max <mountpoint>
```

In the following example command output, the file system of the /dev/vda1 partition is resized. The mount point of the /dev/vda1 partition is the / root directory.

```
[root@ecshost ~]# btrfs filesystem resize max /
```

5. Run the following command to view the results of disk resizing:

```
df -h
```

In the following example command output, the size of the /dev/vda1 partition is 100 GiB, which indicates that the partition is resized.

```
[root@ecshost ~] # df -h

Filesystem Size Used Avail Use% Mounted on
devtmpfs 869M 0 869M 0% /dev

tmpfs 879M 0 879M 0% /dev/shm

tmpfs 879M 492K 878M 1% /run

tmpfs 879M 0 879M 0% /sys/fs/cgroup
/dev/vdal 99G 1.8G 93G 2% /

tmpfs 176M 0 176M 0% /run/user/0
```

Procedure for instances with kernels earlier than 3.6.0

In this example, CentOS 6 is used to describe how to resize the partitions and file systems of the system disk.

1. Change the source address of a YUM repository on CentOS 6.

CentOS 6 has reached its end of life (EOL). To install the software package on CentOS 6 by using YUM, you must change the source address of the YUM repository. For more information, see Change the CentOS 6 source address.

2. Run the following command to install dracut-modules-growroot:

```
yum install -y dracut-modules-growroot
```

? Note If a package manager other than YUM is used, change **yum** in the preceding command based on the package manager.

3. Run the following command to overwrite the existing initramfs file:

```
dracut -f
```

- 4. Check the size of the system disk and the partitions and file systems.
 - Run the following command to check the size of the system disk:

```
fdisk -l
```

In the following example command output, the size of the /dev/vda disk is 100 GiB:

• Run the following command to check the partitions of the system disk and the types of the file systems:

```
df -Th
```

In the following example command output, the size of the /dev/vda1 partition is 20 GiB, and the file system type is ext4:

```
[root@ecshost ~]# df -Th
Filesystem Type Size Used Avail Use% Mounted on
/dev/vda1 ext4 20G 1.1G 18G 6% /
tmpfs tmpfs 7.8G 0 7.8G 0% /dev/shm
```

5. Run the following command to resize a partition:

```
growpart <DeviceName> <PartionNumber>
```

In this command, *<DeviceName>* is the device name of the system disk, and *<PartionNumber>* is the partition number. You must use a space to separate the device name from the partition number.

In the following example command output, the first partition of the system disk is resized:

```
[root@ecshost ~]# growpart /dev/vda 1
CHANGED: partition=1 start=2048 old: size=41940992 end=41943040 new: size=209710462,end
=209712510
```

6. Restart the instance in the ECS console.

Notice The resize operation does not take effect until you restart the instance by using the ECS console or by calling the RebootInstance operation. For more information, see Restart an instance and RebootInstance.

- 7. Connect to the instance.
- 8. Resize the file system.

Run the df -Th command to check the file system type. Then, run one of the following commands to resize the file system based on the file system type.

• If the file system is of an ext type such as ext3 and ext4, run the following command to resize the file system:

```
resize2fs <PartitionName>
```

In the following example command output, the file system of the /dev/vda1 partition is resized:

```
[root@ecshost ~]# resize2fs /dev/vda1
resize2fs 1.41.12 (17-May-2010)
Filesystem at /dev/vda1 is mounted on /; on-line resizing required
old desc_blocks = 2, new_desc_blocks = 7
Performing an on-line resize of /dev/vda1 to 26213807 (4k) blocks.
The filesystem on /dev/vda1 is now 26213807 blocks long.
```

o If the file system type is xfs, run the following command to resize the file system:

```
xfs_growfs <mountpoint>
```

In the following example command output, the file system of the /dev/vda1 partition is resized. The mount point of the /dev/vda1 partition is the / root directory.

```
[root@ecshost ~]# xfs_growfs /
```

9. Run the following command to check the sizes of the disk partitions:

```
df -h
```

In the following example command output, the size of the /dev/vda1 partition is 100 GiB, which indicates that the partition is resized.

Related information

- Resize disks online for Linux instances
- Resize disks offline for Linux instances
- Resize partitions and file systems of Linux data disks

11.2. Resize partitions and file systems of Linux data disks

When you resize a disk of an Elastic Compute Service (ECS) instance, only the storage capacity of the disk is extended. The file systems of the instance are not resized. You can perform the operations in this topic to resize the file systems of the instance and then extend its storage capacity.

Prerequisites

1. A snapshot is created to back up data.

To prevent data loss caused by accidental changes, we recommend that you create snapshots to back up your data. For more information, see Create a snapshot of a disk.

2. A data disk is resized in the ECS console.

If no data disks are resized, perform the operations described in Step 2: Resize the disk in the ECS console to resize a data disk.

3.

Context

In the examples of this topic, the following configurations are used:

- Operating system of the instance: Alibaba Cloud Linux 2.1903 LTS 64-bit public image
- Category of the data disk: ultra disk
- Device name of the data disk: /dev/vdb

Adjust the commands or parameter settings based on the actual operating system and device name of your data disk.

Check the partition format and the file system type

1. Run the following command to check the partition format of the data disk:

```
fdisk -lu /dev/vdb
```

In this example, the disk has a partition named /dev/vdb1.

- \circ If the disk uses the master boot record (MBR) partition format, the value of System is Linux .
- If the disk uses the GUID Partition Table (GPT) partition format, the value of system is GPT.

```
[root@ecshost ~]# fdisk -lu /dev/vdb
Disk /dev/vdb: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x9277b47b
Device Boot Start End Blocks Id System
/dev/vdb1 2048 41943039 20970496 83 Linux
```

2. Run the following command to check the file system type of the partition:

```
blkid /dev/vdb1
```

In this example, the file system type of /dev/vdb1 is ext4.

```
[root@ecshost ~]# blkid /dev/vdb1
/dev/vdb1: UUID="e97bf1e2-fc84-4c11-9652-73******24" TYPE="ext4"
```

Note If the data disk does not have partitions or file systems, or if the data disk has partitions but no file systems, no results are returned.

- 3. Check the status of the file system.
 - o Run the following command to check the status of an ext file system:

```
e2fsck -n /dev/vdb1
```

• Run the following command to check the status of an xfs file system:

```
xfs_repair -n /dev/vdb1
```

• Run the following command to check the status of a btrfs file system:

```
btrfsck /dev/vdb1
```

Different command outputs are returned for different types of file systems.

• The following code shows an example command output for an ext or xfs file system. If the file system works normally, clean is returned. If clean is not returned, troubleshoot the issue.

```
[root@ecshost ~]# e2fsck -n /dev/vdb1
Warning! /dev/vdb1 is mounted.
Warning: skipping journal recovery because doing a read-only filesystem check.
/dev/vdb1: clean, 11/1310720 files, 126322/5242624 blocks
```

o The following code shows an example command output for a btrfs file system. found 114688 bytes used err is 0 indicates that the file system works normally. If an error is returned in the command output, troubleshoot the issue.

```
[root@ecshost ~] # btrfsck /dev/vdb1
Checking filesystem on /dev/vdb1
UUID: 1234b7a7-68ff-4f48-a88c-8943f27f1234
checking extents
checking free space cache
checking fs roots
checking csums
checking root refs
found 114688 bytes used err is 0
total csum bytes: 0
total tree bytes: 114688
total fs tree bytes: 32768
total extent tree bytes: 16384
btree space waste bytes: 109471
file data blocks allocated: 0
referenced 0
```

Choose a method to resize partitions or file systems

Choose a resize method based on the partition format and the file system type.

	Resize method	Scenario
--	---------------	----------

Scenario	Resize method
The data disk has partitions and file systems	 To resize the existing MBR partitions of the data disk, see Option 1: Resize existing MBR partitions. To resize the disk to create more MBR partitions, see Option 2: Create and format MBR partitions. To resize the existing GPT partitions of the data disk, see Option 3: Resize existing GPT partitions. To resize the disk to create more GPT partitions, see Option 4: Create and format GPT partitions.
The new data disk does not have partitions or file systems	After you resize the data disk in the ECS console, perform the operations described in Partition and format a data disk on a Linux instance or Partition and format a data disk larger than 2 TiB in size.
The raw data disk has a file system but no partitions	After you resize the data disk in the ECS console, perform the operations described in Option 5: Resize the file system of a raw data disk.
The data disk is not attached to an instance	After you attach the data disk to an instance, perform the following operations in this topic to resize the data disk.

? Note

- If a data disk contains an MBR partition, the data disk cannot be resized to 2 TiB or larger. To prevent data loss, we recommend that you create a disk larger than 2 TiB in size, format it to a GPT partition, and then copy the data in the MBR partition to the GPT partition. For more information, see Partition and format a data disk larger than 2 TiB in size.
- If data disks cannot be resized due to issues with the resize or format tool, you can upgrade the tool to a later version or re-install the tool.

Option 1: Resize existing MBR partitions

(?) Note To prevent data loss, we recommend that you do not resize partitions and file systems that are mounted to ECS instances. To resize a partition that is mounted to an ECS instance, run the umount command to unmount the partition and then resize it. When the partition resumes normal working, run the mount command to mount it again. Perform operations based on the Linux kernel version:

- If the instance kernel version is earlier than 3.6, unmount the partition, modify the partition table, and then resize the file system.
- If the instance kernel version is 3.6 or later, modify the partition table, notify the kernel to update the partition table, and then resize the file system.

Perform the following operations to resize an existing MBR partition:

1. Modify the partition table.

i. Run the following command to view the partition information and record the start and end sectors of the partition:

```
fdisk -lu /dev/vdb
```

In this example, the start sector number of the /dev/vdb1 partition is 2048 and the end sector number is 41943039.

```
[root@ecshost ~]# fdisk -lu /dev/vdb
Disk /dev/vdb: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x9277b47b
Device Boot Start End Blocks Id System
/dev/vdb1 2048 41943039 20970496 83 Linux
```

ii. View the mount path of the data disk. Unmount the partition based on the returned file path and wait until the partition is unmounted.

Run the following command to view the mount information:

```
mount | grep "/dev/vdb"
```

Run the following command to unmount the /dev/vdb1 partition from the data disk:

```
umount /dev/vdb1
```

Run the following command to view the operation result:

```
mount | grep "/dev/vdb"
```

A command output similar to the following one is returned:

```
[root@ecshost ~]# mount | grep "/dev/vdb"
/dev/vdb1 on /mnt type ext4 (rw,relatime,data=ordered)
[root@ecshost ~]# umount /dev/vdb1
[root@ecshost ~]# mount | grep "/dev/vdb"
```

iii. Run the fdisk command to delete the existing partition.

• Warning If errors occur when you delete a partition, data stored on the partition may be deleted. To prevent data loss, back up important data such as user data in a database before you delete a partition.

- a. Run the fdisk -u /dev/vdb command to partition the data disk.
- b. Enter p to obtain the partition table.
- c. Enter d to delete the partition.
- d. Enter p to verify whether the partition is deleted.
- e. Enter w to save changes and exit.

The following sample code shows how to delete a partition:

```
[root@ecshost ~]# fdisk -u /dev/vdb
Welcome to fdisk (util-linux 2.23.2).
Changes will remain in memory only, until you decide to write them.
Be careful before using the write command.
Command (m for help): p
Disk /dev/vdb: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x9277b47b
Device Boot Start End Blocks Id System
/dev/vdb1 2048 41943039 20970496 83 Linux
Command (m for help): d
Selected partition 1
Partition 1 is deleted
Command (m for help): p
Disk /dev/vdb: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x9277b47b
Device Boot Start End Blocks Id System
Command (m for help): w
The partition table has been altered!
Calling ioctl() to re-read partition table.
Syncing disks.
```

- iv. Run the fdisk command to create a partition.
 - a. Runthe fdisk -u /dev/vdb command to partition the data disk.
 - b. Enter p to obtain the partition table.
 - c. Enter *n* to create a partition.
 - d. Enter p to select the primary partition type.
 - e. Enter <partition number> to select a partition number. In this example, 1 is selected.
 - f. Set the start and end sector numbers for the new partition.

☐ Warning The start sector number of the new partition must be the same as that of the existing partition, and the end sector number must be greater than that of the existing partition. Otherwise, the resize operation fails.

g. Enter w to save changes and exit.

The following sample code shows how to resize a partition. In this example, the /dev/vdb1 partition is resized from 20 GiB to 40 GiB.

```
[root@ecshost ~] # fdisk -u /dev/vdb
Welcome to fdisk (util-linux 2.23.2).
Changes will remain in memory only, until you decide to write them.
Be careful before using the write command.
Command (m for help): p
Disk /dev/vdb: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x9277b47b
Device Boot Start End Blocks Id System
Command (m for help): n
Partition type:
p primary (0 primary, 0 extended, 4 free)
e extended
Select (default p): p
Partition number (1-4, default 1): 1
First sector (2048-83886079, default 2048):
Using default value 2048
Last sector, +sectors or +size{K,M,G} (2048-83886079, default 83886079):
Partition 1 of type Linux and of size 40 GiB is set
Command (m for help): w
The partition table has been altered!
Calling ioctl() to re-read partition table.
Syncing disks.
```

v. Run one of the following commands to notify the kernel to update the partition table:

partprobe /dev/vdb

Note If the -bash: partprobe: command not found error message appears in the CentOS 6 operating system, change the source address of a YUM repository. For more information, see Change the CentOS 6 source address. Run the yum install -y parted command to install the Parted tool. Then, run the partprobe /dev/vdb command again.

partx -u /dev/vdb1

vi. Run the following command to check whether the partition table is added:

lsblk /dev/vdb

vii. Run the following command to check the file system and verify whether the file system is in the clean state:

e2fsck -f /dev/vdb1

Note If the file system is not in the clean state after you run the preceding command, you can run the e2fsck -n /dev/vdb1 command to check the file system.

- 2. Resize the file system.
 - If the file system type is ext such as ext3 and ext4, run the following commands in sequence to resize the file system and remount the partition.

Run the following command to resize the file system:

resize2fs /dev/vdb1

Run the following command to mount the partition to /mnt:

mount /dev/vdb1 /mnt

• If the file system type is xfs, run the following commands in sequence to remount the partition and then resize the file system.

Run the following command to mount the partition to /mnt:

mount /dev/vdb1 /mnt

Run the following command to resize the file system:

xfs growfs /mnt

- Note The new version xfs_growfs identifies the device to resize based on the mount point. Example: xfs_growfs /mnt . You can run the xfs_growfs --help command to check how to use xfs_growfs of different versions.
- If the file system type is btrfs, run the following commands in sequence to remount the partition and then resize the file system.

Run the following command to mount the partition to /mnt:

```
mount /dev/vdb1 /mnt
```

Run the following command to resize the file system:

```
btrfs filesystem resize max /mnt
```

Option 2: Create and format MBR partitions

Perform the following operations to create more MBR partitions:

1. Run the following command to create a partition:

```
fdisk -u /dev/vdb
```

The following sample code shows how to create a partition. In this example, a 20 GiB /dev/vdb2 partition is created.

```
[root@ecshost ~] # fdisk -u /dev/vdb
Welcome to fdisk (util-linux 2.23.2).
Changes will remain in memory only, until you decide to write them.
Be careful before using the write commad.
Command (m for help): p
Disk /dev/vdb: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x2b31a2a3
Device Boot Start End /dev/vdb1 2048 41943039
                                        Blocks Id System
                   2048 41943039 20970496 83 Linux
Command (m for help): n
Partition type:
  p primary (1 primary, 0 extended, 3 free)
     extended
  е
Select (default p): p
Partition number (2-4, default 2): 2
First sector (41943040-83886079, default 41943040):
Using default value 41943040
Last sector, +sectors or +size(K,M,G) (41943040-83886079, default 83886079):
Using default value 83886079
Partition 2 of type Linux and of size 20 GiB is set
Command (m for help): w
The partition table has been altered!
Calling ioctl() to re-read partition table.
Syncing disks.
```

2. Run the following command to view the partition:

```
lsblk /dev/vdb
```

A command output similar to the following one is returned:

- 3. Create a file system.
 - Run the following command to create an ext4 file system:

```
mkfs.ext4 /dev/vdb2
```

• Run the following command to create an xfs file system:

```
mkfs.xfs -f /dev/vdb2
```

• Run the following command to create a btrfs file system:

```
mkfs.btrfs /dev/vdb2
```

4. Run the following command to view the information about the file system:

```
blkid /dev/vdb2
```

A command output similar to the following one is returned:

```
[root@ecshost ~]# blkid /dev/vdb2
/dev/vdb2: UUID="e3f336dc-d534-4fdd-***-b6ff1a55bdbb" TYPE="ext4"
```

5. Run the following command to mount the partition:

```
mount /dev/vdb2 /mnt
```

6. Run the following command to view the space and usage of the data disk:

```
df -h
```

If the partition is mounted, the information about the new file system is displayed, as shown in the following example command output.

```
[root@ecshost ~] # df -h
Filesystem Size Used Avail Use% Mounted on
/dev/vda1 40G 1.6G 36G 5% /
devtmpfs 3.9G 0 3.9G 0% /dev
tmpfs 3.9G 0 3.9G 0% /dev/shm
tmpfs 3.9G 460K 3.9G 1% /run
tmpfs 3.9G 0 3.9G 0% /sys/fs/cgroup
/dev/vdb2 9.8G 37M 9.2G 1% /mnt
tmpfs 783M 0 783M 0% /run/user/0
```

Option 3: Resize existing GPT partitions

Perform the following operations to resize an existing GPT partition:

1. View the mount path of the data disk. Unmount the partition based on the returned file path and wait until the partition is unmounted.

Run the following command to view the mount information:

```
mount | grep "/dev/vdb"
```

Run the following command to unmount the /dev/vdb1 partition from the data disk:

```
umount /dev/vdb1
```

Run the following command to view the operation result:

```
mount | grep "/dev/vdb"
```

A command output similar to the following one is returned:

```
[root@ecshost ~] # mount | grep "/dev/vdb"
/dev/vdb1 on /mnt type ext4 (rw,relatime,data=ordered)
[root@ecshost ~] # umount /dev/vdb1
[root@ecshost ~] # mount | grep "/dev/vdb"
```

- 2. Use Parted to allocate capacity for the existing GPT partition.
 - i. Run the following command to start Parted:

```
parted /dev/vdb
```

To view the instructions on using Parted, run the help command.

ii. Run the following command to view the partition information and record the partition number and start sector number of the existing partition:

```
print
```

If Fix/Ignore/Cancel? and Fix/Ignore? appear, enter Fix.

In this example, the size of the existing partition is 1 TiB, the partition number (the value of $_{
m Nu}$ mber) is 1, and the start sector number (the value of $_{
m Start}$) is $_{
m 1049kB}$.

```
(parted) print
Model: Virtio Block Device (virtblk)
Disk /dev/vdb: 3299GB
Sector size (logical/physical): 512B/512B
Partition Table: gpt
Disk Flags:
Number Start End Size File system Name Flags
1 1049kB 1100GB 1100GB ext4 primary
```

iii. Run the following command to delete the existing partition:

```
rm <Partition number>
```

In this example, the partition number of the existing partition is 1. Run the following command to delete the existing partition:

```
rm 1
```

iv. Run the following command to recreate the primary partition:

```
mkpart primary <Start sector number of the existing partition> <Percentage of the a llocated capacity>
```

In this example, the start sector number of the existing partition is $_{1049kB}$, and the 3 TiB total capacity is allocated to the partition. Run the following command to create the primary partition:

```
mkpart primary 1049kB 100%
```

v. Run the following command to check whether the primary partition is created:

```
print
```

A command output similar to the following one is returned. In the command output, the partition number of the new GPT partition is still 1 but the capacity of the partition is increased to 3 TiB.

```
(parted) rm 1
(parted) mkpart primary 1049kB 100%
(parted) print
Model: Virtio Block Device (virtblk)
Disk /dev/vdb: 3299GB
Sector size (logical/physical): 512B/512B
Partition Table: gpt
Disk Flags:
                       Size
       Start
                End
                                File system Name
                                                       Flags
       1049kB
                3299GB
                       3299GB
                                ext4
                                              primary
```

vi. Run the following command to exit Parted:

```
quit
```

The following code provides an example on how to perform the preceding steps.

```
[root@ecshost ~] # parted /dev/vdb
GNU Parted 3.1
Using /dev/vdb
Welcome to GNU Parted! Type 'help' to view a list of commands.
(parted) print
Error: The backup GPT table is not at the end of the disk, as it should be.
This might mean that another operating system believes the disk is smaller.
Fix, by moving the backup to the end (and removing the old backup)?
Fix/Ignore/Cancel? Fix
Warning: Not all of the space available to /dev/vdb appears to be used, you can
fix the GPT to use all of the space (an extra 4294967296 blocks) or continue
with the current setting?
Fix/Ignore? Fix
Model: Virtio Block Device (virtblk)
Disk /dev/vdb: 3299GB
Sector size (logical/physical): 512B/512B
Partition Table: gpt
Disk Flags:
Number Start End Size File system Name
                                                  Flags
1 1049kB 1100GB 1100GB ext4 primary
(parted) rm 1
(parted) mkpart primary 1049kB 100%
(parted) print
Model: Virtio Block Device (virtblk)
Disk /dev/vdb: 3299GB
Sector size (logical/physical): 512B/512B
Partition Table: gpt
Disk Flags:
Number Start End Size File system Name
1 1049kB 3299GB 3299GB ext4 primary
(parted) quit
Information: You may need to update /etc/fstab.
```

3. Run the following command to check the file system for consistency:

```
fsck -f /dev/vdb1
```

A command output similar to the following one is returned:

```
[root@ecshost ~]# fsck -f /dev/vdb1
fsck from util-linux 2.23.2
e2fsck 1.43.5 (04-Aug-2017)
Pass 1: Checking inodes, blocks, and sizes
Pass 2: Checking directory structure
Pass 3: Checking directory connectivity
Pass 4: Checking reference counts
Pass 5: Checking group summary information
/dev/vdb1: 11/67108864 files (0.0% non-contiguous), 4265369/268434944 blocks
```

- 4. Resize the file system corresponding to the partition and remount the partition.
 - o ext file system such as ext3 and ext4

Run the following command to resize the file system of the new partition:

resize2fs /dev/vdb1

Run the following command to remount the partition:

mount /dev/vdb1 /mnt

o xfs file system

Run the following command to remount the partition:

mount /dev/vdb1 /mnt

Run the following command to resize the file system:

xfs_growfs /mnt

Note The new version xfs_growfs identifies the device to resize based on the mount point. Example: xfs_growfs /mnt . You can run the xfs_growfs --help command to check how to use xfs_growfs of different versions.

o btrfs file system

Run the following command to remount the partition:

mount /dev/vdb1 /mnt

Run the following command to resize the file system:

btrfs filesystem resize max /mnt

Option 4: Create and format GPT partitions

Perform the following operations to create more GPT partitions: In this example, a 32 TiB data disk is used. The disk has a 4.8 TiB /dev/vdb1 partition, and a new /dev/vdb2 partition is to be created.

1. Run the following command to view the information of existing partitions in the data disk:

fdisk -l

A command output similar to the following one is returned:

```
[root@ecshost ~] # fdisk -l
Disk /dev/vda: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x000b1b45
Device Boot Start End Blocks Id System  
/dev/vda1 * 2048 83875364 41936658+ 83 Linux
                                        Blocks Id System
WARNING: fdisk GPT support is currently new, and therefore in an experimental phase. Us
e at your own discretion.
Disk /dev/vdb: 35184.4 GB, 35184372088832 bytes, 68719476736 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: gpt
Disk identifier: BCE92401-F427-45CC-8B0D-B30EDF279C2F
     Start End Size Type
       2048 10307921919 4.8T Microsoft basic mnt
```

- 2. Use Parted to create a partition and allocate capacity for it.
 - i. Run the following command to start Parted:

```
parted /dev/vdb
```

ii. Run the following command to view the disk capacity to be allocated. Record the start and end sectors and the capacity of the existing partition.

```
print free
```

In this example, the start sector number of /dev/vdb1 is 1049KB, the end sector number is 5278GB, and the capacity is 5,278 GiB.

iii. Run the following command to set the start sector and capacity for the new partition:

```
\label{eq:mkpart} \mbox{\tt Percentage of the allocated capacity} > \\
```

In this example, the /dev/vdb2 partition named test is created. The start sector number of the new partition is the end sector number of the existing partition. The new capacity is allocated to the new partition.

```
mkpart test 5278GB 100%
```

iv. Run the following command to check whether the capacity (Size) of the partition is changed:

```
print
```

A command output similar to the following one is returned:

```
(parted) print
Model: Virtio Block Device (virtblk)
Disk /dev/vdb: 35.2TB
Sector size (logical/physical): 512B/512B
Partition Table: gpt
Disk Flags:
Number Start End Size File system Name Flags
1 1049kB 5278GB 5278GB ext4 mnt
2 5278GB 35.2TB 29.9TB test
```

v. Run the following command to exit Parted:

```
quit
```

- 3. Create a file system for the new partition.
 - Run the following command to create an ext4 file system:

```
mkfs.ext4 /dev/vdb2
```

o Run the following command to create an ext3 file system:

```
mkfs.ext3 /dev/vdb2
```

• Run the following command to create an xfs file system:

```
mkfs.xfs -f /dev/vdb2
```

• Run the following command to create a btrfs file system:

```
mkfs.btrfs /dev/vdb2
```

A command output similar to the following one is returned. In this example, an xfs file system is created.

4. Run the following command to view the changes to the partition capacity:

```
fdisk -1
```

A command output similar to the following one is returned:

```
[root@ecshost ~] # fdisk -l
Disk /dev/vda: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x000b1b45
Device Boot Start End Blocks Id System  
/dev/vda1 * 2048 83875364 41936658+ 83 Linux
                                        Blocks Id System
WARNING: fdisk GPT support is currently new, and therefore in an experimental phase. Us
e at your own discretion.
Disk /dev/vdb: 35184.4 GB, 35184372088832 bytes, 68719476736 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: gpt
Disk identifier: BCE92401-F427-45CC-8B0D-B30EDF279C2F
        Start End Size Type
         2048 10307921919 4.8T Microsoft basic mnt
 2 10307921920 68719474687 27.2T Microsoft basic test
```

5. Run the following command to view the types of the file systems:

```
blkid
```

A command output similar to the following one is returned:

```
[root@ecshost ~] # blkid
/dev/vda1: UUID="ed95c595-4813-480e-***-85b1347842e8" TYPE="ext4"
/dev/vdb1: UUID="21e91bbc-7bca-4c08-***-88d5b3a2303d" TYPE="ext4" PARTLABEL="mnt" PART
UUID="576235e0-5e04-4b76-****-741cbc7e98cb"
/dev/vdb2: UUID="a7dcde59-8f0f-4193-***-362a27192fb1" TYPE="xfs" PARTLABEL="test" PART
UUID="464a9fa9-3933-4365-***-c42de62d2864"
```

6. Run the following command to mount the new partition:

```
mount /dev/vdb2 /mnt
```

Option 5: Resize the file system of a raw data disk

If a raw data disk contains a file system but no partitions, perform the following operations to resize the file system:

1. Run the following command to view the type of the file system:

```
df -Th
```

A command output similar to the following one is returned. In this example, the command output indicates that the file system type of /dev/vdb is xfs.

```
[root@ecshost ~] # df -Th

Filesystem Type Size Used Avail Use% Mounted on

devtmpfs devtmpfs 434M 0 434M 0% /dev

tmpfs tmpfs 446M 0 446M 0% /dev/shm

tmpfs tmpfs 446M 524K 446M 1% /run

tmpfs tmpfs 446M 0 446M 0% /sys/fs/cgroup

/dev/vdal ext4 20G 2.5G 17G 14% /

tmpfs tmpfs 90M 0 90M 0% /run/user/0

/dev/vdb xfs 20G 53M 20G 1% /mnt
```

- 2. Resize the file system based on the file system type.
 - o xfs file system

If the file system type is xfs, run the xfs growfs command as the root user:

```
xfs_growfs /mnt
```

In the preceding command, /mnt indicates the mount point of the file system.

? Note The new version xfs_growfs identifies the device to resize based on the mount point. Example: xfs_growfs / mnt . You can run the $xfs_growfs - -help$ command to check how to use xfs_growfs of different versions.

- o ext and btrfs file systems
 - a. Run the following command to view the mount information:

```
mount | grep "/dev/vdb"
```

b. Run the following command to unmount the /dev/vdb data disk:

```
umount /dev/vdb
```

c. Run the following command to view the operation result:

```
mount | grep "/dev/vdb"
```

A command output similar to the following one is returned:

```
[root@ecshost ~]# mount | grep "/dev/vdb"
/dev/vdb on /mnt type ext4 (rw,relatime,data=ordered)
[root@ecshost ~]# umount /dev/vdb
[root@ecshost ~]# mount | grep "/dev/vdb"
```

- d. Run one of the following commands to resize the file system based on the file system type.
 - If the file system type is ext, run the resize2fs command as the root user:

```
resize2fs /dev/vdb
```

• If the file system type is btrfs, run the btrfs command as the root user:

```
btrfs filesystem resize max /mnt
```

In the preceding command, /mnt indicates the mount point of the file system.

e. Run the following command to mount the disk to the mount point:

```
mount /dev/vdb /mnt
```

3. Run the df -h command to view the resize results:

```
df -Th
```

A command output similar to the following one is returned. The command output shows that the file system has a larger capacity, which indicates that the file system is resized.

```
[root@ecshost ~] # df -Th
Filesystem Type Size Used Avail Use% Mounted on
devtmpfs
         devtmpfs 434M 0 434M 0% /dev
          tmpfs 446M 0 446M 0% /dev/shm
tmpfs
          tmpfs 446M 524K 446M 1% /run
tmpfs
          tmpfs 446M 0 446M 0% /sys/fs/cgroup
tmpfs
/dev/vda1
tmpfs
          ext4
                  20G 2.5G 17G 14% /
         tmpfs 90M 0 90M 0% /run/user/0
                   30G 63M 30G 1% /mnt
/dev/vdb
           xfs
```

Related information

- Resize disks online for Linux instances
- Resize disks online for Windows instances
- Resize partitions and file systems of Linux system disks

11.3. Use LVs for Linux

11.3.1. Use LVM to create a logical volume

Logical Volume Manager (LVM) is a Linux mechanism for dynamically managing disks and disk partitions. This topic describes how to use LVM to create a logical volume (LV) on multiple disks of a Linux Elastic Compute Service (ECS) instance.

Prerequisites

- Multiple disks are created and attached to the Linux instance. For more information, see Create a disk and Attach a data disk.
- To prevent data loss caused by accidental changes, we recommend that you create a snapshotconsistent group to back up your data. For more information, see Create a snapshot-consistent group.

Context

In the LVM architecture, a logical layer is created on top of disks and disk partitions to help you manage disk partitions in a more flexible manner. The size of an LV can be dynamically adjusted without losing existing data. The existing LV remains unchanged even if you add new data disks.

Notice

- To prevent data loss, an LV cannot be created on disks that contain data.
- Each disk snapshot can back up only the data of a single disk. If you use a disk snapshot to roll back a disk after LVM is used to partition the disk, data inconsistency occurs. We recommend that you use snapshot-consistent groups to back up data. For more information, see Create a snapshot-consistent group.

Step 1: Create physical volumes (PVs)

1.

2. Run the following command to view information about all disks on the instance:

```
lsblk
```

A command output similar to the following one indicates that you can create a scalable LV on five disks by using LVM.

```
[root@ecs ~]# lsblk
      MAJ:MIN RM SIZE RO TYPE MOUNTPOINT
vda
      253:0
             0 40G 0 disk
└vda1 253:1
              а
                 40G
                      0 nart
      253:16 0
vdb
                 40G
                      0 disk
                      0 disk
vdc
      253:32
                 40G
vdd
      253:48 0 40G
                      0 disk
      253:64
             0 40G
                      0 disk
vdf
      253:80
              0 40G
                     0 disk
```

3. If LVM version 2 (LVM2) is not installed on your instance, run the following command to install LVM2:

```
yum install -y lvm2
```

4. Run the following command to create PVs:

```
pvcreate <Device name of data disk 1> \dots <Device name of data disk N>
```

In this example, run the following command to create PVs for the /dev/vdb, /dev/vdc, /dev/vdd, /dev/vde, and /dev/vdf data disks. Separate the device names of multiple data disks with spaces. In actual operations, replace the device names with the device names of your disks.

```
pvcreate /dev/vdb /dev/vdc /dev/vdd /dev/vde /dev/vdf
```

A command output similar to the following one is returned.

```
[root@ecs ~]# pvcreate /dev/vdb /dev/vdc /dev/vdd /dev/vde /dev/vdf
Physical volume "/dev/vdb" successfully created.
Physical volume "/dev/vdc" successfully created.
Physical volume "/dev/vdd" successfully created.
Physical volume "/dev/vde" successfully created.
Physical volume "/dev/vdf" successfully created.
```

5. Run the following command to view information about the created PVs on the instance:

```
lvmdiskscan | grep LVM
```

A command output similar to the following one is returned.

```
[root@ecs ~]# lvmdiskscan | grep LVM
  /dev/vdb [
                    40.00 GiB]
                                    physical volume
                                    physical volume
                    40.00 GiB]
  /dev/vdc
                    40.00 GiB]
  /dev/vdd
                                   physical volume
                   40.00 GiB] LVM
40.00 GiB] LVM
  /dev/vde
                                   physical volume
  /dev/vdf
                                    physical volume
       physical volume whole disks
       physical volumes
```

Step 2: Create a volume group (VG)

1. Run the following command to create a VG:

```
vgcreate <VG name> <Name of PV 1> ......<Name of PV N>
```

In this example, run the following command to create the lvm_01 VG and add the /dev/vdb, /dev/vdc, /dev/vdd, /dev/vde, and /dev/vdf PVs. Separate multiple PV names with spaces. In actual operations, replace the VG and PV names with your VG and PV names.

```
vgcreate lvm_01 /dev/vdb /dev/vdc /dev/vdd /dev/vde /dev/vdf
```

A command output similar to the following one is returned.

```
[root@ecs ~]# vgcreate lvm_01 /dev/vdb /dev/vdc /dev/vdd /dev/vde /dev/vdf
Volume group "lvm_01" successfully created
```

2. (Optional)Run the following command to add new PVs to the VG:

```
vgextend <VG name> <Name of PV 1> .....<Name of PV N>
```

In this example, run the following command to add the /dev/vdg PV to the lvm_01 VG. If you want to add multiple PVs, separate the PV names with spaces.

```
vgextend lvm_01 /dev/vdg
```

A command output similar to the following one is returned.

```
[root@ecs ~]# vgextend lvm_01 /dev/vdg
Volume group "lvm_01" successfully extended
```

3. Run the following command to view the VG information:

```
vgs
```

A command output similar to the following one is returned.

```
[root@ecs ~]# vgs
VG #PV #LV #SN Attr VSize VFree
lvm_01 5 _0 0 wz--n- 199.98g 199.98g
```

Step 3: Create an LV

1. Run the following command to create an LV:

```
lvcreate [-L <LV size>][ -n <LV name>] <VG name>
```

- ? Note
 - LV size: The LV size must be smaller than the remaining free space of the VG. The unit can be MiB, GiB, or TiB.
 - o LV name: You can specify a name for the LV.
 - VG name: the name of an existing VG.

In this example, run the following command to create a 150 GiB LV:

```
lvcreate -L 150g -n lv01 lvm_01
```

A command output similar to the following one is returned.

```
[root@ecs ~]# lvcreate -L 150g -n lv01 lvm_01
Logical volume "lv01" created.
```

2. Run the following command to view the LV details:

```
lvdisplay
```

A command output similar to the following one is returned.

```
[root@ecs ~]# lvdisplay
   - Logical volume -
 LV Path
                        /dev/lvm_01/lv01
 LV Name
                        lv01
 VG Name
                        lvm_01
 LV UUID
                        3dP9u0-6htd-PPYW-qlfZ-
 LV Write Access
                       read/write
 LV Creation host, time ecs, 2021-06-03 11:37:55 +0800
 LV Status
                       available
 # open
 LV Size
                        150.00 GiB
 Current LE
                        38400
 Segments
 Allocation
                        inherit
 Read ahead sectors
                        auto
  - currently set to
                        8192
 Block device
                        252:0
```

Step 4: Create and mount a file system

1. Run the following command to create a file system on the LV:

```
mkfs.<File system format> <LV path>
```

In these examples, create an ext4 file system and an xfs file system. You can run one of the following commands to create a file system of a specified format that suits your needs.

Create an ext4 file system

```
mkfs.ext4 /dev/lvm_01/lv01
```

A command output similar to the following one is returned.

```
[root@ecs ~]# mkfs.ext4 /dev/lvm_01/lv01
mke2fs 1.42.9 (28-Dec-2013)
Filesystem label=
OS type: Linux
Block size=4096 (log=2)
Fragment size=4096 (log=2)
Stride=0 blocks, Stripe width=0 blocks
9830400 inodes, 39321600 blocks
1966080 blocks (5.00%) reserved for the super user
First data block=0
Maximum filesystem blocks=2187329536
1200 block groups
32768 blocks per group, 32768 fragments per group
8192 inodes per group
Superblock backups stored on blocks:
        32768, 98304, 163840, 229376, 294912, 819200, 884736, 1605632, 2654208,
        4096000, 7962624, 11239424, 20480000, 23887872
Allocating group tables: done
Writing inode tables: done
Creating journal (32768 blocks): done
Writing superblocks and filesystem accounting information: done
```

Create an xfs file system

```
mkfs.xfs /dev/lvm_01/lv01
```

A command output similar to the following one is returned.

```
[root@ecs ~]# mkfs.xfs /dev/lvm_01/lv01
meta-data=/dev/lvm_01/lv01
                                isize=512
                                             agcount=4, agsize=9830400 blks
                                sectsz=512 attr=2, projid32bit=1
                                             finobt=0, sparse=0
                                crc=1
                                bsize=4096 blocks=39321600, imaxpct=25
data
                                            swidth=0 blks
                                sunit=0
naming
        =version 2
                                bsize=4096
                                            ascii-ci=0 ftype=1
                                bsize=4096 blocks=19200, version=2
         =internal log
                                sectsz=512
                                            sunit=0 blks, lazy-count=1
realtime =none
                                extsz=4096 blocks=0, rtextents=0
```

2. Run the following command to create a mount point. Example: /media/lv01.

If you want to use an existing mount point, skip this step.

```
mkdir /media/lv01
```

3. Run the following command to mount the file system:

In this example, set the LV path to /dev/lvm_01/lv01 and set the mount point to /media /lv01. In actual operations, modify the LV path and mount point based on your needs.

```
mount /dev/lvm_01/lv01 /media/lv01
```

4. Run the following command to view the mount information of the LV:

```
df -h
```

A command output similar to the following one is returned.

```
[root@ecs ~]# df -h
                         Size Used Avail Use% Mounted on
Filesystem
/dev/vda1
                                            6% /
                          40G
                               2.0G
                                      36G
                                  0 3.8G
                                            0% /dev
devtmpfs
                         3.8G
tmpfs
                         3.8G
                                  0
                                     3.8G
                                            0% /dev/shm
tmpfs
                         3.8G
                               488K
                                      3.8G
                                            1% /run
tmpfs
                                            0% /sys/fs/cgroup
                         3.8G
                                  0
                                     3.8G
                         768M
                                            0% /run/user/0
tmpfs
                                  0
                                     768M
/dev/mapper/lvm 01-lv01 150G 33M 150G
                                            1% /media/lv01
```

Related information

• Resize an LV by using LVM

11.3.2. Resize an LV by using LVM

This topic describes how to use Logical Volume Manager (LVM) to resize a logical volume (LV) on a Linux Elastic Compute Service (ECS) instance.

Prerequisites

- An LV is created. For more information, see Use LVM to create a logical volume.
- A disk is resized in the ECS console. For more information, see Step 2: Resize the disk in the ECS console. In this example, /dev/vdf is resized by 40 GiB.
- To prevent data loss caused by accidental changes, we recommend that you create a snapshot to back up your data. For more information, see Create a snapshot of a disk.

Procedure

1.

2. Run the following command to view information about the LV that is created on the ECS instance:

```
lvdisplay
```

The following command output indicates that the /dev/lvm_01/lv01 LV is created and has 150 GiB of capacity.

```
[root@ecs ~]# lvdisplay
   -- Logical volume
 LV Path
                         /dev/lvm_01/lv01
 LV Name
                         lv01
 VG Name
                         lvm 01
 LV UUID
                         3dP9u0-6htd-PPYW-q1fZ-p8K
 LV Write Access
                         read/write
 LV Creation host, time ecs, 2021-06-03 11:37:55 +0800
 LV Status
                         available
  # open
 LV Size
                         150.00 GiB
                         38400
  Current LE
  Segments
 Allocation
                         inherit
 Read ahead sectors
                         auto
   currently set to
                         8192
 Block device
                         252:0
```

3. Run the following command to resize the physical volume (PV) corresponding to the LV:

```
pvresize <PV name>
```

In this example, the name of the PV is /dev/vdf. In actual application, replace it with the name of your PV.

```
pvresize /dev/vdf
```

A command output similar to the following one is displayed:

```
[root@ecs ~]# pvresize /dev/vdf
Physical volume "/dev/vdf" changed
1 physical volume(s) resized or updated / 0 physical volume(s) not resized
```

4. Run the following command to view the PV usage:

```
pvs
```

The following command output indicates that the /dev/vdf PV has 80 GiB of free capacity, which consists of the original 40 GiB of capacity and the added 40 GiB of capacity.

```
[root@ecs ~]# pvs
          VG
                Fmt Attr PSize
                               PFree
 /dev/vdb
          lvm_01 lvm2 a-- <40.00g
                                   0
          lvm_01 lvm2 a--
                       <40.00g
 /dev/vdc
                                   0
 /dev/vdd
          lvm_01 lvm2 a-- <40.00g
                                   0
          lvm 01 lvm2 a--
 /dev/vde
                       <40.00g
                                9.98g
```

5. Run the following command to resize the LV:

```
lvextend [-L <LV size>] <LV name>
```

In this example, the capacity of the LV is resized.

```
lvextend -L +80G /dev/lvm_01/lv01
```

The following section describes the variables in this example. You must modify the variables based on your needs.

- +80GB: the amount of capacity to increase or decrease. You can resize the LV only when the volume group (VG) has remaining capacity.
- o /dev/lvm 01/lv01 : the name of the LV.

The following command output indicates that you increase the physical capacity by 80 GiB for the /dev/lvm 01/lv01 LV.

```
[root@ecs ~]# lvextend -L +806 /dev/lvm_01/lv01
Size of logical volume lvm_01/lv01 changed from 150.00 GiB (38400 extents) to 230.00 GiB (58880 extents).
Logical volume lvm_01/lv01 successfully resized.
```

6. Run one of the following commands to resize the file system of the LV.

Run one of the following commands based on the file system type of the LV. In these examples, ext4 and xfs file systems are used.

? Note If you are uncertain about the file system type of the LV, run the of -Th command to query it.

• To resize an ext4 file system, run the following command:

```
resize2fs /dev/lvm_01/lv01
```

• To resize an xfs file system, run the following command:

```
xfs_growfs /dev/lvm_01/lv01
```

7. Run the following command to check whether the file system is resized:

```
df -h
```

In the following command output, the total capacity of the LV is 230 GiB, which indicates that the file system is resized.

```
[root@ecs ~]# df -h
Filesystem
                        Size Used Avail Use% Mounted on
                        40G 2.0G
/dev/vda1
                                   36G
                                          6% /
devtmpfs
                               0 3.8G
                                          0% /dev
                        3.8G
tmpfs
                                          0% /dev/shm
                        3.8G
                                0 3.8G
tmpfs
                        3.8G
                             488K 3.8G
                                          1% /run
                                0 3.8G
                                          0% /sys/fs/cgroup
tmpfs
                        3.8G
                                          0% /run/user/0
                                0
                                   768M
/dev/mapper/lvm_01-lv01 230G
                              33M
                                   230G
                                          1% /media/lv01
```

Related information

• Use LVM to create a logical volume

11.4. Create a RAID array for a Linux instance

This topic describes how to use the mdadm command of the Linux operating system to create a 200-GiB Independent Array of Independent Disks (RAID) array for multiple data disks. In this topic, an Elastic Compute Service (ECS) instance that runs the Ubuntu operating system is used.

Prerequisites

Multiple disks are created and attached to the ECS instance. We recommend that you create disks of the same capacity and category. For more information about how to create and attach a disk, see Create a disk and Attach a data disk.

Context

RAID combines multiple disks into a disk array group. Compared with a single disk, a RAID array offers improved capacity, read/write bandwidth, reliability, and availability.

We recommend that you use the RAID 0 or RAID 1 mode and partition disks in the same size to maximize the use of disk space. We recommend that you do not use the RAID 5 or RAID 6 mode, because parity data in RAID 5 or RAID 6 mode consumes the IOPS of disks and deteriorates performance.

The following table describes the advantages, disadvantages, and application scenarios of the RAID 0 and RAID 1 modes.

Mode Advantage Disadvantage Scenario

Mode	Advantage	Disadvantage	Scenario
RAID 0	Uses the striping technique to allocate I/O loads to different disks. You can increase the throughput by extending disks. The capacity and bandwidth in the array are the sum of the capacity and bandwidth of different disks.	A damaged disk may cause loss of all data due to lack of data redundancy.	Has high requirements for I/O performance and is applicable when data is backed up in other methods or no data backup is required.
RAID 1	Provides higher data redundancy because data is stored in different disks as images. The minimum capacity and bandwidth values of disks are the capacity and bandwidth values of the RAID array.	The write performance is poor because data must be written to multiple disks at the same time.	Focuses more on fault tolerance than I/O performance in key applications.

Procedure

- 1.
- 2. Run the following command to view the information of all disk on the instance:

```
lsblk
```

A command output similar to the following one is returned.

```
root@ecs:~# lsblk
NAME MAJ:MIN RM SIZE RO TYPE MOUNTPOINT
vda
       253:0 0 40G 0 disk
└vda1 253:1 0 40G 0 part /
vdb 253:16 0 40G 0 disk
vdb
vdc
       253:32
               0
                  40G
                       0 disk
       253:48 0 40G
                       0 disk
vdd
vde
       253:64 0 40G
                       0 disk
vdf
       253:80 0 40G 0 disk
```

3. Run the **mdadm** command to create a RAID array named /dev/md0.

You can create an array in RAID 0 or RAID 1 mode based on your actual needs.

- ? Note
 - o In the following commands, /dev/vd[bcdef indicates that the RAID array consists of five disks: /dev/vdb, /dev/vdc, /dev/vdd, /dev/vde, and /dev/vdf. If you want to use other disks, replace them with your disk names.
 - If you are prompted that the mdadm tool is not installed, run the apt-get install md adm command to install the mdadm tool.
- Run the following command to create an array in RAID 0 mode:

```
mdadm --create /dev/md0 --level=0 --raid-devices=5 /dev/vd[bcdef]
```

- --level=0 indicates the array in RAID 0 mode that is used for data striping.
- --raid-devices=5 indicates that the RAID array consists of five disks.

/dev/vd[bcdef] indicates the five disks: /dev/vdb, /dev/vdc, /dev/vdd, /dev/vde, and /dev/vdf.

A command output similar to the following one is returned.

```
root@ecs:~# mdadm --create /dev/md0 --level=0 --raid-devices=5 /dev/vd[bcdef]
mdadm: Defaulting to version 1.2 metadata
mdadm: array /dev/md0 started.
```

• Run the following command to create an array in RAID 1 mode:

```
mdadm --create /dev/md0 --level=1 --raid-devices=5 /dev/vd[bcdef]
```

- --level=1 indicates the array in RAID 1 mode that is used for data mirroring.
- --raid-devices=5 indicates that the RAID array consists of five disks.
- /dev/vd[bcdef] indicates the five disks: /dev/vdb, /dev/vdc, /dev/vdd, /dev/vde, and /dev/vdf.
- 4. Run the following command to view the information of the /dev/md0 RAID array:

```
mdadm --detail /dev/md0
```

A command output similar to the following one is returned.

```
root@ecs:~# mdadm --detail /dev/md0
/dev/md0:
       Version: 1.2
 Creation Time : Sat May 8 15:10:52 2021
    Raid Level : raid0
    Array Size : 209551360 (199.84 GiB 214.58 GB)
  Raid Devices : 5
 Total Devices : 5
   Persistence : Superblock is persistent
   Update Time: Sat May 8 15:10:52 2021
         State : clean
Active Devices : 5
Working Devices : 5
Failed Devices: 0
 Spare Devices : 0
    Chunk Size : 512K
          Name: ecs:0 (local to host ecs)
          UUID : 09873fbc:5172dd8
        Events: 0
                    Minor
                            RaidDevice State
   Number
            Major
                                      active sync
      0
            253
                      16
                                0
                                                     /dev/vdb
            253
                      32
                                1
                                       active sync
                                                     /dev/vdc
                                                     /dev/vdd
            253
                      48
                                2
                                       active sync
      2
            253
                      64
                                       active sync
                                                     /dev/vde
            253
                      80
                                       active sync
                                                     /dev/vdf
```

5. Run the following command to create a file system on the RAID array. In this example, an ext4 file system is created.

You can also create a file system of other types.

```
mkfs.ext4 /dev/md0
```

A command output similar to the following one is returned.

6. Run the following command to create a configuration file that contains the information of the RAID array and configure the RAID array to be automatically reassembled when the ECS instance starts:

```
sudo mdadm --detail --scan | sudo tee -a /etc/mdadm/mdadm.conf
```

- 7. Mount the file system.
 - i. Run the following command to create a mount point. For example, create a mount point named /media/raid0.

```
mkdir /media/raid0
```

- **? Note** You can attach a disk to an existing directory such as /mnt.
- ii. Run the following command to mount the file system. For example, you can mount the /dev/md0 file system to the /media/raid0 mount point.

```
mount /dev/md0 /media/raid0
```

8. Run the following command to view the mount information of the RAID array:

```
df -h
```

The command output indicates that the /dev/md0 file system is mounted to the /media/rad0 mount point.

```
root@ecs:~# df
Filesystem
               Size Used Avail Use% Mounted on
udev
               1.9G
                      0 1.9G
                                  0% /dev
tmpfs
                381M
                     2.9M
                           378M
                                  1% /run
/dev/vda1
                40G
                     2.4G
                            36G
                                  7% /
tmpfs
                1.9G
                        0
                          1.9G
                                  0% /dev/shm
tmpfs
               5.0M
                         0
                          5.0M
                                  0% /run/lock
                1.9G
                                  0% /sys/fs/cgroup
tmpfs
                        а
                          1.9G
                                   0% /run/user/0
                381M
                            381M
/dev/md0
                           187G
                                   1% /media/raid0
                197G
                      60M
```

What's next

To configure the RAID array to be automatically loaded each time the ECS instance starts, perform the following operations on the /etc/fstab configuration file.

1. Run the following command to add the default startup setting to the /etc/fstab configuration file:

echo `blkid /dev/md0 | awk '{print \$2}' | sed 's/\"//g'` /media/raid0 ext4 defaults 0 0 >> /etc/fstab

- o /dev/md0 indicates the name of the disk array.
- /media/raid0 indicates the mount point. If you want to mount the RAID array to a different mount point, you must replace it with the corresponding mount point.
 - **? Note** To start the ECS instance when the RAID array is not mounted, you must add the nof ail setting. Even if an error occurs when you install a disk, the nof ail setting allows the ECS instance to be started. If your instance runs an Ubuntu operating system, you must also add the nobootwait setting.
- 2. Run the following command to mount all file systems in the /etc/fstab configuration file:

mount -a

11.5. Modify the UUID of a disk

If you create a disk from a snapshot of a disk on a Linux instance and attach the created disk to the Linux instance, the universally unique identifier (UUID) of the disk conflicts with that of the disk from which the snapshot is created. This topic describes how to modify the UUID of a disk.

Context

After a disk is created from a snapshot of a disk on a Linux instance, the UUID of the created disk is the same as that of the disk from which the snapshot is created. If you attach the created disk to the Linux instance, the UUID of the disk conflicts with that of the disk from which the snapshot is created. The following issues may occur:

- If you create a disk from a system disk snapshot of a Linux instance and attach the created disk to the Linux instance, Linux may be started from the newly attached data disk rather than the system disk.
- If your disk uses the XFS file system, the mount operation is prohibited due to a UUID conflict. The "mount: wrong fs type, bad option, bad superblock on /dev/vddl," message appears.

Therefore, after you create a disk from a snapshot of a disk on a Linux instance and attach the created disk to the Linux instance in the console, you must log on to the instance to modify the UUID of the created disk before you perform the mount operation. To modify the UUID of your disk, you can run the blkid command to query its file system type and choose one of the following methods based on the command output:

- If the command output is TYPE="ext4", TYPE="ext3", or TYPE="ext2", see Modify the UUID of an ext2, ext3, or ext4 file system.
- If the command output is TYPE="xfs", see Modify the UUID of an XFS file system.

Modify the UUID of an ext2, ext3, or ext4 file system

Note In this example, /dev/vdb1 is used. You must modify the related commands based on your device name.

1.

2. Run the following command to query the UUID of the created disk:

```
blkid
```

The following command output shows that the UUID of the disk created from a snapshot is the same as that of the disk from which the snapshot is created.

```
[root@ecs ~ 1# blkid
/dev/vda1: UUID="dcbdbcd3-f78c-4739-8cc7- 50da3b" TYPE="ext4"
/dev/vdb1: UUID="dcbdbcd3-f78c-4739-8cc7- 50da3b" TYPE="ext4"
[root@ecs ~ 1#
```

3. Run the following command to check the file system:

```
e2fsck -f /dev/vdb1
```

4. Run the following command to generate a new UUID for the created disk:

```
uuidgen | xargs tune2fs /dev/vdb1 -U
```

5. Run the following command to check whether the UUID is modified:

```
blkid
```

The following command output shows that the UUID of /dev/vdb1 is modified.

6. Run the following command to attach the created disk:

```
mount /dev/vdb1 /mnt
```

7. Configure the /etc/fstab file to automatically attach the created disk on startup.

For information about how to configure the /etc/fstab file, see Configure UUIDs in the fstab file to automatically attach data disks.

Modify the UUID of an XFS file system

Note In this example, /dev/vdd1 is used. You must modify the related commands based on your device name.

1.

2. Run the following command to query the UUID of the created disk:

```
blkid
```

The following command output shows that the UUID of the disk created from a snapshot is the same as that of the disk from which the snapshot is created.

```
[root@ecs ~]# blkid
/dev/vda1: UUID="dcbdbcd3-f78c-4739-8cc7- 50da3b" TYPE="ext4"
/dev/vdb1: UUID="56570712-2c72-42c0-9b13- 7cae0b" TYPE="ext4"
/dev/vdc1: UUID="65f0c62a-f980-4a58-8de5- 65b99f" TYPE="xfs" PARTLABEL="primary"
/dev/vdd1: UUID="65f0c62a-f980-4a58-8de5- 65b99f" TYPE="xfs" PARTLABEL="primary"
/dev/vdd1: UUID="65f0c62a-f980-4a58-8de5- 65b99f" TYPE="xfs" PARTLABEL="primary"
/dev/vdd1: UUID="65f0c62a-f980-4a58-8de5- 65b99f" TYPE="xfs" PARTLABEL="primary"
```

3. Run the following command to generate a new UUID for the created disk:

```
xfs_admin -U generate /dev/vdd1
```

4. Run the following command to check whether the UUID is modified:

```
blkid
```

The following command output shows that the UUID of /dev/vdd1 is modified.

```
| Type="ext4" | Continue of the continue of th
```

5. Run the following command to attach the created disk:

```
mount /dev/vdd1 /mnt
```

6. Configure the /etc/fstab file to automatically attach the created disk on startup.

For information about how to configure the /etc/fstab file, see Configure UUIDs in the fstab file to automatically attach data disks.

References

- Create a disk from a snapshot
- Configure UUIDs in the fstab file to automatically attach data disks

11.6. Configure UUIDs in the fstab file to automatically attach data disks

In Linux, you can configure the fstab file to have the file systems of data disks automatically mounted on the startup of an Elastic Compute Service (ECS) instance. If the fstab file is inappropriately configured and the attach sequence of your disks is changed, the ECS instance may not run normally after it is restarted. This topic describes how to configure universally unique identifiers (UUIDs) in the fstab file to automatically mount the system files of data disks. This can handle the preceding restart exception.

Prerequisites

The disks attached to the instance are partitioned and formatted. For more information, see Partition and format a data disk on a Linux instance.

Context

fstab allows you to identify a file system by using a disk partition name such as/dev/vdb1 or by using a UUID. The two methods differ in the following aspects:

- If disk partition names are used to identify file systems in the fstab file, the disk partitions may not be mounted to the original mount points if the attach sequence of the disks is changed. In this case, applications that run on your ECS instance may be affected.
- If UUIDs are used to identify file systems in the fstab file, the disk partitions can still be mounted to

the original mount points if the attach sequence of the disks is changed. Therefore, we recommend that you use UUIDs to identify file systems.

Procedure

1. Connect to the ECS instance.

For information about how to connect to an ECS instance, see Connect to a Linux instance by using a password or key.

2. Run the following command to view information of the disks attached to the instance:

```
fdisk -lu
```

A command output similar to the following one is returned.

```
[root@ecs ~]# fdisk -lu
Disk /dev/vda: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x000c2bef
   Device Boot
dev/vda1 *
                                                                                                                Blocks
41941999+
                                                                                            End
                                                       Start
                                                                                                                                                         System
                                                         2048
                                                                               83886046
                                                                                                                                                        Linux
Disk /dev/vdb: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x41a5a16d
   Device Boot
dev/vdb1
                                                                               End
83886079
                                                                                                                Blocks
41942016
                                                                                                                                                       System
Linux
                                                      Start
2048
Disk /dev/vdc: 42.9 GB, 42949672960 bytes, 83886080 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk label type: dos
Disk identifier: 0x93f147d9
    Device Boot
dev/vdc1
                                                                                            End
                                                                                                                                                        System
                                                                                                                     Blocks
                                                                               83886079
                                                                                                                 41942016
                                                                                                                                                        Linux
  |root@ecs |# _
```

3. Run the following command to query the UUIDs of the disks:

```
blkid
```

A command output similar to the following one is returned.

- 4. Run the following commands to create the mount points for the data disks:
 - Create the /test01 mount point for /dev/vdb1:

```
mkdir /test01
```

• Create the /test02 mount point for /dev/vdc1:

```
mkdir /test02
```

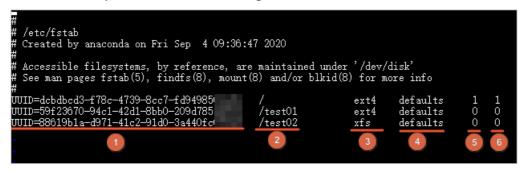
- 5. Add the mount information to the fstab file.
 - i. Run the following command to edit the fstab file:

```
vi /etc/fstab
```

- ii. Press the I key to enter the edit mode.
- iii. Add the following mount information:

```
UUID=59f23670-94c1-42d1-8bb0-209d7854*** /test01 ext4 defaults 0 0 UUID=88619b1a-d971-41c2-91d0-3a440fc0*** /test02 xfs defaults 0 0
```

A command output similar to the following one is returned.



No.	Section	Description
①	<file system=""></file>	The file systems to be mounted to the partitions. We recommend that you use the UUIDs of the file systems. You can run the blkid command to query the UUIDs of the file systems in the partitions.
2	<dir></dir>	The mount points of the file systems. You can create mount points. In this topic, the /test01 and /test02 mount points are created.
3	<type></type>	The types of file systems to be mounted to the partitions. You can run the blkid command to query the types of the file systems.
4	<options></options>	The parameters used for mounting. Typically, the defaults parameter is used. If you want to use multiple parameters, separate them with commas (,). Example: defaults, noatime . For more information about the <options> parameters, see fstab.</options>

No.	Section	Description	
S	<dump></dump>	Indicates whether the dump tool backs up the file systems. O: The dump tool does not back up the file systems. 1: The dump tool backs up the file systems. Typically, the dump tool is not used. In this case, this parameter is set to 0.	
6	<pass></pass>	 The priority in which fsck checks the file systems. 0: The file systems are not checked. 1: The file system corresponding to root directory (/) is checked. 2: All file systems except the one corresponding to root directory (/) are checked. Typically, this parameter is set to 0. 	

- iv. Press the ESC key to exit the edit mode after you complete the preceding configurations.
- v. Enter :wq and press the Enter key to save and exit the file.
- 6. Run the following command to view the fstab file:

```
cat /etc/fstab
```

A command output similar to the following one is returned.

- 7. Run the following command to mount the file systems of the data disk partitions.
 - To mount /dev/vdb1, run the following command:

```
mount /dev/vdb1 /test01
```

• To mount /dev/vdc1, run the following command:

```
mount /dev/vdc1 /test02
```

8. Run the following command to check whether the file systems of the data disk partitions are mounted:

```
df -h
```

A command output similar to the following one is returned:

After the fst ab file is configured, the system attaches the data disks after you restart the ECS instance.

FAQ

What do I do if system startup exceptions occur due to the configuration errors of the /etc/fstab file on Linux instances?

11.7. Shrink a disk

You cannot shrink system or data disks in Elastic Compute Service (ECS). You can use Alibaba Cloud Server Migration Center (SMC) to shrink disks.

Prerequisites

The preparations for the migration are complete. For more information, see Before you begin.

Context

SMC is aimed to balance the cloud-based and offline workloads of Alibaba Cloud users. You can also use SMC to shrink ECS disks.

SMC allows you to create a custom image from an ECS instance or migrate the data of an instance to another instance. When you import the migrated data, you can re-specify the size of a disk to shrink the disk. Before you migrate the data of an instance, take note of the following items:

• If you want to shrink a disk by migrating the data of the instance to which the disk is attached, make sure that the source instance and the destination instance are not the same one. The destination instance must be available and does not contain data or has all data backed up to an image, a snapshot, or a different device.

Warning After the migration task is created, all original data in the destination instance is deleted. If the destination ECS instance contains important data, we recommend that you do not migrate the data of the source instance to the destination instance. In this case, we recommend that you set Resource Type to ECS Image and then create an ECS instance by using a custom image.

• Some properties of the ECS instance are changed because the ECS instance is replaced when SMC is used to shrink the disk. For example, the ID (InstanceId) and public IP address of the instance are changed.

If your source instance resides in a virtual private cloud (VPC), you can convert the public IP address of the instance into an elastic IP address (EIP). This way, you can retain the public IP address. If you have already used EIPs or have less dependence on public IP addresses, we recommend that you use SMC to shrink disks.

Procedure

1. Import the data of the source instance.

You must download the SMC client to the ECS instance whose disks you want to shrink and then use the client to import the source ECS instance to SMC. For more information, see Step 1: Import the information of a migration source.

2. Create and start a migration task.

For information about how to shrink disks by migrating the data of the source instance to the destination instance, see Migrate servers to ECS instances. When you create a migration task, take note of the following items:

- Configure the System Disk and Data Disk parameters in the Target Disk Size section based on your requirements. The values of these parameters cannot be less than the amounts of used space of the system and data disk.
- If you want to create a custom image from the source instance and then create an ECS instance from the custom image, set **Resource Type** to **ECS Image** when you create a migration task.
- 3. Wait until the migration task is completed.
 - If the migration task is in the **Completed** state, the task is completed and you can view the destination ECS instance.
 - If the migration task is in the InError state, the task failed. You can check the logs to troubleshoot the failure. Then, restart the migration task. For information about common errors and solutions, see SMC FAQ.

11.8. Encrypt data stored on ECS resources

Data encryption is suitable for scenarios that require data security and regulatory compliance. To protect data stored on Elastic Compute Service (ECS) resources such as system disks, data disks, and images, you can encrypt the resources. You can use encrypted disks and images to create ECS instances to ensure data privacy and security. This topic describes how to encrypt disks, snapshots, and images and the items that you must take note of during the encryption processes.

Prerequisites

Key Management Service (KMS) is activated in the region where the resources are located. For more information, see Activate KMS.

Context

By default, KMS uses the default service customer master keys (CMKs) to encrypt data. You can also use Bring Your Own Key (BYOK) keys (custom keys imported by using BYOK feature) to encrypt data. You must use a CMK and a data key and the envelope encryption mechanism to encrypt data stored on each disk. For more information, see 加密概述.

When you use CMKs to encrypt data, take note of the items described in the following table.

Each account has a single default service CMK in each region. The default service CMK cannot be deleted or disabled. • The first time that you select a BYOK key to encrypt a disk in the ECS console, click Go to Authorize and follow on-screen tips to attach AliyunECSDiskE ncryptDefaultRole to allow ECS to access your KMS resources. For more information about the role, see RAM overview. • Select Aliyun_AES_256 or Aliyun_SM4 when you create a CMK in the KMS console. CMKs of other types cannot be used to encrypt disks. • Before you delete or disable a BYOK key, make sure that the BYOK key is not associated with disks or that the disk with which the BYOK key is associated is not attached to instances. This prevents disk data loss or instance startup failures. You can call the DescribeDisks operation to query the information of disks with which CMKs are associated. BYOK key BYOK key BYOK key is deleted, it cannot be recovered. If the BYOK key has been used to encrypt data or create data keys, the data or the data keys cannot be decrypted. Before a BYOK key becomes invalid, we recommend that you disable the BYOK key or check whether the BYOK key is associated with cloud resources to prevent data loss.	ltem	Description		
click Go to Authorize and follow on-screen tips to attach ncryptDefaultRole to allow ECS to access your KMS resources. For more information about the role, see RAM overview. Select Aliyun_AES_256 or Aliyun_SM4 when you create a CMK in the KMS console. CMKs of other types cannot be used to encrypt disks. Before you delete or disable a BYOK key, make sure that the BYOK key is not associated with disks or that the disk with which the BYOK key is associated is not attached to instances. This prevents disk data loss or instance startup failures. You can call the DescribeDisks operation to query the information of disks with which CMKs are associated. After a BYOK key is deleted, it cannot be recovered. If the BYOK key has been used to encrypt data or create data keys, the data or the data keys cannot be decrypted. Before a BYOK key becomes invalid, we recommend that you disable the BYOK key or check whether the BYOK key is associated with cloud resources to prevent data loss.	Default service CMK			
Notice A BYOK key becomes invalid when it is deleted or disabled. As a result, the data of the disks, images, and snapshots that were encrypted by using the BYOK key may not be recovered. Note that if BYOK keys become invalid due to your operations, you are responsible for the risk that data stored on disks with which the BYOK keys are associated may not be recovered.	BYOK key	click Go to Authorize and follow on-screen tips to attach alignmecsDiskE noryptDefaultRole to allow ECS to access your KMS resources. For more information about the role, see RAM overview. Select Aliyun_AES_256 or Aliyun_SM4 when you create a CMK in the KMS console. CMKs of other types cannot be used to encrypt disks. Before you delete or disable a BYOK key, make sure that the BYOK key is not associated with disks or that the disk with which the BYOK key is associated is not attached to instances. This prevents disk data loss or instance startup failures. You can call the DescribeDisks operation to query the information of disks with which CMKs are associated. After a BYOK key is deleted, it cannot be recovered. If the BYOK key has been used to encrypt data or create data keys, the data or the data keys cannot be decrypted. Before a BYOK key becomes invalid, we recommend that you disable the BYOK key or check whether the BYOK key is associated with cloud resources to prevent data loss. Notice A BYOK key becomes invalid when it is deleted or disabled. As a result, the data of the disks, images, and snapshots that were encrypted by using the BYOK key may not be recovered. Note that if BYOK keys become invalid due to your operations, you are responsible for the risk that data stored on disks with which the BYOK		

Encrypt a system disk

A system disk is a disk that contains an operating system. System disks can be created only along with instances and share the same lifecycle as the instances. You can encrypt system disks during instance creation only when the requirements described in the following table are met.



Usage notes

ltem	Description
Instance family	The instance family of the instance cannot be ecs.ebmg5, ecs.ebmgn5t, ecs.ebmi3, ecs.sccg5, ecs.scch5, ecs.ebmc4, and ecs.ebmhfg5. For more information, see Instance family.
lmage	The instance uses a public or custom image. The instance cannot use a shared image or an Alibaba Cloud Marketplace image.

Item	Description
Disk category	The disk is an enhanced SSD (ESSD).

Procedure

For information about how to encrypt a system disk, see Encrypt a system disk.

Encrypt a data disk

After a data disk is encrypted, both data in transit and data at rest on the disk are encrypted. You can encrypt data disks when you create instances or when you create the data disks.

Usage notes

If you select **Create from Snapshot** when you encrypt a data disk on the Disk page in the ECS console, you can select **Disk Encryption** only when the following requirements described in the following table are met.

Note The feature is in public preview and supported only in the following zones: Hong Kong Zones B and C and Singapore Zones B and C.

ltem	Description
Instance family	The instance family of the instance cannot be ecs.ebmg5, ecs.ebmgn5t, ecs.ebmi3, ecs.sccg5, ecs.scch5, ecs.ebmc4, and ecs.ebmhfg5. For more information, see Instance family.
lmage	The instance uses a public or custom image. The instance cannot use a shared image or an Alibaba Cloud Marketplace image.
Disk category	The disk is an enhanced SSD (ESSD).

Procedure

For information about how to encrypt a data disk, see Encrypt a data disk.

Encrypt a snapshot

If you create a snapshot from an encrypted disk, the snapshot is encrypted.

Procedure

For information about how to create a snapshot, see Create a snapshot of a disk.

Share an encrypted image

Shared images can be used to deploy ECS instances across accounts. If the encryption feature is enabled for disks of an instance, custom images created from the instance are encrypted. You can share the encrypted custom images to other Alibaba Cloud accounts. These Alibaba Cloud accounts can use the shared images to create instances.

? Note Encrypted custom images can be shared only in the China (Beijing), China (Shanghai), China (Hong Kong), and Singapore (Singapore) regions.

Usage notes

ltem	Description		
	To share encrypted custom images, you must use Resource Access Management (RAM) to create a RAM role named AliyunECSShareEncryptImageDefaultRole and attach the role to the account that shares the custom images. For more information about how to attach the role, see Share encrypted custom images. We recommend that you follow the principle of least privilege and attach only image-related policies when you create and attach the AliyunECSShareEncryptImageDefaultRole role.		
	Sample code:		
Attach a role	<pre>{ "Version": "1", "Statement": [</pre>		

ltem	Description
Use a CMK	We recommend that you create a BYOK key to share encrypted images. This prevents security risks caused by KMS key leakage.
Unshare an encrypted image	 The account is no longer able to query the image by using the ECS console or by calling an API operation. The account is no longer able to create ECS instances or replace system disks by using the image. The system disks of ECS instances that were created from the shared image cannot be re-initialized.
Disable an encrypted image	 The account is no longer able to create ECS instances or replace system disks by using the image. The system disks of ECS instances that were created from the shared image cannot be re-initialized.

Procedure

For information about how to share an encrypted image, see the following topics:

- 1. Share encrypted custom images
- 2. Share or unshare a custom image

12.Best practices for tag design

Increased cloud resources are hard to manage without tags. Tags can be used to manage, group, and search for resources. These resources include personnel, financial costs, and cloud services. This topic describes the best practices for tag design.

Scenarios

Tags are applicable to the following scenarios:

- Management of application publishing procedures
- Resource tracking and tag-based group search and management
- Tag- and group-based automated O&M by using Alibaba Cloud services such as Operation Orchestration Service (OOS), Resource Orchestration Service (ROS), Auto Scaling, and Cloud Assistant
- Tag-based cost management and cost allocation
- Resource- or role-based access control

Principles

You can implement the best practice for tag design based on the following principles:

- Mutual exclusivity
- Collective exhaustion
- Limit ed values
- Considering ramifications of future changes
- Simplified design

Mutual exclusivity

To implement the mutual exclusivity principle, we recommend that an attribute has only a single tag key. For example, if you use the owner tag key to represent the owner attribute, you cannot use other tag keys such as *own* or *belonger* to represent this attribute.

Collective exhaustion

Collective exhaustion indicates that when you plan resources, you must plan tags at the same time and prioritize the tag keys. All resources must have tags that consist of the planned tag keys and the corresponding tag values.

- Each tag key-value pair must be named in a standard format.
- Collective exhaustion is a prerequisite for future tag-based access control, cost tracking, automated O&M, and group search.

Limited values

This principle indicates that excess tag values must be removed and that only core tag values are retained.

Procedures for resource management, access control, automated O&M, and cost allocation can be simplified by implementing this principle. You can also use tags and automation tools under this principle to manage resources. Elastic Compute Service (ECS) allows you to control tags by calling API operations, which makes it easy to automatically manage, search for, and filter resources.

Considering ramifications of future changes

When you plan tags under the limited values principle, you must consider the impact of adding or removing tag values to improve the flexibility of modifying tags.

If you modify tags, tag-based access control, automated O&M, or related billing reports may change. For corporate or personal business, the best practice is to create business-related tag groups to manage resources in technical, business, and security dimensions. When you use automated O&M tools to manage resources and services, you can add automation-specific tags to aid in automated O&M.

Simplified design

Simplified design means that when you plan tags, you must create tag keys that have fixed dimensions to simplify the use of tag keys. By implementing this principle, you can reduce operation errors caused by redundant tag keys.

- You can create business-related tag groups to manage resources in technical, business, and security dimensions.
- When you use automated O&M tools to manage resources and services, you can add automation-specific tags to the resources and services.

Examples of designing tag keys

The following table describes the tag naming examples in the business dimension. We recommend that you use lowercase letters to name tags.

Dimension	Tag key	Tag value
Organization	companydepartmentorganizationteamgroup	Organization-specific names
Business	productbusinessmoduleservice	Business-specific names
Role	roleuser	 network administrator application administrator system administrator opsuser devuser testuser
Purpose	purposeuse	Specific purposes

Dimension	Tag key	Tag value
Project	 From project dimensions: project risk schedule subtask environment From personnel dimensions: sponsor member decisionmaker or owner creator 	Project-related values
Business department (to implement cost allocation and business tracking)	costcenterbusinessunitbizfinancecontact	Department-related values
Owner from the finance dimension (to identify the resource owner)	owner	Names or emails
Customers from the finance dimension (to identify the customers that a specific resource group serves)	Custom or actual values	Customer names
Project from the finance dimension (to identify the projects that are supported by specific resources)	project	Parameter
Order from the finance dimension	order	Order category IDs

References

- Search for resources by tag
- Use OOS to modify a tag value of multiple resources
- Implement automatic resource monitoring by group based on tags
- Create a resource with a specific tag

Related API operations

- TagResources
- ListTagResources
- UntagResources

13.Set the boot mode of custom images to the UEFI mode by calling API operations

This topic describes how to set the boot mode of custom images to the Unified Extensible Firmware Interface (UEFI) mode.

Context

Alibaba Cloud allows you to call API operations to set the boot mode of custom images to the UEFI mode. You may need custom images that support the UEFI boot mode in specific scenarios, such as those in which you use ebmg6a, ebmc6a, or ebmr6a instances. You can call the ImportImage or ModifyImageAttribute operation to set the boot mode of the custom images to the UEFI mode.

Limits

The following limits apply when you set the boot mode of custom images:

- You can set the boot mode of custom images only by calling Elastic Compute Service (ECS) API
 operations.
- You cannot query the boot mode of custom images by using the ECS console or by calling the query API operations such as the Describelmages operation.

Precautions

Take note of the following items:

- When you call an API operation to set the boot mode of custom images to the UEFI mode, the instances that use the custom images start in UEFI mode. Make sure that configurations of the custom images support the UEFI boot mode. Otherwise, the instances cannot start.
- Some instance families that require custom images to use the UEFI boot mode are in invitational or public preview. If exceptions occur when you use the UEFI features and resources, .

Methods of setting the boot mode

You can use one of the following methods to set the boot mode of a custom image to the UEFI mode based on your business needs.

Method	Description
Method 1: Set the UEFI boot mode when you import an image by calling the ImportImage operation.	If you want to import an image that supports the UEFI boot mode to Alibaba Cloud, you can call the ImportImage operation and set BootMode to UEFI . For more information about this operation, see ImportImage.

Method	Description
Method 2: Change the boot mode of a custom image to the UEFI mode by calling the ModifylmageAttribute operation.	If your Alibaba Cloud account has custom images that support the UEFI boot mode, you can call the ModifylmageAttribute operation to modify the attributes of the custom images and set BootMode to UEFI . For more information about this operation, see ModifylmageAttribute.

You can use one of the following methods to call ECS API operations to set the boot mode of custom images to the UEFI mode.

- Alibaba Cloud OpenAPI Explorer. You can use OpenAPI Explorer to call ECS API operations. For more information, see OpenAPI Developer Portal.
- ECS SDKs. You can use ECS SDKs to call ECS API operations. For more information, see SDK overview.
- Alibaba Cloud CLI. You can run the **aliyun** command in the shell command line to interact with Alibaba Cloud services and manage your Alibaba Cloud resources. For more information, see What is Alibaba Cloud CLI? .

14.Best practices for using custom images 14.1. Overview

Alibaba Cloud offers a variety of methods for you to create custom images. You can import a local custom image or create custom images by using Operation Orchestration Service (OOS), snapshots, ECS instances, or Packer. If you frequently use custom images, we recommend that you update them on a regular basis to ensure data security. For example, you can upgrade the OS patches, update the middleware or certificates, or install the latest third-party software for custom images.

Image creation methods

176

The following table compares the image creation methods offered by Alibaba Cloud.

Method	Advantage	Disadvantage	Billing
Use OOS to update custom images	 Official templates, no coding required Online operation, installation-free Password-free logon, secure and reliable Visualized execution process 	None.	You may be charged for resources such as instances, disks, and snapshots. For more information, see Overview.
 Create a custom image from a snapshot Create a custom image from an instance 	Easy operations in the ECS console	 The image creation procedure becomes complex as more pre-installed software is added. Difficult to ensure correct and consistent manual operations. High maintenance cost. 	You may be charged for snapshots. For more information, see Snapshots.
Create a custom image by using Packer	Open source tool that can be used by most cloud service providers	 Installation and maintenance required. Script writing required. 	You may be charged for resources such as instances, disks, and snapshots. For more information, see Overview.

Method	Advantage	Disadvantage	Billing
Create and import a custom image	On-cloud application deployment	Difficult operations. You must know how to use the GUI and configure virtualization platform drivers.	You may be charged for data stored in Object Storage Service (OSS) and snapshots. For more information, see Snapshots and Billing overview in OSS documentation.

Image creation procedure

The preceding image creation methods (except importing a local custom image) depend on the status and application data of an ECS instance at a specific point in time. OOS and Packer automate the image creation procedure by automatically creating a temporary instance and releasing it after the creation, making them more suitable for agile development.

The image creation procedure varies depending on the method you choose:

- If you use OOS to update a custom image, you must select a public template (such as ACS-ECS-UpdateImage) or customize an O&M task template. You can view the process on the YAML, JSON, or Graphical Template tab of OOS in the ECS console.
- If you plan to create a custom image by using a snapshot or an ECS instance, you must use an existing snapshot or a running ECS instance in the ECS console. After the image is created, you must release the temporary ECS instance you created for building the custom image.
- If you use Packer to create a custom image, you must write scripts by using a JSON template such as the image builder.
- If you plan to create and import a local custom image, you must first configure a virtual machine (such as VirtualBox VM) by following the custom image importing instructions, and then use OSS to import the image to ECS. For more information, see Instructions for importing images.

14.2. Use OOS to update custom images

This topic describes how to update custom images by using Operation Orchestration Service (OOS). OOS provides public templates to update images automatically. To create an immediate or scheduled O&M task, you only need to select a source image and specify required parameters such as the Cloud Assistant command in a public template. The O&M task is then automatically executed based on the definitions in the template.

Prerequisites

Context

You can use the ACS-ECS-UpdateImage template to sequentially execute the following tasks to update the image to a new custom image:

- 1. Check whether the name of the new custom image already exists and make sure that the name complies with the naming conventions.
- 2. Create and launch a temporary ECS instance based on parameters such as the instance type, source

image ID, and security group ID that you have configured.

- 3. Check whether the Cloud Assistant client is installed on the temporary ECS instance. If not, install the Cloud Assistant client.
- 4. Run scripts by using Cloud Assistant to update the system environment of the instance.
 - **? Note** OOS calls Cloud Assistant API operations to run shell, bat, or PowerShell scripts to update the system environments of ECS instances. For more information, see Overview.
- 5. Stop the temporary ECS instance.
- 6. Create a custom image from the temporary ECS instance.
- 7. Release the temporary ECS instance.

Procedure

- 1. Log on to the OOS console.
- 2. (Optional)If you are using OOS for the first time, click **Enable Now** to activate OOS.
- 3. In the left-side navigation pane, click Public Templates.

4.

- 5. In the ACS-ECS-UpdateImage section, click Create Execution.
- 6. On the **Create** page, perform the following operations:
 - i. In the Basic Information step, retain the default settings and click Next: Parameter Settings.
 - ii. In the **Parameter Settings** step, specify parameters to automatically create or update custom images. The following table describes these parameters.

Parameter	Description	Example
targetlmageName	The name of the new custom image. The name must comply with the regular expression /^[A-Z a-z0-9\]*\$/, and cannot be the same as an existing image name.	add_testtxt_2019101 0
The ID of the image that you want to update.		
sourcelmageld	Note If you have not created a custom image, you can use a public image, such as centos_7_06_64_20G_alibase_2019 0711.vhd.	m-bp13y4of6mdoqw *****
instanceType	The instance type of the temporary ECS instance to be created. For more information, see Instance family.	ecs.g5.xlarge
securityGroupId	The ID of the security group to which the temporary ECS instance belongs.	sg-bp1azkttqpldxg** ****

Parameter	Description	Example
vSwitchId	The ID of the vSwitch for the temporary ECS instance. The vSwitch and the security group must be in the same VPC.	vsw-bp1s5fnvk4gn2t w*****
ramRoleName	The RAM role of the temporary ECS instance.	TestRAMRole
commandType	 The type of the script that you plan to run by using Cloud Assistant. RunShellScript: shell scripts for Linux instances. RunBatScript: bat scripts for Windows instances. RunPowerShellScript: PowerShell scripts for Windows instances. 	RunShellScript
tags	 Tag Key (Required): the tag key of the image. Tag Value (Optional): the tag value of the image. Attach Resource Tag to Execution: binds resource tags to the OOS template. This option is selected by default. After resource tags are bound to the OOS template, you can use tags to filter the execution results of resources. 	Tag key: ECSTag value: Image
commandContent	The content of the script to be run on the temporary ECS instance.	echo "hello world" >/root/test.txt
Permissions	Optional. Valid values: Use Existing Permissions of Current Account: This is the default value. You have all the permissions granted to your account. Make sure that you have the permissions to call the ECS API operations required to create custom images. Specify RAM Role and Use Permissions Granted to This Role: If a RAM role is specified, OOS performs O&M tasks by assuming that RAM role.	Use Existing Permissions of Current Account

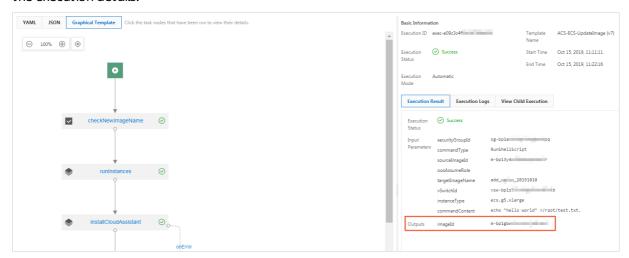
iii. Click Next: OK.

iv. Confirm 0&M task details and high-risk operations. Then, click Confirm and Create.

7. In the left-side navigation pane, click Executions to view the created O&M task.

Result

If the O&M task is created and is in the **Running** state, the custom image is being updated. When **Execution Status** changes to **Success**, the image is updated. You can view the ID of the new image in the execution details.



Note To view the update process, click **Details** of the O&M task to view **Logs**.

Related information

References

•

Introduction to OOS

14.3. Packer: machine images as code 14.3.1. Benefits of using Packer to create custom images

This topic describes the benefits of using Packer to create custom images. Packer is an easy-to-use tool that automates the creation of images. Using packer to create a custom image is as simple as specifying basic image information, the software to be installed, and the relevant configurations in a JSON template.

Prerequisites

Context

Packer is a tool provided by HashiCorp to automate the creation of images. A JSON template is used to ensure that created images are identical each time. Packer simplifies image testing and updating and minimizes image O&M and management costs. For more information, visit the official Packer website.

Operation conditions

This topic highlights the benefits of Packer in DevOps scenarios by comparing the procedure to create a custom image from an Elastic Compute Service (ECS) instance and the procedure to create a custom image by using Packer. The examples provided in this topic are based on the following conditions:

- Region: China (Beijing). For more information, see Regions and zones.
- Operating system: Cent OS 7.3 64-bit. In these examples, the public image cent os_7_03_64_20G_alibase_20170818.vhd is used. You can view the image lists of other operating systems in the ECS console or by calling the Describelmages operation.
- Custom service: Redis.
- Whether to retain temporary resources: No.

Note Paid resources such as instances, public IP addresses, and snapshots will be created in the following procedures. To avoid additional fees, we recommend that you release these resources after custom images are created.

Create a custom image from an ECS instance

This example shows how to create a custom image from an ECS instance in the ECS console.

- 1.
- 2.
- 3.
- 4.
- 5. Create an ECS instance. For more information, see Create an instance by using the wizard.

To minimize charges and simplify the procedure, you can create the instance with the following configurations:

- Billing Method: pay-as-you-go. For more information, see Pay-as-you-go.
- Instance Type: ecs.t5-lc1m1.small. For more information, see Instance family.
- o Public Image: Cent OS 7.3 64-bit.
- o VPC: the default VPC.
- Security Group: the default security group.
- Public IP Address: If your instance does not need to access the Internet, you can choose not to
 assign a public IP address to the instance, and then connect to the instance by using a Virtual
 Network Console (VNC) management terminal. For more information, see Connect to a Linux
 instance by using a password.

6.

- 7. Run the yum install redis.x86 64 -y command to install Redis.
- 8. Return to the homepage of the ECS console and select the China (Beijing) region.
- 9. Create an image from the instance you created. For more information, see Create a custom image from an instance.

10.

- 11. On the Images page, view the progress of image creation.
- 12. (Optional) After the image is created, release the temporary resources created in this procedure,

including the instance. If an Elastic IP Address (EIP) is used, you can also release the EIP.

Create a custom image by using Packer

This example shows how to use Packer to create a custom image. Before you perform the operations, make sure you have installed Packer. For more information, visit Install Packer or see Create a custom image by using Packer. To create a custom image by using Packer, perform the following steps:

1. Create an *alicloud.json* file that contains the following information:

```
"variables": {
 "access_key": "{{env `ALICLOUD_ACCESS_KEY`}}",
 "secret key": "{{env `ALICLOUD SECRET KEY`}}"
},
"builders": [{
 "type": "alicloud-ecs",
 "access_key":"{{user `access_key`}}",
 "secret key":"{{user `secret_key`}}",
 "region": "cn-beijing",
 "image name": "packer basic",
 "source_image":"centos_7_03_64_20G_alibase_20170818.vhd",
 "ssh username":"root",
 "instance_type":"ecs.t5-lc1m1.small",
 "internet_charge_type": "PayByTraffic",
 "io_optimized":"true"
"provisioners": [{
 "type": "shell",
 "inline": [
   "sleep 30",
    "yum install redis.x86 64 -y"
} ]
```

Packer parameters

Parameter	Example	Description
variables{"variabl e1":"value"}	variables{"access _key":"{{env `ALICLOUD_ACCES S_KEY`}}"}	Defines the variables used in image builders. These variables are used to prevent information such as AccessKey pairs from being disclosed. The variables use user-defined values at runtime.
builders{"type":"v alue"}	builders{"type":" alicloud-ecs"}	The image builders defined by Packer. For more information, see builders. Alibaba Cloud supports alicloud-ecs, also known as Alicloud Image Builder, which is used to create custom images in Alibaba Cloud ECS.

Parameter	Example	Description				
provisioners{"typ e":"value"}	provisioners{"typ e":"shell"}	The image provisioners defined by Packer to specify the operations that need to be performed on the temporary instance. The shell provisioner is used in this example. After Packer connects to the Linux instance, Packer runs shell commands to install services based on the shell provisioner configuration. For example, Packer runs the yum install redis.x86_64 -y command to install Redis.				

Alibaba Cloud parameters

Parameter	Data type	Example	Description	Import ance
			The AccessKey ID of your account. For more information, see Obtain an AccessKey pair.	
access_key	String	LT AlnPyXXXXQXX XX	Note To avoid disclosing the AccessKey pair of your Alibaba Cloud account, we recommend that you create a RAM user and use the credentials of the RAM user to create an AccessKey pair. For more information, see Create a RAM user.	High
secret_key	String	CM1ycKrrCekQ0dh XXXXXXXXXI7yav UT	The AccessKey secret of your account.	High
region	String	The region in which to create the custom image. For more information, see Regions and zones.		High
image_name	String	packer_basic	The name of the custom image you want to create. This name must be globally unique in Alibaba Cloud ECS.	
source_image	String	centos_7_03_64_ 20G_alibase_2017 0818.vhd	The ID of the Alibaba Cloud public image based on which to create a custom image. The created custom image uses the same operating system as the public image.	High
instance_type	String	ecs.t5- lc1m1.small	The instance type of the temporary ECS instance used to create the custom image. For more information, see Instance family.	Low

Parameter	Data type	Example	Description	Import ance
internet_charge_t ype	String	PayByT raffic	The billing method for the network usage of the temporary instance. Recommended value: PayByTraffic.	Low
io_optimized	Boolea n	true	Specifies whether the temporary instance is I/O optimized. Recommended value: true.	Low

2. Run the following command to create a custom image:

```
packer build alicloud.json
```

Note It may take a few minutes for the image to be created. After the image is created, it is displayed in the image list of the specified Alibaba Cloud region. You can view the image in the ECS console or by calling the DescribeImages operation.

A log is generated during the image creation process. This log records all the image creation operations, including checking parameters, creating temporary resources, pre-installing software, creating target resources, and releasing temporary resources.

```
alicloud-ecs output will be in this color.
==> alicloud-ecs: Prevalidating image name...
   alicloud-ecs: Found image ID: centos 7 03 64_20G_alibase_20170818.vhd
==> alicloud-ecs: Creating temporary keypair: packer xxx
==> alicloud-ecs: Creating vpc
==> alicloud-ecs: Creating vswitch...
==> alicloud-ecs: Creating security groups...
==> alicloud-ecs: Creating instance.
==> alicloud-ecs: Allocating eip
==> alicloud-ecs: Allocated eip xxx
   alicloud-ecs: Attach keypair packer xxx to instance: i-xxx
==> alicloud-ecs: Starting instance: i-xxx
==> alicloud-ecs: Using ssh communicator to connect: ***
==> alicloud-ecs: Waiting for SSH to become available...
==> alicloud-ecs: Connected to SSH!
==> alicloud-ecs: Provisioning with shell script: /var/folders/k /nv2r4drx3xxxxxxxxxxx
db40000gn/T/packer-shell260049331
   alicloud-ecs: Loaded plugins: fastestmirror
    alicloud-ecs: Determining fastest mirrors
   alicloud-ecs: Resolving Dependencies
   alicloud-ecs: --> Running transaction check
    alicloud-ecs: ---> Package redis.x86 64 0:3.2.12-2.el7 will be installed
    alicloud-ecs: --> Processing Dependency: libjemalloc.so.1()(64bit) for package: red
is-3.2.12-2.el7.x86 64
   alicloud-ecs: --> Running transaction check
    alicloud-ecs: ---> Package jemalloc.x86 64 0:3.6.0-1.el7 will be installed
    alicloud-ecs: --> Finished Dependency Resolution
    alicloud-ecs:
    alicloud-ecs: Dependencies Resolved
    alicloud-ecs:
```

```
alicloud-ecs: ------
  alicloud-ecs: Package
                              Arch
                                            Version
                                                                  Repositor
    Size
  alicloud-ecs: Installing:
                           x86 64
   alicloud-ecs: redis
                                             3.2.12-2.el7
544 k
  alicloud-ecs: Installing for dependencies:
  alicloud-ecs: jemalloc x86 64
                                            3.6.0-1.el7
                                                                  epel
105 k
  alicloud-ecs:
  alicloud-ecs: Transaction Summary
   alicloud-ecs: =======
   alicloud-ecs: Install 1 Package (+1 Dependent package)
   alicloud-ecs:
   alicloud-ecs: Total download size: 648 k
   alicloud-ecs: Installed size: 1.7 M
   alicloud-ecs: Downloading packages:
  alicloud-ecs: -----
   alicloud-ecs: Total
                                                           2.2 MB/s | 648 kB
00:00
  alicloud-ecs: Running transaction check
  alicloud-ecs: Running transaction test
   alicloud-ecs: Transaction test succeeded
   alicloud-ecs: Running transaction
  alicloud-ecs: Installing: jemalloc-3.6.0-1.el7.x86 64
   alicloud-ecs: Installing: redis-3.2.12-2.el7.x86 64
  alicloud-ecs: Verifying: redis-3.2.12-2.el7.x86 64
   alicloud-ecs: Verifying : jemalloc-3.6.0-1.el7.x86 64
2/2
  alicloud-ecs:
   alicloud-ecs: Installed:
   alicloud-ecs: redis.x86_64 0:3.2.12-2.e17
   alicloud-ecs:
   alicloud-ecs: Dependency Installed:
  alicloud-ecs: jemalloc.x86 64 0:3.6.0-1.el7
   alicloud-ecs:
   alicloud-ecs: Complete!
==> alicloud-ecs: Stopping instance: i-xxx
==> alicloud-ecs: Waiting instance stopped: i-xxx
==> alicloud-ecs: Creating image: packer basic
   alicloud-ecs: Detach keypair packer xxx from instance: i-xxx
==> alicloud-ecs: Cleaning up 'EIP'
==> alicloud-ecs: Cleaning up 'instance'
==> alicloud-ecs: Cleaning up 'security group'
==> alicloud-ecs: Cleaning up 'vSwitch'
==> alicloud-ecs: Cleaning up 'VPC'
```

```
==> alicloud-ecs: Deleting temporary keypair...
Build 'alicloud-ecs' finished.
==> Builds finished. The artifacts of successful builds are:
--> alicloud-ecs: Alicloud images were created:
cn-beijing: m-bp67acfmxazb4ph***
```

Related information

References

- Describelmages
- Alicloud Image Builder
- Examples
- Alicloud Image Builder parameters used to implement DevOps

14.3.2. Alicloud Image Builder parameters used to implement DevOps

This topic describes the Alicloud Image Builder parameters that are used to implement DevOps when you use Packer to create Alibaba Cloud ECS custom images.

Parameter used to tag custom images

- Field: tags{"key":"value"}.
- Applicable scenario: If you have multiple custom images, you can tag them for easy management and retrieval. Alicloud Image Builder provides the tags parameter. If you set this parameter when you use Packer to create a custom image, the generated image will be bound with the tag you specified. For more information, see Tag overview.
- Function: When you query images in the ECS console or by calling the Describelmages operation, images are displayed with their tags. You can also filter images by tag. Tags bound to images can be used together with Terraform to standardize enterprise-grade DevOps processes.
- Example: The following configuration file contains the version=v1.0.0 and app=web tags bound to the generated image and relative snapshot:

```
"variables": {
 "access key": "{{env `ALICLOUD ACCESS KEY`}}",
 "secret key": "{{env `ALICLOUD SECRET KEY`}}"
},
"builders": [{
 "type": "alicloud-ecs",
 "access key":"{{user `access key`}}",
 "secret key":"{{user `secret key`}}",
 "region": "cn-beijing",
 "image name": "packer basic",
 "source image": "centos 7 03 64 20G alibase 20170818.vhd",
  "ssh username": "root",
 "instance type": "ecs.t5-lc1m1.small",
 "internet charge type": "PayByTraffic",
 "io optimized":"true",
  "tags": {
   "version": "v1.0.0",
   "app": "web"
 }
} ]
```

Parameter used to control whether to create an image only based on a system disk

- Field: image_ignore_data_disks. Data type: boolean.
- Application scenario: By default, Packer creates images directly from ECS instances. If the ECS
 instances have data disks, the generated images contain data disk snapshots. You can use one of the
 following methods to create an instance that has data disks:
 - Method 1: Set data disk parameters in image_disk_mappings. For more information, see <u>Alicloud Image Builder</u> in *Packer documentation*.
 - Method 2: Select an instance type that comes with data disks by default, such as
 ecs.d1ne.2xlarge. The data disks that instance types come with by default are typically local disks.
 Snapshots cannot be created for local disks. Therefore, you cannot create images directly from
 instances of such instance types.
- Function: If you want to select an instance type that comes with data disks by default, but you do
 not want the image created to contain data disk snapshots, you can add image_ignore_data_disks
 ": "true" in the configuration file. This ensures that the image is created only based on the system disk.

Parameter used to set the snapshot timeout period

- Field: wait_snapshot_ready_timeout. Data type: integer. Default value: 3600. Unit: seconds.
- Applicable scenario: Images are created based on snapshots. The time it takes to create a snapshot for a disk depends on the disk size. The larger a disk is, the more time it takes.
- Function: If a timeout error occurs when you create a snapshot for a large disk, specify a larger value for wait_snapshot_ready_timeout to prolong the snapshot timeout period.

Parameter used to control whether to connect to an instance by using a private IP address

- Field: ssh_private_ip. Data type: boolean.
- Applicable scenario: By default, Packer creates an Elastic IP address (EIP) and associates this EIP with
 the temporary ECS instance when Packer creates a custom image. Then, Packer uses the public IP
 address that corresponds to the EIP to connect to the instance and then installs software or runs
 commands on the instance. If Packer can use private IP addresses to connect to the instance, the
 public IP address is not needed.
- Function: If you set ssh_private_ip to true, Packer does not assign an EIP or a public IP address to the temporary instance but uses private IP addresses to connect to the instance.

Parameter used to control whether to stop an instance

- Field: disable_stop_instance. Data type: boolean.
- Applicable scenario: By default, after Packer runs provisioners, it stops instances and then creates images from the instances. However, in some scenarios (for example, when you run Sysprep in Windows instances), instances must be in the Running state.

For information about the applicable scenarios of Sysprep, see Modify the SID of a Windows ECS instance to create a domain environment.

• Function: After you set disable_stop_instance to true, Packer does not stop instances but assumes that the command provided in the configuration (provisioners) automatically stops the instances.

Parameter used to specify the UserData file path for enabling WinRM

- Field: user_data_file.
- Applicable scenario: For security purposes, the Windows Remote Management (WinRM) function is disabled in Windows images by default. However, Packer must use the WinRM function to connect to a Windows instance and run commands on the instance. You can use the UserData file to enable WinRM when you create a Windows instance.
- Function: You can use the user_data_file parameter to specify the path to the UserData file. For example, you can set the value of user_data_file to examples.ps1.
- Example: In the following sample code, the UserData file is located in the relative path examples/alicloud/basic/winrm_enable_userdata.ps1.

```
"variables": {
 "access key": "{{env `ALICLOUD ACCESS KEY`}}",
  "secret key": "{{env `ALICLOUD SECRET KEY`}}"
},
"builders": [{
 "type": "alicloud-ecs",
  "access key":"{{user `access key`}}",
  "secret key":"{{user `secret key`}}",
 "region": "cn-beijing",
 "image name": "packer test",
  "source image": "win2008r2 64 ent sp1 zh-cn 40G alibase 20181220.vhd",
  "instance type": "ecs.nl.tiny",
 "io optimized": "true",
 "internet charge type": "PayByTraffic",
 "image_force_delete":"true",
  "communicator": "winrm",
 "winrm_port": 5985,
 "winrm username": "Administrator",
  "winrm password": "Test1234",
  "user data file": "examples/alicloud/basic/winrm enable userdata.ps1"
} ],
"provisioners": [{
 "type": "powershell",
  "inline": ["dir c:\\"]
} ]
```

? Note

- o In the preceding sample code:
 - "communicator": "winrm" indicates that you are connected to the instance through WinRM.
 - "winrm port": 5985 indicates that the communication port is port 5985.
 - "winrm_username": "Administrator" indicates that you are connected to the instance as an administrator.
 - "winrm password": "Test1234" indicates that Password Test1234 is used.
- "image_force_delete":"true" indicates that existing images are deleted if they have the same name as the image to be created.

Parameters used to create an image based on an on-premises ISO file and import the image to Alibaba Cloud ECS

- Fields: builders{"type":"qemu"} and post-processors{"type":"alicloud-import"}.
- Applicable scenario: If an on-premises ISO file runs in a non-QEMU-based virtualization environment, you can use Packer to create an image based on the file and import the image to Alibaba Cloud ECS.
- Example: If the on-premises environment is based on QEMU, you can use Packer to create an image and then import the image to Alibaba Cloud ECS. For more information, see Create and import onpremises images by using Packer, which includes two important steps:

- i. Use an on-premises virtualization environment or a builder, such as QEMU Builder, to create an on-premises image.
- ii. Define Alicloud Import Post-Processor to import the generated on-premises image to Alibaba Cloud ECS.

To import an ISO file to Alibaba Cloud ECS, you must first install a virtualization environment on premises and then create an image based on the file. The image must be in a format supported by Alibaba Cloud, such as the QCOW2, VHD, or RAW format. Then, you can import the image to Alibaba Cloud ECS. For more information, see Instructions for importing images.

References

Alicloud Image Builder and Examples.

14.4. Create and import a custom image

If you cannot find the operating system that you need in the Elastic Compute Service (ECS) console when you create an instance, you can create a custom image and import it to the console to create an instance. This topic describes how to create a custom image and import it to Alibaba Cloud.

Prerequisites

- VirtualBox is installed. If not, download it from the VirtualBox official website.
- The network connection is stable.

•

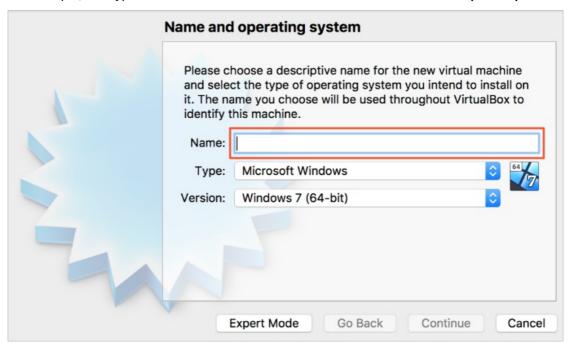
190

- The ISO or binary file of the operating system is installed. Example: Red Hat Enterprise Linux.
- ossbrowser is installed. For more information, see Use ossbrowser.

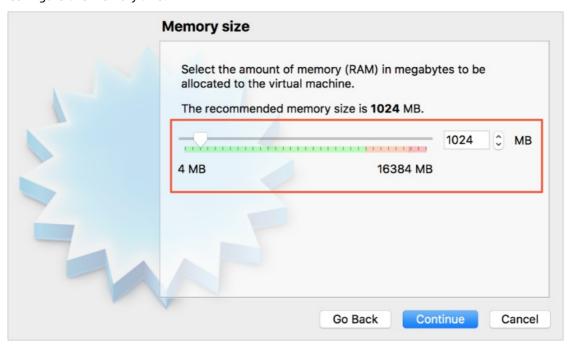
Procedure

- 1. Create a new virtual machine (VM) on VirtualBox.
 - i. Start VirtualBox and click New.

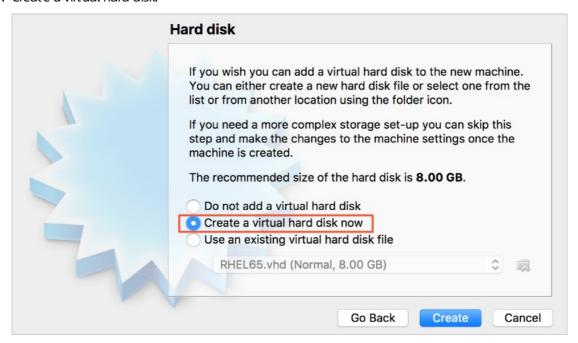
ii. Enter the name of the new VM and select the corresponding operating system and its version. For example, set Type to Microsoft Windows and set Version to Windows 7 (64-bit).



iii. Configure the memory size.



iv. Create a virtual hard disk.



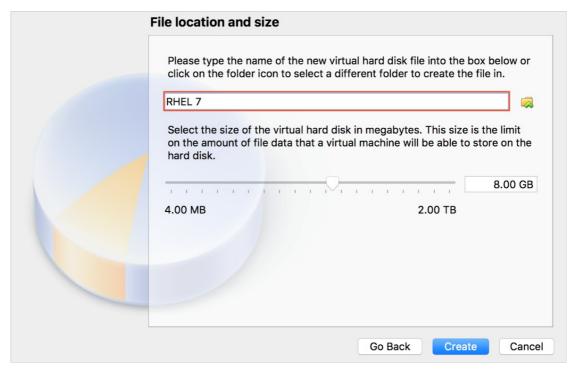
v. Select VHD (Virtual Hard Disk) as the file type of the hard disk. Alibaba Cloud supports the RAW, VHD, and qcow2 formats.



vi. Select Dynamically allocated as the storage type of the hard disk.



vii. Enter the name of the virtual hard disk and click Create.



viii. After the VHD file is created, double-click the new VM to start it.



2. Run the following command to install Kernel-based Virtual Machine (KVM):

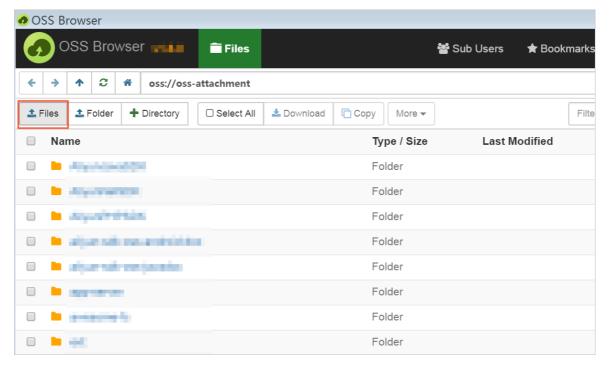
yum install qemu-kvm qemu-img libvirt

3. Run the following command to disable the firewall of the VM:

service firewalld stop

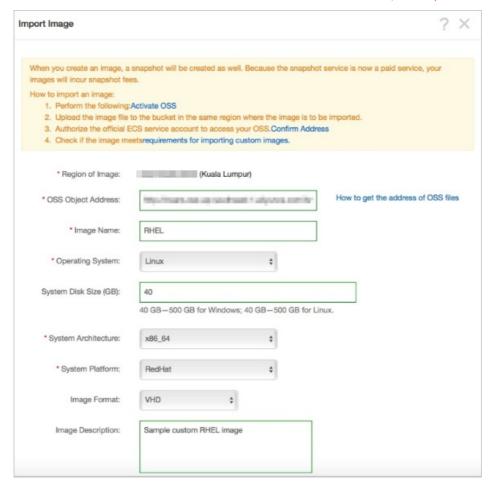
4. Upload the VHD file to Object Storage Service (OSS) and use ossbrowser to upload the file to the region where you want to create an instance.

For more information, see Use ossbrowser.



- 5. Import the custom image.
 - i. Log on to the ECS console.
 - ii. In the left-side navigation pane, choose Instances & Images > Images.
 - iii. In the top navigation bar, select the region where the VHD file is uploaded and click **Import Image**.

For information about how to allow ECS to access OSS resources, see Import custom images.



iv. Configure the parameters and click **OK**.

Note You can log on to the OSS console to obtain the URL of the OSS object that corresponds to the uploaded image file. For more information, see Download objects.

After the custom image is imported, it is displayed in the image list.



15.Monitor

15.1. Use CloudMonitor to monitor websites deployed on ECS instances

Cloud Monitor helps reduce the O&M costs and workloads of cloud services. Cloud Monitor provides real-time operation data that you can use to identify risks in advance, troubleshoot issues, and prevent potential loss. When issues occur, Cloud Monitor immediately sends you an alert notification to help you restore business in a quick manner.

Prerequisites

Before you use CloudMonitor, make sure that the following requirements are met:

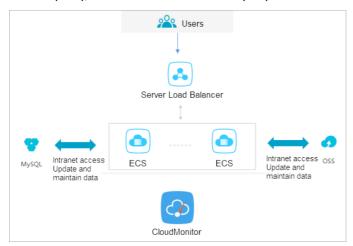
- The CloudMonitor agent is running on the Elastic Compute Service (ECS) instances that you want to monitor and is able to collect metric data. If the CloudMonitor agent is not installed on the instances, manually install it. For more information, see Install the CloudMonitor agent.
- Alert contacts and contact groups are added. We recommend that you add at least two contacts to
 ensure real-time responses to monitoring alerts. For more information about monitoring metrics, see
 Alert service and Overview.

Context

The dashboard feature of CloudMonitor provides system-wide visibility into resource utilization and operational health. In this topic, the CPU utilization, memory usage, and disk usage of ECS instances are separately displayed, and the four metrics of ApsaraDB RDS instances are displayed in two groups.



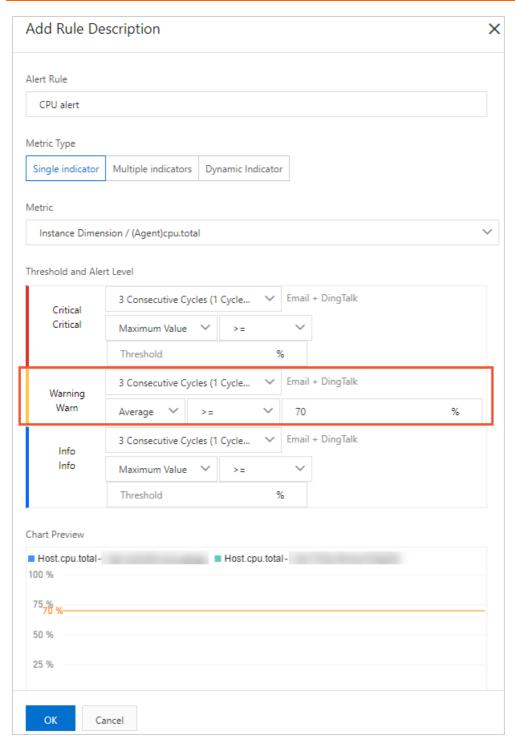
In this topic, a website is used to describe how to configure CloudMonitor. ECS, RDS, Object Storage Service (OSS), and Server Load Balancer (SLB) are used.



Configure alert thresholds and alert rules

We recommend that you configure suitable alert thresholds for monitoring metrics based on your business requirements. A low threshold may lead to frequent triggering of alerts and affect user experience. A high threshold may leave you with insufficient time to respond to events.

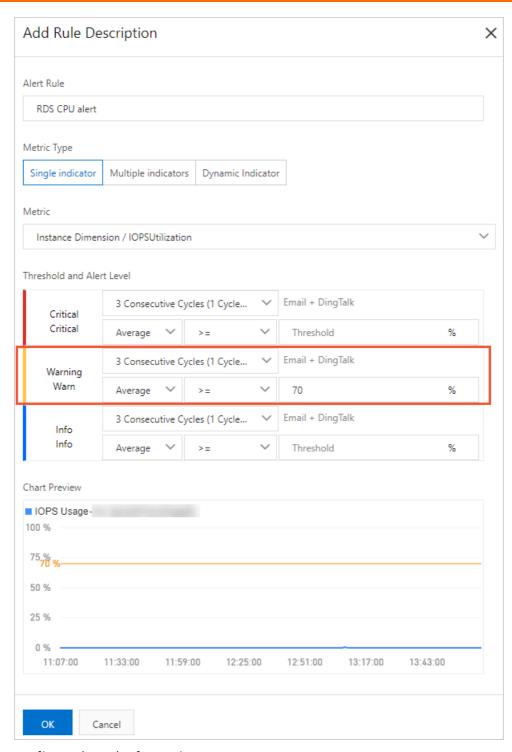
For example, to reserve some processing capacity to ensure the normal operation of the system, you can set the alert threshold for CPU utilization to 70% and set an alert to be triggered when the threshold is exceeded three consecutive times, as shown in the following figure.



If you want to configure alert rules for other metrics, click **Add Alert Rule**. For example, you can perform the following operations:

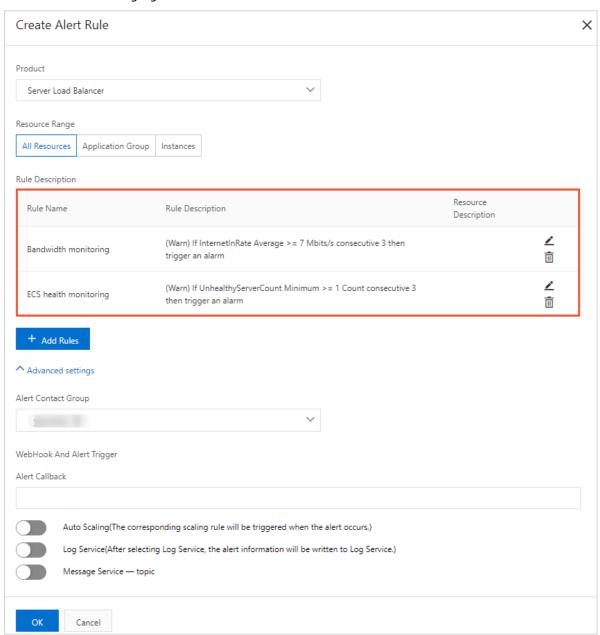
• Configure alert rules for RDS instances

We recommend that you configure alert rules for RDS instances based on your requirements. For example, you can set the alert threshold for the CPU utilization of RDS instances to 70% and set an alert to be triggered when the threshold is exceeded three consecutive times. You can configure alert thresholds for the disk usage, IOPS utilization, and total number of connections based on your requirements. For information about how to view the information about monitoring metrics, see Cloud Service Monitoring.



• Configure alert rules for SLB instances

Before you use CloudMonitor for SLB instances, make sure that health check is enabled for your SLB instances. You can set the alert threshold for the bandwidth value of SLB instances to 7 Mbit/s, as shown in the following figure.



Configure process monitoring

For web applications, you can configure process monitoring to monitor application processes in real time and use monitoring data to troubleshoot issues. For more information, see Configure process monitoring.

Configure site monitoring

Site monitoring is an external monitoring service for ECS instances and is used to simulate real user access scenarios and test the business availability in real time. The monitoring data can also be used to troubleshoot issues.



If the preceding monitoring metrics do not meet your requirements, you can use the custom monitoring feature. For more information, see Overview.

16.Use RAM roles to access other Alibaba Cloud services

This topic describes how to enable applications on ECS instances to access other Alibaba Cloud services by using STS temporary credentials through RAM roles. The examples show how to enable Python to access Object Storage Service (OSS).

Prerequisites

An instance is created. For more information, see Create an instance by using the wizard.

Context

Previously, applications deployed on an ECS instance needed to use AccessKey pairs (AKs) to access other Alibaba Cloud services. An AK allows you to access Alibaba Cloud APIs with full permissions for your account. To facilitate the management of the AK by applications, you must store the AK in the application configuration files or otherwise store the AK in an ECS instance. These operations makes it more complicated to manage the AK and keep it confidential. If you need consistent deployment across regions, the AK is spread along with the images or instances created from the images. In these cases, you must update and redeploy each instance and image individually whenever you make changes to the AK.

You can attach a RAM role to an ECS instance, and use an STS temporary credential to access other cloud services from the applications within the instance. STS temporary credentials are generated and updated automatically. Applications can obtain the STS temporary credentials by using the instance metadata URL. You can use RAM roles and authorization policies to grant ECS instances with different or identical permissions to access other cloud services.

Notice All operations in this topic are performed in OpenAPI Explorer to help you get started with the examples. OpenAPI Explorer obtains a temporary AK of the current account through the information of the logged user, and initiates online resource operations on the current account. You will be charged for creating an instance. Release the instance after the operations are complete.

Procedure

To enable Python to access OSS by using the STS temporary credential through RAM roles, perform the following operations:

- 1. Step 1. Create a RAM role and configure an authorization policy.
- 2. Step 2. Create an ECS instance and attach the RAM role to the instance.
- 3. Step 3. Access the instance metadata URL within the instance to obtain the STS temporary credential.
- 4. Step 4. Use SDK for Python to access OSS by using the STS temporary credential.

Step 1. Create a RAM role and configure an authorization policy.

Perform the following operations to create a RAM role and configure an authorization policy:

1. Create a RAM role.

Call the CreateRole operation to create a RAM role.

- RoleName: Set the name of the RAM role. Set the parameter based on your needs. *EcsRamRoleT est* is used in this example.
- AssumeRolePolicyDocument: Specify a policy as follows to set the created role as a service role and authorize an Alibaba Cloud service (ECS in this example) to assume the role.

2. Create an authorization policy.

Call the CreatePolicy operation to create an authorization policy.

- PolicyName: the name of the authorization policy. EcsRamRolePolicyTest is used in this example.
- PolicyDocument: the content of the authorization policy. The following content is used in the example, indicating that the role has the OSS read-only permission.

- 3. Grant permissions to the role. Call the AttachPolicyToRole operation to grant permissions to the role.
 - PolicyType: Set the parameter to Custom.
 - PolicyName: Set the parameter to the name of the authorization policy created in the second step. *EcsRamRolePolicyTest* is used in this example.
 - RoleName: Set the parameter to the name of the role created in the first step. *EcsRamRoleTest* is used in this example.

Step 2. Create an ECS instance and attach the RAM role to the instance.

You can use one of the following methods to create an ECS instance and attach the RAM role to the instance:

• Attach the RAM role to an existing ECS.

Call the AttachinstanceRamRole operation to attach a RAM role to an existing VPC-type instance. The parameters are configured as follows:

- RegionId: the region ID of the ECS instance.
- RamRoleName: the name of the RAM role. *EcsRamRoleTest* is used in this example.
- InstanceIds: the IDs of VPC-type instances to which you want to attach the RAM role. ["i-bXXXXXXX XX"] is used in this example.
- Create an ECS instance and attach the RAM role to the instance.

To create an ECS instance and attach the RAM role to the instance, perform the following steps:

i. Create an instance.

Call the CreateInstance operation and set parameters based on your actual needs. The following parameters are required:

- RegionId: the region ID of the instance. *cn-hangzhou* is used in this example.
- ImageId: the image of the instance. *centos_7_03_64_40G_alibase_20170503.vhd* is used in this example.
- InstanceType: the instance type of the instance. *ecs.q6.large* is used in this example.
- VSwitchId: the ID of the vSwitch in the VPC to which the instance belongs. RAM roles can be attached to only VPC-type ECS instances. Therefore, the vSwitchId parameter is required.
- RamRoleName: the name of the RAM role. *EcsRamRoleTest* is used in this example.

If you want to authorize a RAM user to create an ECS instance with the specified RAM role attached, in addition to the permission to create an ECS instance, the RAM user must have the PassRole permission. Therefore, you must customize an authorization policy as follows and attach it to the RAM user.

- If you want to configure the RAM user to create an ECS instance, [ECS RAM Action] filed can be replaced with ecs:CreateInstance also grant more permissions to the RAM user based on your actual needs.
- If you want to grant all ECS permissions to the RAM user, [ECS RAM Action] must be replaced with ecs:*

Note For more information about values of [ECS RAM Action], see Authentication rules.

- ii. Configure the password and start the instance.
- iii. Configure the instance to access the Internet in the ECS console or by calling an API operation.

Step 3. Access the instance metadata URL within the instance to obtain the STS temporary credential.

To obtain the STS temporary credential of the instance, perform the following steps.

Notice A new STS temporary credential is generated 30 minutes before the current one expires. Both STS credentials can be used during this period of time.

1.

2. Access http://100.100.100.200/latest/meta-data/ram/security-credentials/EcsRamRoleTest to obtain the STS temporary credential. The last part of the URL is the RAM role name, which must be replaced with the name you specified.

```
Note In this example, the curl command is used to access the preceding URL. If your instance is a Windows instance, see Metadata.
```

The following content shows the sample response:

```
[root@local ~]# curl http://100.100.200/latest/meta-data/ram/security-credentials/E
csRamRoleTest
{
"AccessKeyId" : "STS.J8XXXXXXXXX4",
"AccessKeySecret" : "9PjfXXXXXXXXXXBf2XAW",
"Expiration" : "2017-06-09T09:17:19Z",
"SecurityToken" : "CAIXXXXXXXXXXXXWmBkleCTkyI+",
"LastUpdated" : "2017-06-09T03:17:18Z",
"Code" : "Success"
}
```

Step 4. Use SDK for Python to access OSS by using the STS temporary credential.

In this example, SDK for Python is used to list 10 files in a specified OSS bucket located within the same region as the instance through the STS temporary credential.

Prerequisites:

- The ECS instance is connected.
- The ECS instance is installed with Python. If your instance is a Linux instance, pip is required.
- An OSS bucket is created in the same region as the instance and the name and endpoint of the bucket are obtained. In this example, the bucket name is ramroletest , and the endpoint is oss-c n-hangzhou.aliyuncs.com .

Perform the following steps to use SDK for Python to access OSS:

- 1. Run the pip install oss2 command to install OSS SDK for Python.
- 2. Run the following commands to test whether SDK for Python can be used to list 10 files in the specified OSS bucket:

```
import oss2
from itertools import islice
auth = oss2.StsAuth(<AccessKeyId>, <AccessKeySecret>, <SecurityToken>)
bucket = oss2.Bucket(auth, <Your endpoint>, <Your bucket name>)
for b in islice(oss2.ObjectIterator(bucket), 10):
    print(b.key)
```

where:

- The three parameters of oss2.StsAuth correspond to the AccessKeyId, AccessKeySecret, and SecurityToken values returned in the Step 3. Access the instance metadata URL within the instance to obtain the STS credential section.
- The last two parameters of oss2.Bucket correspond to the name and endpoint of the bucket.

The following content shows the sample response:

17.Networks

17.1. Best practices for configuring public bandwidth

Alibaba Cloud Elastic Compute Service (ECS) provides two billing methods for network usage: pay-by-bandwidth and pay-by-traffic. If the bandwidth needs of your instance fluctuate, you can change its billing method for network usage from pay-by-bandwidth to pay-by-traffic and enforce an upper limit on bandwidth to suit your business needs and reduce costs. This topic describes the billing methods for network usage of ECS instances, how to configure public bandwidth in different scenarios, and how to change billing methods for network usage to increase the upper bandwidth limit.

Prerequisites

Before you change the billing method for network usage of a subscription instance from pay-by-bandwidth to pay-by-traffic, make sure that your account has been granted the privilege to perform configuration downgrades.

Note You can click **Privileges** on the **Overview** page in the ECS console to go to the Privileges and Quotas page and check whether your account has been granted the privilege to perform configuration downgrades.

Context

The following table describes the pay-by-bandwidth and pay-by-traffic billing methods for network usage and the scenarios to which these billing methods are applicable.

Scenario	Billing method for network usage	Bandwidth type	Billing
Scenarios that require stable low-cost bandwidth	Pay-by- bandwidth	For an instance that uses this billing method for network usage, you must specify an outbound bandwidth in Mbit/s. The actual outbound bandwidth used by the instance is capped at your specified bandwidth.	You are charged for the specified bandwidth based on a tiered billing model. The bandwidth fees are built into your instance fees. For more information, see Public bandwidth.

Scenario	Billing method for network usage	Bandwidth type	Billing	
Scenarios that have highly variable bandwidth needs, such as those which have low traffic usage with occasional traffic spikes	Pay-by- traffic	For an instance that uses this billing method for network usage, you must specify a maximum outbound bandwidth (Peak Bandwidth) of up to 100 Mbit/s. To prevent out-of-control fees caused by bursts in traffic, specify a maximum bandwidth for the instance. The actual outbound bandwidth that can be used by the instance is capped at your specified maximum bandwidth.	You are charged for the actual volume (in GB) of outbound traffic to the Internet on a pay-as-you-go basis. These bills are generated every hour on	
		② Note Public bandwidth limits are determined based on instance billing methods and billing methods for network usage. For more information, see 使用限制.	the hour. For more information, see Public bandwidth.	

The following table describes how to configure public bandwidth in different scenarios.

Scenario	Configuration method	Procedure
A public IP address is required to access the Internet.	Assign a public IP address and configure public bandwidth.	 You can assign a public IP address and configure public bandwidth when you create an instance. For more information, see Step 1: Assign a public IP address and configure public bandwidth when you create an ECS instance. If an instance is created but not assigned a public IP address, you can change its public bandwidth configurations to assign a public IP address to the instance. For more information, see Overview of instance configuration changes.

Scenario	Configuration method	Procedure
The current billing method for network usage does not meet your business needs.	Change the billing method for network usage.	For more information, see Change the billing method for network usage.
The specified bandwidth does not suit your business needs.	Modify the public bandwidth configurations.	 For information about how to modify the public bandwidth configurations of a subscription instance, see Modify the bandwidth configurations of subscription instances. For information about how to modify the public bandwidth configurations of a pay-as-you-go instance, see Modify the bandwidth configurations of pay-as-you-go instances.
The bandwidth of an elastic IP address (EIP) does not suit your business needs.	Change the bandwidth of the EIP.	For more information, see Modify the bandwidth of an EIP.

For information about some frequently asked questions about public bandwidth, see Network FAQ.

The following sections describe how to configure public bandwidth when you create an instance and how to change the billing method for network usage to enforce an upper bandwidth limit.

Step 1: Assign a public IP address and configure public bandwidth when you create an ECS instance

Create an ECS instance, assign a public IP address, and configure the pay-by-bandwidth billing method for network usage.

- 1.
- 2.
- 3.
- 4. On the Instances page, click Create Instance.
- 5. Configure parameters in the instance creation wizard.

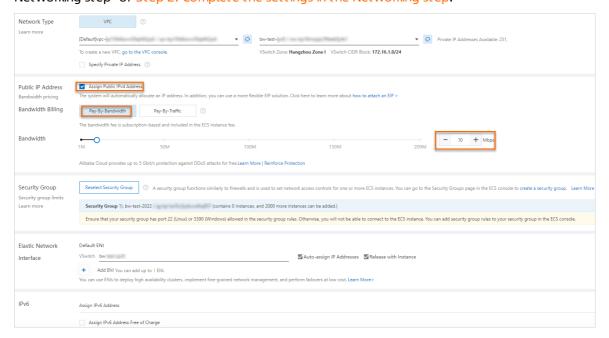
For information about the steps and parameters in the wizard, see Create an instance by using the wizard.

In the Networking step, complete the following settings:

- Network Type: Select VPC.
 - Select an existing virtual private cloud (VPC) and an existing vSwitch from the drop-down lists.
 - You can select **Specify Private IP Address** to specify a private IP address.
- Public IP Address: Select Assign Public IPv4 Address.
- o Bandwidth Billing: Select Pay-By-Bandwidth.
- o Bandwidth: Specify a bandwidth. Unit: Mbit/s. Example: 10. You are charged based on your

specified bandwidth. The actual outbound bandwidth is capped at your specified bandwidth.

- Security Group: Select a security group to control inbound and outbound traffic for the instance.
- ENI: Primary elastic network interfaces (ENIs) cannot be unbound from instances but can be created and released only along with instances. To create a secondary ENI when you create an instance, click the → icon and select a vSwitch with which to associate the secondary ENI.
- o (Optional)IPv6: Assign an IPv6 address to the instance based on your business needs. For more information about these parameters, see the "Step 2: Complete the settings in the Networking step" of Step 2: Complete the settings in the Networking step.



After the instance is created, the specified bandwidth is displayed in the **Specifications** column corresponding to the instance on the Instances page.



Step 2: Change the billing method for network usage from pay-by-bandwidth to pay-by-traffic

If your instance has fluctuating bandwidth needs and requires an upper bandwidth limit, you can change the billing method for network usage of the instance to pay-by-traffic and specify a maximum bandwidth to prevent out-of-control fees caused by bursts in traffic.

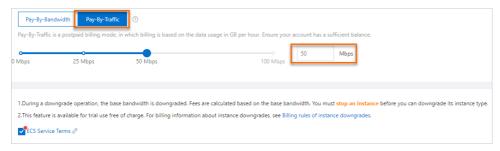
- 1.
- 2.
- 3.
- 4. Find the instance for which you want to change the billing method for network usage. Then, perform different operations to go to the configuration page based on the billing method of the instance.
 - o Subscription

- a. Find the subscription instance for which you want to change the billing method for network usage. Click **Upgrade/Downgrade** in the **Actions** column.
- b. In the dialog box that appears, choose **Downgrade > Bandwidth Configuration**.
- c. Click Continue.
- o Pay-as-you-go

Find the pay-as-you-go instance for which you want to change the billing method for network usage and choose More > Configuration Change > Change Pay-as-you-go Instance Bandwidth in the Actions column.

5. On the Change Bandwidth page, find the Bandwidth section, select **Pay-By-Traffic**, and then specify a maximum bandwidth. Unit: Mbit/s. Example: 50.

You are charged based on the actual traffic volume. The actual outbound bandwidth is capped at the specified maximum bandwidth to prevent out-of-control fees caused by bursts in traffic.



- 6. Read the notes. Read and select ECS Service Terms.
- 7. Confirm the configuration costs, click Confirm in the lower part of the page, and then perform the subsequent operations as instructed on the page.

After the billing method for network usage of the instance is changed, the new configurations take effect immediately. The specified maximum bandwidth is displayed in the **Specifications** column corresponding to the instance on the Instances page. Example: 50 Mbit/s.



Note For information about how to change the billing method for network usage of an instance associated with an EIP, see Modify the bandwidth of an EIP.

17.2. Best practices for testing network performance

This topic describes how to test the packet forwarding rate, network bandwidth, and network latency of an Elastic Compute Service (ECS) instance. In this topic, Netperf and sockperf are used to test the network performance of ECS instances.

Prerequisites

- ECS instances that meet requirements for the test environment are created. For more information about how to create an ECS instance, see Create an instance by using the wizard.
- In the same test environment, all ECS instances reside within the same virtual private cloud (VPC), are connected to the same vSwitch, and belong to the same security group.

Notice

- To prevent data loss, we recommend that you use a tool to test the network performance of new ECS instances that contain no data.
- Instance type specifications are all verified and obtained within a test environment. In actual scenarios, the performance of an instance may vary based on other factors such as instance load and networking model.

Prepare the environment

In this test, ECS instances whose network performance is tested are used as test machines and auxiliary test machines. The packet forwarding rate, network bandwidth, and network latency of these ECS instance are tested. Test machines can be used as clients or servers in tests that use Netperf or sockperf. Auxiliary test machines can also be used as clients or servers in tests that use Netperf or sockperf. Test machines are connected to auxiliary test machines to transmit test configurations.

The following table describes the sample instance types, images, and recommended quantities of instances in different test environments.

Note The packet forwarding rate is not a factor to consider when you test the network bandwidth and network latency of an ECS instance. You can select any instance type for the test.

• Test the packet forwarding rate (less than 6 million pps), network bandwidth, and network latency of an ECS instance

Test item	Test machine	Auxiliary test machine
Instance type	ecs.g7.large	ecs.g7.large
lmage	Alibaba Cloud Linux 2.1903 LTS 64-bit	Alibaba Cloud Linux 2.1903 LTS 64-bit
Quantity	1	1

- For more information about how to test the packet forwarding rate (less than 6 million pps), see the Test the packet forwarding rate (less than 6 million pps) section of this topic.
- For more information about how to test the network bandwidth, see the Test the network bandwidth section of this topic.
- For more information about how to test the network latency, see the Test the network latency section of this topic.
- Test the packet forwarding rate (greater than 6 million pps and less than 20 million pps) of an ECS instance

Test item	Test machine	Auxiliary test machine
Instance type	ecs.g7.16xlarge	ecs.g7.16xlarge
lmage	Alibaba Cloud Linux 2.1903 LTS 64-bit	Alibaba Cloud Linux 2.1903 LTS 64-bit
Quantity	1	3

For more information about how to test the packet forwarding rate (greater than 6 million pps and less than 20 million pps), see the Test the packet forwarding rate (greater than 6 million pps and less than 20 million pps) section of this topic.

• Test the packet forwarding rate (greater than 20 million pps) of an ECS instance

Test item	Test machine	Auxiliary test machine
Instance type	ecs.g7.32xlarge	ecs.g7.32xlarge
lmage	Alibaba Cloud Linux 2.1903 LTS 64-bit	Alibaba Cloud Linux 2.1903 LTS 64-bit
Quantity	1	3

For more information about how to test the packet forwarding rate (greater than 20 million pps), see the Test the packet forwarding rate (greater than 20 million pps) section of this topic.

Test the packet forwarding rate (less than 6 million pps)

- Connect to the test machine and the auxiliary test machine.
 For more information about the methods for connecting to an instance, see Connection methods.
- 2. Run the following command on the test machine and the auxiliary test machine to download Netperf:

```
wget https://benchmark-packages.oss-cn-qingdao.aliyuncs.com/netperf-2.7.0.tar.gz
```

- 3. Run the following command on the test machine and the auxiliary test machine to install Netperf and SAR.
 - i. Run the following command to decompress the Netperf installation package:

```
yum install -y gcc autoconf automake libtool sysstat
tar -zxvf netperf-2.7.0.tar.gz
```

ii. Run the following command to query the version number of GCC:

```
gcc -v 2>&1
```

iii. (Optional)If the GCC version of the test machine and the auxiliary test machine is later than 10, you must run the following command to open the *nettest_omni.c* file and manually delete the declared variables:

```
cd netperf
vim src/nettest_omni.c
```

Manually delete the following declared variables in the *nettest_omni.c* file.

After you delete the previous variables, press the Esc key. Then, enter : wq and press the Enter key to save the modification and exit.

iv. Run the following command to install Netperf and SAR:

```
cd netperf
./configure
make && make install
```

4. Run the following command on the test machine to start 64 netserver services:

```
#!/bin/bash
for j in `seq 64`; do
    netserver -p $[16000+j] > server_$[16000+j].netperf > /dev/null 2>&1 &
done
```

5. Run the following command on the test machine to query the private IP address of the test machine:

```
ifconfig || ip addr
```

```
[root@launch-0322 netperf]# ifconfig || ip addr
eth0: flags=4163<UP.BROADCAST,RUNNING,MULTICAST> mtu 1500
inet 172.16......3 netmask 255.255.240.0 broadcast 172.16.
                      la prefixlen 64 scopeid 0x20<link>
       inet6
       RX errors 0 dropped 0 overruns 0 frame 0
       TX packets 44604 bytes 6561192 (6.2 MiB)
       TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
lo: flags=73<UP,LOOPBACK,RUNNING> mtu 65536
       inet 127.0.0.1 netmask 255.0.0.0
       inet6 ::1 prefixlen 128 scopeid 0x10<host>
       loop txqueuelen 1000 (Local Loopback)
       RX packets 0 bytes 0 (0.0 B)
       RX errors 0 dropped 0 overruns 0 frame 0
       TX packets 0 bytes 0 (0.0 B)
       TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
```

6. Run the following command on the auxiliary test machine to send traffic to the test machine:

```
#!/bin/bash
server_ip=<Private IP address of the test machine>
for j in `seq 64`; do
    port=$[16000+j]
    netperf -H ${server_ip} -l ${run_time:-300} -t UDP_STREAM -p $port -- -m 1 -D > /
dev/null 2>&1 &
done
```

Replace <Private IP address of the test machine> with the private IP address of the test machine obtained in the previous step, as shown in the following figure.

```
netperf]# #!/bin/bash
[root@i
[root@i
                             netperf]# server_ip=172.1
[root@i
                             netperf]#
[root@i
                             netperf]# for j in `seq 64`; do
     port=$[16000+j]
     netperf -H ${server_ip} -l ${run_time:-300} -t UDP_STREAM -p $port -- -m 1 -D > /dev/null 2>81 &
> done
[65] 5819
[66] 5820
[67] 5821
[68] 5822
[69] 5823
```

7. Run the following command on the test machine to test the network traffic:

```
sar -n DEV 1
```

View the values in the rxpck/s column in the test result. The rxpck/s column shows the total number of packets received by the test machine per second. In this example, the number of packets received by the test machine per second is about 900,000, as shown in the following figure.

```
[root@launch-0322 netperf]# sar -n DEV 1
                          64 (launch-0322)
                                                 03/24/2022
                                                                 _x86_64_
                                                                                  (2 CPU)
10:51:19 AM
                        rxpck/s
                                              rxkB/s
                                                        txkB/s
                                                                           txcmp/s rxmcst/s
                IFACE
                                   txpck/s
                                                                 rxcmp/s
10:51:20 AM
                  10
                           0.00
                                      0.00
                                              0.00
                                                          0.00
                                                                    0.00
                                                                              0.00
                                                                                         0.00
10:51:20 AM
                                                                              0.00
                                                                                        0.00
                                                          2.86
                 eth0 898970.00
                                      7.00 37750.35
                                                                    0.00
                        rxpck/s
                                   txpck/s
                                                                 rxcmp/s
                                                                           txcmp/s
                                              rxkB/s
                                                        txkB/s
                                                                                     rxmcst/s
10:51:20 AM
                IFACE
10:51:21 AM
                   lo
                           0.00
                                                                              0.00
                                     0.00
                                               0.00
                                                          0.00
                                                                    0.00
                                                                                         0.00
10:51:21 AM
                 eth0 899982.00
                                      9.00 37792.87
                                                          1.29
                                                                    0.00
                                                                              0.00
                                                                                         0.00
10:51:21 AM
                IFACE
                        rxpck/s
                                   txpck/s
                                              rxkB/s
                                                        txkB/s
                                                                 rxcmp/s
                                                                           txcmp/s rxmcst/s
10:51:22 AM
                   la
                           0.00
                                     0.00
                                                0.00
                                                          0.00
                                                                    0.00
                                                                              0.00
                                                                                         0.00
10:51:22 AM
                 eth0 899955.00
                                     8.00
                                           37791.75
                                                          1.22
                                                                    0.00
                                                                              0.00
                                                                                         0.00
```

Test the packet forwarding rate (greater than 6 million pps and less than 20 million pps)

- 1. Connect to the test machine and the auxiliary test machine.

 For more information about the methods for connecting to an instance, see Connection methods.
- 2. Run the following command on three auxiliary test machines and one test machine to install sockperf:

```
yum install -y sockperf
```

If yum cannot be used to install sockperf, run the following command to install sockperf by using a compiler:

```
yum install -y autoconf automake libtool g++ gcc-c++
cd /opt
wget https://github.com/Mellanox/sockperf/archive/3.6.tar.gz
tar -zxf 3.6.tar.gz
cd sockperf-3.6/
./autogen.sh
./configure
make -j `cat /proc/cpuinfo| grep process | wc -l`
make install
```

3. Run the following command on three auxiliary test machines to send traffic to the test machine:

```
server_ip=<Private IP address of the test machine>
threads=64
msg_size=14
run_time=60
basePort=6666
for((i=0;i<$threads;++i));do
    nohup sockperf tp -i $server_ip --client_port $[${basePort}+${i}] --pps max -m ${ms}
g_size} -t ${run_time} --port $[${basePort}+${i}] 2>&1 &
done
```

Replace <Private IP address of the test machine> with the private IP address of the actual test machine. run time specifies the time period during which traffic is sent.

4. Run the following command on three auxiliary test machines and one test machine to test the network traffic:

```
sar -n DEV 1
```

View the values in the rapek/s column on the test machine. The rapek/s column shows the total number of packets received by the test machine per second. In this example, the number of packets received by the test machine per second is about 12,000,000, as shown in the following figure.

11:56:22	AM	IFACE	rxpck/s	:xpck/s	rxkB/s	txkB/s	rxcmp/s	txcmp/s	rxmcst/s
11:56:23	AM	lo	0.00	0.00	0.00	0.00	0.00	0.00	0.00
11:56:23	AM	eth0	11637644.00	4.00	636433.68	0.92	0.00	0.00	0.00
11:56:23	AM	IFACE	rxpck/s	:xpck/s	rxkB/s	txkB/s	rxcmp/s	txcmp/s	rxmcst/s
11:56:24	AM	lo	0.00	0.00	0.00	0.00	0.00	0.00	0.00
11:56:24	AM	eth0	11649715.00	4.00	637093.79	0.61	0.00	0.00	0.00
11:56:24	AM	IFACE	rxpck/s	:xpck/s	rxkB/s	txkB/s	rxcmp/s	txcmp/s	rxmcst/s
11:56:25	AM	lo	0.00	0.00	0.00	0.00	0.00	0.00	0.00
11:56:25	AM	eth0	11648078.00	5.00	637004.29	0.98	0.00	0.00	0.00
11:56:25	AM	IFACE	rxpck/s	:xpck/s	rxkB/s	txkB/s	rxcmp/s	txcmp/s	rxmcst/s
11:56:26	AM	lo	0.00	0.00	0.00	0.00	0.00	0.00	0.00
11:56:26	AM	eth0	11661760.00	5.00	637752.51	1.04	0.00	0.00	0.00

Test the packet forwarding rate (greater than 20 million pps)

- 1. Connect to the test machine and the auxiliary test machine.

 For more information about the methods for connecting to an instance, see Connection methods.
- 2. Run the following command on three auxiliary test machines and one test machine to install sockperf:

```
yum install -y sockperf
```

If yum cannot be used to install sockperf, run the following command to install sockperf by using a compiler:

```
yum install -y autoconf automake libtool g++ gcc-c++
cd /opt
wget https://github.com/Mellanox/sockperf/archive/3.6.tar.gz
tar -zxf 3.6.tar.gz
cd sockperf-3.6/
./autogen.sh
./configure
make -j `cat /proc/cpuinfo| grep process | wc -l`
make install
```

3. Run the following command on the test machine to bind the interrupts:

```
a=$(cat /proc/interrupts | grep virtio2-input | awk -F ':' '{print $1}')
cpu=0
for irq in $a; do
    echo $cpu >/proc/irq/$irq/smp_affinity_list
    let cpu+=2
done
```

4. Run the following command on three auxiliary test machines to send traffic to the test machine:

```
server_ip=<Private IP address of the test machine>
threads=64
msg_size=14
run_time=60
basePort=6666
for((i=0;i<$threads;++i));do
    nohup sockperf tp -i $server_ip --client_port $[${basePort}+${i}] --pps max -m ${ms}
g_size} -t ${run_time} --port $[${basePort}+${i}] 2>&1 &
done
```

Replace <Private IP address of the test machine> with the private IP address of the actual test machine. run_time specifies the time period during which traffic is sent.

5. Run the following command on three auxiliary test machines and one test machine to test the network traffic:

```
sar -n DEV 1
```

View the values in the rxpck/s column on the test machine. The rxpck/s column shows the total number of packets received by the test machine per second. In this example, the number of packets received by the test machine per second is about 20,000,000, as shown in the following figure.

05:19:12 05:19:13 05:19:13	PM eth0	rxpck/s 20659976.00 0.00	6.00	1129842.55	txkB/s rxcmp/s 1.11 0 0.00 0.00	.00 0	.00 0.00
05:19:13 05:19:14 05:19:14	PM eth0	rxpck/s 20664073.00 0.00	5.00	1130066.52	txkB/s rxcmp/s 0.99 0 0.00 0.00	.00 0	.00 0.00
05:19:14 05:19:15 05:19:15	PM eth0	20658531.00	6.00	1129763.43	txkB/s rxcmp/s 1.16 0 0.00 0.00	.00 0	.00 0.00
05:19:15 05:19:16 05:19:16	PM eth0	20662050.00	6.00	1129955.87	txkB/s rxcmp/s 1.11 0 0.00 0.00	.00 0	.00 0.00
05:19:16 05:19:17 05:19:17	PM eth0	rxpck/s 20660346.00 0.00	5.00	1129862.69	txkB/s rxcmp/s 1.03 0 0.00 0.00	.00 0	.00 0.00
05:19:17 05:19:18 05:19:18	PM eth0	20663060.00	5.00	1130011.11	txkB/s rxcmp/s 1.03 0 0.00 0.00	.00 0	.00 0.00

Test the network bandwidth

- Connect to the test machine and the auxiliary test machine.
 For more information about the methods for connecting to an instance, see Connection methods.
- 2. Run the following command on the test machine and the auxiliary test machine to download Netperf:

```
wget https://benchmark-packages.oss-cn-qingdao.aliyuncs.com/netperf-2.7.0.tar.gz
```

3. Run the following command on the test machine and the auxiliary test machine to install Netperf and SAR:

```
yum install -y gcc autoconf automake libtool sysstat
tar -zxvf netperf-2.7.0.tar.gz
cd netperf
./configure
make && make install
```

4. Run the following command on the test machine to start 64 netserver services:

```
#!/bin/bash
for j in `seq 64`; do
    netserver -p $[16000+j] > server_$[16000+j].netperf > /dev/null 2>&1 &
done
```

5. Run the following command on the test machine to query the private IP address of the test machine:

```
ifconfig || ip addr
```

```
[root@launch-0322 netperf]# ifconfig || ip addr
eth0: flags=4163<UP.BROADCAST,RUNNING,MULTICAST> mtu 1500
       inet 172.16.....3 netmask 255.255.240.0 broadcast 172.16.
                                la prefixlen 64 scopeid 0x20<link>
       inet6
                                  weuelen 1000 (Ethernet)
       ether |
       RX packets 327551824 bytes 14206009941 (13.2 GiB)
       RX errors 0 dropped 0 overruns 0 frame 0
       TX packets 44604 bytes 6561192 (6.2 MiB)
       TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
lo: flags=73<UP,LOOPBACK,RUNNING> mtu 65536
        inet 127.0.0.1 netmask 255.0.0.0
       inet6 ::1 prefixlen 128 scopeid 0x10<host>
       loop txqueuelen 1000 (Local Loopback)
       RX packets 0 bytes 0 (0.0 B)
       RX errors 0 dropped 0 overruns 0 frame 0
       TX packets 0 bytes 0 (0.0 B)
       TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
```

6. Run the following command on the auxiliary test machine to send traffic to the test machine:

```
#!/bin/bash
server_ip=<Private IP address of the test machine>
for j in `seq 64`; do
    port=$[16000+j]
    netperf -H ${server_ip} -l ${run_time:-300} -t UDP_STREAM -p $port -- -m 1 -D > /
dev/null 2>&1 &
done
```

Replace <Private IP address of the test machine> with the private IP address of the test machine obtained in the previous step, as shown in the following figure.

```
netperf]# #!/bin/bash
[root@i
[root@i
                             netperf]# server_ip=172.1
[root@i
                             netperf]#
[root@i
                             netperf]# for j in `seq 64`; do
     port=$[16000+j]
     netperf -H ${server_ip} -l ${run_time:-300} -t UDP_STREAM -p $port -- -m 1 -D > /dev/null 2>81 8
[65] 5819
[66] 5820
[67] 5821
[68]
    5822
[69] 5823
```

7. Run the following command on the test machine to test the network bandwidth:

```
sar -n DEV 1
```

View the values in the $_{\rm rxkB/s}$ column in the test result. The $_{\rm rxkB/s}$ column shows the total number of bytes that are received by the test machine per second. The network bandwidth (unit: Kbit/s) is the value of rxkB/s multiplied by

[root@launch-0322 netperf]# sar -n DEV 1								
District St.		_64 (1	aunch-0322)	03/24	1/2022	_x86_64_	(2	CPU)
10:51:19	AM IFACE	rxpck/s	txpck/s	rxkB/s	txkB/s	rxcmp/s	txcmp/s	rxmcst/s
10:51:20	AM lo	0.00	0.00	0.00	0.00	0.00	0.00	0.00
10:51:20	AM eth0	898970.00	7.00	37750.35	2.86	0.00	0.00	0.00
10:51:20	AM IFACE	rxpck/s	txpck/s	rxkB/s	txkB/s	rxcmp/s	txcmp/s	rxmcst/s
10:51:21	AM lo	0.00	0.00	0.00	0.00	0.00	0.00	0.00
10:51:21	AM eth0	899982.00	9.00	37792.87	1.29	0.00	0.00	0.00
10:51:21	AM IFACE	rxpck/s	txpck/s	rxkB/s	txkB/s	rxcmp/s	txcmp/s	rxmcst/s
10:51:22	AM lo	0.00	0.00	0.00	0.00	0.00	0.00	0.00
10:51:22	AM eth0	899955.00	8.00	37791.75	1.22	0.00	0.00	0.00

Test the network latency

- Connect to the test machine and the auxiliary test machine.
 For more information about the methods for connecting to an instance, see Connection methods.
- 2. Run the following command on the test machine and the auxiliary test machine to install sockperf:

```
yum install -y sockperf
```

If yum cannot be used to install sockperf, run the following command to install sockperf by using a compiler:

```
cd /opt
wget https://github.com/Mellanox/sockperf/archive/3.6.tar.gz
tar -zxf 3.6.tar.gz
cd sockperf-3.6/
./autogen.sh
./configure
make -j `cat /proc/cpuinfo| grep process | wc -l`
make install
```

3. Run the following command on the test machine to start sockperf:

```
sockperf sr --daemonize > /dev/null 2>&1
```

4. Run the following command on the test machine to query the private IP address of the test machine:

```
ifconfig || ip addr
```

```
[root@launch-0322 netperf]# ifconfig || ip addr
eth0: flags=4163<UP.BROADCAST,RUNNING,MULTICAST> mtu 1500
       inet 172.16.....3 netmask 255.255.240.0 broadcast 172.16.
                   la prefixlen 64 scopeid 0x20<link>
ueuelen 1000 (Ethernet)
       inet6
       ether |
       RX packets 327551824 bytes 14206009941 (13.2 GiB)
       RX errors 0 dropped 0 overruns 0 frame 0
       TX packets 44604 bytes 6561192 (6.2 MiB)
       TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
lo: flags=73<UP,LOOPBACK,RUNNING> mtu 65536
        inet 127.0.0.1 netmask 255.0.0.0
       inet6 ::1 prefixlen 128 scopeid 0x10<host>
       loop txqueuelen 1000 (Local Loopback)
       RX packets 0 bytes 0 (0.0 B)
       RX errors 0 dropped 0 overruns 0 frame 0
       TX packets 0 bytes 0 (0.0 B)
       TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
```

5. Run the following command on the auxiliary test machine to send traffic to the test machine:

```
sockperf under-load -i <Private IP address of the test machine> --mps=100000 -t 300 -m 14 --reply-every=50 --full-log=sockperf.out
```

Replace <Private IP address of the test machine> with the private IP address of the test
machine obtained in the previous step.

View the data returned, as shown in the following figure.

- Data returned staring with avg-latency indicates the average latency in microseconds.
- Data returned staring with percentile 99.000 indicates the 99-percentile latency in microseconds.

```
~]# sockperf under-load -i 172.16. --mps=100000 -t 300 -m 14 --reply-every=50 --full-log=sockperf.out
[root@i
sockperf[CLIENT] send on:sockperf: using recvfrom() to block on socket(s)
[ 0] IP = 172.16.119.223 PORT = 11111 # UDP
sockperf: Test end (interrupted by timer)
sockperf: Test ended
sockperf: Test ended
sockperf: [Test ended
sockperf: [Total Run] RunTime-299.999 sec; Warm up time-400 msec; SentMessages-30000953; ReceivedMessages-600018
sockperf: ========= Printing statistics for Server No: 0
sockperf: [Valid Duration] RunTime-299.547 sec; SentMessages-29955801; ReceivedMessages-599117
sockperf: ====> avg=latency=>7.560 (std-dev=38.456)
sockperf: # aropped messages = 0; # auplicated messages = 0; # out-of-order messages = 6
sockperf: Summary: Latency is 37.560 usec
sockperf: Summary: Latency is 37.560 usec
sockperf: Sockperf: ---> dMAX> observations; each percentile contains 5991.17 observations
sockperf: ---> percentile 99.999 = 3299.953
sockperf: ---> percentile 99.990 = 2329.771
sockperf: ---> percentile 99.900 = 180.299
sockperf: ---> percentile 99.900 = 180.299
sockperf:
sockperf:
                  ---> percentile 99.000 =
---> percentile 90.000 =
                                                                  51.570
                                                                   43.321
sockperf:
sockperf:
                  ---> percentile 75.000 = 40.337
---> percentile 50.000 = 36.476
                   ---> percentile 25.000 =
---> <MIN> observation =
                                                                   32.355
25.768
```

18.Automate O&M based on status change events of ECS instances

In addition to the existing system events, CloudMonitor supports the status change events for Elastic Compute Service (ECS). The status change events include interruption notification events that are applied to preemptible instances. A status change event is triggered when the status of an ECS instance changes. The status changes can be caused by operations that you perform in the ECS console and by calling API operations or using SDKs, auto scaling, overdue payment, and system exceptions.

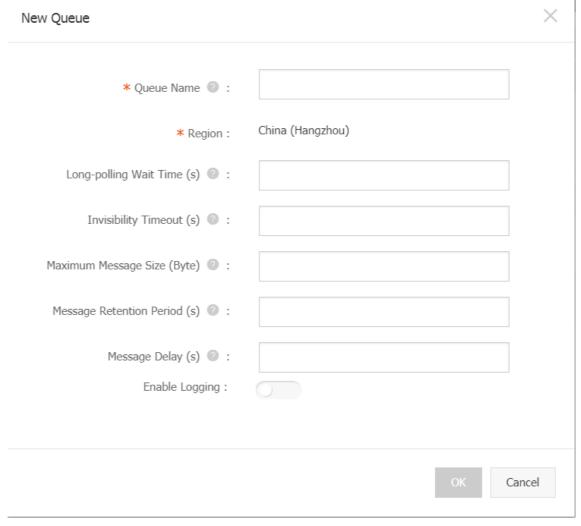
Context

The existing system events for ECS are used to notify you of alerts that require manual operations. The status change events are not about alerts. They are common notifications that are suitable for automated audit and O&M scenarios. CloudMonitor allows you to automatically handle the status change events of ECS instances by using Function Compute or Message Service (MNS).

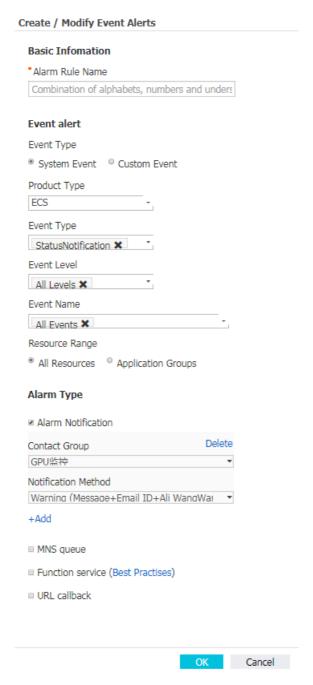
Before you begin

Create an MNS queue

- Create an MNS queue.
 - i. Log on to the MNS console.
 - ii. On the Queues page, select a region and click Create Queue in the upper-right corner.



- iii. In the Create Queue dialog box, enter a queue name, set relevant parameters, and click **OK**. In this example, set the queue name to ecs-cms-event.
- Create an event-triggered alert rule.
 - i. Log on to the CloudMonitor console.
 - ii. In the left-side navigation pane, click Event Monitoring.
 - iii. On the Event Monitoring page, click the Alarm Rules tab. On the Alarm Rules tab, click Create Event Alert.



- iv. In the **Basic Information** section of the Create / Modify Event Alert right-side pane, enter an alert rule name. In this example, enter ecs-test-rule.
- v. In the **Event alert** section, perform the following operations:
 - Set the Event Type parameter to System Event.
 - Set the **Product Type** parameter to **ECS**.
 - Set the Event Type parameter to Status Notification.
 - Set the **Event Name** parameter as needed.
 - Set the Resource Range parameter as needed. If you set the **Resource Range** parameter to **All Resources**, CloudMonitor sends alert notifications for all resource-related events. If you set the **Resource Range** parameter to **Application Groups**, CloudMonitor sends alert notifications for events related to the resources in the specified application group.

vi. In the **Alarm Type** section, perform the following operations:

- Set the Contact Group and Notification Method parameters as needed.
- Select MNS queue and set the Region and Queue parameters as needed. In this example, select the ecs-cms-event queue.
- vii. Click OK.
- Install Python dependencies.

The following code is tested in Python 3.6. You can use other programming languages, such as Java, as needed.

Use Python Package Index (PyPI) to install the following Python dependencies:

- o aliyun-python-sdk-core-v3>=2.12.1
- o aliyun-python-sdk-ecs>=4.16.0
- o aliyun-mns>=1.1.5

Procedure

CloudMonitor sends all status change events of ECS instances to MNS. Then, you can write code to receive messages from MNS and handle the messages.

• Practice 1: Record all creation and release events of ECS instances

You cannot query ECS instances that have been released in the ECS console. If you need to query released ECS instances, you can store status change events of all ECS instances in your own database or logs. When an ECS instance is created, a Pending event is triggered. When an ECS instance is released, a Deleted event is triggered. CloudMonitor records both types of events.

i. Create a Conf file.

Add the following parameters related to MNS in the Conf file:

- endpoint : the endpoint for accessing MNS. You can obtain the endpoint by clicking **Get Endpoint** on the **Queues** page in the MNS console.
- access_key and access_key_secret : the AccessKey ID and AccessKey secret used to access MNS. You can obtain the AccessKey ID and AccessKey secret in the User Management console.
- region_id and queue_name : the region where the MNS queue resides and the name of the MNS queue. You can obtain the region ID and queue name on the Queues page in the MNS console.

```
class Conf:
   endpoint = 'http://<id>.mns.<region>.aliyuncs.com/'
   access_key = '<access_key>'
   access_key_secret = '<access_key_secrect>'
   = 'cn-beijing'
   queue_name = 'test'
   vsever_group_id = '<your_vserver_group_id>'
```

ii. Use the MNS SDK to develop an MNS client for receiving messages from MNS.

```
# -*- coding: utf-8 -*-
import json
from mns.mns exception import MNSExceptionBase
import logging
from mns.account import Account
from . import Conf
class MNSClient(object):
    def init (self):
        self.account = Account(Conf.endpoint, Conf.access key, Conf.access key secre
t)
       self.queue name = Conf.queue name
        self.listeners = dict()
    def regist listener(self, listener, eventname='Instance:StateChange'):
        if eventname in self.listeners.keys():
            self.listeners.get(eventname).append(listener)
        else:
            self.listeners[eventname] = [listener]
    def run(self):
        queue = self.account.get queue(self.queue name)
        while True:
            try:
                message = queue.receive_message(wait_seconds=5)
                event = json.loads(message.message body)
                if event['name'] in self.listeners:
                    for listener in self.listeners.get(event['name']):
                        listener.process(event)
                queue.delete message(receipt handle=message.receipt handle)
            except MNSExceptionBase as e:
                if e.type == 'QueueNotExist':
                    logging.error('Queue %s not exist, please create queue before rec
eive message.', self.queue name)
                else:
                    logging.error('No Message, continue waiting')
class BasicListener(object):
   def process(self, event):
       pass
```

The preceding code is used to receive messages from MNS and delete the messages after the listener is called to consume the messages.

iii. Register a listener to consume events. The following listener generates a log entry after it receives a Pending or Deleted event.

```
# -*- coding: utf-8 -*-
import logging
from .mns_client import BasicListener
class ListenerLog(BasicListener):
    def process(self, event):
        state = event['content']['state']
        resource_id = event['content']['resourceId']
        if state == 'Panding':
            logging.info(f'The instance {resource_id} state is {state}')
        elif state == 'Deleted':
            logging.info(f'The instance {resource_id} state is {state}')
```

Add the following code to the Main function:

```
mns_client = MNSClient()
mns_client.regist_listener(ListenerLog())
mns_client.run()
```

In the production environment, you can store the events in your database or Log Service for subsequent queries and audits.

• Practice 2: Automatically start ECS instances that are shut down

In scenarios where ECS instances may be shut down unexpectedly, you may want to automatically start the ECS instances.

You can reuse the MNS client developed in Practice 1 and create another listener. When the listener receives a Stopped event for an ECS instance, you can run the **start** command on the ECS instance to start it.

```
# -*- coding: utf-8 -*-
import logging
from aliyunsdkecs.request.v20140526 import StartInstanceRequest
from aliyunsdkcore.client import AcsClient
from .mns client import BasicListener
from .config import Conf
class ECSClient(object):
   def init (self, acs client):
       self.client = acs client
    # Start the target ECS instance.
   def start instance(self, instance id):
        logging.info(f'Start instance {instance id} ...')
        request = StartInstanceRequest.StartInstanceRequest()
       request.set accept format('json')
       request.set InstanceId(instance id)
        self.client.do_action_with_exception(request)
class ListenerStart(BasicListener):
   def init (self):
       acs client = AcsClient(Conf.access key, Conf.access key secret, Conf.region id)
        self.ecs client = ECSClient(acs client)
   def process(self, event):
       detail = event['content']
       instance id = detail['resourceId']
        if detail['state'] == 'Stopped':
            self.ecs client.start instance(instance id)
```

In the production environment, you can listen to Starting, Running, or Stopped events after the **start** command is run. Then, you can perform further O&M by using a timer and a counter based on whether the ECS instance is started.

• Practice 3: Automatically remove preemptible instances from SLB before they are released

An interruption notification event is triggered 5 minutes before a preemptible instance is released. During the 5 minutes, you can perform specific operations to prevent your services from being interrupted. For example, you can remove the target preemptible instance from a Server Load Balancer (SLB) instance.

You can reuse the MNS client developed in Practice 1 and create another listener. When the listener receives the interruption notification event for a preemptible instance, you can call the SLB SDK to remove the preemptible instance from an SLB instance.

```
# -*- coding: utf-8 -*-
from aliyunsdkcore.client import AcsClient
from aliyunsdkcore.request import CommonRequest
from .mns client import BasicListener
from .config import Conf
class SLBClient(object):
   def init (self):
        self.client = AcsClient(Conf.access key, Conf.access key secret, Conf.region id)
        self.request = CommonRequest()
        self.request.set method('POST')
        self.request.set accept format('json')
        self.request.set version('2014-05-15')
        self.request.set domain('slb.aliyuncs.com')
        self.request.add query param('RegionId', Conf.region id)
   def remove vserver group backend servers (self, vserver group id, instance id):
        self.request.set_action_name('RemoveVServerGroupBackendServers')
        self.request.add query param('VServerGroupId', vserver group id)
        self.request.add query param('BackendServers',
                                     "[{'ServerId':'" + instance id + "', 'Port':'80', 'Wei
ght':'100'}]")
        response = self.client.do action with exception(self.request)
        return str(response, encoding='utf-8')
class ListenerSLB(BasicListener):
   def init (self, vsever group id):
        self.slb caller = SLBClient()
        self.vsever_group_id = Conf.vsever_group_id
   def process(self, event):
       detail = event['content']
        instance id = detail['instanceId']
        if detail['action'] == 'delete':
            self.slb caller.remove vserver group backend servers(self.vsever group id, in
stance_id)
```

Notice

For interruption notification events, set the event name in the following way: mns_client.re gist_listener(ListenerSLB(Conf.vsever_group_id), 'Instance:PreemptibleInstanceInterruption') .

In the production environment, you can apply for another preemptible instance and add it as a backend server to SLB to ensure the performance of your services.

19.Use the snapshot and image features to migrate instance data

Continuous iterations of Elastic Compute Service (ECS) instances give rise to issues with instances that were created a long time ago. For example, old instances may run short of resources. This affects the O&M efficiency of your cloud-based business. To ensure that you can manage and operate your cloud-based business in an efficient manner, we recommend that you use the snapshot and image features to migrate data from an old instance to a new instance.

Context

Alibaba Cloud snapshot is an agentless data backup feature that backs up or restores data for entire disks. For more information about snapshots, see Snapshot overview. Custom images are created based on snapshots of ECS instances. You can quickly create identical ECS instances from a custom image. For more information about custom images, see Overview.

You can perform the following procedure to migrate instance data by using the snapshot and custom image features:

- Read the notes carefully before you migrate instance data. For more information, see Notes.
- Step 1: Create a custom image from the source instance.
- Step 2: (Optional) Copy the custom image to a different region.
- Step 3: Create a destination instance from the custom image.
- Step 4: Check data in the destination instance.
- Step 5: (Optional) Release or delete the source instance and its resources.

Notes

Before you migrate instance data, take note of the following items:

- Snapshots cannot be created for some instances that use local storage. The procedure described in this topic cannot be used to migrate data to or from such instances.
- Source instances can be located in the classic network or virtual private clouds (VPCs).
- Destination instances can be created only in VPCs.
- Only instance types available within the selected zone can be used to create destination instances.

Note If you want to migrate instance data across zones or regions, we recommend that you plan resources in the source and destination zones or regions in advance.

After you use the snapshot and image features to migrate instance data, disk data on the
destination instance is the same as that on the source instance, but new metadata is generated for
the destination instance and is different from that of the source instance. For more information
about instance metadata, see Overview of ECS instance metadata.

Before you migrate instance data, we recommend that you manually sort out resource associations so that you can update resource associations accordingly in a timely manner after migration. Examples:

- Assume that instances use private IP addresses to communicate within a cluster. After you migrate
 data from one instance to another within the cluster, you must replace the private IP address of
 the source instance with that of the destination instance so that the destination instance can
 communicate with other instances in place of the source instance.
- The licenses of some applications are bound to the media access control (MAC) addresses of
 instances. After you migrate data from one instance to another, the licenses bound to the MAC
 address of the source instance become invalid, and you must rebind the licenses to the MAC
 address of the destination instance.

Step 1: Create a custom image from the source instance

Before you create a custom image from the source instance, make sure that you have understood the items that you must take note of. For more information, see Create a custom image from an instance. When the system creates a custom image from the source instance, the system creates a snapshot for each of the disks attached to the instance. These snapshots are billed based on the storage they use. For more information about the billing of snapshots, see Snapshots.

- 1.
- 2.
- 3.
- 4. Find the source instance and choose More > Disk and Image > Create Custom Image in the Actions column.
- 5. In the Create Custom Image dialog box, configure parameters and click Create.
 - You must specify the Custom Image Name and Custom Image Description parameters. You can optionally specify other parameters based on your business requirements.
- 6.
- 7. On the **Custom Image** tab, find the custom image that you created in the previous step and check its state

If multiple custom images are displayed on the Custom Image tab, you can search for the custom image by name. When the custom image is in the **Available** state, you can proceed to the next step.

Step 2: (Optional) Copy the custom image to a different region

If you want to migrate data from the source instance to a destination instance that is located in a different region, you must first use the image copy feature to copy the custom image that you created in Step 1 to that region. For more information, see Copy a custom image.

After the custom image is copied to the destination region, you must perform the subsequent steps in that region.

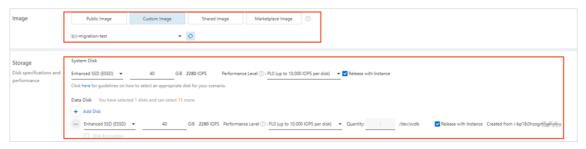
Step 3: Create a destination instance from the custom image

- 1.
- 2. In the top navigation bar, select a region.
 - If you have performed Step 2 to copy the custom image to a different region, switch to that region.
- 3. On the **Custom Image** tab, find the custom image that you created from the source instance. If you have performed Step 2 to copy the custom image, find the image copy.

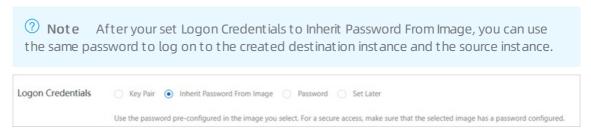
- 4. In the Actions column, click Create Instance.
- 5. On the instance buy page, configure parameters to create an instance to use as the destination instance.

For information about how to create an ECS instance, see Create an instance by using the wizard. When you create the destination instance, take note of the following items:

• In the **Basic Configurations** step, the parameters in the **Image** and **Storage** sections have been automatically specified and do not need to be modified.



• In the System Configurations (Optional) step, set Logon Credentials to Inherit Password From Image.



Step 4: Check data in the destination instance

You must check the data in the destination instance to ensure that services can continue to run smoothly on the destination instance after migration. Examples:

- You must check the data integrity of the destination instance based on the data in the source instance.
- You must compare resource information between the source and destination instances and update the resource associations for the destination instance. For information about how to view instance information, see View instance information or View instance metadata.

Step 5: (Optional) Release or delete the source instance and its resources

After you perform the following operations, you can release or delete the source instance and its resources to avoid unnecessary costs: Verify that data on the destination instance is the same as that on the source instance, update resource associations for the destination instance, and make sure that services can continue to run smoothly on the destination instance.

Notice The instance release, snapshot deletion, and image deletion operations are irreversible. When an instance is released or after a snapshot or image is deleted, data in the instance, snapshot, or image is deleted and cannot be recovered. Make sure that you release or delete resources only after all business data is migrated.

• For information about how to release an instance, see Release an instance.

- For information about how to delete a snapshot, see Delete a snapshot.
- For information about how to delete a custom image, see Delete a custom image.

Note After a custom image is deleted, instances that are using this image cannot have their system disks initialized. If your custom image is free of charge, we recommend that you do not delete it. For information about the billing of images, see Images.

20. Disaster recovery solutions

Disaster recovery solutions help ensure the running stability and security of your services and IT systems by incorporating data backup and disaster recovery. Alibaba Cloud ECS allows you to use snapshots and images to back up data.

Disaster recovery methods

Snapshot backup

Alibaba Cloud ECS allows you to back up system disks and data disks with snapshots. Alibaba Cloud provides the Snapshot 2.0 service, which features a higher snapshot quota and a more flexible automatic task strategy than previous snapshot services, to help reduce the impacts on business I/O. When snapshots are used for data backup, the first snapshot of a disk is a full backup, and subsequent snapshots are incremental backups. Incremental snapshots can be created quickly and have small sizes. The amount of time required for backup depends on the amount of incremental data to be backed up.

Note Snapshots are created on an incremental basis. To improve backup speed, we recommend that you create a new snapshot before deleting the most recent one.



The preceding figure shows how incremental snapshots work. In the figure, Snapshots 1, 2, and 3 represent the first, second, and third snapshot of a disk. The file system checks the disk data block by block. When a snapshot is created, only the blocks with changed data are copied to the snapshot. Alibaba Cloud ECS allows you to configure manual or automatic snapshots of disks. To create automatic snapshots of a disk, you can configure and apply an automatic snapshot policy to the disk. You can specify the hour of the day (on the hour), day of week (Monday through Sunday), and retention period for snapshot creation in the policy. You can customize the retention period to a value from 1 to 65,536 days, or choose to save snapshots permanently.

• Snapshot rollback

When exceptions occur in your system and you want to roll a disk back to a previous state, you can roll the disk back to a created snapshot. For more information, see Roll back a disk by using a snapshot. Note the following points:

- Rollback operations are irreversible. After a rollback is complete, data before the rollback cannot be restored. Exercise caution when you perform this operation.
- When a disk is rolled back, all data created or modified between the current time and the snapshot creation time is lost.
- Image backup

An image works as a copy that stores data from one or more disks. An ECS image may store data from a system disk or from both system and data disks. All image backups are full backups and can only be triggered manually.

Image recovery

You can create a custom image from a snapshot to include the operating system and data environment of the snapshot in the image. Then, you can use the custom image to create multiple instances with the same operating system and data environment. For more information about the configuration of snapshots and images, see Create a snapshot of a disk and Create a custom image from a snapshot.



? Note Custom images cannot be used across regions.

Technical metrics

RTO and RPO are related to the amount of data, typically on an hourly basis.

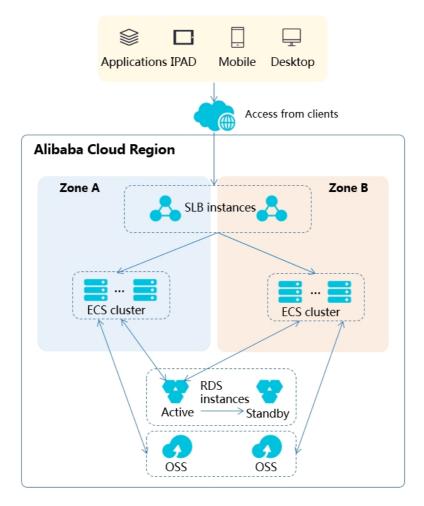
Scenarios

• Backup and recovery

Alibaba Cloud ECS allows you to back up system disks and data disks with snapshots and images. If incorrect data is stored on a disk due to application errors or hackers' malicious access through application vulnerabilities, you can use the snapshot service to restore the disk to a desired state. In addition, Alibaba Cloud ECS allows you to reinitialize disks with images or create ECS instances from cust om images.

Disaster recovery

Alibaba Cloud ECS supports the implementation of disaster recovery architecture. For example, you can buy and use an SLB instance at the frontend of an application, and deploy at least two ECS instances at the backend of the same application. Alternatively, you can use Auto Scaling provided by Alibaba Cloud to perform auto scaling by defining how to use ECS resources. This way, even if one of the ECS instances fails or is overloaded, disaster recovery can be implemented to ensure business continuity. The following figure provides an example in which ECS instances are deployed in data centers in two zones within the same region. All communications are implemented in the Alibaba Cloud Gigabit internal network to ensure fast response and reduce Internet traffic costs.



- SLB: SLB instances are used for load balancing between the two zones. Traffic is distributed to two or more data centers where ECS instance clusters are deployed.
- ECS cluster: ECS instances deployed in the two data centers are equivalent. The failure of a single instance does not affect data layer applications and the ECS control function. If a failure occurs, the system automatically performs hot migration so that other ECS instances can continue to provide services. This can prevent service interruptions caused by a single point of failure or hot migration failures. If hot migration fails, you will receive a notification about the failures based on system events so that you can deploy new nodes in a timely manner.
- Data layer: OSS is deployed at the region level. ECS nodes in data centers in different zones can
 access objects in OSS. For database applications, multi-zone ApsaraDB for RDS service is used.
 Primary nodes can perform read and write operations across zones without conflicting with
 application-layer traffic. In addition, secondary nodes can perform read operations across zones to
 prevent inability of ECS instances to read data in case of failures of the primary nodes.

21.Deploy a highly available architecture 21.1. Deploy a highly available architecture

Highly available architectures provide functions such as service distribution, auto scaling, and multizone deployment. Compared with a single ECS instance, a highly available architecture is more stable and scalable when databases and applications are deployed.

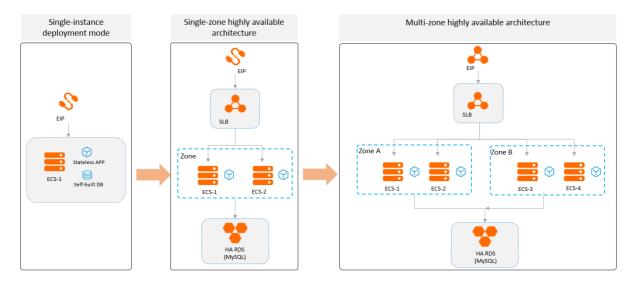
Features

Highly available architectures have the following features:

- A multi-zone, highly available SLB instance distributes traffic to multiple ECS instances, increasing the external service capability of application systems, eliminating single point of failures, and improving the availability of application systems. SLB is used for automatic multi-zone deployment, enhancing disaster recovery capabilities of services.
- You can use custom images to create identical ECS instances, and then add these instances to the
 backend server group of the SLB instance, implementing high availability for your services. SLB can be
 configured with Layer-4 and Layer-7 listeners at the same time, and with multiple algorithms, such as
 round robin, weighted round robin, and weighted least connections, properly allocating computing
 resources to backend ECS instances.
- Relational Database Service (RDS) can be optimized for high concurrency scenarios, ensuring the
 constant stability and high throughput of the system through thread pools, parallel replication, and
 hidden primary keys. CloudDBA provides comprehensive performance monitoring metrics to monitor
 the usage of instances and hardware and slow SQL queries in real time, and gives optimization
 suggestions to help you locate and solve problems.

Deployment process

If you have created an ECS instance and deployed databases and applications on the instance, you can change the single-instance deployment mode to single-zone or multi-zone, highly available architecture. This topic shows you how to use ECS, EIP, SLB, and RDS to deploy a multi-zone, highly available architecture.



- 1. Use a custom image to create multiple identical ECS instances. For more information, see Replicate ECS instances.
- 2. Create an SLB instance and add the ECS instances to the SLB backend server group to mount ECS instances from different zones, achieving the high availability of services. For more information, see Configure an SLB instance.
- 3. Use DTS to migrate a user-created database from an ECS instance to an RDS instance, ensuring that the business database is not interrupted during migration and data is automatically backed up. For more information, see Migrate self-managed databases to ApsaraDB RDS instances that run RDS High-availability Edition.

21.2. Replicate ECS instances

This topic describes how to use a custom image created from a source instance to create three ECS instances for multi-zone disaster recovery. One instance is assigned to the same zone as the source instance, and the other two instances are assigned to a different zone in the same region as the source instance.

Prerequisites

- •
- You have an ECS instance to be replicated.

Procedure

- 1.
- 2. Create three ECS instances from the created custom image.
 - i.
 - ii. On the **Custom Images** tab of the Images page, find the custom image created in the previous step. Click **Create Instance** in the **Actions** column corresponding to the image.

iii. On the **Custom Launch** tab, scroll down to the **Image** section. The custom image you selected is automatically selected. Continue with other settings as prompted and set Quantity to 1.

where:

- Region: Select the same region as the source instance.
- Zone: Select the same zone as the source instance.
- Public IP Address: Clear Assign Public IP Address.

For more information, see Create an instance by using the wizard.

iv. Repeat Step a and Step b. On the Custom Launch tab, scroll down to the Image section. The custom image you selected is automatically selected. Continue with other settings as prompted and set Quantity to 2.

where:

- **Region**: Select the same region as the source instance.
- Zone: Select a zone different from that of the source instance.
- Instance Type: Set Quantity to 2.
- Public IP Address: Clear Assign Public IP Address.

For more information, see Create an instance by using the wizard.

Result

In the left-side navigation pane, choose **Instances & Images > Instances**. On the **Instances** page, the four ECS instances are in the **Running** state. Two instances are located in a zone, and the other two are located in another zone.

What's next

Configure an SLB instance

21.3. Configure an SLB instance

After the ECS instances are replicated, you can create a multi-zone SLB instance in a region that supports multiple zones and bind multi-zone ECS instances to the SLB instance. This task can extend the external service capabilities of application systems, eliminate single point of failures, and improve the availability of application systems. This topic describes how to deploy an SLB instance.

Prerequisites

- Three ECS instances are replicated. For more information, see Replicate ECS instances.
- The web services of the four ECS instances are started and are running normally.

Notice If the web services are not running, the SLB instance and the ECS instances cannot communicate normally.

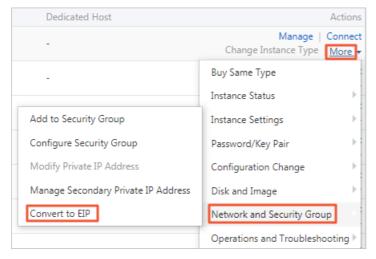
Procedure

1. Create an SLB instance. For more information, see Create a CLB instance.

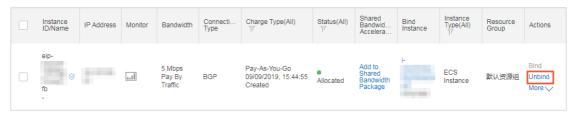
The following settings are used in this topic:

Region: Select the same region as the ECS instances.

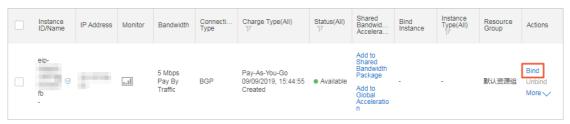
- Incaroni percet the pame region as the Ecombiances
- o Zone Type: Select Multi-zone.
- o Instance Type: Select Internal Network.
- Network Type: Select VPC.
- **Primary Zone** and **Backup Zone**: Configure as needed.
- 2. Convert the public IP address of the source instance into an Elastic IP Address. For more information, see Convert the public IP address of a VPC-type instance to an EIP.
 - Note The IP address of the source instance must remain unchanged so that services are not affected. Therefore, you must first convert the public IP address of the source instance into an Elastic IP Address, unbind the Elastic IP Address from the source instance, and then bind the Elastic IP Address to the multi-zone SLB instance.



- 3. Unbind the Elastic IP Address from the source instance.
 - i. In the IP Address column of the source instance, click the link of the Elastic IP Address.
 - ii. On the Elastic IP Addresses page, click Unbind.



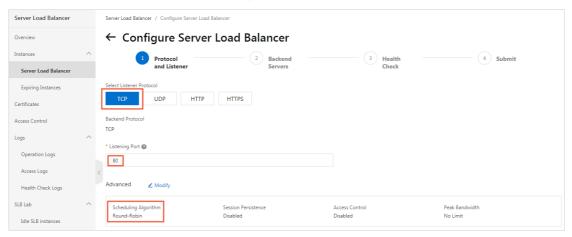
- iii. Click OK. For more information, see Disassociate an EIP from a cloud resource.
- 4. Bind the Elastic IP Address to the SLB instance.
 - i. On the **Elastic IP Addresses** page, find the Elastic IP Address that was unbound from the source instance.



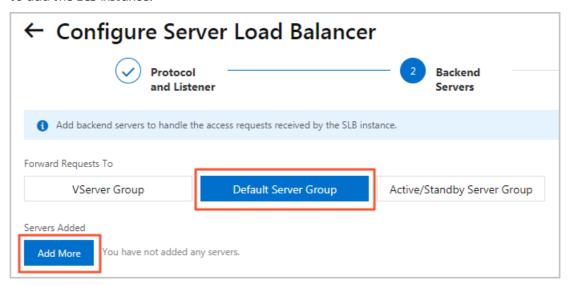
- ii. In the Actions column, click Bind.
- iii. Select SLB Instance for Instance Type, select the SLB instance that you created for SLB Instance, and then click OK. For more information, see Associate with a CLB instance.
- 5. Configure an SLB instance. For more information, see Configure a CLB instance.

Perform the following steps to complete the basic settings:

- i. On the **Protocol and Listener** tab, complete the following configuration:
 - Select Listener Protocol: Select TCP.
 - Listening Port: Enter 80 .
 - Scheduling Algorithm: Set this parameter as needed. In this topic, Scheduling Algorithm is set to Round-Robin (RR).
 - Use the default values for other settings.



ii. Click **Next**. On the **Backend Servers** tab, select **Default Server Group**, and click **Add More** to add the ECS instance.

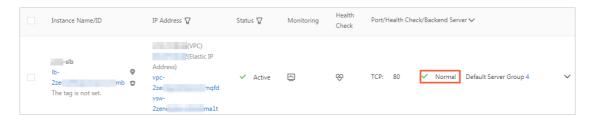


- iii. Select the source instance and the three replicated ECS instances, and click **Next: Set Weight** and **Port**. Set **Port** to 80 and remain the default values for other settings. Click **Next**.
- iv. On the Health Check tab, use the default values and then click Next.
- v. On the Submit tab, verify the information and click Submit.

vi. Click **OK** to go back to the **Server Load Balancer** page, and click ${}^{\square}$.

If the health check is **Normal**, the backend ECS instance is working properly and able to accept requests.

? Note It takes a few minutes to perform the health check. Wait and click the refresh icon to view the status.



Result

In this topic, a static web page is built on each of the four ECS instances to identify each instance. Enter the endpoint of the SLB instance in the browser to test whether the SLB is working properly. Because **Scheduling Algorithm** was set to **Round-Robin (RR)**, requests are sent to each ECS instance in turn.

What's next

Migrate self-managed databases to ApsaraDB RDS instances that run RDS High-availability Edition

21.4. Migrate self-managed databases to ApsaraDB RDS instances that run RDS High-availability Edition

You can migrate databases from Elastic Compute Service (ECS) instances to ApsaraDB RDS instances that run RDS High-availability Edition to ensure high availability, reliability, security, and convenience of database services. This topic describes how to migrate a self-managed database on an ECS instance to an ApsaraDB RDS instance that runs RDS High-availability Edition by using Data Transmission Service (DTS). In this topic, a MySQL database is used.

Prerequisites

- A Server Load Balancer (SLB) instance is configured. For more information, see Configure an SLB instance.
- An ApsaraDB RDS instance that runs RDS High-availability Edition is created and uses the multi-zone deployment method. For more information about how to create an ApsaraDB RDS instance, see Create an ApsaraDB RDS for MySQL instance.

•

Context

DTS supports data migration between homogeneous and heterogeneous data sources. It also supports ETL features such as data mapping at three levels (databases, tables, and columns) and data filtering. You can use DTS for zero-downtime data migration. During data migration, the source database continues to provide services normally, minimizing the impact of data migration on your business. For more information about the database types supported by DTS, see Data migration.

Procedure

- 1. Log on to the DTS console.
- 2. In the left-side navigation pane, click **Data Migration**.

3.

- 4. Configure a migration task.
 - i. Configure a task name.

You can use the default name or customize a name.

ii. Configure the source database.

DTS supports self-managed databases connected over the Internet, VPN Gateway, Express Connect, or Smart Access Gateway. In this topic, the source database is a self-managed database on an ECS instance. For information about migration solutions for other database types, see What is DTS?

iii.

iv.

٧.

vi.

5.

6.

What's next

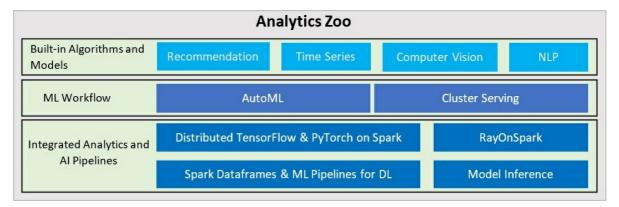
Configure the endpoint, account, and password of the RDS instance in applications to connect to the RDS instance. You can use Data Management (DMS) or the MySQL client to manage the RDS instance. For more information, see Use a database client or the CLI to connect to an ApsaraDB RDS for MySQL instance.

22.Use Analytics Zoo and bfloat16 to accelerate Al applications on ECS instances

This topic describes how to use Analytics Zoo and the Brain Floating Point bit-16 (bfloat16) capability provided by third-generation Intel® Xeon® scalable processors to accelerate the performance of AI applications on Elastic Compute Service (ECS) instances. This topic is applicable to seventh-generation high frequency ECS instances (also called ECS instances with high clock speeds).

Context

- Seventh-generation high frequency ECS instances are built on the third-generation SHENLONG
 architecture and powered by third-generation Intel® Xeon® scalable processors. These instances can
 deliver a compute performance up to 260% higher than the sixth-generation high frequency
 instances. When you use Analytics Zoo on ECS instances, you can leverage advanced pipeline
 features such as the TensorFlow and PyTorch deep learning models optimized by Intel to develop
 deep learning applications.
- Third-generation Intel® Xeon® scalable processors provide industry-leading, workload-optimized platforms by using enhanced Intel® Deep Learning Boost (Intel® DL Boost), which is a built-in AI acceleration feature. Enhanced Intel® DL Boost features the industry's first x86 support for bfloat16 to drive better AI inference and training performance.
 - Third-generation Intel® Xeon® scalable processors can cope with complex AI workloads. Enhanced Intel® DL Boost can increase AI training performance by up to 1.93×, AI inference performance for image classification by up to 1.87×, AI training performance for natural language processing (NLP) by up to 1.7×, and AI inference performance by up to 1.9×. Support for bfloat16 makes it much more efficient to handle AI training workloads for the healthcare, financial services, and retail industries.
- Analytics Zoo is a unified, open source big data analytics and AI platform developed by Intel. It can seamlessly scale AI models such as TensorFlow, Keras, and PyTorch to adapt to distributed big data environments such as Spark, Flink, and Ray. Analytics Zoo provides the following features:
 - End-to-end pipelines for running AI models such as TensorFlow, PyTorch, and OpenVINO on big data platforms. For example, developers can embed TensorFlow or PyTorch code into Spark code for distributed training and inference. Developers can also use the native deep learning models such as TensorFlow, Keras, PyTorch, and BigDL in Spark machine learning (ML) pipelines.
 - High-level ML workflows, such as cluster serving and scalable AutoML for automated ML tasks. Cluster serving is an automatic and distributed inference solution for models such as TensorFlow, PyTorch, and the OpenVINO toolkit. Scalable AutoML is used for time-series predictions.
 - Built-in common models for applications such as recommendation, time series, computer vision, and NLP.



- Bfloat 16 is a numeric format widely used in neural networks.
- ResNet-50 is a residual network that is 50 layers deep. It is widely used in fields such as object classification.

Procedure

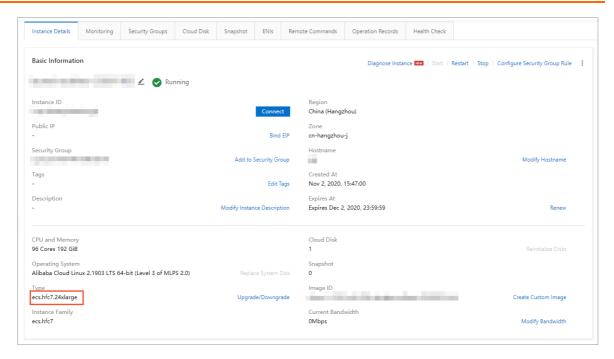
To use Analytics Zoo and bfloat 16 to accelerate AI applications on an ECS instance, perform the following steps:

- 1. Step 1: Create a high frequency ECS instance
- 2. Step 2: Prepare an Analytics Zoo environment that provides enhanced support for bfloat 16 on the ECS instance
- 3. Step 3: Train ResNet-50 models and use bfloat 16 to improve performance on the ECS instance

Step 1: Create a high frequency ECS instance

To create a high frequency ECS instance, perform the following operations:

- 1. Go to the Custom Launch tab of the instance buy page in the ECS console.
- 2. Create a high frequency instance. For more information, see Create an instance by using the wizard. When you configure parameters to create the high frequency instance, you can select the hfc7 or hfg7 instance family. In this example, select the hfc7 instance family. For information about specific instance types, see Instance families with high clock speeds.
- 3. On the Instances page, find the instance that you created in the previous step and click the ID of the instance. Confirm that the instance is of the hfc7 instance type.



Step 2: Prepare an Analytics Zoo environment that provides enhanced support for bfloat16 on the ECS instance

Analytics Zoo provides a precreated Docker image that supports bfloat16. You can use method 1 to obtain this Docker image in Alibaba Cloud ECS. Alternatively, you can use method 2 to obtain a nightly build of Analytics Zoo that supports bfloat16. For information about the related code, see Sample code: How Analytics Zoo uses bfloat16 to accelerate the training of deep learning models.

- Method 1: Obtain the Docker image provided by Analytics Zoo in ECS
 - i. Connect to the ECS instance. For more information, see Connect to an ECS instance.
 - ii. Run the following commands to install and run Docker:

```
yum install docker-io -y
systemctl start docker
```

iii. Run the following command to obtain the Docker image provided by Analytics Zoo that supports bfloat 16:

```
docker pull intelanalytics/analytics-zoo:0.8.1-bigdl_0.10.0-spark_2.4.3-bf16
```

iv. Run the following command to run the Docker container that corresponds to the Docker image:

```
docker run -itd --name az1 --net=host --privileged intelanalytics/analytics-zoo:0.8.
1-bigdl_0.10.0-spark_2.4.3-bf16
```

v. Run the following command to access the container:

```
docker exec -it az1 bash
```

- Method 2: Use a nightly build of Analytics Zoo that supports bfloat16 to manually create an Analytics Zoo environment
 - i. Connect to the ECS instance. For more information, see Connect to an ECS instance.
 - ii. Run the following commands to download and decompress the latest Analytics Zoo prebuilt

release and nightly build package:

```
wget https://oss.sonatype.org/content/repositories/snapshots/com/intel/analytics/zoo/analytics-zoo-bigdl_0.11.1-spark_2.4.3/0.9.0-SNAPSHOT/analytics-zoo-bigdl_0.11.1-spark_2.4.3-0.9.0-20201026.210040-51-dist-all.zip unzip analytics-zoo-bigdl_0.11.1-spark_2.4.3-0.9.0-{datetime}-dist-all.zip -d analytics-zoo
```

iii. Run the following command to install Git:

```
yum -y install git
```

iv. Run the following commands to download TensorFlow source code:

```
git clone https://github.com/Intel-tensorflow/tensorflow.git git checkout v1.15.0up1
```

v. Run the following commands to compile TensorFlow:

```
bazel build --cxxopt=-D_GLIBCXX_USE_CXX11_ABI=0 --copt=-O3 --copt=-Wformat
--copt=-Wformat-security --copt=-fstack-protector --copt=-fPIC
--copt=-fpic --linkopt=-znoexecstack --linkopt=-zrelro
--linkopt=-znow --linkopt=-fstack-protector --config=mkl --define
build_with_mkl_dnn_v1_only=true --copt=-DENABLE_INTEL_MKL_BFLOAT16
--copt=-march=native
//tensorflow/tools/lib_package:libtensorflow_jni.tar.gz
//tensorflow/java:libtensorflow.jar
//tensorflow/java:libtensorflow-src.jar
//tensorflow/tools/lib_package:libtensorflow_proto.zip
```

vi. Run the following commands to organize the library files required by Analytics Zoo:

```
cd bazel-bin/tensorflow/tools/lib_package
mkdir linux-x86_64
tar -xzvf libtensorflow_jni.tar.gz -C linux_x86-64
rm libtensorflow_framework.so
rm libtensorflow_framework.so.1
mv libtensorflow_framework.so.1.15.0 libtensorflow_framework-zoo.so
cp ../../../_solib_k8/_U@mkl_Ulinux_S_S_Cmkl_Ulibs_Ulinux___Uexternal_Smkl_Ulinux_Slib/* ./
```

vii. Run the following command to update the Analytics Zoo Jar folder:

```
cd ~/analytics-zoo/lib/
cp ~/tensorflow/bazel-bin/tensorflow/tools/lib_package/linux-x86_64 ./
jar -ufanalytics-zoo-bigdl_0.11.1-spark_2.4.3-0.9.0-SNAPSHOT-jar-with-dependencies.ja
r linux-x86_64/*
```

Step 3: Train ResNet-50 models and use bfloat16 to improve performance on the ECS instance

1. Run the following command to access the Analytic Zoo Docker container:

```
docker exec -it azl bash
```

2. Run the following commands to configure Spark and modify the /opt/work/spark-2.4.3/conf/spar

k-defaults.conf file:

```
spark.authenticate=false
spark.ui.killEnabled=true
spark.eventLog.enabled=true
spark.history.ui.port=18080
spark.eventLog.dir=file:///var/log/spark/spark-events
spark.history.fs.logDirectory=file:///var/log/spark/spark-events
spark.shuffle.service.port=7337
spark.master=spark://$(hostname):7077
```

3. Run the following commands to start the Spark master:

```
cd /opt/work/spark-2.4.3
./sbin/start-master.sh
```

4. Run the numactl command to start eight Spark workers and bind each Spark worker to 12 vCPUs. Then, create the following script in the /opt/work/spark-2.4.3/bin directory:

```
numactl -C 0-11 ./spark-class org.apache.spark.deploy.worker.Worker spark://$(hostname):7077 &
numactl -C 12-23 ./spark-class org.apache.spark.deploy.worker.Worker spark://$(hostname):7077 &
numactl -C 24-35 ./spark-class org.apache.spark.deploy.worker.Worker spark://$(hostname):7077 &
numactl -C 36-47 ./spark-class org.apache.spark.deploy.worker.Worker spark://$(hostname):7077 &
numactl -C 48-59 ./spark-class org.apache.spark.deploy.worker.Worker spark://$(hostname):7077 &
numactl -C 60-71 ./spark-class org.apache.spark.deploy.worker.Worker spark://$(hostname):7077 &
numactl -C 72-83 ./spark-class org.apache.spark.deploy.worker.Worker spark://$(hostname):7077 &
numactl -C 84-95 ./spark-class org.apache.spark.deploy.worker.Worker spark://$(hostname):7077 &
```

5. Run the following commands to check how many Spark workers have been started. If eight Spark workers have been started, 8 is displayed in the command output.

```
jps | grep Worker | wc -l
```

6. Run the following commands to download ResNet-50 sample code from Github:

```
git clone https://github.com/yangw1234/models-1.git git checkout branch-1.6.1-zoo
```

7. Run the run.sh script in the *models-1/models/image_recognition/tensorflow/resnet50v1_5/trainin g/mlperf_resnet* directory and enable bfloat16-based training by adding the --use_bfloat16 option. If the --use_bfloat16 option is not added, float32 is used for training by default.

```
# Register the model as a source root
export PYTHONPATH="$(pwd):${PYTHONPATH}"
export KMP_BLOCKTIME=0
# 8 instances
export OMP_NUM_THREADS=6
export KMP_AFFINITY=granularity=fine,compact,1,0
export KMP_SETTINGS=1
export ANALYTICS_ZOO_HOME=/opt/work/analytics-zoo/dist
export SPARK_HOME=/opt/work/spark-2.4.3
bash $ANALYTICS_ZOO_HOME/bin/spark-submit-python-with-zoo.sh --master
spark://$(hostname):7077 \
--executor-cores 1 --total-executor-cores 8 --driver-memory 20g --executor-memory 18g \
--conf spark.network.timeout=10000000 --conf spark.executor.heartbeatInterval=100000 \
imagenet_main.py 1 --model_dir ./logs --batch_size 128 --version 1 \
--resnet_size 50 --train_epochs 90 --data_dir /opt/ILSVRC2012/ --use_bfloat16
```

The following table describes the training test results.

ResNet-50 model training	FP32	BF16	Performance improvement achieved by bfloat16 over float32
Throughput(images/se c)	119.636	212.315	1.775

Sample code: How Analytics Zoo uses bfloat16 to accelerate the training of deep learning models

The following code provides an example of how to use bfloat 16 to accelerate the training of deep learning models such as ResNet-50 models. The code is for reference only. It is already included in the Docker image provided by Analytics Zoo, and does not require manual operations.

1. Use the following code to convert input images into the bfloat16 format:

```
if use_bfloat16 == True:
dtype = tf.bfloat16
features = tf.cast(features, dtype)
```

2. Use the following code to compile custom_dtype_getter:

3. Use the following code to create a variable scope and build a model under the variable scope:

4. Use the following code to cast logits to the float 32 format and then calculate loss to ensure numerical stability:

```
logits = tf.cast(logits, tf.float32)
```

5. Use Analytics TFPark for distributed training. For more information, see TFPark Overview.