

Elastic Compute Service

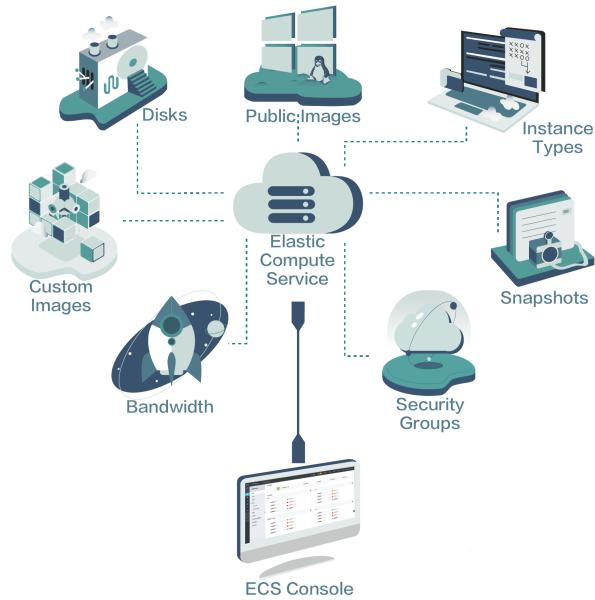
Product Introduction

Product Introduction

Elastic Compute Service (ECS) is a type of computing service that features elastic processing capabilities. ECS has a simpler and more efficient management mode than physical servers. You can create instances, change the operating system, and add or release any number of ECS instances at any time to fit your business needs.

An ECS instance is a virtual computing environment that includes CPU, memory, and other basic computing components. An instance is the core component of ECS and is the actual operating entity offered by Alibaba Cloud. Other resources, such as disks, images, and snapshots, can only be used in conjunction with an ECS instance.

The following figure illustrates the concept of an ECS instance. You can use ECS console to configure the instance type, disks, operating system, and other affiliated resources for your ECS instance.



Advantages

ECS has the following advantages over traditional servers and virtual hosts:

Stability

ECS has 99.95% service availability and 99.9999999% data reliability. It also supports switchover, data snapshot backup and rollback, and system performance alarms.

Disaster recovery

Each data segment has multiple copies, which guarantees rapid restoration when one data segment is physically damaged.

Security

ECS supports security groups, Anti-DDoS, multi-user isolation, and password cracking defense.

Multiline access

ECS is based on the optimal routing algorithm of the Border Gateway Protocol (BGP). Multiline BGP data centers ensure smooth and balanced access throughout the geographic region. Backbone data centers ensure high output bandwidth and dedicated bandwidth.

Low cost

Large one-time payments are not required. Flexible payment options and Pay-As-You-Go let you cope with business changes.

Controllability

As an ECS user, you have the permission of a super administrator. This allows you to completely control the operating system of ECS instances, resolve system problems using the management terminal, and perform operations such as environment deployment and software installation.

Ease of use

A variety of operating systems and applications are supported. Images can be deployed with the click of a button. You can quickly replicate the environment to multiple ECS instances for easy scaling. You can also create ECS instances in batches using custom images and disk snapshots.

API

API invocation management allows configuration of access to one or multiple servers with the security group feature, making development more convenient.

Features

ECS supports the following features:

Flexible instance configuration

Supports multiple instance generations, three instance type families, and dozens of instance types (ranging from 1-core 1 GiB to 56-core 480 GiB).

Multiple regions and zones

Allows instance creation in all regions, some of which have multiple zones.

Abundant image resources

Provides various image resources, including public images, custom images, and shared images, allowing quick operating system deployment and applications without installation.

Numerous operating systems

Supports multiple Windows and Linux operating systems.

Multiple storage methods

Provides three types of data storage disks (Basic Cloud Disks, Ultra Cloud Disks, and SSD Cloud Disks) and I/O-optimized instances.

Robust network and security

- Supports two network types (Classic Network and VPC), allowing network management in different dimensions.
- Supports two types of IP addresses (public and private IP addresses), allowing for Intranet interconnection and Internet access.
- Allows free activation of Alibaba Cloud Security products and provides network monitoring.

Convenient management

Provides multiple management methods, including the console, VNC, and APIs, ensuring complete control.

Flexible payment

Provides flexible payment methods (Subscription and Pay-As-You-Go).

Compared with Internet Data Centers (IDCs) and server vendors, Alibaba Cloud adopts more stringent IDC standards, server access standards, and O&M standards to ensure data reliability and high availability of cloud computing infrastructure and cloud servers.

In addition, each region of Alibaba Cloud consists of multiple zones. For greater fault tolerance, you can build active/standby or active/active services in multiple zones. For a finance-oriented solution with three IDCs in two regions, you can build fault tolerant systems in multiple regions and zones. Those services include disaster tolerance and backup, which are supported by mature solutions of Alibaba Cloud. Services can be switched smoothly within Alibaba Cloud framework. For more information, see [E-Commerce Solutions](#). Alibaba Cloud industry solutions support a variety of services, such as finance, E-commerce, and video services. Alibaba Cloud provides you with the following support services:

- Products and services for availability improvement, including cloud servers, Server Load Balancer, multi-backup databases, and Data Transformation Services (DTS).
- Industry partners and ecosystem partners that help you build a more advanced and stable architecture and ensure service continuity.
- Diverse training services that enable you to connect with high availability from the business end to the underlying basic service end.

Users of cloud computing are most concerned about security and stability. Alibaba Cloud has recently passed a host of international information security certifications, including ISO 27001 and MTCS, which demand strict confidentiality of user data and user information as well as user privacy protection. Alibaba Cloud VPC is the prime choice for providing your cloud computing services.

Alibaba Cloud VPC offers more business possibilities. You only need to perform simple configuration to connect your business environment to global IDCs, making your business more flexible, stable, and extensible.

Alibaba Cloud VPC can connect your IDC through a leased line to build a hybrid cloud architecture. You can build more flexible business with the robust networking derived from Alibaba Cloud's various hybrid cloud solutions and network products. A superior business ecosystem is possible based on Alibaba Cloud's ecosystem.

Alibaba Cloud VPC is more stable and secure.

Stable: After you build your business on VPC, you can update your network architecture and obtain new network functions on a daily basis as the network infrastructure evolves constantly, allowing your business to run steadily. You can divide, configure, and manage your network on VPC according to your need.

Secure: VPC features traffic isolation and attack isolation protect your services from endless attack traffic on the Internet. After you build your business on VPC, the first line of defense is established.

VPC provides a stable, secure, fast-deliverable, self-managed, and controllable network environment. The capability and architecture of VPC hybrid cloud bring the technical advantages of cloud computing to traditional industries as well as industries and enterprises not engaged in cloud computing.

Regions

The following table lists the regions, corresponding cities, and Region IDs.

Regions in Mainland China

Region	China North 1	China North 2	China North 3	China North 5	China East 1	China East 2	China South 1
City	Qingdao	Beijing	Zhangjiakou	Huhehaote	Hangzhou	Shanghai	Shenzhen
RegionId	cn-qingdao	cn-beijing	cn-zhangjiakou	cn-huhehaote	cn-hangzhou	cn-shanghai	cn-shenzhen

International regions

Region	Hong Kong	Asia Pacific SE 1	Asia Pacific SE 2	Asia Pacific SE 3	Asia Pacific NE 1	US West 1	US East 1	Germany 1	Middle East 1
City	Hong Kong	Singapore	Sydney	Kuala Lumpur	Tokyo	Silicon Valley	Virginia	Frankfurt	Dubai
RegionId	cn-hongkong	ap-south-east-1	ap-south-east-2	ap-south-east-3	ap-north-east-1	us-west-1	us-east-1	eu-central-1	me-east-1



Zones

Zones are physical areas with independent power grids and networks in one region. The network latency for ECS instances within the same zone is shorter.

Intranet communication can take place between ECS instances in different zones of the same region, and fault isolation can be performed between zones. Whether ECS instances can be deployed in the

same zone depends on the requirements for disaster recovery capabilities and network latency.

- If your applications require high disaster recovery capabilities, we suggest you deploy your ECS instances in different zones of the same region;
- If your applications require low network latency between instances, we suggest you create your ECS instances in the same zone.

How to select a region

Regions in mainland China

Regions in mainland China include China East 1, China East 2, China North 1, China North 2, China North 3, China North 5, and China South 1.

They offer BGP backbone network lines covering all provinces and municipalities in mainland China and providing stable and fast access within mainland China.

They are similar to each other in terms of infrastructure, BGP network quality, service quality, ECS operation, and configuration. In general cases, we recommend you select a data center closest to your end users to further speed up user access.

International regions

Alibaba Cloud's international regions are data centers outside mainland China. They offer access at international bandwidth, targeting areas outside mainland China. Access to these regions from mainland China may experience high latency. Therefore, they are not recommended for access from mainland China.

Asia Pacific

Hong Kong

The data center in Hong Kong offers access at international bandwidth, covering Hong Kong and Southeast Asia.

If you have business operation in Hong Kong and Southeast Asia, you can select this region.

Asia Pacific SE 1

The data center in Asia Pacific SE 1 is provided by Alibaba Cloud's partner, SingTel, a dominant operator in Southeast Asia. With highly reliable business expertise and maturity, the company is well positioned to serve users across the region.

If you have business operation in Hong Kong and Southeast Asia, you can select this region.

Asia Pacific SE 2

The data center in Asia Pacific SE 2 is located in Sydney.

If you have business operation in Australia, select the Asia Pacific SE 2 region.

Asia Pacific SE 3

The data center in Asia Pacific SE 3 is located in Kuala Lumpur, Malaysia. It provides services of high speed and low latency for users in Southeast Asia countries, including Malaysia, Singapore, Thailand, and their neighboring countries.

- Asia Pacific NE 1**

The data center in Asia Pacific NE 1 is located in Tokyo, Japan.

If you have business operation in Japan, Northeast Asia, and South Korea, select the Asia Pacific NE 1 region.

North America and South America

US West 1

The data center in US West 1 is located in Silicon Valley. It is directly connected to the backbone networks of multiple American operators through BGP lines. In addition to the United States, this data center extends its reach to South America and Continental Europe.

If you have business operation in America and Continental Europe, select this region.

US East 1

The data center in US East 1 is located in Virginia of the United States.

If you have business operation in America and Continental Europe, select this region.

Germany 1

The data center in Germany 1 is located in Frankfurt.

If you have business operation in Continental Europe, select the Germany 1 region.

Middle East 1

The data center in Middle East 1 is located in Dubai.

If you have business operation in Middle East, select the Middle East 1 region.

Intranet communication between Alibaba Cloud products across regions

Intranet communication between Alibaba Cloud products that are not in the same region is not

supported, which means:

- ECS instances in different regions cannot communicate with each other on the intranet.
- ECS instances and other products in different regions, such as ApsaraDB for RDS and OSS instances, cannot communicate with each other on the intranet.
- Server Load Balancer cannot be deployed for ECS instances in various regions.

About business license record filing

If you must file your business license for record, pay attention to the following points:

- If your company is located in Beijing, select the **China North 2** region for the ECS instance you bought.
- If your company is located in Guangdong, select the **China South 1** region for the ECS instance you bought.

Notice: In mainland China, the approval requirements for business record filing vary depending on provincial communication management bureaus. For more information, see the information published on the business record filing website of your local bureau.

ECS is a highly flexible solution. It can be used independently as a simple web server, or used with other Alibaba Cloud products, such as OSS and CDN, to provide advanced solutions.

ECS can be used in applications such as:

Official corporate websites and simple web applications

In the initial stage, corporate websites have low traffic volumes and require only low-configuration ECS instances to run applications, databases, storage files, and other resources. As your business expands, you can upgrade the ECS configuration and increase the number of ECS instances at any time. You no longer need to worry about insufficient resources during peak traffic.

Multimedia and large-traffic apps or websites

ECS can be used with OSS to store static images, videos, and downloaded packages, reducing storage fees. In addition, ECS can be used with CDN or Server Load Balancer to greatly reduce user access waiting time, reduce bandwidth fees, and improve availability.

Databases

A high-configuration I/O-optimized ECS instance can be used with an SSD cloud disk to support high I/O concurrency with higher data reliability. Alternatively, multiple lower-configuration I/O-optimized ECS instances can be used with Server Load Balancer to deliver a high-availability architecture.

Apps or websites with large traffic fluctuations

Some applications may encounter large traffic fluctuations within a short period. When ECS is used with Auto Scaling, the number of ECS instances is automatically adjusted based on traffic. This feature allows you to meet resource requirements while maintaining a low cost. ECS can be used with Server Load Balancer to deliver a high availability architecture.

An ECS instance is the minimal unit that can provide computing services for your business. It provides computing capabilities at a certain specification.

The availability of instance type families and their types varies according to the regions and the amount of resources. Go to the [purchase page](#) to check the available instance types.

ECS instances are categorized into multiple specification types, which are called type families, based on the business and usage scenarios. In the same business scenario, you can select various type families. Each type family contains multiple instance types based on the CPU and memory specifications.

We define two basic attributes for an ECS instance type: the specifications of the CPU and the memory, including CPU model and clock speed. However, the attributes of a disk, an image, and the network service of an ECS instance must be defined simultaneously for the specific service form of the instance to be determined.

According to the release history and the business scenarios, Alibaba Cloud ECS instances are categorized into the following type families:

- The latest type families for various business scenarios, including:
 - Type families for enterprise-class computing on the x86-architecture, including:
 - General-purpose instance type families, including:
 - g5, general-purpose type family
 - sn2ne, general-purpose type family with enhanced network performance
 - sn2, general purpose type family
 - Compute instance type families, including:
 - c5, compute instance type family
 - sn1ne, compute optimized type family with enhanced network performance
 - sn1, compute optimized type family
 - Memory instance type families, including:
 - r5, memory instance type family
 - se1ne, memory optimized type family with enhanced network performance
 - se1, memory optimized type family
 - Big data instance type families, including:

- d1ne, big data type family with enhanced network performance
 - d1, big data type family
 - Instance type families with ephemeral storage, including:
 - i2, type family with ephemeral SSD disks
 - i1, type family with ephemeral SSD disks
 - Instance type families with high clock speed, including:
 - hfc5, compute optimized type family with high clock speed
 - hfg5, general-purpose type family with high clock speed
 - c4, cm4, and ce4, compute optimized type family with high clock speed
 - Type families for enterprise-class heterogeneous computing, including:
 - gn5, compute optimized type family with GPU
 - gn4, compute optimized type family with GPU
 - ga1, visualization compute type family with GPU
 - f1, compute optimized type family with FPGA
- Type families of previous generations for beginners, computing on the x86-architecture

The latest type families

All the ECS instances of the latest type families are I/O-optimized. They support the following disk types:

- SSD cloud disks
- Ultra cloud disks

Instances of the latest type families are categorized into the following type families based on the business scenarios.

Type families for enterprise-class computing on the x86-architecture

Type families for enterprise-class computing must be featured with stable performance and resource dedication. Within the enterprise-class instances, each vCPU core is supported by one Intel Xeon CPU core through hyper-threading.

g5, general-purpose type family

Features

- vCPU : Memory = 1:4
- Ultra high packet forwarding rate
- 2.5 GHz Intel Xeon Platinum 8163 (Skylake) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)

- Ideal for:
- Scenarios of receiving and transmitting a large volume of packets, such as video bullet screen and retransmission of telecommunication services
- Enterprise-class applications of various types and sizes
- Medium and small database systems, cache, and search clusters
- Data analysis and computing
- Computing clusters, and data processing depending on memory

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.g5.large	2	8.0	N/A	1.0	30	2
ecs.g5.xlarge	4	16.0	N/A	1.5	50	2
ecs.g5.2xlarge	8	32.0	N/A	2.5	80**	2
ecs.g5.4xlarge	16	64.0	N/A	5.0	100**	4
ecs.g5.6xlarge	24	96.0	N/A	7.5	150**	6
ecs.g5.8xlarge	32	128.0	N/A	10.0	200**	8
ecs.g5.16xlarge	64	256.0	N/A	20.0	400**	16

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

sn2ne, general purpose type family with enhanced network performance

Features

- vCPU : Memory = 1:4
- Ultra high packet forwarding rate
- 2.5 GHz Intel Xeon, E5-2682 v4 (Broadwell), or E5-2680 v3 (Haswell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Scenarios of receiving and transmitting a large volume of packets, such as video

- bullet screen and retransmission of telecommunication services
- Enterprise-class applications of various types and sizes
 - Medium and small database systems, cache, and search clusters
 - Data analysis and computing
 - Computing clusters, and data processing depending on memory

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.sn2ne.large	2	8.0	N/A	1.0	30	2
ecs.sn2ne.xlarge	4	16.0	N/A	1.5	50	2
ecs.sn2ne.2xlarge	8	32.0	N/A	2.0	100**	4
ecs.sn2ne.4xlarge	16	64.0	N/A	3.0	160**	4
ecs.sn2ne.8xlarge	32	128.0	N/A	6.0	250**	8
ecs.sn2ne.14xlarge	56	224.0	N/A	10.0	450**	16

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You can change the configurations of an instance between any two type families of sn2, sn2ne, sn1, sn1ne, se1, and se1ne, and within the same instance type family.

sn2, general purpose type family

Features

- vCPU : Memory = 1:4
- 2.5 GHz Intel Xeon, E5-2682 v4 (Broadwell), or E5-2680 v3 (Haswell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Enterprise-class applications of various types and sizes
 - Medium and small database systems, cache, and search clusters

- Data analysis and computing
- Computing clusters, and data processing depending on memory

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.sn2.medium	2	8.0	N/A	0.5	10	1
ecs.sn2.large	4	16.0	N/A	0.8	20	1
ecs.sn2.xlarge	8	32.0	N/A	1.5	40	1
ecs.sn2.3xlarge	16	64.0	N/A	3.0	50****	2
ecs.sn2.7xlarge	32	128.0	N/A	6.0	80***	3
ecs.sn2.13xlarge	56	224.0	N/A	10.0	120**	4

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Four queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Three queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

**** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You can change the configurations of an instance between any two type families of sn2, sn2ne, sn1, sn1ne, se1, and se1ne, and within the same instance type family.

c5, compute instance type family

Features

- vCPU : Memory = 1:2
- Ultra high packet forwarding rate
- 2.5 GHz Intel Xeon Platinum 8163 (Skylake) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:

- Scenarios of receiving and transmitting a large volume of packets, such as video bullet screen and retransmission of telecommunication services
- Web front-end servers
- Front ends of Massively Multiplayer Online (MMO) games
- Data analysis, batch compute, and video coding
- High performance science and engineering applications

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.c5.large	2	4.0	N/A	1.0	30	2
ecs.c5.xlarge	4	8.0	N/A	1.5	50	2
ecs.c5.2xlarge	8	16.0	N/A	2.5	80**	2
ecs.c5.4xlarge	16	32.0	N/A	5.0	100**	4
ecs.c5.6xlarge	24	48.0	N/A	7.5	150**	6
ecs.c5.8xlarge	32	64.0	N/A	10.0	200**	8
ecs.c5.16xlarge	64	128.0	N/A	20.0	400**	16

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

sn1ne, compute optimized type family with enhanced network performance

Features

- vCPU : Memory = 1:2
- Ultra high packet forwarding rate
- 2.5 GHz Intel Xeon, E5-2682 v4 (Broadwell), or E5-2680 v3 (Haswell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Scenarios of receiving and transmitting a large volume of packets, such as video bullet screen and retransmission of telecommunication services

- Web front-end servers
- Front ends of Massively Multiplayer Online (MMO) games
- Data analysis, batch compute, and video coding
- High performance science and engineering applications

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.sn1ne.large	2	4.0	N/A	1.0	30	2
ecs.sn1ne.xlarge	4	8.0	N/A	1.5	50	2
ecs.sn1ne.2xlarge	8	16.0	N/A	2.0	100**	4
ecs.sn1ne.4xlarge	16	32.0	N/A	3.0	160**	4
ecs.sn1ne.8xlarge	32	64.0	N/A	6.0	250**	8

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You can change the configurations of an instance between any two type families of sn2, sn2ne, sn1, sn1ne, se1, and se1ne, and within the same instance type family.

sn1, compute optimized type family

Features

- vCPU : Memory = 1:2
- 2.5 GHz Intel Xeon, E5-2682 v4 (Broadwell), or E5-2680 v3 (Haswell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Web front-end servers
 - Front ends of Massively Multiplayer Online (MMO) games
 - Data analysis, batch compute, and video coding
 - High performance science and engineering applications

Instance types

Instance	vCPU	Memory	Ephemeral	Intranet	Packet	NIC

type		(GiB)	al storage (GiB)	bandwidt h (Gbit/s)	forwardin g rate (10 thousand PPS)*	queues
ecs.sn1.m edium	2	4.0	N/A	0.5	10	1
ecs.sn1.la rge	4	8.0	N/A	0.8	20	1
ecs.sn1.xl arge	8	16.0	N/A	1.5	40	1
ecs.sn1.3x large	16	32.0	N/A	3.0	50**	2
ecs.sn1.7x large	32	64.0	N/A	6.0	80***	3

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Three queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You can change the configurations of an instance between any two type families of sn2, sn2ne, sn1, sn1ne, se1, and se1ne, and within the same instance type family.

r5, memory instance type family

Features

- Ultra high packet forwarding rate
- 2.5 GHz Intel Xeon Platinum 8163 (Skylake) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Scenarios of receiving and transmitting a large volume of packets, such as video bullet screen and retransmission of telecommunication services
 - High performance databases and memory databases
 - Data analysis and mining, and distributed memory cache
 - Hadoop, Spark, and other enterprise-class applications that require large volume of memory

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidt h (Gbit/s)	Packet forwardin g rate (10 thousand PPS)	NIC queues

					PPS)*	
ecs.r5.large	2	16.0	N/A	1.0	30	2
ecs.r5.xlarge	4	32.0	N/A	1.5	50	2
ecs.r5.2xlarge	8	48.0	N/A	2.5	80**	2
ecs.r5.4xlarge	16	64.0	N/A	5.0	100**	4
ecs.r5.6xlarge	24	128.0	N/A	7.5	150**	6
ecs.r5.8xlarge	32	256.0	N/A	10.0	200**	8
ecs.r5.16xlarge	64	512.0	N/A	20.0	400**	16
ecs.r5.22xlarge	88	704.0	N/A	30.0	450**	22

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

se1ne, memory optimized type family with enhanced network performance

Features

- vCPU : Memory = 1:8
- Ultra high packet receive and forwarding rate
- 2.5 GHz Intel Xeon, E5-2682 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Scenarios of receiving and transmitting a large volume of packets, such as video bullet screen and retransmission of telecommunication services
 - High performance databases and memory databases
 - Data analysis and mining, and distributed memory cache
 - Hadoop, Spark, and other enterprise-class applications that require large volume of memory

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand)	NIC queues

					PPS)*	
ecs.se1ne.large	2	16.0	N/A	1.0	30	2
ecs.se1ne.xlarge	4	32.0	N/A	1.5	50	2
ecs.se1ne.2xlarge	8	64.0	N/A	2.0	100**	4
ecs.se1ne.4xlarge	16	128.0	N/A	3.0	160**	4
ecs.se1ne.8xlarge	32	256.0	N/A	6.0	250**	8
ecs.se1ne.14xlarge	56	480.0	N/A	10.0	450**	16

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You can change the configurations of an instance between any two type families of sn2, sn2ne, sn1, sn1ne, se1, and se1ne, and within the same instance type family.

se1, memory optimized type family

Features

- vCPU : Memory = 1:8
- 2.5 GHz Intel Xeon, E5-2682 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - High performance databases and memory databases
 - Data analysis and mining, and distributed memory cache
 - Hadoop, Spark, and other enterprise-class applications that require large volume of memory

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.se1.large	2	16.0	N/A	0.5	10	1
ecs.se1.xl	4	32.0	N/A	0.8	20	1

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.se1.2xlarge	8	64.0	N/A	1.5	40	1
ecs.se1.4xlarge	16	128.0	N/A	3.0	50****	2
ecs.se1.8xlarge	32	256.0	N/A	6.0	80***	3
ecs.se1.14xlarge	56	480.0	N/A	10.0	120**	4

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Four queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Three queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

**** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You can change the configurations of an instance between any two type families of sn2, sn2ne, sn1, sn1ne, se1, and se1ne, and within the same instance type family.

d1ne, big data type family with enhanced network performance

Features

- High-volume ephemeral SATA HDD disks with high I/O throughput and a maximum of 35 Gbit/s of bandwidth for a single instance
- vCPU : Memory = 1:4, designed for big data scenarios
- 2.5 GHz Intel Xeon E5-2682 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Hadoop MapReduce, HDFS, Hive, HBase, and so on
 - Spark in-memory computing, MLlib, and so on
 - For those enterprises that require big data computing and storage analysis, such as enterprises in Internet and finance industries, to store and compute massive data
 - Elasticsearch, logs, and so on

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues

ecs.d1ne.2xlarge	8	32.0	4 * 5500	6.0	100**	4
ecs.d1ne.4xlarge	16	64.0	8 * 5500	12.0	160**	4
ecs.d1ne.6xlarge	24	96.0	12 * 5500	16.0	200**	6
ecs.d1ne.8xlarge	32	128.0	16 * 5500	20.0	250**	8
ecs.d1ne.14xlarge	56	224.0	28 * 5500	35.0	450**	14

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You cannot change configurations of d1ne instances.

For more information of d1ne type families, see [FAQ on d1 and d1ne](#).

d1, big data type family

Features

- High-volume ephemeral SATA HDD disks with high I/O throughput and a maximum of 17 Gbit/s of intranet bandwidth for a single instance
- vCPU : Memory = 1:4, designed for big data scenarios
- 2.5 GHz Intel Xeon E5-2682 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Hadoop MapReduce, HDFS, Hive, HBase, and so on
 - Spark in-memory computing, MLlib, and so on
 - For those enterprises that require big data computing and storage analysis, such as enterprises in Internet and finance industries, to store and compute massive data
 - Elasticsearch, logs, and so on

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.d1.2xlarge	8	32.0	4 * 5500	3.0	30	1
ecs.d1.4xl	16	64.0	8 * 5500	6.0	60****	2

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 Gbit/s)	NIC queues
ecs.d1.6xlarge	24	96.0	12 * 5500	8.0	80****	2
ecs.d1-c8d3.8xlarge	32	128.0	12 * 5500	10.0	100***	4
ecs.d1.8xlarge	32	128.0	16 * 5500	10.0	100***	4
ecs.d1-c14d3.14xlarge	56	160.0	12 * 5500	17.0	180**	6
ecs.d1.14xlarge	56	224.0	28 * 5500	17.0	180**	6

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Six queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Four queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

**** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You cannot change configurations of d1 instances.

For more information of d1 type family, see [FAQ on d1 and d1ne](#).

i2, type family with ephemeral SSD disks

Features

- High-performance ephemeral NVMe SSD disks: supporting high IOPS and I/O throughput and low latency.
- vCPU : Memory = 1:8, designed for high performance databases
- 2.5 GHz Intel Xeon Platinum 8163 (Skylake) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - OLTP and high performance relational databases
 - NoSQL databases, such as Cassandra and MongoDB
 - Search applications, such as Elasticsearch

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 Gbit/s)	NIC queues

					thousand PPS)*	
ecs.i2.xlarge	4	32.0	1 * 894	1.0	50	2
ecs.i1.2xlarge	8	64.0	1 * 1788	2.0	100**	2
ecs.i1.4xlarge	16	128.0	2 * 1788	3.0	150**	4
ecs.i1.8xlarge	32	256.0	4 * 1788	6.0	200**	8
ecs.i1.16xlarge	64	512.0	8 * 1788	10.0	400**	16

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You cannot change configurations of i2 instances.

i1, type family with ephemeral SSD disks

Features

- High-performance ephemeral NVMe SSD disks: supporting high IOPS and I/O throughput and low latency.
- vCPU : Memory = 1:4, designed for high performance databases
- 2.5 GHz Intel Xeon E5-2682 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - OLTP and high performance relational databases
 - NoSQL databases, such as Cassandra and MongoDB
 - Search applications, such as Elasticsearch

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.i1.xlarge	4	16.0	2 * 104	0.8	20	1
ecs.i1.2xlarge	8	32.0	2 * 208	1.5	40	1

ecs.i1.4xlarge	16	64.0	2 * 416	3.0	50****	2
ecs.i1-c5d1.4xlarge	16	64.0	2 * 1456	3.0	40****	2
ecs.i1.8xlarge	32	128.0	2 * 832	6.0	80***	3
ecs.i1-c10d1.8xlarge	32	128.0	2 * 1456	6.0	80***	3
ecs.i1.14xlarge	56	224.0	2 * 1456	10.0	120**	4

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Four queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Three queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

**** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You cannot change configurations of i1 instances.

hfc5, compute optimized type family with high clock speed

Features

- 3.1 GHz Intel Xeon Gold 6149 (Skylake) processors
- vCPU : Memory = 1:2
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - High performance Web front-end servers
 - High performance science and engineering applications
 - Massively Multiplayer Online (MMO) games and video coding

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.hfc5.large	2	4.0	N/A	1.0	30	2
ecs.hfc5.x	4	8.0	N/A	2.5	50	2

large						
ecs.hfc5.2 xlarge	8	16.0	N/A	5.0	100	2
ecs.hfc5.4 xlarge	16	32.0	N/A	8.0	160	4
ecs.hfc5.6 xlarge	24	48.0	N/A	12.0	240	6
ecs.hfc5.8 xlarge	32	64.0	N/A	16.0	320	8

* For more information about PPS testing, see [Test network performance](#). Multiple NIC queues must be enabled for testing the instances with more than four vCPU cores.

You can change the configurations of an instance within hfc5, and between hfc5 and hfg5 families.

hfg5, general-purpose type family with high clock speed

Features

- 3.1 GHz Intel Xeon Gold 6149 (Skylake) processors
- vCPU : Memory = 1:4, except that of the instance type with 56 vCPU cores
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - High performance Web front-end servers
 - High performance science and engineering applications
 - Massively Multiplayer Online (MMO) games and video coding

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.hfg5.large	2	8.0	N/A	1.0	30	2
ecs.hfg5.xlarge	4	16.0	N/A	2.5	50	2
ecs.hfg5.2xlarge	8	32.0	N/A	5.0	100	2
ecs.hfg5.4xlarge	16	64.0	N/A	8.0	160	4
ecs.hfg5.6	24	96.0	N/A	12.0	240	6

xlarge						
ecs.hfg5.8xlarge	32	128.0	N/A	16.0	320	8
ecs.hfg5.14xlarge	56	160.0	N/A	28.0	450	14

* For more information about PPS testing, see [Test network performance](#). Multiple NIC queues must be enabled for testing the instances with more than four vCPU cores.

You can change the configurations of an instance within hfg5, and between hfc5 and hfg5 families.

c4, cm4, and ce4, compute optimized type family with high clock speed

Features

- 3.2 GHz Intel Xeon E5-2667 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - High performance Web front-end servers
 - High performance science and engineering applications
 - Massively Multiplayer Online (MMO) games and video coding

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.c4.xlarge	4	8.0	N/A	1.5	20	1
ecs.c4.2xlarge	8	16.0	N/A	3.0	40	1
ecs.c4.4xlarge	16	32.0	N/A	6.0	80***	2
ecs.cm4.xlarge	4	16.0	N/A	1.5	20	1
ecs.cm4.2xlarge	8	32.0	N/A	3.0	40	1
ecs.cm4.4xlarge	16	64.0	N/A	6.0	80***	2
ecs.cm4.6xlarge	24	96.0	N/A	10.0	120**	4

ecs.ce4.xlarge	4	32.0	N/A	1.5	20	1
----------------	---	------	-----	-----	----	---

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Four queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You can change the configurations of an instance within c4, cm4, and ce4.

Type families for enterprise-class heterogeneous computing

gn5, compute optimized type family with GPU

Features

- NVIDIA P100 GPU processors
- No fixed ratio of vCPU to memory
- High performance ephemeral NVMe SSD disks
- 2.5 GHz Intel Xeon E5-2682 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Deep learning
 - Scientific computing, such as computational fluid dynamics, computational finance, genomics, and environmental analysis
 - High performance computing, rendering, multi-media coding and decoding, and other server-side GPU compute workloads

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	GPU	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.gn5 - c4g1.xlarge	4	30.0	440	1 * NVIDIA P100	3.0	30	1
ecs.gn5 - c8g1.2xlarge	8	60.0	440	1 * NVIDIA P100	3.0	40	1

ecs.gn5 - c4g1.2xl arge	8	60.0	880	2 * NVIDIA P100	5.0	100***	2
ecs.gn5 - c8g1.4xl arge	16	120.0	880	2 * NVIDIA P100	5.0	100***	4
ecs.gn5 - c28g1.7 xlarge	28	112.0	440	1 * NVIDIA P100	5.0	100**	8
ecs.gn5 - c8g1.8xl arge	32	240.0	1760	4 * NVIDIA P100	10.0	200**	8
ecs.gn5 - c28g1.1 4xlarge	56	224.0	880	2 * NVIDIA P100	10.0	200**	14
ecs.gn5 - c8g1.14 xlarge	54	480.0	3520	8 * NVIDIA P100	25.0	400**	14

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: No more than (vCPU core number/4) queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You cannot change configurations of gn5 instances.

gn4, compute optimized type family with GPU

Features

- NVIDIA M40 GPU processors
- No fixed ratio of CPU to memory
- 2.5 GHz Intel Xeon E5-2682 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Deep learning
 - Scientific computing, such as computational fluid dynamics, computational finance, genomics, and environmental analysis

- High performance computing, rendering, multi-media coding and decoding, and other server-side GPU compute workloads

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	GPU	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.gn4-c4g1.xlarge	4	30.0	N/A	1 * NVIDIA M40	3.0	30	1
ecs.gn4-c8g1.2xlarge	8	60.0	N/A	1 * NVIDIA M40	3.0	40	1
ecs.gn4.8xlarge	32	48.0	N/A	1 * NVIDIA M40	6.0	80***	3
ecs.gn4-c4g1.2xlarge	8	60.0	N/A	2 * NVIDIA M40	5.0	50	1
ecs.gn4-c8g1.4xlarge	16	60.0	N/A	2 * NVIDIA M40	5.0	50****	1
ecs.gn4.14xlarge	56	96.0	N/A	2 * NVIDIA M40	10.0	120**	4

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Four queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Three queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

**** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

See [Create a gn4 instance](#) in the [ECS User Guide](#).

You can change the configurations of an instance within the gn4 family.

ga1, visualization compute type family with GPU

Features

- AMD S7150 GPU processors
- vCPU : Memory = 1:2.5
- High performance ephemeral NVMe SSD disks
- 2.5 GHz Intel Xeon E5-2682 v4 (Broadwell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Rendering, multimedia coding and decoding
 - Machine learning, high-performance computing, and high performance databases
 - Other server-end business scenarios that require powerful concurrent floating-point compute capabilities

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	GPU	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.ga1.2xlarge	8	20.0	1 * 175	0.5 * AMD S7150	1.5	30	1
ecs.ga1.4xlarge	16	40.0	1 * 350	1 * AMD S7150	3.0	50****	2
ecs.ga1.8xlarge	32	80.0	1 * 700	2 * AMD S7150	6.0	80***	3
ecs.ga1.14xlarge	56	160.0	1 * 1400	4 * AMD S7150	10.0	120**	4

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Four queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

*** Testing conditions: Three queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

**** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You cannot change configurations of ga1 instances.

f1, compute optimized type family with FPGA

Features

- Intel Arria 10 GX 1150 FPGA
- vCPU : Memory = 1:7.5
- 2.5 GHz Intel Xeon E5-2682 v4 (Broadwell) processors
- High performance ephemeral NVMe SSD disks
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- Ideal for:
 - Deep learning and reasoning
 - Genomics research and finance analysis
 - Computational workloads, such as real-time video processing and security

Instance types

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	FPGA	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.f1-c8f1.2xlarge	8	60.0	440	Intel Arria 10 GX 1150	3.0	40	4
ecs.f1-c8f1.4xlarge	16	120.0	880	Intel Arria 10 GX 1150 * 2	4.0	200**	4

* For more information about PPS testing, see [Test network performance](#).

** Testing conditions: Two queues for NICs are enabled and CentOS 7.3 is used. If you want to adjust multiple-queue for NIC, restart the instance.

You cannot change configurations of f1 instances.

Type families of previous generations for beginners, computing on the x86-architecture

xn4/n4/mn4/e4, shared instance type families

Features

- 2.5 GHz Intel Xeon E5-2682 v4 (Broadwell) processors
- The latest DDR4 memory
- No fixed ratio of CPU to memory

Instance types

Type family	Features	vCPU : Memory	Idea for
xn4	Compact shared instances	1:1	<ul style="list-style-type: none"> - Front ends of Web applications - Light load applications and microservices - Applications for development or testing environments
n4	General shared instances	1:2	<ul style="list-style-type: none"> - Websites and Web applications - Development environment, building servers, code repositories, microservices, and testing and staging environments - Lightweight enterprise applications

mn4	Balanced shared instances	1:4	- Websites and Web applications - Lightweight databases and cache - Integrated applications and lightweight enterprise services
e4	Memory shared instances	1:8	- Applications that require large volume of memory - Lightweight databases and cache

You can change the configurations of an instance between any two type families of xn4, n4, mn4, and e4, and within the same instance type family.

xn4

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.xn4.small	1	1.0	N/A	0.5	5	1

* For more information about PPS testing, see [Test network performance](#).

n4

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues

					PPS)*	
ecs.n4.small	1	2.0	N/A	0.5	5	1
ecs.n4.large	2	4.0	N/A	0.5	10	1
ecs.n4.xlarge	4	8.0	N/A	0.8	15	1
ecs.n4.2xlarge	8	16.0	N/A	1.2	30	1
ecs.n4.4xlarge	16	32.0	N/A	2.5	40	1
ecs.n4.8xlarge	32	64.0	N/A	5.0	50	1

* For more information about PPS testing, see [Test network performance](#).

mn4

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.mn4.small	1	4.0	N/A	0.5	5	1
ecs.mn4.large	2	8.0	N/A	0.5	10	1
ecs.mn4.xlarge	4	16.0	N/A	0.8	15	1
ecs.mn4.2xlarge	8	32.0	N/A	1.2	30	1
ecs.mn4.4xlarge	16	64.0	N/A	2.5	40	1

* For more information about PPS testing, see [Test network performance](#).

e4

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)	Intranet bandwidth (Gbit/s)	Packet forwarding rate (10 thousand PPS)*	NIC queues
ecs.e4.small	1	8.0	N/A	0.5	5	1

* For more information about PPS testing, see [Test network performance](#).

n1/n2/e3, shared instance type families

Features

- 2.5 GHz Intel Xeon E5-2680 v3 (Haswell) processors
- The network performance of an instance matching the computing type (the more advanced the computing type, the more powerful the network performance)
- I/O-optimized
- Supporting the following disk types:
 - SSD cloud disks
 - Ultra cloud disks

Instance types

Type family	Features	vCPU : Memory	Idea for
n1	General shared instances	1:2	<ul style="list-style-type: none"> - Small and medium-sized web servers - Batch processing - Distributed analysis - Advertisement services
n2	Balanced shared instances	1:4	<ul style="list-style-type: none"> - Medium-sized Web servers - Batch processing - Distributed analysis - Advertisement services - Hadoop clusters
e3	Memory shared instances	1:8	<ul style="list-style-type: none"> - Cache, Redis - Search

			<ul style="list-style-type: none"> - Memory databases - Databases with high I/O, for example, Oracle and MongoDB - Hadoop clusters - Computing scenarios that involve massive data processing
--	--	--	---

n1

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)
ecs.n1.tiny	1	1.0	N/A
ecs.n1.small	1	2.0	N/A
ecs.n1.medium	2	4.0	N/A
ecs.n1.large	4	8.0	N/A
ecs.n1.xlarge	8	16.0	N/A
ecs.n1.3xlarge	16	32.0	N/A
ecs.n1.7xlarge	32	64.0	N/A

n2

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)
ecs.n2.small	1	4.0	N/A
ecs.n2.medium	2	8.0	N/A
ecs.n2.large	4	16.0	N/A
ecs.n2.xlarge	8	32.0	N/A
ecs.n2.3xlarge	16	64.0	N/A
ecs.n2.7xlarge	32	128.0	N/A

e3

Instance type	vCPU	Memory (GiB)	Ephemeral storage (GiB)
ecs.e3.small	1	8.0	N/A
ecs.e3.medium	2	16.0	N/A
ecs.e3.large	4	32.0	N/A
ecs.e3.xlarge	8	64.0	N/A
ecs.e3.3xlarge	16	128.0	N/A

You can change the configurations among the three shared instance type families (n1, n2, and e3), and within the same instance type family.

If you are using t1, s1, s2, s3, m1, m2, c1, or c2, see [Generation I instance types](#).

Instances

An ECS instance is a virtual computing environment that includes CPU, memory, operating system, bandwidth, disks, and other basic computing components. An ECS instance is an independent virtual machine, and is the core element of ECS. Other resources, such as disks, IPs, images, and snapshots can only be used in conjunction with an ECS instance.

The life cycle of an instance begins when you create the instance and ends when the instance is released, either after a monthly or yearly subscription expires, when you manually release a Pay-As-You-Go instance, or because of an outstanding payment.

The life cycle of an instance involves several inherent instance states, as listed in the following table.

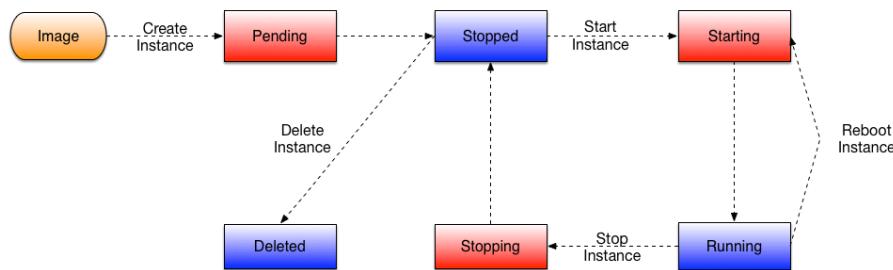
State	Type of state	Description	Corresponding API state
Preparing*	Intermediate state	After an instance is created, it remains in this state before running.	Pending
Created*	Stable state	An instance is in this state when it has been created and is awaiting start up.	Stopped
Starting*	Intermediate state	An instance is in this	Starting

		state after it is started or restarted in the console or using API until it is running.	
Running	Stable state	The instance is operating properly and can accommodate your business needs.	Running
Stopping*	Intermediate state	An instance is in this state after the stop operation is performed in the console or through API but before it actually stops.	Stopping
Stopped	Stable state	The instance has been stopped properly. In this state, the instance cannot accommodate external services.	Stopped
Re-initializing*	Intermediate state	An instance is in this state after the system disk and/or data disk is re-initialized in the console or using API until it is running.	Stopped
Replacing System Disk	Intermediate state	An instance is in this state after the operating system is replaced or another such operation is performed in the console or using API until it is running.	Stopped
Expired	Stable state	<p>The yearly/monthly instance subscription has expired because it has not been properly renewed. The Pay-As-You-Go instances have expired because of overdue payments.</p> <p>Note: After expiration, both the yearly/monthly and Pay-As-You-Go</p>	Stopped

		instances will continue running for 15 days, and data will be retained for an extra 15 days, after which the instances will be released and the data will be removed permanently.	
--	--	---	--

* If an instance remains in the Preparing, Starting, Stopping, Re-initializing, or Replacing System Disk state for a long time, it may encounter an exception.

The following illustration describes the transition between API instance states.



Burstable instances (also called t5 instances) can handle sudden rise of requirements for CPU performance. Each t5 instance provides a baseline CPU performance. The instance type determines the rate at which CPU credits are obtained. When your CPU usage is below the baseline performance, the corresponding CPU credits accumulate, than being consumed. When your CPU usage exceeds the baseline performance, the instance consumes the accumulated CPU credits to meet the elevated performance requirements. t5 instances seamlessly increase your CPU performance, without affecting the instance environment or applications.

t5 instances are ideal for scenarios where you usually do not require high CPU performance, but occasionally require a high computing performance, such as lightweight web servers, development and testing environments, and a low or mid-performance databases.

How t5 instances work

You must have a basic understanding of the following concepts, before you use t5 instances:

Baseline CPU performance

The instance type of any t5 instance determines its baseline CPU performance, which means each vCPU core of an instance has a maximum usage for normal workloads. For example,

when a t5-1c1m2.small instance is used for normal workloads, the maximum CPU usage is 10%.

CPU credits

Each t5 instance obtains CPU credits at a fixed distribution rate, which is determined by the baseline CPU performance. A CPU credit is a measurement unit of computing performance, which is determined by the number of vCPU core, CPU usage, and work time:

- 1 CPU credit = 1 vCPU core at 100% usage for 1 minute
- 1 CPU credit = 1 vCPU core at 50% usage for 2 minutes
- 1 CPU credit = 2 vCPU cores at 25% usage for 2 minutes

If one vCPU core runs at 100% usage for one hour (60 minutes), it consumes 60 CPU credits.

CPU credit distribution rate

The CPU credit distribution rate is the number of CPU credits that a t5 instance obtains per minute. It is determined by the baseline CPU performance. You can use the following formula to determine the CPU credit distribution rate by the baseline CPU performance:

$$\text{CPU credit distribution rate} = (60 \text{ CPU credits} * \text{Baseline CPU performance}) / 60 \text{ minutes}$$

Example: A t5 instance of the t5-1c1m2.small type provides a baseline CPU performance of 10%, so the CPU credit distribution rate is 0.1 CPU point per minute, or six CPU credits per hour.

Initial CPU credits

Every time you create a new t5 instance, 30 CPU credits are immediately allocated to the instance, which is called as initial CPU credits. Instances are only allocated with initial CPU credits at the time of creation. When an instance begins to consume CPU credits, the initial CPU credits are used first.

CPU credit accumulation

When the CPU usage of a t5 instance is lower than or equal to the baseline CPU performance, the instance accumulates CPU credits, but does not consume them. You can view the CPU credits of a t5 instance in the ECS console.

CPU credit consumption

When workload of a t5 instance bursts to a level above the baseline CPU performance, the instance consumes accumulated CPU credits to raise the CPU usage to meet your business requirements.

When you want to use one vCPU at 100% usage for one minute, the number of the consumed CPU credits can be calculated by using the following formula:

CPU credits consumed per minute = 1 CPU point * (100% - Baseline CPU performance)

Example: A t5 instance of the t5-lc1m2.small type provides a baseline CPU performance of 10%. When this instance is used at 100% usage for 1 minute, it consumes 0.9 CPU credits.

When the CPU credit accumulation rate exceeds the consumption rate, the instance continues accumulating CPU credits. If the consumption rate exceeds the accumulation rate, the total instance points decrease in number. Once accumulated, the CPU credits are saved only for 24 hours. The credit points are no longer valid after 24 hours.

However, if the instance is stopped, the existing CPU credits are still valid and the credit accumulation continues. And, when you restart the instance, the instance continues to accumulate CPU credits.

If the instance runs out-of-service due to overdue payments, the CPU credits still remain valid, but the CPU credit accumulation stops. After the instance is relaunched, the accumulation resumes automatically.

Instance types

t5 instances use Intel Xeon processors. The instance types are shown in the following table. In this table:

- CPU credits/hour is the total number of CPU credits allocated to all vCPU cores of a t5 instance per hour.
- Average baseline CPU performance is the average baseline CPU performance of each vCPU core for a t5 instance.

Instance type	vCPU	CPU credits/hour	Avg baseline CPU performance	Memory (GiB)
t5-lc2m1.nano	1	6	10%	0.5
t5-lc1m1.small	1	6	10%	1
t5-lc1m2.small	1	6	10%	2
t5-lc1m2.large	2	12	10%	4
t5-lc1m4.large	2	12	10%	8
t5-c1m1.large	2	18	15%	2
t5-c1m2.large	2	18	15%	4
t5-c1m4.large	2	18	15%	8
t5-c1m1.xlarge	4	36	15%	4
t5-c1m2.xlarge	4	36	15%	8
t5-c1m4.xlarge	4	36	15%	16
t5-	8	72	15%	8

c1m1.2xlarge				
t5-c1m2.2xlarge	8	72	15%	16
t5-c1m4.2xlarge	8	72	15%	32

Here, we use t5-c1m1.xlarge as an example to explain the t5 instance configuration:

Each vCPU core has an average baseline computing performance of 15%. Therefore, the total baseline computing performance of a t5 instance of the t5-c1m1.xlarge type is 60%, which means:

- If the instance only uses one vCPU core, this core has a baseline computing performance of 60%.
- If the instance only uses two vCPU cores, each core is allocated with a baseline computing performance of 30%.
- If the instance only uses three vCPU cores, each core is allocated with a baseline computing performance of 20%.
- If the instance uses all four vCPU cores, each core is allocated with a baseline computing performance of 15%.

One instance is allocated with 36 CPU credits per hour, which means that each vCPU core is allocated with nine CPU credits per hour.

Billing method

t5 instances support the following billing methods: Pay-As-You-Go and Subscription. For more information on the billing methods, see [Purchase ECS instances](#).

Create an instance

See [Create an ECS instance](#) to create a t5 instance. When creating a t5 instance, consider the following settings:

- Region: Now only Zone E of China North 2 (Beijing) supports t5 instances.
- Network type: Only VPC is supported.
- Image and instance type: The minimum t5 instance memory specification of 512 MiB only supports Linux. To create a Windows instance, the minimum memory is 1 GiB. For more information on image selection, see [How to select a system image](#).

Manage t5 instances

View CPU usage

In the ECS console, go to the Monitoring Information section of the **Instance Details** page to view the instance CPU usage. You can also remotely connect to the ECS instance to view CPU usage.

The CPU usage information displayed in the ECS console and in the instance may differ:

- When a t5 instance has CPU credits available, the CPU usage information is same for both the methods.
- When a t5 instance does not have available CPU credits:
 - If the vCPU usage is lower than the baseline CPU performance, the CPU usage information is the same for both methods.
 - If the vCPU usage exceeds or is equal to the baseline CPU performance, the CPU usage information displayed in the console shows the baseline CPU performance, but the information in the instance shows the actual workload.

View CPU usage in the ECS console

To view CPU usage in the ECS console, follow these steps:

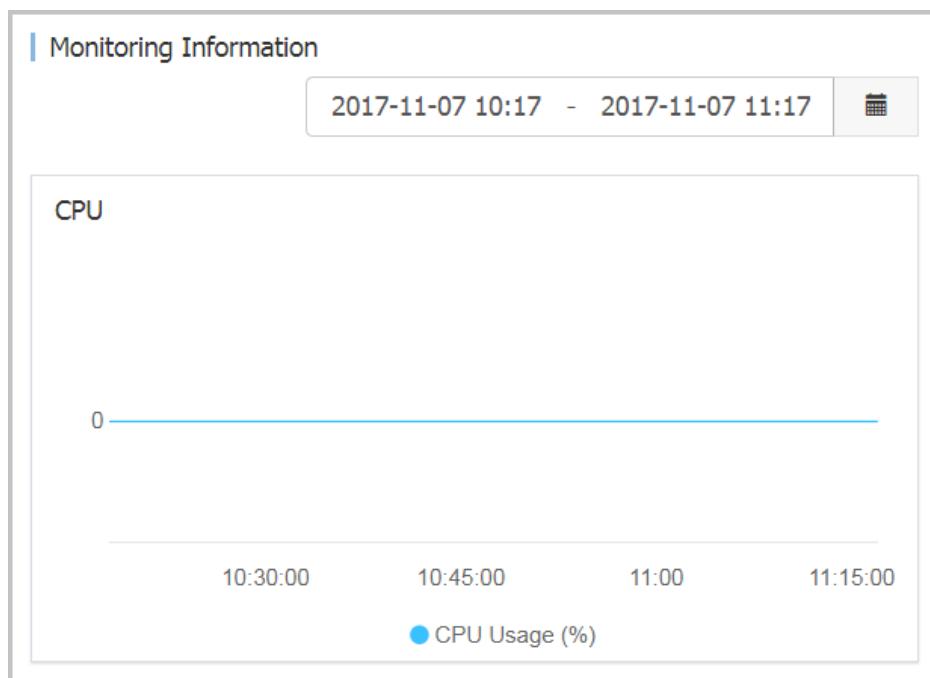
Log on to the ECS console.

On the left-side navigation pane, click Instances.

Select a region.

Find a t5 instance, and click the instance ID or in the **Actions** column, click **Manage**.

In the **Monitoring Information** section, view CPU usage information.



Remotely connect to the instance to view CPU usage

The methods vary as per the operating system:

- Windows: Connect to the instance and view the information in the **Task Manager**.
- Linux: Connect to the instance and run the top command to view the CPU usage.

Change configurations

In the ECS console, if you see that the CPU usage is at the baseline level of CPU performance for an extended period of time or it never exceeds the baseline level, your current instance type is either insufficient for your needs or exceeds your needs. In these cases, consider changing the instance type.

If your instance billing method is Subscription, you can upgrade the instance type. For more information, see [Upgrade configurations](#).

If your instance billing method is Pay-As-You-Go, you can change the instance type.

View CPU credits

Log on to the ECS console and go to the [Instance Details](#) page to view the accumulated and consumed CPU credits of a t5 instance.

Disks

An ECS disk can be used jointly or separately to meet the requirements of different application scenarios. ECS disks are categorized into ephemeral SSD disks and cloud disks. Compared with ephemeral SSD disks, cloud disks are more reliable because they use a triplicate distributed system to provide block-level data storage for ECS instances, ensuring 99.9999999% data reliability. Cloud disks are categorized as one of the following:

SSD cloud disks

Ideal for I/O intensive applications, and provide stable and high random IOPS performance.

Ultra cloud disks

Ideal for application scenarios with medium I/O load and provide a storage performance of up to 3000 random IOPS for ECS instances.

Basic cloud disks

Ideal for application scenarios with low I/O load and provide an I/O performance of several hundred IOPS for ECS instances.

Note: For detailed instructions on attaching a disk, refer to *Attach a data disk* from *User Guide of ECS*.

Disk comparison

The following table lists the features and typical application scenarios of different types of cloud disks.

Item	SSD cloud disk	Ultra cloud disk	Basic cloud disk
Maximum capacity	32,768 GB	32,768 GB	2,000 GB
Maximum IOPS	20,000 To calculate the maximum IOPS of a disk: $IOPS = \min\{1200 + 30 * capacity, 20000\}$	3,000 To calculate the maximum IOPS of a disk: $IOPS = \min\{1000 + 6 * capacity, 3000\}$	Several hundreds
Maximum throughput	300 MBps To calculate the maximum throughput of a disk: $Throughput = \min\{80 + 0.5 * capacity, 300\} MBps$	80 MBps To calculate the maximum throughput of a disk: $Throughput = \min\{50 + 0.1 * capacity, 80\} MBps$	30 MBps

Access latency	0.5–2 ms	1–3 ms	5–10 ms
Data reliability	99.9999999%	99.9999999%	99.9999999%
API name	cloud_ssd	cloud_efficiency	cloud
Price*	\$0.15 USD/GB/month	\$0.08 USD/GB/month	\$0.05 USD/GB/month
Typical application scenarios	<ul style="list-style-type: none"> - I/O intensive applications - Medium/Large relational databases - NoSQL databases 	<ul style="list-style-type: none"> - Medium/Small databases - Large-scale development and testing - Web server logs 	Applications with infrequent access or low I/O load

* Prices shown are for the US West region. For more information, see [ECS Price](#) for details.

For more information about ephemeral SSD disks, see [ephemeral SSD disks](#).

Test disk performance

You can use fio to test the performance of a cloud disk.

Warning:

Testing naked disks can obtain more accurate performance data, but will damage the file structure. Make sure that you back up your data before testing. We recommend that you use a new ECS instance without data on the disks to test the disks by using fio.

Make sure that disks are 4K aligned before performing the following tests:

Test random writes IOPS

```
fio -direct=1 -iodepth=128 -rw=randwrite -ioengine=libaio -bs=4k -size=1G -numjobs=1 -runtime=1000 -group_reporting -filename=iotest -name=Rand_Write_Testing
```

Test random reads IOPS

```
fio -direct=1 -iodepth=128 -rw=randread -ioengine=libaio -bs=4k -size=1G -numjobs=1 -runtime=1000 -group_reporting -filename=iotest -name=Rand_Read_Testing
```

Test writes throughput

```
fio -direct=1 -iodepth=64 -rw=write -ioengine=libaio -bs=64k -size=1G -numjobs=1 -runtime=1000 -group_reporting -filename=iotest -name=Write_PPS_Testing
```

Test reads throughput

```
fio -direct=1 -iodepth=64 -rw=read -ioengine=libaio -bs=64k -size=1G -numjobs=1 -runtime=1000 -group_reporting -filename=iotest -name=Read_PPS_Testing
```

Take the command for testing random reads IOPS as an example to describe the meaning of the parameters of a fio command, as shown in the following table.

Parameter	Description
-direct=1	Ignore I/O buffer when testing. Data is written directly.
-rw=randwrite	Read and write policies. Available options: randread (random read), randwrite(random write), read(sequential read), write(sequential write), and randrw (random read and write).
-ioengine=libaio	Use libaio as the testing method (Linux AIO, Asynchronous I/O). Usually there are two ways for an application to use I/O: synchronous and asynchronous. Synchronous I/O only sends out one I/O request each time, and returns only after the kernel is completed. In this case, the iodepth is always less than 1 for a single job, but can be resolved by multiple concurrent jobs. Usually 16–32 concurrent jobs can fill up the iodepth. Asynchronous method uses libaio to submit a batch of I/O requests each time, thus reduces interaction times, and makes interaction more effective.
-bs=4k	The size of each block for one I/O is 4k. If not specified, the default value 4k is used. When IOPS is tested, we recommend that you set the bs to a small value, such as 4k in this example command. When throughput is tested, we recommend that you set the bs to a big value, such as 1024k in the IOPS tests.
-size=1G	The size of the testing file is 1 GB.
-numjobs=1	The number of testing jobs is 1.
-runtime=1000	Testing time is 1000 seconds. If not specified, the test will go on with the value specified for -size, and write data in -bs each time.
-group_reporting	The display mode of showing the testing results. Group_reporting means sums up statistics of each job, instead of showing statistics by different jobs.
-filename=iotest	The output path and name of the test files. Testing naked disks can obtain more accurate performance data, but will damage the file

	structure. Make sure that you back up your data before testing.
-name=Rand_Write_Testing	The name of the testing task.

Take the output of a random reads IOPS test on an SSD cloud disk of 800 GB capacity as an example, we describe how to explain the test report.

```
Rand_Read_Testing: (g=0): rw=randread, bs=4K-4K/4K-4K/4K-4K, ioengine=libaio, iodepth=128
fio-2.2.8
Starting 1 process
Jobs: 1 (f=1): [r(1)] [21.4% done] [80000KB/0KB/0KB /s] [20.0K/0/0 iops] [eta 00Jobs: 1 (f=1): [r(1)] [28.6% done]
[80000KB/0KB/0KB /s] [20.0K/0/0 iops] [eta 00Jobs: 1 (f=1): [r(1)] [35.7% done] [80000KB/0KB/0KB /s] [20.0K/0/0
iops] [eta 00Jobs: 1 (f=1): [r(1)] [42.9% done] [80004KB/0KB/0KB /s] [20.1K/0/0 iops] [eta 00Jobs: 1 (f=1): [r(1)]
[50.0% done] [80004KB/0KB/0KB /s] [20.1K/0/0 iops] [eta 00Jobs: 1 (f=1): [r(1)] [57.1% done] [80000KB/0KB/0KB /s]
[20.0K/0/0 iops] [eta 00Jobs: 1 (f=1): [r(1)] [64.3% done] [80144KB/0KB/0KB /s] [20.4K/0/0 iops] [eta 00Jobs: 1 (f=1):
[r(1)] [71.4% done] [80388KB/0KB/0KB /s] [20.1K/0/0 iops] [eta 00Jobs: 1 (f=1): [r(1)] [78.6% done]
[80232KB/0KB/0KB /s] [20.6K/0/0 iops] [eta 00Jobs: 1 (f=1): [r(1)] [85.7% done] [80260KB/0KB/0KB /s] [20.7K/0/0
iops] [eta 00Jobs: 1 (f=1): [r(1)] [92.9% done] [80016KB/0KB/0KB /s] [20.4K/0/0 iops] [eta 00Jobs: 1 (f=1): [r(1)]
[100.0% done] [80576KB/0KB/0KB /s] [20.2K/0/0 iops] [eta 00m:00s]
Rand_Read_Testing: (groupid=0, jobs=1): err= 0: pid=9845: Tue Sep 26 20:21:01 2017
read : io=1024.0MB, bw=80505KB/s, iops=20126, runt= 13025msec
slat (usec): min=1, max=674, avg= 4.09, stdev= 6.11
clat (usec): min=172, max=82992, avg=6353.90, stdev=19137.18
lat (usec): min=175, max=82994, avg=6358.28, stdev=19137.16
clat percentiles (usec):
| 1.00th=[ 454], 5.00th=[ 668], 10.00th=[ 812], 20.00th=[ 996],
| 30.00th=[ 1128], 40.00th=[ 1256], 50.00th=[ 1368], 60.00th=[ 1480],
| 70.00th=[ 1624], 80.00th=[ 1816], 90.00th=[ 2192], 95.00th=[79360],
| 99.00th=[81408], 99.50th=[81408], 99.90th=[82432], 99.95th=[82432],
| 99.99th=[82432]
bw (KB /s): min=79530, max=81840, per=99.45%, avg=80064.69, stdev=463.90
lat (usec) : 250=0.04%, 500=1.49%, 750=6.08%, 1000=12.81%
lat (msec) : 2=65.86%, 4=6.84%, 10=0.49%, 20=0.04%, 100=6.35%
cpu : usr=3.19%, sys=10.95%, ctx=23746, majf=0, minf=160
IO depths : 1=0.1%, 2=0.1%, 4=0.1%, 8=0.1%, 16=0.1%, 32=0.1%, >=64=100.0%
submit : 0=0.0%, 4=100.0%, 8=0.0%, 16=0.0%, 32=0.0%, 64=0.0%, >=64=0.0%
complete : 0=0.0%, 4=100.0%, 8=0.0%, 16=0.0%, 32=0.0%, 64=0.0%, >=64=0.1%
issued : total=r=262144/w=0/d=0, short=r=0/w=0/d=0, drop=r=0/w=0/d=0
latency : target=0, window=0, percentile=100.00%, depth=128

Run status group 0 (all jobs):
READ: io=1024.0MB, aggrb=80504KB/s, minb=80504KB/s, maxb=80504KB/s, mint=13025msec, maxt=13025msec

Disk stats (read/write):
vdb: ios=258422/0, merge=0/0, ticks=1625844/0, in_queue=1625990, util=99.30%
```

In the result, you must pay the most attention to the following code line:

```
read : io=1024.0MB, bw=80505KB/s, iops=20126, runt= 13025msec
```

The line means that fio did a total of 1 GB of IO at 80 MB/s for a total of 20126 IOPS (at the default 4k

block size), and ran for 13 seconds.

According to the formula, the maximum IOPS of an SSD cloud disk of 800 GB capacity is: IOPS = min{1200+30 * capacity, 20000} = IOPS = min{1200+30 * 800, 20000} = 20000

The calculated IOPS is similar to the value in the output report, 20126.

Alibaba Cloud disks can be categorized into SSD cloud disks, ultra cloud disks, and basic cloud disks. This document describes the features and application scenarios of each disk category.

SSD cloud disks

Features

SSD cloud disks use a distributed, triplicate mechanism to provide high-performance storage with stable and high random I/O and high data reliability. They provide the following features:

High random I/O performance

The maximum random read/write IOPS is 20,000. The base is 1,200 IOPS, and each GB of capacity provides 30 random IOPS. For example, a 100 GB SSD cloud disk can provide 4,200 IOPS, and a 334 GB SSD cloud disk can provide 11,220 IOPS.

High throughput

The maximum throughput is 300 MBps. The throughput of an SSD cloud disk can be calculated using the formula $\min\{80 + 0.5 * \text{capacity}, 300\}$ MBps.

High data reliability

SSD cloud disks adopt a distributed triplicate mechanism to provide 99.999999% data reliability.

Large storage capacity

A single SSD cloud disk provides up to 32,768 GB storage space.

Independent attaching

SSD cloud disks can be attached to any ECS instance in the same zone.

Note:

Expected IOPS performance can be achieved only when the SSD cloud disk is attached to an I/O-optimized instance. An SSD cloud disk that is attached to a non I/O-optimized instance cannot achieve the expected IOPS performance.

Performance baselines

The performance of SSD cloud disks varies by the data block size. A smaller data block causes higher IOPS and smaller throughput.

Block size	Maximum IOPS	Maximum throughput
4 KB or 8 KB	20,000	Small, far below 300 MBps
16 KB	About 17,200	
32 KB	About 9,600	Almost 300 MBps
64 KB	About 4,800	

Use cases

SSD cloud disks have stable and high random I/O performance, and high data reliability. They are idea for the following scenarios:

- PostgreSQL, MySQL, Oracle, SQL Server, and other medium/large relational database applications.
- Medium to large development and testing environments with high requirements for data reliability.

Ultra cloud disks

Features

Ultra cloud disks adopt the hybrid media of SSD and HDD as the storage media. They provide the following features:

High random I/O performance

The maximum random read/write IOPS is 3,000. The random read/write IOPS is initially 1,000 and increases by 6 IOPS for each GB. For example, a 250 GB ultra cloud disk features 2,500 random read/write IOPS.

High throughput

The maximum throughput is 80 MBps. The throughput is initially 50 MBps and increases by

0.1 MBps for each GB. For example, a 250 GB ultra cloud disk features a throughput of 75 MBps.

High data reliability

Ultra cloud disks adopt a distributed, triplicate mechanism to provide 99.999999% data reliability.

Large storage capacity

A single ultra cloud disk provides up to 32,768 GB storage space.

Independent attaching

Ultra cloud disks can be attached to any ECS instance in the same zone.

Use cases

Ultra cloud disks are idea for the following scenarios:

- MySQL, SQL Server, PostgreSQL, and other small or medium relational database applications.
- Medium or large development and testing environments with high requirements for data reliability and intermediate performance.

Basic cloud disks

Features

Basic cloud disks adopt HDDs as the storage medium and use a distributed, triplicate mechanism to provide high data reliability. They provide the following features:

High random I/O performance

The maximum random read/write IOPS is of several hundreds.

High throughput

The maximum throughput is 30 MBps–40 MBps.

High data reliability

Disks adopt a distributed triplicate mechanism provides 99.999999% data reliability.

Large storage capacity

A single basic cloud disk provides up to 2,000 GB storage space.

Independent attaching

Basic cloud disks can be attached to any ECS instance in the same zone.

Use cases

Basic cloud disks are idea for the following scenarios:

- Applicable to scenarios in which data is not frequently accessed, or which have low I/O loads.
If an application requires higher I/O performance, we recommend that you use an SSD cloud disk.
- Application environments that require low costs and have random I/O read/write.

Introduction to triplicate technology

The Alibaba Cloud Distributed File System provides stable, efficient, and reliable random data access capabilities for ECS.

Chunks

When ECS users perform read and write operations onto virtual disks, the operations are translated into corresponding processes on the files stored in the Alibaba Cloud Distributed File System. Alibaba Cloud provides a flat storage space, in which the linear addresses are divided into chunks, also referred to as slices. Alibaba Cloud employs a certain strategy to create three copies for each chunk and stores these copies on different nodes, ensuring the reliability of user data.

Principles of triplicate technology

The Alibaba Cloud data storage system consists of three roles: Master, Chunk Server, and Client. The write operation of an ECS user goes through several conversions and is executed by the Client. The procedure is as follows:

1. The Client calculates the chunk corresponding to a given writing operation.
2. The Client sends a request to the Master for the storage location of the three copies of the chunk.
3. The Client sends writing requests to the three Chunk Servers according to the results returned from the Master.

4. The Client returns a message to the user indicating whether the operation was successful.

The distribution strategy of the Master is decided based on an overall consideration of the following:

- Disk use conditions of all Chunk Servers in the cluster.
- Distribution of the Chunk Servers under different kinds of switch racks.
- The power supply.
- The instrument load.

This strategy ensures that all the copies of a Chunk are distributed on different Chunk Servers on different racks. This can effectively prevent data unavailability caused by the failure of a Chunk Server or rack.

Data protection mechanism

When some data nodes are corrupted, or some hard drives on a certain data node fail, the number of valid copies of some Chunks in the cluster will be less than three. If this occurs, the Master initiates the copy mechanism to copy data between Chunk Servers, making three valid copies of all Chunks in the cluster.

In sum, for the data on the cloud disk, all user operations and data addition or modification will be synchronized to the three copies. This mode ensures the reliability and consistency of user data.

To prevent data losses caused by virus infection or cyber-attacks, we recommend that you use the triplicate technology with other protection methods, such as taking snapshots.

Network and security

Intranet

Currently, Alibaba Cloud servers communicate through the intranet. They use a gigabit of shared bandwidth for non I/O optimized instances, and 10 gigabits of shared bandwidth for I/O Optimized instances, with no special restrictions. However, because this is a shared network, the bandwidth speed may fluctuate.

If you need to transmit data between two ECS instances in the same region, you should use an intranet connection. Intranet connections can also be used to connect RDS, Server Load Balancer, and OSS instances. The internet speed of these instances is based on a gigabit shared bandwidth environment. At present, you can also use a direct intranet connection to link RDS, Server Load Balancer, and OSS instances with ECS instances in the same region.

For ECS instances in the intranet:

For instances of Classic network:

- Intranet communication is by default used only for instances in the same security group of the same account in the same region.
- An intranet communication can also be used for instances in the same security group of the same account and region but of different zones, even if the intranet IP addresses are in different network segments.
- For intranet communication between instances in the same region but of different accounts, you can use security groups. For more information, see [Application scenarios of security group from ECS User Guide](#).

For instances of VPC network:

- Intranet communication is by default used only for instances in the same security group of the same account and same VPC network in the same region.
- An intranet communication can also be used for instances of the same account and region but of different VPC networks only if you use ExpressConnect to authorize their intranet communication. For more information, see [Application scenarios from Product Introduction to ExpressConnect](#).

The intranet IP addresses of instances cannot be modified or changed.

Intranet and Internet addresses of instances do not support virtual IP (VIP) configuration.

Instances of different network types cannot communicate with each other in intranet.

IP addresses for Classic network

IP addresses are an important means for users to access ECS instances, and for ECS instances to provide external services. Currently, classic IP addresses are uniformly distributed by Alibaba Cloud. They are divided into public and private IP addresses.

Private IP addresses

An instance is allocated with a private network card and bound to a specific private IP address. Private IP addresses are required and cannot be modified.

If a private IP address is changed independently in an operating system, communication in the private network will be interrupted.

Communication traffic through private IP addresses between instances in the same region is free. Private IP addresses can be used in the following scenarios:

- Load balancing of the Server Load Balancer
- Intranet mutual access between ECS instances
- Intranet mutual access between an ECS instance and another cloud service (such as OSS and RDS)

Public IP addresses

Each instance is by default configured with a public network interface card. Unlike private IP addresses, public IP addresses are optional. If you select a public network bandwidth greater than 0 Mbps when purchasing an instance, a public IP address will be allocated during creation of the instance.

Regardless of your selected billing method, you must select a public network bandwidth limit. The bandwidth limit you select will determine the limit of the outgoing bandwidth for the public network card.

Public network traffic will be charged. Public IP addresses can be used in the following scenarios:

- Mutual access between an ECS instance and the Internet
- Mutual access between an ECS instance and another cloud service

Multicast and Broadcast

ECS does not support multicast or broadcast.

Security groups

A security group is a logical group that groups instances in the same region with the same security requirements and mutual trust. Each instance belongs to at least one security group, which must be specified at the time of creation. Instances in the same security group can communicate through the network, but instances in different security groups by default cannot communicate through an intranet. However, mutual access can be authorized between two security groups.

A security group is a virtual firewall that provides stateful packet inspection (SPI). Security groups are used to set network access control for one or more ECSs. As an important means of security isolation, security groups are used to divide security domains on the cloud.

Security group restrictions

A single security group cannot contain more than 1,000 instances. If you require intranet mutual

access between more than 1,000 instances, you can allocate them to different security groups and permit mutual access through mutual authorization.

- Each instance can join up to five security groups.
- Each user can have up to 100 security groups.
- Adjusting security groups will not affect the continuity of user service.
- Security groups are stateful. If an outbound packet is permitted, inbound packets corresponding to this connection will also be permitted.
- Security groups have two network types: classic network and Virtual Private Cloud (VPC).
 - Classic Network type instances can join security groups on classic networks in the same region.
 - VPC type instances can join security groups on the same VPC.

Security group rules

Security group rules can be set to permit or forbid ECS instances associated with security groups to access a public network or an intranet from inbound and outbound directions.

You can authorize or delete security group rules at any time. Security group rules you have changed will automatically apply to ECS instances associated with security groups.

When setting security group rules, make sure security group rules are simple. If you associate an ECS instance with multiple security groups, up to hundreds of rules may apply to the instance, which may cause connection errors when you access the instance.

Security group rule restrictions

Each security group can have a maximum of 100 security group rules.

Alibaba Cloud offers two authentication methods for remote logon to ECS instances:

- Password logon: A standard authentication method using the administrator password. It applies to both Windows instances and Linux instances.
- SSH key pair logon: This method only applies to Linux instances. If you are running Linux, it is recommended that you choose this authentication method to protect your ECS instance's security.

An SSH key pair is a pair of keys generated through an encryption algorithm: one key is intentionally available, known as the **public key**, and the other key is kept confidential, known as the **private key**.

If you have placed the public key in a Linux instance, you can use the private key to log on to the instance through using SSH commands or related tools from local computer or another instance, without the need to enter a password.

Benefits

SSH key pairs provides the following benefits:

High security:

- The security strength of a key pair is much higher than that of user passwords. A key pair can hinder brute force password-cracking attacks.
- It is impossible to deduce the private key even if the public key is maliciously acquired.

Ease-of-use: You can log on to the instance remotely through simple configuration in the console and on the local client. No password is required for the logon next time. If you need to maintain multiple ECS instances in batch, this logon method is recommended.

Alibaba Cloud SSH key pairs

To generate key pairs, you can use either of the following methods:

- Use Alibaba Cloud to generate key pairs. Alibaba Cloud uses 2048-bit RSA keys by default.
- Import the public key of a key pair that has been generated by another key pair generation tool. The key pair must be one of the following types:

- rsa
- dsa
- ssh-rsa
- ssh-dss
- ecdsa
- ssh-rsa-cert-v00@openssh.com
- ssh-dss-cert-v00@openssh.com
- ssh-rsa-cert-v01@openssh.com
- ssh-dss-cert-v01@openssh.com
- ecdsa-sha2-nistp256-cert-v01@openssh.com
- ecdsa-sha2-nistp384-cert-v01@openssh.com
- ecdsa-sha2-nistp521-cert-v01@openssh.com

If you generate your key pair using Alibaba Cloud, you must download the private key immediately after the creation and keep it secure. If you do not have the private key, you will not be able to log on to the ECS instance that is bound to this key pair.

You can allocate a key pair to an instance when you create the Linux instance, or you can allocate it after the instance is created.

If you use an SSH key pair to log on to a Linux instance, password authentication is disabled by default to improve instance security.

Limits for using SSH keys are as follows:

- Only Linux instances are supported. Windows instances are not supported.
- An account can have a maximum of 500 key pairs in a region.
- A Linux instance can only bind one SSH key pair. If your instance has already bound a key pair, the new key pair will replace the original key pair.
- Within the lifecycle of a Linux instance, you can re-bind the SSH key pair and instance. After re-binding, the key pair will take effect without the need to restart the instance.
- All instances of any instance type family, except those I/O optimized instances of Generation I, support SSH key pairs.

Operate an SSH key pair

- If you do not have an SSH key pair, you can create an SSH key pair.
- If you have had an SSH key pair generated by another tool, you can import an SSH key pair.
- If you no longer require an SSH key pair, you can delete an SSH key pair.
- If you want to enable or disable SSH key pair authentication for logging on to a Linux ECS instance, you can bind or unbind an SSH key pair.
- Allocate an SSH key pair when creating an ECS instance.
- Log on to an instance using an SSH key pair.

If you use a single CPU, chances of many network interruptions increase. For processing, you can route NIC interruptions in the ECS instances to different CPUs. In the network PPS and network bandwidth tests, a solution that uses two queues instead of one queue can enhance the performance by 50% to 100%. A solution that uses four queues can bring significant increase in the performance.

ECS instance types supporting multi-queue

See [Instance generations and type families](#) to find instance types supporting multi-queue and the number of queues that are supported.

Images supporting multi-queue

The public images officially provided by Alibaba Cloud shown in the following table support multi-queue. Whether an image supports multi-queue is not related to the memory address width of the operating system.

Image	Notes
Windows 2012 R2	Unavailable. You may be invited to test this feature in the future.
Windows 2016	Unavailable. You may be invited to test this feature in the future.

CentOS 6.8/6.9/7.2/7.3/7.4	None
Ubuntu 14.04/16.04	None
Debian 8.9	None
SUSE Linux Enterprise Server 12 SP1	None
SUSE Linux Enterprise Server 12 SP2	Available soon

Configure multi-queue support for NICs on a Linux ECS instance

We recommend that you use one of the latest Linux distributions, such as CentOS 7.2, to configure multi-queue for the NICs.

Here we take CentOS 7.2 as an example to illustrate how to configure multi-queue for the NIC. In this example, we want to configure two queues, and the NIC name is eth0.

- To check whether the NIC supports multi-queue, run the command: ethtool -l eth0.
- To enable multi-queue for the NIC, run the command: ethtool -L eth0 combined 2.

If you are using more than one NIC, configure each NIC.

```
[root@localhost ~]# ethtool -l eth0
Channel parameters for eth0:
Pre-set maximums:
RX: 0
TX: 0
Other: 0
Combined: 2 # This line indicates that a maximum of two queues can be configured
Current hardware settings:
RX: 0
TX: 0
Other: 0
Combined: 1 #It indicates that one queue is currently taking effect

[root@localhost ~]# ethtool -L eth0 combined 2 # It sets eth0 to use two queues currently
```

We recommend that you enable the irqbalance service so that the system can automatically adjust the allocation of the NIC interrupts on multiple CPU cores. Run the command: systemctl start irqbalance. This feature is enabled by default in CentOS 7.2.

If the network performance improvement is not as high as you expected after the multi-queue feature is enabled, you can enable the RPS feature. See the following Shell script.

```
#!/bin/bash
cpu_num=$(grep -c processor /proc/cpuinfo)
```

```
quotient=$((cpu_num/8))
if [ $quotient -gt 2 ]; then
quotient=2
elif [ $quotient -lt 1 ]; then
quotient=1
fi
for i in $(seq $quotient)
do
cpuset="${cpuset}f"
done

for rps_file in $(ls /sys/class/net/eth*/queues/rx-*/rps_cpus)
do
echo $cpuset > $rps_file
done
```

Configure multi-queue support for NICs on a Windows ECS instance

Note:

We are inviting Windows users to test the performance improvement.

Windows systems see improved network performance after using multi-queue for NICs, but the improvement is not as good as seen in Linux systems.

If you are using a Windows instance, you must install the driver to use the multi-queue feature for NICs.

To install the driver for Windows systems, follow these steps:

Open a ticket to request and download the driver installation package.

Unzip the driver installation package. For Windows 2012/2016 systems, use the driver in the **Win8/amd64** folder.

Upgrade the NIC driver:

- i. Select **Device Manager > Network adapters**.
- ii. Right click **Red Hat VirtIO Ethernet Adapter** and select **Update Driver...**
- iii. Select the **Win8/admin64** directory of the driver directory that you unzipped, and update the driver.

After the driver is upgraded, we recommend that you restart the Windows system.

The multi-queue feature for NICs is ready to use.

Images

An image is a running environment template for ECS instances. It generally includes an operating system and preinstalled software. You can use an image to create an ECS instance or change the system disk of an ECS instance.

ECS allows you to easily obtain an image in the following ways:

- Choosing a public image officially provided by Alibaba Cloud (multiple Windows and Linux versions are available).
- Creating a custom image based on an existing ECS instance.
- Choosing an image shared by another Alibaba Cloud account.

You can import an offline image file into an ECS cluster to generate a custom image.

You can also copy a custom image to another region to maintain a consistent environment and application deployment across multiple regions.

Snapshots

A snapshot is a copy of data on a disk at a certain point in time. Scheduled creation of disk snapshots ensures continuous operation of your business. Snapshot is a simple and efficient data protection method, and is recommended for the following scenarios:

Routine backup of system and data disks

You can back up business-critical data at regular intervals using snapshots to prevent data loss from misoperations, attacks, and viruses.

OS replacement

Before important operations such as upgrading application software or migrating business data, you need to create one or more snapshots. In case of any issues occurring during the upgrade or migration, you can restore timely to normal status using the snapshots.

Use of multiple copies of production data

You can take snapshots of production data to provide close-to-real-time production data for

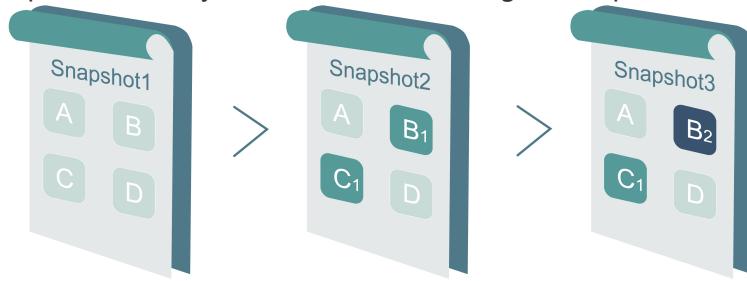
data mining, report queries, and developing and testing applications. You can also take snapshots to reuse data on a disk as basic data for another disk.

Restoring data stored on cloud disks

Cloud disks provide a secure storage method to ensure that your stored content will not be lost. However, if the data stored on a cloud disk is incorrect due to an application error, or the data is maliciously tampered by hackers through an application vulnerability, snapshots ensure that your data can be restored to the desired state.

Incremental snapshot mechanism

Snapshots on Alibaba Cloud are taken using an incremental method. In this method, two snapshots are compared and only the data that has changed is copied, as shown in the following image:



In the preceding figure, Snapshot 1, Snapshot 2, and Snapshot 3 are the first, second, and third snapshots of a disk. The file system checks the disk data by blocks. When a snapshot is created, only the blocks with changed data will be copied to the snapshot. In this example:

- In Snapshot 1, all data on the disk is copied since it is the first disk snapshot.
- Snapshot 2 only copies the changed data blocks B₁ and C₁. Data blocks, A and D, are referenced from Snapshot 1.
- Snapshot 3 copies the changed data block B₂ but references data blocks, A, D, from Snapshot 1, and references C₁ from Snapshot 2.
- When you roll back the disk to Snapshot 3, blocks A, B₂, C₁, and D are copied to the disk, to replicate Snapshot 3.
- When you delete Snapshot 2, block B₁ will be deleted, but C₁ will remain because blocks that are referenced by other snapshots cannot be deleted. When you roll back to Snapshot 3, block C₁ will be recovered.

When the disk needs to be restored to the status at the time of Snapshot 3, you can perform snapshot rollback to copy data blocks A, B₂, C₁, and D to the disk.

If Snapshot 2 is deleted, data block B₁ in the snapshot will be deleted but data block C₁ will not be deleted. In this way, when the disk is restored to the status at the time of Snapshot 3, data block C₁

can also be restored.

Snapshot creation time varies depending on actual volume. For a frame of reference, it typically takes several minutes to manually create a 40 GB snapshot.

Snapshots are stored on the Object Storage Service (OSS), but they are invisible to users and will not be computed in the OSS space occupied by the users' buckets. Snapshot operations can only be performed through the ECS console or APIs.

ECS Snapshot 2.0

Built on original basic snapshot features, ECS Snapshot 2.0 data backup service provides a higher snapshot quota and more flexible automatic task policies, further reducing its impact on business IO. The features of ECS Snapshot 2.0 are described in the following table.

Feature	Original snapshot specifications	Snapshot 2.0 specifications	User benefit
Snapshot quota	(Number of disks)*6+6	64 snapshots for each disk	Longer protection circle Smaller protection granularity
Automatic task policy	Hardcoded, triggered once daily, and unmodifiable	Customizable weekly snapshot day, time of day, and snapshot retention period Query-able disk quantity and related details associated with an automatic snapshot policy	More flexible protection policy
Implementation principle	COW (Copy-on-write)	ROW (Redirect-on-write)	Mitigated performance impact of the snapshot task on business IO write

The implementation of ECS Snapshot 2.0 features is described in the following table.

Feature	Implementation
Snapshot quota	Snapshot backup of a data disk for non-core businesses occurs at 00:00 every day. This backup data is retained for over 2 months. Snapshot backup of a data disk for core businesses occurs every 4 hours. This backup data is retained for over 10 days.
Automatic task policy	A user can take snapshots on the hour and for several times in a day. A user can choose any day as the recurring

	day for taking weekly snapshots. A user can specify the snapshot retention period or choose to retain it permanently (When the maximum number of automatic snapshots has been reached, the oldest automatic snapshot will be deleted).
Implementation principle	The implementation principle is not made visible to users, allowing snapshots to be taken at any time of day without affecting user experience.

ECS Snapshot 2.0 vs. traditional storage products

Alibaba Cloud ECS Snapshot 2.0 has many advantages compared with the snapshot feature of traditional storage products, as described in the following table.

Comparison item	ECS Snapshot 2.0	Snapshot feature of traditional storage products
Capacity limit	Unlimited capacity, meeting data protection needs for extra-large businesses.	Capacity limited by initial storage device capacity, merely meeting data protection needs for a few core services.
Scalability	One-click auto scaling, allowing you to scale up and down according to their business scale, in mere seconds.	Poor scalability, restrained by factors such as production and storage performance, available capacity, and vendor support capabilities. Scaling typically takes 1 ~ 2 weeks.
Cost	Billed based on the actual amount of data changed in your business and snapshot size.	Large, inefficient upfront investment involving software licenses, reserved space, and upgrade and maintenance expenses.
Usability	24x7 online post-sales support.	Complex operations, greatly restrained by vendor support capabilities.