

Alibaba Cloud Elastic Compute Service

Best Practices

Issue: 20190217

Legal disclaimer

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








1. You shall download and obtain this document from the Alibaba Cloud website or other Alibaba Cloud-authorized channels, and use this document for your own legal business activities only. The content of this document is considered confidential information of Alibaba Cloud. You shall strictly abide by the confidentiality obligations. No part of this document shall be disclosed or provided to any third party for use without the prior written consent of Alibaba Cloud.
2. No part of this document shall be excerpted, translated, reproduced, transmitted, or disseminated by any organization, company, or individual in any form or by any means without the prior written consent of Alibaba Cloud.
3. The content of this document may be changed due to product version upgrades, adjustments, or other reasons. Alibaba Cloud reserves the right to modify the content of this document without notice and the updated versions of this document will be occasionally released through Alibaba Cloud-authorized channels. You shall pay attention to the version changes of this document as they occur and download and obtain the most up-to-date version of this document from Alibaba Cloud-authorized channels.
4. This document serves only as a reference guide for your use of Alibaba Cloud products and services. Alibaba Cloud provides the document in the context that Alibaba Cloud products and services are provided on an "as is", "with all faults" and "as available" basis. Alibaba Cloud makes every effort to provide relevant operational guidance based on existing technologies. However, Alibaba Cloud hereby makes a clear statement that it in no way guarantees the accuracy, integrity, applicability, and reliability of the content of this document, either explicitly or implicitly. Alibaba Cloud shall not bear any liability for any errors or financial losses incurred by any organizations, companies, or individuals arising from their download, use, or trust in this document. Alibaba Cloud shall not, under any circumstances, bear responsibility for any indirect, consequential, exemplary, incidental, special, or punitive damages, including lost profits arising from the use

or trust in this document, even if Alibaba Cloud has been notified of the possibility of such a loss.

5. By law, all the content of the Alibaba Cloud website, including but not limited to works, products, images, archives, information, materials, website architecture, website graphic layout, and webpage design, are intellectual property of Alibaba Cloud and/or its affiliates. This intellectual property includes, but is not limited to, trademark rights, patent rights, copyrights, and trade secrets. No part of the Alibaba Cloud website, product programs, or content shall be used, modified, reproduced, publicly transmitted, changed, disseminated, distributed, or published without the prior written consent of Alibaba Cloud and/or its affiliates. The names owned by Alibaba Cloud shall not be used, published, or reproduced for marketing, advertising, promotion, or other purposes without the prior written consent of Alibaba Cloud. The names owned by Alibaba Cloud include, but are not limited to, "Alibaba Cloud", "Aliyun", "HiChina", and other brands of Alibaba Cloud and/or its affiliates, which appear separately or in combination, as well as the auxiliary signs and patterns of the preceding brands, or anything similar to the company names, trade names, trademarks, product or service names, domain names, patterns, logos, marks, signs, or special descriptions that third parties identify as Alibaba Cloud and/or its affiliates).
6. Please contact Alibaba Cloud directly if you discover any errors in this document.

Generic conventions

Table -1: Style conventions

Style	Description	Example
	This warning information indicates a situation that will cause major system changes, faults, physical injuries, and other adverse results.	 Danger: Resetting will result in the loss of user configuration data.
	This warning information indicates a situation that may cause major system changes, faults, physical injuries, and other adverse results.	 Warning: Restarting will cause business interruption. About 10 minutes are required to restore business.
	This indicates warning information, supplementary instructions, and other content that the user must understand.	 Notice: Take the necessary precautions to save exported data containing sensitive information.
	This indicates supplemental instructions, best practices, tips, and other content that is good to know for the user.	 Note: You can use Ctrl + A to select all files.
>	Multi-level menu cascade.	Settings > Network > Set network type
Bold	It is used for buttons, menus, page names, and other UI elements.	Click OK.
Courier font	It is used for commands.	Run the <code>cd /d C:/windows</code> command to enter the Windows system folder.
<i>Italics</i>	It is used for parameters and variables.	<code>bae log list --instanceid Instance_ID</code>
[] or [a b]	It indicates that it is an optional value, and only one item can be selected.	<code>ipconfig [-all -t]</code>

Style	Description	Example
<code>{}</code> or <code>{a b}</code>	It indicates that it is a required value, and only one item can be selected.	<code>switch {stand slave}</code>

Contents

Legal disclaimer.....	I
Generic conventions.....	I
1 Security.....	1
1.1 Best practices of the security group (part 1).....	1
1.2 Best practices of the security group (part 2).....	4
1.3 Best practices of the security group (part 3).....	10
1.4 Best practices of ECS data security.....	14
1.5 How to configure instances to access each other in classic networks.....	17
1.6 Modify the default remote access port.....	22
1.7 Use logs in Windows instances.....	28
1.8 Overview and best practices of Windows Firewall with Advanced Security.....	30
1.9 Isolation of instances within a security group.....	45
1.10 Security group quintuple rules.....	48
2 Disaster recovery solutions.....	52
3 Data recovery.....	57
3.1 How to restore the data that is deleted by mistake.....	57
3.2 Data restoration in Linux instances.....	61
3.3 Data restoration in Windows instances.....	68
4 Configuration preference.....	72
4.1 Time setting: Synchronize NTP servers for Windows instances.....	72
4.2 Time setting: NTP servers and other public services.....	74
4.3 Configure language settings for multiple instances.....	75
4.4 Time setting: Synchronize NTP servers and change time zone for Linux instances.....	78
5 Monitor.....	81
5.1 Use CloudMonitor to monitor ECS instances.....	81
5.2 Automatically manage instances.....	83
6 User-defined data.....	90
6.1 User-defined yum sources, NTP services and DNS services.....	90
6.2 Create a new account with the root user privilege.....	92
7 GPU instances.....	95
7.1 Deploy an NGC on gn5 instances.....	95
7.2 Install a GRID driver on a gn5/gn5i/gn6v instance.....	99
8 FaaS instances best practices.....	109
8.1 Use RTL compiler on an f1 instance.....	109
8.2 Use OpenCL on an f1 instance.....	112
8.3 Best practices for OpenCL on an f3 instance.....	117

8.4 Best practices for RTL design on an f3 instance.....	125
8.5 faascmd tool.....	131
8.5.1 faascmd overview.....	131
8.5.2 Install faascmd.....	131
8.5.3 Configure faascmd.....	132
8.5.4 Use faascmd.....	133
8.5.5 FAQ.....	140
9 Access other Cloud Product APIs by the Instance RAM Role.....	146
10 Shrink disk.....	152
11 Terraform.....	155
11.1 What is Terraform?.....	155
11.2 Install and configure Terraform.....	156
11.3 Create an ECS instance.....	157
11.4 Create multiple ECS instances.....	160
11.5 Deploy a Web cluster.....	162

1 Security

1.1 Best practices of the security group (part 1)

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

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

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

- [Security group restrictions](#)
- [Default security group rules](#)
- [Set the inbound access of a security group](#)
- [Set the outbound access of a security group](#)

General suggestions for security group practices

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

- The most important rule: A security group should be used as a whitelist.
- The "minimum authorization" principle should be observed when you configure the inbound or outbound rules for applications. For example, you can allow a specific port (such as port 80).
- It is not recommended to use one security group to manage all applications, because requirements must be different at different layers.
- For distributed applications, different security groups should be used for different application types. For example, you should use different security groups for the Web, Service, Database and Cache layers to apply different inbound/outbound rules and permissions.
- There is no need to set a separate security group for every instance, as this would unnecessarily add to management costs.
- VPC should be preferred.

- Do not assign Internet addresses to resources that require no Internet access.
- Keep the rules of each security group as concise as possible. A single instance can join up to five security groups, and a security group can contain up to 100 security group rules, so an instance may be subject to hundreds of security group rules at the same time. You can aggregate all the assigned security rules to determine whether inbound or outbound traffic is permitted or not. However, overly complicated rules for a single security group can increase management complexity. For this reason, it is recommended to keep the rules of each security group as concise as possible.
- The ECS console allows you to clone a security group and security group rules. If you want to modify an active security group and its rules, you should clone the security group and modify the cloned security group, avoiding any impacts on online applications.

**Note:**

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

Set inbound access rules of security groups

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

Do not use the 0.0.0.0/0 inbound rule

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

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

Disable unneeded inbound rules

If your current inbound rules include 0.0.0.0/0, review the ports and services that must be exposed for your applications. If you do not want some ports to directly

provide services for external applications, add denial rules for them. For example, if you have installed MySQL database services on the server, port 3306 should not be exposed to the Internet by default. You can add a denial rule, as shown below. Set the priority value to 100, which is the lowest priority.

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

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

Authorize another security group for inbound access

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

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

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

Authorize another CIDR for inbound access

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

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

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

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

Steps and instructions for changing security group rules

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



Note:

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

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

For detailed instructions, see [How to configure intercommunication among instances in the classic network](#).

1.2 Best practices of the security group (part 2)

This document introduces the following:

- [Authorize](#) and [revoke](#) security groups.
- [Join](#) and [leave](#) security groups.

Alibaba Cloud provides two types of networks, namely classic networks and VPC networks. They support different security group rules:

- For classic networks, you can set the following rules: intranet inbound, intranet outbound, Internet inbound and Internet outbound.
- For VPC networks, you can set the following rules: intranet inbound and intranet outbound.

Basic knowledge of intranet communication for security groups

Firstly, learn about the following points about intranet communication for security groups:

- By default, only the ECS instances in the same security group can access each other. In other words, the instances of the same account in different security groups are inaccessible to each other on the intranet. This applies to both classic and VPC networks. Therefore, the ECS instances in classic networks are secure over the intranet.
- If you have two ECS instances in different security groups, and you want them to be inaccessible to each other over the intranet but they are actually accessible, you should check the intranet rule settings of your security group. If the intranet rules include the following items, you are recommended to reconfigure them.
 - Allow all ports;
 - The authorized object is a CIDR segment (SourceCidrIp): `0.0.0.0/0` or `10.0.0.0/8`. For classic networks, the above rules can expose your intranet to external access.
- If you want to implement network intercommunication among the resources of different security groups, you should adopt security group authorization. For intranet access, you are recommended to adopt the source security group authorization, instead of CIDR segment authorization.

Attributes of security rules

Security rules mainly describe different access permissions with the following attributes:

- Policy: authorization policies. The parameter value can be *accept* or *drop*.
- Priority: priority levels. The priority levels are sorted by creation time in descending order. The rule priority ranges from 1 to 100. The default value is 1, which is the highest priority. A greater value indicates a lower priority.
- NicType: network type. In security group authorization (namely SourceGroupId is specified while SourceCidrIp is not), you must specify NicType as *intranet*.
- Description:
 - IpProtocol: IP protocol. Values: *tcp*, *udp*, *icmp*, *gre* or *all*. The value "all" indicates all the protocols.

- **PortRange**: the range of port numbers related to the IP protocol:
 - When the value of **IpProtocol** is *tcp* or *udp*, the port range is 1-65535. The format must be "starting port number/ending port number". For example, "1/200" indicates that the port range is 1-200. If the input value is "200/1", an error will be reported when the interface is called.
 - When the value of **IpProtocol** is *icmp*, *gre* or *all*, the port range is -1/-1, indicating no restriction on ports.
- If security group authorization is adopted, the **SourceGroupId** (namely the source security group ID) should be specified. In this case, you can choose to set **SourceGroupOwnerAccount** based on whether it is cross-account authorization. **SourceGroupOwnerAccount** indicates the account to which the source security group belongs.
- If CIDR authorization is adopted, **SourceCidrIp** should be specified. **SourceCidrIp** is the source IP address segment, which must be in the CIDR format.

Create a rule to authorize inbound requests

When you create a security group in the console or through the API, the default inbound rule is *deny all*, that is, all the inbound requests are rejected by default. This does not apply to all the situations, so you need to configure inbound rules accordingly.

If you need to enable port 80 on the Internet to provide HTTP services for external applications, do not impose any restrictions on IP network segments but set it to *0.0.0.0/0* in order to allow all inbound requests. For this purpose, you can refer to the following properties where console parameters are outside of brackets and OpenAPI parameters are within brackets (no difference is made if both parameters are the same).

- **NIC Type (NicType)**: Internet (internet). For VPCs, simply enter intranet to implement Internet access through EIP.
- **Action (Policy)**: allow (accept).
- **Rule Direction (NicType)**: inbound.
- **Protocol Type (IpProtocol)**: TCP (tcp).
- **Port Range (PortRange)**: 80/80.
- **Authorized Objects (SourceCidrIp)**: 0.0.0.0/0.
- **Priority (Priority)**: 1.

**Note:**

These recommended values apply only to the Internet. For intranet requests, you are not recommended to use CIDR network segments. Please refer to [Do not use CIDR or IP authorization for intranet security group rules of classic networks](#).

Create a rule to deny inbound requests

To deny inbound requests, you only need to configure a denial policy with a low priority. In this way, you can configure another rule with a higher priority to overwrite this rule when needed. For example, the following explains how to deny access to port 6379.

- NIC Type (NicType): Intranet (intranet).
- Action (Policy): forbid (drop).
- Rule Direction (NicType): inbound.
- Protocol Type (IpProtocol): TCP (tcp).
- Port Range (PortRange): 6379/6379.
- Authorized Objects (SourceCidrIp): 0.0.0.0/0.
- Priority (Priority): 100.

Do not use CIDR or IP authorization for intranet security group rules of classic networks

For ECS instances in classic networks, no intranet inbound rules are enabled by default. Always exercise caution for intranet authorization.

**Note:**

For the sake of security, it is not recommended to enable any authorization that is based on CIDR network segments.

For elastic computing, intranet IP addresses change frequently and the network segment to which the IP addresses map varies dynamically. For this reason, you are only recommended to authorize intranet access through security groups in classic networks.

For example, if you build a Redis cluster in the sg-redis security group and only permit certain computers (such as those in sg-web) to access the servers of this Redis cluster, you do not need to configure CIDR. Instead, you only need to add an inbound rule to specify relevant security group IDs.

- NIC Type (NicType): Intranet (intranet).

- Action (Policy): allow (accept).
- Rule Direction (NicType): inbound.
- Protocol Type (IpProtocol): TCP (tcp).
- Port Range (PortRange): 6379/6379.
- Authorized Objects (SourceGroupId): sg-web.
- Priority (Priority): 1.

For instances in a VPC, if you have planned an IP address range through multiple VSwitches, you can use the CIDR settings as the security group inbound rules. However, if your VPC network segment is ambiguous, you are recommended to prioritize security groups for inbound rules.

Add ECS instances requiring intercommunication into the same security group

A single ECS instance can join up to five security groups, and the ECS instances in the same security group can intercommunicate over the intranet. If you have created multiple security groups during planning and directly setting multiple security rules is too complicated, you can create a security group and add the instances that require intranet communication to it.

Different security groups may have different network types. More importantly, an ECS instance in a classic network can only join a security group created for classic networks. An ECS instance in a VPC can only join a security group created for the same VPC.

Additionally, you are not recommended to add all the ECS instances into the same security group as this will make the configuration of security group rules quite messy. For a large or medium-sized application, each server group has a different role and it is important to plan inbound and outbound requests in a rational manner.

In the console, you can add an instance to a security group by following the description in [Join a security group](#).

If you are quite familiar with Alibaba Cloud OpenAPI, you can perform batch operations through OpenAPI. For more information, see [Manage ECS instances elastically by using OpenAPI](#). The corresponding Python snippets are as follows.

```
def join_sg(sg_id, instance_id):
    request = JoinSecurityGroupRequest()
    request.set_InstanceId(instance_id)
    request.set_SecurityGroupId(sg_id)
    response = _send_request(request)
    return response
```

```
# send open api request
def _send_request(request):
    request.set_accept_format('json')
    try:
        response_str = clt.do_action(request)
        logging.info(response_str)
        response_detail = json.loads(response_str)
        return response_detail
    except Exception as e:
        logging.error(e)
```

Remove an ECS instance from a security group

If an ECS instance is added to an inappropriate security group, your services may be exposed or blocked. In this case, you can remove the ECS instance from the security group. Before the removal, however, you must ensure that your ECS instance has been added to another security group.



Note:

You are recommended to perform sufficient tests before the removal as this may cause intercommunication failure between the instance and other instances in the current security group.

The corresponding Python snippets are as follows:

```
def leave_sg(sg_id, instance_id):
    request = LeaveSecurityGroupRequest()
    request.set_InstanceId(instance_id)
    request.set_SecurityGroupId(sg_id)
    response = _send_request(request)
    return response
# send open api request
def _send_request(request):
    request.set_accept_format('json')
    try:
        response_str = clt.do_action(request)
        logging.info(response_str)
        response_detail = json.loads(response_str)
        return response_detail
    except Exception as e:
        logging.error(e)
```

Define reasonable names and tags for security groups

Reasonable names and descriptions for security groups help you quickly identify the meanings of complicated rule combinations. You can change security group names and descriptions as needed.

Also, you can set tags for security groups. You can manage your own security groups by grouping them with tags. To [set tags](#), you can directly configure them in the console or by using APIs.

Delete undesired security groups

Security rules of security groups are like whitelist and blacklist items. Therefore, you are recommended to delete unnecessary security groups to prevent unexpected problems caused by adding an ECS instance to those groups by mistake.

1.3 Best practices of the security group (part 3)

In practice, all instances may be placed in the same security group, thus reducing the configuration workload in the initial period. In the long run, however, interactions of the business systems will become complicated and uncontrollable. When you modify a security group, you will be unable to clearly identify the impact scope of adding or removing a rule.

Rational planning and differentiation of security groups makes it easy to adjust your systems, sort out the services provided by the applications, and arrange applications at different layers. We recommend that you plan different security groups and set different security group rules for different businesses.

Distinguish between different security groups

- Use different security groups for ECS instances on the Internet and those on the intranet

ECS instances that provide Internet services, either through exposure of some ports for external access (such as 80 and 443) or through provision of port forwarding rules (for example, instances configured with forwarding rules for Internet IP address, EIP address or NAT ports), will expose their applications to the Internet.

For the two scenarios above, the relevant security groups should adopt the strictest rules. We recommend that Internet access should be rejected first. Specifically, all ports and protocols should be disabled by default except the ports needed to provide external services, such as 80 and 443. As the security group only contains the ECS instances that provide Internet access, it is easier to adjust the security group rules.

For a group of ECS instances that provide Internet access, their responsibilities should be clear and simple, avoiding offering other external services on the same instances. For MySQL, Redis, and more, for example, it is recommended to install

such services on ECS instances that disable Internet access, and then enable access to them through security group authorization.

Assume you have an ECS instance that provides Internet access, which is in the security group SG_CURRENT as the instances of other applications. You can make changes by performing the steps below.

1. Sort out the ports and protocols exposed by the current Internet services, such as 80 and 443.
2. Create a new security group such as SG_WEB and add corresponding ports and rules.



Note:

Action: Allow; Protocol Type: All; Port Range: 80/80; Authorization Objects: 0.0.0.0/0; Action: Allow; Protocol Type: All; Port Range: 443/443; Authorization Objects: 0.0.0.0/0.

3. Select the security group SG_CURRENT and add a rule for security group authorization, that is, allowing the resources in SG_WEB to access the resources in SG_CURRENT.



Note:

Action: Allow; Protocol Type: All; Port Range: -1/-1; Authorization Objects: SG_WEB; Priority: Choose from [1-100] according to actual conditions.

4. Add ECS_WEB_1 to the new security group. It is an instance that needs to switch its security group.
 - a. In the ECS console, select Security groups.
 - b. Select SG_WEB > Manage Instances > Add Instance. Add the instance ECS_WEB_1 to the new security group SG_WEB. Make sure ECS_WEB_1 works normally.
5. Remove the instance ECS_WEB_1 from the original security group.
 - a. In the ECS console, select Security Groups.
 - b. Select SG_WEB > Manage Instances > Add Instance. Select ECS_WEB_1 and remove it from SG_CURRENT. Verify that the traffic and network are in normal condition.

- c. If errors occur, add ECS_WEB_1 back to the security group SG_CURRENT.

Check whether the ports of SG_WEB are exposed as expected, and then make adjustments accordingly.

- 6. Make other changes to the security group.

- Use different security groups for different applications

In production environments, different operating systems generally do not belong to the same application group to provide load balancing services. Providing different services means that exposed ports are different from rejected ports. Therefore, it is recommended that instances with different operating systems belong to different security groups.

For example, TCP port 22 may be exposed for implementing SSH in Linux, while TCP port 3389 may be exposed for implementing remote desktop connection in Windows.

In addition, for instances that have the same type of images but provide different services, it is recommended to put them into different security groups if they do not need to access each other over the intranet. This facilitates decoupling and future changes to security group rules as the rules can be as simple as possible.

When planning and adding new applications, you should reasonably organize the security groups apart from dividing different VSwitches to configure subnets.

You can use network segments and security groups to distinguish yourself as the service provider or consumer.

For specific change procedures, see the operations above.

- Use different security groups for production environments and testing environments

To better isolate systems, you may build multiple testing environments and one online environment during actual development. For better network isolation, you need to configure different security policies for different environments, preventing changes to the testing environment from being synchronized to the online environment, which may affect the stability of online services.

By creating different security groups, you can restrict the access domains of applications and avoid interoperability between the production environment and testing environment. Also, you can create different security groups for different

test environments, thus avoiding interference between test environments and improving development efficiency.

Only assign Internet addresses to subnets or instances that require Internet access

Whether it is a classic network or a VPC, rational allocation of Internet addresses facilitates Internet management of the system and reduces the risk of attack. For VPCs, we recommend that you place the IP segments of instances requiring Internet access onto several dedicated VSwitches (subnet CIDR) when creating a VSwitch. This facilitates auditing and differentiation and helps avoid accidental Internet access.

Most distributed applications have different layers and groups. For ECS instances that offer no Internet access, try your best not to provide Internet addresses for them. If there are multiple instances that provide Internet access, we recommend you to configure the [Server Load Balancer](#) to distribute traffic of Internet services, thus improving system availability and avoiding a single point of failure.

For ECS instances that require no Internet access, try your best not to assign Internet addresses to them. In VPCs, when your ECS instances need to access the Internet, we recommend you to use the [NAT gateway](#) to provide Internet proxy services for ECS instances without Internet addresses in the VPC. By simply configuring the corresponding SNAT rules, you can enable a specific CIDR segment or subnet to access the Internet. For specific configurations, see [SNAT](#). In this way, exposure of services to the Internet can be avoided after Elastic IP (EIP) addresses are allocated when only outbound access is required.

Minimum principle

A security group should work as a whitelist. Therefore, try your best to open and expose as few ports as possible, and allocate as few Internet addresses as possible. Although allocating Internet addresses or binding EIPs makes it easy to access online instances for troubleshooting, it exposes the entire instance to the Internet after all. A safer policy is to manage IP addresses by using the Jump Server.

Use the Jump Server

As the Jump Server has much higher permissions, relevant operations should be well recorded and audited through tools. In addition, it is recommended to choose a dedicated VSwitch for the Jump Server in VPCs, providing the corresponding EIP or NAT port forwarding tables to it.

First, create a dedicated security group SG_BRIDGE by enabling the corresponding port such as TCP 22 in Linux or RDP 3389 in Windows. To restrict the inbound access, you can limit the authorization objects to the Internet egress ports of your company, lowering the probability of being scanned and accessed.

After that, you can add the Jumper Server instance to that security group. In order for that Jumper Server to access other appropriate instances, you can configure appropriate group authorization. For example, add a rule for SG_CURRENT, allowing SG_BRIDGE to access certain ports and protocols.

When you use the Jumper Server for SSH communication, it is recommended to use the *SSH key pair* for logon, instead of the password.

In summary, reasonable planning of security groups makes it easy for you to expand the applications and makes your system more secure.

1.4 Best practices of ECS data security

This document introduces how to implement data security for ECS instances from the O&M perspective.

Intended audience

This document applies to individuals and enterprises that are new to Alibaba Cloud.

Contents

- Back up data regularly
- Design security domains properly
- Set security group rules
- Set logon passwords
- Server port security
- Application vulnerability protection
- Security information collection

Back up data regularly

As the foundation of disaster tolerance, data backup is intended to reduce the risk of data loss due to system failures, operation errors, and security problems. ECS instances come with the snapshot backup function. Correctly using the snapshot function can satisfy the data backup requirements for most users. It is recommended to customize your own backup policy according to actual business needs. You can

select [Create Snapshot](#) or [Create Automatic Snapshot Policy](#), and [apply the policy to specific disks](#).

It is recommended to take automatic snapshots on a daily basis, and store them for at least seven days. Good backup habits contribute to rapid data recovery and minimizing losses in the case of failure.

Design security domains properly

Developed upon the Software Defined Network (SDN) technology, VPCs allow you to build private networks that separate servers of different security levels in your enterprise, preventing servers from impacting each other over an interconnected network.

It is recommended that you [create VPC](#), and set the IP address range, network segments, route tables, and gateways. You can store important data in an intranet that is totally isolated from the Internet. You can use Elastic IP (EIP) addresses or the Jumper Server to manage data in daily O&M.

Set security group rules

As an important means for security isolation, security groups are used to set network access control for one or more ECS instances. With security groups, you can set firewall policies at the instance level, filtering active and passive access of an instance at the network layer. Specifically, you can restrict inbound and outbound access on a port, and authorize access to IP addresses, reducing attacks and enhancing instance security.

For example, the remote port is 22 by default in Linux, which should not be open to the Internet directly. You can set up a security group to control the Internet access to an ECS instance, such as authorizing fixed IP addresses to access the instance. To learn more about security groups, see [Application cases](#). If you have higher requirements, you can also use third-party VPN products to encrypt the logon data. For more software, visit [Alibaba Cloud Market](#).

Set logon passwords

Weak passwords have been a major cause of data leakage, as they are one of the most common vulnerabilities and can be exploited very easily. It is recommended that the server password should contain at least eight characters, and should be complicated enough by including uppercase and lowercase letters, numbers and special characters. In addition, you should change the password regularly.

Server port security

As long as servers provide Internet services, the corresponding ports will be exposed to the Internet. From the perspective of security management, more open ports mean more system risks. It is recommended to open as few ports as necessary to the Internet. Common ports should be changed to custom ports (port 30000 or greater), and access control should be implemented on the service ports.

For example, it is recommended to restrict database services to the intranet and prevent access from the Internet. If it is necessary to access the database directly from the Internet, you need to change the connection port from 3306 to a greater port , and authorize the relevant IP addresses according to the business needs.

Application vulnerability protection

Application vulnerabilities are security defects that can be exploited by hackers to illegally access data from Web applications, cache, database and storage. Common application vulnerabilities include SQL injection, XSS attacks, Web shells, backdoor, command injection, illegal HTTP requests, common Web server vulnerability attacks, unauthorized access to core files, path traversal, and more. These vulnerabilities are different from system vulnerabilities, and are difficult to fix. If application security cannot be guaranteed during the initial design, servers may be attacked due to such vulnerabilities. Therefore, it is recommended to install a [Web Application Firewall \(WAF\)](#) to prevent various attacks, thus ensuring website security and availability.

Security information collection

In today's Internet security field, both security engineers and hackers are racing against the clock. As a security service based on big data, [Alibaba Cloud Security Situational Awareness](#) can fully, rapidly and accurately capture and analyze factors that may lead to security situation changes in large cloud computing environments. After that, Situational Awareness associates the current threats with the past ones to perform big data analysis, so as to predict potential security events in the future and provide a systematic solution.

Therefore, in addition to daily O&M, technicians should obtain as much information as possible to improve warning capability. In this way, quick recovery can be made possible in the case of security problems and ECS data security can be truly guaranteed.

1.5 How to configure instances to access each other in classic networks

A security group is an instance-level firewall. To ensure the instance security, the minimum authorization principle should be observed for the setting of security group rules. This document introduces four safe methods of enabling intranet intercommunication for instances.

Method 1. Authorize access to a single IP address

- **Application scenario:** intercommunication of a small number of instances over the intranet.
- **Advantage:** Authorizing access to IP addresses makes the security group rules clear and easy to understand.
- **Disadvantage:** When a great number of instances need to access each other over the intranet, it is limited by the quota of 100 security group rules. In addition, the maintenance workload will be high.
- **Configuration:**
 1. Select the instance that requires intercommunication, and click Security Groups.
 2. Select the expected security group and click Add Rules.
 3. Click Ingress and then click Add Security Group Rule.
 4. Add security group rules as instructed below:
 - **Action:** Allow.
 - **Protocol Type:** Select the protocol type as needed.
 - **Port Range:** Set the port range as needed. The format is “start port number/end port number” .
 - **Authorization Type:** CIDR.
 - **Authorization Objects:** Enter the expected intranet IP address for intranet intercommunication. The format must be *a.b.c.d/32*. Where, the subnet mask must be /32.

Add Security Group Rule

NIC:

Internal Network

Rule Direction:

Ingress

Action:

Allow

Protocol Type:

Customized TCP

* Port Range:

Example: 22/22 or 3389/3388

Priority:

1

Authorization Type:

CIDR

* Authorization Objects:

Example: 10.0.0.0/32

Tutorial

Description:

It can be 2 to 256 characters in length and cannot start with http:// or https://.

OK

Cancel

Method 2. Join the same security group

- **Application scenario:** If your application architecture is relatively simple, you can add all the instances to the same security group. Such instances need no special rules as they can access each other over the intranet by default.
- **Advantage:** Security group rules are clear and easy to understand.
- **Disadvantage:** It is only applicable to simple application network architecture. When the network architecture is adjusted, the authorization method should be modified accordingly.

Method 3. Bind instances with a security group that is created solely for intercommunication

- **Application scenario:** You can bind expected instances to a dedicated security group for intercommunication. This method applies to the network architecture with multiple layers of applications.
- **Advantage:** This method is easy to implement, and allows you to quickly establish interconnection between instances. It is applicable to complicated network architecture.
- **Disadvantage:** The instances need to be bound to multiple security groups and the security group rules are hard to comprehend.
- **Configuration:**
 1. Create a new security group with the name of “security group for intercommunication” . No rules are required for the new security group.
 2. Add the expected instances to the newly created “security group for intercommunication” . The instances will be interconnected over the intranet as this is a default feature for instances in the same security group.

Method 4. Security group authorization

- **Application scenario:** If your network architecture is complicated, and the applications deployed on different instances have different service roles, you can select security group authorization.
- **Advantage:** The security group rules are clear and easy to understand. Besides, intercommunication can be implemented across accounts.
- **Disadvantage:** You need to configure a lot of security group rules.
- **Configuration:**
 1. Select the expected instance, and enter the Security Groups page.
 2. Select the expected security group and click Add Rules.
 3. Click Ingress, and then click Add Security Group Rule.
 4. Add security group rules as described below:
 - **Action:** Allow.
 - **Protocol Type:** Select the protocol type as needed.
 - **Port Range:** Set it as needed.
 - **Authorization Type:** Security Group.
 - **Authorization Objects:**

- **Allow Current Account:** Based on your networking requirements, select the security group IDs of the peer instances for intranet intercommunication in Authorized Objects.
- **Allow Other Accounts:** Enter the security group IDs of the peer instances in Authorized Objects. Enter the peer account ID in Account ID. You can query it in Account Management > Security Settings.

Add Security Group Rule

NIC:

Internal Network

Rule Direction:

Ingress

Action:

Allow

Protocol Type:

Customized TCP

* Port Range:

Example: 22/22 or 3389/338

Priority:

1

Authorization Type:

Security Group

☒ Allow Current Account

☐ Allow Other Accounts

Authorization Objects:

Select Security Group

Description:

It can be 2 to 256 characters in length and cannot start with http:// or https://.

OK

Cancel

Add Security Group Rule

NIC:

Internal Network

Rule Direction:

Ingress

Action:

Allow

Protocol Type:

Customized TCP

* Port Range:

Example: 22/22 or 3389/338

Priority:

1

Authorization Type:

Security Group

☐ Allow Current Account

☒ Allow Other Accounts

Authorization Objects:

sg-xxxxxxxxxxxxxxxxxxxxxxxx

Account ID:

xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx

Enter an account ID. To query your account ID, go to [Account Center](#)

Description:

It can be 2 to 256 characters in length and cannot start with http:// or https://.

OK

Cancel

Suggestions

If too much access is granted by the security group in the early stage, it is recommended to reduce the authorization scope with the following procedure.



In the figure, Delete 0.0.0.0 means to delete the original security group rule that allows the inbound access from the 0.0.0.0/0 address segment.

If the security group rules are changed improperly, the communications between your instances may be affected. Please back up the security group rules you want to change before changing the settings for timely recovery upon intercommunication problems.

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

- Web layer security group: Open port 80.
- Application layer security group: Open port 8080.
- DB layer security group: Open port 3306.

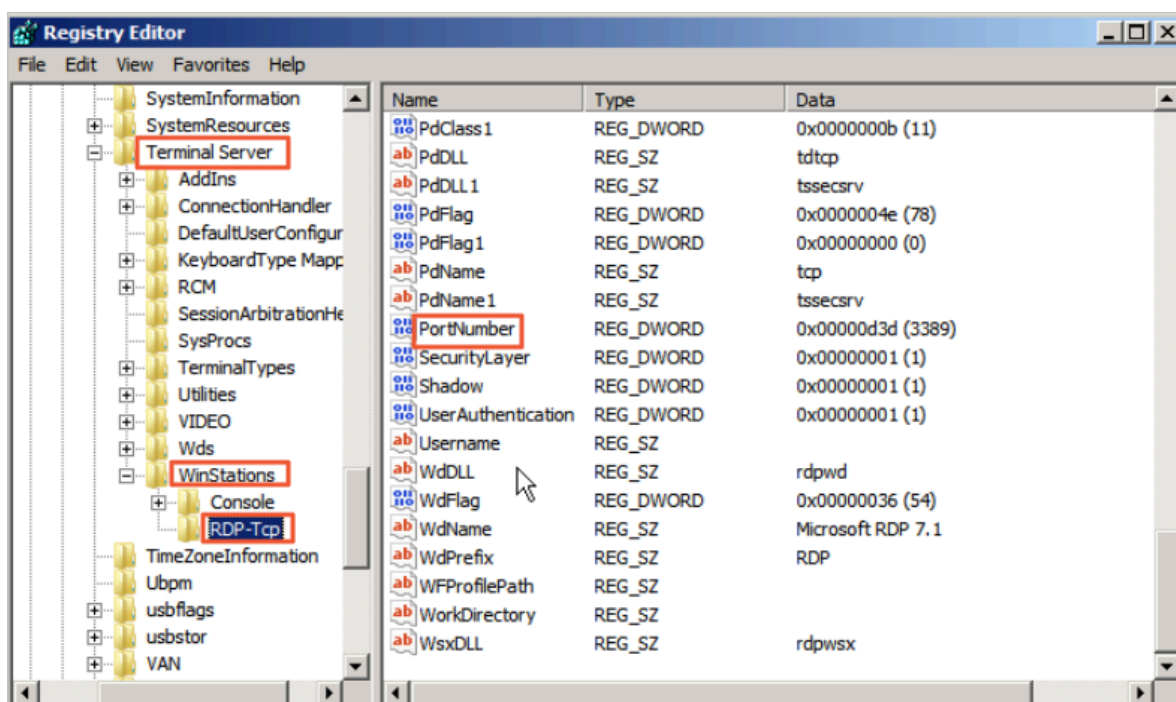
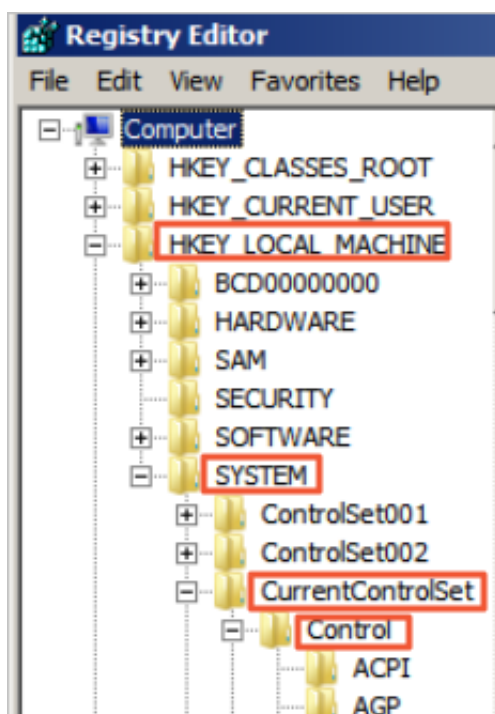
1.6 Modify the default remote access port

This topic describes how to modify the remote port of a Windows or Linux instance.

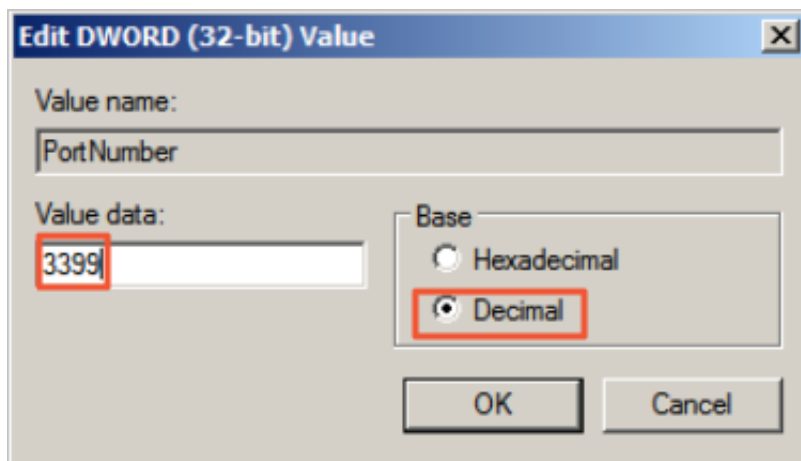
Modify the default remote port of a Windows instance

This section describes how to modify the remote port of a Windows instance running Windows Server 2008.

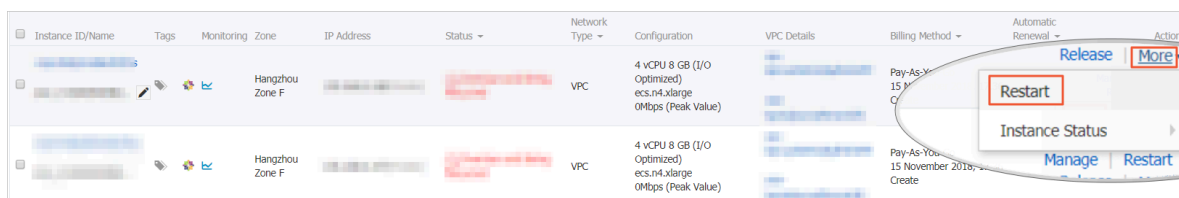
1. [Connect to the Windows instance](#).
2. Run `regedit.exe` to open Registry Editor.
3. On the left-side navigation pane of the Registry Editor, find `HKEY_LOCAL_MACHINE\System\CurrentControlSet\Control\Terminal Server\WinStations\RDP-Tcp\PortNumber`.



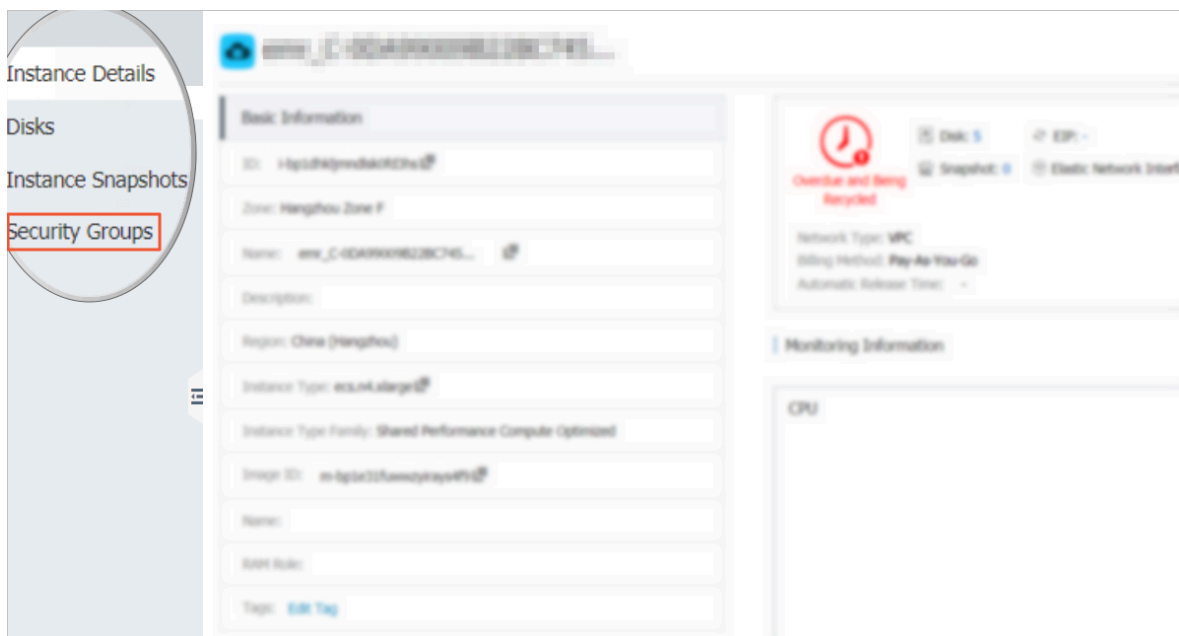
4. In the dialog box, select Decimal as Base, and then type a number in the Value data field as the new remote port number, which is 3399 in this example. Click OK.



5. (Optional) If you have enabled firewall, open the new port on the firewall.
6. Log on to the [ECS console](#), find the instance, and then select More > Restart.



7. After the instance is restarted, click the Manage of the instance to enter the Instance Details page. Click Security Groups.



8. On the Security Groups page, click Add Rules.
9. On the Security Group Rules page, click Add Security Group Rule. Add a new security group rule to allow access to the new remote port. For more information about adding security group rules, see [Add security group rules](#).

Add Security Group Rule ? Add security group rules

NIC:

Internal Network ▼

Rule Direction:

Ingress ▼

Action:

Allow ▼

Protocol Type:

Customized TCP ▼

Port Range:

3399/3399

i

Priority:

1

i

Authorization Type:

IPv4 CIDR Block ▼

* Authorization Objects:

Example: 10.x.y.z/32. You can specify up to 10 authorization objects separated with commas (,.)

i Tutorial

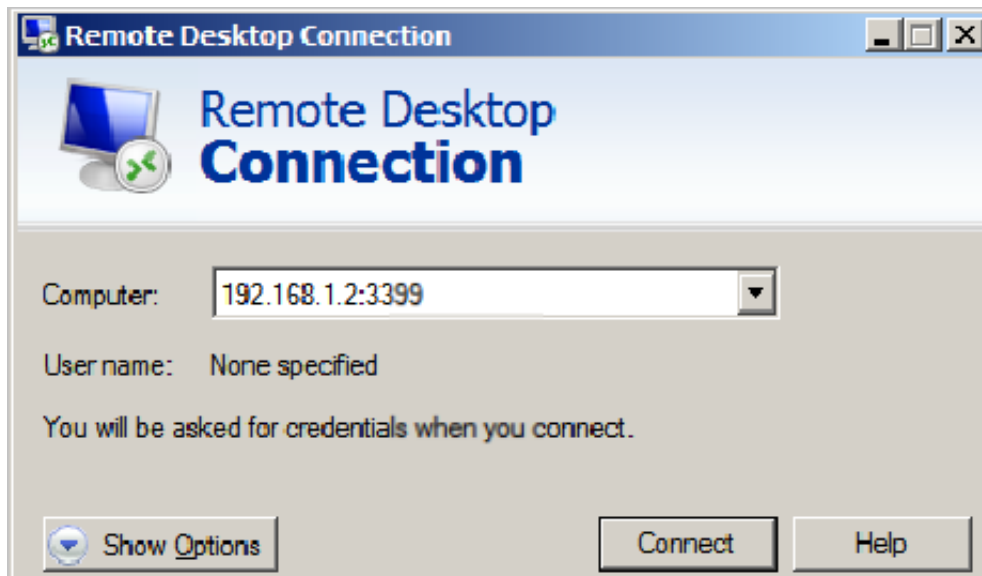
Description:

It can be 2 to 256 characters in length and cannot start with http:// or https://.

OK

Cancel

10.Connect to the instance by accessing the IP address ending with the new port number. For example, 192.168.1.2:3399 in this example.

**Note:**

Only the default port 3389 can be used for access by Mac remote desktop users.

Modify the default remote port of a Linux instance

This section describes how to modify the remote port of a Linux instance running CentOS 6.8.

**Note:**

Do not modify the 22 port directly, first add the new default remote port. Set two ports first and delete one after the test succeeds. It ensures that you can use port 22 to debug any problems if you cannot connect the instance through the new port.

1. [Connect to the Linux instance.](#)
2. Run `vim /etc/ssh/sshd_config`.
3. Press the "I" key on the keyboard to enter the Edit mode. Add new remote service port (for example, Port 1022). Enter *Port 1022* under *Port 22*.
4. Press "ESC" and enter `: wq` to exit the editing.
5. Restart the instance by executing the following command. You can then log on to the Linux instance through 22 port and 1022 port.

```
/etc/init.d/ssh restart
```

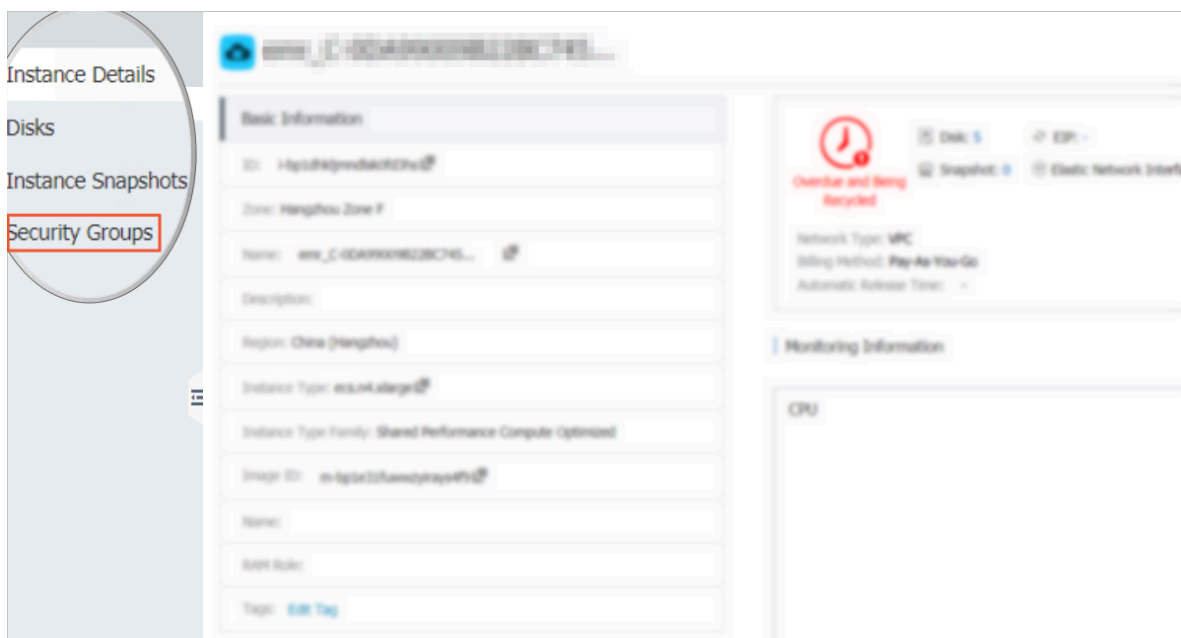
6. (Optional) Configure the firewall. When you use Linux versions earlier than CentOS 7 and has enabled firewall iptables, note that iptables do not intercept access by default. If you configured iptables rules, run `iptables -A INPUT -p tcp --dport`

1022 -j ACCEPT to configure the firewall. Then perform `service iptables restart` to restart the firewall.

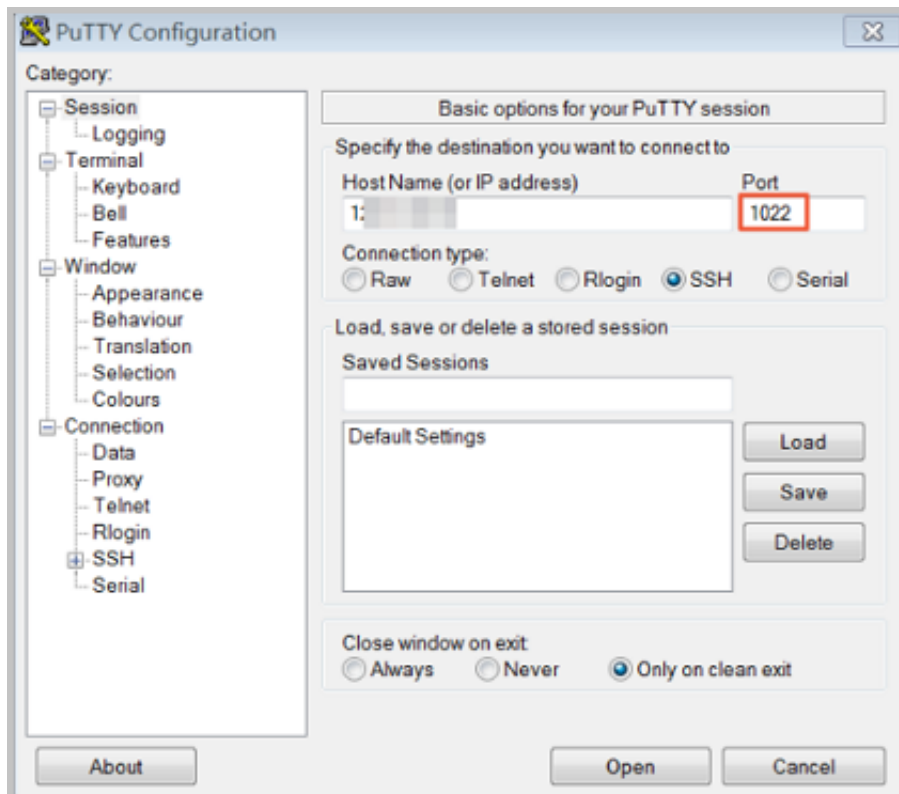
**Note:**

Firewalld is installed by default on CentOS 7 and later versions. If you have enabled `firewalld.service`, open TCP port 1022 by running the command `firewall-cmd --add-port=1022/tcp --permanent`. If success is returned, TCP port 1022 is opened.

7. Log on to the [ECS console](#), find the instance, and then select Manage.
8. Enter the Instance Details page. Click Security Groups.



9. On the Security Groups page, click Add Rules.
10. On the Security Group Rules page, click Add Security Group Rule. Add a new security group rule to allow access to the new remote port. For more information about adding security group rules, see [Add security group rules](#).
11. Use the SSH tool to connect to the new port to test if the default remote port is modified successfully. Enter the new port number in Port when logging on to the instance, which is 1022 in this example.



12. Once you successfully connect the instance through port 1022, run `vim /etc/ssh/sshd_config` again to remove port 22.
13. Run `/etc/init.d/sshd` to restart the instance and the default remote port is successfully modified. Connect to the instance by accessing the IP address ending with the new port number.

1.7 Use logs in Windows instances

Logs are records of hardware and software in the system, and system error information. They can also be used to monitor system events. When a server intrusion or system (application) error occurs, administrators can quickly locate the problems by using logs and solve the problems quickly, which improves work efficiency and server security substantially. Windows logs can be mainly divided into four categories: system logs, application logs, security logs, and applications and services logs. In this example, we use Windows Server 2008 R2 to introduce the use and analysis of the four categories of logs.

Open the Event Viewer

Follow these steps to open Event Viewer: Open the Run window, type `eventvwr`, and then click OK to open the Event Viewer.

Then, you can view the following four categories of logs in Event Viewer.



Note:

You can find the solutions to any error event ID that you can find in these logs in Microsoft knowledge base.

- **System Logs**

System logs include events recorded by Windows system components. For example, system logs record failures that occur when loading drivers or other system components during startup.

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

- **Application logs**

Application logs include events recorded by applications or programs. For example, a database application can record file errors in application logs.

The types of events recorded are determined by developers.

- **Security logs**

Security logs include events such as valid and invalid logon attempts, and resource usage related events such as creation, opening, or deletion of files or other objects.

Administrators can specify the types of events recorded in security logs. For example, if logon has been set to be audited, logon attempts are recorded in security logs.

- **Application and service logs**

Application and service logs are a new type of event logs. These logs store events from a single application or component, rather than events that may affect the global system.

Modify log path and back up logs

Logs are stored on the system disk by default. The maximum log size is 20 MB by default, and the earliest events are overwritten when 20 MB is exceeded. You can modify the maximum log size according to your needs.

Follow these steps to modify the log path and back up logs:

1. In the left-side navigation pane of Event Viewer, click Windows Logs.

2. Right click a log name, such as Application and click Properties.
3. In the Log Properties dialog box, you can modify the following settings:
 - Log path
 - Maximum log size
 - Operations executed when maximum event log size is reached

1.8 Overview and best practices of Windows Firewall with Advanced Security

This article introduces Windows Firewall with Advanced Security (WFAS), its application scenarios, and common operations.

Overview

As an important part of the hierarchical security model, WFAS was launched after Windows NT6.0 by Microsoft. WFAS blocks unauthorized traffic that flows in or out of local computers by providing bi-directional filtering based on the current connection status. WFAS also uses Network Location Awareness (NLA) to apply the corresponding firewall profile to the computer based on its current connection status. The security rules of Windows Firewall and Internet Protocol Security (IPsec) are configured in the Microsoft Management Console (MMC) snap-in, and WFAS is also an important part of the network isolation policy.

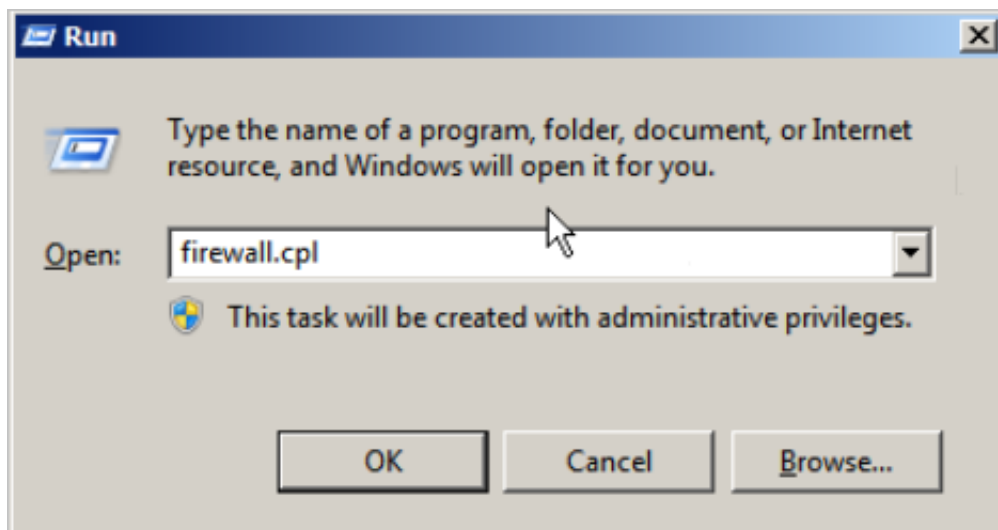
Application Scenario

More and more O&M personnel are reporting that servers are attacked and passwords are cracked, which in most cases, are due to the “backdoor” left open to “intruders”. Intruders scan open ports on your computers and penetrate them through vulnerable ports, for example, the remote port 3389 in Windows and the remote port 22 in Linux. Now that we know where the problem is, we can take the effective countermeasure. Specifically, we can close these “backdoors” by modifying the default remote ports and restricting remote access. So how do we restrict remote access? Now let's demonstrate how to restrict the remote desktop connection by taking an ECS instance (Windows Server 2008 R2) for example.

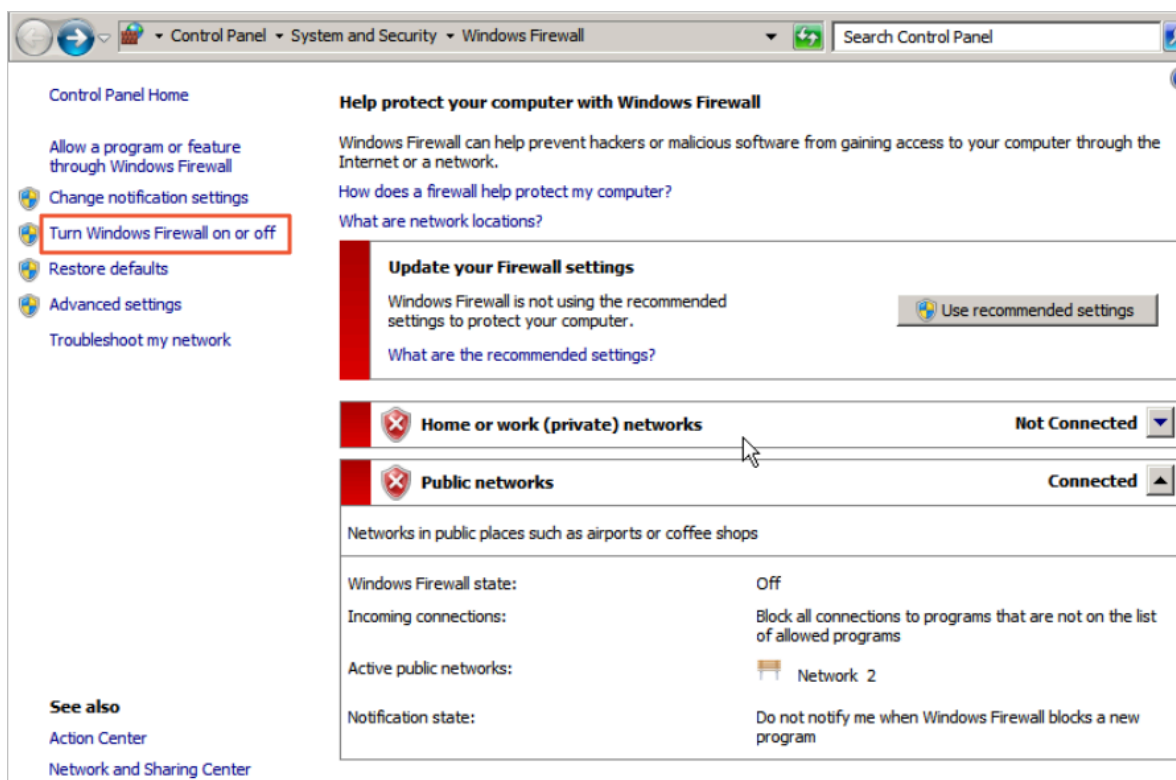
Procedure

1. View the Windows Firewall status

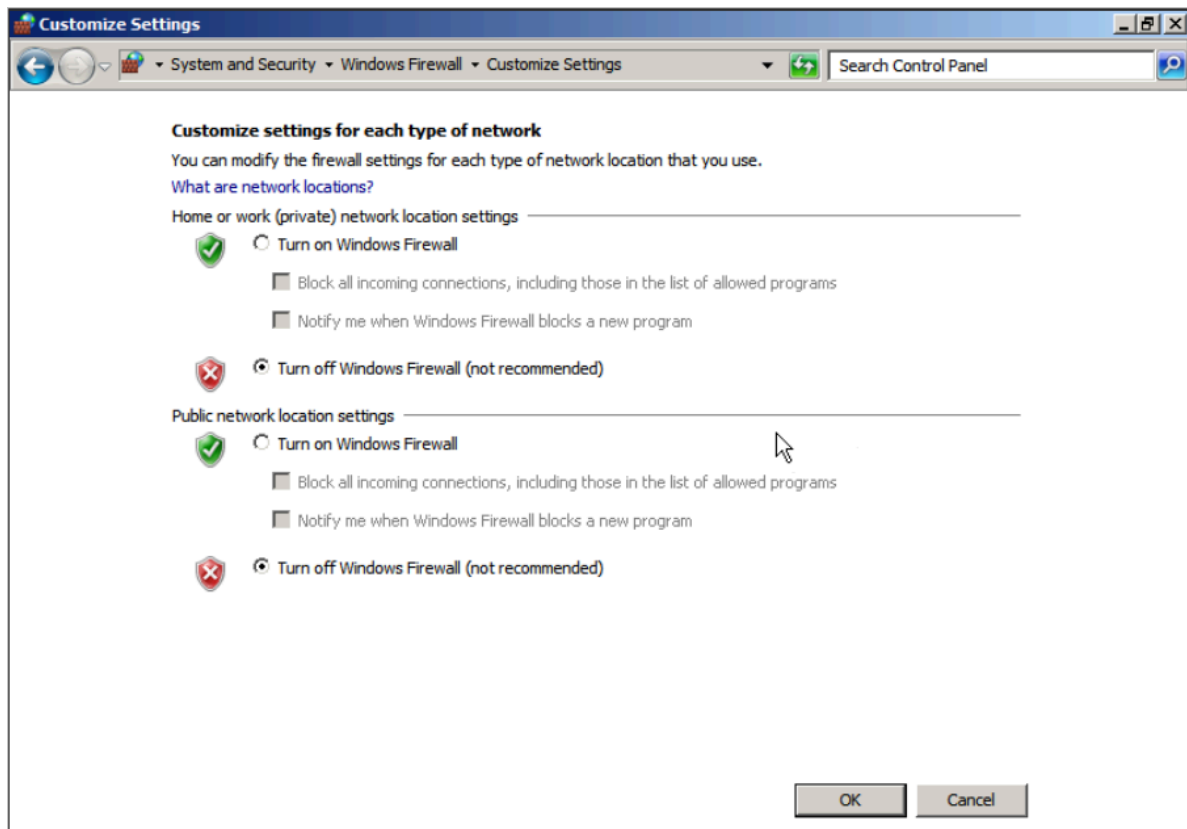
Windows Firewall of the ECS instance is disabled by default. You can press Win+R to open the Run window, enter *firewall.cpl*, and then press Enter to open the Windows Firewall console, as shown below.



Enable or disable Windows Firewall.

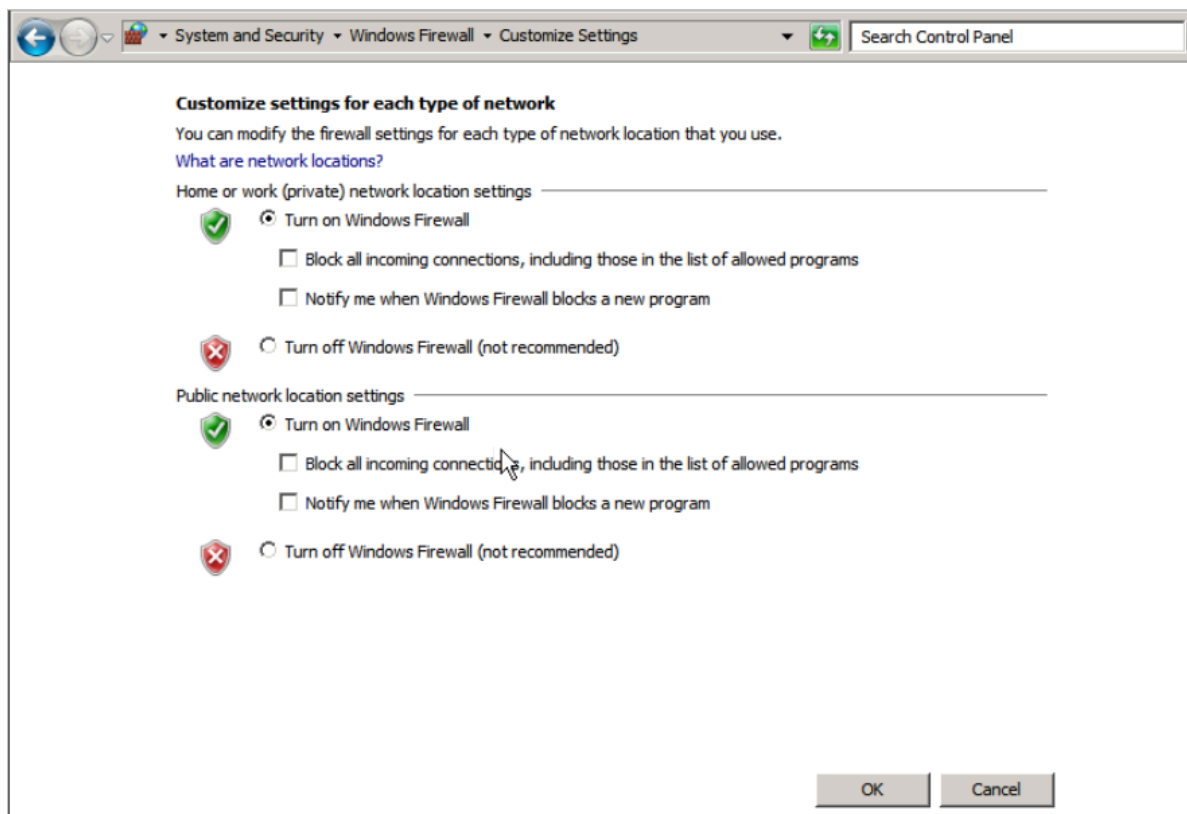


As shown below, Windows Firewall is disabled by default.

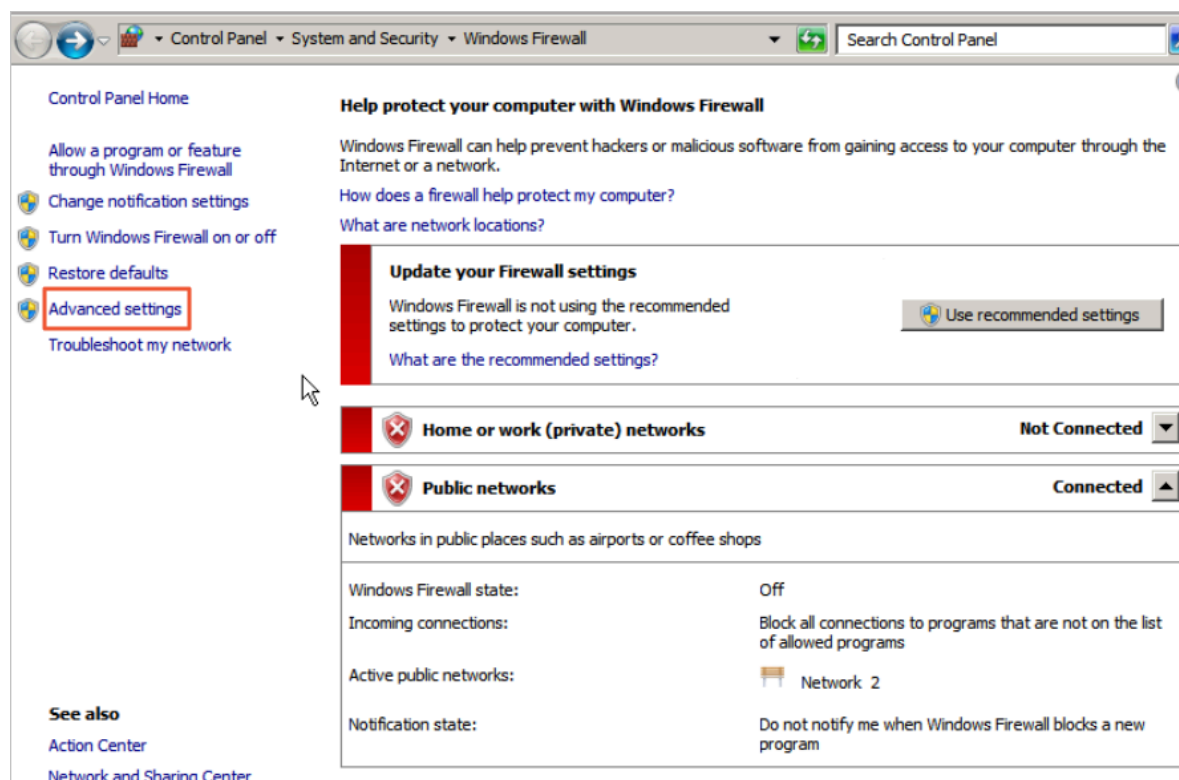


2. Enable Windows Firewall

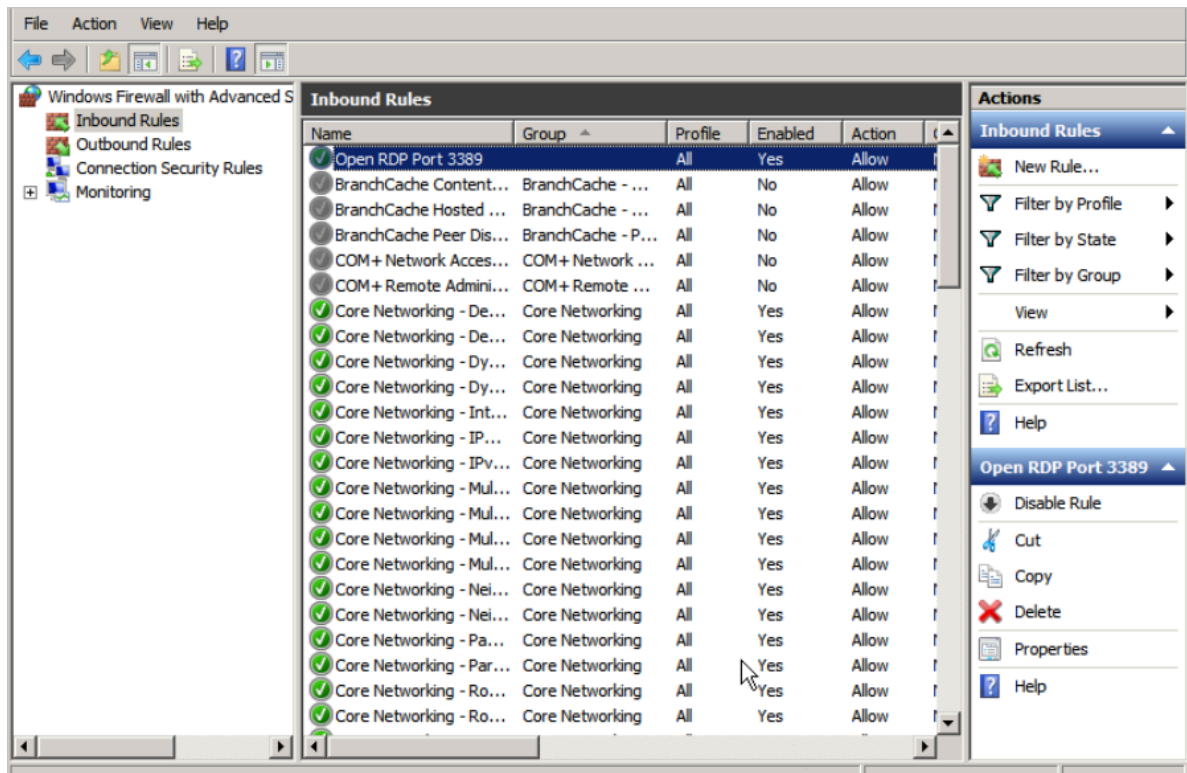
Enable Windows Firewall through the previous steps, as shown below.



Before enabling Windows Firewall, make sure the remote port is open in the inbound rules, or you cannot establish the remote connection even yourself. WFAS, however, opens RDP port 3389 in its inbound rules by default. Select Advanced settings.

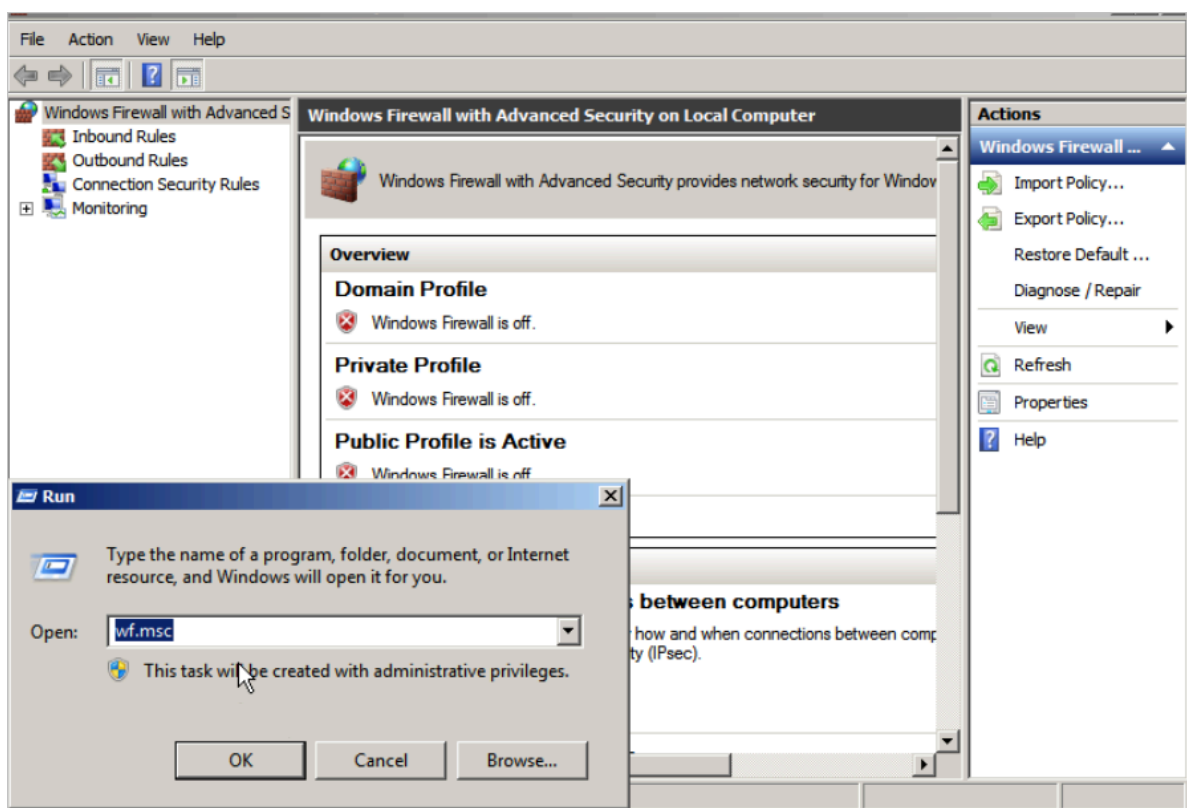


Select Inbound Rules. We can see that the Open RDP Port 3389 rule is enabled by default.

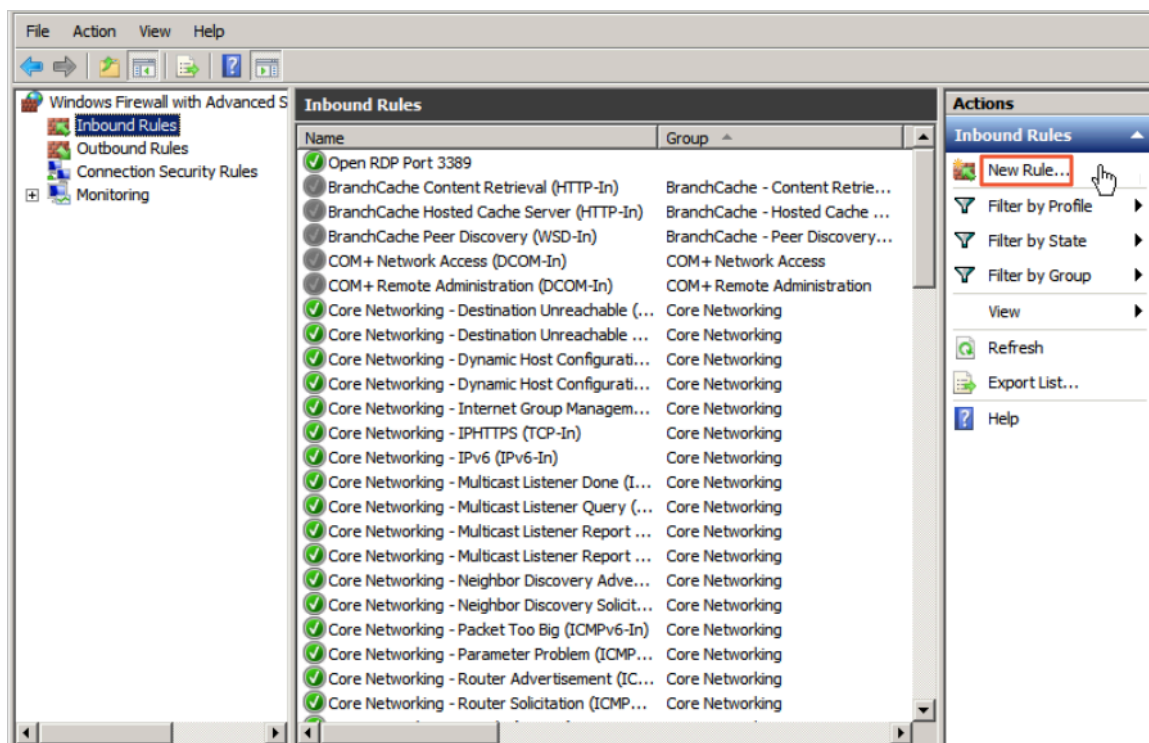


3. Configure WFAS

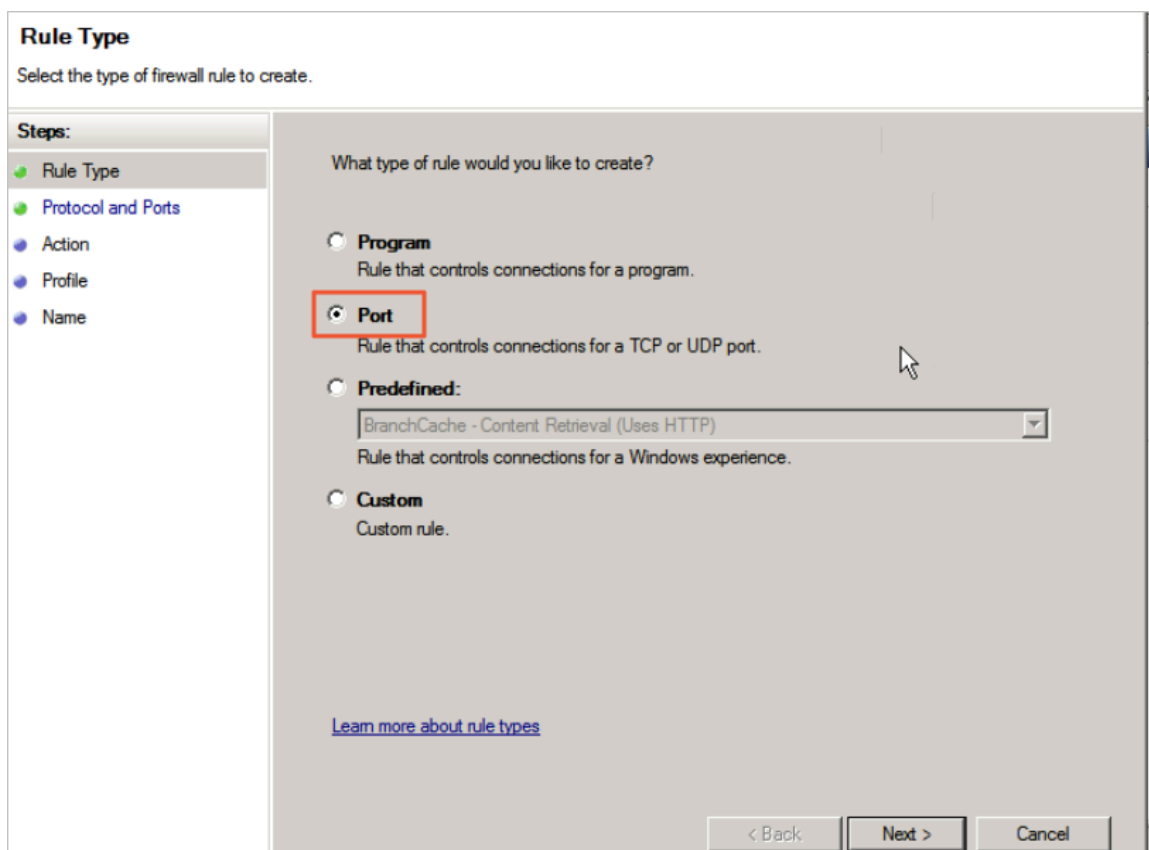
Press Win+R to open the Run window, enter *wf.msc*, and then press Enter to open the WFAS window, as shown below.



a. Create an inbound rule manually



In the New Inbound Rule Wizard window, select Port and click Next.



Select TCP and set Specific Local Ports to 3389.

Protocol and Ports

Specify the protocols and ports to which this rule applies.

Steps:

- Rule Type
- Protocol and Ports
- Action
- Profile
- Name

Does this rule apply to TCP or UDP?

☒ **TCP**

☐ **UDP**

Does this rule apply to all local ports or specific local ports?

☐ **All local ports**

☒ **Specific local ports:**

Example: 80, 443, 5000-5010

[Learn more about protocol and ports](#)

< Back Next > Cancel

Click Next and select Allow the connection.

Action

Specify the action to be taken when a connection matches the conditions specified in the rule.

Steps:

- Rule Type
- Protocol and Ports
- Action
- Profile
- Name

What action should be taken when a connection matches the specified conditions?

☒ **Allow the connection**

This includes connections that are protected with IPsec as well as those are not.

☐ **Allow the connection if it is secure**

This includes only connections that have been authenticated by using IPsec. Connections will be secured using the settings in IPsec properties and rules in the Connection Security Rule node.

☐ **Block the connection**

[Learn more about actions](#)

< Back Next > Cancel

Click Next and keep the default configurations.

Profile
Specify the profiles for which this rule applies.

Steps:

- Rule Type
- Protocol and Ports
- Action
- Profile
- Name

When does this rule apply?

- ☒ **Domain**
Applies when a computer is connected to its corporate domain.
- ☒ **Private**
Applies when a computer is connected to a private network location.
- ☒ **Public**
Applies when a computer is connected to a public network location.

[Learn more about profiles](#)

< Back Next > Cancel

Click Next and enter the rule name (for example, "RemoteDesktop"), and click Finish.

Name

Specify the name and description of this rule.

Steps:

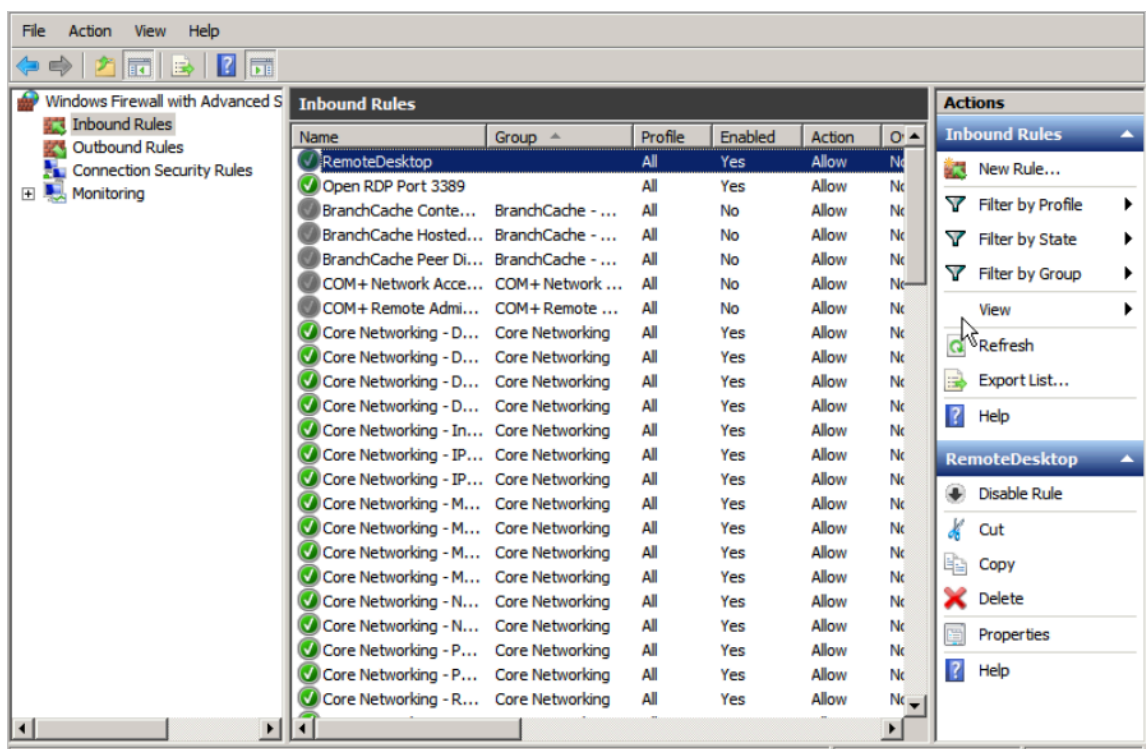
- Rule Type
- Protocol and Ports
- Action
- Profile
- Name**

Name: RemoteDesktop

Description (optional):

< Back Finish Cancel

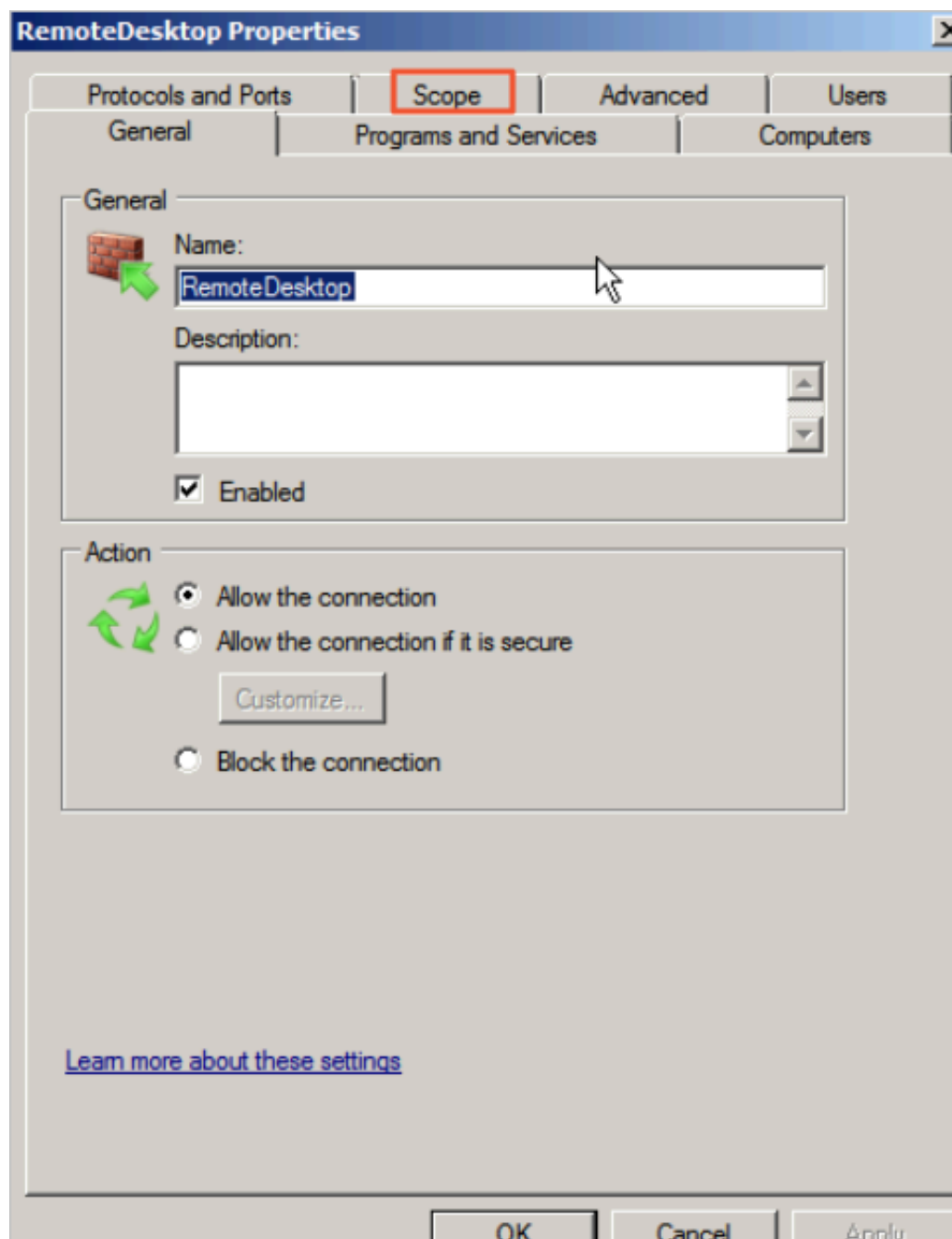
The new rule is shown in the Inbound Rule list.



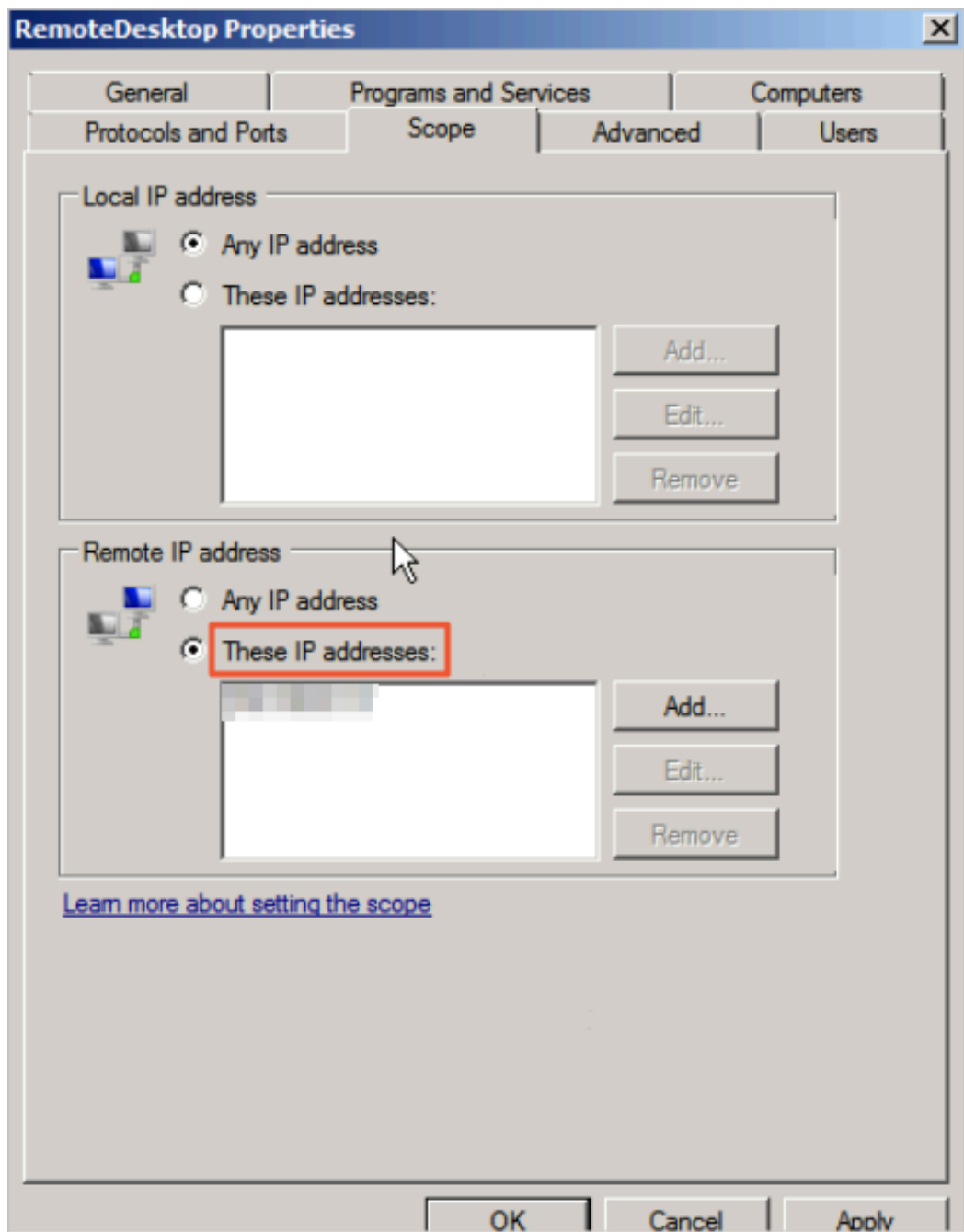
With the above steps, the remote port is added to WFAS, but access restriction is still not implemented. Let's implement it now.

b. Configure the IP address scope

Right-click the just created inbound rule, and select Properties in the context menu. In the displayed dialog box, click the Scope tab. Then add the remote IP addresses that can access this ECS instance. Note that once the IP address settings here are enabled, other IP addresses will be unable to access this ECS instance.

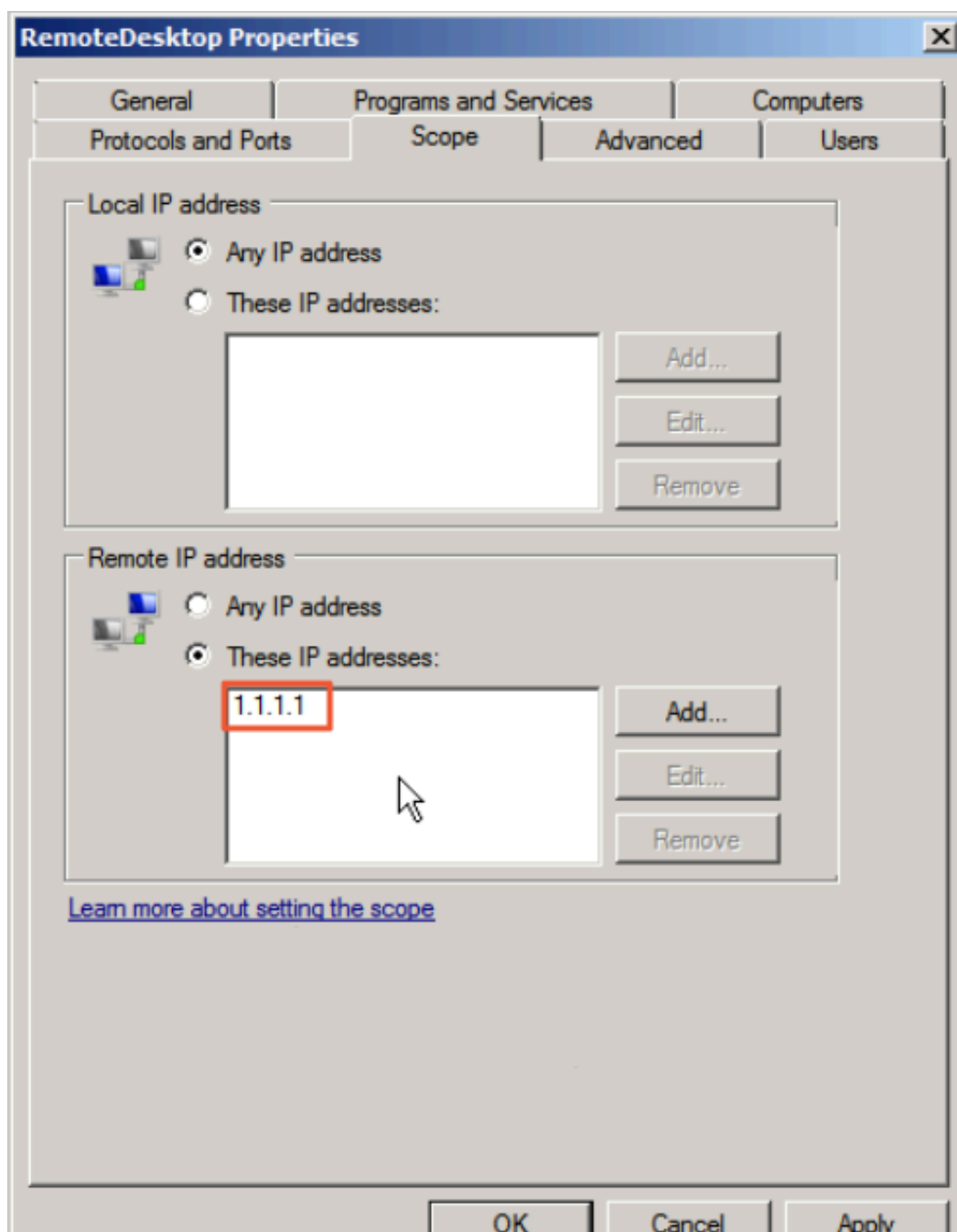


Add remote IP addresses.

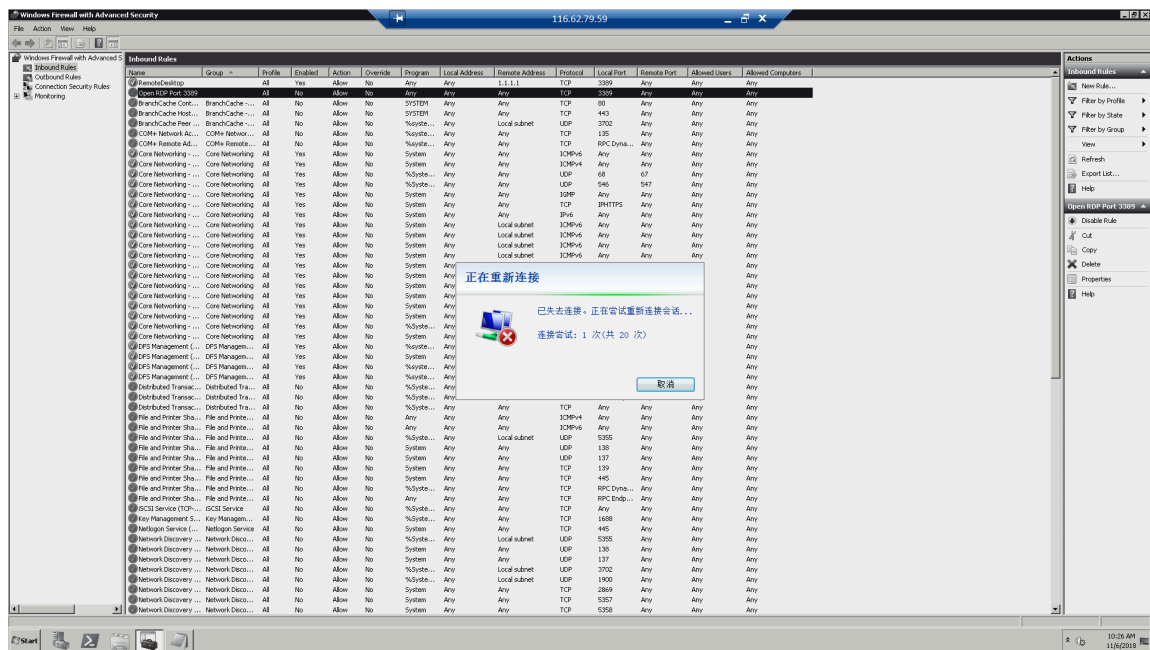


c. Validate the IP address scope

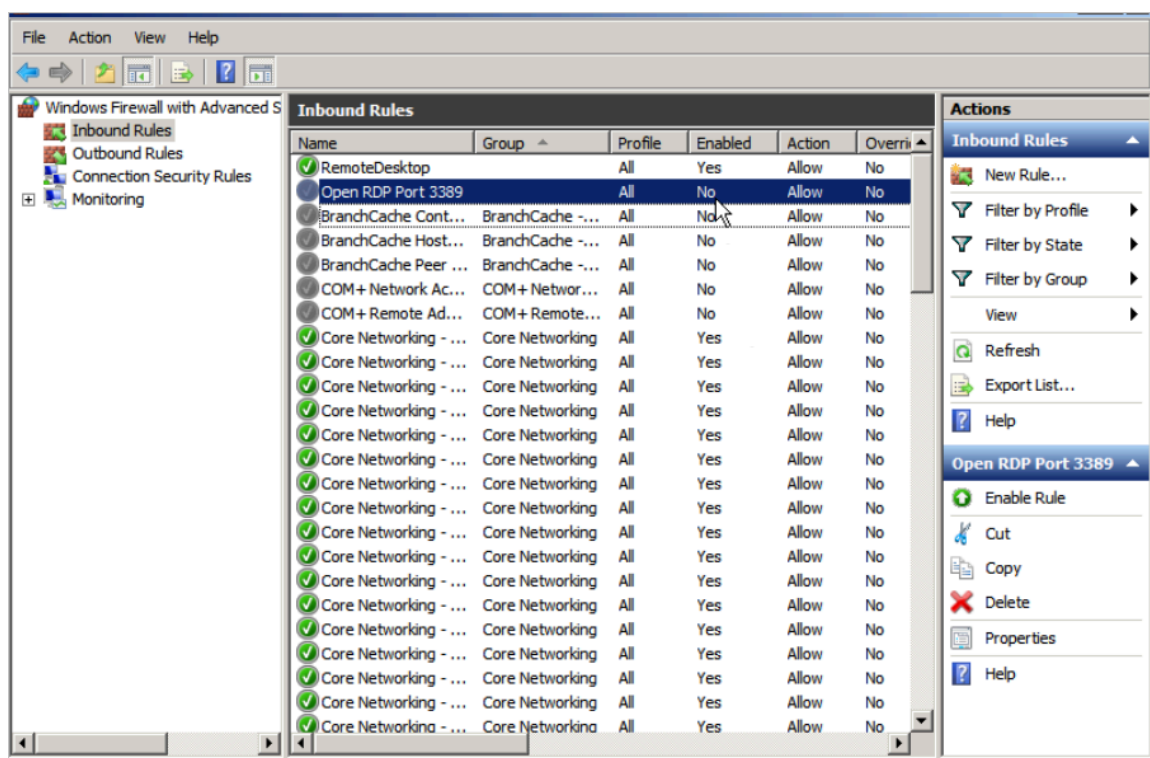
Let's add an IP address arbitrarily in the Remote IP address box and see what happens to the remote connection.



The remote connection is down.



If the remote connection is still up, we can just disable the Open RDP Port 3389 rule.



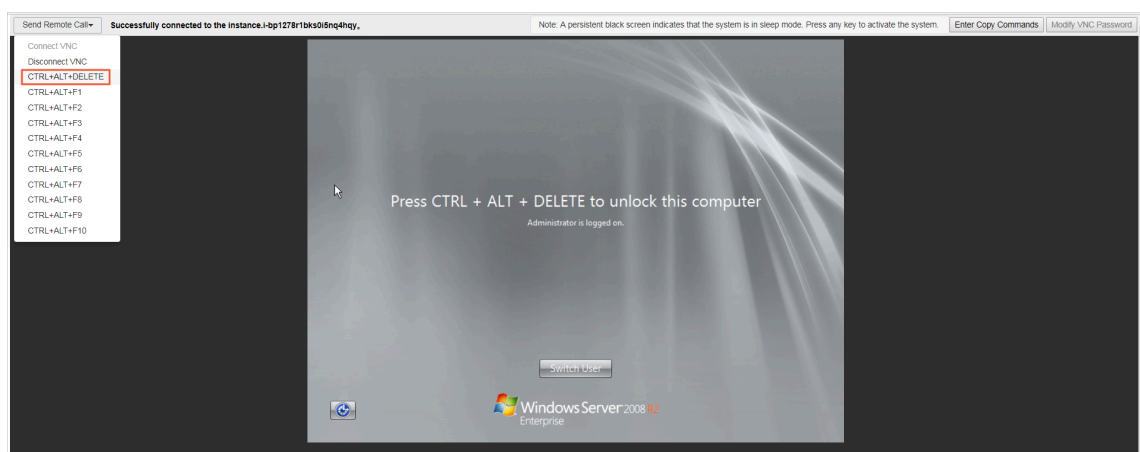
If the remote connection is down, it means that the IP address scope has taken effect. However, we cannot connect to the ECS instance ourselves now. What should we do? We now can turn to the ECS console. Log on to the ECS console, and replace the remote IP address previously configured in the Scope tab with

our own address (enter the Internet address unless your work environment is connected to Alibaba Cloud). You can connect to the ECS instance again now.

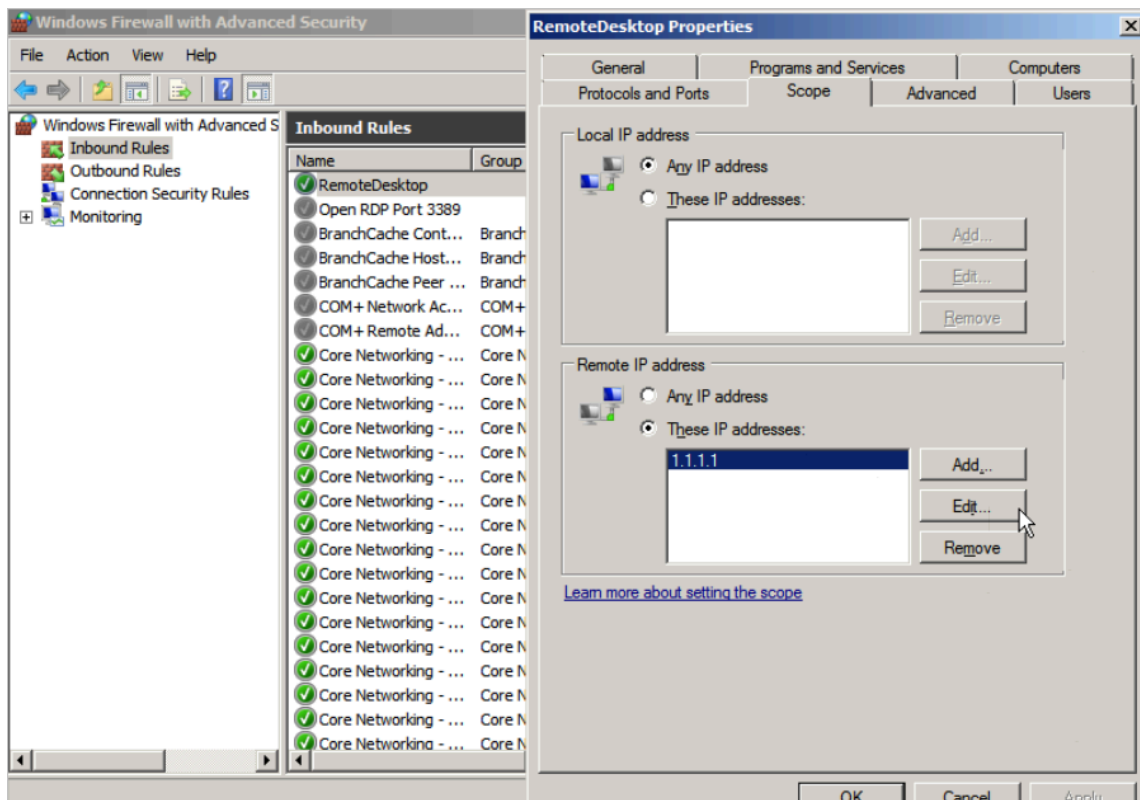
Enter the ECS console, find the corresponding instance, and then connect to it.



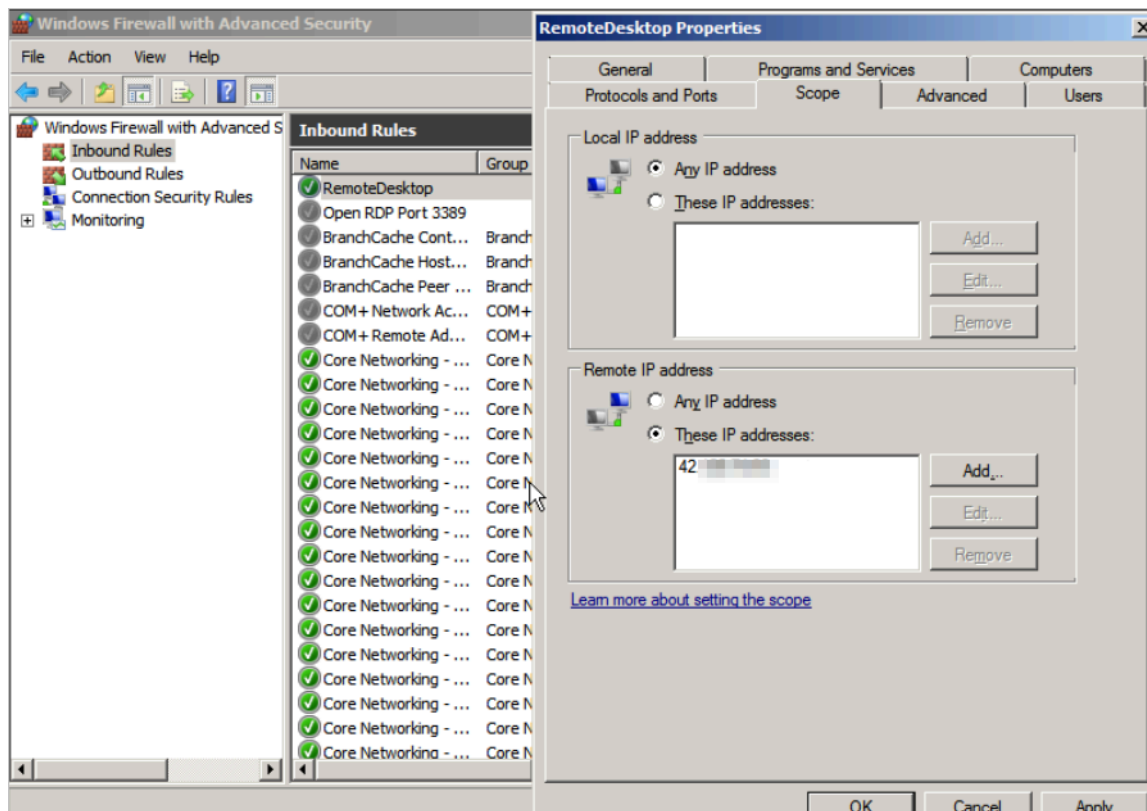
Log on to the ECS instance.



Modify the remote IP address in the Scope tab of the RemoteDesktop rule in the same way. Specifically, replace 1.1.1.1 with our own IP address.



Now we can connect to the ECS instance normally after adding our IP address. If you do not know your Internet address, you can [click here](#) to view it.



The above steps implement remote access restriction on an ECS instance through WFAS. For other services and ports, restrictions can be implemented in the same way, for example, disabling ports 135, 137, 138, and 445 that are not used frequently, limiting access to FTP and related services, and more, thus maximizing the protection of ECS instances.

Command line operations

1. Export the firewall configurations to a file.

```
netsh advfirewall export c:\adv.pol
```

2. Import the firewall configuration file to the system.

```
netsh advfirewall import c:\adv.pol
```

3. Restore the default firewall settings.

```
Netsh advfirewall reset
```

4. Disable the firewall.

```
netsh advfirewall set allprofiles state off
```

5. Enable the firewall.

```
netsh advfirewall set allprofiles state on
```

6. Configure to block inbound traffic and allow outbound traffic by default in all configuration files.

```
netsh advfirewall set allprofiles firewallpolicy blockinbound,  
allowoutbound
```

7. Delete the rule named “ftp” .

```
netsh advfirewall firewall delete rule name=ftp
```

8. Delete all inbound rules for local port 80.

```
netsh advfirewall firewall delete rule name=all protocol=tcp  
localport=80
```

9. Add the RemoteDesktop rule to allow port 3389.

```
netsh advfirewall firewall add rule name=RemoteDesktop (TCP-In-3389  
) protocol=TCP dir=in localport=3389 action=allow
```

References

[How to restrict the access of ports/IP addresses/applications using Windows 2008/2012 Firewall](#)

More open source software are available at [Alibaba Cloud Marketplace](#)

1.9 Isolation of instances within a security group

A security group is a virtual firewall that provides Stateful Packet Inspection (SPI) and packet filtering. It contains instances in the same region with the same security

requirements and mutual trust. Alibaba Cloud provides various access control policies to allow you isolate instances within a security group.

Intra-group isolation rules

- Network isolation in a security group is implemented between network interfaces, not between instances. If multiple Elastic Network Interfaces (ENIs) are bound to an instance, you need to set isolation rules for each ENI.
- Instances in a security group can access each other by default, which is not changed by the isolation rules.

Intra-group isolation rules are user-defined access control policies, and are invalid for the default security groups and new security groups. The default access control policy for a security group is: instances in the same security group can access each other over the intranet, while instances in different security groups cannot.

- Intra-group isolation rules have the lowest priority.

To isolate instances in a security group, make sure no intercommunication rules apply to them except for the isolation rules. In the following cases, instances can still access each other even though intra-group isolation rules are set:

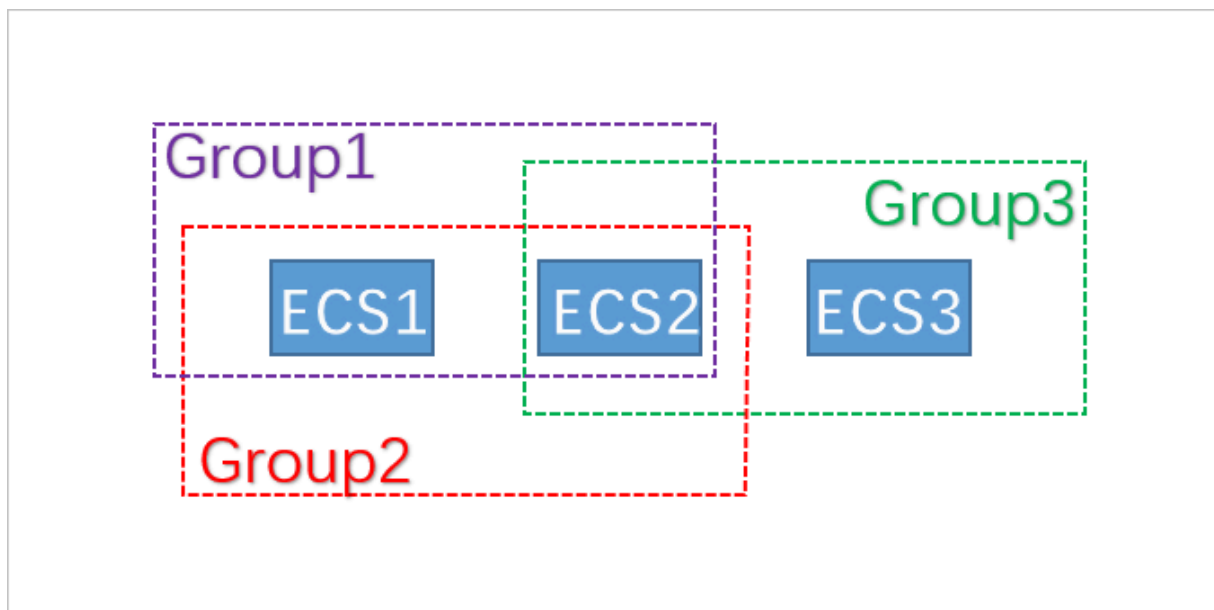
- Intra-group isolation rules are set in a security group, while an Access Control List (ACL) that permits intra-group communication between instances is set at the same time.
- Intra-group isolation rules are set in a security group, while intra-group intercommunication is configured at the same time.
- Intra-group isolation rules only apply to the instances in the current security group
-

Modify the access control policy

You can use the [ModifySecurityGroupPolicy](#) interface to modify the access control policy within a security group.

Case analysis

The following figure shows the relationship between three instances and their security groups.



In this example, Group1, Group2, and Group3 are three different security groups. ECS1, ECS2, and ECS3 are three different ECS instances. ECS1 and ECS2 belong to Group1 and Group2. ECS2 and ECS3 belong to Group3.

The intra-group intercommunication policies of the three security groups are as follows:

Security group	Intra-group intercommunication policy	Instances included
Group1	Isolated	ECS1 and ECS2
Group2	Interconnected	ECS1 and ECS2
Group3	Interconnected	ECS2 and ECS3

The communication status between instances is as follows:

Instance	Interconnected or isolated?	Reason
ECS1 and ECS2	Interconnected	ECS1 and ECS2 belong to both Group1 and Group2. The policy of Group1 is "isolated", while that of Group2 is "interconnected". As intra-group isolation has the lowest priority, ECS1 and ECS2 are interconnected.
ECS2 and ECS3	Interconnected	Both ECS2 and ECS3 belong to Group3. The policy of Group3 is "interconnected", so ECS2 and ECS3 are interconnected.

Instance	Interconnected or isolated?	Reason
ECS1 and ECS3	Isolated	ECS1 and ECS3 belong to different security groups . Instances in different security groups are not interconnected by default. To permit access between instances in two security groups, you can authorize security groups through security group rules.

1.10 Security group quintuple rules

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

Background information

Previously, security group rules have the following characteristics:

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

In most scenarios, these types of security group rules simplify the setup process, but possess the following disadvantages:

- The source port range of an ingress rule is not restricted. That is, all source ports are permitted by default.
- The destination IP address of an ingress rule is not restricted. That is, all IP addresses in a security group are permitted by default.
- The source port range of an egress rule is not restricted. That is, all source ports are permitted by default.
- The source IP address of an egress rule is not restricted. That is, all IP addresses in a security group are permitted by default.

Definition of a quintuple rule

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

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

The following shows an example quintuple rule:

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

The example egress rule indicates that 172.16.1.0/32 is prohibited from accessing 10.0.0.1/32 from port 22 through TCP.

Scenarios

- Some platform products are connected to the solutions of third-party vendors to provide them with network services. To prevent these products from illegally accessing users' ECS instances, it is needed to set quintuple rules in the security group to control the inbound and outbound traffic more precisely.
- If your instances are isolated within a security group due to settings, and you want to precisely control the access between several ECS instances in the group, you can set security group quintuple rules to meet your needs.

How to configure quintuple rules

You can use OpenAPI to set quintuple rules.

- To add a security group ingress rule, see [AuthorizeSecurityGroup](#).
- To add a security group egress rule, see [AuthorizeSecurityGroupEgress](#).
- To delete a security group ingress rule, see [RevokeSecurityGroup](#).
- To delete a security group egress rule, see [RevokeSecurityGroupEgress](#).

Parameters

The following table describes the parameters.

Parameter	Meaning in ingress rules	Meaning in egress rules
SecurityGroupId	The ID of the security group to which the current ingress rule belongs (that is, the ID of the destination security group).	The ID of the security group to which the current egress rule belongs (that is, the ID of the source security group).
DestCidrIp	Destination IP address range; optional. <ul style="list-style-type: none"> • If DestCidrIp is specified, you can control the destination IP address range in an ingress rule more precisely. • If DestCidrIp is not specified, the IP address range in an ingress rule includes all the IP addresses in the security group indicated by the SecurityGroupId. 	Destination IP addresses. Either DestGroupId or DestCidrIp must be specified. If both are specified, DestCidrIp takes priority.
PortRange	Destination port range; required.	Destination port range; required.
DestGroupId	Input not allowed. The destination security group ID must be a SecurityGroupId.	The destination security group ID. Either DestGroupId or DestCidrIp must be specified. If both are specified, DestCidrIp takes priority.
SourceGroupId	The source security group ID. Either SourceGroupId or SourceCidrIp must be specified. If both are specified, SourceCidrIp takes priority.	Input not allowed. The source security group ID in an egress rule must be a SecurityGroupId.

Parameter	Meaning in ingress rules	Meaning in egress rules
SourceCidrIp	Source IP address range. Either SourceGroupId or SourceCidrIp must be specified. If both are specified, SourceCidrIp takes a higher priority.	Source IP address range; optional. <ul style="list-style-type: none">· If SourceCidrIp is specified, you can control the source IP address range in an egress rule more precisely.· If SourceCidrIp is not specified, the source IP addresses in an egress rule include all the IP addresses in the security group indicated by the SecurityGroupId.
SourcePort Range	Source port range; optional. If it is not specified, source ports are not restricted.	Source port range; optional. If it is not specified, source ports are not restricted.

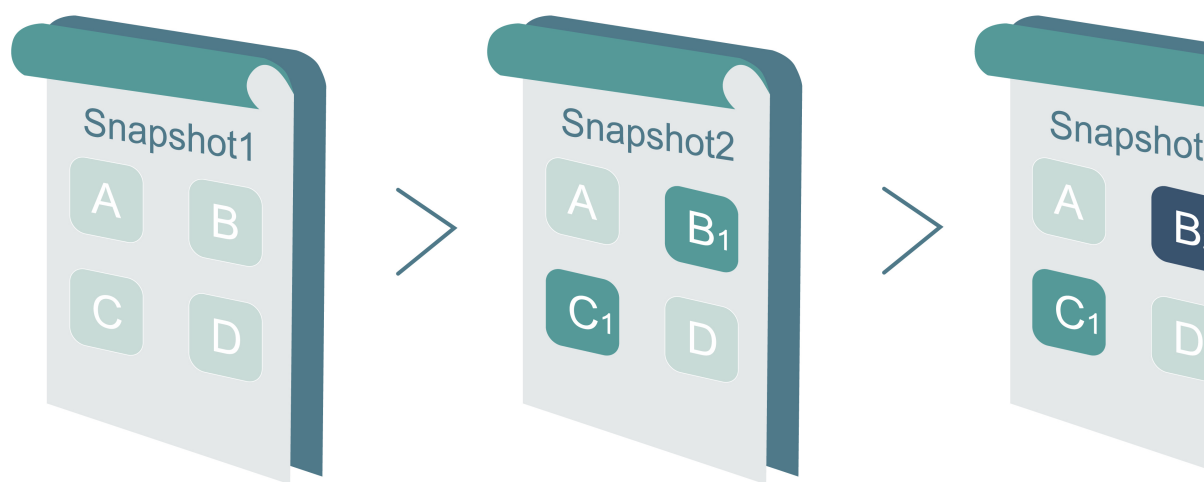
2 Disaster recovery solutions

Disaster recovery solutions help guarantee the running stability and data security of your IT system. Specifically, the solutions incorporate data backup and disaster recovery of systems and applications. Alibaba Cloud ECS allows you to use snapshots and images for data backup.

Disaster recovery methods

- Snapshot backup

Alibaba Cloud ECS allows you to back up system disks and data disks with snapshots. Currently, Alibaba Cloud provides the Snapshot 2.0 service, which features a higher snapshot quota and a more flexible automatic task strategy than previous snapshot services, helping to reduce impact on business I/O. When snapshots are used for data backup, the first backup is a full backup, followed by incremental backups. The backup duration depends on the amount of data to be backed up.



As shown 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 are copied to the snapshot. Alibaba Cloud ECS allows you to configure manual or automatic snapshot backup. With automatic backup, you can specify the time of day (24 options, on the hour), recurring day of week (Monday through Sunday), and retention time for snapshot creation. The retention time is customizable, and you can set a value from 1 to 65,536 days or choose to save snapshots permanently.

- Snapshot rollback

If exceptions occur in your system and you need to roll back a disk to a previous state, you can [roll back the disk](#) so long as it has a corresponding snapshot created.

Note:

- Rolling back a disk is an irreversible action. After disk rollback is completed, data cannot be restored. Exercise caution when performing this action.
- After a disk is rolled back, data will be irretrievably lost from the creation time of the snapshot to the current time.

- Image backup

An image file is equivalent to a replica file that contains all the data from one or more disks (a system disk or both the system disk and data disks). All image backups are full backups and can only be triggered manually.

- Image recovery

You can create custom images from snapshots to include the operating system and data environment in the image. The custom images can then be used to create multiple instances with the same operating system and data environment. For the configuration of snapshots and images, see [Snapshots](#) and [Images](#).



Note:

Custom images cannot be used across regions.

Technical metrics

RTO and RPO: related to the amount of data, usually at an hourly level.

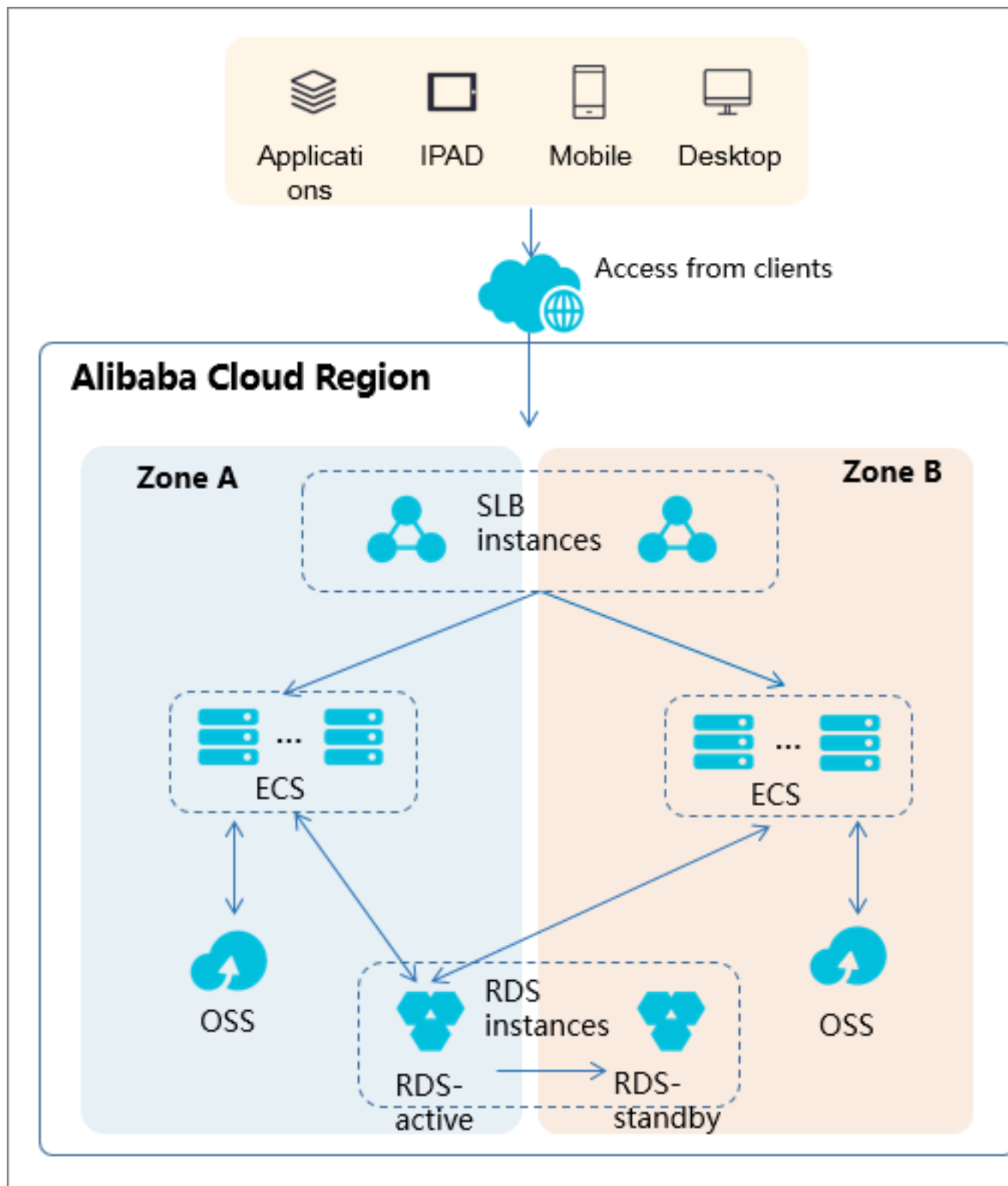
Scenarios

- Backup and restoration

Alibaba Cloud ECS allows you to back up system disks and data disks with snapshots and images. If incorrect data is stored on a disk due to data errors caused by application errors, or hackers exploiting application vulnerabilities for malicious access, you can use the snapshot service to restore the disk to a desired state. In addition, Alibaba Cloud ECS allows you to reinitialize disks with images or purchase new ECS instances with a custom image.

- Disaster recovery application

Alibaba Cloud ECS supports the implementation of disaster recovery architecture. For example, you can buy and use a Server Load Balancer (SLB) at the front end of an application, and deploy at least two ECS instances at the back end of the same application. Alternatively, you can implement an Auto Scaling solution using the auto scaling technology provided by Alibaba Cloud by defining how to use the ECS resources. In this way, even if one of the ECS instances fails or resources are overused, business continuity will not be interrupted, thus realizing disaster recovery. Take the deployment of ECS instances in two Internet Data Centers (IDCs) in the same city for example:



- A cluster of ECS instances is deployed in both IDCs. At the access side, SLBs are used for load balancing between the two IDCs.
- The Region Master nodes in both IDCs are identical and operate in active/standby mode. The failure of one node does not affect the ECS control function.

- To switch the control node of ECS instances in the case of IDC failure, the middleware domain name is associated anew as it is used for controlling the cluster. If the IDC of the control node experiences problems, the middleware domain name needs to be associated with the control node of the other IDC.

3 Data recovery

3.1 How to restore the data that is deleted by mistake

By taking CentOS 7 for example, this document introduces how to use Extundelete, an open source tool, to quickly restore accidentally deleted data.

Overview

In practice, data may be deleted accidentally. In this case, how to restore the data quickly and effectively? Alibaba Cloud offers several ways to restore data, for example :

- Roll back a *snapshot* or *custom image* through the ECS console.
- Purchase several ECS instances to implement *load balancing* and high availability for your services.
- Use *Object Storage Service (OSS)* to store a massive amount of data such as web pages, images, and videos.

There are a variety of open source data recovery tools for Linux, such as debugfs, R-Linux, ext3grep, Extundelete, and more. Of them, ext3grep and Extundelete are generally used. Both tools adopt the same recovery techniques, just that Extundelete is more powerful.

Extundelete is a Linux-based open source data recovery software. When using Linux instances, you can install this tool conveniently to quickly restore the data deleted accidentally as no Recycle Bin is available in Linux.

Extundelete can locate the position of an inode block by combining the inode information and logs so as to search for and restore the desired data. This powerful tool supports the disk-wide restoration of ext3/ext4 dual-format partitions.

Once data is deleted accidentally, firstly you need to unmount the disk or disk partition that contains the deleted data. This is because after a file is deleted, only the inode pointers of the file are zeroed while the actual file is still stored on the disk. If the disk is mounted in read/write mode, data blocks of the deleted file may be reallocated by the operating system. Once the data blocks are overwritten by new data, the original data will be lost completely, and cannot be restored by any means.

Therefore, mounting a disk in read-only mode can reduce the risk of overwriting the data in blocks, thus improving the chances of restoring the data successfully.

**Note:**

During the online restoration process, do not install Extundelete on the disk that has the deleted file. Otherwise, the data to be restored might be overwritten. Keep in mind to back up the disk by taking a snapshot before any operations.

Intended audience

- Users who accidentally delete files on a disk and no write operations have been performed on the disk after the deletion.
- Users whose websites have low traffic and who have few ECS instances.

Procedure

Software release: e2fsprogs-devel e2fsprogs gcc-c++ make (compiler and more)
Extundelete-0.2.4.

**Note:**

The libext2fs 1.39 or above is required for the normal operation of Extundelete. For ext4 support, however, make sure e2fsprogs 1.41 or higher is provided (you can run the command `dumpe2fs` to check the version output).

The above releases are available when this document is being written. Your downloads may be different.

- Deploy Extundelete

```
wget http://zy-res.oss-cn-hangzhou.aliyuncs.com/server/extundelete-0.2.4.tar.bz2
yum -y install bzip2 e2fsprogs-devel e2fsprogs gcc-c++ make
#Install related dependencies and libraries
tar -xvzf extundelete-0.2.4.tar.bz2
cd extundelete-0.2.4
program directory #Enter the
```



```
./configure
successfully as shown below
```

#Installed

```
extundelete-0.2.4/src/Makefile.am
extundelete-0.2.4/configure.ac
extundelete-0.2.4/depcomp
extundelete-0.2.4/Makefile.in
extundelete-0.2.4/Makefile.am
[root@iZy930wmhyutC2Z ~]# cd extundelete-0.2.4
[root@iZy930wmhyutC2Z extundelete-0.2.4]# ./configure
Configuring extundelete 0.2.4
Writing generated files to disk
[root@iZy930wmhyutC2Z extundelete-0.2.4]# █
```

```
make && make install
```

At this point, the `src` directory appears. It contains an Extundelete executable file and a corresponding path. As shown below, the default file is installed in `usr/local/bin`, and the following demo is made in the `usr/local/bin` directory.

- Delete a file and use Extundelete to restore it

1. Check the available disks and partitions of your ECS instance, then format and partition the `/dev/vdb` partition. For more information about formatting and partitioning, see [Format and mount a data disk](#).

```
fdisk -l
```

```
Disk label type: dos
Disk identifier: 0x0000efd2

   Device Boot      Start         End      Blocks    Id  System
/dev/vda1  *           2048     83886079     41942016    83  Linux

Disk /dev/vdb: 21.5 GB, 21474836480 bytes, 41943040 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
```

2. Mount the partitioned disk under the `/zhuyun` directory, and then create a file named `hello`.

```
mkdir /zhuyun
directory.
mount /dev/vdb1 /zhuyun
the zhuyun directory.
```

#Create the zhuyun

#Mount the disk under

```
echo test > hello #Create a test file.
```

3. Run the `md5sum` command to generate the MD5 value of the file and note it down. You can compare the MD5 values of the file before and after the deletion to verify its integrity.

```
md5sum hello
```

```
[root@iZbp13micdqi2364umm8aZ zhuyun]# md5sum hello
d8e8fca2dc0f896fd7cb4cb0031ba249 hello
```

4. Delete the `hello` file.

```
rm -rf hello
cd ~
fuser -k /zhuyun #Terminate the process tree
that uses a certain partition (skip this if you are sure that no
resources are occupied).
```

5. Unmount the data disk

```
umount /dev/vdb1 #Before using any file
restoration tool, unmount or mount the partitions to be restored
in read-only mode to prevent their data from being overwritten.
```

6. Use `Extundelete` to restore the file.

```
extundelete --inode 2 /dev/vdb1 #Query the contents in a
certain inode. Using "2" means to search the entire partition. To
search a directory, just specify the inode and directory. Now
you can see the deleted file and inode.
```

```
Direct blocks: 127754, 4, 0, 0, 1, 9252, 0, 0, 0, 0, 0, 0
Indirect block: 0
Double indirect block: 0
Triple indirect block: 0
```

File name	Inode number	Deleted status
.	2	
..	2	
lost+found	11	
hello	12	Deleted

```
/usr/local/bin/extundelete --restore-inode 12 /dev/vdb1 #
Restore the deleted file.
```

At this point, the `RECOVERED_FILES` directory appears under the directory where the command is executed. Check whether the file is restored.

```
[root@iZbp13micdqi2364umm8aZ /]# ll RECOVERED_FILES/
total 4
-rw-r--r-- 1 root root 5 Mar  8 14:20 hello
```

Check the MD5 values of the files before and after deletion. If they are the same, restoration is successful.

```
--restore-inode 12          # --restore-inode Restore by  
the specified inode.        # --restore-inode Restore by  
--extundelete --restore-all # --restore-all   Restore all.
```

3.2 Data restoration in Linux instances

When solving problems related to disks, you may frequently encounter the loss of data disk partitions. This article describes common data partition loss problems and corresponding solutions in Linux, and provides common mistakes and best practices for cloud disks to avoid possible risks of data loss.

Before restoring data, you must create snapshots for data disks that lose partitions. If problems occur during the restoration process, you can roll back data disks to the status before restoration.

Prerequisites

Before restoring data, you must create snapshots for data disks that lose partitions. If problems occur during the restoration process, you can roll back data disks to the status before restoration.

Introduction to disk management tools

You can select one of the following tools to fix the disk partition and restore the data in a Linux instance:

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

Handle data disk partition loss and data restoration in Linux

After you restart a Linux instance, you may encounter data disk partition loss or data loss issues. This may be because you have not set the partitions to be mounted automatically on startup of the instance in the `etc/fstab` file. In this case, you

can manually mount the data disk partition first. If the system prompts partition table loss when you manually mount the data disk, you can try to solve the problem through the following three methods: [Restore partitions by using fdisk](#), [Restore partitions by using testdisk](#), or [Restore data by using testdisk](#).

- Restore partitions by using fdisk

Default values usually apply to the starting and ending sectors of the partition when you partition a data disk. You can then directly use fdisk to restore the partition. For more information about this tool, see [Linux Format and mount a data disk](#).

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

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

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

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

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

If the preceding operations do not help, you can try testdisk for the restoration.

- Restore partitions by using testdisk

Here we suppose the cloud disk device is named `/dev/xvdb`. Follow these steps to restore the partitions by using testdisk:

1. Run `testdisk /dev/xvdb` (replace the device name as appropriate), and then select Proceed (default value) and press the Enter key.

```

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

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

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

>[Proceed] [ Quit ]

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

```

2. Select the partition table type for scanning: *Intel* by default. If your data disk uses the GPT format, select *EFI GPT*.

```

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

Disk /dev/xvdb - 5368 MB / 5120 MiB

Please select the partition table type, press Enter when done.
>[Intel] Intel/PC partition
[EFI GPT] EFI GPT partition map (Mac i386, some x86_64...)
[Humax] Humax partition table
[Mac] Apple partition map
[None] Non partitioned media
[Sun] Sun Solaris partition
[XBox] Xbox partition
[Return] Return to disk selection

Note: DO NOT select 'None' for media with only a single partition. It's very
rare for a disk to be 'Non-partitioned'.

```

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

```

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

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

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

```

4. If you cannot see any partition, select *Quick Search* and then press the Enter key for a quick search.

```

Disk /dev/xvdb - 5368 MB / 5120 MiB - CHS 652 255 63
Current partition structure:
    Partition                Start          End      Size in sectors
No partition is bootable

*=Primary bootable P=Primary L=Logical E=Extended D=Deleted
>[Quick Search]
Trv to locate partition

```

The partition information is displayed in the returned result, as shown in the following figure.

```

Disk /dev/xvdb - 5368 MB / 5120 MiB - CHS 652 255 63
    Partition                Start          End      Size in sectors
>* Linux                    0 32 33      652 180 40    10483712

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

```

5. Select the partition and press the Enter key.
6. Select *Write* to save the partition.



Note:

Select *Deeper Search* to continue searching if the expected partition is not listed.

```

Disk /dev/xvdb - 5368 MB / 5120 MiB - CHS 652 255 63
    Partition                Start          End      Size in sectors
1 * Linux                    0 32 33      652 180 40    10483712

[ Quit ] [Deeper search] >[ write ]
write partition structure to disk

```

7. Press the Y key to save the partition.

```

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

write partition table, confirm ? (Y/N)

```

8. Run `partprobe /dev/xvdb` (replace the device name as appropriate) to refresh the partition table manually.

9. Mount the partition again and view the data in the data disk.

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

- Restore data by using testdisk

In some cases, you can use testdisk to scan and locate the disk partition, but you cannot save the partition. In this case, you can try to restore files directly. Follow these steps:

1. Find the partition following Step 1 to Step 4 described in [Restore partitions by using testdisk](#).
2. List files by pressing the P key. The returned result is shown in the following figure.

```
* Linux
Directory /
0 32 33 652 180 40 10483712
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 .
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 ..
drwx----- 0 0 16384 21-Feb-2017 11:56 lost+found
-rw-r--r-- 0 0 1701 21-Feb-2017 11:57 install_edsd.sh
-rw-r--r-- 0 0 5848 21-Feb-2017 11:57 install.sh
>-rw-r--r-- 0 0 12136 21-Feb-2017 11:57 ip.qz
-rw-r--r-- 0 0 0 21-Feb-2017 11:57 test
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 123.sh
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 configclient
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 data
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 diamond
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 logs

Next
Use Right to change directory, h to hide deleted files
q to quit, : to select the current file, a to select all files
C to copy the selected files. c to copy the current file
```

3. Select the files to restore, and press the C key.
4. Select a directory. In this example, the file is restored and copied to the `/home` directory.

```

Please select a destination where /ip.gz will be copied.
Keys: Arrow keys to select another directory
      C when the destination is correct
      Q to quit
Directory /
drwxr-xr-x 0 0 4096 11-Jan-2017 09:32 .
drwxr-xr-x 0 0 4096 11-Jan-2017 09:32 ..
dr-xr-xr-x 0 0 4096 25-Jul-2016 16:23 boot
drwxr-xr-x 0 0 2940 21-Feb-2017 12:30 dev
drwxr-xr-x 0 0 4096 21-Feb-2017 12:12 etc
>drwxr-xr-x 0 0 4096 16-Feb-2017 11:48 home
drwx----- 0 0 16384 12-May-2016 19:58 lost+found
drwxr-xr-x 0 0 4096 12-Aug-2015 22:22 media
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 mnt
drwxr-xr-x 0 0 4096 12-Aug-2015 22:22 opt
dr-xr-xr-x 0 0 0 16-Feb-2017 21:35 proc
dr-xr-x--- 0 0 4096 21-Feb-2017 11:57 root
drwxr-xr-x 0 0 560 21-Feb-2017 12:12 run
drwxr-xr-x 0 0 4096 12-Aug-2015 22:22 srv
dr-xr-xr-x 0 0 0 16-Feb-2017 21:35 sys
drwxrwxrwt 0 0 4096 21-Feb-2017 12:34 tmp
drwxr-xr-x 0 0 4096 16-Feb-2017 11:48 usr
drwxr-xr-x 0 0 4096 16-Feb-2017 21:35 var
lrwxrwxrwx 0 0 7 3-May-2016 13:48 bin
lrwxrwxrwx 0 0 7 3-May-2016 13:48 lib
lrwxrwxrwx 0 0 9 3-May-2016 13:48 lib64
lrwxrwxrwx 0 0 8 3-May-2016 13:48 sbin

```

If you see Copy done! 1 ok, 0 failed, it indicates that copy was successful, as shown in the following figure.

```

* Linux 0 32 33 652 180 40 10483712
directory /
Copy done! 1 ok, 0 failed
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 .
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 ..
drwx----- 0 0 16384 21-Feb-2017 11:56 lost+found
-rw-r--r-- 0 0 1701 21-Feb-2017 11:57 install_edsd.sh
-rw-r--r-- 0 0 5848 21-Feb-2017 11:57 install.sh
>-rw-r--r-- 0 0 12136 21-Feb-2017 11:57 ip.gz
-rw-r--r-- 0 0 0 21-Feb-2017 11:57 test
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 123.sh
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 configclient
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 data
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 diamond
drwxr-xr-x 0 0 4096 21-Feb-2017 11:57 logs

```

5. Switch to the `/home` directory to view details. If you can see files, it indicates that files have been restored successfully.

```

[root@Aliyun /]# ls /home/
admin ip.gz
[root@Aliyun /]#

```


Common mistakes and best practices

Data is users' core asset. Many users establish websites and databases (MYSQL/MongoDB/Redis) on ECS. Huge risks to the users' services may occur when data is lost. Common mistakes and best practices are summarized as follows.

- Common mistakes

The bottom layer of Alibaba Cloud block-level storage is based on [triplicate technology](#). Therefore, some users consider that no risk of data loss in the operating system exists. It is actually a misunderstanding. The three copies of data stored in the bottom layer provide physical layer protection for data disks. However, if problems occur to the cloud disk logic in the system, such as viruses, accidental data deletion, and file system damage, the data may still be lost. To guarantee data security, you have to use technologies such as Snapshot and backup.

- Best practices

Data disk partition restoration and data restoration are the final solutions for solving data loss problems, but it is never guaranteed. We strongly recommend that you follow the best practices to perform auto or manual snapshot on data and run different backup schemes to maximize your data security.

- Enable automatic snapshots

Automatic snapshots are enabled for the system disk and data disk based on actual service conditions. Note that automatic snapshot may be released when the system disk is changed, the instance is expired, or the disk is manually released.

You log on to the ECS console to change the attributes of the disks to enable snapshot release with the disk. Disable snapshot release with the disk if you want to retain the snapshots.

For more information, see [FAQ about automatic snapshots](#).

- Create manual snapshots

Create snapshots manually before any important or risky operations such as:

- Upgrade the kernel
- Upgrade or change of applications
- Restoration of disk data

You must create snapshots for disks before restoring them. After the snapshots are completed, you can perform other operations.

- OSS, offline, or offsite backup

You can back up important data by means of OSS, offline, or offsite backup based on actual conditions.

3.3 Data restoration in Windows instances

When solving problems related to disks, you may frequently encounter the loss of data disk partitions. This article describes common data partition loss problems and corresponding solutions in Windows, and provides common mistakes and best practices for cloud disks to avoid possible risks of data loss.

Prerequisites

Before restoring data, you must create snapshots for data disks that lose partitions. If problems occur during the restoration process, you can roll back data disks to the status before restoration.

Introduction to disk management tools

In Windows instances, you can select either of the following tools for restoring data disk data:

- **Disk Management:** A tool provided by Windows for partitioning and formatting the disk.
- **Data restoration software:** Generally, they are commercial software, and can be downloaded from the providers' official websites. They are mainly used for restoring data in an abnormal file system.

Status of the disk is Foreign and no partitions are displayed

In the Disk Management of Windows, the disk is in the Foreign status and displays no partitions.

Solution:

Right click the Foreign disk, select Import Foreign Disks, and then click OK.

Status of the disk is Offline and no partitions are displayed

In the Disk Management of Windows, the disk is in the Offline status and displays no partitions.

Solution:

Right click the Offline disk (for example, Disk 1), select Online, and then click OK.

No drive letter assigned

In the Disk Management of Windows, you can view data disk information, but no drive letter is allocated to the data disk.

Solution:

Right click primary partition of the disk (for example, Disk 1), click Change drive letter and paths, and then complete operations by prompt.

Error occurred during storage enumeration

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

**Note:**

Some versions may report Error occurred during enumeration of volumes. They are the same.

Solution:

1. Start Windows PowerShell.
2. Run `winrm quickconfig` for restoring. When “Make these changes [y/n]?” is displayed on the interface, you must type `y` to run the command.

After the restoration, you can have the data disks in the Disk Management.

Data disk is in RAW format

In some special circumstances, the disk in Windows is in RAW format.

If the file system of a disk is unrecognizable to Windows, it is displayed as a RAW disk. This usually occurs when the partition table or boot sector that records the type or location of the file system is lost or damaged. Common causes are listed as follows:

- Safely remove hardware is not used when disconnecting the external disk.
- Disk problems caused by power outages or unexpected shutdown.
- Hardware layer failure may also cause information loss of the disk partition.
- Bottom layer drivers or disk-related applications. For example, DiskProbe can be used to directly modify the disk table structure.

- Computer viruses.

For more information about how to fix these problems, see [Dskprobe Overview](#) document.

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

Common mistakes and best practices

Data is users' core asset. Many users establish websites and databases (MySQL/MongoDB/Redis) on ECS. Huge risks to the users' services may occur when data is lost. Common mistakes and best practices are summarized as follows.

- Common mistakes

The bottom layer of Alibaba Cloud block-level storage is based on [triplicate technology](#). Therefore, some users consider that no risk of data loss in the operating system exists. It is actually a misunderstanding. The three copies of data stored in the bottom layer provide physical layer protection for data disks. However, if problems occur to the cloud disk logic in the system, such as viruses, accidental data deletion, and file system damage, the data may still be lost. To guarantee data security, you have to use technologies such as Snapshot and backup.

- Best practices

Data disk partition restoration and data restoration are the final solutions for solving data loss problems, but it is never guaranteed. We strongly recommend that you follow the best practices to perform auto or manual snapshot on data and run different backup schemes to maximize your data security.

- Enable automatic snapshots

Automatic snapshots are enabled for the system disk and data disk based on actual service conditions. Note that automatic snapshot may be released when the system disk is changed, the instance is expired, or the disk is manually released.

You log on to the ECS console to change the attributes of the disks to enable snapshot release with the disk. Disable snapshot release with the disk if you want to retain the snapshots.

For more information, see [FAQ about automatic snapshots](#).

- **Create manual snapshots**

Create snapshots manually before any important or risky operations such as:

- Upgrade the kernel
- Upgrade or change of applications
- Restoration of disk data

You must create snapshots for disks before restoring them. After the snapshots are completed, you can perform other operations.

- **OSS, offline, or offsite backup**

You can back up important data by means of OSS, offline, or offsite backup based on actual conditions.

4 Configuration preference

4.1 Time setting: Synchronize NTP servers for Windows instances

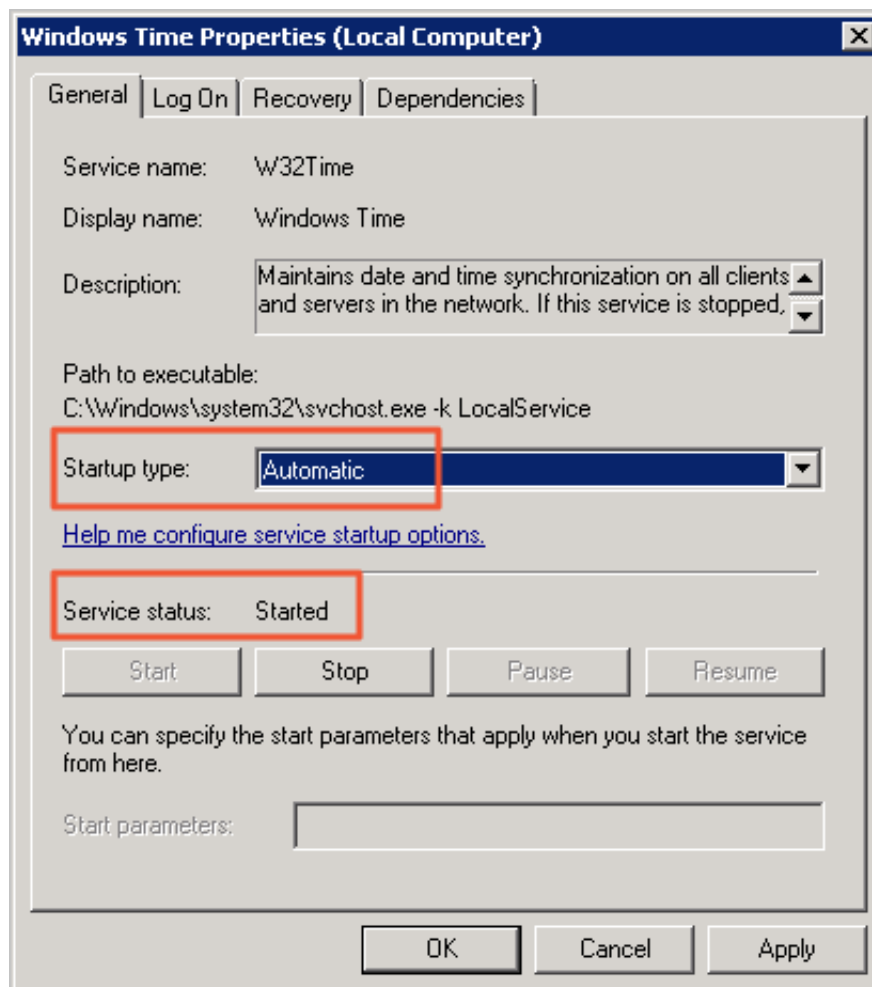
Network Time Protocol (NTP) is a networking protocol for clock synchronization between computer systems over networks. For highly time-sensitive applications (such as those in the communication industry), clock variation between different computers may lead to serious data inconsistencies. You can use the NTP service to synchronize clocks of all servers within the network. The current default time zone for Alibaba Cloud ECS instances across all regions is CST (China Standard Time).

This article describes how use the NTP service to synchronize the clock of a Windows ECS instance running Windows Server 2008 R2 Enterprise Edition x64.

Windows Time service is enabled by default on Windows Server. You must enable the NTP service in the instance to make sure that the NTP service can normally synchronize time after successful NTP service configuration. To check and enable the NTP service, follow these steps:

1. [Connect to a Windows instance](#). Select Start > All Programs > Accessories > Run to open the Run dialog box, and run `services.msc`.
2. In the Services window, double click the Windows Time service.
3. In the Windows Time Properties (Local Computer) dialog box, follow these steps:
 - a. Set Startup type to Automatic.
 - b. Check if the Service status is Started. If not, click Start.

After completing the settings, click Apply, and then click OK.



Modify the default NTP server address

time.windows.com is used as the default NTP server in Windows Server, but synchronization errors may frequently occur due to network issues. When using a Windows instance, you can replace the default NTP server with the intranet NTP server provided by Alibaba Cloud. For more information, see [Internet and intranet NTP servers](#). To modify the default NTP server address, follow these steps:

1. [Connect to a Windows instance](#).
2. In the notification area of the task bar, click Date and Time, and then click Change date and time settings.
3. In the Date and Time dialog box, click the Internet Time tab, and then click Change settings.
4. In the Internet Time Settings dialog box, select Synchronize with an Internet time server, type an Alibaba Cloud intranet NTP server address (for detailed list, see [Internet and intranet NTP servers](#)), and then click Update now.

You are prompted if the synchronization is successful.

Modify NTP synchronization interval

The default NTP synchronization interval is 5 minutes. To modify the NTP synchronization interval, follow these steps:

1. [Connect to a Windows instance](#).
2. Select Start > All Programs > Accessories > Run to open the Run dialog box, and run `regedit`.
3. On the left-side navigation pane of the Registry Editor, find `HKEY_LOCAL_MACHINE/SYSTEM/CurrentControlSet/services/W32Time/TimeProviders/NtpClient`, and then double click `SpecialPollInterval`.
4. In the Edit DWORD (32-bit) Value dialog box, select Decimal as the Base, and then type the Value data as needed. The number you typed is the synchronization interval you need. Unit: seconds.

4.2 Time setting: NTP servers and other public services

Alibaba Cloud ECS provides standard intranet NTP servers, which you can access from your instances. We also provide external NTP services for instances that need the Internet access.

Intranet and Internet NTP servers

To counterbalance the leap seconds in our world, ECS provides free of charge, highly accurate, and reliable NTP service for both classic network- and VPC-Connected instances. Among the NTP servers, the `ntp.cloud.aliyuncs.com` achieves nearly zero difference of atomic reference by synchronizing with satellite services. See the following table for the NTP servers provided by Alibaba Cloud ECS.

Classic network intranet	VPC intranet	Internet
<code>ntp.cloud.aliyuncs.com</code>		<code>ntp1.aliyun.com</code>
<code>ntp1.cloud.aliyuncs.com</code>	<code>ntp7.cloud.aliyuncs.com</code>	<code>ntp2.aliyun.com</code>
<code>ntp2.cloud.aliyuncs.com</code>	<code>ntp8.cloud.aliyuncs.com</code>	<code>ntp3.aliyun.com</code>
<code>ntp3.cloud.aliyuncs.com</code>	<code>ntp9.cloud.aliyuncs.com</code>	<code>ntp4.aliyun.com</code>
<code>ntp4.cloud.aliyuncs.com</code>	<code>ntp10.cloud.aliyuncs.com</code>	<code>ntp5.aliyun.com</code>
<code>ntp5.cloud.aliyuncs.com</code>	<code>ntp11.cloud.aliyuncs.com</code>	<code>ntp6.aliyun.com</code>

Classic network intranet	VPC intranet	Internet
ntp6.cloud.aliyuncs.com	ntp12.cloud.aliyuncs.com	ntp7.aliyun.com

Other public services of Alibaba Cloud ECS

See the following list for some public services provided by Alibaba Cloud ECS.

Public service	Description
Public DNS: 223.5.5.5 / 223.6.6.6	Domain name: http://www.alidns.com
Open source images: http://mirrors.aliyun.com	Update frequency: The image files are updated at everyday 02:00–04:00 (UTC+8:00), including a lot of Linux distributions and open source applications.

4.3 Configure language settings for multiple instances

This tutorial takes German as an example. The German language package is downloaded from Windows Update. A custom image is then created that uses the German language and German keyboard settings. You can then use the custom image to create as many instances as required.

Context

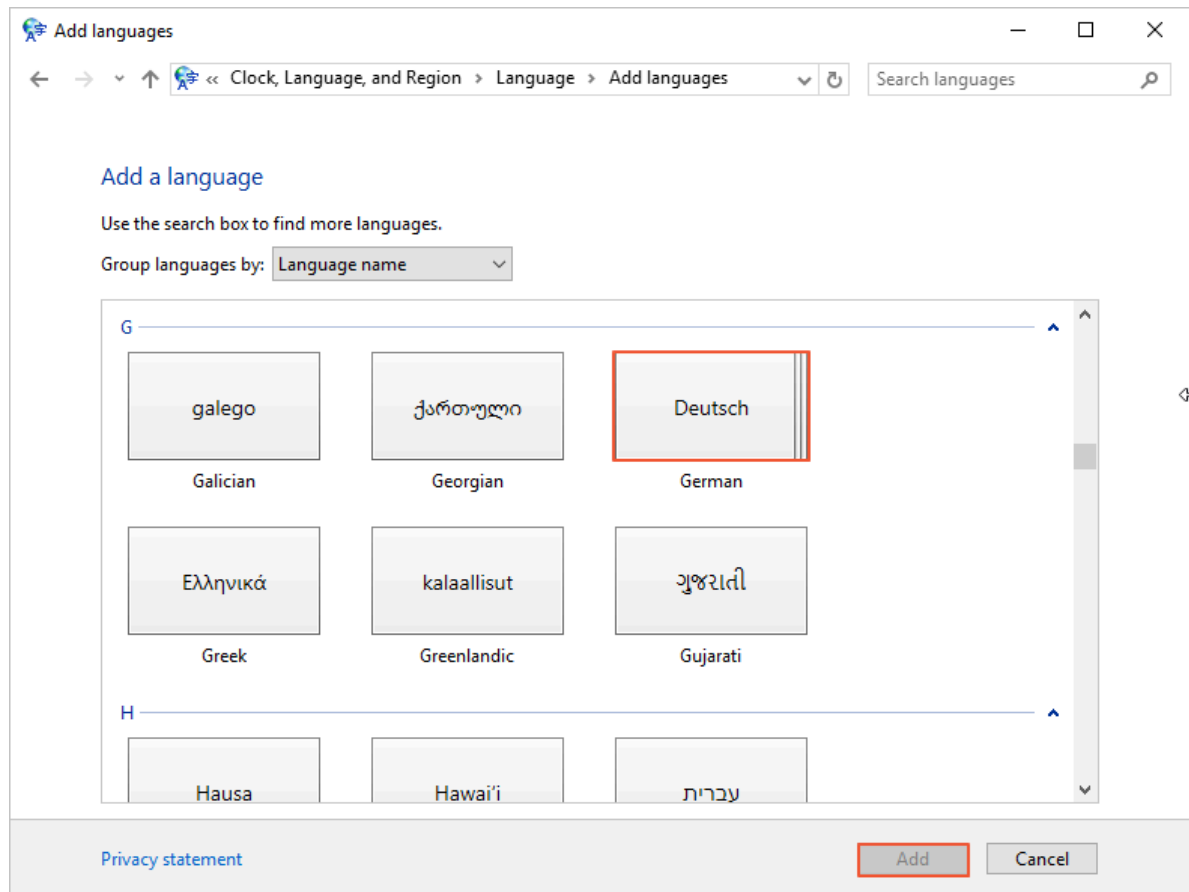
Currently, Alibaba Cloud ECS provides only Chinese and English editions of Windows Server images. If you want to use other language editions, such as Arabic, German, or Russian, you can follow this tutorial to set up and deploy your ECS instances.

Procedure

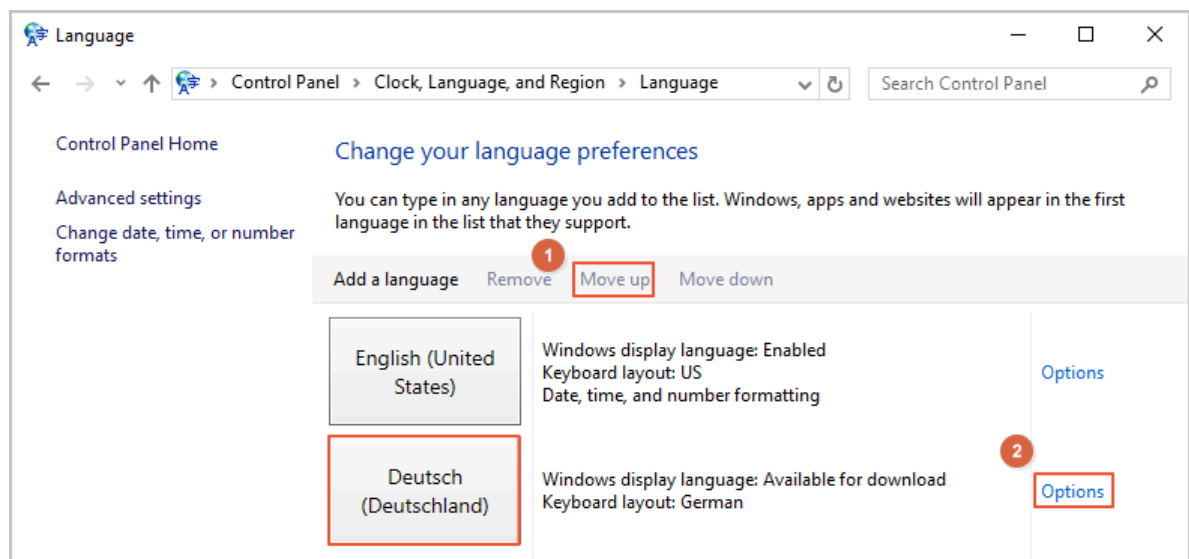
1. [Connect to the Windows instance.](#)
2. Open the PowerShell module.
3. Run the following commands to disable WSUS temporarily.

```
Set-ItemProperty -Path 'HKLM:\SOFTWARE\Policies\Microsoft\Windows\WindowsUpdate\AU' -Name UseWUService -Value 0
Restart-Service -Name wuauclt
```

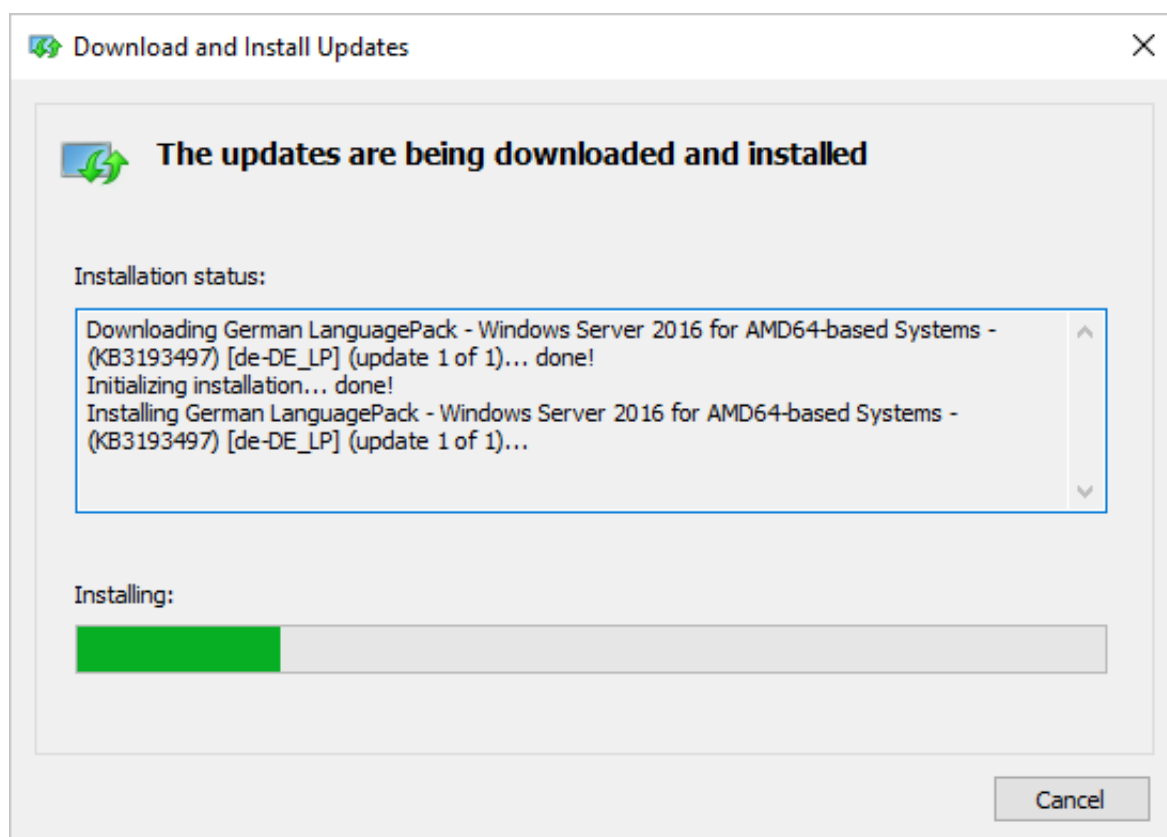
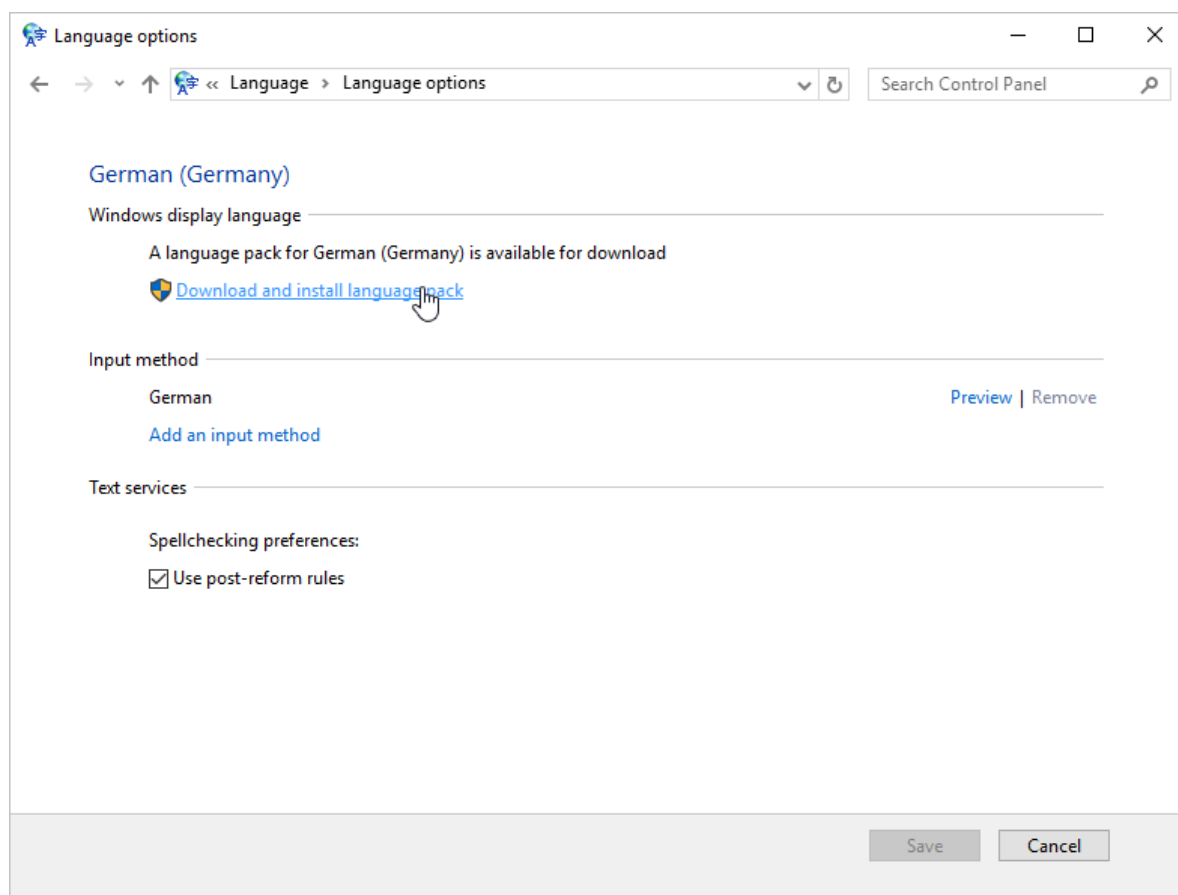
4. Find the Control Panel, click Clock, Language, and Region > Language > Add a language.
5. In the Add languages dialog box, select a language, for example, Deutsch (German) > Deutsch (Deutschland), and click Add.



6. Select the language, such as Deutsch (Deutschland), and click Move up to change the language priority.
7. Click Options next to the selected language to check online for language updates.



8. Wait for about 3 minutes while the instance checks for updates. Once the update is available for download, click Download and install language pack and wait until the installation is complete.



9. *Restart your instance*, and the display language is changed on next login.

10. [Connect to the Windows instance](#) again. The display language is now Deutsch (German).

11. Open the PowerShell ISE module and run the following commands to turn WSUS back on.

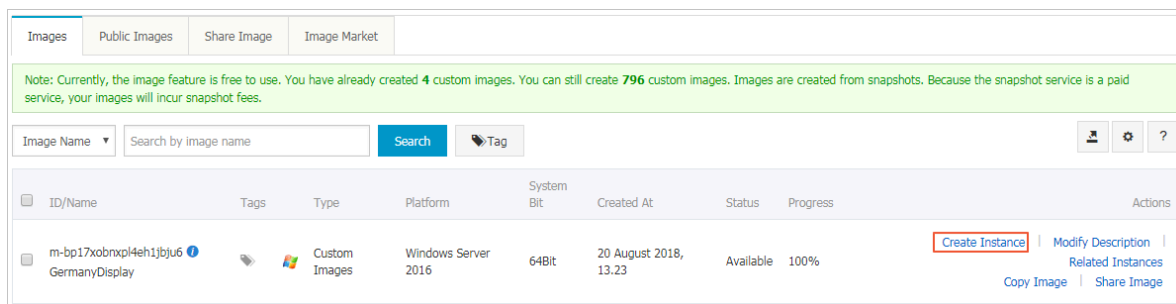
```
Set-ItemProperty -Path 'HKLM:\SOFTWARE\Policies\Microsoft\Windows\WindowsUpdate\AU' -Name UseWUService -Value 1
Restart-Service -Name wuauserv
```

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

What's next

Create multiple instances with the same language settings

1. Log on to the [ECS console](#).
2. and [create a custom image](#) by using the Windows instance with the new display language.
3. [Create a specified number of instances from the custom image](#).



4.4 Time setting: Synchronize NTP servers and change time zone for Linux instances

The current default time zone for Alibaba Cloud ECS instances across all regions is CST (China Standard Time). In addition, the NTP (Network Time Protocol) service guarantees that your instances are synchronized with the standard time. Follow these steps in this topic to change the time zone for your ECS instances and configure your NTP service.

Context

Synchronizing time and the time zone is crucial for Elastic Compute Service (ECS) instances, for example, an inaccurate time may have a significant impact on business when updating your database. To avoid both business disruptions running on your instances and networking request errors, you must configure one or more instances

in the same time zone, such as `Asia/Shanghai` or `America/Los Angeles`. Take CentOS 6.5 as an example to demonstrate how to change the time zone by modifying configuration file.

**Note:**

After you change the time zone for an instance, always run `hwclock -w` to update the real-time clock (RTC) of the instance.

Procedure

1. [Connect](#) to the Linux instance.

**Note:**

Only a root user can open and edit time zone configuration files, so we use the `sudo` command here.

2. Run `sudo rm /etc/localtime` to delete the local time in the instance.
3. Run `sudo vi /etc/sysconfig/clock` to edit the configuration file `/etc/sysconfig/clock`.
4. Enter `i` to add the time zone and city. For example, add `Zone=Asia/Shanghai`. Press `Esc` to exit the edit and enter `:wq` to save and exit.

Optional. Run `ls /usr/share/zoneinfo` to query the list of available time zones. For example, `Shanghai` is one of them.

5. Run `sudo ln -sf /usr/share/zoneinfo/XXXX/XXXXXXX /etc/localtime` to update the time zone change, for example, run `sudo ln -sf /usr/share/zoneinfo/Asia/Shanghai /etc/localtime`.
6. Run `hwclock -w` to update the RTC.
7. Run `sudo reboot` to restart the instance.
8. Run `date -R` to check whether the new time zone is effective or not. If not, repeat the preceding steps.

What's next

The Linux instance offers the `ntpdate` and the `ntpd` two approaches of synchronizing the NTP service. The `ntpdate` can be used to force an immediate update and the `ntpd` offers a systematic approach. The `ntpdate` service can be used for new instances,

whereas `ntpd` is recommended for instances that run your business. Both standard and custom NTP service configurations are provided in this section. For more information about the NTP service, see [internal and public NTP server](#).

Prerequisites

The communication port of the NTP service is UDP 123. Before configuring the service, make sure that you enabled the UDP port 123. You can use `netstat -nupl` in the Linux instance to make sure whether the UDP port 123 is enabled or not. For more information, see [add a security group rule](#).

Set up standard NTP service

1. [Connect](#) to the Linux instance.
2. Run `sudo service ntpd start` to enable the NTP service.
3. Run `chkconfig ntpd on` to enable the NTP service.
4. Run `ntpstat` to check whether the NTP service is enabled or not.
5. Optional. Run `ntpq -p` to view a list of NTP service peers. Run `sudo chkconfig --list ntpd` to view the NTP service running level.

Set up custom NTP service

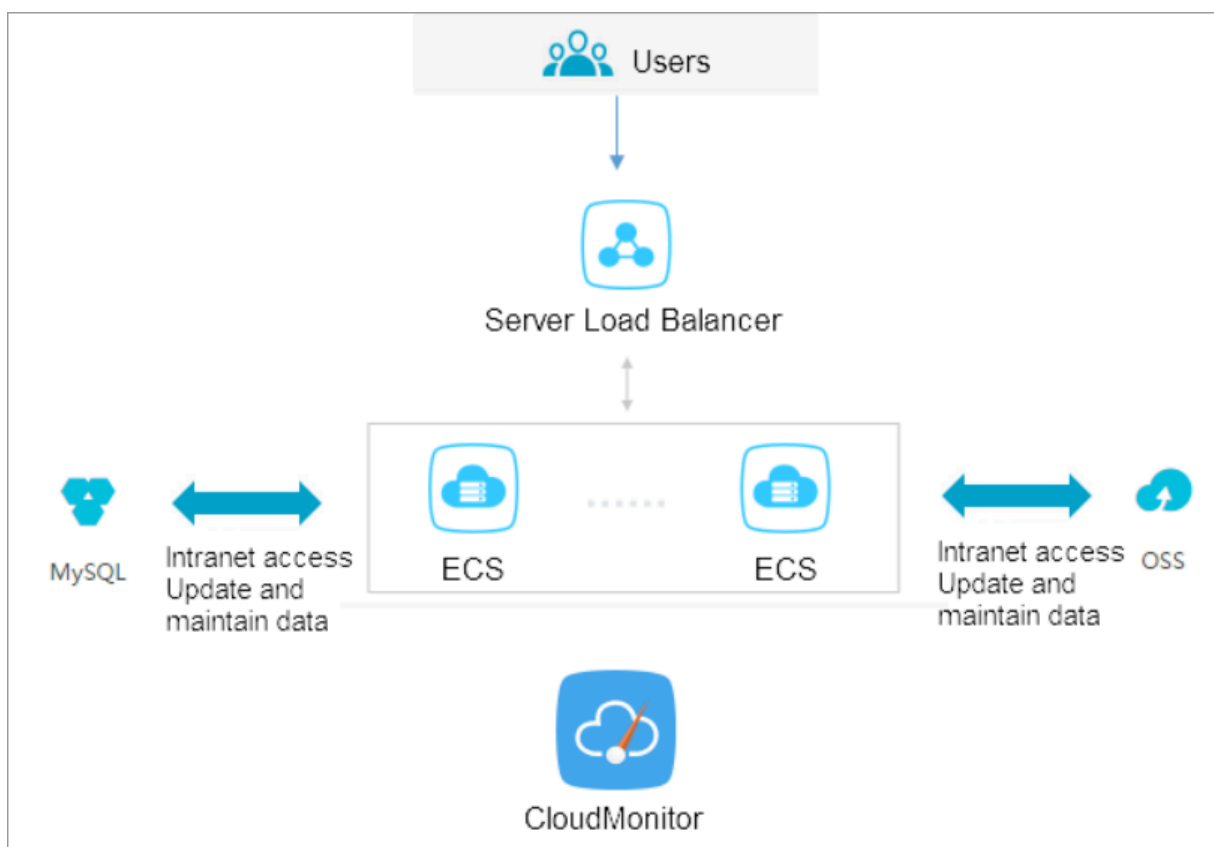
1. [Connect](#) to the Linux instance.
2. Run `sudo vi /etc/ntp.conf` to edit the NTP service configuration files.
3. After finding the information about `ntp server XXXX iburst`, enter `i` and start editing the file. NTP servers that are not currently needed can be hidden by adding a pound (#) at the beginning of the lines.
4. Add a new line of NTP server information in the format of `server XXXX iburst`, and the XXXX is the custom NTP endpoint. For more information, see [Internet and intranet NTP servers](#). After editing, press Esc and enter `:wq` to save and exit.
5. Run `sudo service ntpd start` to enable the customized NTP service.
6. Run `chkconfig ntpd on` to enable the NTP service.
7. Run `ntpstat` to check whether the NTP service is enabled or not.

5 Monitor

5.1 Use CloudMonitor to monitor ECS instances

Many businesses are moving to cloud computing because it is cost-effective, and saves customers of heavy lifting. This can be greatly attributed to the leverage of monitoring. Monitoring service provides real-time operation data for you to identify risks in advance, avoid potential loss, and troubleshoot as quickly as possible.

This article takes a website for example (the website architecture is shown as follows) to illustrate how to configure CloudMonitor. The example website uses Alibaba Cloud services such as ECS, RDS, OSS, and Server Load Balancer.



Prerequisites

Before you begin, you must complete the following operations:

- Make sure that your ECS monitoring agents are functional to collect metric data. Otherwise, you must install the agent manually. For more information, see [How to install CloudMonitor agent](#).

- [Add alarm contacts and contact groups](#). We recommend that you add at least two contacts to make sure real-time responses to monitoring alarms. For more information about metrics, see [Cloud service overview and alarm overview](#).
- With CloudMonitor Dashboard, you can gain system-wide visibility into resource utilization and operational health. You can select a metrics dimension. You can choose per-instance metrics dimension if you only have several instances.

Otherwise, you can choose ECS groups dimension or user dimension, and choose the average value.

Setting alarm threshold

We recommend that you set the alarm threshold according to your business status. A much lower threshold may trigger alarm too often and render monitoring meaningless, while a much higher threshold may leave you with no time to respond to a major event.

Set alarm rules

Take CPU utilization as an example. We have to reserve some processing capacity to guarantee the normal function, so you can set the threshold to 70% and to trigger an alarm when the threshold is exceeded by three times in a row, as shown in the following figure.

If you have to set alarm rules for other metrics, click Add Alarm Rule.

2 Set Alarm Rules

Alarm Type : **Threshold Value Alarm** Event Alarm

Alarm Rule : CPU Alarm

Rule Describe : (ECS) CPU Usage 5mins Average >= 70 %

[+Add Alarm Rule](#)

Mute for : 24h

Triggered when threshold is exceeded for : 3

Effective Period : 00:00 To: 23:59

Set process monitoring

For Web applications, you can [add monitoring for process](#) . For more information, see [Process monitoring](#).

Set site monitoring

Site monitoring is at the network access layer to test the availability.

Set RDS monitoring

We recommend that you set the RDS CPU utilization alarm threshold to 70% and to trigger an alarm when the threshold is exceeded by three times in a row. You can set the disk utilization , IOPS utilization, total connections and other [metrics](#) as needed.

Set Server Load Balancer monitoring

Before you begin, make sure that you have enabled health check for your Server Load Balancer instance.

You can use Custom monitoring metrics if the metrics you need are not covered.

5.2 Automatically manage instances

ECS instances maintenance aims to keep ECS instances in the best state and guarantee the troubleshooting efficiency. However, manual maintenance involves a huge amount of time and effort. To address this issue, you can use cloud assistant for automation and batch processing of daily maintenance tasks. This topic illustrates how to automatically maintain ECS instances by invoking cloud assistant commands on ECS instances.

Context

Cloud assistant supports the following three command types.

Command type	Parameter	Description
Shell script	RunShellScript	A shell script that is running on running Linux instances.
PowerShell script	RunPowerShellScript	A PowerShell script that is running on running Windows instances.
Bat script	RunBatScript	A Bat script that is running on running Windows instances.

Prerequisites

- You must make sure that the network type of the target ECS instances is [VPC](#).
- The target ECS instances must be in the Running (Running) status.
- The target ECS instances must have the Cloud Assistant client installed in advance. For more information, see [Cloud Assistant Client](#).
- To perform a PowerShell command, you must make sure that the target Windows instances has the PowerShell feature configured.
- You can get the latest version of Alibaba Cloud CLI from [GitHub](#).
- You must make sure that you have installed [Alibaba Cloud CLI \(Command-Line Interface\)](#).
- You must [have your SDK upgraded](#).

The following example illustrates how to use APIs in Alibaba Cloud CLI to use Cloud Assistant. For example, we want to run the `echo 123` command on Linux instances.

Procedure

1. In the CMD, PowerShell, or Shell of a local computer, run `aliyuncli ecs CreateCommand --CommandContent ZWNobyAxMjM= --Type RunShellScript --Name test --Description test` to [create a shell script](#) (CreateCommand). The Command ID information is returned after successful creation.



Note:

- The `ZWNobyAxMjM=` in `CommandContent` is the Base64 code of the `echo 123` command. For more information about Base64 encoding or decoding, see [Wikipedia - Base64](#).
- If the operating system of the target ECS instances are Windows, change `type` to `RunBatScript` or `RunPowershellScript`.
- After the script is created, `CommandId` is returned.

```
C:\Windows\System32>aliyuncli ecs CreateCommand --CommandContent ZWNobyAxMjM= --
Type RunShellScript --Name test --Description test
=====
|                               CreateCommand                               |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|                               |                               |                               |
| CommandId                     | RequestId                     |                               |
|                               |                               |                               |
+-----+-----+-----+-----+-----+-----+-----+
| c-f0902c0972984e31aaf2129fd48a9c6d | 34E84CD7-723B-47D6-8568-1FCC8604ED4E |
+-----+-----+-----+-----+-----+-----+-----+
C:\Windows\System32>
```

2. Run `aliyuncli ecs InvokeCommand --InstanceId.1 your-vm-instance-id1 --InstanceId.2 your-vm-instance-id2 --CommandId your-command-id --Timed false` to [run the command](#) (InvokeCommand).

**Note:**

- The `InstanceIds` indicates your ECS instances IDs. Up to 100 ECS instances are supported each time.
- The `Timed` indicates whether the task is a periodical one or not. `--Timed True` indicates that the task is a periodical one, while `--Timed False` indicates the opposite.
- When your task is a periodical one and the `Timed` parameter value is `True`, you must specify the interval value in the `Frequency` parameter. For example, `0 */20 * * * *` indicates that the interval value is 20 minutes. For more information, see [expressions](#).
- A shared `InvokeId` is returned for all target ECS instances. You can use the `InvokeId` to check the invocation status of the command.

3. Optional. Run `aliyuncli ecs DescribeInvocations --InstanceId your-vm-instance-id --InvokeId your-invoke-id` to [query the invocation status](#) (DescribeInvocations). Specifically, the `InvokeId` is the invocation ID returned in [step 2](#) during command invocation on the ECS instances.

When the returned `InvokeStatus` value is `Finished`, it indicates that the command process is complete, but not necessarily as effective as expected. You must check the `Output` parameter in [DescribeInvocationResults](#) to get the specific invocation result.

4. (Optional). Run `aliyuncli ecs DescribeInvocationResults --InstanceId your-vm-instance-id --InvokeId your-invoke-id` to [check the results of the invocation](#) (DescribeInvocationResults). Specifically, the `InvokeId` is the invocation ID returned in [step 2](#) during command invocation on the ECS instances.

Result

When [creating a command](#) (CreateCommand), you can set the following request parameters for the command.

Command Property	Parameter	Description
Execution directory	WorkingDir	<p>Specifies the path in an ECS instance where the command is performed. Default value:</p> <ul style="list-style-type: none"> For Linux instances: /root. For Windows instances: In the path where the cloud assistant client process is located, such as C:\ProgramData\aliyun\assist\\$(version).
Timeout period	TimeOut	<p>Modifies the invocation timeout value of a command on ECS instances. The unit is seconds.</p> <p>When your command fails for some reason, the invocation may time out, and the cloud assistant client forces to terminate the command process afterwards. The parameter value must be greater than or equal to 60. If the value is smaller than 60, the timeout value is 60 seconds by default.</p> <p>Default value: 3600</p> <ul style="list-style-type: none"> One-time invocation: <ul style="list-style-type: none"> After invocation timeout, the command invocation status (DescribeInvocationResults) for the specified ECS instances becomes Failed. Periodical invocation: <ul style="list-style-type: none"> The timeout value of periodical invocation is effective for every invocation record. After one invocation operation timed out, the status for the invocation record (DescribeInvocationResults) becomes Failed. The timeout status of last invocation does not affect the next invocation.

Sample of Python SDK to use cloud assistant

You can also use the cloud assistant by using the [Alibaba Cloud SDK](#). For more information about how to configure Alibaba Cloud SDK, see [for Alibaba Cloud users](#). The following is the Python SDK code to use cloud assistant.

```
# coding=utf-8
# if the python sdk is not install using 'sudo pip install aliyun-python-sdk-ecs'
# if the python sdk is install using 'sudo pip install --upgrade aliyun-python-sdk-ecs'
# make sure the sdk version is 2.1.2, you can use command 'pip show aliyun-python-sdk-ecs' to check
```

```

import json
import logging
import os
import time
import datetime
import base64
from aliyunsdkcore import client

from aliyunsdkecs.request.v20140526. CreateCommandRequest import
CreateCommandRequest
from aliyunsdkecs.request.v20140526. InvokeCommandRequest import
InvokeCommandRequest
from aliyunsdkecs.request.v20140526. DescribeInvocationResultsRequest
import DescribeInvocationResultsRequest

# configuration the log output formatter, if you want to save the
output to file,
# append ",filename='ecs_invoke.log'" after datefmt.
logging.basicConfig(level=logging.INFO,
                    format='%(asctime)s %(filename)s[line:%(lineno)d]
                    %(levelname)s %(message)s',
                    datefmt='%a, %d %b %Y %H:%M:%S',filename='
aliyun_assist_openapi_test.log', filemode='w')
#access_key = 'Your Access Key Id'
#access_key_secret = 'Your Access Key Secret'
#region_name = 'cn-shanghai'
#zone_id = 'cn-shanghai-b'

access_key = 'LTAIXXXXXXXXXXXXXX'
access_key_secret = '4dZXXXXXXXXXXXXXXXXXXXXXXXXXXXXX'
region_name = 'cn-hangzhou'
zone_id = 'cn-hangzhou-f'

clt = client.AcsClient(access_key, access_key_secret, region_name)

def create_command(command_content, type, name, description):
    request = CreateCommandRequest()
    request.set_CommandContent(command_content)
    request.set_Type(type)
    request.set_Name(name)
    request.set_Description(description)
    response = _send_request(request)
    if response is None:
        return None
    command_id = response.get('CommandId')
    return command_id;

def invoke_command(instance_id, command_id, timed, cronat):
    request = InvokeCommandRequest()
    request.set_Timed(timed)
    InstanceIds = [instance_id]
    request.set_InstanceIds(InstanceIds)
    request.set_CommandId(command_id)
    request.set_Frequency(cronat)
    response = _send_request(request)
    invoke_id = response.get('InvokeId')
    return invoke_id;

def get_task_output_by_id(instance_id, invoke_id):
    logging.info("Check instance %s invoke_id is %s", instance_id,
    invoke_id)
    request = DescribeInvocationResultsRequest()
    request.set_InstanceId(instance_id)
    request.set_InvokeId(invoke_id)

```

```

        response = _send_request(request)
        invoke_detail = None
        output = None
        if response is not None:
            result_list = response.get('Invocation').get('Invocation
Results').get('InvocationResult')
            for item in result_list:
                invoke_detail = item
                output = base64.b64decode(item.get('Output'))
                break;
            return output;

def execute_command(instance_id):
    command_str = 'yum check-update'
    command_id = create_command(base64.b64encode(command_str), '
RunShellScript', 'test', 'test')
    if(command_id is None):
        logging.info('create command failed')
        return

    invoke_id = invoke_command(instance_id, command_id, 'false', '')
    if(invoke_id is None):
        logging.info('invoke command failed')
        return

    time.sleep(15)

    output = get_task_output_by_id(instance_id, invoke_id)
    if(output is None):
        logging.info('get result failed')
        return

    logging.info("output: %s is \n", output)

# send open api request
def _send_request(request):
    request.set_accept_format('json')
    try:
        response_str = clt.do_action(request)
        logging.info(response_str)
        response_detail = json.loads(response_str)
        return response_detail
    except Exception as e:
        logging.error(e)

if __name__ == '__main__':
    execute_command('i-bp17zhpbXXXXXXXXXXXXX')

```

References

The preceding examples demonstrate how to auto manage ECS instances maintenance by using Alibaba Cloud CLI and cloud assistant APIs [CreateCommand](#), [InvokeCommand](#), [DescribeInvocations](#), and [DescribeInvocationResults](#). You can also use other APIs of the cloud assistant:

- [StopInvocation](#): Stops a scheduled command process.
- [ModifyCommand](#): Modifies the content of a command.
- [DescribeCommands](#): Queries the available commands.

- *DeleteCommand*: Deletes a command.

6 User-defined data

6.1 User-defined yum sources, NTP services and DNS services

User-defined scripts are a type of script provided by Alibaba Cloud for users to customize the startup behaviors of ECS instances. For more information, see [User-defined data](#).

This example uses a Linux instance to demonstrate how to use a user-defined script to configure your own yum repository, NTP service, and DNS service when creating a Linux instance. User-defined scripts also enable you to configure NTP service and DNS service for a Windows instance.

Scenarios

When a Linux instance is started, Alibaba Cloud automatically configures a pre-defined yum repository, NTP service, and DNS service for the instance. However, if you want to have your own yum repository, NTP service, and DNS service, use user-defined scripts to implement them.

- If you are using a custom yum repository, Alibaba Cloud does not provide support for it.
- If you are using a custom NTP service, Alibaba Cloud does not provide time service.

Procedure

To customize your yum repository, NTP service, and DNS service for a Linux instance when creating it, follow these steps:

1. Log on to the [ECS console](#) and create an instance. Configure the instance as follows:
 - Network Type: Select VPC.
 - Instance Type: Select an I/O-optimized instance.
 - Operating System: Select CentOS 7.2 in Public Image tab.
2. Enter the following script in the User Data box on the instance creation page.

```
#!/bin/sh
# Modify DNS
echo "nameserver 8.8.8.8" | tee /etc/resolv.conf
# Modify yum repo and update
rm -rf /etc/yum.repos.d/*
```



```
touch myrepo.repo
echo "[base]" | tee /etc/yum.repos.d/myrepo.repo
echo "name=myrepo" | tee -a /etc/yum.repos.d/myrepo.repo
echo "baseurl=http://mirror.centos.org/centos" | tee -a /etc/yum.
repos.d/myrepo.repo
echo "gpgcheck=0" | tee -a /etc/yum.repos.d/myrepo.repo
echo "enabled=1" | tee -a /etc/yum.repos.d/myrepo.repo
yum update -y
# Modify NTP Server
echo "server ntp1.aliyun.com" | tee /etc/ntp.conf
systemctl restart ntpd.service
```

**Note:**

- The first line must be `#!/bin/sh`, with no leading space.
- Do not add unnecessary spaces or carriage return characters in the full text.
- You can customize URLs of your own DNS server, NTP Server, and yum repository based on the instance situations.
- The preceding content applies to CentOS 7.2. If you are using other images, modify the scripts as needed.
- You can also define the yum repository in the scripts of the `cloud config` type, but it is not recommended because it is not flexible enough to get adapted to Alibaba Cloud that may pre-configure some yum repository. Scripts of `script` type is recommended for changing the yum repository.

3. Complete the security settings as needed.
4. After you complete the configuration, click Buy Now and activate the instance following the instructions on the page.

After the instance is created, you can connect to the instance to view the implementation details, as shown in the following figure.

```

[root@iZwz99v9qbmmk2dswgnzg8Z yum.repos.d]# cat /etc/resolv.conf
nameserver 8.8.8.8
[root@iZwz99v9qbmmk2dswgnzg8Z yum.repos.d]# ping www.baidu.com
PING www.a.shifen.com (103.235.46.39) 56(84) bytes of data:
64 bytes from 103.235.46.39: icmp_seq=1 ttl=48 time=73.3 ms
^C64 bytes from 103.235.46.39: icmp_seq=2 ttl=48 time=74.8 ms

--- www.a.shifen.com ping statistics ---
2 packets transmitted, 2 received, 0% packet loss, time 100lms
rtt min/avg/max/mdev = 73.393/74.113/74.833/0.720 ms
[root@iZwz99v9qbmmk2dswgnzg8Z yum.repos.d]# cat /etc/ntp.conf
server ntp1.aliyun.com
[root@iZwz99v9qbmmk2dswgnzg8Z yum.repos.d]# systemctl status ntpd.service
● ntpd.service - Network Time Service
   Loaded: loaded (/usr/lib/systemd/system/ntpd.service; enabled; vendor preset: disabled)
   Active: active (running) since Mon 2017-03-13 11:08:11 CST; 1min 58s ago
     Process: 6235 ExecStart=/usr/sbin/ntpd -u ntp:ntp $OPTIONS (code=exited, status=0/SUCCESS)
    Main PID: 6237 (ntpd)
      CGroup: /system.slice/ntpd.service
              └─6237 /usr/sbin/ntpd -u ntp:ntp -g

Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: 0.0.0.0 c01d 0d kern kernel time sync enabled
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: ntp_io: estimated max descriptors: 1024, initial socket boundary: 16
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: Listen and drop on 0 v4wildcard 0.0.0.0 UDP 123
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: Listen and drop on 1 v6wildcard :: UDP 123
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: Listen normally on 2 lo 127.0.0.1 UDP 123
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: Listen normally on 3 eth0 172.18.48.114 UDP 123
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: Listening on routing socket on fd #20 for interface updates
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: 0.0.0.0 c016 06 restart
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: 0.0.0.0 c012 02 freq_set kernel 0.000 PPM
Mar 13 11:08:11 iZwz99v9qbmmk2dswgnzg8Z ntpd[6237]: 0.0.0.0 c011 01 freq_not_set
[root@iZwz99v9qbmmk2dswgnzg8Z yum.repos.d]# cat /etc/yum.repos.d/myrepo.repo
[base]
name=myrepo
baseurl=http://mirror.centos.org/centos
gpgcheck=0
enabled=1
[root@iZwz99v9qbmmk2dswgnzg8Z yum.repos.d]#

```

The preceding figure shows that you have successfully customized the DNS service, the NTP service, and the yum repository.

6.2 Create a new account with the root user privilege

User-defined scripts are a type of script provided by Alibaba Cloud to enable users to customize the startup behavior of ECS instances. For details, see [User-defined data](#).

This example uses a Linux instance to demonstrate how to use a user-defined script to create a new account, with the root user privilege, when creating a Linux instance. User-defined scripts can also be used to create a new account with the administrator privilege for a Windows instance.

Scenarios

Use user-defined scripts of instances if you want to achieve the following results when creating a Linux ECS instance:

- Disable the default root account that comes with a Linux ECS instance. You can use the script to customize how to disable the root user and how many root user privileges are disabled.
- Create a new account with the root user privilege and customize the account name.
- Use only SSH key pairs, but not user passwords, for remote logon to manage the instance by using the new account with the root user privilege.

- If this new account is required to perform operations that can only be done by a user with root user privilege, the `sudo` command can be used without a password for privilege escalation.

Procedure

To create a new account with the root user privilege, follow these steps:

1. [Create a Linux instance](#). Configure the instance as follows:

- Network Type: Select VPC.
- Instance Type: Select an I/O-optimized instance.
- Operating System: Select CentOS 7.2 in Public Image tab.

2. Enter the following script in the User Data box on the instance creation page:

```
#!/bin/sh
useradd test
echo "test    ALL=(ALL)        NOPASSWD:ALL" | tee -a /etc/sudoers
mkdir /home/test/.ssh
touch /home/test/.ssh/authorized_keys
echo "ssh-rsa AAAAB3NzaC1yc2EAAAABJQAAAQEAhGqEh/rGbIMCGItF
VtYpsXPQrCaunGJKZVIWtINrGZwusLc290qDZ93KCeb8o6X1Iby1Wm+psZY8THE+/
BsXq0M0HzfkQZD2vXuhRb4x1lZ98JHskX+0jnbjqYGY+Brgai9BvKDXTTsyJtCYU
nEKxvcK+d1ZwxbNuk2QZ0ryHESDbSaczlNFgFQEDxhCrvko+zWLjTVnomVUDhdMP2g6f
Z0tgFVwkJFV0bE7oob3N0Vcrx2TyhfcAjA4M2/Ry7U2MFADDC+EVkpoVDm0SOT/
hYJgaVM1xMDlSeE7kzX7yZbJLR1XAWV1xzZkNclY5w1kPnW8qMYuSwphXzt4gsF0w==
rsa-key-20170217" | tee -a /home/test/.ssh/authorized_keys
```



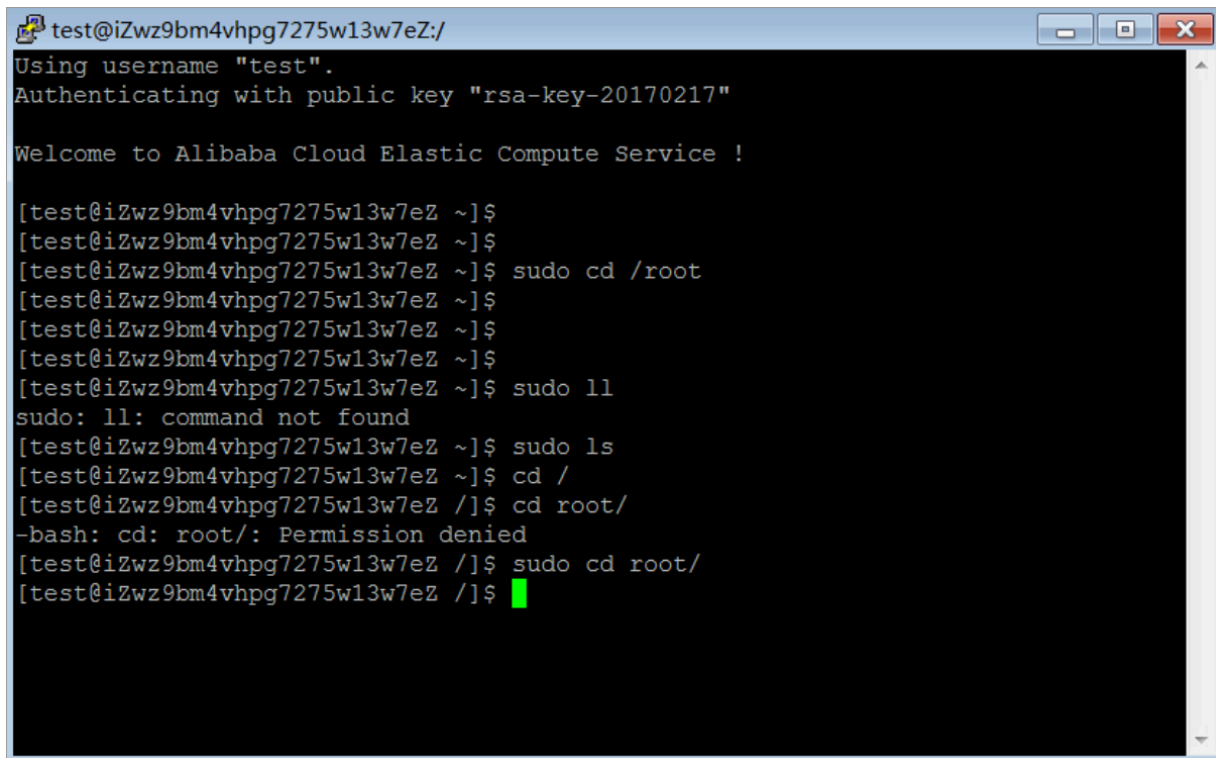
Note:

- The first line must be `#!/bin/sh`, with no leading space.
- Do not enter unnecessary spaces or carriage return characters in the text.
- The last line is your public key. You can define it.
- You can add other configuration in the script, as you need.
- The example script only applies to CentOS 7.2. If you are using other images, customize the script according to the operating system types.

3. Select post-creation settings in security settings.

4. After you finish the configuration, click **Buy Now** and activate the instance by following instructions on the page.

After the instance is created, you can use the new test user to connect to the instance using an SSH private key. You can also escalate the permission level using the `sudo` command and run operations that require the root user privilege, as shown in the following figure.



```
test@iZwz9bm4vhpg7275w13w7eZ:/
Using username "test".
Authenticating with public key "rsa-key-20170217"

Welcome to Alibaba Cloud Elastic Compute Service !

[test@iZwz9bm4vhpg7275w13w7eZ ~]$
[test@iZwz9bm4vhpg7275w13w7eZ ~]$
[test@iZwz9bm4vhpg7275w13w7eZ ~]$ sudo cd /root
[test@iZwz9bm4vhpg7275w13w7eZ ~]$
[test@iZwz9bm4vhpg7275w13w7eZ ~]$
[test@iZwz9bm4vhpg7275w13w7eZ ~]$
[test@iZwz9bm4vhpg7275w13w7eZ ~]$ sudo ll
sudo: ll: command not found
[test@iZwz9bm4vhpg7275w13w7eZ ~]$ sudo ls
[test@iZwz9bm4vhpg7275w13w7eZ ~]$ cd /
[test@iZwz9bm4vhpg7275w13w7eZ /]$ cd root/
-bash: cd: root/: Permission denied
[test@iZwz9bm4vhpg7275w13w7eZ /]$ sudo cd root/
[test@iZwz9bm4vhpg7275w13w7eZ /]$
```

7 GPU instances

7.1 Deploy an NGC on gn5 instances

As a deep learning ecosystem from NVIDIA, NVIDIA GPU CLOUD (NGC) allows developers to access the deep learning software stack free of charge and is fit for creating a deep learning development environment.

At present, NGC has been fully deployed in the gn5 instances. Moreover, the image market also provides NGC container images optimized for NVIDIA Pascal GPU . By deploying NGC container images from the image market, developers can build an NGC container environment conveniently, and access optimized deep learning frameworks instantly, thus reducing the product development and business deployment time considerably. Other benefits include pre-installation of the development environment, support for optimized algorithm frameworks, and continuous updates.

The [NGC website](#) provides images of different versions of the current mainstream deep learning frameworks (such as Caffe, Caffe2, CNTK, MxNet, TensorFlow, Theano, and Torch). You can select the desired image to build the environment. By taking the TensorFlow deep learning framework for example, this article describes how to build an NGC environment on gn5 instances.

Before building a TensorFlow environment, you must do the following:

- Sign up with Alibaba Cloud and finish real-name registration.
- Log on to the [NGC website](#) and create your NGC account.
- Log on to the [NGC website](#), get the NGC API Key and save it locally. The NGC API Key will be verified when you log on to the NGC container environment.

Procedure

1. Create a gn5 instance by referring to [create an ECS instance](#). Pay attention to the following configurations:
 - Region: Only China North 1, China North 2, China North 5, China East 1, China East 2, China South 1, Hong Kong, Asia Pacific SE 1 (Singapore), Asia Pacific SE 2 (Sydney), US West 1 (Silicon Valley), US East 1 (Virginia), and Germany 1 (Frankfurt) are available.

- **Instance:** Select a gn5 instance type.
- **Image:** Select Marketplace Image. In the displayed dialog box, search for NVIDIA GPU Cloud VM Image, and then click Continue.
- **Network Billing Method:** Select Assign Public IP.

**Note:**

If you do not assign a public IP address here, you can bind an EIP address after the instance is created successfully.

- **Security Group:** Select a security group. Access to TCP port 22 must be allowed in the security group. If your instance needs to support HTTPS or *DIGITS 6*, access to TCP port 443 (for HTTPS) or TCP port 5000 (for DIGITS 6) must be allowed.

After the ECS instance is created successfully, *log on to the ECS console* and note down the public IP address of the instance.

2. **Connect to the ECS instance:** Based on the logon credentials selected during instance creation, you can *connect to an ECS instance by using a password* or *connect to an ECS instance by using an SSH key pair*.
3. Enter the NGC API Key obtained from the NGC website, and then press the Enter key to log on to the NGC container environment.

```

? MobaXterm 8.4 ?
(SSH client, X-server and networking tools)

> SSH session to [redacted]
? SSH compression : ✓
? SSH-browser      : ✓
? X11-forwarding   : ✓ (remote display is forwarded through SSH)
? DISPLAY          : ✓ (automatically set on remote server)

> For more info, ctrl+click on help or visit our website

Welcome to Ubuntu 16.04.4 LTS (GNU/Linux 4.4.0-116-generic x86_64)

* Documentation:  https://help.ubuntu.com
* Management:    https://landscape.canonical.com
* Support:       https://ubuntu.com/advantage

Welcome to the NVIDIA GPU Cloud Virtual Machine. This environment is provided
to enable you to easily run the Deep Learning containers from the NGC Registry.
All of the documentation for how to use NGC and this VM are found at
http://docs.nvidia.com/deeplearning/ngc

Welcome to Alibaba Cloud Elastic Compute Service !

/usr/bin/xauth:  file /root/.Xauthority does not exist

Please enter your NGC APIkey to login to the NGC Registry:

```

4. Run `nvidia-smi`. You can view the information about the current GPU, including the GPU model, the driver version, and more, as shown below.

```

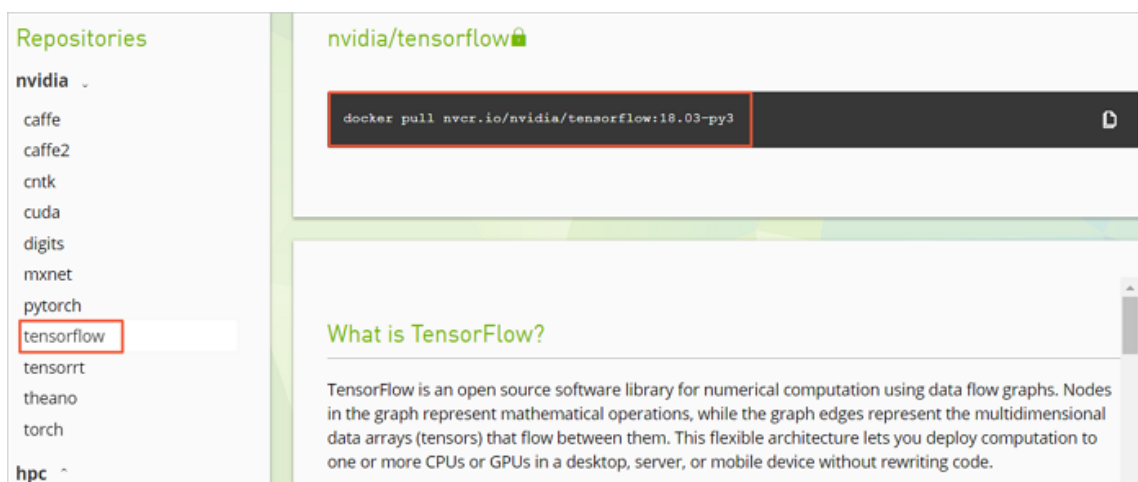
root@--:~# nvidia-smi
Thu Mar 29 20:50:01 2018

+-----+
| NVIDIA-SMI 384.111                Driver Version: 384.111          |
+-----+-----+
| GPU Name           Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
+-----+-----+
| 0   Tesla P100-PCIE...    Off   | 00000000:00:08.0 Off  |             0        |
| N/A   29C    P0       27W / 250W |      0MiB / 16276MiB |           0%      Default |
+-----+-----+

+-----+
| Processes:                                                       GPU Memory |
|  GPU       PID    Type    Process name                     Usage    |
+-----+-----+
| No running processes found                                     |
+-----+

```

5. Follow the steps below to build the TensorFlow environment:
 - a. Log on to the [NGC website](#), go to the TensorFlow image page, and then get the `docker pull` command.



b. Download the TensorFlow image.

```
docker pull nvcr.io/nvidia/tensorflow:18.03-py3
```

c. View the downloaded image.

```
docker image ls
```

d. Run the container to deploy the TensorFlow development environment.

```
nvidia-docker run --rm -it nvcr.io/nvidia/tensorflow:18.03-py3
```

```
root@ :~# nvidia-docker run --rm -it nvcr.io/nvidia/tensorflow:18.03-py3
=====
== TensorFlow ==
=====

NVIDIA Release 18.03 (build 349854)

Container image Copyright (c) 2018, NVIDIA CORPORATION. All rights reserved.
Copyright 2017 The TensorFlow Authors. All rights reserved.

Various files include modifications (c) NVIDIA CORPORATION. All rights reserved.
NVIDIA modifications are covered by the license terms that apply to the underlying project or file.
```

6. Test TensorFlow by using one of the following methods:

- Simple test of TensorFlow.

```
$python
```

```
>>> import tensorflow as tf
>>> hello = tf.constant('Hello, TensorFlow!')
>>> sess = tf.Session()
>>> sess.run(hello)
```

If TensorFlow loads the GPU device correctly, the result is as shown below.


```

root@----- # python
Python 3.5.2 (default, Nov 23 2017, 16:37:01)
[GCC 5.4.0 20160609] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> import tensorflow as tf
>>> hello = tf.constant('Hello, TensorFlow!')
>>> sess = tf.Session()
2018-03-30 03:37:53.682157: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:892] s
be at least one NUMA node, so returning NUMA node zero
2018-03-30 03:37:53.682544: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1030] Foun
name: Tesla P100-PCI-E-16GB major: 6 minor: 0 memoryClockRate(GHz): 1.3285
pciBusID: 0000:00:08:0
totalMemory: 15.89GiB freeMemory: 15.60GiB
2018-03-30 03:37:53.682583: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1120] Crea
16GB, pci bus id: 0000:00:08:0, compute capability: 6.0)
>>> sess.run(hello)
b'Hello, TensorFlow!'
>>>

```

- Download the TensorFlow model and test TensorFlow.

```

git clone https://github.com/tensorflow/models.git
cd models/tutorials/image/alexnet
python alexnet_benchmark.py --batch_size 128 --num_batches 100

```

The running status is as shown below.

```

conv1 [128, 56, 56, 64]
pool1 [128, 27, 27, 64]
conv2 [128, 27, 27, 192]
pool2 [128, 13, 13, 192]
conv3 [128, 13, 13, 384]
conv4 [128, 13, 13, 256]
conv5 [128, 13, 13, 256]
pool5 [128, 6, 6, 256]
2018-03-30 03:40:13.357785: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:892] successful NUMA node read from SysFS
be at least one NUMA node, so returning NUMA node zero
2018-03-30 03:40:13.358297: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1030] Found device 0 with properties:
name: Tesla P100-PCI-E-16GB major: 6 minor: 0 memoryClockRate(GHz): 1.3285
pciBusID: 0000:00:08:0
totalMemory: 15.89GiB freeMemory: 15.60GiB
2018-03-30 03:40:13.358245: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1120] Creating TensorFlow device (/device:GPU:
16GB, pci bus id: 0000:00:08:0, compute capability: 6.0)
2018-03-30 03:40:15.916471: step 0, duration = 0.038
2018-03-30 03:40:16.299169: step 10, duration = 0.038
2018-03-30 03:40:16.682881: step 20, duration = 0.038
2018-03-30 03:40:17.065379: step 30, duration = 0.038
2018-03-30 03:40:17.448118: step 40, duration = 0.038
2018-03-30 03:40:17.830372: step 50, duration = 0.038
2018-03-30 03:40:18.213018: step 60, duration = 0.038
2018-03-30 03:40:18.595734: step 70, duration = 0.038
2018-03-30 03:40:18.978311: step 80, duration = 0.038
2018-03-30 03:40:19.361063: step 90, duration = 0.038
2018-03-30 03:40:19.705396: Forward across 100 steps, 0.038 +/- 0.000 sec / batch
2018-03-30 03:40:21.164735: step 0, duration = 0.090
2018-03-30 03:40:22.062778: step 10, duration = 0.090
2018-03-30 03:40:22.962202: step 20, duration = 0.090
2018-03-30 03:40:23.860856: step 30, duration = 0.090
2018-03-30 03:40:24.758891: step 40, duration = 0.090
2018-03-30 03:40:25.657170: step 50, duration = 0.090
2018-03-30 03:40:26.555194: step 60, duration = 0.090
2018-03-30 03:40:27.452843: step 70, duration = 0.090
2018-03-30 03:40:28.351092: step 80, duration = 0.090
2018-03-30 03:40:29.249606: step 90, duration = 0.090
2018-03-30 03:40:30.058089: Forward-backward across 100 steps, 0.090 +/- 0.000 sec / batch

```

7. Save the changes made to the TensorFlow image. Otherwise, the configuration will be lost the next time you log on.

7.2 Install a GRID driver on a gn5/gn5i/gn6v instance

If your GPU instance (available in the gn5, gn5i and gn6v families) requires OpenGL, you must install the GRID driver on the instance. The NVIDIA GRID license granted to the NVIDIA GPU (such as P100, P4 and V100) of gn5, gn5i and gn6v instances cannot

meet the graphics requirements of OpenGL. However, you can use the trial version of the GRID driver to meet the requirements.

This article explains how to install the GRID driver and deploy a desktop environment on a Linux GPU instance running Ubuntu 16.04 or CentOS 7.3.

Ubuntu 16.04

This section describes how to install the GRID driver on a GPU instance running Ubuntu 16.04 64-bit.

Prerequisites

- You have created a gn5, gn5i or gn6v instance. For more information, see [create a compute optimized instance with GPUs](#). Make sure that the instance can access the Internet.



Note:

We recommend that you use a public image rather than an image from the marketplace that is pre-installed with a NVIDIA driver. Otherwise, you have to disable the Nouveau driver after the instance is created. To disable the Nouveau driver, create a file named `nouveau.conf` in the directory of `/etc/modprobe.d` and add `blacklist nouveau` into the file.

- You have installed a VNC application on your local machine. In this example, VNC Viewer is used.

Install a GRID driver

To install the GRID driver, follow these steps:

1. [Connect to the Linux instance](#).
2. Run the following commands in sequence to upgrade the system and install the KDE.

```
apt-get update
apt-get upgrade
apt-get install kubuntu-desktop
```

3. Run `reboot` to restart the system.
4. [Connect to the Linux instance](#) again, and then run the following commands to download and decompress the NVIDIA GRID driver package.

The NVIDIA GRID driver package contains the drivers for various operating systems. For Linux OS, select NVIDIA-Linux-x86_64-390.57-grid.run.

```
wget https://nvidia-driver.oss-cn-huhehaote.aliyuncs.com/NVIDIA-Linux-x86_64-390.57-grid.run
```

5. Run the following commands in sequence, then follow the prompts to install the NVIDIA GRID driver.

```
chmod 777 NVIDIA-Linux-x86_64-390.57-grid.run
./NVIDIA-Linux-x86_64-390.57-grid.run
```

6. Run `nvidia-smi` to verify the installation.

If the following message appears, the driver is installed successfully.

```
root@ Z:~# nvidia-smi
Wed Jul  4 15:34:55 2018

+-----+
| NVIDIA-SMI 390.57                  Driver Version: 390.57          |
+-----+-----+
| GPU   Name           Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
+-----+-----+
|  0 Tesla P100-PCIE...    On      | 00000000:00:08.0 Off |                    |
| N/A   27C    P0      25W / 250W |  0MiB / 16280MiB |      0%      Default |
+-----+-----+

+-----+
| Processes:                         GPU Memory                       |
|  GPU       PID    Type    Process name                       Usage                          |
+-----+-----+
| No running processes found         |
+-----+
```

7. Add License Server to activate the License:

- a. Switch to the `/etc/nvidia` directory: `cd /etc/nvidia`.
- b. Create a file named `gridd.conf`: `cp gridd.conf.template gridd.conf`.
- c. Add the following lines about License Server to the `gridd.conf` file.

```
ServerAddress=116.62.19.179
ServerPort=7070
FeatureType=2
EnableUI=TRUE
```

8. Run the command to install `x11vnc`.

```
apt-get install x11vnc
```

9. Run `lspci | grep NVIDIA` to check GPU BusID.

In this example, the GPU BusID is `00:07:0`.

10. Configure the X Server environment and restart the system:

- Run `nvidia-xconfig --enable-all-gpus --separate-x-screens`.
- Edit `/etc/X11/xorg.conf`: Add your GPU BusID to the Section "Device". In this example, BusID `"PCI:0:7:0"` is added.

```
Section "Device"
    Identifier      "Device0"
    Driver          "nvidia"
    VendorName      "NVIDIA Corporation"
    BoardName       "Tesla P4"
    BusID           "PCI:0:7:0"
EndSection
```

- Run `reboot` to restart the system.

Verify installation

To verify the installation of the GRID driver, follow these steps:

- Run the following command to install the GLX application.

```
apt-get install mesa-utils
```

- Run `startx` to start X Server.



Note:

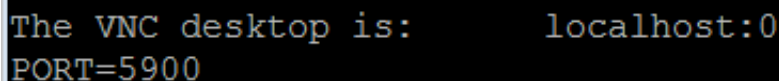
- If the `startx` command cannot be found, run `apt-get install xinit` to install it.
- Running `startx` may result in the `hostname: Name or service not known` error. This error has no effect on starting X Server. Run `hostname` to obtain the host name of the instance, and then modify the `/etc/hosts` file by replacing the `hostname`, which is preceded by `127.0.0.1`, with the actual host name of your instance.

```
root@iZ... Z:~# startx
hostname: Name or service not known
xauth: (stdin):1: bad display name "iZ... Z:1" in "add" command
```

3. Start a new terminal session of the SSH client and run the following command to start x11vnc.

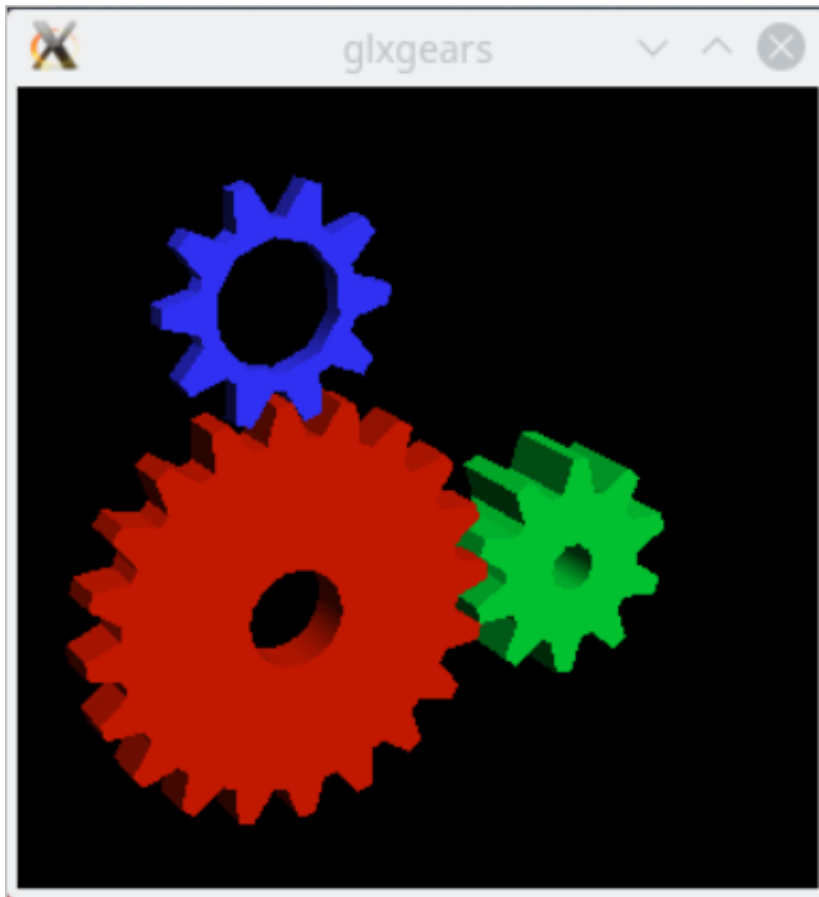
```
x11vnc -display :1
```

If the following message appears, x11vnc has been restarted successfully. Now you can connect to the instance by using a VNC application. In this example, VNC Viewer is used.

A terminal window with a black background and green text. It displays the message: "The VNC desktop is: localhost:0" on the first line and "PORT=5900" on the second line.

```
The VNC desktop is: localhost:0
PORT=5900
```

4. Log on to the ECS console, and [add security group rules](#) in the security group to allow inbound traffic from the Internet on the TCP 5900 port.
5. On the local machine, start VNC Viewer and type in Public IP address of the instance:5900 to connect to the instance and enter the KDE desktop.
6. Run `glxinfo` to view the configurations supported by the current GRID driver:
 - a. Start a new terminal session of the SSH client.
 - b. Run `export DISPLAY=:1`.
 - c. Run `glxinfo -t` to list the configurations supported by the current GRID driver.
7. Run `glxgears` to test the GRID driver:
 - a. On the KDE desktop, right-click the desktop and select Run Command.
 - b. Run `glxgears` to start the testing application. If the following figure appears, the GRID driver works normally.



CentOS 7

This section describes how to install the GRID driver on a GPU instance running CentOS 7.3 64-bit.

Prerequisites

- You have created a gn5, gn5i or gn6v instance. For more information, see [Create a compute optimized instance with GPUs](#). Make sure that the instance can access the Internet.
- You have installed a VNC application on your local machine. In this example, VNC Viewer is used.

Install a GRID driver

To install the GRID driver, follow these steps:

1. [Connect to the Linux instance](#).
2. Run the following commands in sequence to upgrade the system and install the KDE.

```
yum update
yum install kernel-devel
```

```
yum groupinstall "KDE Plasma Workspaces"
```

3. Run `reboot` to restart the system.
4. [Connect to the Linux instance](#) again, and then run the following commands to download and decompress the NVIDIA GRID driver package.

The NVIDIA GRID driver package contains the drivers for various operating systems. For Linux OS, select `NVIDIA-Linux-x86_64-390.57-grid.run`.

```
wget https://nvidia-driver.oss-cn-huhehaote.aliyuncs.com/NVIDIA-Linux-x86_64-390.57-grid.run
```

5. Disable the nouveau driver:
 - a. Run `vim /etc/modprobe.d/blacklist.conf`, and add `blacklist nouveau` to the file.
 - b. Run `vim /lib/modprobe.d/dist-blacklist.conf` and add the following lines.

```
blacklist nouveau
options nouveau modeset=0
```

- c. Run `mv /boot/initramfs-$(uname -r).img /boot/initramfs-$(uname -r)-nouveau.img`.
 - d. Run `dracut /boot/initramfs-$(uname -r).img $(uname -r)`.

6. Run `reboot` to restart the system.
7. Run the following commands in sequence, then follow the prompts to install the NVIDIA GRID driver.

```
chmod 777 NVIDIA-Linux-x86_64-390.57-grid.run
. /NVIDIA-Linux-x86_64-390.57-grid.run
```

8. Run `nvidia-smi` to verify the installation.

If the following message appears, the driver is installed successfully.


```
[root@i-xxxxx ~]# nvidia-smi
Wed Jul  4 16:30:12 2018

+-----+
| NVIDIA-SMI 390.57                  Driver Version: 390.57          |
+-----+-----+
| GPU   Name           Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap|  Memory-Usage | GPU-Util  Compute M. |
+-----+-----+
|  0   Tesla P100-PCIE...    On   | 00000000:00:08.0 Off |                    0 |
| N/A   28C    P0      25W / 250W |  0MiB / 16280MiB |      0%      Default |
+-----+-----+

+-----+
| Processes:                         GPU Memory |
|  GPU       PID    Type    Process name      Usage   |
+-----+-----+
| No running processes found              |
+-----+
```

9. Add License Server to activate the License:

- Run `cd /etc/nvidia` to switch to the `/etc/nvidia` directory.
- Run `cp gridd.conf.template gridd.conf` to create a file named `gridd.conf`.
- Add the following lines about License Server to the `gridd.conf` file.

```
ServerAddress=116.62.19.179
ServerPort=7070
FeatureType=2
EnableUI=TRUE
```

10.Run the following command to install x11vnc.

```
yum install x11vnc
```

11.Run `lspci | grep NVIDIA` to check GPU BusID.

In this example, the GPU BusID is `00:07:0`.

12.Configure the X Server environment:

- Run `nvidia-xconfig --enable-all-gpus --separate-x-screens`.
- Edit `/etc/X11/xorg.conf`: Add your GPU BusID to the Section "Device". In this example, BusID "PCI:0:7:0" is added

```
Section "Device"
    Identifier      "Device0"
    Driver          "nvidia"
    VendorName      "NVIDIA Corporation"
    BoardName       "Tesla P4"
    BusID           "PCI:0:7:0"
EndSection
```


13. Run `reboot` to restart the system.

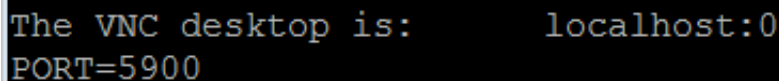
Verify installation

To verify the installation of the GRID driver, follow these steps:

1. Run `startx` to start X Server.
2. Start a new terminal session of the SSH client and run the command to start `x11vnc`.

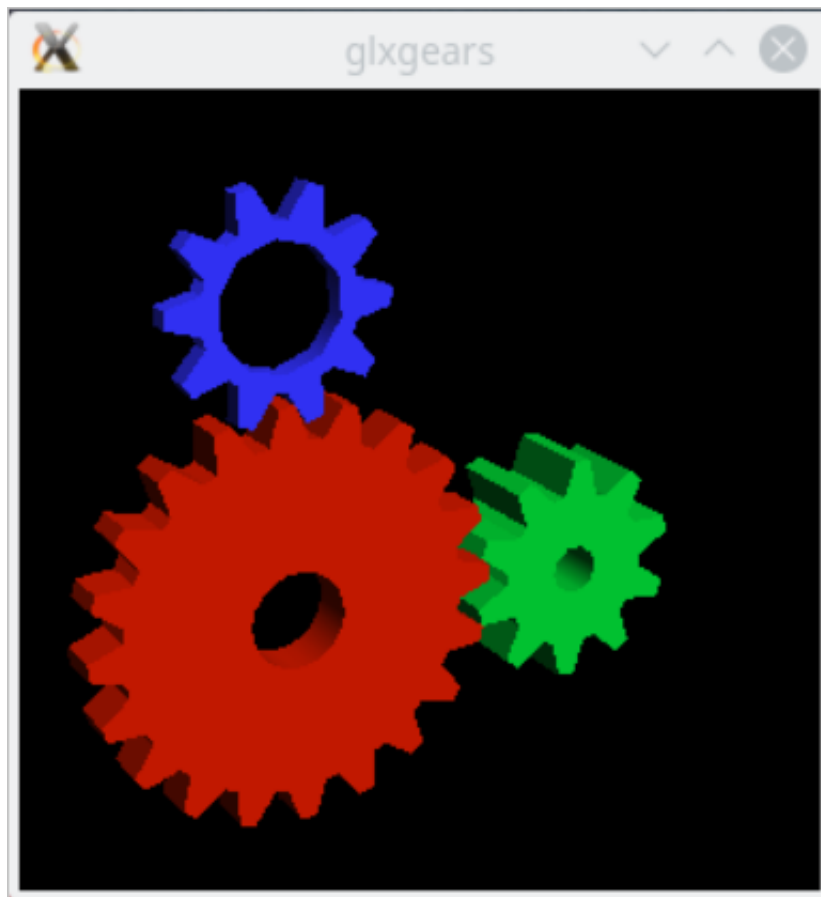
```
x11vnc -display :0
```

If the following message appears, `x11vnc` has been restarted successfully. Now you can connect to the instance by using a VNC application. In this example, VNC Viewer is used.



```
The VNC desktop is:      localhost:0  
PORT=5900
```

3. Log on to the ECS console, and [add security group rules](#) in the security group to allow inbound traffic from the Internet on TCP 5900 port.
4. On the local machine, start VNC Viewer and type in `Public IP address of the instance:5900` to connect to the instance and enter the KDE desktop.
5. Run `glxinfo` to view the configurations supported by the current GRID driver:
 - a. Start a new terminal session of the SSH client.
 - b. Run `export DISPLAY=:0`.
 - c. Run `glxinfo -t` to list the configurations supported by the current GRID driver.
6. Run `glxgears` to test the GRID driver:
 - a. On the VNC Viewer, right-click the desktop and select Run Command.
 - b. Run `glxgears` to start the testing application. If the following image appears, the GRID driver works normally.



8 FaaS instances best practices

8.1 Use RTL compiler on an f1 instance

This article describes how to use Register Transfer Level (RTL) compiler on an f1 instance.



Note:

- All the operations described in this article must be performed by one account in the same region.
- We strongly recommend that you use an f1 instance as a RAM user. To avoid unwanted operations, you must authorize the RAM user to perform required actions only. You must create a role for the RAM user and grant temporary permissions to the role to access the OSS buckets. If you want to encrypt the IP address, grant the RAM user to use Key Management Service (KMS). If you want the RAM user to check permissions, authorize the RAM user to view the resources of an account.

Prerequisites

- Create an f1 instance and add a security group rule to allow Internet access to SSH Port 22 of the instance.



Note:

Only the image we share with you can be used on an f1 instance. For more information, see [Create an f1 instance](#).

- Log on to the [ECS console](#) to obtain the instance ID.
- Activate OSS and [create an OSS bucket](#) to upload your files. The OSS bucket and the f1 instance must be owned by one account and operated in the same region.
- For encryption, activate [Key Management Service \(KMS\)](#).
- To operate FPGA as a RAM user, do the following in advance:
 - [Create a RAM and grant permissions](#).
 - [Create a RAM and grant permissions](#).
 - Use the AccessKey to complete the authentication.

Procedure

To use RTL compiler on an f1 instance, follow these steps.

Step 1. Connect to the f1 instance

Connect to your f1 instance.

Step 2. Configure the basic environment

Run the script to configure the basic environment.

```
source /opt/dcp1_1/script/f1_env_set.sh
```

Step 3. Compile the project

Run the following commands to compile the project.

```
cd /opt/dcp1_1/hw/samples/dma_afu
afu_synth_setup --source hw/rtl/filelist.txt build_synth
cd build_synth/
run.sh
```



Note:

It takes a long time to compile the project.

Step 4. Create an image

To create an image, follow these steps:

1. Run the following commands to initialize `faascmd`.

```
# If needed, add the environment variable and grant permission to
run the commands.
export PATH=$PATH:/opt/dcp1_1/script/
chmod +x /opt/dcp1_1/script/faascmd
# Replace hereIsMySecretId with your AccessKey ID. Replace
hereIsMySecretKey with your AccessKey Secret. faascmd config --id=
hereIsMySecretId --key=hereIsMySecretKey
faascmd config --id=hereIsYourSecretId --key=hereIsYourSecretKey
# Replace hereIsYourBucket with the OSS bucket name in the China
East 1 region.
faascmd auth --bucket=hereIsYourBucket
```

2. Make sure you are at the `/opt/dcp1_1/hw/samples/dma_afu` directory, and run the command to upload the gbs file.

```
faascmd upload_object --object=dma_afu.gbs --file=dma_afu.gbs
```

3. Run the command to create an image.

```
# Replace hereIsYourImageName with your image name.
faascmd create_image --object=dma_afu.gbs --fpgatype=intel --name=
hereIsYourImageName --tags=hereIsYourImageTag --encrypted=false --
shell=V1.1
```

Step 5. Download the image

To download the image, follow these steps:

1. Run the `faascmd list_images` command to check whether the image is created.

If `"State": "success"` exists in the returned result, it means the image is created.

Record the `FpgaImageUUID`. Record the `FpgaImageUUID`.

```
[root@izop: ~]# faascmd list_images
{"FpgaImages": [{"FpgaImage": [{"Name": "Image_1_dma_afu", "Tags": "ImageTag_1_dma_afu", "ShellUUID": "V0.11", "Description": "None", "FpgaImageUUID": "inteld98db1d1-0238", "State": "success", "CreateTime": "Fri Jan 26 2018 10:15:59 GMT+0800 (CST)", "Encrypted": "false", "UpdateTime": "Fri Jan 26 2018 10:17:08 GMT+0800 (CST)"}]}]}
```

2. Run the command to obtain FPGA ID.

```
# Replace hereIsYourInstanceId with your f1 instance ID.
faascmd list_instances --instanceId=hereIsYourInstanceId
```

Record `FpgaUUID` in the returned result.

```
[root@izb: ~]# faascmd list_instances --instanceId=i-bp15n6gzu...
{"Instances": [{"Instance": [{"ShellUUID": "V0.11", "FpgaType": "intel", "FpgaUUID": "0x6c92bf4786940500", "InstanceId": "i-bp15n6gzu...", "Dev": "iceBDF": "05:00.0", "FpgaStatus": "valid"}]}]}
```

3. Run the command to download the image to your f1 instance.

```
# Replace hereIsYourInstanceId with your f1 instance ID. Replace
hereIsFpgaUUID with your FpgaUUID. Replace hereIsImageUUID with your
FpgaImageUUID.
faascmd download_image --instanceId=hereIsYourInstanceId --fpgauuid=
hereIsFpgaUUID --fpgatype=intel --imageuuid=hereIsImageUUID --
imagetype=afu --shell=V0.11
```

4. Run the command to check whether the image is downloaded.

```
# Replace hereIsYourInstanceId with your f1 instance ID. Replace
hereIsFpgaUUID with your FpgaUUID.
faascmd fpga_status --instanceId=hereIsYourInstanceId --fpgauuid=
hereIsFpgaUUID
```

If `"TaskStatus": "operating"` exists in the returned result, and the displayed `FpgaImageUUID` is identical with your recorded `FpgaImageUUID`, the image is downloaded.

```
[root@ ~]# faascmd fpga_status --instanceId=i-bp15n6gzu... --fpgauuid=0x6c92bf4786940500
{"shellUUID": "V0.11", "FpgaImageUUID": "inteld98db1d1-0238", "FpgaUUID": "0x6c92bf4786940500", "InstanceId": "i-bp15n6gzu...", "CreateTime": "Fri Jan 26 2018 10:40:41 GMT+0800 (CST)", "TaskStatus": "operating", "Encrypted": "false"}
0.291(s) elapsed
```

Step 6. Test

Run the commands one by one for test.

```
cd /opt/dcp1_1/hw/samples/dma_afu/sw
make
sudo LD_LIBRARY_PATH=/opt/dcp1_1/hw/samples/dma_afu/sw:$LD_LIBRARY_PATH ./fpga_dma_test 0
```

If the following result is returned, the test is completed.

```
[root@iZ...Z sw]# ./fpga_dma_test use_ase=0
Running test in HW mode
Buffer Verification Success!
Buffer Verification Success!
Running DDR sweep test
Allocated test buffer
Fill test buffer
DDR Sweep Host to FPGA
Measured bandwidth = 5726.623061 Megabytes/sec
Clear buffer
DDR Sweep FPGA to Host
Measured bandwidth = 4473.924267 Megabytes/sec
Verifying buffer..
Buffer Verification Success!
```



Note:

If the Huge pages feature is not enabled, run the following command to enable it.

```
sudo bash -c "echo 20 > /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages"
```

8.2 Use OpenCL on an f1 instance

This article introduces how to use Open Computing Language (OpenCL) to create an image file, and then download the image to an FPGA chip.



Note:

- All the operations described in this article must be performed by one account in the same region.
- We strongly recommend that you use an f1 instance as a RAM user. To avoid unwanted operations, you must authorize the RAM user to perform required actions only. You must create a role for the RAM user and grant temporary

permissions to the role to access the OSS buckets. If you want to encrypt the IP address, grant the RAM user to use Key Management Service (KMS). If you want the RAM user to check permissions, authorize the RAM user to view the resources of an account. Before you begin, complete the following:

Prerequisites

- Create an f1 instance and add a security group rule to allow Internet access to SSH Port 22 of the instance.



Note:

Only the image we share with you can be used on an f1 instance. For more information, see [Create an f1 instance](#).

- Log on to the [ECS console](#) to obtain the instance ID.
- [Create an OSS bucket](#) to upload your custom bitstream files. The OSS bucket and the f1 instance must be owned by one account and in the same region.
- To encrypt your bitstream, activate Key Management Service (KMS).
- To operate an f1 instance as a RAM user, you must do the following operations:
 - [Create a RAM user](#) and [grant permissions](#).
 - [Create a RAM role](#) and [grant permissions](#).
 - Create an AccessKey.

Procedure

To configure the environment of FPGA Server Example, follow these steps.

Step 1. Connect to your f1 instance

[Connect to the Linux instance](#).

Step 2. Install the basic environment

Run the following script to install the base environment.

```
source /opt/dcp1_1/script/f1_env_set.sh
```

Step 3. Download the OpenCL Example

Follow these steps to download the official opencl example.

1. Create the `/opt/tmp` directory, and change the current directory to it.

```
mkdir -p /opt/tmp
```

```
cd /opt/tmp
```

Now, you are at the `/opt/tmp` directory.

```
[root@iZt1z1z1z1z1Z tmp]# pwd
/opt/tmp
```

2. Run the commands one by one to download and decompress the OpenCL Example file.

```
wget https://www.altera.com/content/dam/altera-www/global/en_US/
others/support/examples/download/exm_openc1_matrix_mult_x64_linux.
tgz
tar -zxvf exm_openc1_matrix_mult_x64_linux.tgz
```

The following figure displays the directory after decompression.

```
[root@iZt1z1z1z1z1Z tmp]# tree -L 1
.
├── common
├── exm_openc1_matrix_mult_x64_linux.tgz
└── matrix_mult

2 directories, 1 file
```

3. Change the current directory to the `matrix_mult` directory and run the command for compilation.

```
cd matrix_mult
aoc -v -g --report ./device/matrix_mult.cl
```

The process of compilation takes several hours. You can open a new console, and run the `top` command to monitor processes and system resource usage on the instance and view the status of the compilation process.

Step 4. Upload the configuration file to the OSS bucket

Follow these steps to upload the configuration file.

1. Run the commands to initialize the `faascmd`.

```
# If needed, add the environment variable and grant the permission
to run the commands
export PATH=$PATH:/opt/dcp1_1/script/
chmod +x /opt/dcp1_1/script/faascmd
# Replace hereIsYourSecretId with your AccessKey ID. Replace
hereIsYourSecretKey with your AccessKey Secret
faascmd config --id=hereIsYourSecretId --key=hereIsYourSecretKey
```



```
# Replace hereIsYourBucket with the bucket name of your OSS in the
Region China East 1.
faascmd auth --bucket=hereIsYourBucket
```

2. Change the current directory to the `matrix_mult/output_files` directory, and upload the configuration file.

```
cd matrix_mult/output_files # Now you are accessing/opt/tmp/
matrix_mult/matrix_mult/output_files
faascmd upload_object --object=afu_fit.gbs --file=afu_fit.gbs
```

3. Use `gbs` to create an FPGA image.

```
# Replace hereIsYourImageName with your image name. Replace
hereIsYourImageTag with your image tag.
faascmd create_image --object=dma_afu.gbs --fpgatype=intel --name=
hereIsYourImageName --tags=hereIsYourImageTag --encrypted=false --
shell=V1.1
```

4. Run the `faascmd list_images` command to check whether the image is created. In the returned result, if `"State": "success"` is displayed, it means the image is created. Record the `FpgaImageUUID`.

```
[root@f1p. ...]# faascmd list_images
{"FpgaImages":[{"fpgaImage":{"Name":"Image_1_dma_afu","Tags":"ImageTag_1_dma_afu","ShellUUID":"V0.11","Description":"None","FpgaImageUUID":"intel98db1d1-023...8","State":"success","CreateTime":"Fri Jan 26 2018 10:15:59 GMT+0800 (CST)","Encrypted":"false","UpdateTime":"Fri Jan 26 2018 10:17:08 GMT...
```

Step 5. Download the image to your f1 instance

To download the image to your f1 instance, follow these steps:

1. Run the command to obtain FPGA ID.

```
# Replace hereIsYourInstanceId with your f1 instance ID.
faascmd list_instances --instanceId=hereIsYourInstanceId
```

Returned results sample: Record `FpgaUUID` in the returned result.

```
[root@f1p. ...]# faascmd list_instances --instanceId=i-bp15n6gzu...
{"Instances":[{"Instance":{"ShellUUID":"V0.11","FpgaType":"intel","FpgaUUID":"0xe...","InstanceId":"i-bp15n6gzu...","DeviceBDF":"05:00.0","FpgaStatus":"valid"}}]}
```

2. Run the command to download the image to your f1 instance.

```
# Replace hereIsYourInstanceId with your f1 instance ID. Replace
hereIsFpgaUUID with your FPGA UUID. Replace hereIsImageUUID with
your image UUID.
faascmd download_image --instanceId=hereIsYourInstanceId --fpgauuid=
hereIsFpgaUUID --fpgatype=intel --imageuuid=hereIsImageUUID --
imagetype=afu --shell=V0.11
```

3. Run the command to check whether the image is downloaded.

```
# Replace hereIsYourInstanceId with your f1 instance ID. Replace
hereIsFpgaUUID with your FPGA UUID.
```

```
faascmd fpga_status --fpgauid=hereIsFpgaUUID --instanceId=
hereIsYourInstanceID
```

If "TaskStatus": "operating" exists in the returned result, it means the image is downloaded.

```
[root@iZbpXXXXXZ ~]# faascmd fpga_status --instanceId=i-bp1ite6wvjlcsl3e6s --fpgauid=0x40500
{"shellUUID":"V0.11","FpgaImageUUID":"inteld98db18","FpgaUUID":"0x40500","InstanceId":"i-bp1ite6wvjlcsl3e6s","CreateTime":"Fri Jan 26 2018 10:40:41 GMT+0800 (CST)","TaskStatus":"operating","Encrypted":"false"}
0.291(s) elapsed
```

Step 6. Download the FPGA image to an FPGA chip

To download the FPGA image to an FPGA chip, follow these steps:

1. Open the console in Step 1. If it is closed, repeat Step 1.
2. Run the following command to configure the runtime environment for OpenCL.

```
sh /opt/dcp1_1/opencl/opencl_bsp/linux64/libexec/setup_permissions.
sh
```

3. Run the command to go back to the parent directory.

```
cd .. /.. # Now, you are at the /opt/tmp/matrix_mult directory
```

4. Run the command to compile.

```
make
# Output the environment configuration
export CL_CONTEXT_COMPILER_MODE_ALTERA=3
cp matrix_mult.aocx ./bin/matrix_mult.aocx
cd bin
host matrix_mult.aocx
```

If the following result is returned, it means the configuration is successful. Note that the last line must be `Verification: PASS`.

```
[root@iZbpXXXXXZ bin]# ./host matrix_mult.aocx
Matrix sizes:
  A: 2048 x 1024
  B: 1024 x 1024
  C: 2048 x 1024
Initializing OpenCL
Platform: Intel(R) FPGA SDK for OpenCL(TM)
Using 1 device(s)
  skx_fpga_dcp_ddr : SKX DCP FPGA OpenCL BSP (acl0)
Using AOCX: matrix_mult.aocx
Generating input matrices
Launching for device 0 (global size: 1024, 2048)
Time: 40.415 ms
Kernel time (device 0): 40.355 ms
Throughput: 106.27 GFLOPS
Computing reference output
Verifying
```

Verification: PASS

8.3 Best practices for OpenCL on an f3 instance

This topic describes how to use Open Computing Language (OpenCL) to create an image, and then download it to an FPGA chip in an f3 instance.



Note:

- All the operations described in this topic must be performed by one account in the same region.
- We recommend that you use an f3 instance as a RAM user. You must create a role for the RAM user and grant the role temporary permissions to access the specified OSS buckets.

Prerequisites

- [Create an f3 instance](#).



Note:

- Only the image we share with you can be used on an f3 instance.
 - Select Assign public IP when creating an instance, so that the instance can access the Internet.
 - The security group of the f3 instance has added the rule for allowing access to the SSH port 22.
- Log on to the ECS console and obtain the instance ID of your f3 instance.
 - Create an OSS bucket in the same region as your f3 instance by using the same account. For more information, see [Sign up for OSS](#) and [Create a bucket](#).
 - To operate FPGA as a RAM user, do the following in advance:
 - [Create a RAM user](#) and [grant permissions](#).
 - [Create a RAM role](#) and [grant permissions](#).
 - Obtain the AccessKey ID and AccessKey Secret.

Procedure

To create an image and download it to an FPGA chip on an f3 instance by using OpenCL, follow these steps.

Step 1. Set up the environment

To set up the environment, follow these steps:

1. [Connect to the f3 instance.](#)



Note:

The subsequent compilation process may take a few hours. We recommend that you log on through screen or nohub, so as to avoid forced logout due to an SSH timeout.

2. Run the command to install Screen.

```
yum install screen -y
```

3. Run the command to enter Screen.

```
screen -S f3openc1
```

4. Run the command to set up the environment.

```
source /root/xbinst_oem/f3_env_setup.sh xocl # Run the command each  
time you open a new terminal window
```



Note:

- Configuring the environment involves installing the xocl driver, setting the vivado environment variable, checking the vivado license, detecting the aliyun-f3 sdaccel platform, configuring 2018.2 runtime, and detecting the faascmd version.
- If you want to run an emulation of sdaccel, do not run the above command to configure the environment. Instead, you only need to configure the environment variable for vivado separately.
- We recommend that you use Makefile for emulation.

Step 2. Compile a binary file

- Example 1: vadd

To compile the vadd binary file, follow these steps:

1. Copy the `example` directory.

```
cp -rf /opt/Xilinx/SDx/2018.2/examples . /
```

2. Enter the `vadd` directory.

```
cd examples/vadd/
```

3. Run the command `cat sdaccel.mk | grep "XDEVICE="` to view the value of `XDEVICE`. Make sure its configuration is `XDEVICE=xilinx_aliyun-f3_dynamic_5_0`.

4. Follow these steps to modify the `common.mk` file.

a. Run the `vim ../common/common.mk` command to open the file.

b. At the end of the code line 61, add the compilation parameter `--xp param:`

`compiler.acceleratorBinaryContent=dcg` (the parameter may be in the line 60-62, depending on your file). The modified code is:

```
CLCC_OPT += $(CLCC_OPT_LEVEL) ${DEVICE_REPO_OPT} --platform ${XDEVICE} -o ${XCLBIN} ${KERNEL_DEFS} ${KERNEL_INCS} --xp param:compiler.acceleratorBinaryContent=dcg
```



Note:

Given that you must submit a DCP file to the compilation server, you need to add the parameter `--xp param:compiler.acceleratorBinaryContent=dcg`, so that Xilinx® OpenCL™ Compiler (xocc) generates a DCP file (instead of a bit file) after the placement and routing is complete.

5. Run the command to compile the program.

```
make -f sdaccel.mk xbin_hw
```

If the following information is displayed, the compilation of the binary file has started. This process may take several hours.

```
[root@ ~]# cd vadd/ && make -f sdaccel.mk xbin_hw
make SDA_FLOW=hw xbin -f sdaccel.mk
make[1]: Entering directory '/root/xilinx_example/examples/vadd'
xocc -c -t hw --platform xilinx_aliyun-f3_dynamic_5_0 --xp param:compiler.acceleratorBinaryContent=dcg -s --kernel krnl_vadd krnl_vadd.cl -o bin_vadd_hw.xo
***** xocc v2018.2 (64-bit)
**** SW Build 2258646 on Thu Jun 14 20:02:38 MDT 2018
** Copyright 1986-2018 Xilinx, Inc. All Rights Reserved.
Attempting to get a license: ap_opencl
Feature available: ap_opencl
INFO: [XOCC 60-585] Compiling for hardware target
Running SDx Rule Check Server on port:39076
INFO: [XOCC 60-895] Target platform: /opt/Xilinx/SDx/2018.2/platforms/xilinx_aliyun-f3_dynamic_5_0/xilinx_aliyun-f3_dynamic_5_0.xpfm
INFO: [XOCC 60-423] Target device: xilinx_aliyun-f3_dynamic_5_0
```

• Example 2: `kernel_global_bandwidth`

Follow these steps to compile the `kernel_global_bandwidth` binary file:

1. Clone xilinx 2018.2 example.

```
git clone https://github.com/Xilinx/SDAccel_Examples.git
cd SDAccel_Examples/
git checkout 2018.2
```



Note:

The git branch must be the 2018.2 version.

2. Run the `cd getting_started/kernel_to_gmem/kernel_global_bandwidth/` command to enter the directory.

3. Follow these steps to modify the Makefile file.

a. Run the `vim Makefile` command to open the file.

b. Set `DEVICES=xilinx_aliyun-f3_dynamic_5_0`.

c. In the code line 33, add the compilation parameter `--xp param:compiler`.

`acceleratorBinaryContent=dcg`. The modified code is:

```
CLFLAGS +=--xp "param:compiler.acceleratorBinaryContent=dcg" --xp
xp "param:compiler.preserveHlsOutput=1" --xp "param:compiler
.generateExtraRunData=true" --max_memory_ports bandwidth -
DNDDR_BANKS=$(ddr_banks)
```

4. Run the command to compile the program.

```
make TARGET=hw
```

If the following information is displayed, the compilation of the binary file has started. This process may take several hours.

```
[root@ip-10.10.10.10 ~]# make TARGET=hw
mkdir -p ./xclbin
/opt/Xilinx/SDx/2018.2/bin/xccp -I /opt/Xilinx/SDx/2018.2/runtime/include/1_2/ -I /opt/Xilinx/SDx/2018.2/Vivado_HLS/include/ -O0 -g -Wall -fmessage-length=0 -std=c++14 -DNDDR_BANKS=4 -I../
../libs/xcl2 src/kernel_global_bandwidth.cpp ../../libs/xcl2/xcl2.cpp -o 'kernel_global' -lOpenCL -lpthread -lrt -lstdc++ -L/opt/Xilinx/SDx/2018.2/runtime/lib/x86_64
src/kernel_global_bandwidth.cpp: In function 'int main(int, char**)':
src/kernel_global_bandwidth.cpp:260:89: warning: value computed is not used [-Wunused-value]
    nsduration = OCL_CHECK(err, event.getProfilingInfo(CL_PROFILING_COMMAND_END>(&err)) - OCL_CHECK(err, event.getProfilingInfo(CL_PROFILING_COMMAND_START>(&err)));
                                                                    ^
mkdir -p ./xclbin
/opt/Xilinx/SDx/2018.2/bin/xocc -t hw --platform xilinx_aliyun-f3_dynamic_5_0 --save-temps --xp "param:compiler.acceleratorBinaryContent=dcg" --xp "param:compiler.preserveHlsOutput=1" --xp
"param:compiler.generateExtraRunData=true" --max_memory_ports bandwidth -DNDDR_BANKS=4 -c -k bandwidth -I'src' -o'xclbin/bandwidth.hw.xilinx_aliyun-f3_dynamic_5_0.xo' 'src/kernel.cl'
***** xocc v2018.2 (64-bit)
***** SW Build 2258646 on Thu Jun 14 20:02:38 MDT 2018
** Copyright 1986-2018 Xilinx, Inc. All Rights Reserved.
```

Step 3. Check the packaging script

Run the command to check whether the packaging script exists or not.

```
file /root/xbinst_oem/sdaccel_package.sh
```

If the returned message contains cannot open (No such file or directory), the file does not exist. You need to download the script by running the following command.

```
wget http://fpga-tools.oss-cn-shanghai.aliyuncs.com/sdaccel_package.sh
```

Step 4. Create an image

To create an image, follow these steps:

1. Run the command to set up the OSS environment.

```
faascmd config --id=hereIsMySecretId --key=hereIsMySecretKey #
Replace hereIsMySecretId, hereIsMySecretKey with your AccessKeyID,
AccessKeySecret
faascmd auth --bucket=hereIsMyBucket # Replace hereIsMyBucket with
your bucket name
```

2. Run the `ls` command to obtain the file suffixed by `.xclbin`.

```
[root@... vadd]# ls
bin_vadd_hw.xclbin      krnl_vadd.cl  vadd.cpp
description.json        README.md     vadd.h
Export_Compliance_Notice.md sdaccel.mk    _xocc_krnl_vadd_bin_vadd_hw.dir
```

3. Run the command to package the binary file.

```
/root/xbinst_oem/sdaccel_package.sh -xclbin=/opt/Xilinx/SDx/2017.4.
op/examples/vadd/bin_vadd_hw.xclbin
```

After the packaging is completed, you can find a package file in the same directory , as shown in the following figure.

```
[root@... vadd]# ls
17_10_28-021904-primary.bit      krnl_vadd.cl
SDAccel_Kernel.tar.gz           README.md
17_10_28-021904-xclbin.xml       sdaccel.mk
bin_vadd_hw.xclbin               to_aliyun
description.json                 vadd.cpp
Export_Compliance_Notice.md      vadd.h
header.bin                       _xocc_krnl_vadd_bin_vadd_hw.dir
```

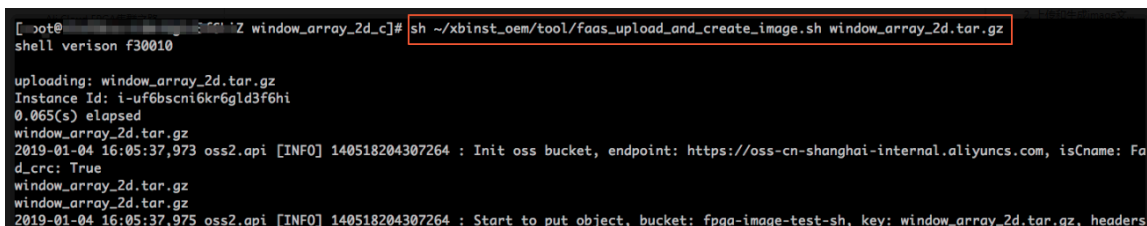
Step 5. Download the image

You can use a scripted process or step-by-step process to upload the package file and download the FPGA image.

- Scripted process: Only applicable to f3 instances with one FPGA chip.

- Run the following commands to upload the package and generate the image file.

```
sh /root/xbinst_oem/tool/faas_upload_and_create_image.sh <bit.tar.gz - the package to upload>
```



```
[root@i-uf6bscni6kr6gld3f6hi ~]# sh /root/xbinst_oem/tool/faas_upload_and_create_image.sh window_array_2d.tar.gz
shell verison f30010
uploading: window_array_2d.tar.gz
Instance Id: i-uf6bscni6kr6gld3f6hi
0.065(s) elapsed
window_array_2d.tar.gz
2019-01-04 16:05:37,973 oss2.api [INFO] 140518204307264 : Init oss bucket, endpoint: https://oss-cn-shanghai-internal.aliyuncs.com, isCName: Fa
d_crc: True
window_array_2d.tar.gz
2019-01-04 16:05:37,975 oss2.api [INFO] 140518204307264 : Start to put object, bucket: fpga-image-test-sh, key: window_array_2d.tar.gz, headers
```

- Download the image file.

```
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz -
package name> <0/1> # The last number <0/1> stands for the FPGA
serial No. in the instance
```

0 indicates the first FPGA of the f3 instance. For single-FPGA instances, the FPGA serial No. is always 0. For instances with multiple FPGAs, such as an instance with four FPGAs, the serial No. are 0, 1, 2 and 3.

To download the same image to multiple FPGAs, add the serial No. to the end.

For example, run the command to download the same image to four FPGA chips:

```
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz -
package name> 0
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz -
package name> 1
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz -
package name> 2
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz -
package name> 3
```

- Step-by-step process: [Use the faascmd tool](#) to perform operations.

- Run the command to upload the package to your OSS bucket. Then, upload gbs in your OSS bucket to the OSS bucket in the FaaS administrative unit.

```
faascmd upload_object --object=bit.tar.gz --file=bit.tar.gz
```



```
faascmd create_image --object=bit.tar.gz --fpgatype=xilinx --name=
hereIsFPGAImageName --tags=hereIsFPGAImageTag --encrypted=false --
shell=hereIsShellVersionOfFPGA
```

```
[root@iz-... ~]# faascmd upload_object --object=rion.zj_test_SDAccel_Kernel.tar.gz --file=18_05_03-222718_SDAccel_Kernel.tar
.gz
rion.zj_test_SDAccel_Kernel.tar.gz
18_05_03-222718_SDAccel_Kernel.tar.gz
4.735(s) elapsed
```

```
[root@iz-... ~]# faascmd create_image --object=rion.zj_test_SDAccel_Kernel.tar.gz --fpgatype=xilinx --name=rion.zj_xilinx_f3
_test --tags=hereIsFPGAImageTag --encrypted=false --shell=f30001
{"Name": "rion.zj_xilinx_f3_test", "CreateTime": "Fri May 04 2018 20:24:21 GMT+0800 (CST)", "ShellUUID": "f30001", "Description": "None", "FpgaImageUU
ID": "xilinx1", "State": "5", "State": "queued"}
```

2. Run the command to view if the FPGA image is downloadable.

```
faascmd list_images
```

If the returned message shows `State:compiling`, the FPGA image is being compiled. If the returned message shows `State:success`, the FPGA image is ready for downloading. Find `FpgaImageUUID` and note it down.

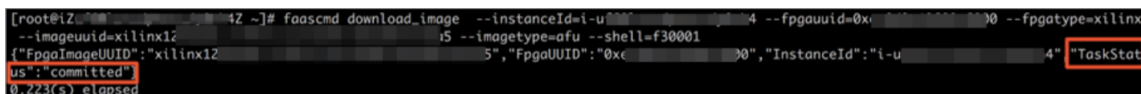
```
[root@... ~]# faascmd list_images
{
  "FpgaImages": {
    "fpgaImage": [
      {
        "CreateTime": "Fri Jan 04 2019 16:05:43 GMT+0800 (CST)",
        "Description": "None",
        "Encrypted": "false",
        "FpgaImageUUID": "xilinx8858a3c1-...",
        "Name": "window_array_2d.tar.gz",
        "ShellUUID": "f30010",
        "State": "compiling",
        "Tags": "hereIsFPGAImageTag",
        "UpdateTime": "Fri Jan 04 2019 16:05:44 GMT+0800 (CST)"
      },
      {
        "CreateTime": "Thu Jan 03 2019 15:58:58 GMT+0800 (CST)",
        "Description": "None",
        "Encrypted": "false",
        "FpgaImageUUID": "xilinx6cbd48c1-...",
        "Name": "vadd.tar.gz",
        "ShellUUID": "f30010",
        "State": "success",
        "Tags": "hereIsFPGAImageTag",
        "UpdateTime": "Thu Jan 03 2019 16:32:32 GMT+0800 (CST)"
      }
    ]
  }
}
```

3. Run the following command. In the returned message, find and note down `FpgaUUID`.

```
faascmd list_instances --instanceId=hereIsYourInstanceId # Replace
hereIsYourInstanceId with the f3 instance ID
```

4. Run the command to download the FPGA image.

```
faascmd download_image --instanceId=hereIsYourInstanceId --
fpgauid=hereIsFpgaUUID --fpgatype=xilinx --imageuuid=hereIsImag
eUUID --imagetype=afu --shell=hereIsShellVersionOfFpga
# Replace hereIsYourInstanceId with the f3 instance ID, hereIsFpga
UUID with the FpgaUUID, and hereIsImageUUID with the FpgaImageUUID
```



```
[root@iz... ~]# faascmd download_image --instanceId=i-u... 4 --fpgauid=0x... 30 --fpgatype=xilinx
--imageuuid=xilinx12... 5 --imagetype=afu --shell=f30001
{"FpgaImageUUID":"xilinx12... 5","FpgaUUID":"0x... 30","InstanceID":"i-u... 4","TaskStat
us":"committed"}
0.223(s) elapsed
```

5. Run the command to view if the image is downloaded successfully.

```
faascmd fpga_status --fpgauid=hereIsFpgaUUID --instanceId=
hereIsYourInstanceId # Replace hereIsFpgaUUID with the obtained
FpgaUUID, and hereIsYourInstanceId with the f3 instance ID
```

Below is an example of the returned message. If the FpgaImageUUID in the message is the same as the FpgaImageUUID you note down and the message shows "TaskStatus": "valid", the image is downloaded successfully.



```
[root@iz... ~]# faascmd fpga_status --fpgauid=0xe... 30 --instanceId=i-u... 4
{"shellUUID":"f30001","FpgaImageUUID":"xilinx1... 5","FpgaUUID":"0xe... 30","InstanceID":"i-u... 4",
"CreateTime":"Fri May 04 2018 21:25:53 GMT+0800 (CST)","TaskStatus":"valid","Encrypted":"false"}
0.263(s) elapsed
```

Step 6: Run the Host program

To run the Host program, follow these steps:

1. Run the following command to configure the environment.

```
source /root/xbinst_oem/f3_env_setup.sh xocl # Run the command each
time you open a new terminal window
```

2. Configure the sdaccel.ini file.

In the directory where the Host binary file is located, run the `vim sdaccel.ini` command to create the sdaccel.ini file and enter the following content.

```
[Debug]
profile=true
[Runtime]
runtime_log = "run.log"
hal_log = hal.log
ert=false
kds=false
```

3. Run the Host.

- For vadd, run the command:

```
make -f sdaccel.mk host
```

```
./vadd bin_vadd_hw.xclbin
```

- For `kernel_global_bandwidth`, run the command:

```
./kernel_global
```

If `Test Passed` is returned, the test is successful.

Other common commands

This section introduces some common commands for f3 instances.

Task	Command
View the help document	<code>make -f ./sdaccel.mk help</code>
Run software emulation	<code>make -f ./sdaccel.mk run_cpu_em</code>
Run hardware emulation	<code>make -f ./sdaccel.mk run_hw_em</code>
Compile the host code only	<code>make -f ./sdaccel.mk host</code>
Compile and generate files for downloading	<code>make -f sdaccel.mk xbin_hw</code>
Clean a work directory	<code>make -f sdaccel.mk clean</code>
Forcibly clean a work directory	<code>make -f sdaccel.mk cleanall</code>



Note:

- During emulation, follow the Xilinx emulation process. You do not need to set up the `f3_env_setup` environment.
- The SDAccel runtime and SDAccel development platform are available in the official f3 images provided by Alibaba Cloud. You can also download them at [SDAccel runtime](#) and [SDAccel development platform](#).

8.4 Best practices for RTL design on an f3 instance

This topic describes how to implement the Register Transfer Level (RTL) design on an f3 instance.



Note:

- All the operations described in this topic must be performed by one account in the same region.

- We recommend that you use an f3 instance as a RAM user. To avoid unwanted operations, you must authorize the RAM user to perform required actions only. To use the FaaS service, you need to authorize the FaaS service account to access the OSS bucket that you specify. Therefore, you need to create the service role faasRole in the RAM console, and grant it the faasPolicy permission. If you want to encrypt IP addresses by using the Key Management Service (KMS), you must authorize the KMS-related permissions in faasPolicy.

Prerequisites

- [Create an f3 instance](#) and add a security group rule to allow Internet access to SSH port 22 of the instance.
- Log on to the [ECS console](#) to obtain the instance ID on the details page of the f3 instance.
- [Create an OSS bucket](#) in China East 2 (Shanghai) for the FaaS service.



Note:

The bucket will provide read and write access to the FaaS administrative account. We recommend that you do not store objects that are not related to FaaS.

- To operate an f3 instance as a RAM user, do the following:
 - [Create a RAM user](#) and [grant permissions](#).
 - [Create a RAM role](#) and [grant permissions](#).
 - Create the AccessKey ID and AccessKey Secret.

Procedure

1. [Connect to your f3 instance](#).



Note:

It takes two or three hours to compile the project. We recommend that you use nohup or VNC to connect to the instance to avoid unexpected disconnection.

2. Download and decompress the [RTL reference design](#).
3. Configure the f3 environment.

- If the driver is `xdma`, run the following command to configure the environment:

```
source /root/xbinst_oem/F3_env_setup.sh xdma #Run this command  
each time you open a new terminal window
```

- If the driver is `xocl`, run the following command to configure the environment:

```
source /root/xbinst_oem/F3_env_setup.sh xocl #Run this command  
each time you open a new terminal window
```

**Note:**

Configuring the environment mainly includes mounting the `xdma` or `xocl` driver, setting the `vivado` environment variable, checking the `vivado` license, detecting the `aliyun-f3` `sdaccel` platform, configuring 2018.2 runtime, and detecting the `faascmd` version.

4. Specify an OSS bucket.

```
faascmd config --id=hereIsYourSecretId --key=hereIsYourSecretKey #  
Replace hereIsYourSecretId and hereIsYourSecretKey with your RAM  
user AccessKey  
faascmd auth --bucket=hereIsYourBucket #Replace hereIsYourBucket  
with your OSS bucket name
```

5. Run the following commands to compile the RTL project:

```
cd <decompressed directory>/hw/ # Enter the decompressed hw  
directory  
sh compiling.sh
```

**Note:**

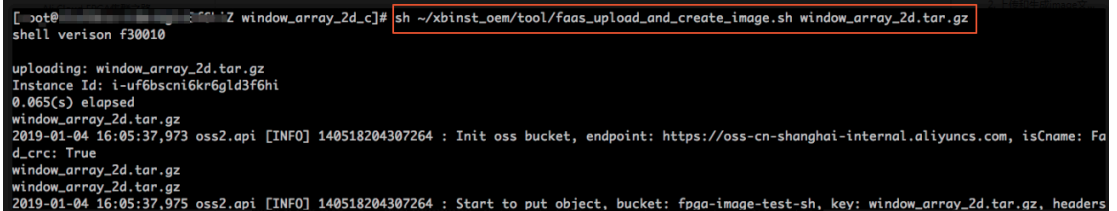
It takes two or three hours to compile the project.

6. Upload the Netlist files and download the FPGA image. You can use the scripted process or the step-by-step process to finish this task.

- Scripted process: Applicable to the `f3` instances with a single FPGA chip.

- a. Run the following commands to upload the package and generate the image file:

```
sh /root/xbinst_oem/tool/faas_upload_and_create_image.sh <bit.tar.gz - the package to upload>
```



```
[root@iZuf6bscni6kr6gld3f6hi ~]# sh ~/xbinst_oem/tool/faas_upload_and_create_image.sh window_array_2d.tar.gz
shell version f30010
uploading: window_array_2d.tar.gz
Instance Id: i-uf6bscni6kr6gld3f6hi
0.065(s) elapsed
window_array_2d.tar.gz
2019-01-04 16:05:37.973 oss2.api [INFO] 140518204307264 : Init oss bucket, endpoint: https://oss-cn-shanghai-internal.aliyuncs.com, isCName: False, crc: True
window_array_2d.tar.gz
window_array_2d.tar.gz
2019-01-04 16:05:37.975 oss2.api [INFO] 140518204307264 : Start to put object, bucket: fpga-image-test-sh, key: window_array_2d.tar.gz, headers
```

- b. Download the image file.

```
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz - the package filename> <0/1> # The last number <0/1> stands for the FPGA serial No. of the instance
```

0 indicates the first FPGA of the f3 instance. For single-FPGA instances, the FPGA serial No. is always 0. For instances with multiple FPGAs, such as an instance with four FPGAs, the serial No. are 0, 1, 2 and 3.

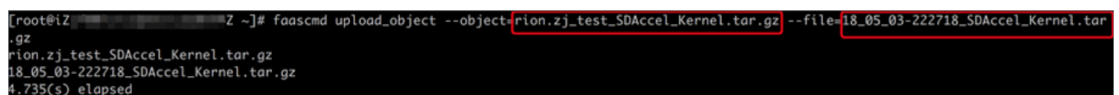
To download the same image to multiple FPGAs, add the serial No. to the end of the command. For example, to download the same image to four FPGAs, use the following commands:

```
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz - package filename> 0
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz - package filename> 1
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz - package filename> 2
sh /root/xbinst_oem/tool/faas_download_image.sh <bit.tar.gz - package filename> 3
```

- Step-by-step process: [Use the faascmd tool](#) to perform the operations.

- a. Run the following commands to upload the package to your OSS bucket, and then upload gbs in your OSS bucket to the OSS bucket of the FaaS unit:

```
faascmd upload_object --object=bit.tar.gz --file=bit.tar.gz
faascmd create_image --object=bit.tar.gz --fpgatype=xilinx --name=hereIsFPGAImageName --tags=hereIsFPGAImageTag --encrypted=false --shell=hereIsShellVersionOfFPGA
```



```
[root@iZuf6bscni6kr6gld3f6hi ~]# faascmd upload_object --object=rion.zj_test_SDAccel_Kernel.tar.gz --file=18_05_03-222718_SDAccel_Kernel.tar.gz
rion.zj_test_SDAccel_Kernel.tar.gz
18_05_03-222718_SDAccel_Kernel.tar.gz
4.735(s) elapsed
```

```
[root@iz-2 ~]# faascmd create_image --object=rion.zj_test_SDAccel_Kernel.tar.gz --fpgatype=xilinx --name=rion.zj_xilinx_f3_test --tags=hereIsFPGAImageTag --encrypted=false --shell=f30001
{"Name": "rion.zj_xilinx_f3_test", "CreateTime": "Fri May 04 2018 20:24:21 GMT+0800 (CST)", "ShellUUID": "f30001", "Description": "None", "FpgaImageUUID": "xilinx1.5", "State": "queued"}
0.221(s) elapsed
```

- b. Run the following command to check if the FPGA image is ready for downloading:

```
faascmd list_images
```

If the returned message shows `State:compiling`, the FPGA image is being compiled, and you still need to wait. If the returned message shows `State:success`, the FPGA image is ready for downloading. Find the `FpgaImageUUID` and note it down.

```
[root@ ~]# faascmd list_images
{
  "FpgaImages": {
    "fpgaImage": [
      {
        "CreateTime": "Fri Jan 04 2019 16:05:43 GMT+0800 (CST)",
        "Description": "None",
        "Encrypted": "false",
        "FpgaImageUUID": "xilinx8858a3c1-...",
        "Name": "window_array_2d.tar.gz",
        "ShellUUID": "f30010",
        "State": "compiling",
        "Tags": "hereIsFPGAImageTag",
        "UpdateTime": "Fri Jan 04 2019 16:05:44 GMT+0800 (CST)"
      },
      {
        "CreateTime": "Thu Jan 03 2019 15:58:58 GMT+0800 (CST)",
        "Description": "None",
        "Encrypted": "false",
        "FpgaImageUUID": "xilinx6cbd48c1-...",
        "Name": "vadd.tar.gz",
        "ShellUUID": "f30010",
        "State": "success",
        "Tags": "hereIsFPGAImageTag",
        "UpdateTime": "Thu Jan 03 2019 16:32:32 GMT+0800 (CST)"
      }
    ]
  }
}
```

- c. Run the following command. In the returned message, note down the `FpgaUUID`.

```
faascmd list_instances --instanceId=hereIsYourInstanceId #
Replace hereIsYourInstanceId with the f3 instance ID
```

- d. Run the following command to download the FPGA image:

```
faascmd download_image --instanceId=hereIsYourInstanceId
--fpgauuid=hereIsFpgaUUID --fpgatype=xilinx --imageuuid=
hereIsImageUUID --imagetype=afu --shell=hereIsShellVersionOf
Fpga
```



```
# Replace hereIsYourInstanceId with the f3 instance ID,
hereIsFpgaUUID with the obtained FpgaUUID, and hereIsImageUUID
with the obtained FpgaImageUUID
```

```
[root@iz... ~]# faascmd download_image --instanceId=i-u... --fpgauid=0xe... --fpgatype=xilinx
--imageuid=xilinx12... --imagetype=afu --shell=f30001
{"FpgaImageUUID":"xilinx12...","FpgaUUID":"0xe...","InstanceId":"i-u..."}
"TaskStatus":"committed"
0.223(s) elapsed
```

- e. Run the following command to check whether the image has been successfully downloaded:

```
faascmd fpga_status --fpgauuid=hereIsFpgaUUID --instanceId=
hereIsYourInstanceId # Replace hereIsFpgaUUID with the obtained
FpgaUUID, and hereIsYourInstanceId with the f3 instance ID
```

The following is an example of the returned message. If the FpgaImageUUID in the message is identical to the FpgaImageUUID you note down, and the message shows "TaskStatus":"valid", the image has been successfully downloaded.

```
[root@iz... ~]# faascmd fpga_status --fpgauid=0xe... --instanceId=i-u...
{"shellUUID":"f30001","FpgaImageUUID":"xilinx1...","FpgaUUID":"0xe...","InstanceId":"i-u..."}
{"CreateTime":"Fri May 04 2018 21:25:53 GMT+0800 (CST)","TaskStatus":"valid","Encrypted":"false"}
0.263(s) elapsed
```

FAQ

How do I view the details of errors that occur during image upload?

If your project reports errors during image upload, such as compilation errors, you can view the error details in two ways:

- Check `faas_compiling.log`. When the upload script `faas_upload_and_create_image.sh` is used, `faas_compiling.log` is automatically downloaded and printed onto the terminal if compilation fails.
- Run the command to view the log file: `sh /root/xbinst_oem/tool/faas_check_log.sh <bit.tar.gz - package uploaded previously>`

How do I reload the image?

To reload the image, follow these steps:

1. Uninstall the driver.

- If you have installed the `xdma` driver, run the command `sudo rmmod xdma` in the instance to uninstall it.
- If you have installed the `xocl` driver, run the command `sudo rmmod xocl` in the instance to uninstall it.

2. Download the image in either of the two ways :

- Use the script.

```
sh faas_download_image.sh bit.tar.gz <0/1> #The last number stands  
for the FPGA serial No. of the instance
```

- Use faascmd.

```
faascmd download_image --instanceId=hereIsYourInstanceId --  
fpgauuid=hereIsFpgaUUID --fpgatype=xilinx --imageuuid=hereIsImag  
eUUID --imagetype=afu --shell=hereIsShellVersionOfFpga
```

3. Install the driver.

- To install the `xdma` driver, run the following command:

```
sudo depmod  
sudo modprobe xdma
```

- To install the `xocl` driver, run the following command:

```
sudo depmod  
sudo modprobe xocl
```

8.5 faascmd tool

8.5.1 faascmd overview

`faascmd` is a command-line tool provided by the Alibaba Cloud FPGA cloud server. It is a script that is developed based on the Python SDK.

You can use `faascmd` to:

- Perform authorization and related operations.
- Manage and operate FPGA images.
- View and upload objects.
- Obtain information about FPGA instances.

8.5.2 Install faascmd

This topic describes how to download and install `faascmd`.

Preparations

- Perform the following steps on the instance for which you want to run `faascmd`:

1. Run the following command to check that the Python version is 2.7.x.

```
python -V
```

```
[root@testhost script]# python -V
Python 2.7.5
```

2. Install the Python module by running the following commands:

```
pip -q install oss2
pip -q install aliyun-python-sdk-core
pip -q install aliyun-python-sdk-faas
pip -q install aliyun-python-sdk-ram
```

3. Run the following command to check that the aliyun-python-sdk-core version is 2.11.0 or later.

```
cat /usr/lib/python2.7/sitepackages/aliyun-sdk-core/__init__.py
```

```
[root@testhost python2.7]# cat /usr/lib/python2.7/site-packages/aliyun-sdk-core/__init__.py
version = "2.11.0" [root@testhost python2.7]#
```



Note:

If the version is earlier than 2.11.0, run `pip install --upgrade aliyun-python-sdk-core` to upgrade aliyun-python-sdk-core to the latest version.

- Obtain the AccessKeyID and AccessKeySecret of the RAM user.

Procedure

1. Log on to your instance and run `wget http://fpga-tools.oss-cn-shanghai.aliyuncs.com/faascmd` in the current or any other directory to download faascmd.



Note:

When you [configure faascmd](#), you need to add the absolute path of the directory where faascmd is installed to the PATH variable.

2. Add executable permissions to faascmd by running the following command:

```
chmod +x faascmd
```

8.5.3 Configure faascmd

Before using faascmd, you need to configure the related environment variable and the AccessKey of the RAM user.

Procedure

1. Log on to your instance and configure the PATH environment variable by running the following command:

```
export PATH=$PATH:<path where faascmd is located>
```

2. Configure the AccessKey (that is, the AccessKeyId and AccessKeySecret) by running the following command:

```
faascmd config --id=<yourAccessKeyID> --key=<yourAccessKeySecret>
```

```
[root@testhost script]# faascmd config --id= --key=
Your configuration is saved into /root/.faascredentials .
[root@testhost script]#
```

8.5.4 Use faascmd

This topic describes how to use faascmd commands.

Prerequisite

You have [configured faascmd](#) before using it.

Syntax description

- All commands and parameters provided by faascmd are case-sensitive.
- There must be no space before and after equal signs (=) in the parameters of faascmd commands.

Authorize users

The `faascmd auth` command is used to authorize the faas admin user to access the users' OSS buckets.

Prerequisites

1. You have created an OSS bucket for FaaS to upload the originally compiled DCP file.
2. You have created a folder named `compiling_logs` in the FaaS OSS bucket.

Command format

```
faascmd auth --bucket=<yourFaasOSSBucketName>
```

Code example

```
[root@testhost script]# faascmd auth --bucket=juliabucket
faasRole has existed!
RAMSECTION has existed!
OSSSECTION has existed!
RoleArn: acs:ram::[REDACTED]:role/faasrole
Create role success
faasPolicy has not existed! Create it Now!
Create policy success
Attach policy to role success
0.459(s) elapsed
```

**Note:**

If an Alibaba Cloud account has multiple RAM user accounts, we recommend that the RAM user accounts share an OSS bucket to prevent authorization policies from being repeatedly modified or overwritten.

View authorization policies

The `faascmd list_policy` command is used to view whether the specified OSS bucket has been added to the corresponding authorization policy (faasPolicy).

Command format

```
faascmd list_policy
```

Code example

```
[root@testhost script]# faascmd list_policy
VersionId : v1   CreateTime : 2018-11-09T03:22:01Z   IsDefaultVersion : True
{
  "Statement": [
    {
      "Action": "ecs:DescribeInstances",
      "Effect": "Allow",
      "Resource": "acs:ecs:*:*:*"
    },
  ],
}
```

**Note:**

You need to check whether your OSS bucket and OSS bucket/compiling_logs appear in the policy information.

Delete authorization policies

The `faascmd delete_policy` command is used to delete authorization policies (faasPolicy).

Command format

```
faascmd delete_policy
```

Code example

```
[root@testhost script]# faascmd delete_policy
Detach faasPolicy from faasRole successfully!!!
Delete the faasPolicy successfully!!!
0.306(s) elapsed
```



Note:

If an Alibaba Cloud account has multiple RAM user accounts, we recommend that you delete the target policy in the RAM console to prevent incorrect authorization policy deletion.

View all objects under an OSS bucket

The `faascmd list_objects` command is used to view all objects under your OSS bucket.

Command format

```
faascmd list_objects
```

Code example

```
[root@testhost script]# faascmd list_objects
compiling_logs/
juliabucket
juliafile
0.081(s) elapsed
[root@testhost script]# faascmd list_objects |grep "julia"
0.082(s) elapsed
juliabucket
juliafile
```



Note:

You can use this command with the `grep` command to filter for the files you want, for example, `faascmd list_objects | grep "xxx"`.

Upload original compilation files

The `faascmd upload_object` command is used to upload the original files that are compiled on your local PC to a specified OSS bucket.

Command format

```
faascmd upload_object --object=<newFileNameinOSSBucket> --file= <
your_file_path>/fileNameYouWantToUpload
```

Code example

```
[root@testhost script]# faascmd upload_object --object=juliaOSSFile1 --file=julia_test.tar
juliaOSSFile1
julia_test.tar
0.091(s) elapsed
[root@testhost script]# faascmd upload_object --object=juliaOSSFile2 --file=/opt/dcp1_0/testfile.tar
juliaOSSFile2
/opt/dcp1_0/testfile.tar
0.089(s) elapsed
```



Note:

- No path is needed if the target files are stored in the current directory.
- Locally compiled original files provided by Intel FPGA are in .gbs format and those provided by Xilinx FPGA are compressed as packages in .tar format after script processing.

Download objects from an OSS bucket

The `faascmd get_object` command is used to download a specified object from an OSS bucket.

Command format

```
faascmd get_object --object=<yourObjectName> --file=<your_local_path>/
<yourFileName>
```

Code example

```
[root@ ~]# faascmd get_object --object=juliaOSSFile3 --file=vivadol.log
2018-12-04 10:09:47,342 oss2.api [INFO] 140410558318400 : Init oss bucket, endpoint: https://oss-cn-hangzhou-internal.aliyuncs.com, isCname: False, connect_timeout: None, app_name: , enabled_crc: True
juliaOSSFile3
vivadol.log
2018-12-04 10:09:47,344 oss2.api [INFO] 140410558318400 : Start to get object to file, bucket: juliabucket, key: juliaOSSFile3, file path: vivadol.log
2018-12-04 10:09:47,344 oss2.api [INFO] 140410558318400 : Start to get object, bucket: juliabucket, key: juliaOSSFile3, range: , headers: {}, params: {}
2018-12-04 10:09:47,456 oss2.api [INFO] 140410558318400 : Get object done, req_id: SC05E1EB874F5A9B75E1728B, status code: 200
0.117(s) elapsed
```



Note:

If no path is provided, the objects are downloaded to the current folder by default.

Create FPGA images

The `faascmd create_image` command is used to submit FPGA image creation requests. If the request succeeds, `fpga_imageuuid` is returned.

Command format

```
faascmd create_image --object=<yourObjectName>
--fpgatype=<intel/xilinx> --encrypted=<true/false>
--kmskey=<key/mandatory if encrypted is true>
--shell=<Shell Version/mandatory> --name=<name/optional>
--description=<description/optional> --tags=<tags/optional>
```

Code example

```
[root@testhost script]# faascmd create_image --object=juliasbucket --fpgatype=intel --encrypted=false --shell=V1.1
{"Name": "None", "CreateTime": "Fri Nov 09 2018 11:42:47 GMT+0800 (CST)", "ShellUUID": "V1.1", "Description": "None", "FpgaImageUUID": "
0.250(s) elapsed
```

View FPGA images

The `faascmd list_images` command is used to view information about all the FPGA images you have created.

Command format

```
faascmd list_images
```

Code example

```
[root@testhost script]# faascmd list_images
{
  "FpgaImages": {
    "fpgaImage": [
      {
        "CreateTime": "Fri Nov 09 2018 11:42:47 GMT+0800 (CST)",
        "Description": "None",
        "Encrypted": "false",
        "FpgaImageUUID": " ",
        "Name": "None",
        "ShellUUID": "V1.1",
        "State": "success",
        "Tags": "None",
        "UpdateTime": "Fri Nov 09 2018 11:43:53 GMT+0800 (CST)"
      }
    ]
  }
}
0.076(s) elapsed
```



Note:

A maximum of 10 FPGA images can be reserved for each RAM user account.

Delete FPGA images

The `faascmd delete_image` command is used to delete FPGA images.

Command format

```
faascmd delete_image --imageuuid=<yourImageuuid>
```

Code example

```
[root@testhost script]# faascmd delete_image --imageuuid=  
{\"Status\":200,\"FpgaImageUUID\":\"j\", \"Message\":\"delete succeed!\"}  
0.143(s) elapsed
```

Download FPGA images

The `faascmd download_image` command is used to submit FPGA image download requests.

Command format

```
faascmd download_image --instanceId=<yourInstanceId>  
--fpgauuid=<yourfpgauuid> --fpgatype=<intel/xilinx>  
--imageuuid=<yourImageuuid> --imagetype=<afu>  
--shell=<yourImageShellVersion>
```

Code example

```
faascmd download_image --instanceId=XXXXX --fpgauuid=XXXX --fpgatype=  
intel --imageuuid=XXXX
```

View the FPGA image download status

The `faascmd fpga_status` command is used to view the status of the current FPGA board card and the FPGA image download status.

Command format

```
faascmd fpga_status --fpgauuid=<fpgauuid> --instanceId=<instanceId>
```

Code example

```
[root@testhost script]# faascmd fpga_status --fpgauuid=  
--instanceId=  
{\"shellUUID\":\"V1.0\",\"FpgaImageUUID\":\"\", \"FpgaUUID\":\"\",  
askStatus\":\"invalid\",\"Encrypted\":\"false\"}  
0.310(s) elapsed
```

Publish FPGA images

The `faascmd publish_image` command is used to submit FPGA image publishing requests.

Command format

```
faascmd publish_image --imageuuid=<yourImageuuid> --imageid=<yourFPGAImageid>
```



Note:

- imageuuid is the ID of the image you are going to publish to the cloud marketplace. You can view the image ID by running the `faascmd list_images` command.
- imageid is the FPGA image ID. You can view the ID on the instance details page in the ECS console.

View FPGA instance information

The `faascmd list_instances` command is used to obtain basic information about an FPGA instance, including the instance ID, FPGA board card information, and shell version.

Command format

```
faascmd list_instances --instanceId=<yourInstanceId>
```

Code example

```
[root@testhost script]# faascmd list_instances --instanceId=
{
  "Instances": {
    "instance": [
      {
        "DeviceBDF": "05:00.0",
        "FpgaStatus": "invalid",
        "FpgaType": "intel",
        "FpgaUUID": " ",
        "InstanceId": " ",
        "ShellUUID": "V1.1"
      }
    ]
  }
}
0.275(s) elapsed
```

8.5.5 FAQ

This topic lists common FAQs relating to the faascmd tool and provides corresponding solutions.

FAQ

- What do I do if an error indicating "Name Error:global name'ID' is not defined." is reported?

Cause: faascmd cannot obtain your AccessKeyId or AccessKeySecret.

Solution: Run the `faascmd config` command. Then, the information about the AccessKeyId and AccessKeySecret you have entered will be saved in the `/root/.faascredentials` file.

- What do I do if an error indicating "HTTP Status:403 Error:RoleAccessError. You have no right to assume this role." is reported?

Cause: faascmd cannot obtain information about the role ARN or the obtained ARN does not belong to the same account as the existing AccessKeyId and AccessKeySecret.

Solution: Check whether the following information is contained in the `/root/.faascredentials` file:

```
[FaaScredentials]
accessid=xxxxxxxxxx
accesskey=xxxxxxxxxxxxxxxxxxxxxxxxxx
[Role]
role=acs:ram::1234567890123456:role/xxxxxx
[OSS]
bucket=xxxx
```



Note:

- If the preceding information already exists, check whether the role ARN and the AccessKeyId/AccessKeySecret belong to the same account.
 - If the preceding information does not exist, run `faascmd auth bucket=xxxx` to grant permissions.
- What do I do if an error indicating "HTTP Status: 404 Error: EntityNotExist. Role Error. The specified Role not exists." is reported?
- Cause:** There is no faasrole role in your account.
- Solution:** Log on to the RAM console to check whether a faasrole role exists.

- If no faasrole role exists, run the `faascmd config` and `faascmd auth` commands to create such a role and grant permissions to it.
- If a faasrole role already exists, open a ticket.
- What do I do if an error indicating "SDK.InvalidRegionId. Can not find endpoint to access." is reported?

Cause: faascmd cannot obtain the endpoint address of FaaS.

Solution: Perform the following steps check whether faascmd configurations meet the specified requirements:

- Run the `python -V` command to check whether the Python version is 2.7.x.
- Run the `which python` command to check whether the default installation path of Python is `/usr/bin/python`.
- Run the `cat /usr/lib/python2.7/site-packages/aliyun-sdk-core/__init__.py` command to check whether the aliyun-sdk-core version is 2.11.0 or later.



Note:

If the aliyun-sdk-core version is earlier than 2.11.0, you need to run the `pip install --upgrade aliyun-python-sdk-core` command to upgrade the aliyun-sdk-core to the latest version.

- What do I do if an error indicating "HTTP Status:404 Error:SHELL NOT MATCH The image Shell is not match with fpga Shell! Request ID:D7D1AB1E-8682-4091-8129-C17D54FD10D4" is returned when I download an image?

Cause: The shell versions of the target FPGA image and the specified FPGA do not match.

Solution: Perform the following steps:

- Run the `faascmd list_instances --instance=xxx` command to check the shell version of the current FPGA.
- Run the `faascmd list_images` command to check the shell version of the specified FPGA image.



Note:

- If the two shell versions are different, you need to create an FPGA image whose shell version is the same as that of the FPGA, and then download the image.
- If the two shell versions are consistent, open a ticket.

- What do I do if an error indicating "HTTP Status:503 Error:ANOTHER TASK RUNNING. Another task is running,user is allowed to take this task half an hour Request ID: 5FCB6F75-8572-4840-9BDC-87C57174F26D" is returned when I download an image?

Cause: The FPGA is stuck in operating state due to unexpected failure or interruption of the download request you have submitted.

Solution: Wait for 10 minutes until the download task ends, and then resubmit an image download request.



Note:

If the problem persists, open a ticket.

- What do I do if the image status is failed when I run the `faascmd list_images` command?

Solution: Obtain the compilation logs for troubleshooting by running the following command:

```
faascmd list_objects|grep vivado
faascmd get_object --object=<yourObjectName> --file=<your_local_path>/vivado.log #The path is optional. The compilation logs are downloaded to the current folder by default.
```

Common error codes

faascmd command	API name	Error message	Error description	Error code
Applicable to all commands	Applicable to all APIs	PARAMETER INVALIDATE	The input parameter is incorrect.	400
Applicable to all commands	Applicable to all APIs	InternalError	There is an internal error. Please open a ticket.	500
auth	auth	NoPermisson	You do not have the permission to access a specific open API.	403

faascmd command	API name	Error message	Error description	Error code
create_image	CreateFpgaImage	IMAGE NUMBER EXCEED	There cannot be more than 10 images in the image list . Please delete unnecessary images and try again.	401
		FREQUENCY ERROR	The interval for submitting image requests is 30 minutes.	503
		SHELL NOT SUPPORT	The input shell version is not supported. Please verify that the shell version is correct.	404
		EntityNotExist. RoleError	The current account has no faasrole role.	404
		RoleAccess Error	The role ARN is empty, or the role ARN and the AccessKeyId/AccessKeySecret do not belong to the same account.	403
		InvalidAccessKeyIdError	The AccessKeyId/AccessKeySecret is invalid.	401
		Forbidden. KeyNotFoundError	The specified KMS key cannot be found. Please log on to the KMS console and check whether the input KeyId exists .	503
		AccessDeniedError	The faas admin account is not authorized to access the current bucket.	
		OSS OBJECT NOT FOUND	The specified OSS bucket/object does not exist or is inaccessible.	404
delete_image	DeleteFpgaImage	IMAGE NOT FOUND	The specified FPGA image cannot be found.	400
list_instances	DescribeFpgaInstances	NOT AUTHORIZED	The specified instance does not exist or does not belong to the current account.	401
		RoleAccess Error	The role ARN is empty, or the role ARN and the AccessKeyId/AccessKeySecret do not belong to the same account.	403

faascmd command	API name	Error message	Error description	Error code
		INSTANCE INVALIDATE	The specified instance is not an FPGA instance. If the specified instance is an FPGA instance, please open a ticket.	404
fpga_status	DescribeLoadTaskStatus	NOT AUTHORIZED	The specified instanceId cannot be found. Please check the input parameter.	401
		FPGA NOT FOUND	The specified fpgauid cannot be found. Please check the input parameter.	404
download_image	LoadFpgaImage	ANOTHER TASK RUNNING	The image download task you submitted is still in operating state.	503
		IMAGE ACCESS ERROR	The specified image does not belong to the current account.	401
		YOU HAVE NO ACCESS TO THIS INSTANCE	The specified instance does not belong to the current account.	401
		IMAGE NOT FOUND	The specified FPGA image cannot be found.	404
		FPGA NOT FOUND	The specified FPGA cannot be found.	404
		SHELL NOT MATCH	The image and the specified FPGA do not match in shell version.	404
		RoleAccess Error	The role ARN is empty, or the role ARN and the AccessKeyId/AccessKeySecret do not belong to the same cloud account.	403
		Image not in success state	The specified image is not in success state. Only images in success state can be downloaded.	404
publish_image	PublishFpgaImage	FPGA IMAGE STATE ERROR	The specified image is not in success state.	404

faascmd command	API name	Error message	Error description	Error code
		FPGA IMAGE NOT FOUND	The specified image cannot be found or does not belong to the current account.	404

9 Access other Cloud Product APIs by the Instance RAM Role

Previously, applications deployed on an ECS Instance usually needed to use AccessKey ID and AccessKey Secret (AK) to access APIs of other Alibaba Cloud products. AK is the key to accessing Alibaba Cloud APIs and has all of the permissions of the corresponding accounts. To help applications manage the AK, you have to save AK in the configuration files of the application or save it in an ECS instance by using other methods, which makes it more complicated to manage the AK and reduces its confidentiality. What's more, if you need concurrent deployment across regions, the AK is diffused along with the images or instances created by the image, which makes you have to update and re-deploy the instances and images one by one when changing the AK.

Now with the help of the instance [RAM role](#), you can assign a RAM role to an ECS instance. The applications on the instance can then access APIs of other cloud products with the STS credential. The STS credential is automatically generated and updated by the system, and the applications can use the specified [meta data](#) URL to obtain the STS credential without special management. Meanwhile, you can modify the RAM role and the authorization policy to grant different or identical access permissions to an instance to different Alibaba Cloud products.

This article introduces how to create an ECS instance that plays a RAM role and how to enable applications on the ECS instance to access other Alibaba Cloud products with the STS credential.



Note:

To make it easy for you to get started with the example in this article, all of the operations in the document are done in [OpenAPI Explorer](#). OpenAPI Explorer obtains the temporary AK of the current account through the logged user information, and initiates online resource operation to the current account. Please execute operations carefully. Creating an instance will incur charges. Please release the instance soon after completing the operation.

Procedure

To enable python on an instance to access an OSS bucket under the same account by using the instance RAM role, follow these steps:

Step 1. Create a RAM role and attach it to an authorization policy.

Step 2. Create an ECS instance playing the RAM role to create.

Step 3. Within the instance, access the metadata URL to obtain the STS credential.

Step 4. Use Python to access OSS using the STS credential.

Step 1. Create a RAM role and attach it to an authorization policy

Use the `CreateRole` API to

1. create a RAM role. The required request parameters are:

- **RoleName:** Specify a name for the role. *EcsRamRoleTest* is used in this example.
- **AssumeRolePolicyDocument:** Specify a policy as follows, which indicates that the role to be created is a service role and an Alibaba Cloud product (ECS in this example) is assigned to play this role.

```
{
  "Statement": [
    {
      "Action": "sts:AssumeRole",
      "Effect": "Allow",
      "Principal": {
        "Service ": [
          "ecs.aliyuncs.com"
        ]
      }
    }
  ],
  "Version": "1"
}
```

2. Use the `CreatePolicy` API to create an authorization policy. The required request parameters are:

- **PolicyName:** Specify a name for the authorization policy. *EcsRamRolePolicyTest* is used in this example.
- **PolicyDocument:** Specify a policy as follows, which indicates that the role has OSS read-only permission.

```
{
  "Statement": [
    {
      "Action": [
        "oss:Get*",

```

```
"oss:List*"
],
"Effect": "Allow",
"Resource " : "*"
}
],
"Version": "1"
}
```

3. Use the `AttachPolicyToRole` API to attach the authorization policy to the role. The required request parameters are:

- **PolicyType:** Set it to *Custom*.
- **PolicyName:** Use the policy name specified in step 2. Use *EcsRamRolePolicyTest* in this example.
- **RoleName:** Use the role name specified in step 1. Use *EcsRamRoleTest* in this example.

Step 2. You can use either method to create an ECS instance playing the RAM role:

Attach a RAM role to an existing VPC-Connected ECS instance.

- Create a VPC-Connected ECS instance with the RAM role
- Attach a RAM role to an existing VPC-Connected ECS instance

Create a VPC-Connected ECS instance with the RAM role

Use the `AttachInstanceRamRole` API to attach a RAM role to an existing VPC-Connected ECS instance. The parameters are as follows:

- **RegionId:** The ID of the region where the instance is located.
- **RamRoleName:** The name of a RAM role. In this example, *EcsRamRoleTest* is used. In this example, *EcsRamRoleTest*.
- **InstanceIds:** The IDs of VPC-Connected ECS instances that you want to attach the RAM role to, in the format of ["i-bXXXXXXXX"] for one instance, or ["i-bXXXXXX" , "i-cXXXXXX" , ["i-bXXXXXXXX"]] for multiple instances.

Create a VPC-Connected ECS instance with the RAM role

You must have a VPC network before creating an ECS instance with the RAM role.

1. To create a VPC-Connected ECS instance with the RAM role, follow these steps: Use the `CreateInstance` API to create an ECS instance. The required request parameters are:

- **RegionId**:The region of the instance. In this example, `cn-hangzhou` is used. In this example, `cn-hangzhou` is used.
- **ImageId**:The image of the instance. In this example, `centos_7_03_64_40G_alibase_20170503.vhd` is used. In this example, `centos_7_03_64_40G_alibase_20170503.vhd` is used.
- **InstanceType**:The type of the instance. In this example, `ecs.xn4.small` is used.
- **VSwitchId**:The virtual switch of the VPC network where the instance is located. Because the instance RAM role only supports VPC network, **VSwitchId** is required.
- **RamRoleName**:The name of RAM Role. In this example, `EcsRamRoleTest` is used.

If you want to authorize a sub account to create an ECS instance playing the specified RAM role, besides the permission to create an ECS instance, the sub account must have the `PassRole` permission. Therefore, you must customize an authorization policy as follows and attach it to the sub account. If the action is creating an ECS instance only, set [ECS RAM Action] to `ecs:CreateInstance`. If you want to grant all ECS action permissions to the sub account, set [ECS RAM Action] to `ecs:*`.

```
{
  "Statement": [
    {
      "ecs: [ECS RAM Action]",
      "Resource": "*",
      "Effect": "Allow"
    },
    {
      "Action": "ram:PassRole",
      "Resource": "*",
      "Effect": "Allow"
    }
  ],
  "Version": "1"
}
```

2. Set the password and start the instance.
3. Set the ECS instance to access the Internet by using API or in the ECS console.

Step 3: Access the metadata URL within the instance to obtain the STS credential

To obtain the STS credential of the instance, follow these steps:



Note:

A new STS credential is generated 30 minutes before the current one expires. Both STS credentials can be used during this period of time.

1. [Connect to the instance](#).
2. Access the following URL to obtain the STS credential. `http://100.100.100.200/latest/meta-data/ram/security-credentials/EcsRamRoleTest` The last part of the URL is the RAM role name, which must be replaced with the one you create. The last part of the path is the RAM role name which should be replaced by one you create.

**Note:**

In this example, use the curly command to access the above curl. In this example, we run the curl command to access the URL. If you are using a Windows ECS instance, see [Use metadata of an instance](#) in ECS the User Guide to obtain the STS credential.

The return parameters are as follows.

```
[root@local ~]# curl http://100.100.100.200/latest/meta-data/ram/
security-credentials/EcsRamRoleTest
{
  "AccessKeyId" : "XXXXXXXXXX",
  "AccessKeySecret" : "XXXXXXXXXX",
  "Expiration" : "2017-06-09T09:17:19Z",
  "SecurityToken" : "CAIXXXXXXXXXXXXwmBkIeCTkyI+",
  "LastUpdated" : "2017-10-31T23:20:01Z",
  "Code" : "Success"
}
```

Step 4: Use Python SDK to access OSS with the STS credential

In this example, with the STS credential, we use Python to list 10 files in an OSS bucket that is in the same region with the instance.

Prerequisites

You have remotely connected to the ECS instance.

Python has been installed on the ECS instance. If you are using a Linux ECS instance, pip must be installed.

A bucket has been created in the region of the instance, and the bucket name and the Endpoint have been acquired. In this example, the bucket name is `ramroletest`, and the endpoint is `oss-cn-hangzhou.aliyuncs.com`.

Procedure

To use Python to access the OSS bucket, follow these steps:

1. Run the command `pip install oss2` to install OSS Python SDK.

2. Run the following commands to test, of which:

- The three parameters in `oss2. StsAuth` correspond respectively to `AccessKeyId`, `AccessKeySecret` and `SecurityToken` returned by the above URL.
- The last two parameters in `oss2. Bucket` are the `bucketcodeph` name and the endpoint.

```
import oss2
from itertools import islice
auth = oss2. StsAuth(<AccessKeyId>, <AccessKeySecret>, <SecurityToken>)
bucket = oss2. bucket = oss2.Bucket(auth, <your Endpoint>, <your Bucket name>)
for b in islice(oss2. ObjectIterator(bucket), 10):
    print(b.key)
```

Output results are as follows:

```
[root@local ~]# python
Python 2.7.5 (default, Nov 6 2016, 00:28:07)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-11)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
>>> import oss2
>>> from itertools import islice
>>> auth = oss2. StsAuth("STS.J8XXXXXXXXXX4", "9PjfXXXXXXXXXBf2XAW",
"CAIXXXXXXXXXXXwmBkleCTkyI+")
>>> bucket = oss2. Bucket(auth, "oss-cn-hangzhou.aliyuncs.com", "
ramroletest")
>>> for b in islice(oss2. ObjectIterator(bucket), 10):
...     print(b.key)
...
ramroletest.txt
test.sh
```

10 Shrink disk

Currently, Elastic Compute Service (ECS) does not support system disk or data disk shrink. If you want to shrink your disk volumes, try [Alibaba Cloud Migration Tool](#) instead.

Though Cloud Migration Tool is designed to balance the cloud-based and offline workloads of Alibaba Cloud users, you can use it to shrink ECS disk volumes.

The tool creates a custom image based on your ECS instance. During this process, it re-specifies the size of the disk to shrink it. Apart from replacing the target object with an ECS instance, the tools for cloud migration and disk volume shrinking are identical, in terms of [both operation and limitations](#). Because the ECS instance is already virtual, it is more convenient to use and the chances of reporting errors is reduced.

However, using this tool may change some attributes of the ECS instance. For example, instance ID (`InstanceId`) and public IP. If your instance is a [VPC-Connected](#) instance, you can reserve the public IP address by [converting public IP address to EIP address](#). We recommend that users using [Alibaba Cloud Elastic IP \(EIP\)](#) and users with less dependency on public IP use this approach to shrink the disk.

Prerequisites

- When the disk is mounted on a Linux instance, you must first install rsync, a remote data synchronization tool.
 - CentOS Instance: Run `yum install rsync -y`.
 - Ubuntu Instance: Run `apt-get install rsync -y`.
 - Debian Instance: Run `apt-get install rsync -y`.
 - Other distributions: Please visit the official website to find the relevant installation documents.
- You must [create an AccessKey](#) in the console first, which is used to output it into the configuration file `user_config.json`.



Note:

To prevent data leakage due to excessive permissions for AccessKey, we recommend that you [create a RAM sub-account](#) and use this account to [create an AccessKey](#).

- For other prerequisites and limitations, see [migrate to Alibaba Cloud by using Cloud Migration Tool](#).

Procedure

1. [Connect](#) to the target ECS instance by using the administrator/root account.
2. [Download](#) the Alibaba Cloud Migration Tool zip file.
3. Unzip the Cloud Migration Tool. Enter the corresponding operating system and version of the client file directory to find the configuration file `user_config.json`.
4. See customize [user_config.json](#) to complete the configuration.

See the following figure for the configuration file in a Linux instance.

```
"access_id": "",
"secret_key": "",
"region_id": "",
"image_name": "",
"system_disk_size": ,
"platform": "",
"architecture": "",
"data_disks": [],
"bandwidth_limit": 0
```

The most important parameters to configure for shrinking a disk volume are as follows:

- `system_disk_size`: Set this parameter to the expected system disk size in GB. The value cannot be less than the actual size of the system disk.
- `data_disks`: Set this parameter to the expected data disk size in GB. The value cannot be less than the actual size of the data disk.



Note:

- When a Linux instance comes with a data disk, the `data_disks` parameter is required even if you do not want to shrink the data disk volume. If it is not configured, Cloud Migration Tool copies data from the data disk to the system disk by default.
- When a Windows instance comes with a data disk, the `data_disks` parameter is optional if you do not want to shrink the size of the data disk.

5. Run the program `go2aliyun_client.exe`:
 - Windows instance: Right-click `go2aliyun_client.exe` and select Run as administrator.
 - Linux instance:

- a. Run `chmod +x go2aliyun_client` to give the client executable permissions.
- b. Run `./ go2aliyun_client` to run the client.

6. Wait for the running results:

- If `Goto Aliyun Finished!` is displayed, go to the [ECS console](#) and check the custom image after shrinking. If the custom image has been generated, you can release the original instance and use the custom image to [create an ECS instance](#). After you create a new instance, the disk volume shrinking process is complete.
- If `Goto Aliyun Not Finished!` is displayed, check the log files in the same directory for [troubleshooting](#). After fixing any problems, run Cloud Migration Tool again to resume volume shrinking. The tool continues the most recent migration progress and does not start over.

References

- For a detailed introduction to Cloud Migration Tool, see [what is Alibaba Cloud Migration Tool](#).
- For instructions on how to use Cloud Migration Tool, see [migrate to Alibaba Cloud by using Cloud Migration Tool](#).

11 Terraform

11.1 What is Terraform?

Terraform is an open source tool for securely and efficiently provisioning and managing cloud infrastructure.

[HashiCorp Terraform](#) is an automated IT infrastructure orchestration tool that can use codes to manage and maintain IT resources. The Command Line Interface (CLI) of Terraform provides a simple mechanism, which is used for deploying and versioning configuration files on Alibaba Cloud or any other supported cloud.

Terraform writes the infrastructure, for example, virtual machines, storage accounts, and network interfaces in the configuration file that describes the cloud resource topology. The Command Line Interface (CLI) of Terraform provides a simple mechanism, which is used for deploying and versioning configuration files on Alibaba Cloud or any other supported cloud.

Terraform is a highly scalable tool that supports new infrastructure through providers. You can use Terraform to create, modify, or delete multiple resources, such as ECS, VPC, RDS, and SLB.

Benefits

- Multiple-cloud infrastructure deployment

Terraform applies to multi-cloud scenarios, where similar infrastructure is deployed on Alibaba Cloud, other cloud providers, or local data centers.

Developers can use the same tool and configuration file to simultaneously manage the resources of different cloud providers.

- Automated infrastructure management

Terraform can create configuration file templates to define, provision, and configure ECS resources in a repeatable and predictable manner, reducing deployment and management errors resulting from human intervention. In addition, Terraform can deploy the same template multiple times to create the same development, test, and production environment.

- Infrastructure as code

With Terraform, you can use codes to manage and maintain resources. It allows you to store the infrastructure status, so that you can track the changes in different components of the system (infrastructure as code) and share these configurations with others.

- Reduced development costs

You can reduce costs by creating on-demand development and deployment environments. In addition, you can evaluate such environments before making system changes.

Application scenarios

Terraform is a well-proven and open-source automated operation and maintenance tool for managing cloud infrastructure, creating images, and supporting multi-cloud business scenarios.

For the application scenarios of Terraform, see [Terraform details](#).

Use Terraform

Terraform allows you to use a [simple template language](#) to easily define, preview, and deploy cloud infrastructure on Alibaba Cloud. The steps for Terraform to provision resources in ECS are described as follows:

1. Install Terraform.
2. Configure Terraform.
3. Use Terraform to create one or more ECS instances.

More information

- [Terraform Alibaba provider](#)
- [Terraform Alibaba github](#)
- [Terraform Registry Alibaba Modules](#)

11.2 Install and configure Terraform

This article describes how to install Terraform.

Procedure

1. Go to the [Terraform website](#) to download the package for your operating system.
2. Decompress the package to `/usr/local/bin`.

If you decompress the executable file to another directory, define a global path for the file as follows:

- **Linux:** See [How to define a global path on Linux](#).
- **Windows:** See [How to define a global path on Windows](#).
- **Mac:** See [How to define a global path on Mac](#).

3. Run the `terraform` command to verify the path.

A list of available Terraform options is displayed as follows, indicating that the installation is completed.

```
username:~$ terraform
Usage: terraform [-version] [-help] <command> [args]
```

4. For higher flexibility and security of rights management, it is recommended that you create and authorize a RAM user.

- a. Log on to the [RAM console](#).
- b. Create a RAM user named *Terraform* and create an AccessKey for the user. For detailed steps, see [Create a RAM user](#).
- c. Authorize the RAM user. In this example, the user *Terraform* is granted the `AliyunECSFullAccess` and `AliyunVPCFullAccess` permissions. For detailed steps, see [Attach policies to a RAM user](#).

5. Create an environment variable to store identity authentication information.

```
export ALICLOUD_ACCESS_KEY="LTAIUrZCw3*****"
export ALICLOUD_SECRET_KEY="zfwWAMWIAiooj14GQ2*****"
export ALICLOUD_REGION="cn-beijing"
```

11.3 Create an ECS instance

This article describes how to create an ECS instance by using Terraform.

Procedure

1. Create a VPC and a switch.
 - a. Create the *terraform.tf* file, enter the following, and save it to the current execution directory.

```
resource "alicloud_vpc" "vpc" {
  name      = "tf_test_foo"
  cidr_block = "172.16.0.0/12"
}
```

```
resource "alicloud_vswitch" "vsw" {
  vpc_id          = "${alicloud_vpc.vpc.id}"
  cidr_block      = "172.16.0.0/21"
  availability_zone = "cn-beijing-b"
}
```

b. Run `terraform apply` to start the creation.

c. Run `terraform show` to view the VPC and VSwitch that have been created.

You can also log on to the VPC console to view the attributes of the VPC and VSwitch.

2. Create a security group and apply it to the VPC created in the previous step.

a. In the file `terraform.tf`, add the following:

```
resource "alicloud_security_group" "default" {
  name = "default"
  vpc_id = "${alicloud_vpc.vpc.id}"
}

resource "alicloud_security_group_rule" "allow_all_tcp" {
  type          = "ingress"
  ip_protocol   = "tcp"
  nic_type      = "internet"
  policy        = "accept"
  port_range    = "1/65535"
  priority      = 1
  security_group_id = "${alicloud_security_group.default.id}"
  cidr_ip       = "0.0.0.0/0"
}
```

b. Run `terraform apply` to start the creation.

c. Run `terraform show` to view the security group and security group rules that have been created.

You can also log on to the ECS console to view the security group and security group rules.

3. Creates an ECS instance.

a. In the file `terraform.tf`, add the following:

```
resource "alicloud_instance" "instance" {
  # cn-beijing
  availability_zone = "cn-beijing-b"
  security_groups = ["${alicloud_security_group.default.*.id}"]

  # series III
  instance_type      = "ecs.n2.small"
  system_disk_category = "cloud_efficiency"
  image_id           = "ubuntu_140405_64_40G_cloudinit_20161115.vhd"
  instance_name      = "test_foo"
  vswitch_id         = "${alicloud_vswitch.vsw.id}"
  internet_max_bandwidth_out = 10
}
```

```
password = "<replace_with_your_password>"
}
```

**Note:**

- In the above example, `internet_max_bandwidth_out = 10` is specified. Therefore, the instance is assigned a public IP automatically.
- For a detailed explanation of the parameters, see the [Alibaba Cloud parameter descriptions](#).

- Run `terraform apply` to start the creation.
- Run `terraform show` to view the ECS instance that has been created.
- Run `ssh root@<publicip>` and enter the password to access the ECS instance.

```
provider "alicloud" {}

resource "alicloud_vpc" "vpc" {
  name      = "tf_test_foo"
  cidr_block = "172.16.0.0/12"
}

resource "alicloud_vswitch" "vsw" {
  vpc_id      = "${alicloud_vpc.vpc.id}"
  cidr_block  = "172.16.0.0/21"
  availability_zone = "cn-beijing-b"
}

resource "alicloud_security_group" "default" {
  name = "default"
  vpc_id = "${alicloud_vpc.vpc.id}"
}

resource "alicloud_instance" "instance" {
  # cn-beijing
  availability_zone = "cn-beijing-b"
  security_groups = ["${alicloud_security_group.default. *.id}"]

  # series III
  instance_type      = "ecs.n2.small"
  system_disk_category = "cloud_efficiency"
  image_id           = "ubuntu_140405_64_40G_cloudinit_20161115.vhd"
  instance_name      = "test_foo"
  vswitch_id         = "${alicloud_vswitch.vsw.id}"
  internet_max_bandwidth_out = 10
}

resource "alicloud_security_group_rule" "allow_all_tcp" {
  type      = "ingress"
  ip_protocol = "tcp"
  nic_type  = "intranet"
}
```

```
policy          = "accept"
port_range      = "1/65535"
priority        = 1
security_group_id = "${alicloud_security_group.default.id}"
cidr_ip         = "0.0.0.0/0"
}
```

11.4 Create multiple ECS instances

This article describes how to create multiple ECS instances in batches by using Terraform.

Procedure

1. Create a VPC and a VSwitch.

- a. Create the `terraform.tf` file, enter the following, and save it to the current execution directory.

```
resource "alicloud_vpc" "vpc" {
  name      = "tf_test_foo"
  cidr_block = "172.16.0.0/12"
}

resource "alicloud_vswitch" "vsw" {
  vpc_id      = "${alicloud_vpc.vpc.id}"
  cidr_block   = "172.16.0.0/21"
  availability_zone = "cn-beijing-b"
}
```

- b. Run `terraform apply` to start the creation.
- c. Run `terraform show` to view the VPC and VSwitch that have been created.

You can also log on to the VPC console to view the attributes of the VPC and VSwitch.

2. Create a security group and apply it to the VPC created in the previous step.

- a. In the file `terraform.tf`, add the following:

```
resource "alicloud_security_group" "default" {
  name = "default"
  vpc_id = "${alicloud_vpc.vpc.id}"
}

resource "alicloud_security_group_rule" "allow_all_tcp" {
  type          = "ingress"
  ip_protocol    = "tcp"
  nic_type       = "internet"
  policy         = "accept"
  port_range     = "1/65535"
  priority       = 1
  security_group_id = "${alicloud_security_group.default.id}"
}
```

```
cidr_ip      = "0.0.0.0/0"
}
```

- b. Run `terraform apply` to start the creation.
- c. Run `terraform show` to view the security group and security group rules that have been created.

You can also log on to the ECS console to view the security group and security group rules.

3. Use the Module to create multiple ECS instances. In this example, three ECS instances are created.

- a. In the file `terraform.tf`, add the following:

```
module "tf-instances" {
  source = "alibaba/ecs-instance/alibabacloud"
  vswitch_id = "${alicloud_vswitch.vsw.id}"
  group_ids = ["${alicloud_security_group.default.*.id}"]
  availability_zone = "cn-beijing-b"
  disk_category = "cloud_ssd"
  disk_name = "my_module_disk"
  disk_size = "50"
  number_of_disks = 7

  instance_name = "my_module_instances_"
  host_name = "sample"
  internet_charge_type = "PayByTraffic"
  number_of_instances = "3"
  password="User@123"
}
```



Note:

- In the above example, `internet_max_bandwidth_out = 10` is specified. Therefore, the instances are assigned public IP addresses automatically.
- For a detailed explanation of the parameters, see the [Parameter descriptions](#).

- b. Run `terraform apply` to start the creation.
- c. Run `terraform show` to view the ECS instances that have been created.
- d. Run `ssh root@<publicip>` and enter the password to access the ECS instances.

```
provider "alicloud" {}

resource "alicloud_vpc" "vpc" {
  name      = "tf_test_foo"
  cidr_block = "172.16.0.0/12"
}

resource "alicloud_vswitch" "vsw" {
  vpc_id = "${alicloud_vpc.vpc.id}"
}
```

```
cidr_block      = "172.16.0.0/21"
availability_zone = "cn-beijing-b"
}

resource "alicloud_security_group" "default" {
  name = "default"
  vpc_id = "${alicloud_vpc.vpc.id}"
}

resource "alicloud_security_group_rule" "allow_all_tcp" {
  type          = "ingress"
  ip_protocol   = "tcp"
  nic_type      = "intranet"
  policy        = "accept"
  port_range    = "1/65535"
  priority      = 1
  security_group_id = "${alicloud_security_group.default.id}"
  cidr_ip       = "0.0.0.0/0"
}

module "tf-instances" {
  source = "alibaba/ecs-instance/alicloud"
  vswitch_id = "${alicloud_vswitch.vsw.id}"
  group_ids = ["${alicloud_security_group.default.*.id}"]
  availability_zone = "cn-beijing-b"
  disk_category = "cloud_ssd"
  disk_name = "my_module_disk"
  disk_size = "50"
  number_of_disks = 7

  instance_name = "my_module_instances_"
  host_name = "sample"
  internet_charge_type = "PayByTraffic"
  number_of_instances = "3"
  password="User@123"
}
```

11.5 Deploy a Web cluster

When you deploy a website or application, you need to deploy a series of nodes, and allow them to scale up or down automatically according to the number of visits or resource usage. SLB distributes requests to respective nodes. This article describes how to deploy a Web cluster by using Terraform.

Context

In this example, the entire application is deployed in one zone, and the "Hello, World" web page can be accessed only through port 8080.

Procedure

1. Create a VPC and a VSwitch.

- a. Create the `terraform.tf` file, enter the following, and save it to the current execution directory.

```
resource "alicloud_vpc" "vpc" {
  name      = "tf_test_foo"
  cidr_block = "172.16.0.0/12"
}

resource "alicloud_vswitch" "vsw" {
  vpc_id          = "${alicloud_vpc.vpc.id}"
  cidr_block      = "172.16.0.0/21"
  availability_zone = "cn-beijing-b"
}
```

- b. Run `terraform apply` to start the creation.
- c. Run `terraform show` to view the VPC and VSwitch that have been created.

You can also log on to the VPC console to view the attributes of the VPC and VSwitch.

2. Create a security group and apply it to the VPC created in the previous step.

- a. In the file `terraform.tf`, add the following:

```
resource "alicloud_security_group" "default" {
  name = "default"
  vpc_id = "${alicloud_vpc.vpc.id}"
}

resource "alicloud_security_group_rule" "allow_all_tcp" {
  type          = "ingress"
  ip_protocol   = "tcp"
  nic_type      = "internet"
  policy        = "accept"
  port_range    = "1/65535"
  priority      = 1
  security_group_id = "${alicloud_security_group.default.id}"
  cidr_ip       = "0.0.0.0/0"
}
```

- b. Run `terraform apply` to start the creation.
- c. Run `terraform show` to view the security group and security group rules that have been created.

You can also log on to the ECS console to view the security group and security group rules.

3. Create a Server Load Balancer (SLB) instance and assign a public IP address to it. In this example, the SLB instance is configured with a mapping from front end port 80 to back end port 8080. In addition, it is configured to output the public IP address for subsequent testing.

a. Create the file `slb.tf` and add the following.

```
resource "alicloud_slb" "slb" {
  name          = "test-slb-tf"
  vswitch_id    = "${alicloud_vswitch.vsw.id}"
  internet      = true
}
resource "alicloud_slb_listener" "http" {
  load_balancer_id = "${alicloud_slb.slb.id}"
  backend_port     = 8080
  frontend_port    = 80
  bandwidth       = 10
  protocol         = "http"
  sticky_session   = "on"
  sticky_session_type = "insert"
  cookie           = "testslblistenercookie"
  cookie_timeout   = 86400
  health_check     = "on"
  health_check_type = "http"
  health_check_connect_port = 8080
}

output "slb_public_ip" {
  value = "${alicloud_slb.slb.address}"
}
```

b. Run `terraform apply` to start the creation.

c. Run `terraform show` to view the SLB instance that has been created.

You can also log on to the SLB console to view the new SLB instance.

4. Creates an Auto Scaling solution.

In this example, the following resources are created:

- **Scaling group:** specifies the minimum number of ECS instances as 2 in the template, and the maximum number as 10. Meanwhile, bind the scaling group to the newly created SLB instance. Due to the configuration requirements of the scaling group, SLB must have a listener configured accordingly. As a result, the order of deployment is specified with the `depends_on` attribute in the template.
- **Scaling group configuration:** specifies the specific configuration of the ECS instance in the template. Generate a "hello World" web page in the initialization configuration (user-data), and provide services on port 8080. To simplify operations, in this example, the virtual machine is assigned a public IP address, and `force_delete=true` is set for subsequent deletion of the environment.
- **Scaling rules:** define specific scaling rules.

a. Create the file `ess.tf` and add the following to it:

```
resource "alicloud_ess_scaling_group" "scaling" {
```

```

    min_size = 2
    max_size = 10
    scaling_group_name = "tf-scaling"
    vswitch_ids=["${alicloud_vswitch.vsw. *.id}"]
    loadbalancer_ids = ["${alicloud_slb.slb. *.id}"]
    removal_policies = ["OldestInstance", "NewestInstance"]
    depends_on = ["alicloud_slb_listener.http"]
  }

  resource "alicloud_ess_scaling_configuration" "config" {
    scaling_group_id = "${alicloud_ess_scaling_group.scaling.id}"
    image_id = "ubuntu_140405_64_40G_cloudinit_20161115.vhd"
    instance_type = "ecs.n2.small"
    security_group_id = "${alicloud_security_group.default.id}"
    active=true
    enable=true
    user_data = "#! /bin/bash\nnecho \"Hello, World\" > index.html\n\nnohup busybox httpd -f -p 8080&"
    internet_max_bandwidth_in=10
    internet_max_bandwidth_out= 10
    internet_charge_type = "PayByTraffic"
    force_delete= true
  }

  resource "alicloud_ess_scaling_rule" "rule" {
    scaling_group_id = "${alicloud_ess_scaling_group.scaling.id}"
    adjustment_type = "TotalCapacity"
    adjustment_value = 2
    cooldown = 60
  }

```

b. Run `terraform apply` to start the creation.

After it is created successfully, the public IP address of the SLB is generated.

c. In about two minutes, Auto Scaling will automatically create the ECS instance.

d. Enter the command `curl http://<slb public ip>` for verification.

If you see `Hello, World`, you have successfully accessed the web page provided by the ECS instance through an SLB instance.

5. Run `terraform destroy` to delete the test environment. Once confirmed, the entire deployed environment will be deleted.

Terraform makes it easy to remove and redeploy an environment. If you want to redeploy, just run `terraform apply`.