

Server Load Balancer

Quick Start

Quick Start

This section provides a complete tutorial on using Server Load Balancer. An Internet-facing Server Load Balancer instance will be created to distribute the received HTTP requests to the backend servers.

Note: Before creating a Server Load Balancer instance, you need to plan and design your load balancing service, such as the instance type, instance region, and more. For details, see [Plan and prepare](#).

The tutorial includes the following tasks:

Create ECS instances.

Server Load Balancer is a complementary service for ECS multi-machine solutions, and must be used in conjunction with ECS. In this tutorial, two ECS instances are created to process the distributed traffic.

Install web pages.

Create required applications on the ECS instances. In this tutorial, a static web page is created to test the load balancing service.

Create a Server Load Balancer instance.

A Server Load Balancer instance is a running entity of Server Load Balancer. In this tutorial, an Internet-facing Server Load Balancer instance is created.

Configure the Server Load Balancer instance.

After creating a Server Load Balancer instance, you have to add at least one listener, and multiple ECS instances as backend servers. In this tutorial, a TCP listener is added, and the ECS instances created in task 1 are used as backend servers.

Delete the Server Load Balancer instance.

If you no longer need Server Load Balancer, delete it to avoid additional charges.

Plan the region of the Server Load Balancer instance

Alibaba Cloud provides the Server Load Balancer service in various regions.

To provide more stable and reliable load balancing services, Server Load Balancer has deployed multiple zones in most regions for better disaster tolerance. Additionally, to improve high availability in different regions, you can deploy the Server Load Balancer instances in multiple regions and use DNS to resolve the domain name to the IP addresses of the Server Load Balancer instances.

When selecting the region, note the following considerations:

To reduce latency and increase the download speed, we recommend choosing a region that is physically closest to where your customers are located.

Server Load Balancer does not support the cross-region development. Ensure that the region is the same for the Server Load Balancer and the backend ECS instances.

Plan the instance type (Internet or intranet)

Choose the instance type according to your business. After you create a Server Load Balancer instance, a private or public IP is allocated. You can resolve a domain name to the IP to provide services.

The Internet Server Load Balancer instance only has a public IP, which is accessible from the Internet.

If you choose the Internet type, you need to consider the billing method:

By traffic: Suitable for an application with obvious traffic changes.

By bandwidth: Suitable for an application with relatively stable bandwidth.

The intranet Server Load Balancer instance only has a private IP, which is accessible only from classic network or VPC, and cannot be accessed by the Internet.

Plan the listening protocol

Server Load Balancer supports layer-4 (TCP and UDP) and layer-7 (HTTP and HTTPS) listening.

Layer-4 listener distributes connection requests directly to the backend servers without modifying the HTTP headers. After the request arrives at the listener, the Server Load Balancer server uses the backend protocol port configured in the listener to create a TCP connection with the backend ECS.

Layer-7 listener is an implementation of reverse proxy. After the request arrives at the listener, the Server Load Balancer server uses the TCP connection to transfer the data packets instead of transferring the data packets directly to the backend ECS.

Prepare the backend servers

Before using Server Load Balancer, you have to create ECS instances and build corresponding applications. Then, add the ECS instances to a Server Load Balancer instance as the backend servers to process the distributed requests.

ECS region

Ensure the region is the same for the ECS instance and Server Load Balancer instance. Also, we recommend deploying the ECS instances in different zones to improve availability.

ECS configurations

Additional configuration is not required after deploying the applications. However, if you are going to create a layer-4 listener, and the ECS instances use the Linux operating system, ensure the values of the following parameters in the `net.ipv4.conf` file are 0:

```
net.ipv4.conf.default.rp_filter = 0
net.ipv4.conf.all.rp_filter = 0
net.ipv4.conf.eth0.rp_filter = 0
```

ECS amount

There is no restriction on the number of ECS instances added to a Server Load Balancer instance. To improve the stability and efficiency of the service, we recommend adding the ECS instances responsible for different tasks or provide different services to different Server Load Balancer instances.

Before using Server Load Balancer, you have to create at least two ECS instances and build corresponding applications. Then, add them to the Server Load Balancer instance to process the distributed client requests.

Follow the instruction in this document to create two ECS instances, ECS01 and ECS02.

Procedure

Log on to the ECS console.

In the left-side navigation pane, click **Instances** and then click **Create Instance**.

On the buy page, configure the ECS instance.

The following are ECS settings used in this tutorial. For more information, see [Create Linux ECS instances](#).

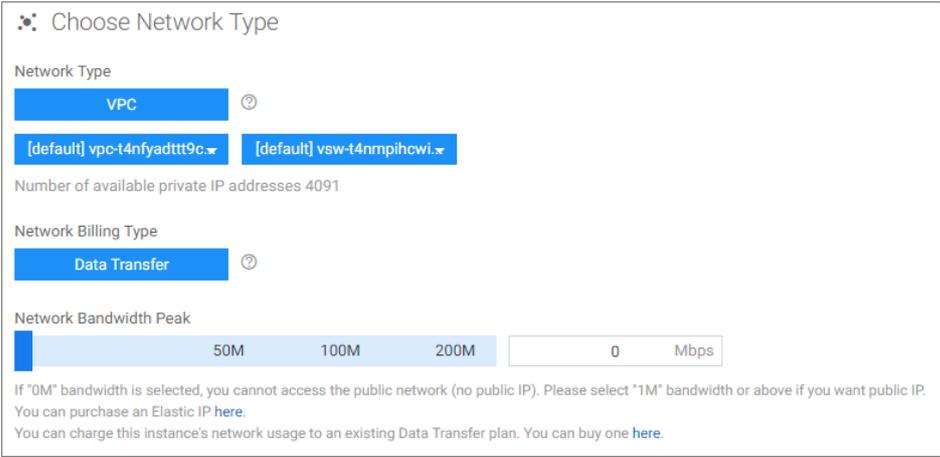
Region: In this tutorial, select **China East 1**.

Note: Server Load Balancer does not support cross-region deployment. The region must be the same for the Server Load Balancer instance and the ECS instance.

Network Type: In this tutorial, select **VPC**. Use the default VPC and VSwitch.

Operating System: In this tutorial, select **Ubuntu 16.04 64 bit**.

Number of Instances: In this tutorial, select **2**. The system will simultaneously create two ECS instances with identical settings.



Choose Network Type

Network Type

VPC ⓘ

[default] vpc-t4nfyadttt9c ⌵ [default] vsw-t4nmpihcwi ⌵

Number of available private IP addresses 4091

Network Billing Type

Data Transfer ⓘ

Network Bandwidth Peak

50M 100M 200M 0 Mbps

If "0M" bandwidth is selected, you cannot access the public network (no public IP). Please select "1M" bandwidth or above if you want public IP. You can purchase an Elastic IP [here](#). You can charge this instance's network usage to an existing Data Transfer plan. You can buy one [here](#).

Click **Buy Now** and complete payment.

Go back to the **ECS Instance List** page and click **China East 1**.

Hover the mouse pointer over the instance name and click the displayed pencil icon to change the instance name to ECS01 and ECS02, separately.

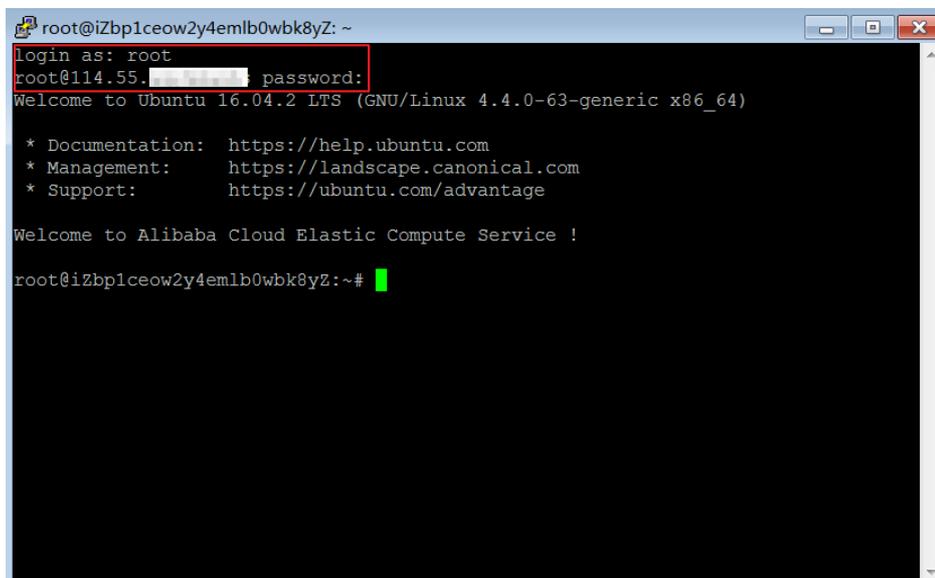
Instance ID/Name	IP Address	Status(All)	Network Type(All)	Billing Method(All)	Action
i-hc1t-ecs01	172.17.0.1 (Private IP Address)	Running	VPC	Pay-As-You-Go 17-07-23 17:23 created	Manage Connect More
i-hc1t-ecs02	172.17.0.2 (Private IP Address)	Running	VPC	Pay-As-You-Go 17-07-23 17:23 created	Manage Connect More

After you create the ECS instances, you need to deploy applications. In this tutorial, two static web pages will be deployed on the ECS instances using Apache.

Note: We use the default settings of Apache and only modify the content of the index file. Additionally, two Elastic IPs are bound to the ECS instances for easy management. For more information, see [Bind an EIP](#).

Procedure

Log on to the ECS instance.



```
root@iZbp1ceow2y4emlb0wbk8yZ: ~  
login as: root  
root@114.55.114.114 password:  
Welcome to Ubuntu 16.04.2 LTS (GNU/Linux 4.4.0-63-generic x86_64)  
  
* Documentation:  https://help.ubuntu.com  
* Management:    https://landscape.canonical.com  
* Support:       https://ubuntu.com/advantage  
  
Welcome to Alibaba Cloud Elastic Compute Service !  
root@iZbp1ceow2y4emlb0wbk8yZ:~#
```

Enter the following command to install Apache.

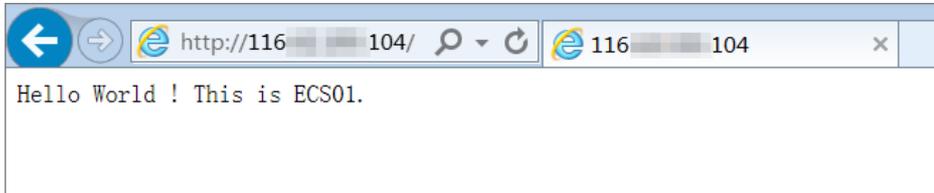
```
sudo apt-get install apache2
```

Enter the following command to modify the content of the index.html file.

```
cd /var/www/html
```

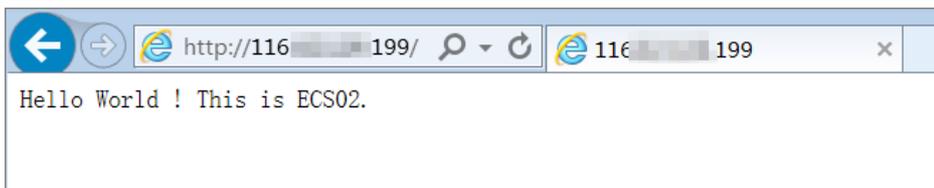
```
echo "Hello World! This is ECS01." > index.html
```

After modifying the content, enter the Elastic IP of the ECS instance in the web browser, you will see the following content.



Repeat the previous steps to create a web page on the another ECS instance and change the content to Hello World! This is ECS02..

After modifying the content, enter the Elastic IP of the ECS instance in the web browser, you will see the following content.



Before using Server Load Balancer, create a Server Load Balancer instance. You can add multiple listeners and backend servers to the Server Load Balancer instance.

Follow this document to create an Internet-facing Server Load Balancer instance. After creating the instance, a public IP is allocated to it and resolve a domain name to this IP.

Procedure

Log on to the Server Load Balancer console.

On the **Instances** page, click **Create Server Load Balancer**.

Configure the Server Load Balancer instance.

The configurations for the Server Load Balancer instance in this tutorial are as follows. For more information, see [Server Load Balancer configurations](#).

Region: Server Load Balancer does not support cross-region deployment. The region must be the same for the Server Load Balancer instance and the ECS instances. In this tutorial, select **China East 1**, which is the region of the ECS instance.

Zone type: Server Load Balancer has deployed multiple zones in most regions for better disaster tolerance. If Server Load Balancer service is unavailable in the primary zone, it will switch to a backup zone to restore service (within 30 seconds). Then, it will automatically switch back to the primary zone when the service is restored.

In this tutorial, select **China East 1 Zone B** as the primary zone and **China East 1 Zone D** as the backup zone.

Instance type: Select **Internet**.

The screenshot shows the configuration page for a Server Load Balancer instance. The page is titled "Server Load Balancer" and is divided into three main sections: "Basic Configuration", "Network and Instance Type", and "Purchase Plan".

Basic Configuration:

- Region:** A grid of region options is shown. "China East 1" is selected and highlighted in blue. Other regions include Singapore, Hong Kong, US East 1, US West 1, Asia Pacific NE 1, China East 2, China North 1, China North 2, China South 1, Northern China 3, Europe Central 1, Middle East 1, and Asia Pacific SE 2.
- Zone type:** "Multi-zone" is selected.
- Primary zone:** "China East 1 Zone B" is selected from a dropdown menu.
- Backup zone:** "China East 1 Zone D" is selected from a dropdown menu.

Network and Instance Type:

- Instance type:** "Internet" is selected from two options: "Internet" and "Intranet".
- Bandwidth:** "By traffic" is selected.

Purchase Plan:

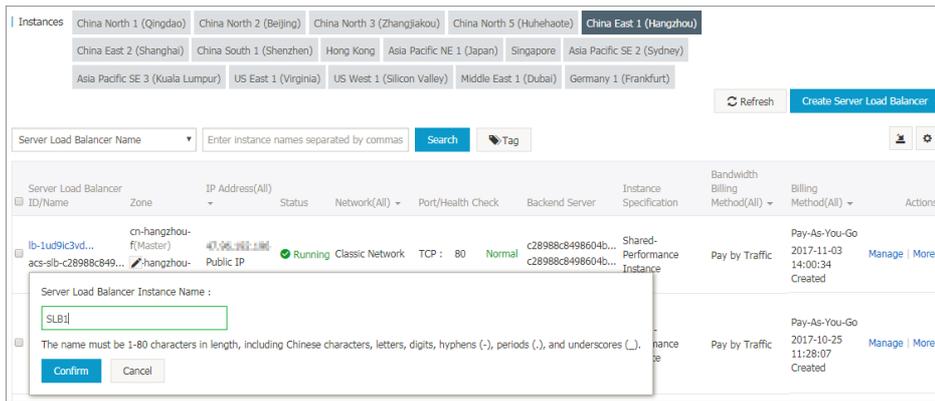
- Quantity:** "1" is entered in a text box.
- A note below the quantity field states: "You currently have 5 instances. You can create 25 more instances".

Click **Buy Now**.

Go back to the **Instances** page, find the created instance.

Hover the mouse pointer over the instance ID and then click the pencil icon.

Enter the name **SLB1** and click **Confirm**.



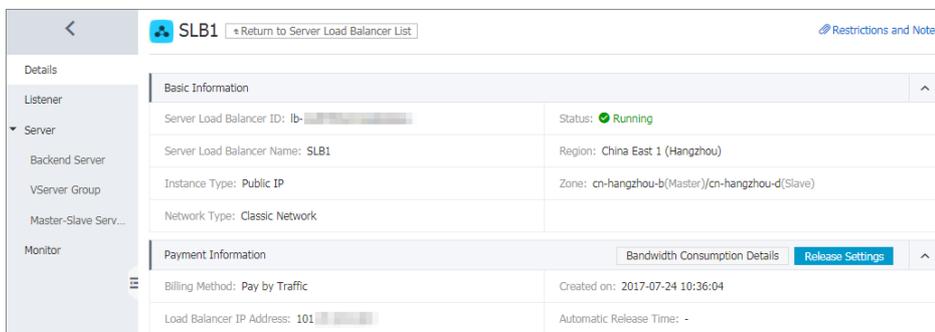
After creating a Server Load Balancer instance, you have to add at least one listener and a group of backend servers to it.

In this tutorial, we will create a TCP listener and add the ECS instances that have deployed web pages as the backend servers.

Procedure

Log on to the Server Load Balancer console.

On the **Instance Management** page, click the ID of the target Server Load Balancer instance.



On the left-side navigation pane, click **Listener** and then click **Add Listener**.

Configure the listener as follows and use the default settings for other options:

Frontend Protocol [Port]: The front-end protocol and port of the Server Load Balancer system that is used to receive and distribute connection requests. The port number cannot be the same in a Server Load Balancer instance.

In this tutorial, select the **TCP** protocol with port number **80**.

Backend Protocol [Port]: The port number that is opened on the ECS instances to receive the distributed requests. The port number can be the same in a Server Load Balancer instance.

In this tutorial, set to **80**.

Peak Bandwidth: You can set a peak bandwidth to limit the service capabilities that the application of the ECS instance can provide.

In this tutorial, no need to set the peak bandwidth because the instance is payed by traffic.

Scheduling Algorithm: Server Load Balancer supports the following scheduling algorithms. In this tutorial, the round-robin method is used.

Round robin: Requests are distributed evenly across the group of the backend ECS servers sequentially.

Weighted round robin (WRR): You can set a weight for each backend server. Servers with higher weights receive more requests than those with less weights.

Weighted least connections (WLC): In addition to the weight set to each backend ECS server, the number of connections to the client is also considered. The servers with a higher weight value will receive a larger percentage of live connections at any one time. If the weights are the same, the system directs network connections to the server with the least number of established connections.

Add Listener

1.Listener Configuration 2.Health Check 3.Success

Frontend Protocol [Port]* TCP : 80
You can enter any port number from 1-65535.

Backend Protocol [Port]* TCP : 80
You can enter any port number from 1-65535.

Peak Bandwidth: Unlimited [Configure](#)
You can set a peak bandwidth from 1-5000M. By default, the instances charged by traffic do not have peak bandwidth limit.

Scheduling Algorithm: Round Robin

Use Server Group:

Automatically Activate Listener after Creation: Activated

[Expand Advanced Options](#)

[Next Step](#) [Cancel](#)

Click **Next Step** to configure health check settings. Select the **TCP** mode and keep other settings as default, click **Confirm**.

Through health check on backend ECS instances, Server Load Balancer can automatically block abnormal ECS instances and distribute requests to them again when they become normal.

Add Listener
✕

1.Listener Configuration
2.Health Check
3.Success

Health Check Mode: TCP HTTP

Health Check Port:

If no port number is specified, the backend server port will be used for health checks by default.

Collapse Advanced Options

Response Timeout Duration: Second(s)

Max timeout for each health check request. Enter a value from 1-300 seconds, and the default value is 5 seconds.

Health Check Interval: Second(s)

Interval between health checks. Enter a value from 1-50 seconds, and the default value is 2 seconds.

Unhealthy Threshold: 2 3 4 5 6 7 8 9 10

The number of consecutive health check failures on the ECS servers (from success to failure).

Healthy Threshold: 2 3 4 5 6 7 8 9 10

The number of consecutive health check successes on the ECS servers (from failure to success).

Previous Step
Confirm
Cancel

Click **Confirm** to complete the configuration.

On the left-side navigation pane, click **Server** > **Backend Server**.

On the **Load Balancer Server Pool** page, click the **Servers Not Added** tab and select the previously created ECS instances, then click **Add in Batch**.

SLB1 - Return to Server Load Balancer List
Restrictions and Notes

Load Balancer Server Pool Region: China East 1 (Hangzhou) Zone: cn-hangzhou-b (Master) / cn-hangzhou-d (Slave)

Servers Added
Servers Not Added

Instance Name

ECS Instance ID/Name	Zone	Public/Internal IP Address	Status(AI)	Network Type(AI)	Server Load Balancer	Action
i-4h3...	cn-hangzhou-f	192.168.1.1 (Public) 192.168.1.2 (Private)	Running	Public	-	Add
i-4h3...	cn-hangzhou-f	192.168.1.1 (Public) 192.168.1.2 (Private)	Running	Public	-	Add
i-4h3...	cn-hangzhou-f	192.168.1.1 (Public) 192.168.1.2 (Private)	Running	Public	-	Add
i-4h3...	cn-hangzhou-e	192.168.1.1 (Public) 192.168.1.2 (Private)	Running	Public	-	Add
i-4h3...	cn-hangzhou-e	192.168.1.1 (Public) 192.168.1.2 (Private)	Stopped	Public	-	Add
<input type="button" value="Add in Batch"/>						

Total: 5 Item(s) | Per Page: 20 Item(s) | | |

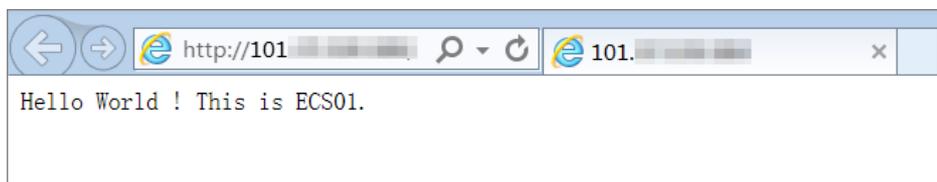
In the **Add a Backend Server** dialog, use the default weight value and click **Confirm**.

The higher the weight, the more requests are received.

Go back to the **Instance Management** page, click **Refresh**. When the health check is **Normal**, you can send requests to the Server Load Balancer instance.

Server Load Balancer ID/Name	Zone	IP Address	Status	Network	Port/Health Check	Backend Server	Instance Spec	Bandwidth	Billing Method	Billing Method	Action
lb-1uff19bc...	cn-hangzhou-5f7nvalv1	101.111.111.111	Running	Classic Network	TCP: 80 Normal	ECS02 ECS01	performance shared instance	Pay by Traffic	Pay-As-You-Go	2017-07-24 10:36:04 Created	Manage More

In the web browser, enter the IP address of the Server Load Balancer instance to test the service.



When you no longer need Server Load Balancer, delete the corresponding instance to avoid additional charges. When deleting the Server Load Balancer instance, the backend ECS will not be deleted or affected.

Note: After the Server Load Balancer instance is released, the backend ECS instances are still running. If you want to release the ECS instances, see **Release an instance**.

Procedure

Log on to the ECS console.

On the **Instances** page, select the region where the instance is located.

Select the target instance and click **Release Instance**.

In the **Release Instance** dialog, select **Release Now** or **Timed Release**.

If you select **Timed Release**, select the time to release the instance.

Click **Next** and click **OK** to finish.