**Solution**

# High-Availability AIGC Applications with Open-Source Models

| | |
|---|---|
| **Issue** | 1.0 |
| **Date** | 2023-08-10 |

# Contents

# 1 Solution Overview

## Scenarios

This solution helps you use Stable Diffusion to build high-availability Artificial Intelligence-Generated Content (AIGC) web applications on Huawei Cloud Elastic Cloud Server (ECS). Stable Diffusion is a latent text-to-image diffusion model capable of generating photo-realistic images given any text and images input.

## Solution Architecture

This solution helps you use Stable Diffusion to build high-availability AIGC web applications on Huawei Cloud Elastic Cloud Server (ECS). The following figure shows the architecture of this solution.

**Figure 1-1** Solution architecture



The following resources are required for deploying this solution:

- Two GPU-accelerated Linux ECSs, which will be used for running AIGC applications
- Three Elastic IP addresses (EIPs), which will be bound to the two Linux ECSs and an **Elastic Load Balance (ELB)**, respectively, for internal and external communication
- An ELB, which will be used to distribute traffic across availability zones (AZs)
- An Object Storage Service (OBS) bucket, which will be used to store generated image files
- Stable Diffusion web UI, inotify-tools, and OBS obsutil, which will be installed on each Linux ECS to automatically upload the images saved on the web UI

## Advantages

- High availability

  ECSs are deployed across AZs for multi-AZ disaster recovery (DR) and automatic, quick failover.

- Open source and custom development

  This solution is open-source and free for commercial use. You can also make custom development based on source code.

- Easy deployment

  You can easily deploy this solution with just a few clicks.

## Constraints

- Before deploying this solution, you need to have created a Huawei Cloud account and completed real-name authentication. You also need to ensure that the account is not frozen and has sufficient balance to pay for the resources required. You can estimate the total price according to the resource planning and costs tables.

# 2 Resource Planning and Costs

This solution deploys the services listed in the following table. The costs are only estimates and may differ from the final prices. For details, see **Price Calculator**.

**Table 2-1** Resource planning and costs (pay-per-use)

| Huawei Cloud Service | Example Configuration | Estimated Monthly Cost |
|---|---|---|
| Elastic Cloud Server (ECS) | <ul><li>Pay-per-use: $1.01 USD</li><li>Region: AP-Singapore</li><li>Billing Mode: Pay-per-use</li><li>Selected specifications: pi2.2xlarge.4 \| 8 vCPUs \| 32 GB</li><li>Accelerator: 1 x NVIDIA T4 / 1 x 16 GB</li><li>Image: Ubuntu 20.04 server 64bit with Tesla Driver 460.73.01 and CUDA 11.2</li><li>System Disk: High I/O \| 100 GB</li><li>Quantity: 2</li></ul> | $1.01 USD x 2 x 24 x 30 = $1,454.40 USD |
| Elastic IP (EIP) | <ul><li>Pay-per-use: $5.88 USD</li><li>Region: AP-Singapore</li><li>Billing Mode: Pay-per-use</li><li>Routing Type: Dynamic BGP</li><li>Billed By: Traffic</li><li>Traffic: 20 GB</li><li>IP Required Duration: 720 hours</li><li>EIP Quantity: 2</li></ul> | $5.88 USD x 2 = $11.76 USD |

| Huawei Cloud Service | Example Configuration | Estimated Monthly Cost |
|---|---|---|
| Elastic IP (EIP) | <ul><li>Pay-per-use: $0.13 USD/5 Mbit/s/hour</li><li>Region: AP-Singapore</li><li>Billing Mode: Pay-per-use</li><li>Product Type: Dedicated</li><li>Routing Type: Dynamic BGP</li><li>Billed By: Bandwidth</li><li>Bandwidth: 5 Mbit/s</li><li>Quantity: 1</li></ul> | $0.13 USD x 24 x 30 = $93.60 USD |
| Elastic Load Balance (ELB) | <ul><li>Pay-per-use: $0.05 USD</li><li>Region: AP-Singapore</li><li>Billing Mode: Pay-per-use</li><li>Type: Shared load balancer</li><li>Required Duration: 1 hour</li></ul> | $0.05 USD x 24 x 30 = $36.00 USD |
| Object Storage Service (OBS) | <ul><li>Region: AP-Singapore</li><li>Product: FunctionGraph</li><li>Request pricing tier:<br>≤ 1 million requests: $0 USD per 1 million requests<br>> 1 million requests: $0.2 USD per 1 million requests</li><li>Traffic pricing tier:<br>≤ 400,000 GB-seconds: $0 USD per GB-second<br>> 400,000 GB-seconds: $0.00001667 USD per GB-second</li></ul> | The OBS cost covers the storage and request cost as well as traffic cost. For details, see the monthly bill. |
| **Total** | - | **$1,595.76 USD + OBS price** |

**Table 2-2** Resource planning and costs (yearly/monthly)

| Huawei Cloud Service | Example Configuration | Estimated Monthly Cost |
|---|---|---|
| Elastic Cloud Server (ECS) | <ul><li>Region: AP-Singapore</li><li>Billing Mode: Yearly/Monthly</li><li>Selected specifications: pi2.2xlarge.4 \| 8 vCPUs \| 32 GB</li><li>Accelerator: 1 x NVIDIA T4 / 1 x 16 GB</li><li>Image: Ubuntu 20.04 server 64bit with Tesla Driver 460.73.01 and CUDA 11.2</li><li>System Disk: High I/O \| 100 GB</li><li>Quantity: 2</li></ul> | $549.30 USD x 2 = $1,098.60 USD |
| Elastic IP (EIP) | <ul><li>Pay-per-use: $5.88 USD</li><li>Region: AP-Singapore</li><li>Billing Mode: Pay-per-use</li><li>Routing Type: Dynamic BGP</li><li>Billed By: Traffic</li><li>Traffic: 20 GB</li><li>IP Required Duration: 720 hours</li><li>EIP Quantity: 2</li></ul> | $5.88 USD x 2 = $11.76 USD |
| Elastic IP (EIP) | <ul><li>Region: AP-Singapore</li><li>Billing Mode: Yearly/Monthly</li><li>Product Type: Dedicated</li><li>Routing Type: Dynamic BGP</li><li>Billed By: Bandwidth</li><li>Bandwidth: 5 Mbit/s</li></ul> | $57.00 USD |
| Elastic Load Balance (ELB) | <ul><li>Pay-per-use: $0.05 USD</li><li>Region: AP-Singapore</li><li>Billing Mode: Pay-per-use</li><li>Type: Shared load balancer</li><li>Required Duration: 1 hour</li></ul> | $0.05 USD x 24 x 30 = $36.00 USD |

| Huawei Cloud Service | Example Configuration | Estimated Monthly Cost |
|---|---|---|
| Object Storage Service (OBS) | <ul><li>Region: AP-Singapore</li><li>Product: FunctionGraph</li><li>Request pricing tier: ≤ 1 million requests: $0 USD per 1 million requests<br>> 1 million requests: $0.2 USD per 1 million requests</li><li>Traffic pricing tier: ≤ 400,000 GB-seconds: $0 USD per GB-second<br>> 400,000 GB-seconds: $0.00001667 USD per GB-second</li></ul> | The OBS cost covers the storage and request cost as well as traffic cost. For details, see the monthly bill. |
| **Total** | - | **$1,203.36 USD + OBS price** |

# 3 Procedure

## 3.1 Preparations

### Creating the rf_admin_trust Agency

**Step 1**  Log in to the **Huawei Cloud console**, hover the mouse pointer over the account name in the upper right corner, and choose **Identity and Access Management**.
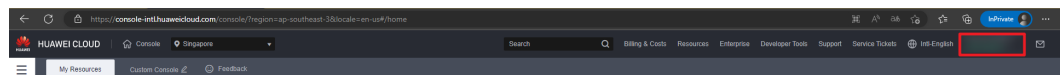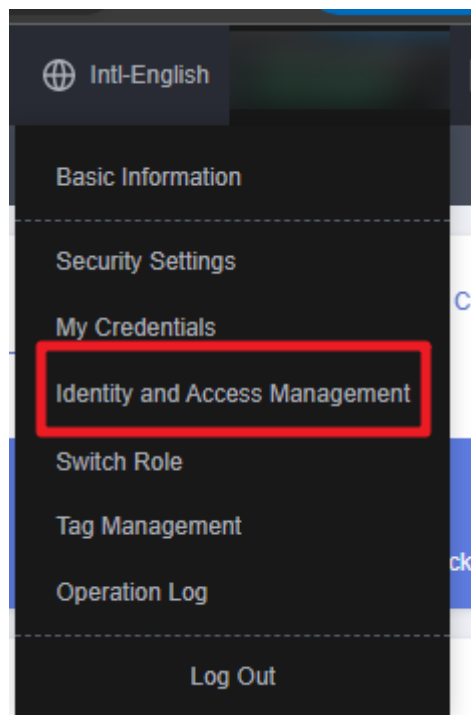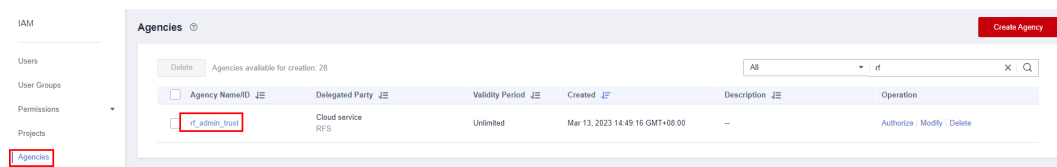
**Figure 3-1** Huawei Cloud console

**Figure 3-2** Identity and Access Management



**Step 2** Choose **Agencies** and then search for the **rf_admin_trust** agency in the agency list.

**Figure 3-3** Agency list



- If the agency is found, skip the following steps.
- If the agency is not found, perform the following steps to create it.

**Step 3** Click **Create Agency** in the upper right corner of the page. On the displayed page, enter **rf_admin_trust** for **Agency Name**, select **Cloud service** for **Agency Type**, select **RFS** for **Cloud Service**, and click **Next**.

**Figure 3-4** Create Agency



**Step 4**  Search for **Tenant Administrator**, select it in the search results, and click **Next**.

**Figure 3-5** Selecting a policy/role



**Step 5**  Select **All resources** and click **OK**.

**Figure 3-6** Selecting a scope



**Step 6**  If **rf_admin_trust** is displayed in the agency list, the agency has been created.
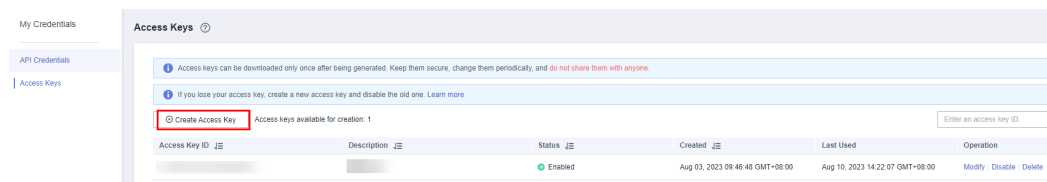
---

**Figure 3-7** Agency list



**----End**

## Obtaining an Access Key (AK/SK)

Before deploying this solution, you need to obtain the AK/SK from the Huawei Cloud console and then configure parameters listed in **Step 3**.

On the Huawei Cloud console, hover the mouse pointer over the account name in the upper right corner and choose **My Credentials**. On the **Access Keys** page, create an access key and download it. For details, see **How Do I Obtain an Access Key (AK/SK)?**

**Figure 3-8** Creating an access key



# 3.2 Quick Deployment

This section describes how to quickly deploy this solution.

**Table 3-1** Parameters required for deploying this solution

| Parameter | Type | Mandatory | Description | Default Value |
|---|---|---|---|---|
| vpc_name | String | Yes | Virtual Private Cloud (VPC) name. This template uses a newly created VPC and the VPC name must be unique. The value can contain 1 to 54 characters, including letters, digits, underscores (_), hyphens (-). and periods (.). | high-availability-aigc-applications-demo |

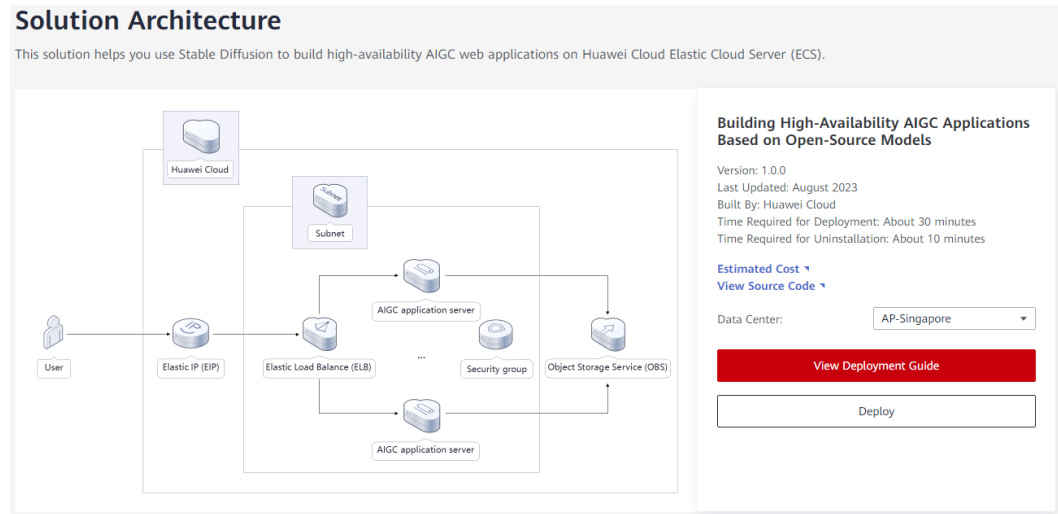| Parameter | Type | Mandatory | Description | Default Value |
|---|---|---|---|---|
| security_group_name | String | Yes | Security group name. This template uses a newly created security group. The value can contain 1 to 64 characters, including letters, digits, underscores (_), hyphens (-), and periods (.). | high-availability-aigc-applications-demo |
| ecs_name | String | Yes | ECS name, which must be unique. The name format is {ecs_name}-digit. It can contain 1 to 60 characters, including letters, digits, underscores (_), hyphens (-). and periods (.). | high-availability-aigc-applications-demo |
| image_bucket_name | String | Yes | OBS bucket name, which is globally unique. The bucket is used to store automatically uploaded images saved on the web UI. The bucket name can contain 3 to 63 characters, including lowercase letters, digits, hyphens (-), and periods (.). Do not start or end with a hyphen (-) or a period (.). | None |
| ecs_count | String | Yes | Number of ECSs, which is greater than or equal to 1. The maximum number of ECSs is determined by the user quota listed in **My Quotas**. | 2 |
| ecs_flavor | String | Yes | ECS flavor. This template uses a GPU-accelerated flavor. For details about flavors, see **A Summary List of x86 ECS Specifications**. | pi2.2xlarge.4 |

| Paramete r | Type | Mandator y | Description | Default Value |
|---|---|---|---|---|
| ecs_passw ord | String | Yes | Initial password of the ECS. After the ECS is created, reset the password by referring to **Step 1**. It can contain 8 to 26 characters, including at least three of the following character types: uppercase letters, lowercase letters, digits, and special characters (! @$%^-_=+[{()}]:,./?~#*). For Windows ECSs, the password cannot contain the username, the username spelled backwards, or more than two consecutive characters in the username. The default administrator account is root. | None |
| elb_name | String | Yes | ELB name, which can contain 1 to 64 characters, including letters, digits, underscores (_), hyphens (-), and periods (.). | high- availability- aigc- applications- demo |
| eip_band width_size | Number | Yes | EIP bandwidth, which is billed by traffic. Value range: 1-300 Mbit/s | 300 Mbit/s |
| charging_ mode | String | Yes | Billing mode. By default, expenses are automatically deducted. The value can be **prePaid** (yearly/ monthly) or **postPaid** (pay-per-use). | postPaid |
| charge_pe riod_unit | String | Yes | Unit of a subscription period. This parameter is valid and mandatory only when **charging_mode** is set to **prePaid**. The value can be **month** or **year**. | month |

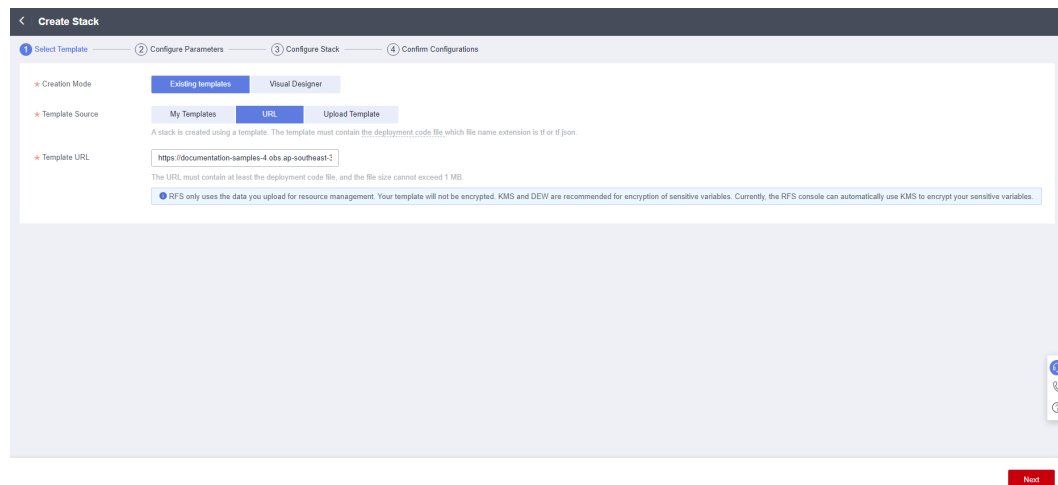| Parameter | Type | Mandatory | Description | Default Value |
|---|---|---|---|---|
| charge_period | Number | Yes | Subscription period. This parameter is valid and mandatory only when **charging_mode** is set to **prePaid**. Value range:<br><br>If **charge_period_unit** is set to **month**, the value range is from **1** to **9**.<br><br>If **charge_period_unit** is set to **year**, the value range is from **1** to **3**. | 1 |
| access_key_id | String | Yes | Access Key ID (AK), which is used to verify the identity of a user attempting to upload generated images to the OBS bucket. For details about how to obtain the AK, see **Obtaining an Access Key (AK/SK)**. | None |
| secret_access_key | String | Yes | Secret Access Key (SK), which is used to sign requests. It must be used together with the AK to authenticate image upload requests. For details about how to obtain the SK, see **Obtaining an Access Key (AK/SK)**. | None |

**Step 1** Log in to **Practical Application of Huawei Cloud Solutions** and select **Building High-Availability AIGC Applications Based on Open-Source Models**.
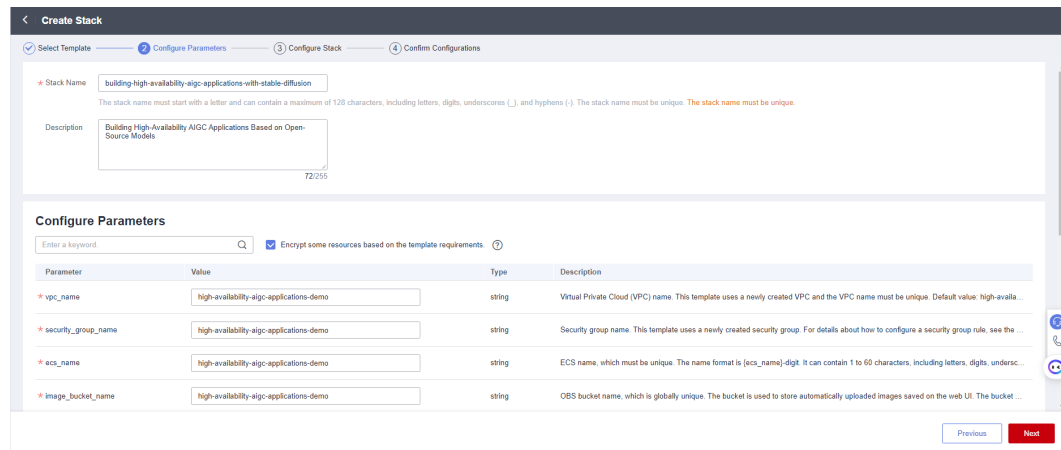
**Figure 3-9** Selecting a solution



**Step 2** Click **Deploy** to switch to the **Create Stack** page.

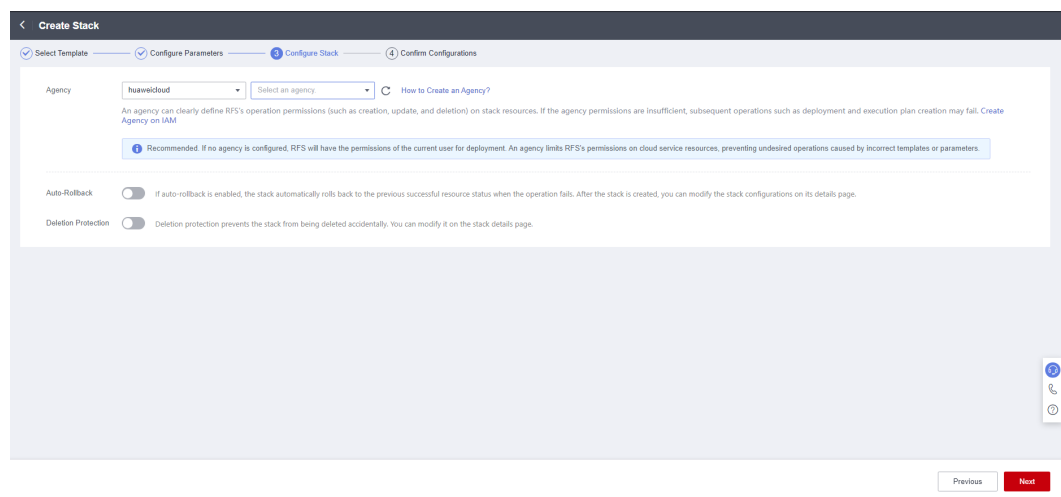**Figure 3-10** Creating a stack



**Step 3** Click **Next**. On the displayed page, set parameters by referring to **Table 3-1** and click **Next**.
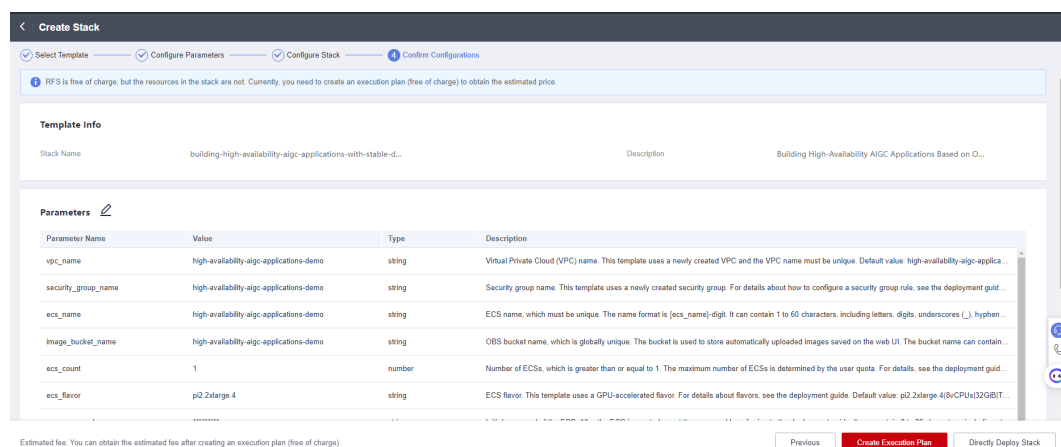
**Figure 3-11** Configuring parameters



**Step 4** (Optional) On the **Configure Stack** page, select **rf_admin_trust** from the agency drop-down list and click **Next**.

**Figure 3-12** Configuring a stack



**Step 5** On the **Confirm Configurations** page, click **Create Execution Plan**.

**Figure 3-13** Creating an execution plan

**Step 6** In the displayed **Create Execution Plan** dialog box, enter an execution plan name and click **OK**.

**Figure 3-14** Creating an execution plan



**Step 7** On the **Execution Plans** tab, click **Deploy**. In the displayed dialog box, click **Execute**.
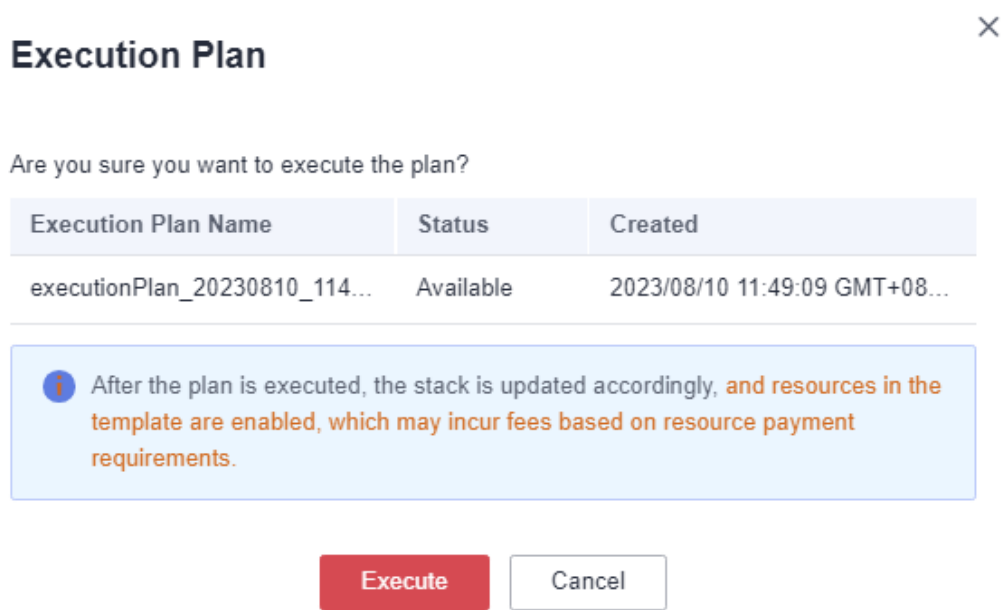
**Figure 3-15** Deploying the execution plan

**Figure 3-16** Confirming the deployment



**Step 8** Wait until the deployment is successful and click the **Events** tab to view details.

**Figure 3-17** Successful deployment



**Step 9** Refresh the page and view the web UI access description on the **Outputs** tab.

**Figure 3-18** Outputs



**----End**

# 3.3 Getting Started

## (Optional) Modifying Security Group Rules

> **NOTICE**
>
> By default, this solution creates the security group rule that uses the ping command to test ECS connectivity. To remotely log in to an ECS in the security group, you need to add an inbound rule, for example, by setting the login port to 3389 and adding a whitelist IP address.

A security group is a collection of access control rules for cloud resources, such as cloud servers, containers, and databases, to control inbound and outbound traffic. Cloud resources associated with the same security group have the same security requirements and are mutually trusted within a VPC.

You can modify security group rules, for example, by adding, modifying, or deleting a TCP port.

- Adding a security group rule: **Add an inbound rule** and enable a TCP port if needed.

- Modifying a security group rule: Inappropriate security group settings can be a serious security risk. You can **modify security group rules** to improve the network security of ECSs.

- Deleting a security group rule: If the source or destination IP address of an inbound or outbound security group rule changes, or a port needs to be disabled, you can **delete the security group rule**.
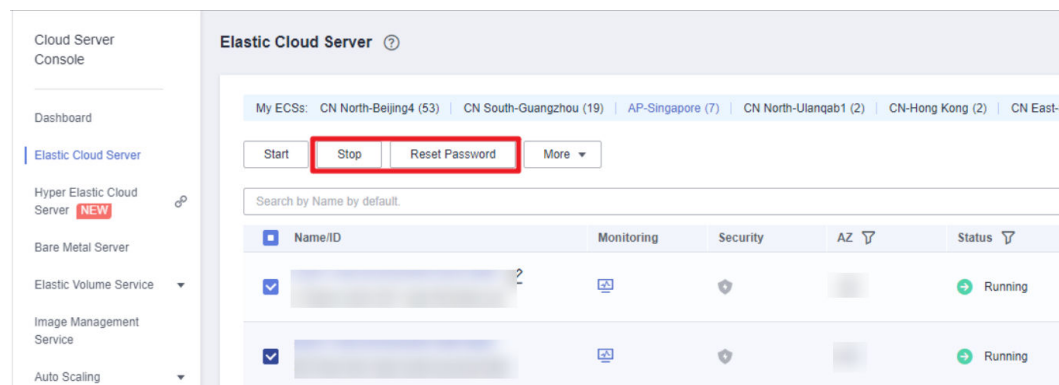
## (Optional) Configuring a Domain Name for an Application

Configure domain name resolution to resolve the website domain name to the IP address displayed in the figure in **Step 9**. In this way, the website can be accessed over its domain name. For details about DNS resolution, see **Public Domain Name Resolution**.

## Using the AIGC Web UI Application

**Step 1**  (Optional) Log in to the **ECS console**, select the created ECSs and click **Stop** above the ECS list. After the ECSs are stopped, click **Reset Password**, enter a new password, and click **OK**. The new password will take effect after the ECSs are started.

**Figure 3-19** Resetting the password



**Step 2**  Log in to the **ELB console** and choose **Backend Server Groups** from the left navigation pane. Click the target backend server group name to view its details. On the **Backend Servers** tab, check whether servers are in the **Healthy** state. (Service initialization will be completed 20 minutes after this solution is deployed based on default settings. All backend servers will be healthy on port 7860 then.)
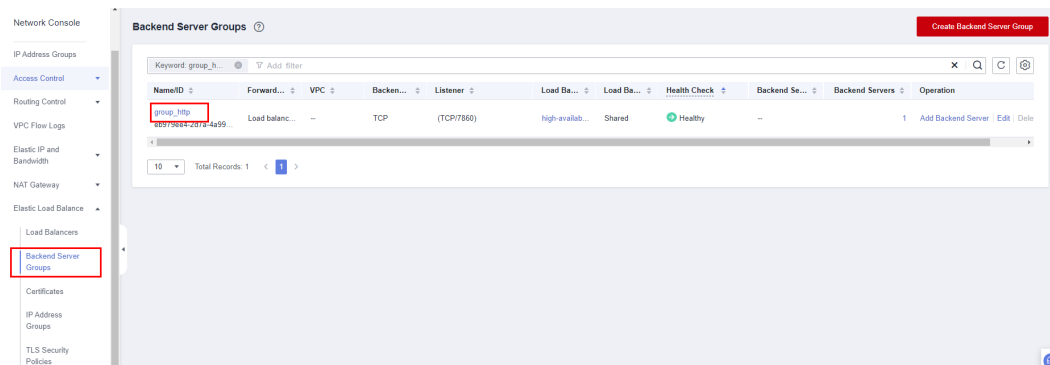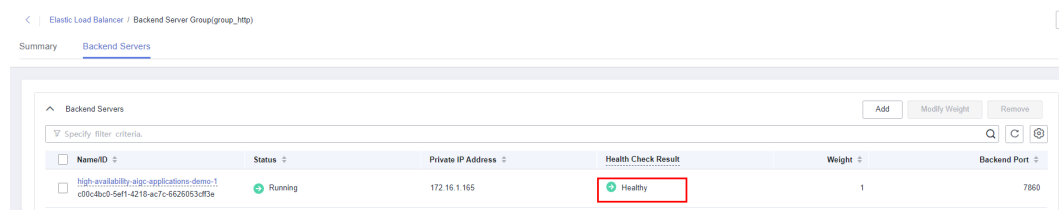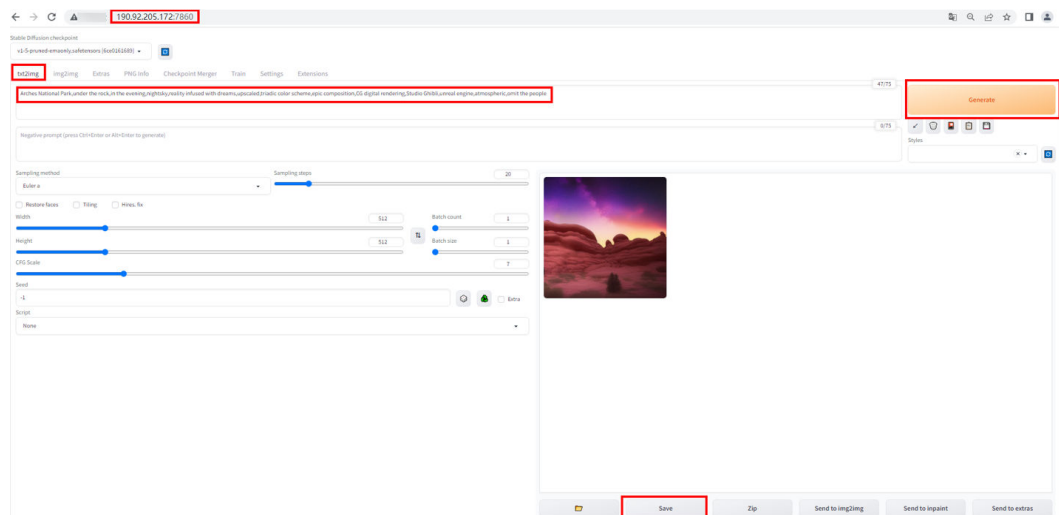
**Figure 3-20** Backend Server Groups



**Figure 3-21** Backend server status



**Step 3** Use the URL in **Step 9** to access the AIGC web UI. Click **txt2img**, input text in the box, and click **Generate**. After the image is generated, click **Save**. For details about how to use the Stable Diffusion web UI, see **stable-diffusion-webui** or search for tutorials on the Internet. This solution creates the user **aigc** with the default password **aigc@123**.

**Figure 3-22** AIGC web UI



Example text prompt:
Arches National Park,under the rock,in the evening,nightsky,reality infused with dreams,upscaled,triadic color scheme,epic composition,CG digital rendering,Studio Ghibli,unreal engine,atmospheric,omit the people

**Step 4** On the **OBS console**, click the OBS bucket created in **Step 3** to view the saved images. You can also click **Share** to share the images. For more OBS functions, see **Object Management**.
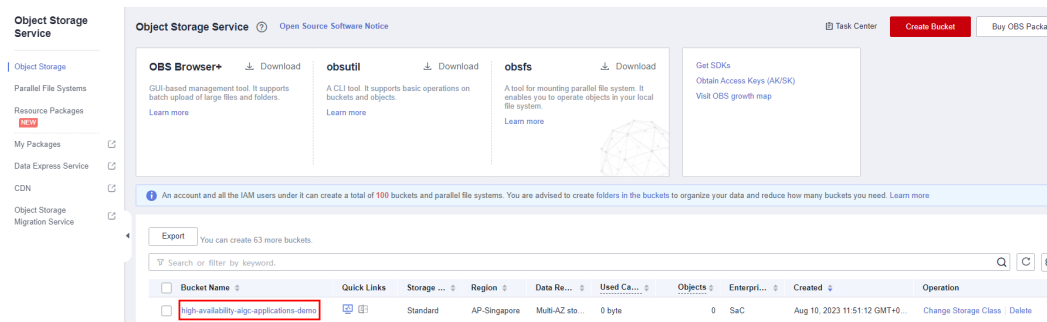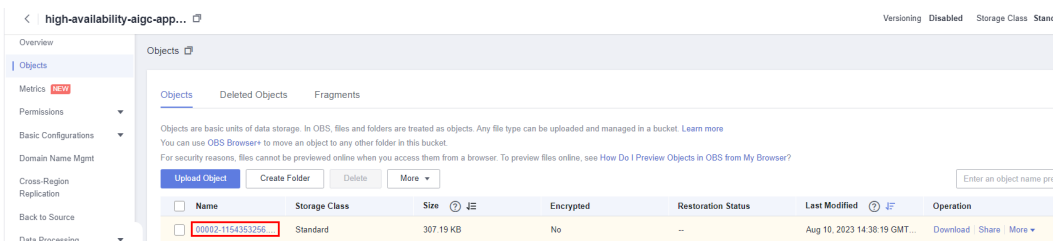
**Figure 3-23** OBS bucket list



**Figure 3-24** Viewing the saved images



📖 **NOTE**

This solution has added inotify-tools and OBS obsutil to run automatically at startup in
ECSs, so the images you saved on the web UI can be automatically uploaded to the OBS
bucket. You can also right-click on the web UI and choose **Save as** to save the images
locally. AIGC can run automatically at ECS startup.
Example command for starting AIGC:
Start in the foreground:
cd /home/aigc && sudo -u aigc bash -c "source /home/aigc/webui.sh --listen --port 7860 --api --
enable-insecure-extension-access"
Start in the background:
cd /home/aigc && sudo -u aigc bash -c "source /home/aigc/webui.sh --listen --port 7860 --api --
enable-insecure-extension-access &">> /home/aigc/aigc-applications.log (change the storage path as
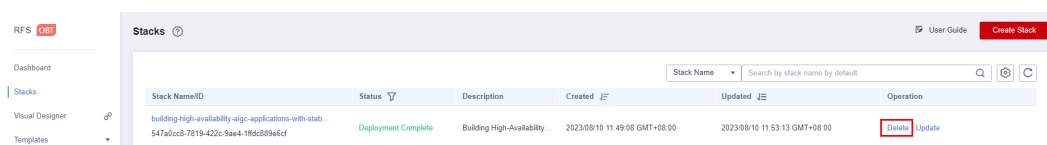needed)

**----End**

# 3.4 Quick Uninstallation

**NOTICE**

If there is data stored in the OBS bucket, the solution cannot be uninstalled. To
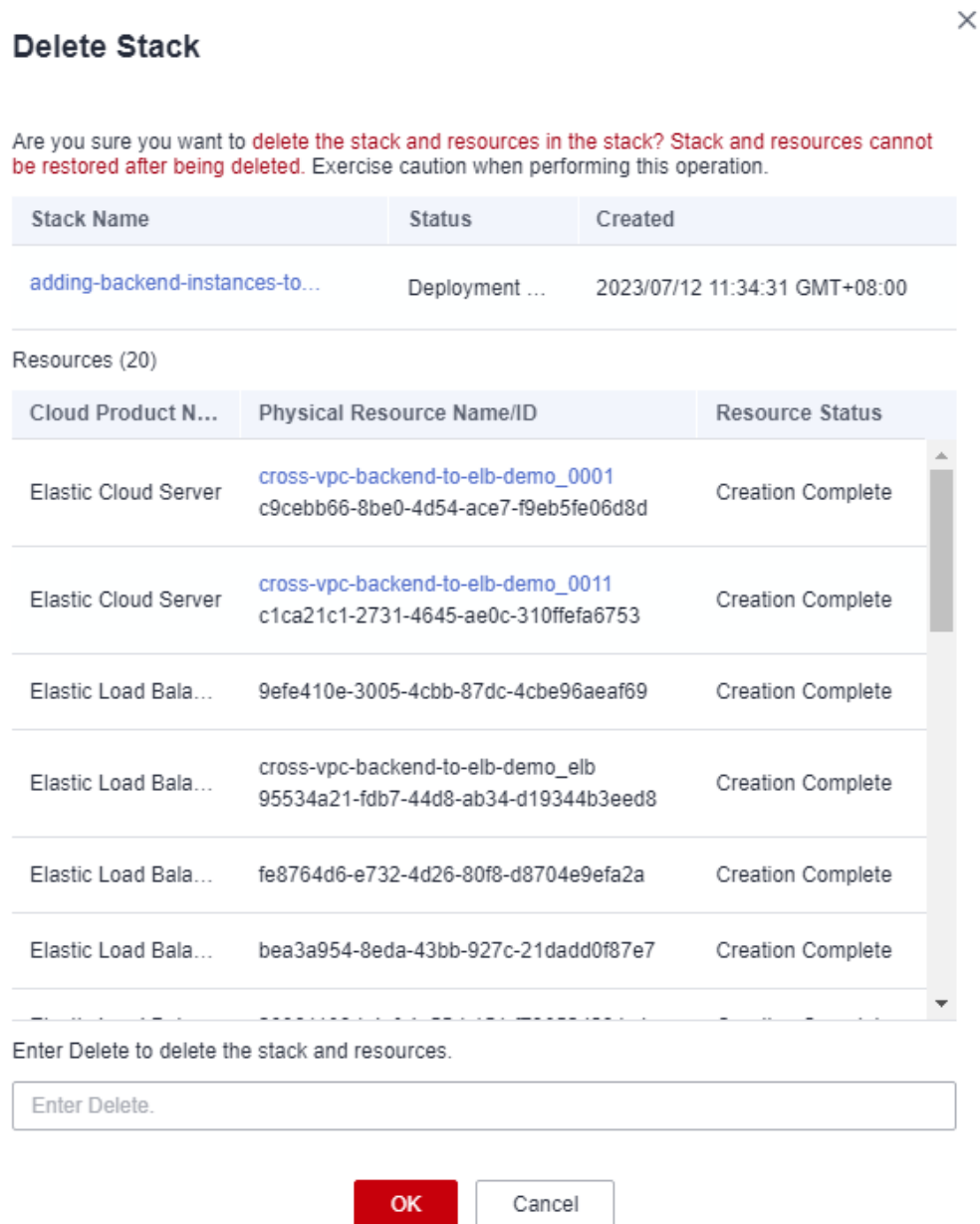uninstall this solution, back up the data and empty the bucket first.

**Step 1** Click **Delete** in the row of the solution stack.

**Figure 3-25** Deleting the stack

**Step 2** Enter **Delete** and click **OK**.

**Figure 3-26** Confirming the uninstallation



----**End**

# 4 Appendix

## Terms

Basic concepts and cloud service introduction

- Elastic Cloud Server (ECS): a scalable and on-demand cloud server. It helps you to efficiently set up reliable, secure, and flexible application environments, ensuring stable service running and improving O&M efficiency.

- Elastic Load Balance (ELB): automatically distributes incoming traffic across multiple servers to balance their workloads, increasing service capabilities and fault tolerance of your applications.

- Elastic IP (EIP): enables your cloud resources to communicate with the Internet using static public IP addresses and scalable bandwidths. EIPs can be bound to or unbound from ECSs, BMSs, virtual IP addresses, load balancers, and NAT gateways.

- Virtual Private Cloud (VPC): an isolated and private virtual network environment. You can configure IP address segments, subnets, and security groups, assign EIPs, and allocate bandwidths in a VPC.

- Object Storage Service (OBS): a secure, highly reliable object storage service that allows you to inexpensively store any amount of data.

- Security group: a collection of access control rules for ECSs that have the same security protection requirements and are mutually trusted in a VPC. You can define inbound and outbound rules to control traffic to and from the ECSs in a security group.

- inotify-tools: a command-line tool in Linux to monitor file system changes and trigger corresponding operations.

# 5 Change History

| Released On | Description |
|---|---|
| 2023-08-10 | This issue is the first official release. |