

Deploying Cloudera In Multiple GCP Availability Zones (Preview)

Date published: 2025-01-22

Date modified: 2025-02-25

Legal Notice

© Cludera Inc. 2024. All rights reserved.

The documentation is and contains Cludera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cludera software may be found within the documentation accompanying each component in a particular release.

Cludera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 ("ASLv2"), the Affero General Public License version 3 (AGPLv3), or other license terms.

Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cludera software product page for more information on Cludera software. For more information on Cludera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cludera reserves the right to change any products at any time, and without notice. Cludera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cludera.

Cludera, Cludera Altus, HUE, Impala, Cludera Impala, and other Cludera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners. Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER'S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided 'as is' without warranty or support. Further, Cludera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.

Contents

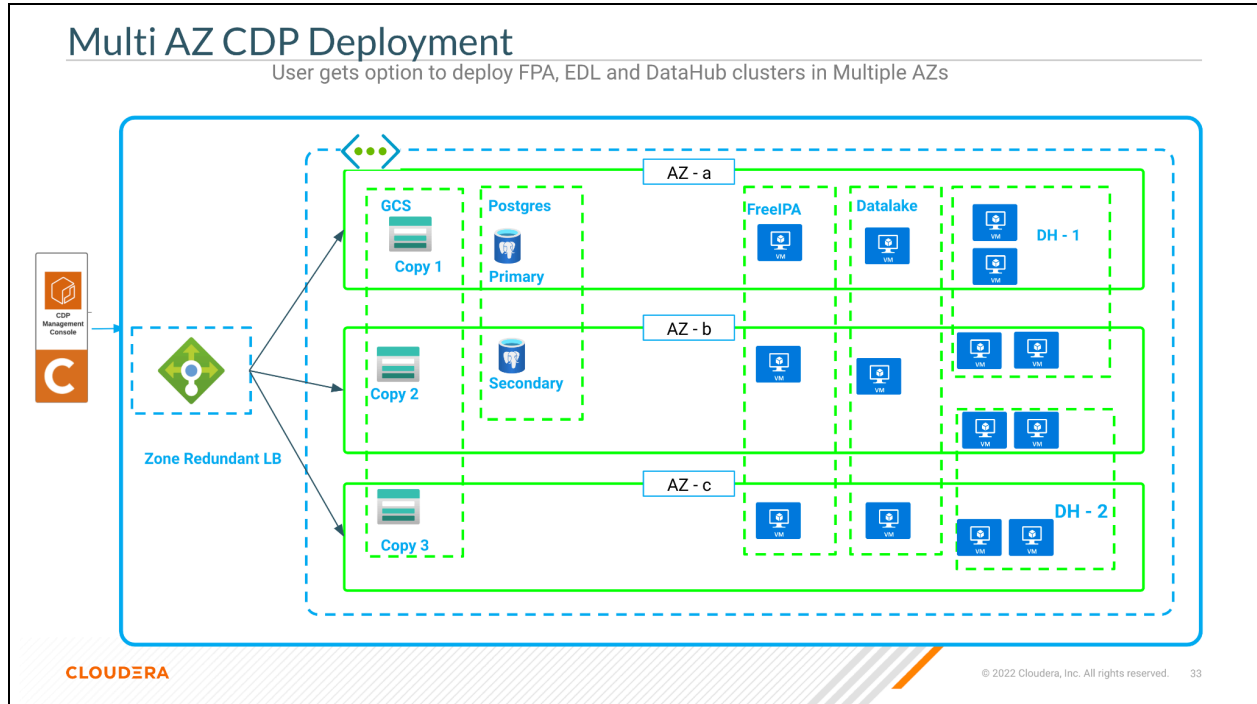
Legal Notice	2
Contents	3
Use cases	6
Limitations	7
Register a multi-AZ environment	7
Cloudera UI	7
CDP CLI	8
Modify an environment to add AZs	9
Cloudera UI	9
CDP CLI	9
Create a multi-AZ Cloudera Data Hub cluster	10
Prerequisites	10
Cloudera UI	10
CDP CLI	10

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided 'as is' without warranty or support. Further, Cloudera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.

You can optionally choose to deploy Data Lake, FreeIPA, and Cloudera Data Hub clusters across multiple availability zones (multi-AZ). With multi-AZ support, newly created GCP environments, enterprise Data Lakes and Cloudera Data Hub clusters using HA templates can be deployed across multiple availability zones of the selected GCP region. This provides fault tolerance during the extreme event of an availability zone outage.

Note: Only Enterprise Data Lakes can use multi-AZ; other Data Lake shapes do not support it. Enterprise Data Lake is only available in Cloudera Runtime 7.2.17 and newer; therefore, multi-AZ for Cloudera on GCP is only available in Cloudera Runtime 7.2.17 and newer. You need to contact Cloudera Support to have this feature enabled.

Each GCP region has multiple availability zones, which act as failure domains, preventing small outages from affecting entire regions. If you choose to deploy your Cloudera environment (FreeIPA and Data Lake) and Cloudera Data Hub clusters across multiple availability zones, each of these components is spread across the configured availability zones, providing high availability and fault tolerance. This is illustrated in the following diagram:



With the multi-AZ option enabled, your services are deployed in the following way:

- GCP environments are always created with the configured number of FreeIPA servers, deployed on virtual machines spread across the selected availability zones.
- In a GCP Enterprise Data Lake each host group is configured so that virtual machines of all critical services are spread across the selected availability zones.
- In Cloudera Data Hub clusters, virtual machines of each host group are evenly spread across the selected availability zones, following a round-robin logic.

When a zone failure happens and a cluster needs to be repaired, the replacement VMs are always provisioned in the same subnet and availability zone as the old ones since the detached disks can only be reattached to a VM in the same availability zone. This means that if there is an availability zone outage, cluster repair is not possible.

By default, if you do not enable multi-AZ, All GCP resources are placed in a fixed Availability Zone and the zone is either provided by the customer or selected randomly.

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided 'as is' without warranty or support. Further, Cloudera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.

When creating Cloudera Data Hub clusters via CDP CLI, you have the option to specify the AZ, which, in addition to allowing you to select the AZs that should be used, allows you to set up AZ targeting, where all nodes of the cluster are placed on the same AZ. This enables creating disaster recovery scenarios, where a primary and secondary cluster are running in different AZs. If an AZ outage occurs and the primary cluster is lost, it is guaranteed that the secondary cluster is not impacted.

Use cases

A multi-AZ Data Lake and FreeIPA constitute a resilient environment that provides a solid basis for multi-AZ Cloudera Data Hub clusters and Cloudera data services. Cloudera Data Hub clusters and Cloudera data services depend on the FreeIPA instance in the Data Lake to provide DNS resolution. Deploying FreeIPA across multiple availability zones ensures that critical DNS resolution is available in the event of an availability zone outage. Furthermore, a medium duty or enterprise Data Lake provides high availability, and additional compute and memory resources for key SDX services and is recommended for production workloads.

Deploying your Cloudera Data Hub clusters across multiple availability zones is key if your mission-critical applications depend on HBase and Kafka. Multiple availability zone deployment for operational workloads is considered best practice by the cloud vendors. It ensures that your applications can continue to run in the event of an availability zone outage.

When an entire availability zone fails, HBase automatically rebalances regions among the remaining instances in the cluster to maintain availability. The write-ahead log (WAL), which is replicated across the configured availability zones is automatically replayed by the newly assigned region servers in other availability zones to ensure writes to the database are not lost.

When using the multi availability zone feature, Cloudera ensures that Kafka replicates partitions across brokers in different availability zones. During an availability zone failure this ensures that no data is lost and applications can continue to access the data they need. Cruise Control, which is deployed alongside every Kafka cluster in Cloudera on cloud, detects that topics need to be rebalanced to the remaining brokers. Once the availability zone is back online, you can repair your Kafka cluster, restoring the initial broker distribution across availability zones. Afterwards Cruise Control kicks in and ensures that all topic partitions are balanced across the cluster.

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided 'as is' without warranty or support. Further, Cloudera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.

Limitations

The following limitations apply when deploying a multi-AZ Cloudera:

- When an AZ is down, you cannot create a new Cloudera Data Hub cluster, and create or activate Cloudera data services within the environment. Existing workloads will continue to work.
- When an AZ is down, you cannot resize, stop, or restart Cloudera Data Hub clusters.
- Non-AZ environments or clusters cannot be converted to multi-AZ.

Register a multi-AZ environment

You can register a multi-AZ GCP environment via Cloudera UI or CDP CLI. You may choose to enable multi-AZ for Data Lake only or for FreeIPA only. There is no requirement to enable both.

Cloudera UI

Register your environment as usual, just make sure to do the following:

1. On the Data Access and Data Lake Scaling page:
 - a. Select to use the Enterprise Data Lake.
 - b. On the same page, scroll down and in the bottom of the page enable the **Advanced Options**.
 - c. In the **Network and Availability** section enable the **Enable Multiple Availability Zones for Data Lake** toggle button in order to enable multi-AZ for Data Lake. The option is disabled by default. The option only appears when the Enterprise Data Lake is selected.
2. On the **Region, Networking, and Security** page:
 - a. Scroll down and in the bottom of the page enable the **Advanced Options**.
 - b. In the **Network and Availability** section enable the **Enable Multiple Availability Zones for Data Lake** toggle button in order to enable multi-AZ for FreeIPA. The option is disabled by default.
3. Finish registering your environment as usual.

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided 'as is' without warranty or support. Further, Cloudera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.

CDP CLI

Use the following CDP CLI commands to register an environment with a multi-AZ Data Lake and FreeIPA:

1. Register an GCP environment using the `cdp environments create-gcp-environment` command and include `multiAz=true` in the `--free-ipa` parameter as shown in this example:

Unset

```
cdp environments create-gcp-environment \
  --environment-name test-env \
  ...
  --free-ipa instanceCountByGroup=3,multiAz=true \
```

If you do not include the `multiAz=true`, the fixed AZ will be used.

You can also optionally include the `--availability-zones` parameter to select the specific availability zones that should be used. Valid value for availability zone is `<region>-<zone>`. `<zone>` is a letter like `a`, `b`, `c`, etc. Each region has a different number of availability zones. For example, the `us-west2` region has zones `us-west2-a`, `us-west2-b`, and `us-west2-c`. If this parameter is not provided, all AZs are used. For example:

Unset

```
cdp environments create-gcp-environment \
  --environment-name test-env \
  ...
  --free-ipa instanceCountByGroup=3,multiAz=true \
  --availability-zones us-west2-a,us-west2-b
```

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided 'as is' without warranty or support. Further, Cloudera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.

2. Set IDBroker mappings as usual using the `cdp environments set-id-broker-mappings` command.
3. Create a Data Lake using the `cdp datalake create-gcp-datalake` command and adding the `--multi-az` parameter. For example:

Unset

```
cdp datalake create-gcp-datalake \
  --datalake-name test-dl \
  --environment-name test-env \
  ...
  --scale ENTERPRISE \
  --runtime 7.2.17 \
  --multi-az
```

Modify an environment to add AZs

You can modify an environment to add more AZs to it. In this case, added AZs will only be used for new Cloudera Data Hub cluster clusters. Existing clusters will continue to use the original AZ.

Cloudera UI

1. In the Cloudera Management Console, navigate to environment details > **Summary**.
2. Scroll down to the **Region** section.
3. Under **Select GCP availability zones**, select the availability zones:
4. Click **Save**.

CDP CLI

Use the CLI command for updating AZs of an existing environment:

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided 'as is' without warranty or support. Further, Cloudera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.

Unset

```
cdp environments update-gcp-availability-zones \
--environment-name <ENV-NAME> \
--availability-zones <LIST-OF-AZs>

cdp environments update-gcp-availability-zones \
--environment-name test-env \
--availability-zones us-west2-a us-west2-b
```

Create a multi-AZ Cloudera Data Hub cluster

You can create multi-AZ Cloudera Data Hub clusters within any existing environment. Detailed steps are provided below.

Prerequisites

You can create a multi-AZ Cloudera Data Hub cluster in a multi-AZ environment only. If you are trying to create a multi-AZ Cloudera Data Hub cluster in an environment that uses the default AZ distribution, you need to first edit that environment and add AZs to it.

Cloudera UI

To enable multi-AZ when creating a Cloudera Data Hub cluster on GCP, navigate to the **Advanced Options > Network And Availability** and in the “GCP Availability Zones” section click the toggle button next to **Enable using multiple availability zones**.

CDP CLI

You can create a multi-AZ Cloudera Data Hub cluster by adding the `--multi-az` option to the Cloudera Data Hub cluster creation command.

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided ‘as is’ without warranty or support. Further, Cloudera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.

In the `--instance-groups` parameter, you can optionally include the `availabilityZones` to select the specific availability zones that should be used. If this parameter is not provided, all three AZs are used. For example:

Unset

```
cdp datahub create-gcp-cluster \
  --cluster-name test-cluster1 \
  --environment-name test-env \
  --cluster-template-name "7.2.17 - Data Engineering: Apache
Spark, Apache Hive, Apache Oozie" \
  --multi-az \
  --instance-groups
nodeCount=1,instanceGroupName=compute,instanceGroupType=CORE,inst
anceType=Standard_D5_v2,rootVolumeSize=100,attachedVolumeConfigur
ation=\[\{volumeSize=100,volumeCount=0,volumeType=StandardSSD_LRS
\}\],recoveryMode=MANUAL,availabilityZones=\[us-west2-a,us-west2-
b\]
nodeCount=0,instanceGroupName=gateway,instanceGroupType=CORE,inst
anceType=Standard_D8_v3,rootVolumeSize=100,attachedVolumeConfigur
ation=\[\{volumeSize=100,volumeCount=1,volumeType=StandardSSD_LRS
\}\],recoveryMode=MANUAL,availabilityZones=\[us-west2-b,us-west2-
c\]
nodeCount=1,instanceGroupName=master,instanceGroupType=GATEWAY,in
stanceType=Standard_D16_v3,rootVolumeSize=100,attachedVolumeConfi
guration=\[\{volumeSize=100,volumeCount=1,volumeType=StandardSSD_
LRS\}\],recoveryMode=MANUAL,availabilityZones=\[us-west2-a,us-wes
t2-b,us-west2-c\]
nodeCount=3,instanceGroupName=worker,instanceGroupType=CORE,insta
nceType=Standard_D5_v2,rootVolumeSize=100,attachedVolumeConfigura
tion=\[\{volumeSize=100,volumeCount=1,volumeType=StandardSSD_LRS\
}\],recoveryMode=MANUAL,availabilityZones=\[us-west2-a,us-west2-c
\]
```

This document has been released as part of a technical preview for features described herein. Technical preview components are provided as a convenience to our customers for their evaluation and trial usage. These components are provided 'as is' without warranty or support. Further, Cloudera assumes no liability for the usage of technical preview components, which should be used by customers at their own risk.