# How to enable Cloudera AI Inference Service (CAII) on Sandbox and Workshop tenants

## About this Page 🔗

The CAII service has multiple dependencies and so deploying it is a multi step process.  This guide will walk through the steps to set up the environment/datalake, enable compute clusters for the environment, create a compute cluster for CAII, deploy a CAI Workbench, deploy a model registry, and finally deploy the CAII service and create a model endpoint.

Documentation on CAII Service can be found here: Ⓒ Cloudera AI Inference service Overview

> ⚠️ ***Please note:*** using the CAII service requires significant resources - 4 total kubernetes clusters plus GPUs for model serving - be mindful of the cost associated with deployment and utilize autoscaling and appropriate CPU and GPU node shapes for a given demo/workshop/poc to help limit cost. Please also stop/terminate the environment when not in use

## Instructions 🔗

## 1. Create Environment and Datalake 🔗

It's recommended that you use terraform automation for speed and simplicity. The terraform quickstart can be found here: ⭘ GitHub - cloudera-labs/cdp-tf-quickstarts
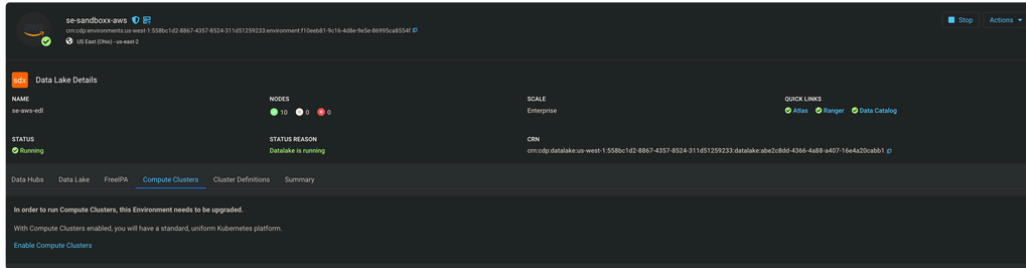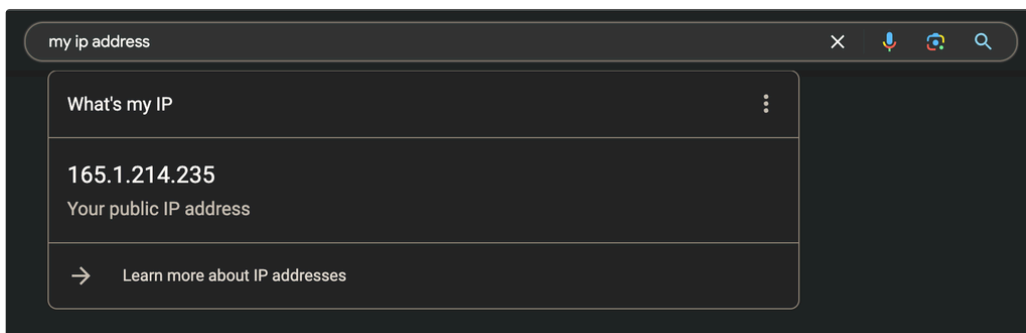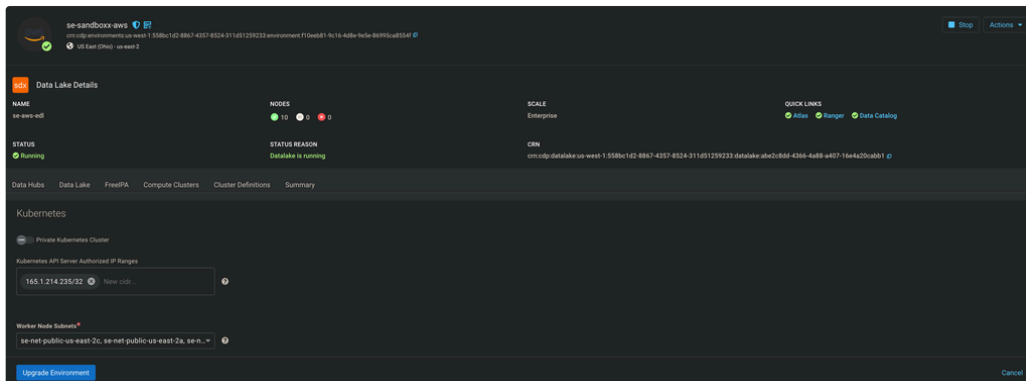
## 2. Initialize Compute Clusters 🔗

Once the environment is available you have to initialize the environment to use compute clusters. This operation will migrate the environment from v1 to v2 and create a default compute cluster (aka containerized datalake).

Option 1: UI

Navigate to the environment you deployed in the previous step and go to the `Compute Clusters` tab



Select `Enable Compute Clusters` on the screen. You will be prompted for a few options. We recommend keeping the `Private Kubernetes Cluster` option disabled. We also recommend using public subnets in the `Worker Node Subnets` field. Last, for the `Kubernetes API Server Authorized IP Ranges` enter the cidr for the VPN endpoint that you are connected to. If you are unsure of your endpoint IP you can search `My IP Address` on Chrome and it will show you the IP. Just add a `/32` to the end and it will work. Note that this IP whitelist is for access to the kubernetes cluster for administrative purposes. Once everything is set select `Upgrade Environment`





Option 2: CLI

Run the below command from the CDP CLI. This command is available in the public cli release - more details can be found here
🔗 initialize-aws-compute-cluster — CDP CLI 0.9.137 documentation

```
1  cdp environments initialize-aws-compute-cluster --cli-input-json file://convert-v2-env.json
```

Example convert-v2-env.json:

```json
{
    "environmentName": "avk-upgrade-cdp-env",
    "computeClusterConfiguration": {
        "privateCluster": false,
        "kubeApiAuthorizedIpRanges": [
            "165.1.214.236/16",
            "134.238.0.0/16",
            "208.127.0.0/16"
        ],
        "workerNodeSubnets": [
            "subnet-054399ed657492b02",
            "subnet-094a93473bca14029",
            "subnet-025cf012f047c7392"
        ]
    }
}
```
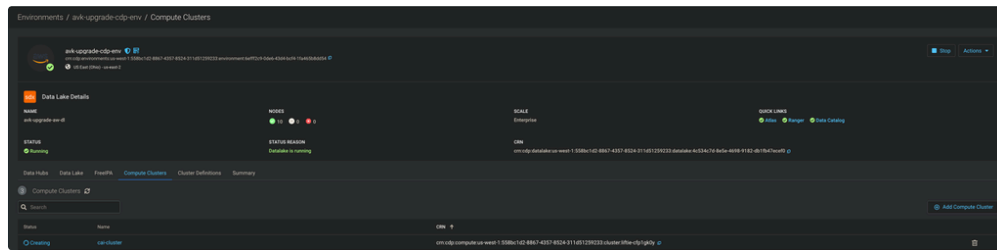
*Note:* Make sure you use the correct environmentName and workerNodeSunets for your environment created in step 1.  Also make sure to add the IP for the VPN endpoint you are connected to in `kubeApiAuthorizedIpRanges` so that you can connect to the kubernetes cluster.





## 3. Create New Compute Cluster for CAII Service 🔗

The CAII Service requires a dedicated compute cluster.  This can be created directly from the UI.

## 4. Create CAI Workbench and Model Registry 🔗

CAII also requires a CAI Workbench and Model Registry to use. These can be created in parallel to the compute cluster.





## 5. Create the CAII Service App 🔗

Once the Workbench, Registry, and Compute Cluster are ready you can deploy the CAII service app. This can currently be done through the UI or CLI

Option 1: UI

Go to `Cloudera AI` → `AI Inference Services` and select `Create New AI Inference Service`



In the configuration page enter a name for the Inference Service, select the correct environment and compute cluster created in Step 3, select the CPU and GPU instance shapes and select `Enable Public IP Address for Load Balancer`

**_Note:_** To create a model serving app you must have the EnvironmentAdmin and MLAdmin roles on the environment

**_Note:_** We recommend using the g5 gpu shape for the CAII service



Option 2: CLI

Run the below command from the CLI:

```
1  cdp ml create-ml-serving-app --cli-input-json file://create-serving-app-input.json
```



Example create-serving-app-input.json:

```
1   {
2       "appName": "avk-serving-app",
3       "environmentCrn": "crn:cdp:environments:us-west-1:558bc1d2-8867-4357-8524-
    311d51259233:environment:6efff2c9-0de6-43d4-bcf4-1fa465b8dd54",
4       "clusterCrn": "crn:cdp:compute:us-west-1:558bc1d2-8867-4357-8524-311d51259233:cluster:liftie-cfp1gk0y",
5       "provisionK8sRequest": {
6           "instanceGroups": [
7               {
8                   "instanceType": "m5.4xlarge",
9                   "instanceTier": "",
10                  "instanceCount": 1,
11                  "name": "",
12                  "ingressRules": [
13                      ""
14                  ],
15                  "rootVolume": {
16                      "size": 256
```

```
17                },
18                "autoscaling": {
19                    "minInstances": 0,
20                    "maxInstances": 2,
21                    "enabled": true
22                }
23            },
24            {
25                "instanceType": "g5.12xlarge",
26                "instanceCount": 1,
27                "rootVolume": {
28                "size": 256
29                },
30                "autoscaling": {
31                    "minInstances": 0,
32                    "maxInstances": 1,
33                    "enabled": true
34                }
35            }
36        ],
37        "environmentCrn": "crn:cdp:environments:us-west-1:558bc1d2-8867-4357-8524-
   311d51259233:environment:6efff2c9-0de6-43d4-bcf4-1fa465b8dd54",
38        "tags": [
39            {
40                "key": "owner",
41                "value": "akahan"
42            }
43        ],
44        "network": {
45            "plugin": "",
46            "topology": {
47                "subnets": [
48                    ""
49                ]
50            }
51        }
52    },
53    "usePublicLoadBalancer": true,
54    "skipValidation": true,
55    "loadBalancerIPWhitelists": [
56        ""
57    ],
58    "subnetsForLoadBalancers": [
59        ""
60    ],
61    "staticSubdomain": "avk-serving"
62 }
```

**Note:** Make sure to update the environmentCrn and clusterCrn with the environment created in step 1 and the compute cluster created in step 3

**Note:** To create a model serving app you must have the EnvironmentAdmin and MLAdmin roles on the environment

**Note:** We recommend using the g5 gpu shape for the CAII service

## 5.a Check CAII Service App Status 🔗

You can check the status of the CAII Service app install using command:

```
1   cdp ml list-ml-serving-apps
```



## 6. Deploy a Model Endpoint 🔗

When the serving app is complete, you can create an endpoint from a model imported from the Model Hub or from a model created in the CAI Workbench. The below is an example deployment of the wine-deploy model.

**Note:** To create a model endpoint you must have the EnvironmentUser and MLUser roles on the environment you are using

# 7. Cleanup of CAII and supporting Services 🔗

To delete the CAII service and supporting resources you must delete all of the CAI resources first before you can terminate the environment/datalake.  Failure to do so will result in orphaned cloud resources and continue to accrue cost after the environment is gone.

## 7.a Delete the CAII Service 🔗

To delete the CAII service you first need to get the service CRN. This can currently only be done with the cli.  You can do this with the command:

```
1  cdp ml list-ml-serving-apps
```



Then you can delete the serving app with the command:

```
1  cdp ml delete-ml-serving-app --app-crn <serving app crn>
```
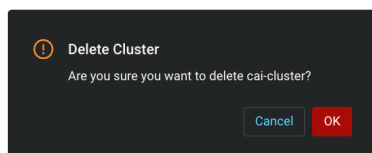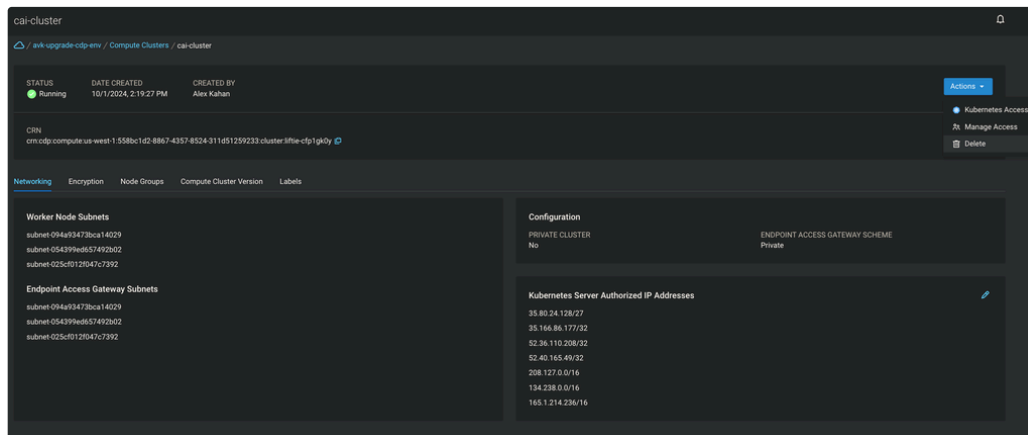


## 7.b Delete the CAI Workbench and Model Registry 🔗

While the CAII service app is deleting you can also delete the Workbench and Model Registry.  This can be done through the UI.

## 7.c Delete the CAII Compute Cluster 🔗

When the CAII service app is deleted you then have to do delete the compute cluster that was created in step 3.  This can be done through the UI

## 7.d Destroy the Environment, Datalake, and Cloud Prereqs 🔗

The last step is to delete the environment, datalake, and cloud resources.  If you used terraform to create the environment you should also use terraform to destroy it.  This will ensure that all resources on the cloud side are deleted.  If you created the environment any other way, make sure that all resources are deleted on the cloud side.

> ℹ️ Highlight important information in a panel like this one. To edit this panel's color or style, select one of the options in the menu.

📋 Related articles

📄 How to Enable Dataviz On Base With Knox SSO

📄 How To Connect Octopai To Hive/Impala on Base

📄 How to Configure Octopai Client on Platform9

📄 How to Configure Octopai Client on GCP

📄 How to Configure Octopai Client on AWS