# Step by Step install Helm chart
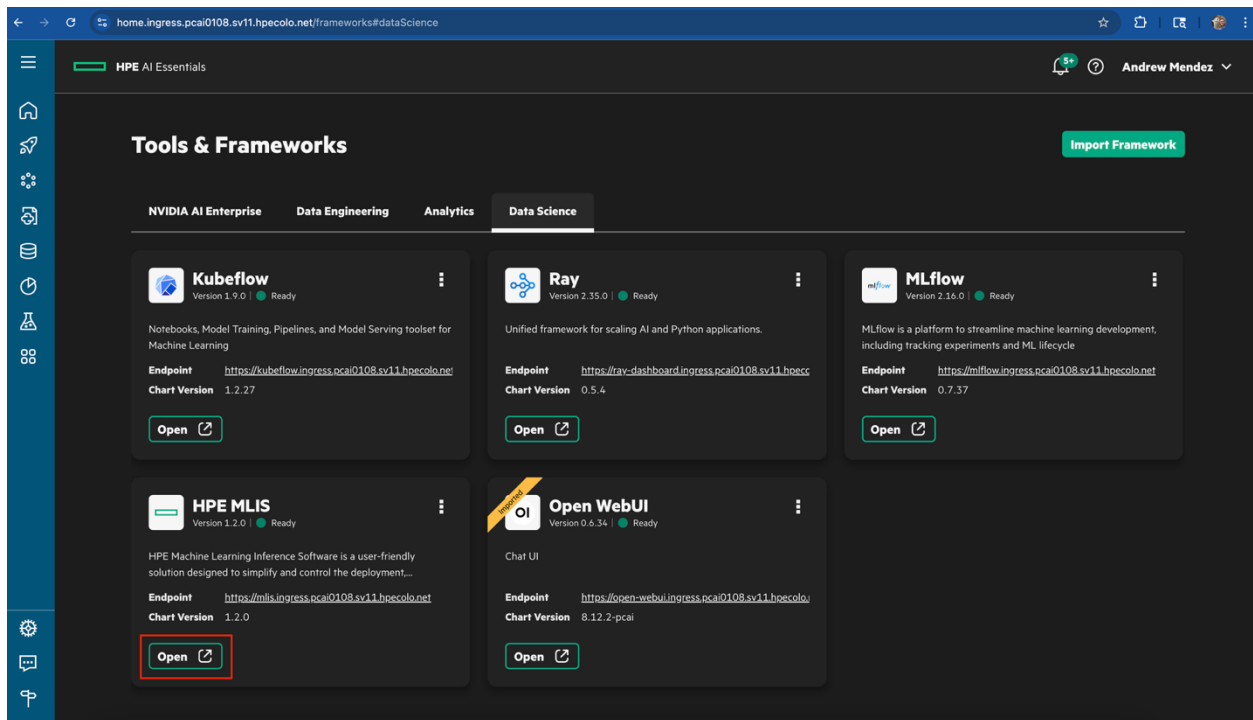
## Install on AIE 1.6

Part 1: Deploy Vista model on MLIS

Go to MLIS



Next create registry select add new registry

## Registries

Registries store your models and code.

Add new registry

| Registry name | | Last modified ^ | Type |
|---|---|---|---|
| huggingface-registry | ⋯ | 10 hours ago | openllm |
| local-s3-bentotest | ⋯ | 4 days ago | s3 |
| NGC | ⋯ | 2 months ago | ngc |

HPE MLIS

- Deployments
- Packaged models
- Registries
- API Tokens

mlis.ingress.pcai0108.sv11.hpecolo.net/ui/registries

Fill in the information needed to create a registry, this will hold your nvidia enterprise API key

# Add new registry

A registry stores information needed to access your models. Learn how to setup NGC registry.

**Name** ⍰

| name |
|---|

**Type** ⍰

| NGC ⌄ |
|---|

**API key** ⍰

| dummy 👁 |
|---|

**Org name** ⍰

| nim |
|---|

**Team name (optional)** ⍰

| nvidia |
|---|

**Endpoint (optional)** ⍰

| https://api.ngc.nvidia.com |
|---|

Cancel        **Create registry**

Next create a packaged model, name your model vista



Select your registry and list of supported models will populate, select vista model docker image

# Add new packaged model

A model is required for an inference deployment. Learn how to setup a model.

| Your model | **Storage** | Resources | Advanced (optional) |

**Registry** ⓘ

⁞  NGC                                                                    ngc ⌄

**NGC Supported Models** ⓘ

▢  vista3d                                                               vtest ⌄

**Image** ⓘ

nvcr.io/nim/nvidia/vista3d:latest

**Path (optional)** ⓘ

path

Cancel    Back    Next

Next set resources to the following screenshot. Make sure to manually change gpu to 1 and 1, this model only needs 1 l40s GPU

# Add new packaged model

A model is required for an inference deployment. Learn how to setup a model.

| Your model | Storage | **Resources** | Advanced (optional) |

ⓘ Requested resources are the minimum your packaged model needs to operate. You can set limits to handle spikes to manage additional traffic without affecting other nodes.

**Resource Template** ⓘ

| 🎞 gpu-small | ⌄ |

**CPU** ⓘ

| 2 | → | 6 |

**Memory** ⓘ

| 20Gi | → | 40Gi |

**GPU** ⓘ

| 1 | → | 1 |

Cancel　　Back　　Next

Next in the advanced settings, Add environment variable

DOMAIN_WHITELIST
[".*","http://.*","https://.*","http://.*:.*","https://.*:.*","file:///.*","*"]

# Add new packaged model

A model is required for an inference deployment. Learn how to setup a model.

| Your model | Storage | Resources | **Advanced (optional)** |
|---|---|---|---|

ⓘ   The following configuration values are optional. **Learn more.**

**Environment Variables** ⓘ

| DOMAIN_WHITELIST | [".*","http://.*","https://.*","http://.*:.*","https://.* |
|---|---|

**Add new**

**Arguments** ⓘ

```
ex: --arg --foo
```

Cancel    Back    **Create model**

Now lets create a deployment, name your deployment same as packaged model

# Create new deployment

A deployment is a running instance of a packaged model. Learn how to setup a deployment.

| **Deployment** | Packaged Model | Infrastructure | Scaling | Advanced (optional) |
|---|---|---|---|---|

**Deployment Name** ⓘ

```
vista-two
```

Cancel    **Next**

Select packaged model you just created.

# Create new deployment

A deployment is a running instance of a packaged model. Learn how to setup a deployment.

Deployment | **Packaged Model** | Infrastructure | Scaling | Advanced (optional)

## Which packaged model do you want to serve? ⑦

Select packaged model... ⌄

| Type to select a model | |
|---|---|
| ⬡ Qwen2.5-VL-32B-Instruct-AWQ | |
| ⬡ bento-taxi | |
| ⬡ bge-cpu | |
| ⬡ bge-large-en-v1.5 | |
| ⬡ chatterbox-tts | 4 versions |
| ⬡ kokoro-fastapi-cpu | 3 versions |
| ⬡ kokoro-fastapi-gpu | |
| ⬡ llama-3-1 | |
| ⬡ qwen3-8b | 2 versions |
| ⬡ vista | |
| ⬡ whisper-v3-turbo | 2 versions |

next  set auto scaling to fixed-1

# Create new deployment

A deployment is a running instance of a packaged model. Learn how to setup a deployment.

Deployment    Packaged Model    Infrastructure    **Scaling**    Advanced (optional)

Auto scaling targets template ⓘ

select an auto scaling template...    ⌄

| | |
|---|---|
| ◎ | fixed-1 |
| ◎ | fixed-2 |
| ◎ | scale-0-to-1-concurrency-3 |
| ◎ | scale-1-to-8-concurrency-3 |
| ◎ | scale-0-to-4-rps-10 |
| ◎ | scale-0-to-8-rps-20 |
| ◎ | scale-1-to-4-rps-10 |
| ◎ | custom |

Select Done, and wait for the model to deploy, this will take a few minutes.

When its deployed, copy URL, example shown here

| resumai-llm-server | ... | Ready | Serving | hugo-boulet-7c022924 | qwen3-8b v2 | https://resumai-llm-server-predictor-hugo-boulet-7c022924.ingress.pcai0108.sv11.hpecolo.net |
|---|---|---|---|---|---|---|

Next, lets make an API key

Select deployed deployment, go to users and click add new user api token

## Deployments

Deployments host the model and infrastructure that makes everything happen.

| Deployment name | | Status | Latest event | Namespace | Packaged model |
|---|---|---|---|---|---|
| bento-taxi | ⋯ | Paused | Paused | tanguy-pomas-1d2af612 | bento-taxi v1 |
| bge-cpu | ⋯ | Paused | Paused | tanguy-pomas-1d2af612 | bge-cpu v1 |
| bge-large-en-v1-5 | ⋯ | Paused | Paused | tanguy-pomas-1d2af612 | bge-large-en-v1.5 v1 |
| chatterbox-tts | ⋯ | Paused | Paused | tanguy-pomas-1d2af612 | chatterbox-tts v4 |
| kokoro-fastapi-gpu | ⋯ | Paused | Paused | tanguy-pomas-1d2af612 | kokoro-fastapi-gpu v1 |
| qwen3-8b | ⋯ | Paused | Paused | isabelle-steinh-74bc67b1 | qwen3-8b v2 |
| resumai-llm-server | ⋯ | Ready | Serving | hugo-boulet-7c022924 | qwen3-8b v2 |
| vista | ⋯ | Paused | Paused | andrew-mendez-fa786398 | vista v1 |
| whisper-v3-turbo | ⋯ | Paused | Paused | tanguy-pomas-1d2af612 | whisper-v3-turbo v2 |

**resumai-llm-server**

Ready

General  Timeline  Advanced  **Users**

| User/Token | Status/Expiration | Actions |
|---|---|---|
| There are no user tokens. | | |

Add new user API token

Can set the role to whatever, I usually do  admin

# Create new token

Access tokens enable you to control who can use protected deployments.

**Which deployment do you want to create a token for?** ⓘ

resumai-llm-server ⌄

**Select 1 or more users**

1 user selected ⌄

**Selected users**

admin ✕

**Description of this token** ⓘ

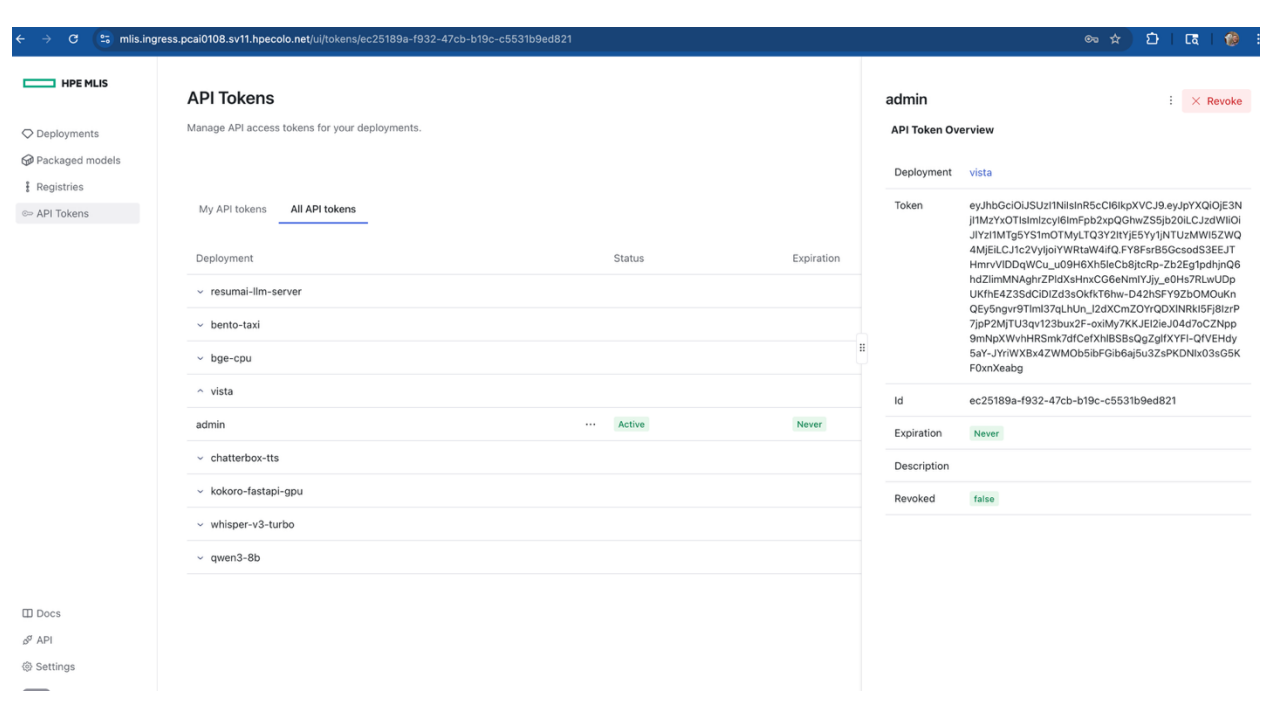Provide a description of this token

**When should this expire?** ⓘ

YYYY-MM-DDTHH:MM:SSZ

Quick selects:  30 days ,  60 days ,  90 days ,  120 days ,  Never

Cancel    Create

When the API is created, Select API Tokens, select all API tokens, and copy the API token shown on the right



You are ready to deploy the helm chart

Pre-req; you will need the tgz file of the helm chart ready to upload.


Part 2: Install Helm chart

GO to Tools & Frameworks > Import Framework


Add name, description, and logo

Drag .tgz of helm chart in UI, you should see this next

Add namespace



Now in the helm chart UI, only change the VISTA3D_SERVER and the VISTA3D_API_KEY

```
frontend:
  image: mendeza/vista3d-frontend-helm:v1.0.1 # "mendeza/vista3d-frontend:v1.0.8"
  imagePullPolicy: Always
  port: 8501
  env:
    VISTA3D_SERVER: "<REPLACE_ME>.${DOMAIN_NAME}"
    IMAGE_SERVER: "https://vista3d-image-server.${DOMAIN_NAME}"
    VISTA3D_IMAGE_SERVER_URL: "https://vista3d-image-server.${DOMAIN_NAME}"
    EXTERNAL_IMAGE_SERVER: "https://vista3d-image-server.${DOMAIN_NAME}"
    VISTA3D_API_KEY: ""    # <-- set this manually
```

Wait until its deployed, if its deployed successfully, you should see Open Button.