How to add models the pre-deployed NIM model to MLIS

## Create NGC Registry

Ideally to access NIM models you will need to have a NGC key.

But these NIM models under Tools & Frameworks >> NVIDIA AI Enterprise, are already available on the system for use.

There is an in-built registry which is valid only for these models which come along with the system.

Thus, you do not need to register for NGC key on MLIS for these models.

Do verify that there is an entry under Registries as show below.



## Packaged Model

Go to the NVIDIA AI Enterprise Tab.

Select the model you want to deploy.

Click on **View Details** for the model.

Make note of the parameters marked in red as they will be required while packaging the model for being hosted on MLIS.



Now back to the MLIS page.

**Name:** Enter a name which helps you to identify the model easily, preferably model name.

**Description:** Favourably add the actual model's name. So, its easy for reference later.

In the Storage tab,

**Registry:** Select the **NGC** registry.

As soon as you select that the **NGC Supported Models will pe pre-populated with values:** Select the model you want to deploy from the drop down.

As soon as you select the model from the drop down, the fields **Image** and **Path** will be auto filled. The values should match the values from the model's **View Details** page captured earlier.

## 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** ⓘ

| ⬡  meta/llama3-8b-instruct | ve4 ⌄ |

**Image** ⓘ

| 10.19.111.55/ezmeral-common/nvcr.io/nim/meta/llama3-8b-instruct:1.0.3 |

**Path (optional)** ⓘ

| pvc://models-pvc |

Cancel    Back    **Next**

**Resource Template:** Choose appropriate template or define a custom one based on the model requirements. Refer the documentation for help.

# 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-tiny | ⌄ |

**CPU** ⓘ

| 1 | → | 1 |

**Memory** ⓘ

| 40Gi | → | 60Gi |

**GPU** ⓘ

| 1 | → | 1 |

Cancel    Back    **Next**

---

**Advanced (optional)**

**Environment Variables:** Add the required ones. Refer documentation

AIOLI_PROGRESS_DEADLINE: 3600s

The above indicates deadline for downloading the model: default value is 30 mins. Thus, it is most commonly set since for some models this can go beyond that time frame.

## 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 ⓘ

| AIOLI_PROGRESS_DEADLINE | 3600s |

[Add new]

### Arguments ⓘ

| ex: --arg --foo |

[Cancel] [Back] [**Create model**]

---

**HPE MLIS**

♡ Deployments
⊕ Packaged models
⫶ Registries
⇔ API Tokens

### Packaged models

Models are the things that you want to get out to the real world.

[🔎 Search]   [Add new model]

| Model name | | | Status | Last modified ^ | Description | Registry used | Path |
|---|---|---|---|---|---|---|---|
| llama-3-8b-nim | v1 | ⋯ | Staged | 55 seconds ago | meta/llama-3.1-8b-instruct | NGC | pvc://models-pvc |
| llama3-70b | v1 | ⋯ | Staged | 19 days ago | | HF | openllm://meta-llama/Meta-Llama-3-70B-Instruct |
| llama3-8b | v1 | ⋯ | Staged | 19 days ago | | HF | openllm://meta-llama/Meta-Llama-3-8B-Instruct |
| ⌄ blip2-serving | v6 | ⋯ | Staged | 2 months ago | | none | |
| blip2-ac-aws | v1 | ⋯ | Staged | 2 months ago | | aws-source-images | s3://poc-mercedes-gp/model/v1 |

## Deployments

**Deployment Name:** Provide a name of your choice. just ensure it needs to be all in lowercase and you can only use special characters "." and "-".

## 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** ⓘ

llama-3-8b-nim

Cancel   **Next**

**Packaged Model**: Select the packaged model created in previous step.

## 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?** ⓘ

◻ llama-3-8b-nim   ⌄

Cancel   Back   **Next**

**Scaling:** Choose option as applicable, for testing purpose we go with 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** ⓘ

⊚  fixed-1                                                                      ⌄

**Minimum instances** ⓘ                          **Maximum instances** ⓘ

1                                                  1

**Auto scaling target** ⓘ

rps ⌄    5

Cancel    Back    **Next**

**Advanced:** Fill in as applicable (normally left blank)

## 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)**

ⓘ   The following configuration values are optional. Learn more.

**Environment Variables** ⓘ
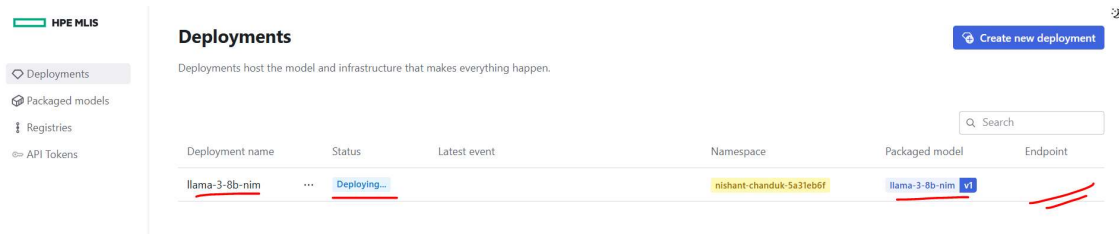
| field | value |

**Add new**
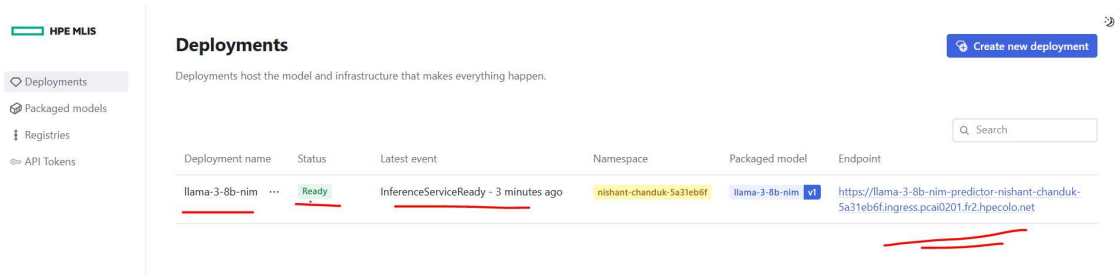
**Arguments** ⓘ

ex: --arg --foo

Cancel    Back    **Done**

The above snapshot shows the model is getting deployed which will take some time.

Once the model is in **Ready** state you would see a valid endpoint.



Now how to we access the model anywhere within or outside the system.
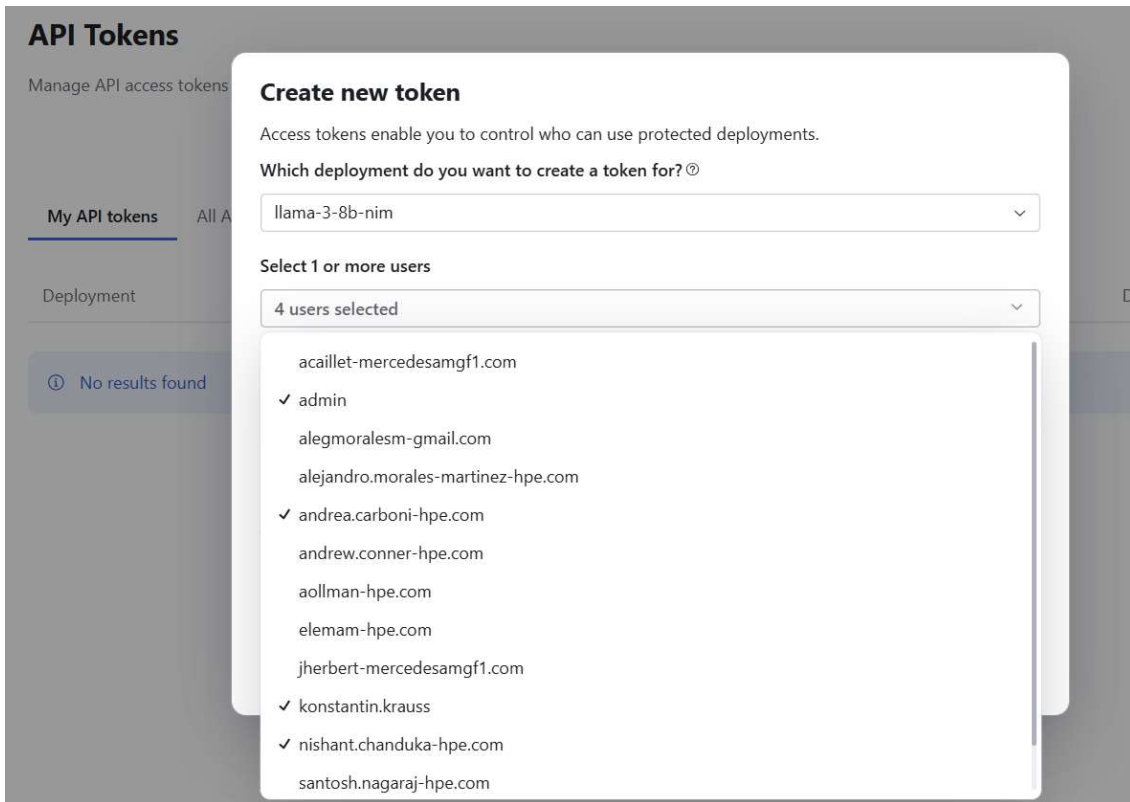
Go to **API Tokens**

Click on **Create new API access Token**



The dropdown menu **Which deployment do you want to create a token for?** Will list the models you see in the **Deployments** tab only.

Select your model.

**Select 1 or more users:** Select the users you want to be able to use the token.

**API Tokens**

Manage API access tokens

My API tokens   All A

Deployment

ⓘ No results found

**Create new token**

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

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

llama-3-8b-nim                                              ⌄

**Select 1 or more users**

4 users selected                                           ⌄

acaillet-mercedesamgf1.com
✓ admin
alegmoralesm-gmail.com
alejandro.morales-martinez-hpe.com
✓ andrea.carboni-hpe.com
andrew.conner-hpe.com
aollman-hpe.com
elemam-hpe.com
jherbert-mercedesamgf1.com
✓ konstantin.krauss
✓ nishant.chanduka-hpe.com
santosh.nagaraj-hpe.com

Add **Description of this token**

Select **When should this expire?**

# Create new token

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

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

> llama-3-8b-nim ⌄

**Select 1 or more users**

> 4 users selected ⌄

**Selected users**

> admin ✕    andrea.carboni-hpe.com ✕    konstantin.krauss ✕    nishant.chanduka-hpe.com ✕

**Description of this token** ⓘ

> access llama-3-8b NIM model

**When should this expire?** ⓘ

> 2025-05-28T05:12:04.900Z

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

Cancel    Create

---

Copy the Created token