Ismail Syed          Specialist Leader EMEA - Data Science, Vector & ML          Oracle

# Deploying LLM from AI Quick Actions Catalog

## Referenced Documentation

https://docs.oracle.com/en-us/iaas/data-science/using/ai-quick-actions-model-deploy.htm
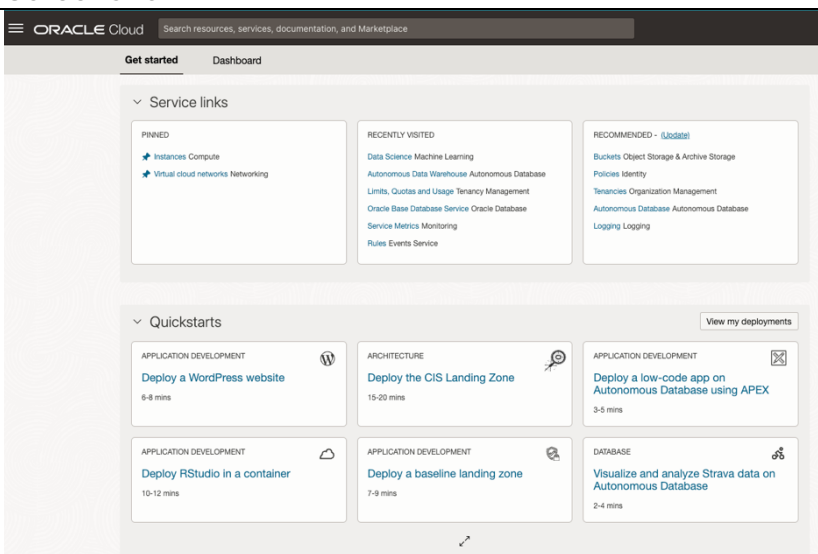
## Description

You can create a Model Deployment from the foundation models with the tag Ready to Deploy in the Model Explorer, or with fine-tuned models. When you create a Model Deployment in AI Quick Actions, you're creating an OCI Data Science Model Deployment, which is a managed resource in the OCI Data Science Service. You can deploy the model as HTTP endpoints in OCI.
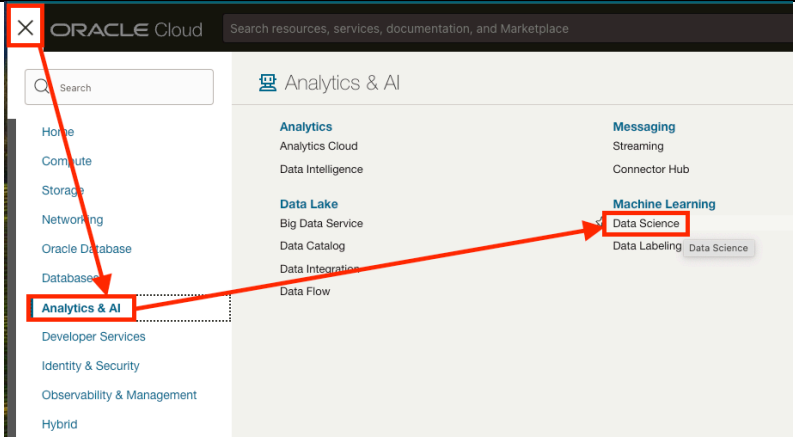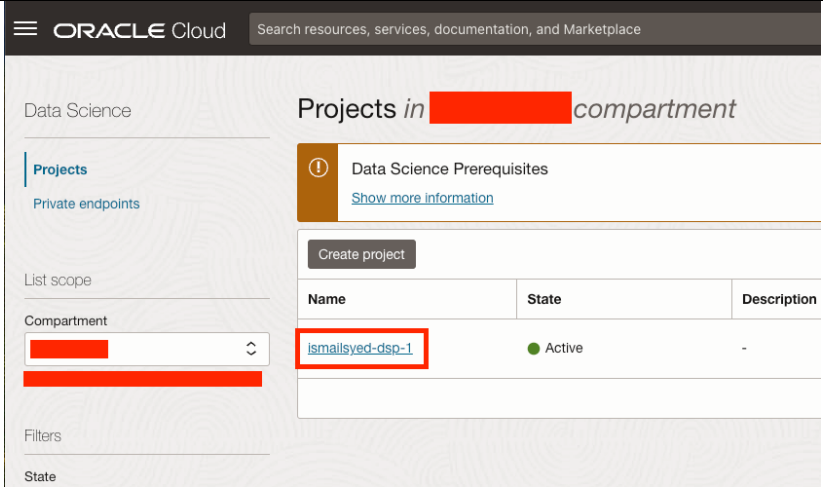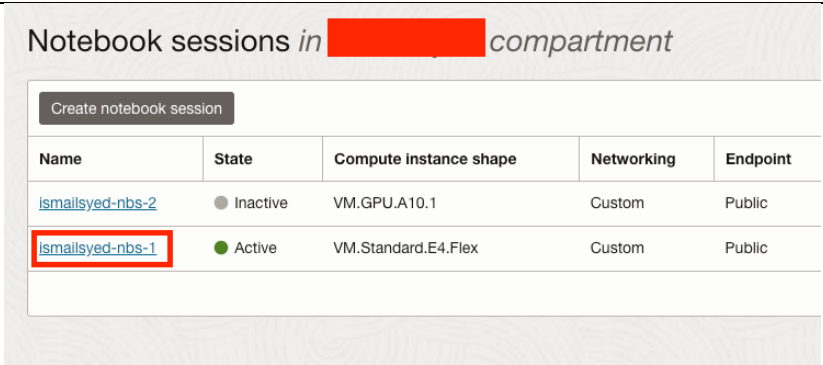
## Pre-Requisites

- Implement the required policies - https://docs.oracle.com/en-us/iaas/data-science/using/ai-quick-actions-set-up.htm
- Ensure you have your OCI Data Science GPU service limits raised for the GPU Shapes you plan to use. This can be done from OCI Console.
- Provisioned OCI Data Science Project and Notebook Session (Must be deactivated and reactivated if created before the policies where implemented).
- OCI Log Group & Log Created (Optional)

## Guide

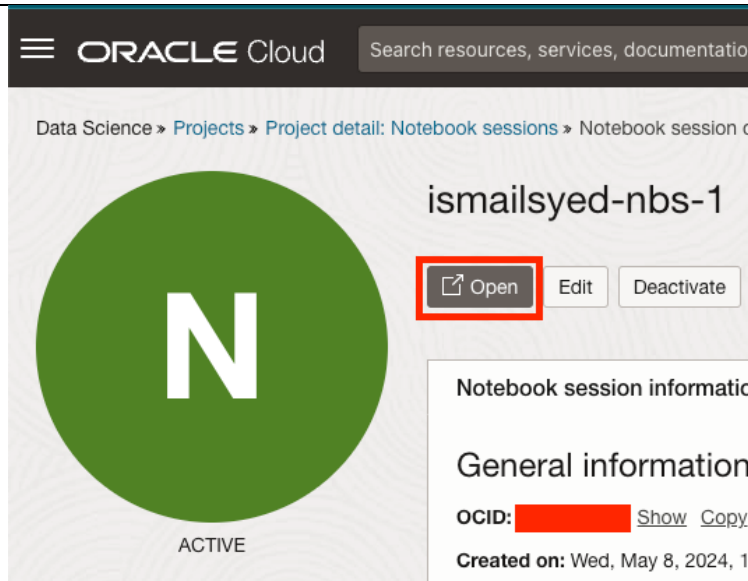| Step | Screenshot |
|---|---|
| Login to the Cloud Console.<br><br>cloud.oracle.com |  |

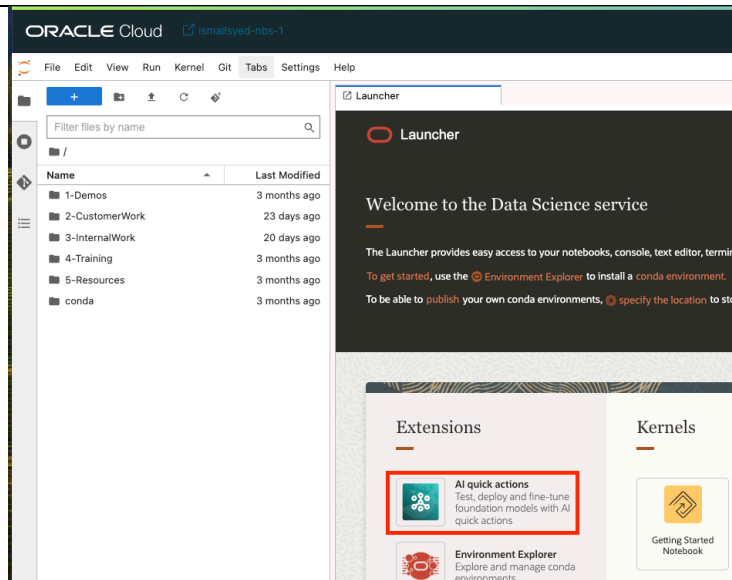| | |
|---|---|
| Navigate to your Data Science Projects.<br><br>**OCI Menu > Analytics & AI > Data Science.** |  |
| Open up your existing Data Science Project. |  |
| **Click on** your existing Data Science Notebook Session.<br><br>Note – This does not have to be a GPU Shape. |  |

***Click on Open.***

This will open up your Data Science Notebook Session.

***You will have to reauthenticate.***



If the policies within the pre-requisites have been implemented correctly you should be able to open up the AI Quick Actions Extension within the Launcher.

***Click AI quick actions.***



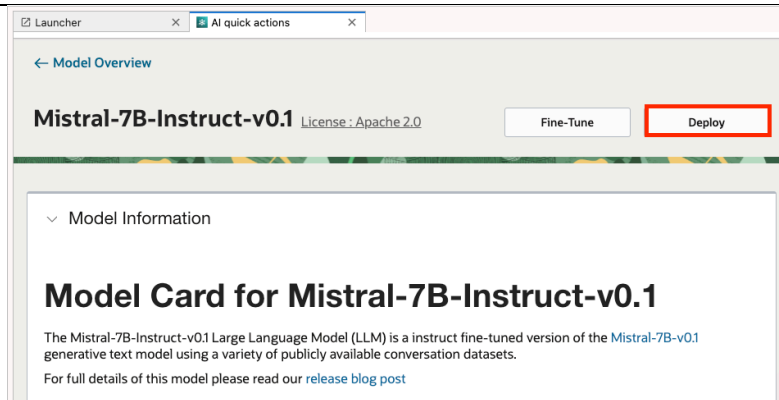The AI Quick Actions Catalog will be displayed.

In this tutorial we will deploy the ***Mistral-7B-Instruct-v0.1 model***.

***Select this model.***

| | |
|---|---|
| We will then be taken to the Model card which gives you a bit of information about the Model.<br><br>***Click on Deploy.*** |  |
| ***Enter:***<br><br>Deployment Name<br><br>Compute Shape – leave the recommend shape selected (VM.GPU.A10.1)<br><br>As an optional step, select your Log Group and Log to store your logging information.<br><br>***Click Deploy*** |  |
| The Model will then enter the ***Creating*** Lifecycle state.<br><br>You can continue to do other tasks within OCI Data Science while the deployment is happening. |  |

| | |
|---|---|
| Once the deployment is complete the lifecycle state will update to ***Active***. |  |
| If you scroll down, you will see a sandpit environment to start testing your model along with editing the Model Parameters on the right-hand side. |  |
| Head back to the OCI Console within our Data Project.<br><br>***Navigate to Models.***<br><br>Here you will see our new LLM in the Model Catalog. |  |
| We can then ***navigate to Model Deployments.***<br><br>Here we can see our deployed Model.<br><br>Select our deployed model. |  |

| | |
|---|---|
| ***Navigate to Invoking your Model under Resources.***<br><br>Here we can view the deployed Model Endpoint along with some sample code to make inference against the model. |  |
| Assuming you have the necessary policies and dynamic groups in place to use a resource principal.<br><br>Please refer to the following Notebook which provides some sample code.<br><br>***00-invoke-deployed-model.ipynb***<br><br>Remember to replace the Model Endpoint with your model along with editing the JSON Body with the parameters you would like sent to the LLM. |  |
| If you are unable to use Resource Principals, then use the code provided under the ***'Invoking your Model'*** section on your Model Deployment Page. ||