# Vector Similarity Search with Redis Enterprise Cloud on AWS

## Learning objectives

In this workshop, you will learn
- How to configure Redis Enterprise as a vector database?
- How to invoke embedding and LLM models programmatically using the Sagemaker Jupiter notebook
- How does vector similarity search work with Redis using a sample dataset?
- How do you query vectors stored as Hash/JSON in Redis using search in Redis?

## Pre-requisites

Following are some of the skillsets required to complete this training/workshop successfully:
1. Comfortable with AWS management console and cloud formation templates
2. Good in Python as the exercise and coding will be done in Python only
3. General idea of how Machine Learning works like training, testing and deployment of models
4. No prior knowledge of Redis Enterprise is required

## Target Personas

Developers, Data architects, Devops, SAs, Team lead and other technical stakeholders who have coding skills

## Duration

1-1.5 hours

## Tasks

The entire workshop is divided into several labs the details of which are mentioned below and also present in the following Jupyter notebook:

https://github.com/bestarch/redis-vss-getting-started/blob/main/vector_similarity_with_redis.ipyn
b


# Infrastructure set-up

## Provision Sagemaker Jupiter instance

Log in using either the root user or the IAM user. In case you are login in using an IAM user,
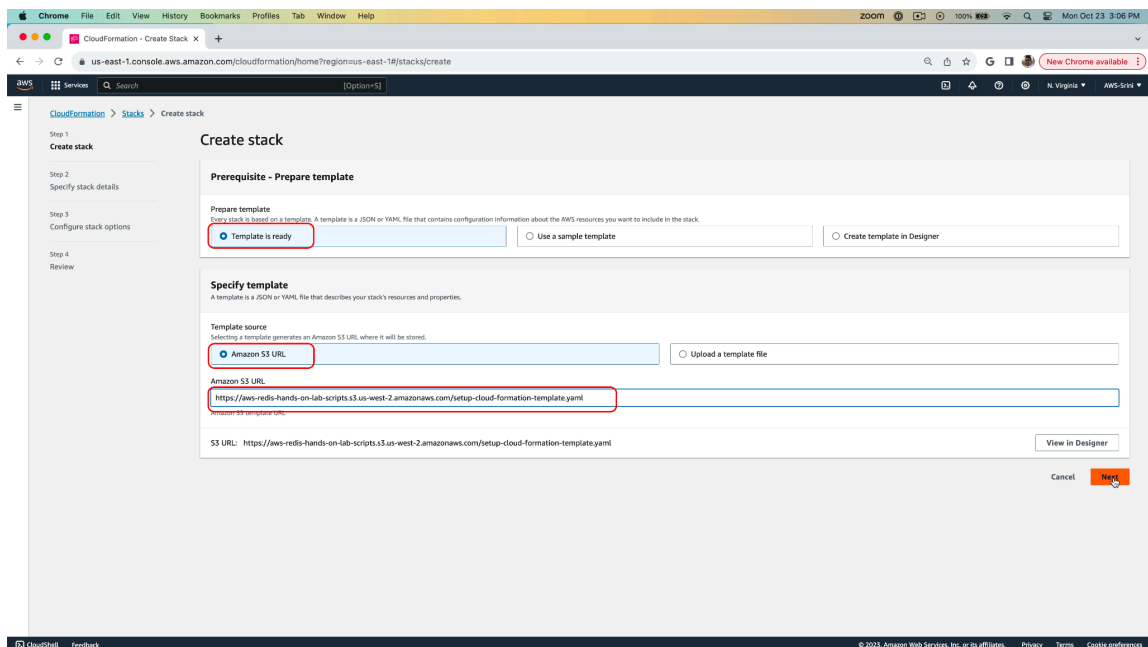make sure the following permissions are enabled for you:
- AdministratorAccess
- AmazonEC2FullAccess
- AmazonElasticContainerRegistryPublicFullAccess
- CloudWatchFullAccess
- IAMUserChangePassword
- PowerUserAccess
- AWSCloudFormationFullAccess
- AmazonSageMakerFullAccess

Before starting the lab exercises, we will first create a Sagemaker notebook instance. We will
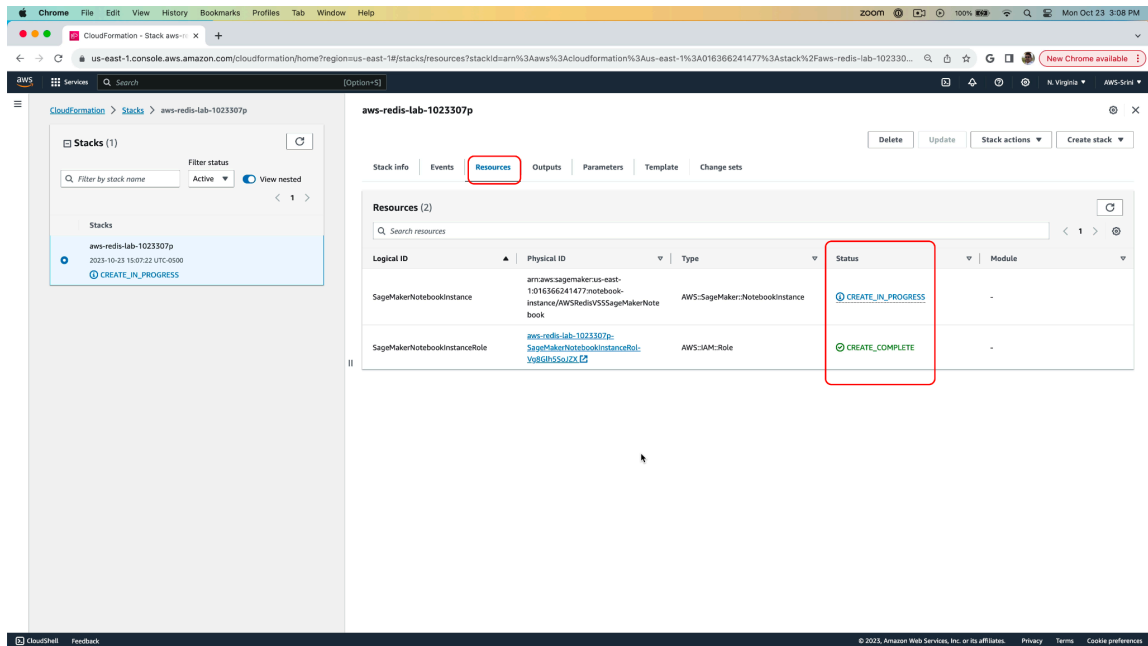use the CloudFormation template to do this.
Follow these steps:
1. Open the AWS management console and navigate to the CloudFormation service.
2. Create a new stack from the following template:
https://aws-cloud-formation-template-sagemaker.s3.ap-south-1.amazonaws.com/cloud-f
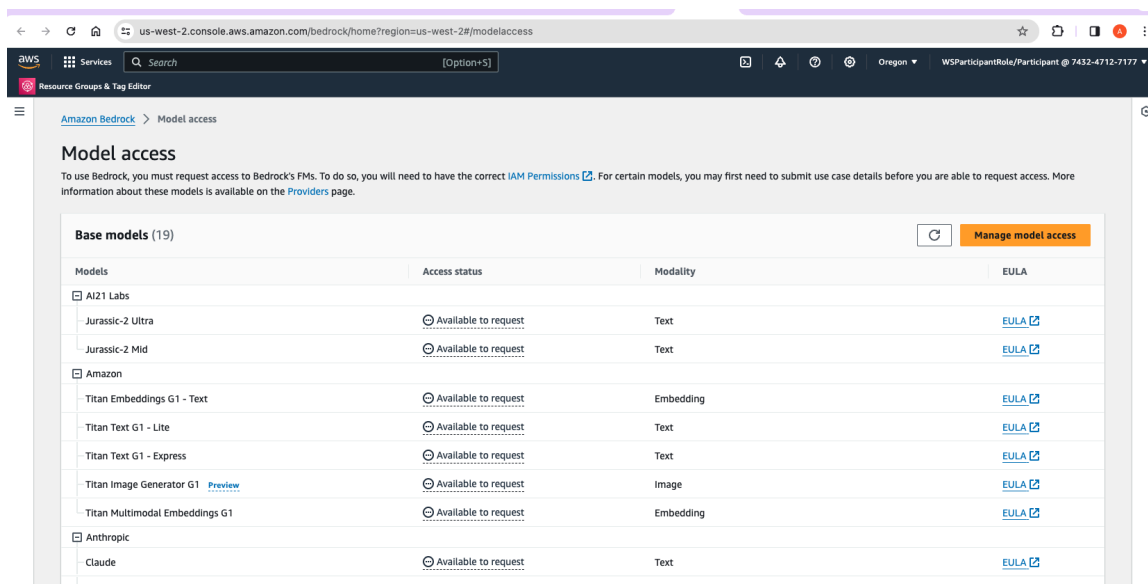ormation-template.yaml

3. Wait for a few minutes before the resources are ready.



4. Now we have created 2 resources:
   - A Sagemaker notebook instance
   - Suitable IAM roles (to access Sagemaker & Bedrock APIs) needed to complete our exercise using the notebook instance created above
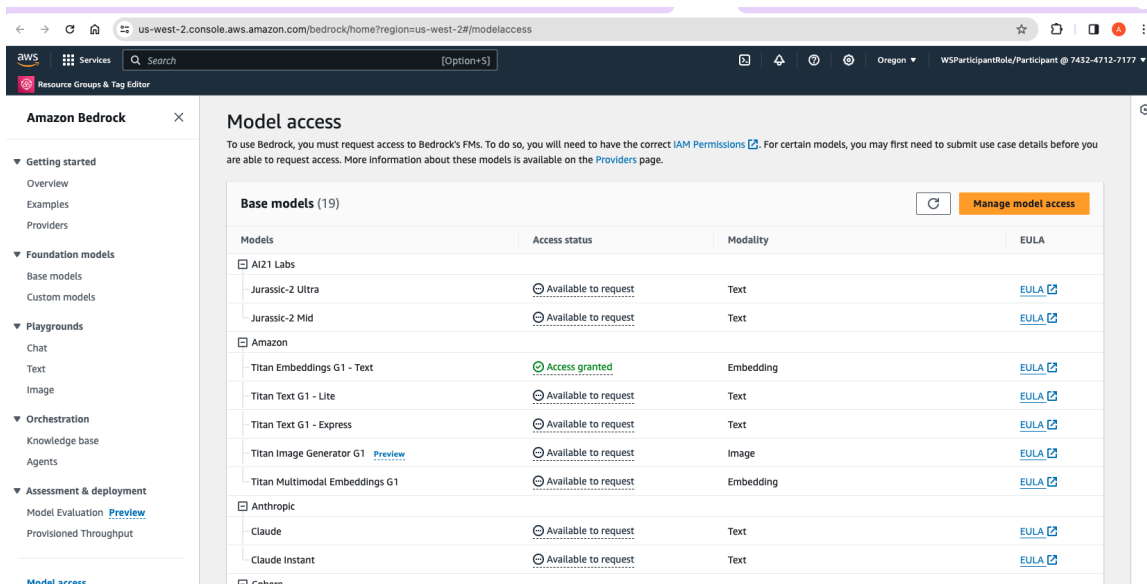
   But still, you can't use Amazon Bedrock FM at this time. To access these FMs you need to request them from your AWS management console.
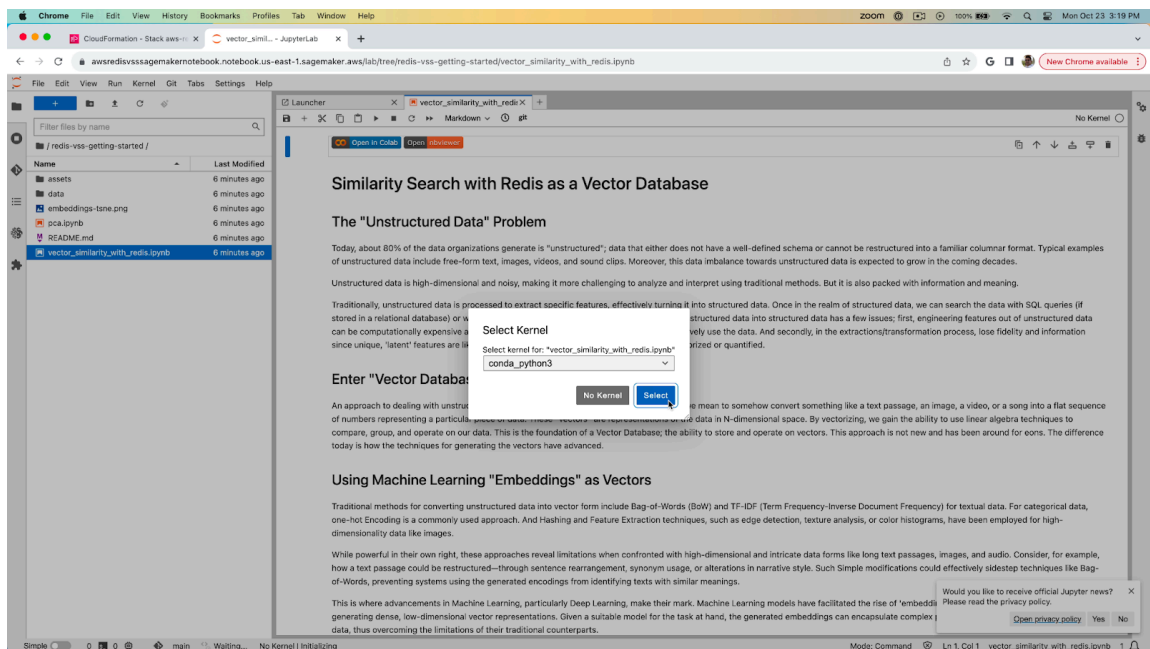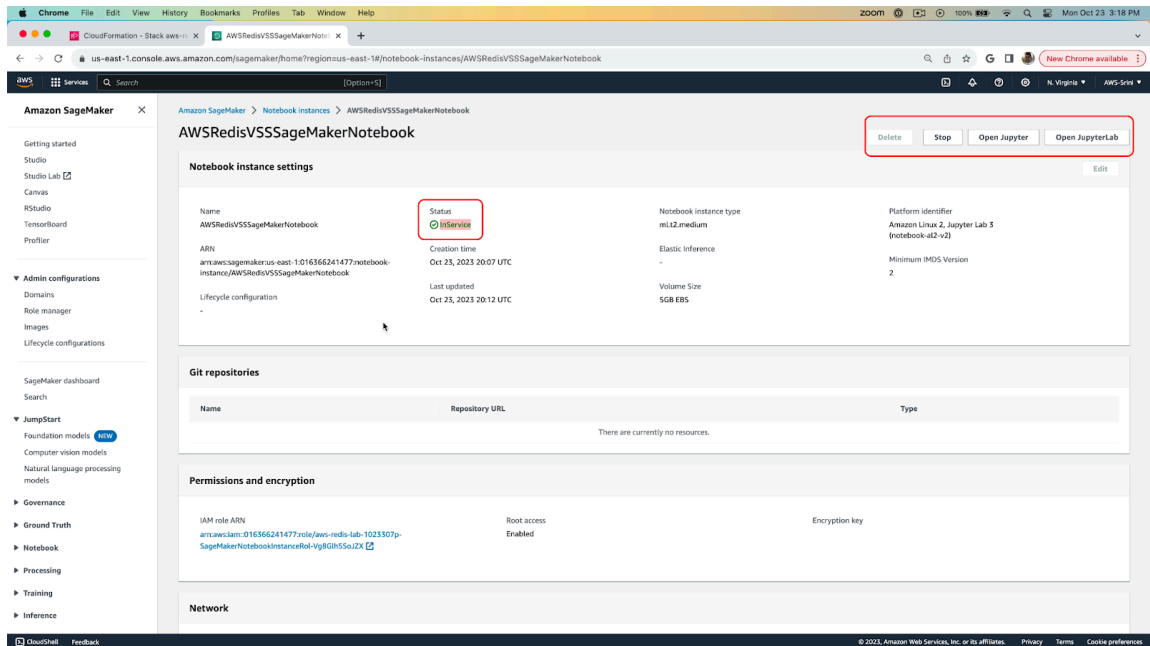
We will use the following FM for our vector embeddings:
`'amazon.titan-embed-text-v1'`

5. Goto Bedrock service from the AWS management console and click on the "Base Model" link present in the left navigation menu

6. Locate the 'Amazon Titan Embeddings G1 - Text' model and raise a request for this model by clicking on the 'Available to request' link.

7. This takes a few seconds and if everything is right, this model will be immediately available. Now we will use this model in our coding exercise.



8. Now that all the prerequisites are met, let's navigate to the notebook instance we created above using the Sagemaker service

# Lab 1: Load "Bikes" Dataset

To investigate Vector Similarity, we'll use a subset of the Bikes dataset, a relatively simple synthetic dataset.

The dataset has 11 bicycle records in a JSON file named `bikes.json` and includes the fields `model`, `brand`, `price`, `type`, `specs`, and `description`.
The `description` field is particularly interesting for our purposes since it consists of a free-form textual description of a bicycle.

For hands-on exercise, navigate to the 'Lab 1' section of the Jupyter notebook.

## Lab 2: Store the 'bikes' dataset in Redis

At this point, we have loaded our **'bikes'** dataset in memory.
Now, let's configure **Redis** and store this dataset as a collection of JSON documents. For this purpose, we will use the **Redis Enterprise Cloud** database.
In the next section, we will leverage the same Redis database to store vector embeddings as well.

***Note***: *In case Redis Enterprise Cloud is not available, use Redis Stack open-source.*
This can be done by installing Redis Stack Server using:
*!curl -fsSL*
*https://packages.redis.io/redis-stack/redis-stack-server-6.2.6-v7.focal.*
*x86_64.tar.gz -o redis-stack-server.tar.gz*
*!tar -xvf redis-stack-server.tar.gz*
*!./redis-stack-server-6.2.6-v7/bin/redis-stack-server --daemonize yes*
*oss = True*

For hands-on exercise, navigate to the 'Lab 2' section of the Jupyter notebook.

## Lab 3: Vectorize the 'description' field inside JSON documents

Now, let's vectorize the **'description'** field present inside the JSON documents and finally load these vector embeddings inside the same document.
We will first collect all the Redis keys for the bikes.
We'll use the keys as a parameter to the `JSON.MGET` command, along with the JSONPath expression `$.description` to collect the descriptions in a list. We will iterate through this list to create the embeddings.

To create the vector embeddings we can leverage:

1. **Recommended**: Use AWS Bedrock API and invoke the Titan embedding model
2. **Alternative**: Use the SentenceTransformers framework to generate embeddings for the bike's descriptions. Sentence-BERT (SBERT) is a BERT model modification that produces consistent and contextually rich sentence embeddings. SBERT improves tasks like semantic

search and text grouping by allowing for efficient and meaningful comparison of sentence-level semantic similarity.

Both of these options are discussed in detail in the notebook.
For hands-on exercise, navigate to the '**Lab 3**' section of the Jupyter notebook.

## Lab 4: Execute queries & verify the data in Redis

At this point, we should have created all the vector embeddings and loaded them into our Redis database.
Next, we will perform the following tasks in this lab:
- Create an index to make the bikes searchable
- Write queries to search the documents
- Semantic searching using VSS query, hybrid query & range query

For hands-on exercise, navigate to the '**Lab 4**' section of the Jupyter notebook.

## Final: Wrap-up

In this lab exercise, we learned how Redis Enterprise provides powerful search capabilities over structured and unstructured data. Redis support for vector data can enrich and enhance the user's search experience. Although we focused on generating embeddings for unstructured data, the vector similarity approach can equally be employed with structure data, as long as a suitable vector generation technique is used.

The references below can help you learn more about Redis search capabilities:

- https://redis.com/solutions/use-cases/vector-database/
- https://redis.io/docs/stack/search/
- https://redis.io/docs/stack/search/indexing_json/