

Table of contents

1. Ba	ckground	1
2. Arc	chitecture	1
2.1.	High level apps integration architecture	1
2.2.	. How it works	2
3. Sol	lutions component	2
	Amazon Bedrock	
	. Amazon Aurora PostgreSQL	
	. AWS Lambda	
3.4.	. Amazon API Gateway	3
3.5.	Amazon S3	3
4. AW	VS Account Configuration	3
	AWS workload account	
5. AW	VS workload configuration	4
5.1.	Foundation model	4
5.2.	Bedrock foundation model availability	4
	5.2.1. Activate Foundation model	5
5.3.	. CloudFormation	7
	5.3.1. CloudFormation parameters	7
	5.3.2. CloudFormation execution	8
5.4.	. Database configuration	10
	5.4.1. Query editor setup	10
	5.4.2. SQL script extension on PostgreSQL	11
5.5.	Knowledge base configuration	12
	5.5.1. Create knowledge base	12
	5.5.2. Setup knowledge base integration	18
5.6.	. Lambda configuration	20
	5.6.1. Lambda Layer Boto3	
	5.6.2. Lambda Layer JWT	
	5.6.3. Lambda token configuration	
	5.6.4. Lambda prompt configuration	25
6. Ge	nAl test API	29
6.1.	Get token	29
6.2.	Prompt test	30



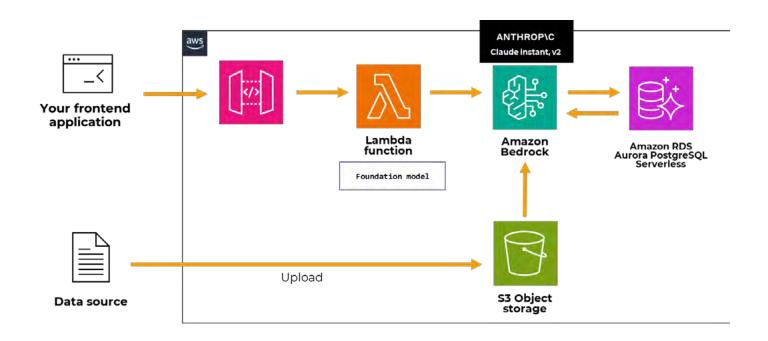
1. Background

AWS Bedrock offers a transformative approach to utilizing generative AI (GenAI) in business and technology, providing a robust, scalable, and cost-effective platform for innovation. The power of AWS Bedrock lies in its ability to streamline and democratize access to cutting-edge GenAl models, making sophisticated Al capabilities available to a broader range of industries and developers. This accessibility fosters innovation, as organizations can leverage advanced AI without the need for extensive infrastructure or expertise. By integrating GenAl into AWS Bedrock, businesses can automate complex processes, enhance customer experiences, and drive decision-making with data-driven insights, thereby staying competitive in an increasingly Al-driven market.

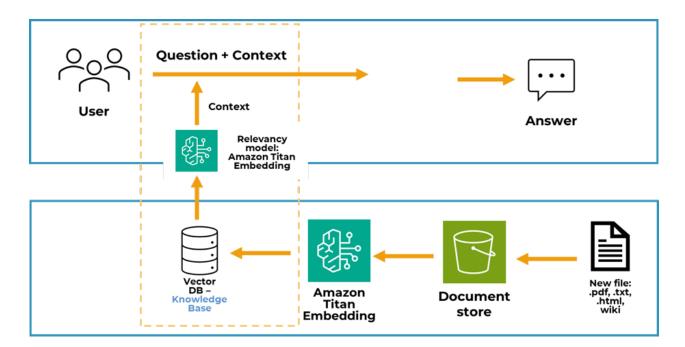
Moreover, AWS Bedrock ensures reliability and security, critical factors when handling sensitive data and deploying AI at scale. With AWS's proven track record in cloud services, Bedrock provides a seamless and secure environment for developing, training, and deploying AI models. This integration allows businesses to focus on innovation rather than infrastructure management. Additionally, AWS Bedrock's flexibility supports various AI use cases, from natural language processing to image generation, enabling organizations to tailor AI solutions to their specific needs. Embracing AWS Bedrock for GenAI thus not only enhances operational efficiency but also drives strategic growth by unlocking the full potential of artificial intelligence.

2. Architecture

2.1. High level apps integration architecture



2.2. How it works



- a. A document is broken up into chunks of text. The chunks are passed to Titan Embeddings to be converted to vectors. The vectors are then saved to the vector database.
- b. The user submits a question.
- c. The question is converted to a vector using Amazon Titan Embeddings, then matched to the closest vectors in the vector database.
- d. The combined content from the matching vectors + the original question is then passed to the large language model to get the best answer.

3. Solutions component

3.1. Amazon Bedrock

AWS Bedrock is a fully managed service from Amazon Web Services (AWS) that enables users to build and scale generative AI applications with foundation models (FMs) from AI21 Labs, Anthropic, and Stability AI. It simplifies access to these advanced models, offering APIs to integrate them into various applications. AWS Bedrock is designed to support tasks such as text generation, image creation, and data summarization, and it allows customization and fine-tuning of these models according to specific needs.

3.2. Amazon Aurora PostgreSQL

AWS Aurora PostgreSQL is a fully managed relational database service by Amazon Web Services (AWS) that is compatible with PostgreSQL. It combines the performance and availability of high-end commercial databases with the simplicity and cost-effectiveness of open-source databases. Aurora PostgreSQL provides features such as automated backups, replication, and seamless scaling of storage and compute resources. It is designed for enterprise-grade performance and reliability, supporting complex queries and large-scale database applications while ensuring high availability with multi-AZ (Availability Zone) deployments. Aurora also integrates with other AWS services, offering a robust solution for applications requiring a powerful and scalable PostgreSQL-compatible database.



3.3. AWS Lambda

AWS Lambda is a serverless compute service provided by Amazon Web Services (AWS) that allows users to run code without provisioning or managing servers. It automatically scales applications by running code in response to events, such as changes in data, shifts in system state, or user actions. With AWS Lambda, users pay only for the compute time consumed, with billing based on the number of requests and the duration of code execution. It supports various programming languages, including Node.js, Python, Java, and Go, enabling developers to build and deploy applications quickly and efficiently while focusing solely on writing code, without the overhead of managing infrastructure.

3.4. Amazon API Gateway

AWS API Gateway is a fully managed service that allows developers to create, publish, maintain, monitor, and secure APIs at any scale. It acts as a "front door" for applications to access data, business logic, or functionality from backend services such as workloads running on Amazon EC2, AWS Lambda, or any web application. API Gateway handles all the tasks associated with accepting and processing up to hundreds of thousands of concurrent API calls, including traffic management, authorization and access control, monitoring, and API version management. It supports RESTful APIs and WebSocket APIs, enabling real-time two-way communication applications. By using API Gateway, developers can create robust, secure APIs that are easy to maintain and scalable to meet the demands of any application.

3.5. Amazon S3

Amazon Simple Storage Service (AWS S3) is a scalable, high-speed, web-based cloud storage service designed for online data and file storage. Offered by Amazon Web Services (AWS), S3 provides secure, durable, and highly available object storage, capable of storing and retrieving any amount of data from anywhere on the web. It is designed for 99.9999999999 (11 9's) durability and offers comprehensive security and compliance capabilities. AWS S3 is used for a wide range of applications, including backup and restore, data archiving, big data analytics, content distribution, and hosting static websites. Users pay only for the storage they use, making it a cost-effective solution for businesses of all sizes.

4. AWS Account Configuration

4.1. AWS workload account

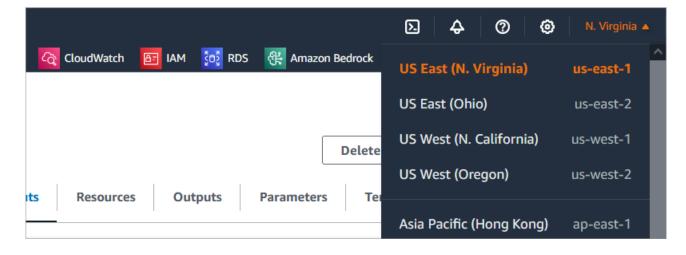
- 1. Visit the AWS Website Go to the AWS website and click on the "Create an AWS Account" button.
- **2. Sign Up** Provide your email address, choose a password, and enter an AWS account name. Click on "Continue."
- **3. Contact Information** Enter your contact information, including your name, address, and phone number. You will also need to choose an account type (Personal or Professional).
- **4. Payment Information** Enter your credit or debit card information. AWS requires this to verify your identity and for billing purposes. Your card will not be charged unless you use paid services.
- **5. Identity Verification** AWS will verify your identity via a phone call or text message. Enter the verification code you receive.
- **6. Select a Support Plan** Choose a support plan that fits your needs. The Basic Support plan is free and sufficient for getting started.
- 7. Complete the Setup Review the details and complete the setup process.

After completing these steps, you will receive a confirmation email. You can then log in to the AWS Management Console and start using AWS services.

5. AWS workload configuration

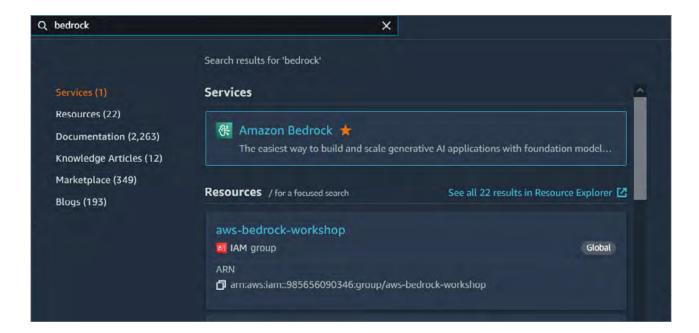
5.1. Foundation model

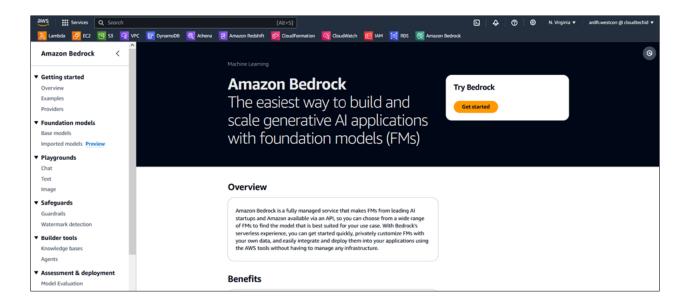
Enabling an AWS Bedrock foundation model involves leveraging AWS's managed service to access and deploy powerful generative AI models. To get started, you must sign in to your AWS Management Console and navigate to the AWS Bedrock service. From there, you can select and activate a foundation model from providers like AI21 Labs, Anthropic, or Stability AI, tailored to your specific use case. AWS Bedrock offers APIs to seamlessly integrate these models into your applications, allowing for tasks such as natural language processing, image generation, and more.



5.2. Bedrock foundation model availability

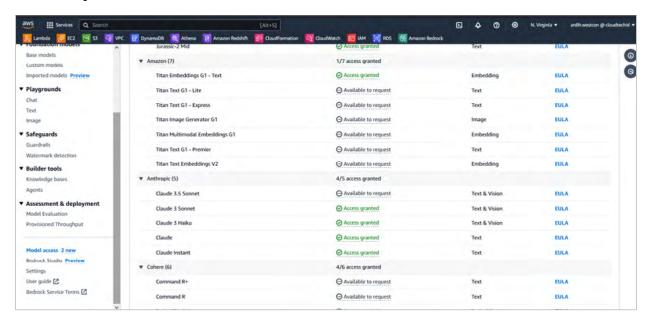
- 1. Use your AWS credentials to log in to the console. Make sure to check your region. The region should be N. Virginia
- 2. In the AWS Management Console, type "Bedrock" in the search bar and select it from the services list. Check Amazon Bedrock Foundation Model





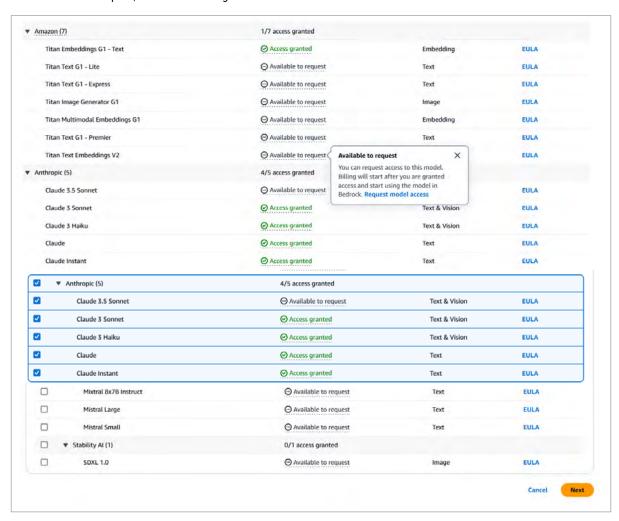
5.2.1. Activate Foundation model

In the AWS Bedrock dashboard, you will find options to view and manage available foundation models. These models come from providers like AI21 Labs, Anthropic, and Stability Al.





3. In the AWS Bedrock dashboard, you will find options to view and manage available foundation models. These models come from providers like AI21 Labs, Anthropic, and Stability Al.





5.3. CloudFormation

Deploying an AWS CloudFormation stack involves several steps to define and launch AWS resources using a CloudFormation template.

5.3.1. CloudFormation parameters

1) Stack name

CloudFormation stack name.

2) VpcCidrBlock

VPC CIDR block, default value is 10.1.0.0/16. Equivalent to 65,536 IP.

3) PublicSubnetCidrBlock

Public subnet that interacts with internet gateway. Provided 3 subnets for public subnet. Default value:

- a. 10.1.10.0/24 → PublicSubnetACidrBlock
- **b. 10.1.11.0/24** → PublicSubnetBCidrBlock
- c. 10.1.12.0/24 → PublicSubnetCCidrBlock

4) PrivateSubnetCidrBlock

Private subnet that interacts without internet gateway. Provided 3 subnets for private subnet. Default value:

- a. 10.1.0.0/24 → PrivateSubnetACidrBlock
- **b. 10.1.1.0/24** → PrivateSubnetBCidrBlock
- c. 10.1.2.0/24 → PrivateSubnetCCidrBlock

5) VPCName

VPC name for your project.

6) ProjectName

Project name that affect your service name. For example, your s3 bucket name format would be **<ProjectName>-lambda-layer**

7) AppCode

Application code to get token from your Lambda function

8) AppSecret

Application secret to interact with Lambda application JWT. Information about JWT can be found in this link Introduction to **JSON Web Token**.

9) DBMasterUsername

User name for your vector database in Aurora PostgreSQL

10) DBmasterUserPassword

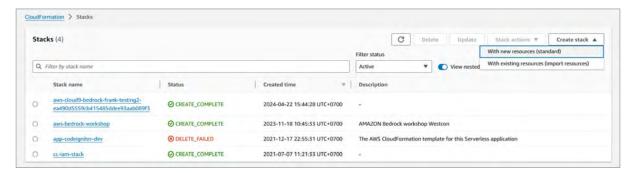
Password for your vector database in Aurora PostgreSQL



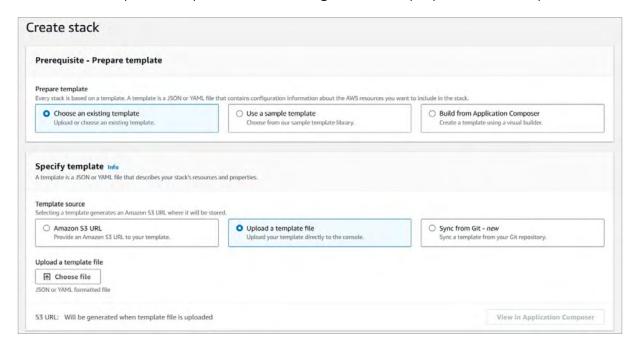


5.3.2. CloudFormation execution

1) Log in to the AWS Management Console, then navigate to the CloudFormation service.



Click on "Create stack" and choose either "With new resources 2) (standard)" or "With existing resources (import resources)."



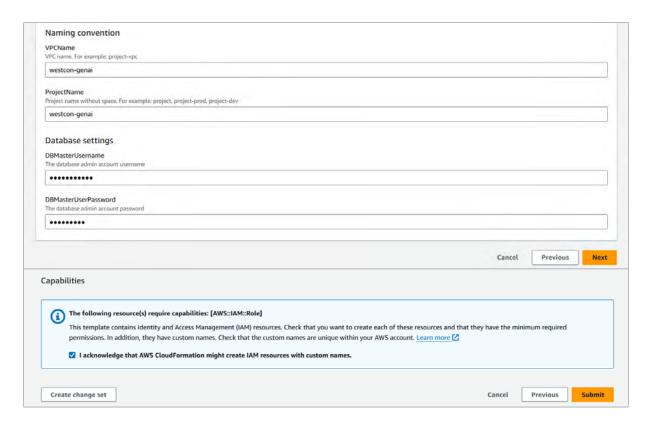
Choose the source of your template 3)



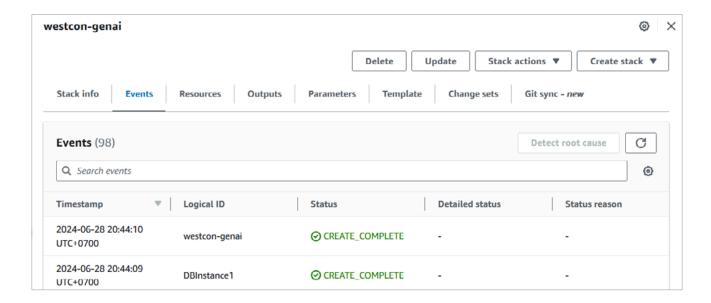
- Configure stack details
 - a) Stack Name: Provide a unique name for your stack.
 - b) Parameters: If your template includes parameters, enter the required values.







- 5) Review all the configurations and template details. Acknowledge that AWS CloudFormation might create IAM resources if your template includes IAM roles or policies. Then, click "Create stack"
- 6) CloudFormation will start creating the resources defined in your template. You can monitor the progress in the "Events" tab of your stack in the CloudFormation console
- 7) Once the stack creation is complete, the status will change to "CREATE_COMPLETE." You can now see the created resources in the AWS Management Console.

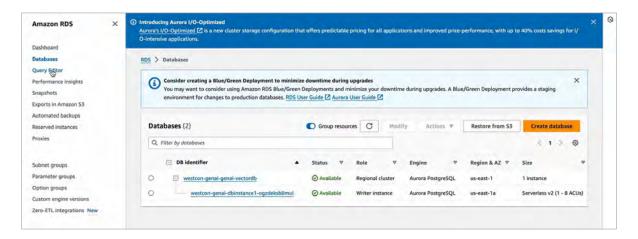




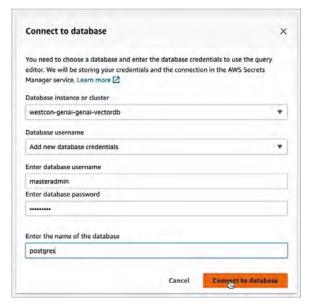
5.4. Database configuration

Query editor setup 5.4.1.

1. In the AWS Management Console, type "RDS" in the search bar and select it from the services list.

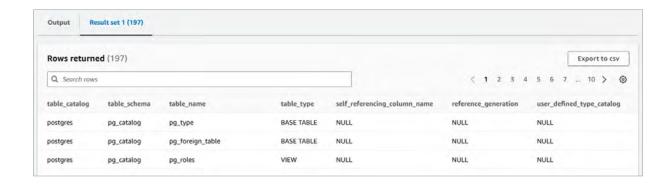


2. Setup database credentials. Ensure you have the database username and password for your RDS instance, as you will need these credentials to connect to the Query Editor.









5.4.2. SQL script extension on PostgreSQL

1. Here are the query command to create vector extension.

```
CREATE EXTENSION IF NOT EXISTS vector;

SELECT extversion FROM pg_extension WHERE extname='vector';

CREATE SCHEMA bedrock_integration;

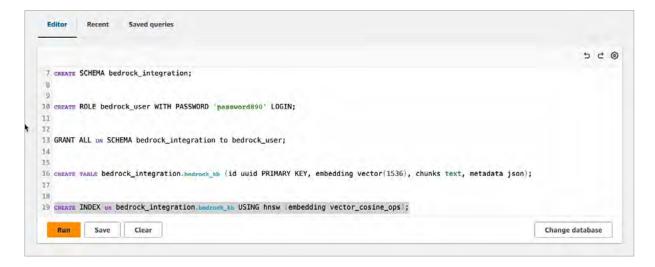
CREATE ROLE bedrock_user WITH PASSWORD <yourPassword> LOGIN;

GRANT ALL ON SCHEMA bedrock_integration to bedrock_user;

CREATE TABLE bedrock_integration.bedrock_kb (id uuid PRIMARY KEY, embedding vector(1536), chunks text, metadata json);

CREATE INDEX on bedrock_integration.bedrock_kb USING hnsw (embedding vector_cosine_ops);
```

2. Block each line of the code and hit Run.

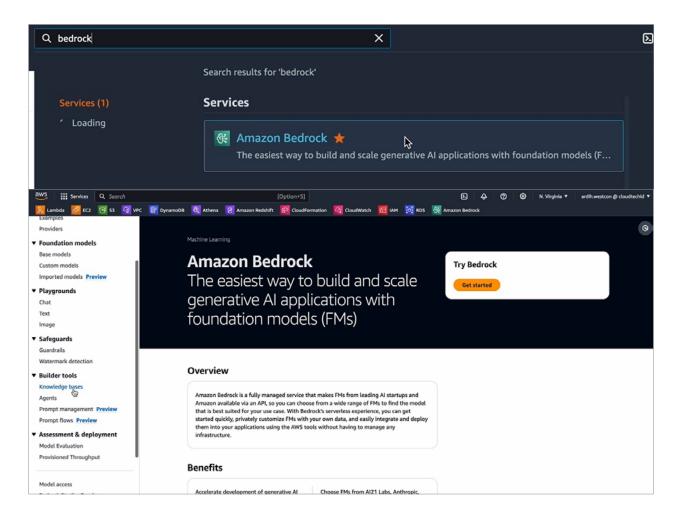


5.5. Knowledge base configuration

In AWS Bedrock, a knowledge base refers to a centralized repository of information that can be utilized by generative AI models to enhance their performance and accuracy. This repository stores structured and unstructured data, such as documents, FAQs, manuals, and other relevant content, which the AI models can access and learn from. By leveraging a knowledge base, AWS Bedrock can provide more accurate and contextually relevant responses, improve decision-making processes, and streamline information retrieval. This capability is particularly useful for applications requiring detailed and specific knowledge, such as customer support, technical documentation, and content creation, enabling more intelligent and informed Al-driven interactions.

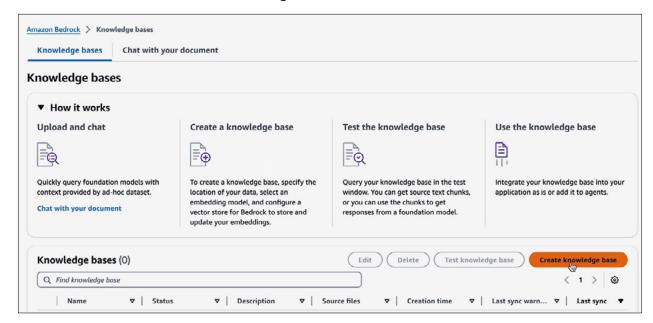
5.5.1. Create knowledge base

1. In the AWS Management Console, search for "Bedrock" and select it from the services list.

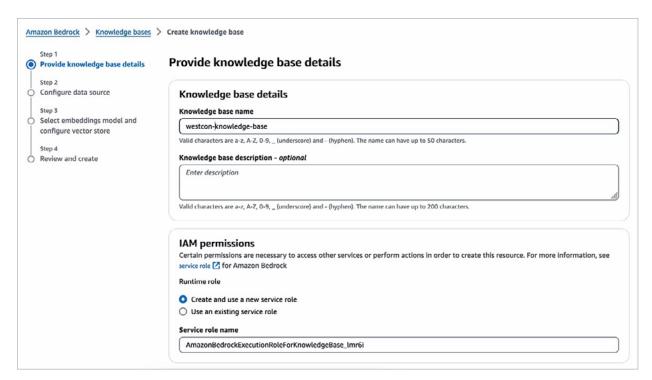




2. Create a New Knowledge Base



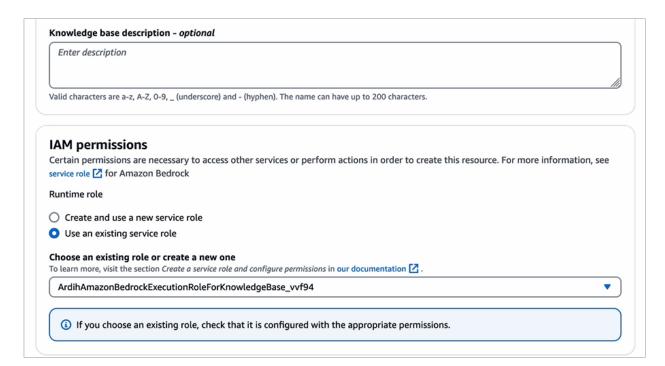
3. Provide a name and description for your Knowledge Base.



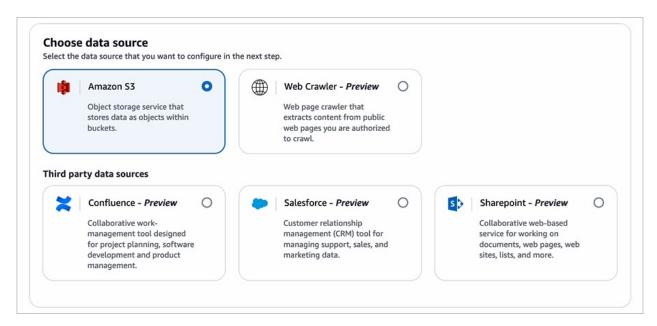




4. Create IAM permission or use existing role.

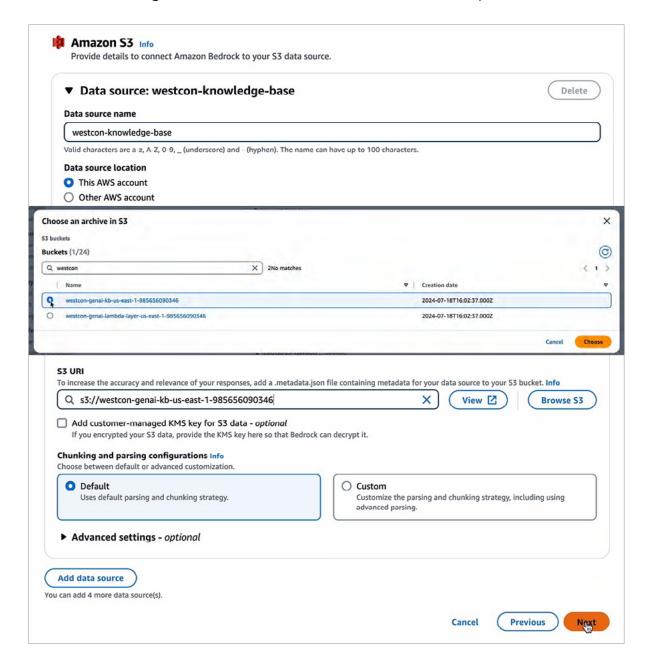


5. Choose Amazon S3 for data source.

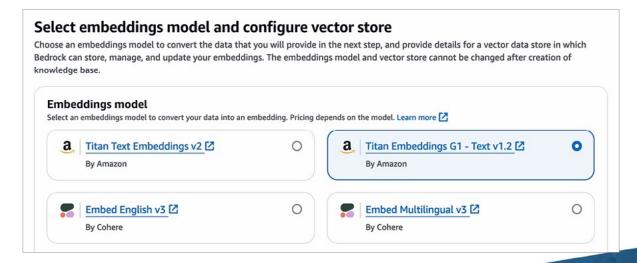




6. Navigate to S3 bucket from CloudFormation template.

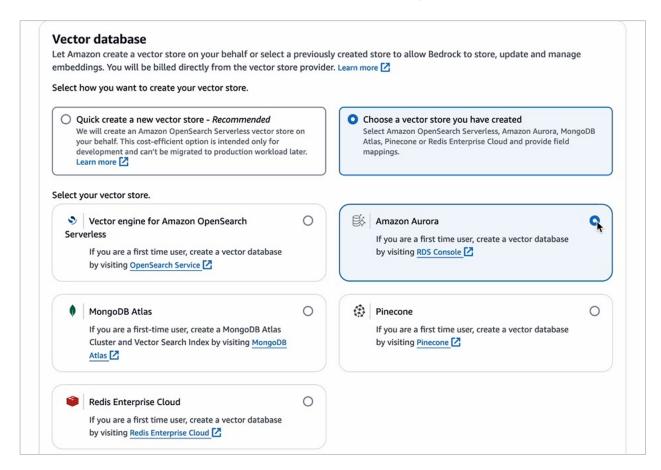


7. Use Titan Embeddings V1.2 for Embeddings model

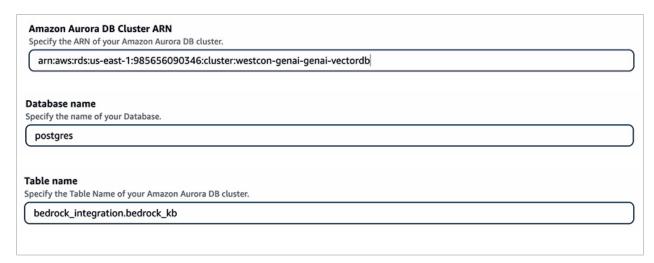




8. Choose vector store that we have created, and select for Amazon Aurora.



9. Fill Aurora DB ARN. Database name would be postgres and table name be seen from query editor in this link







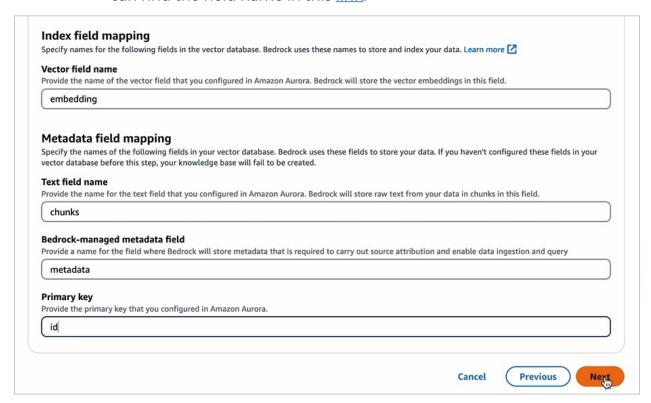
10. To find Secret ARN, you can go to AWS Secret Manager page.



11. Copy and paste the ARN to secret ARN field.



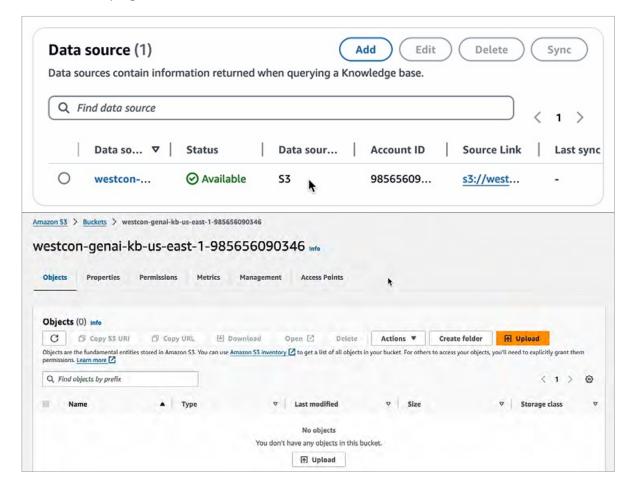
12. Index field mapping would like this according to query editor step. You can find the field name in this link.



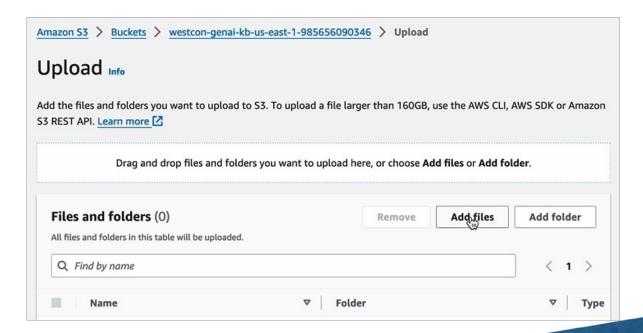
5.5.2. Setup knowledge base integration

After knowledge base creation, navigate to data source and you should find data source that we have setup.

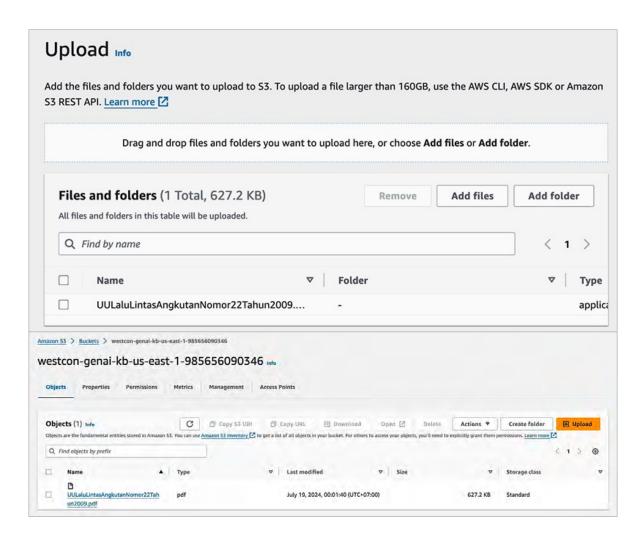
1. Select data source, and hit the source link. This will navigate you to S3 page.



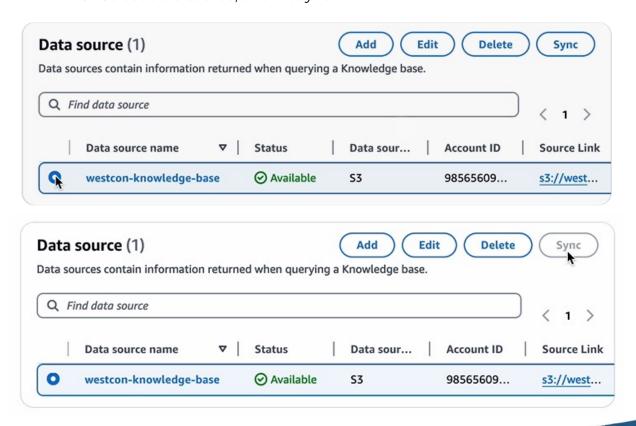
2. Try to upload one or more documents to test.



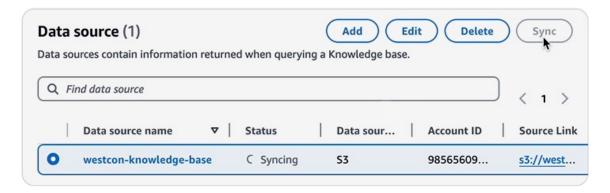




3. Select data source, and hit sync.



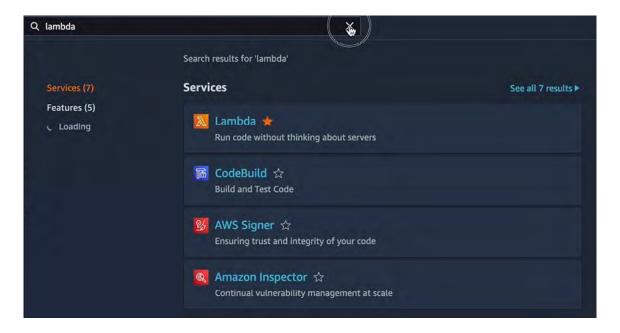
4) Wait until syncing process is completed.



5.6. Lambda configuration

Before we setup AWS Lambda Function, we need to add 2 Lambda Layer to work, Boto3 and PyJWT. AWS Lambda layers are used to simplify and streamline the management of dependencies and shared code across multiple Lambda functions. By creating a layer, you can package libraries, custom runtime environments, and other dependencies separately from your function code, which allows for cleaner and more modular code management. This separation enhances reusability and efficiency, as the same layer can be shared across different functions without duplicating code. Layers also make it easier to update dependencies independently of the function code, reducing deployment time and minimizing the risk of errors. This approach not only optimizes the development process but also ensures better version control and consistency across serverless applications.

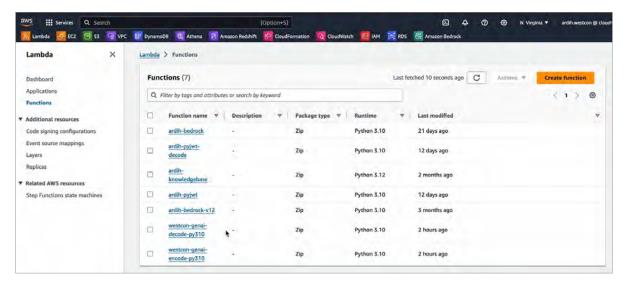
1. Navigate to Lambda page in the AWS console.



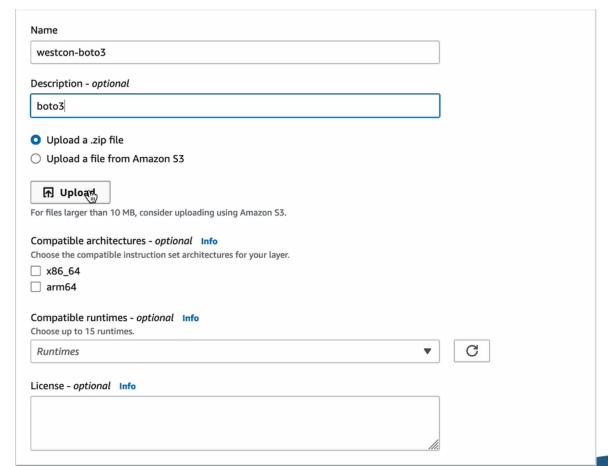
5.6.1. Lambda Layer Boto3

Boto3 is the Amazon Web Services (AWS) Software Development Kit (SDK) for Python, enabling developers to interact with AWS services programmatically. It simplifies the process of integrating Python applications with AWS by providing an easy-to-use, consistent interface for managing and utilizing AWS resources like S3, EC2, DynamoDB, Lambda, and more.

1. Expand the left panel and search for Layers menu. Click Create layer.



2. Upload the zip file you created or specify an S3 URL if you've uploaded the zip to an S3 bucket. Fill the field with your preferred name that indicate library name. For this instance, **Westcon-boto3** is the name of layer.

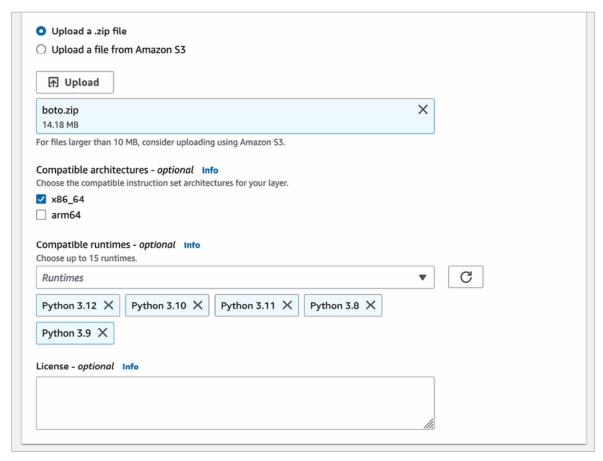




3. Choose boto.zip



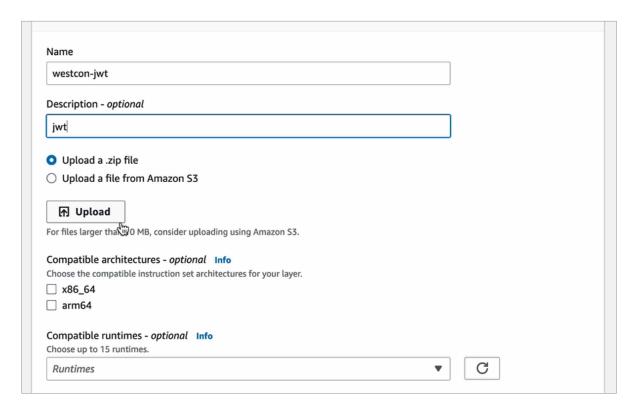
4. Choose x86_64 for compatible architectures, and select all compatible runtime from Python 3.8 to Python 3.12



5.6.2. Lambda Layer JWT

PyJWT is a Python library used for encoding and decoding JSON Web Tokens (JWT). JWTs are a compact, URL-safe means of representing claims to be transferred between two parties, commonly used for authentication and information exchange.

1. Upload the zip file you created or specify an S3 URL if you've uploaded the zip to an S3 bucket. Fill the field with your preferred name that indicate library name. For this instance, Westcon-jwt is the name of layer.



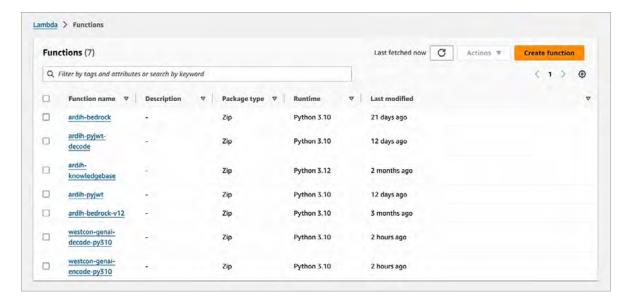
2. Choose jwt.zip



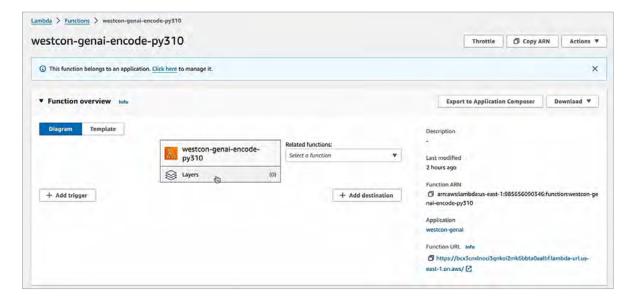
3. Choose <u>x86_64</u> for compatible architectures, and select all compatible runtime from Python 3.8 to Python 3.12

5.6.3. Lambda token configuration

Log in to the AWS Management Console and navigate to the AWS Lambda service. In the Lambda console, select "Functions" from the left navigation pane.



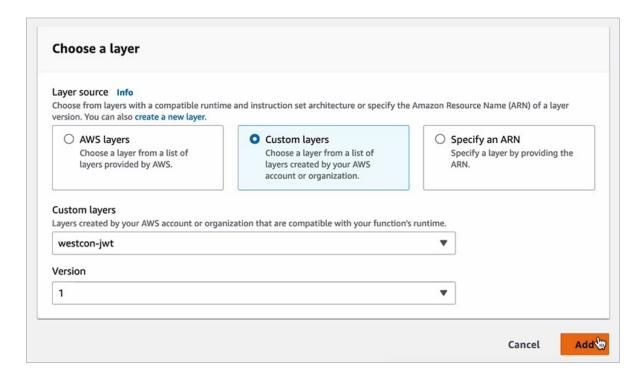
2. Select layers.



Add new layer

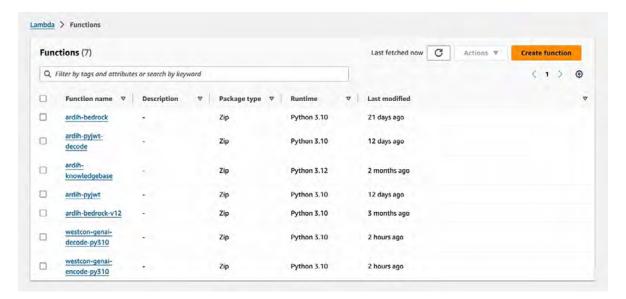


4. Choose custom layer. Select custom layer named Westcon-jwt and version 1.



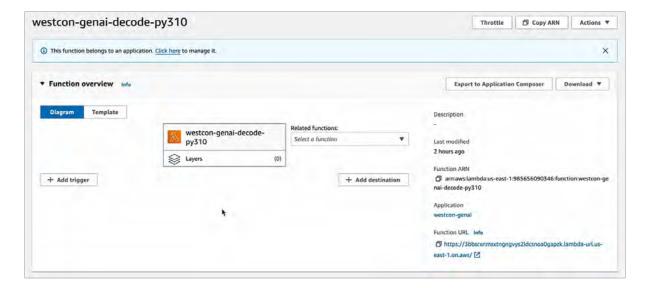
Lambda prompt configuration

1. Log in to the AWS Management Console and navigate to the AWS Lambda service. In the Lambda console, select "Functions" from the left navigation pane.





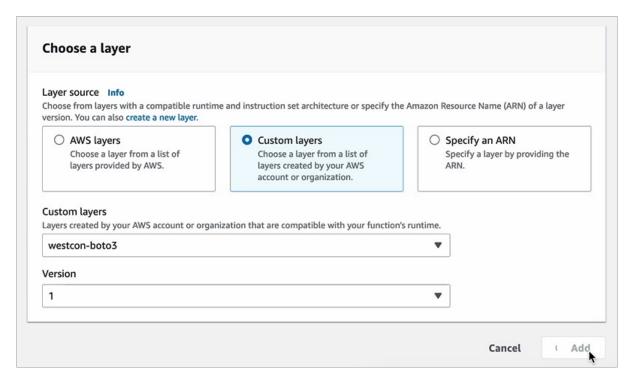
2. Select layers.



3. Add new layer

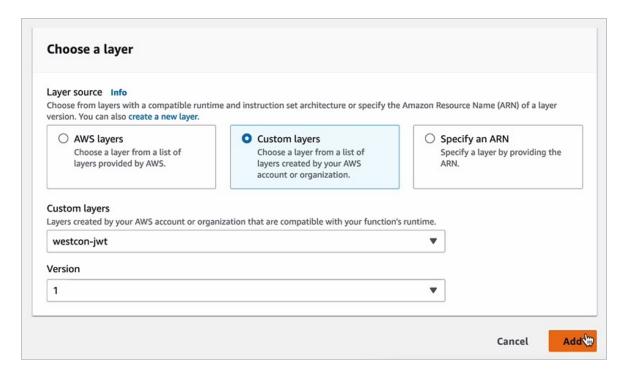


4. Add custom layer named **Westcon-boto3** version 1.

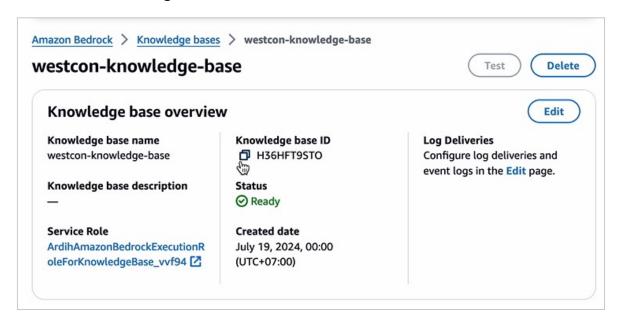




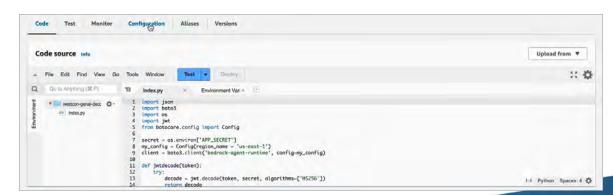
5. Repeat the process for Westcon-jwt.



6. Go to Amazon Bedrock page, select the knowledge base and copy the knowledge base ID.



7. Back to Lambda page, click on **Configuration** tab.

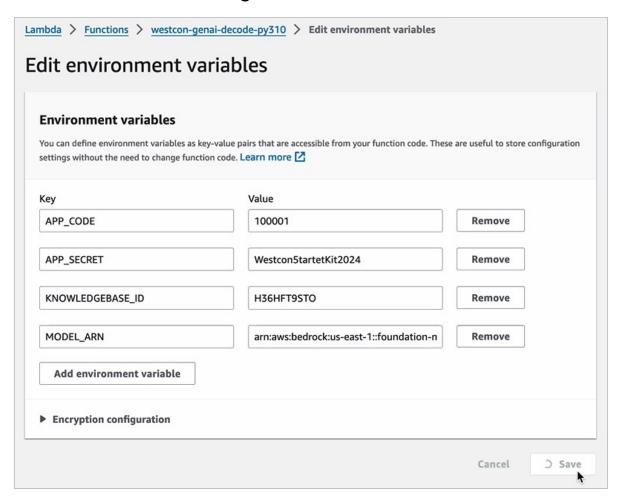




8. Scroll down to **Environment variables**, and edit.



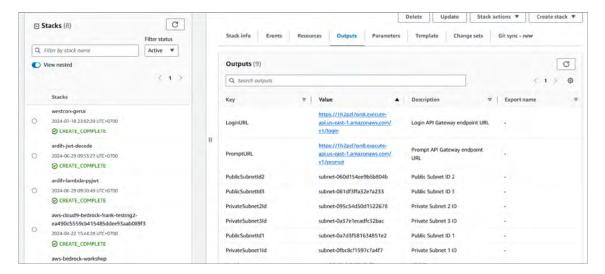
9. Paste the **Knowledge base ID** and hit **Save**.



6. GenAl test API

6.1. Get token

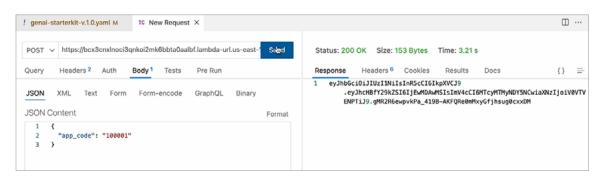
1. Open CloudFormation page and open **Outputs** tab.



- 2. Copy LoginURL value
- 3. Open Lambda page and copy APP_CODE



4. Paste URL and APP_CODE to your API tester tool such as Postman, Insomnia, etc.



5. You can follow this JSON structure to POST data

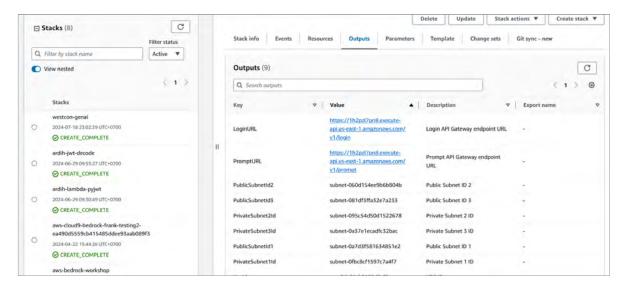
JSON POST structure

```
{
    "app_code": "<YourAppCode>"
}
```

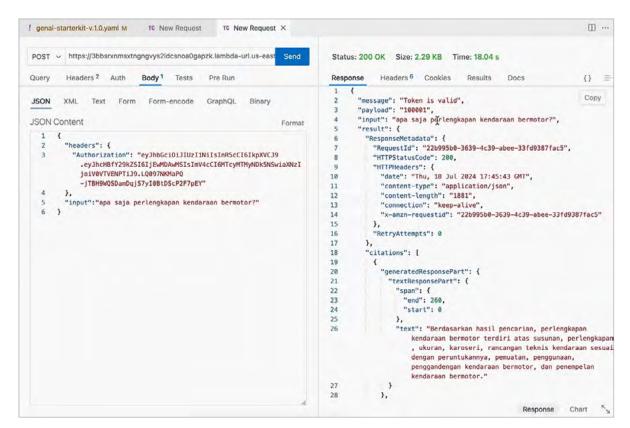
JSON JWT response

6.2. Prompt test

1. Open CloudFormation page and open **Outputs** tab.



- 2. Copy PromptURL value
- 3. Copy PromptURL and token from previous operation



4. Follow this JSON structure to POST data

```
JSON POST structure
{
    "headers": {
        "Authorization": "<token>"
     },
     "input": "<your-question>"
}
```

```
JSON response structure
 "message": "<token-message>",
 "payload": "<YourAppCode>",
 "input": "<question>",
 "result": {
  "ResponseMetadata": {
   "RequestId": "08b8e586-9940-4e4c-88b9-b5cce47fc797",
   "HTTPStatusCode": 200,
   "HTTPHeaders": {
    "date": "Fri, 19 Jul 2024 08:17:19 GMT",
    "content-type": "application/json",
    "content-length": "3689",
    "connection": "keep-alive",
    "x-amzn-requestid": "08b8e586-9940-4e4c-88b9-b5cce47fc797"
   "RetryAttempts": 0
  },
  "citations": [
    "generatedResponsePart": {
     "textResponsePart": {
      "span": {
       "end": 381,
       "start": 0
      },
      "text": "<text>"
     }
    },
    "retrievedReferences": [
      "content": {
       "text": "<text>"
      },
      "location": {
       "s3Location": {
        "uri": "s3://<s3-object-path>"
       },
        "type": "S3"
     },
      "content": {
       "text": "<text>"
      },
      "location": {
        "s3Location": {
        "uri": "s3://<s3-object-path>"
        "type": "S3"
  ],
  "output": {
   "text": "<answer>"
  "sessionId": "e7bba22d-0749-4402-9b10-3b96cee90655"
 }
}
```



Have a question?

Contact us

NZ Cloud Sales: +64 9 477 7211 cloudsales.nz@westcon.com

AU Cloud Sales: +61 2 8412 1212 cloudsales.au@westcon.com

SG Cloud Sales: +65 6424 0570 cloudsales.sg@westcon.com

ID Cloud Sales: +62 21 8062 1470 cloudsales.id@westcon.com

VN Cloud Sales: +84 2836203338 cloudsales.vn@westcon.com