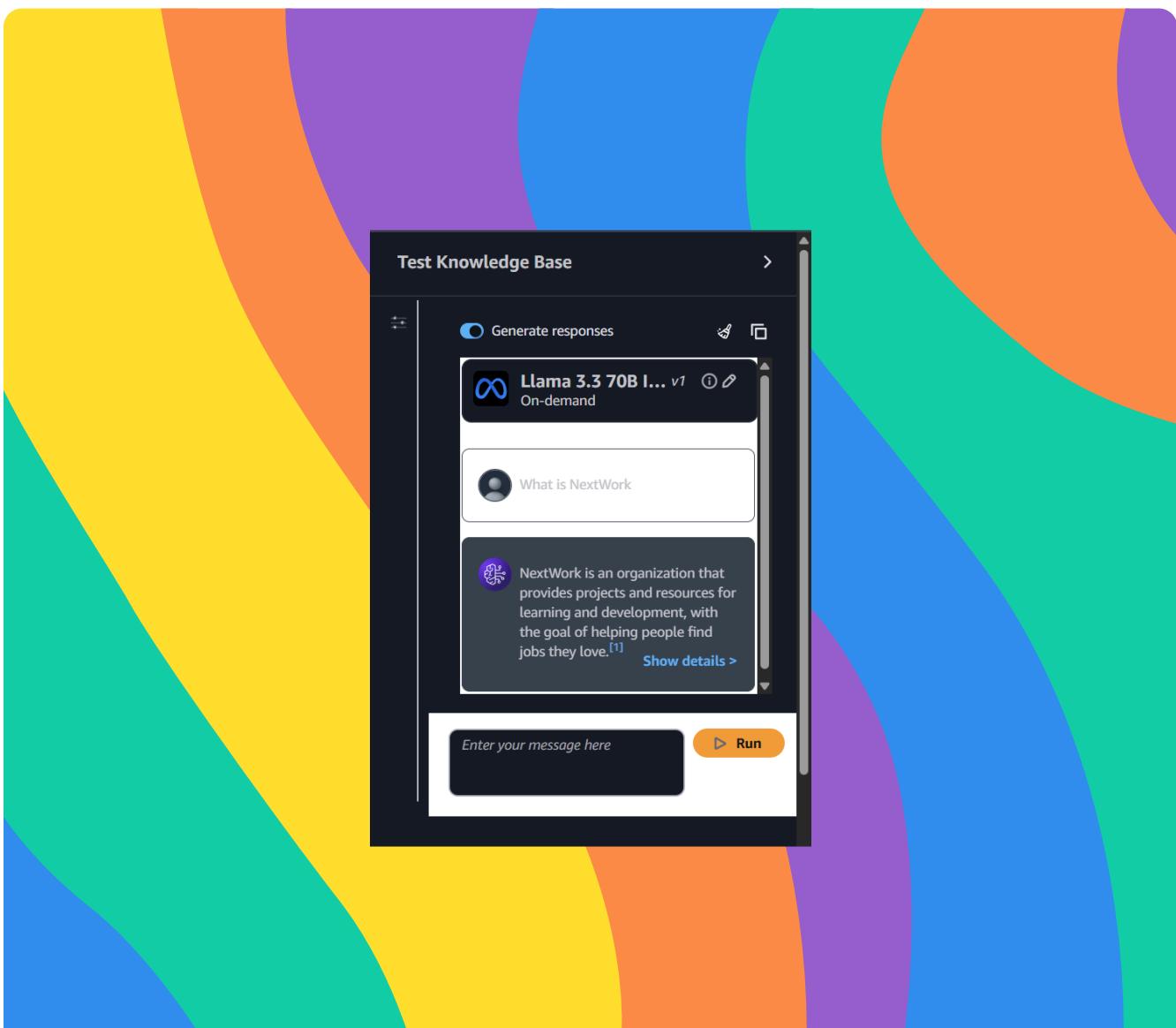


nextwork.org

Set Up a RAG Chatbot in Bedrock



jyothesh karnam





Introducing Today's Project!

RAG (Retrieval Augmented Generation) is an AI technique that lets you train an AI model on your own personal documents. In this project, I will demonstrate RAG by setting up a RAG chatbot in Amazon Bedrock.

Tools and concepts

Services I used in this project were Amazon Bedrock, S3, and OpenSearch Serverless. Key concepts I learnt include knowledge bases, requesting access to AI models, how chatbots generate responses (i.e. AI models + knowledge base), vector stores.

Project reflection

This project took me approximately two hours including project demo time. The most challenging part was the error with my AI models (and understanding on-demand vs pre-provisioned inference). It was most rewarding to level up my chatbot's responses!

I did this project today to learn more about Bedrock and RAG. This project definitely met my goals – awesome to learn both over a hands-on project!



Understanding Amazon Bedrock

Amazon Bedrock is an AWS service that makes it easy to build generative AI applications. This is because Bedrock is like an AI model marketplace that lets us find, use and test models from different providers in a central place. I am using Bedrock in this project to create a knowledge base.

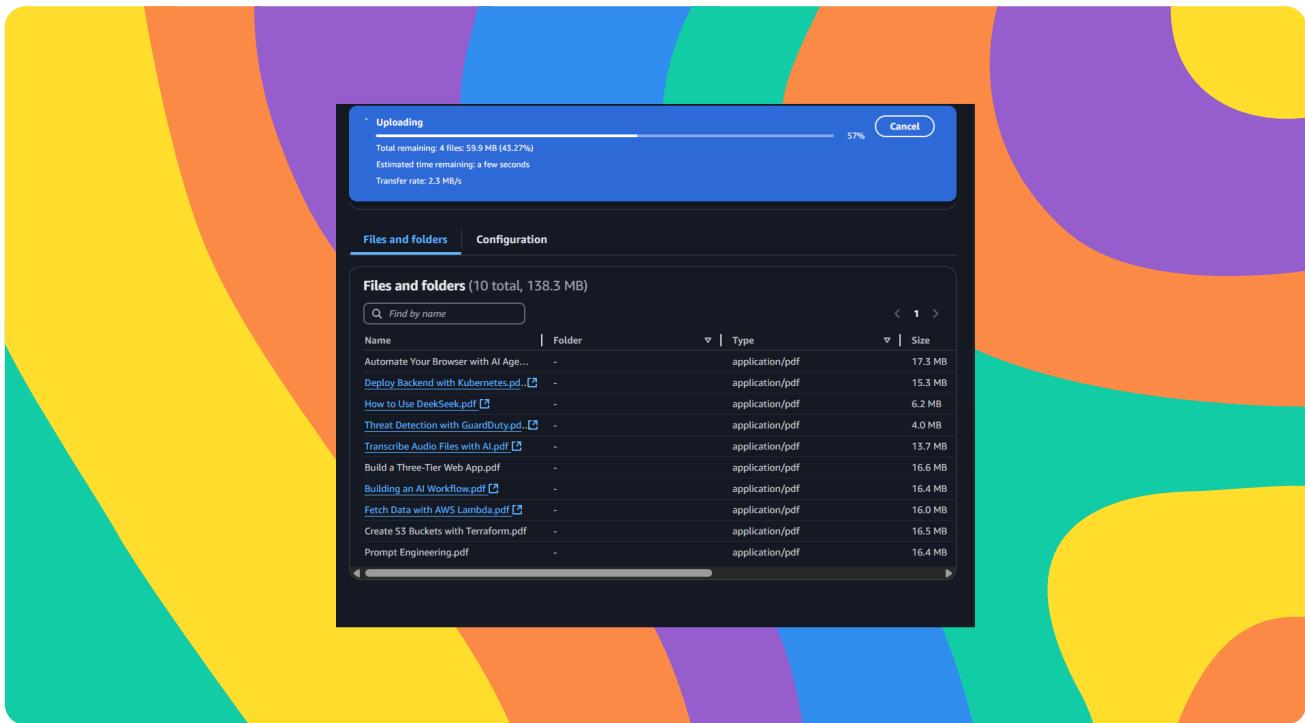
My Knowledge Base is connected to S3 because S3 is going to be the storage/source for my Knowledge Base's raw documents. S3 is AWS's storage service, where you can store all kinds of objects (e.g. videos, documents, audio) in the same bucket.

In an S3 bucket, I uploaded the documents that will make up my chatbot's knowledge. My S3 bucket is in the same region as my Knowledge Base because Bedrock is a regional service – data must live in the same region as the Bedrock resource (Kbase).



jyothesh karnam
NextWork Student

NextWork.org





My Knowledge Base Setup

My Knowledge Base uses a vector store, which means a search engine/database that stores data based on their semantic meaning! When I query my Knowledge Base, OpenSearch will find the relevant chunks of data to the query, and pass it to Bedrock.

Embeddings are vector representations of the semantic meaning of a chunk of text. The embedding model I'm using is Titan Text Embeddings v2 because it's fast, accurate and a lot more affordable!

Chunking is the process of splitting large documents into smaller, manageable pieces called chunks, which makes it easier for the model to retrieve and understand relevant information during a query. In my Knowledge Base, chunks are set to a specific size and overlap value to ensure context is preserved across splits, improving the accuracy of responses.



jyothesh karnam
NextWork Student

[NextWork.org](https://www.nextwork.org)

Preparing service decisions in Amazon OpenSearch Service. This process may take several minutes to complete.

Amazon Bedrock | [Troubleshooting](#) > Create knowledge base with vector store

Step 1: Provide knowledge base details

Knowledge base name: Knowledge base description: Service role: Knowledge base type: Data source type: Log tail index:

Step 2: Configure data source

Data source: s3://bucket-next-network-rag-bedrock

| | | |
|--|---|--|
| Data source name: s3://bucket-next-network-rag-bedrock | Account ID: 432345678901 (This account) | S3 URI: s3://next-network-rag-bedrock/ |
| Customer-managed AWS Key for S3: | AWS key for transient data storage: | Checking strategy (Default) |
| Pending at stage: DEFAULT | Lambda function: | S3 bucket for Lambda function: |
| Data deletion policy: DELETE | | |

Step 3: Configure data storage and processing

Embeddings model: Embedding type: Vector dimension:

Vector store: Quick connector type: DocumentDB connection:

S3 URI:

[Cancel](#) [Purchase](#) [Create Knowledge Base](#)



jyothesh karnam
NextWork Student

NextWork.org

AI Models

AI models are important for my chatbot because they're the translator of my Knowledge Base's search results into human-like text. Without AI models, my chatbot would only respond with chunks of text from my documents – which isn't the best experience!

To access an AI model in Bedrock, I had to visit the "Model Access" page and request access explicitly! AWS needs explicit access because some AI model providers have extra forms/rules if you wanted to use them, and AWS needs to check availability.

The screenshot shows the "Base models (18)" section of the AWS Bedrock Model Access interface. The interface has a dark theme with colored sections (yellow, orange, purple, blue, green) on the sides. The main content area displays a table of models categorized by provider:

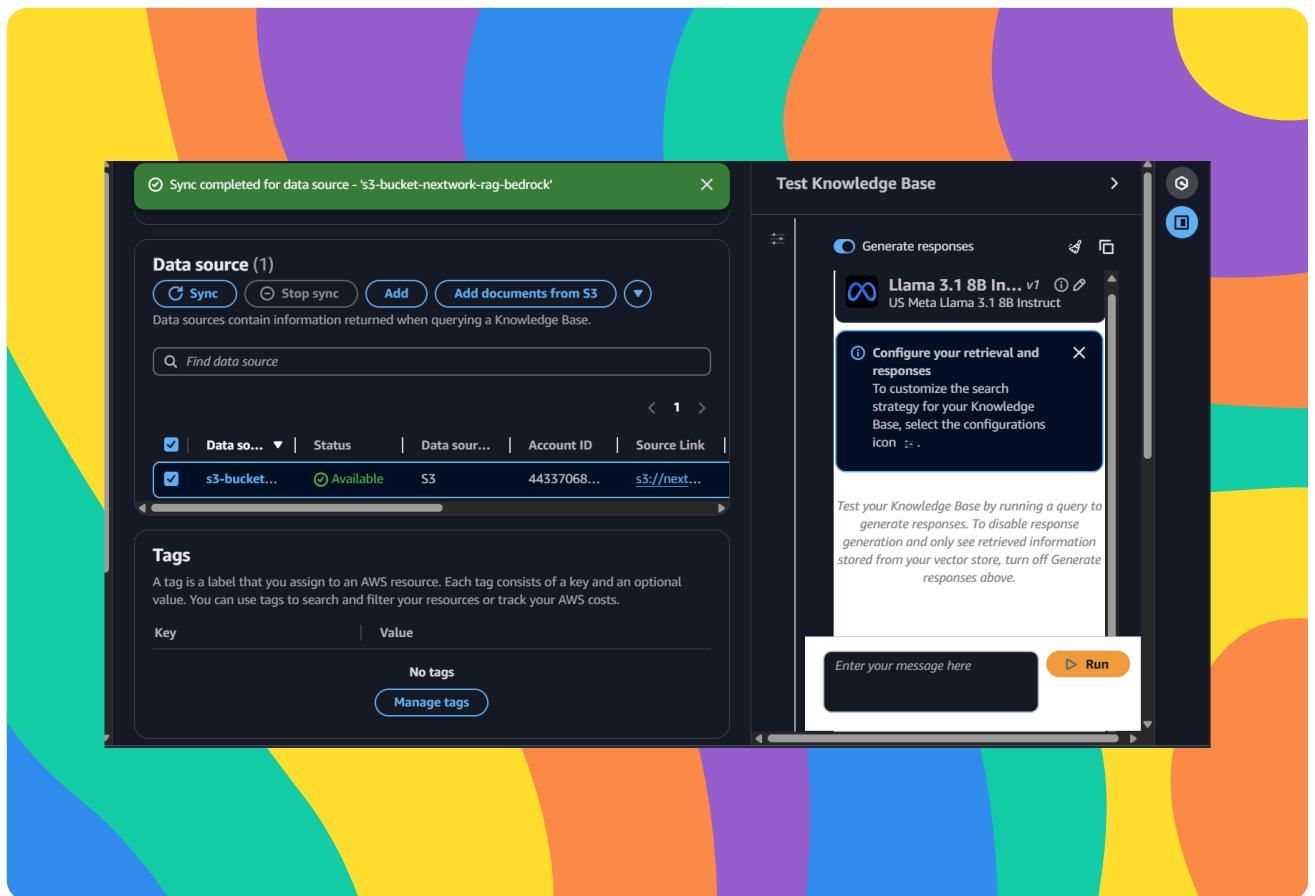
| Models | Access status | Modality | EULA |
|-------------------------------|----------------------|---------------|------|
| Amazon (4) | 1/4 access granted | | |
| Titan Text Embeddings V2 | Access granted | Embedding | EULA |
| Nova Pro | Available to request | Text & Vision | EULA |
| Nova Lite | Available to request | Text & Vision | EULA |
| Nova Micro | Available to request | Text | EULA |
| Anthropic (5) | 0/5 access granted | | |
| Claude 3.7 Sonnet | Available to request | Text & Vision | EULA |
| Claude 3.5 Haiku | Available to request | Text | EULA |
| Claude 3.5 Sonnet v2 | Available to request | Text & Vision | EULA |
| Claude 3.5 Sonnet | Available to request | Text & Vision | EULA |
| Claude 3 Haiku | Available to request | Text & Vision | EULA |
| DeepSeek (1) | 0/1 access granted | | |
| DeepSeek-R1 | Available to request | Text | EULA |
| Meta (8) | 2/8 access granted | | |
| Llama 3.5 70B Instruct | Access granted | Text | EULA |
| Llama 3.2 1B Instruct | Available to request | Text | EULA |
| Llama 3.2 3B Instruct | Available to request | Text | EULA |
| Llama 3.2 11B Vision Instruct | Available to request | Text & Vision | EULA |
| Llama 3.2 90B Vision Instruct | Available to request | Text & Vision | EULA |
| Llama 3.1 405B Instruct | Available to request | Text | EULA |
| Llama 3.1 70B Instruct | Available to request | Text | EULA |
| Llama 3.1 8B Instruct | Access granted | Text | EULA |



Syncing the Knowledge Base

Even though I've already connected my S3 bucket when creating the Knowledge Base, I still need to sync my data, which will actually move the data from S3 and into my Knowledge Base + OpenSearch Serverless. Currently, the Knowledge Base is empty!

The sync process involves three key steps: Ingesting (i.e. Bedrock takes the data from S3), Processing (i.e. Bedrock chunks and embeds the data) and Storing (i.e. Bedrock stores the processed data in my vector store, OpenSearch Serverless).





jyothesh karnam
NextWork Student

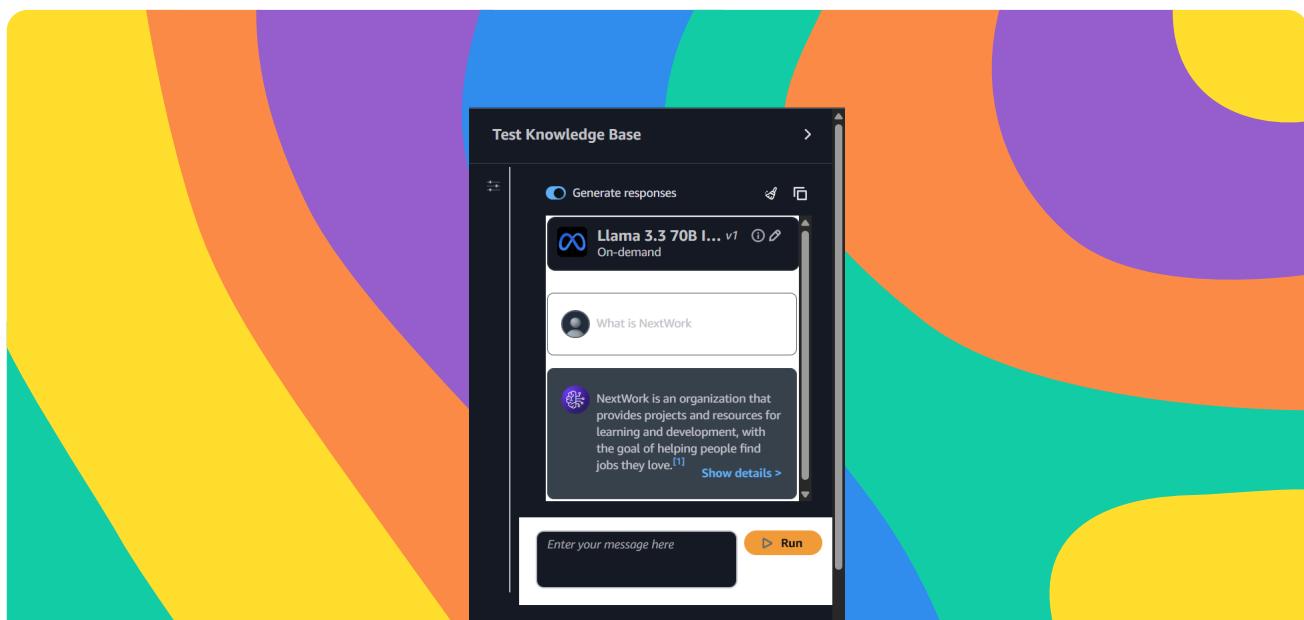
NextWork.org

Testing My Chatbot

I initially tried to test my chatbot using Llama 3.1 8B as the AI model, but it triggered an error – it was not available on demand. I had to switch to Llama 3.3 70B because it was offered on-demand by AWS (since it's a newer, efficient model).

When I asked about topics unrelated to my Knowledge Base's data, my chatbot responds that it cannot help me with this request. This proves that the chatbot only knows the information that I gave it – it won't know even common knowledge outside.

I can also turn off the Generate Responses setting to see the raw chunks of data directly from my Knowledge Base. When I tested this, my chatbot gave a list of 5 paragraphs to answer a question, whereas the AI model will convert it to a sentence.





Upgrading My Chatbot

In a project extension, I looked into increasing the number of source chunks, which are the number of chunks of text that the Knowledge Base will return for a query. This will improve my chatbot's responses by giving my AI model more data to use.

I also added a custom prompt that tells my chatbot to always respond by mentioning skills that the student's learnt, even if the student didn't ask for it directly. I also gave context that the database contains documentation for projects.

After adjusting the source chunks and the generation prompt, my chatbot's response became so much more detailed! The chatbot would talk about skills I've learnt, and used 10 or more examples within the response.

The screenshot shows a chatbot interface with a colorful background. A user query "What projects are in this database?" is shown in a white input field. Below it, a dark gray card displays a response from a bot icon. The response details various projects including Threat Detection with GuardDuty, Building an AI Workflow, a three-tier architecture project, Deploy Backend with Kubernetes, and a project on DeepSeek. It also mentions using n8n, OpenAI, Google Calendar, CloudFormation, S3, CloudShell, DynamoDB, API Gateway, Lambda, and CloudFront, along with projects on setting up a local AI workflow and deploying a backend with Kubernetes. A "Show details >" link is at the bottom right of the card.

What projects are in this database?

The projects in this database include Threat Detection with GuardDuty, Building an AI Workflow, a three-tier architecture project, Deploy Backend with Kubernetes, and a project on DeepSeek, which involves setting up a local Large Language Model.^[1]^[2]^[3]^[4]^[5] Additionally, there are projects that involve using n8n, OpenAI, Google Calendar, CloudFormation, S3, CloudShell, DynamoDB, API Gateway, Lambda, and CloudFront, as well as a project on setting up a local AI workflow and a project on deploying a backend with Kubernetes.^[6]^[7]^[8]^[9]

Show details >



NextWork.org

Everyone should be in a job they love.

Check out nextwork.org for
more projects

