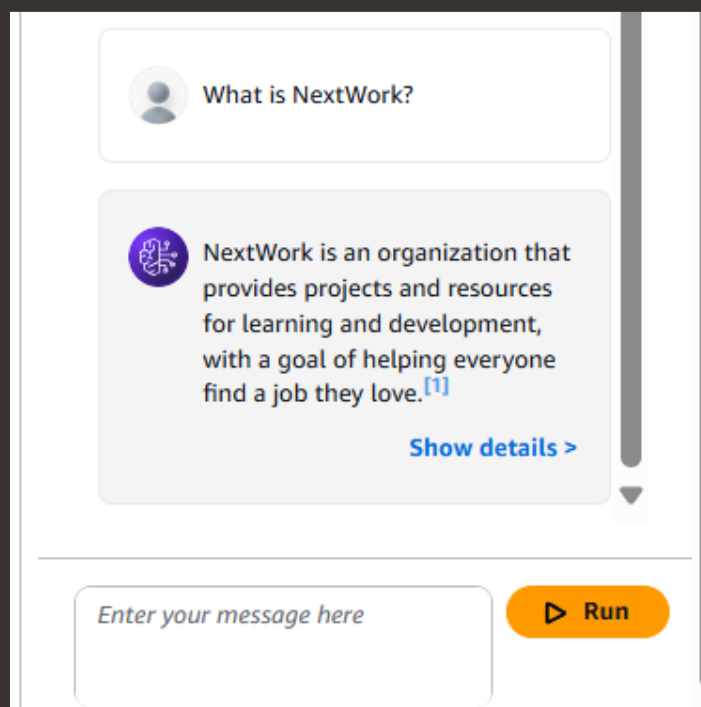# Set Up a RAG Chatbot in Bedrock

DI   Dineshraj Dhanapathy

# Introducing Today's Project!

RAG (Retrieval Augmented Generation) is a method that combines retrieved data with AI to produce informed answers. In this project, I will demonstrate RAG by merging live search results with AI output.

## Tools and concepts

Services I used were Amazon Bedrock, S3, and OpenSearch Serverless. Key concepts I learnt include chunking, embeddings, vector stores, and the importance of syncing data for effective AI-powered

## Project reflection

This project took me approximately 2 hours. The most challenging part was setting up the Knowledge Base and troubleshooting model responses. It was most rewarding to see the chatbot provide accurate, real-time answers.

I did this project today to deepen my understanding of AI models and enhance my skills in integrating knowledge bases. The project met my goals by successfully building a functional, responsive AI chatbot.

# Understanding Amazon Bedrock

Amazon Bedrock is a managed service offering scalable foundation models and tools to build generative AI applications. I'm using Bedrock in this project to integrate and deploy AI-powered solutions.

My Knowledge Base is connected to S3 because S3 is a secure, scalable, and durable storage solution that hosts our data and enables efficient retrieval and processing for AI models for fast access.

In an S3 bucket, I uploaded documents including FAQs, manuals, and guides for our chatbot's Knowledge Base. My S3 bucket is in the same region as my KB because I'm using Ohio East 2 for low latency

| Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|
| ☐ Automate Your Browser with AI Agents.pdf | pdf | February 13, 2025, 02:48:12 (UTC+05:30) | 17.3 MB | Standard |
| ☐ Build a Three-Tier Web App.pdf | pdf | February 13, 2025, 02:48:12 (UTC+05:30) | 16.6 MB | Standard |
| ☐ Building an AI Workflow.pdf | pdf | February 13, 2025, 02:48:43 (UTC+05:30) | 16.4 MB | Standard |
| ☐ Create S3 Buckets with Terraform.pdf | pdf | February 13, 2025, 02:48:12 (UTC+05:30) | 16.5 MB | Standard |
| ☐ Deploy Backend with Kubernetes.pdf | pdf | February 13, 2025, 03:04:20 (UTC+05:30) | 15.3 MB | Standard |

# My Knowledge Base Setup

My Knowledge Base uses a vector store, which means it efficiently stores and retrieves embeddings for fast AI responses. When I query my Knowledge Base, OpenSearch will quickly find relevant document chunks

Embeddings are numerical representations of text that help AI understand context and meaning. The embedding model I'm using is Titan Text Embeddings v2 because it provides accurate, efficient semantic search.

Chunking is the process of breaking large text into smaller parts for efficient AI processing. In my Knowledge Base, chunks are set to 300 tokens, ensuring the chatbot retrieves precise and relevant info.

## Review and create

### Step 1: Provide details

Edit

#### Knowledge Base details

| | | |
|---|---|---|
| **Knowledge Base name** | **Knowledge Base description** | **Service role** |
| nextwork-rag-documentation | This Knowledge base stores all documentation at NextWork. | AmazonBedrockExecutionRoleForKnowledgeBase_st1hm |
| **Knowledge base type** | **Data source type** | **Log Deliveries** |
| Knowledge base use vector store | S3 | — |

### Step 2: Setup up data source

Edit

#### Data source: s3-bucket-nextwork-rag-bedrock

| | | |
|---|---|---|
| **Data source name** | **Account ID** | **S3 URI** |
| s3-bucket-nextwork-rag-bedrock | 466742534146 (this account) | s3://nextwork-rag-bedrock-dd ↗ |
| **Customer-managed KMS Key for S3** | **KMS key for transient data storage** | **Chunking strategy** |

# AI Models

AI models are important for my chatbot because they generate intelligent, context-aware responses. Without AI models, my chatbot would only retrieve static text without understanding user intent.

To get access to AI models in Bedrock, I had to request model access in the AWS console. AWS needs explicit access because it ensures compliance, security, and proper usage of foundation models.
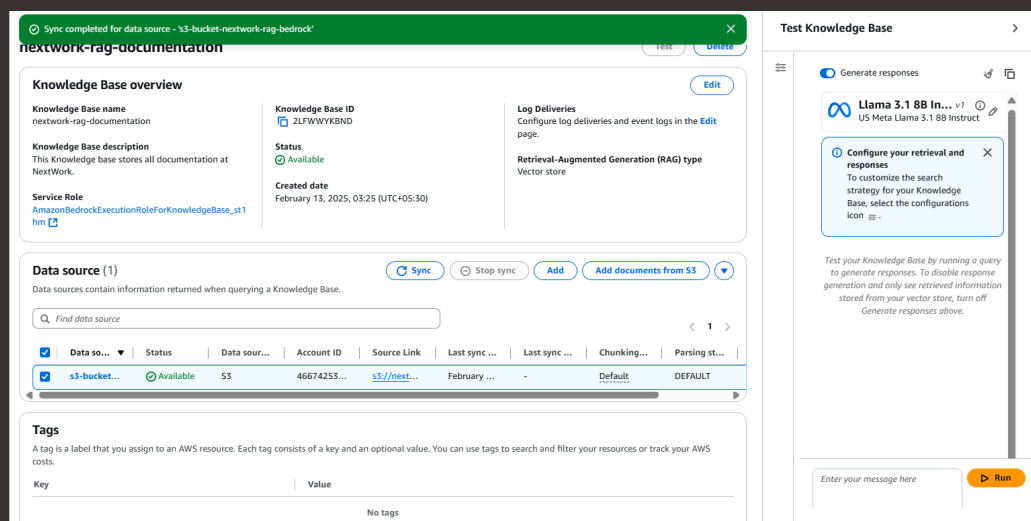
# Syncing the Knowledge Base

Even though I already connected my S3 bucket when creating the Knowledge Base, I still need to sync because it ensures the latest data is indexed, updated, and ready for accurate retrieval by the AI.

The sync process involves three steps: Ingesting, where Bedrock retrieves data from S3; Processing, where it chunks and embeds the data; and Storing, where Bedrock stores the processed data in OpenSearch Serverless.

# Testing My Chatbot

I initially tried to test my chatbot using Llama 3.1 8B as the AI model, but it lacked the depth needed for complex queries. I had to switch to Llama 3.3 70B because it provides more accurate, nuanced responses.

When I asked about topics unrelated to my data, my chatbot struggled to provide relevant answers. This proves that the chatbot relies on the knowledge base, and without related data, its accuracy drops.

You can also turn off the Generate Responses setting to prevent the AI from automatically generating answers. This allows you to manually adjust or review the responses before they are delivered, giving more control over the interaction.

What is NextWork?

NextWork is an organization that provides projects and resources for learning and development, with a goal of helping everyone find a job they love.[1]

Show details >

Enter your message here

▷ Run