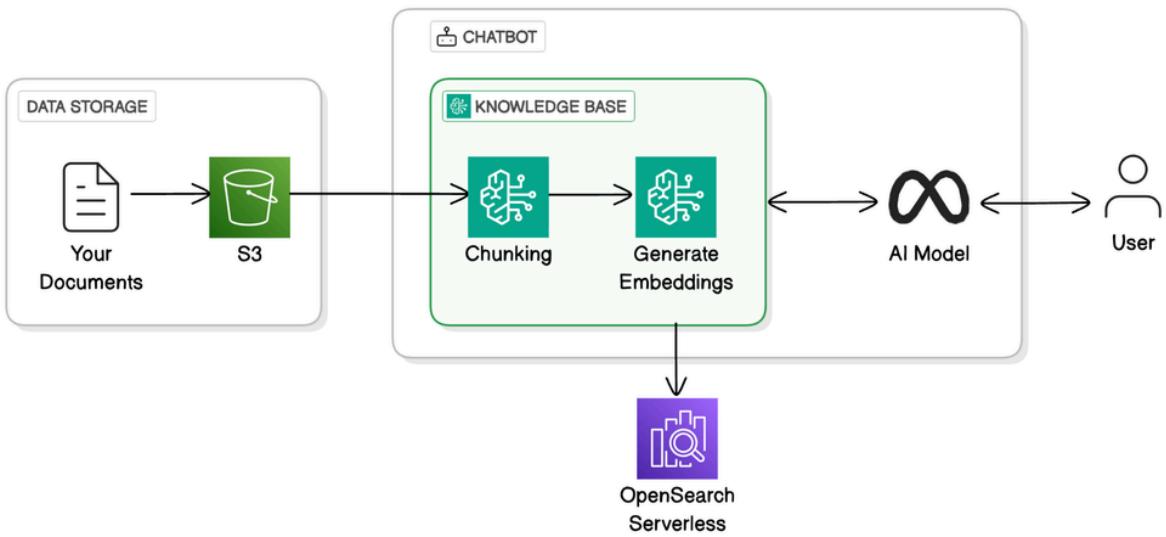




Kareshma
Rajaananthapadmanaban

Set Up a RAG Chatbot in Bedrock

Build an AI chatbot that learns from your data with RAG and Amazon Bedrock





Set Up a RAG Chatbot in Bedrock



Kareshma Rajaananthapadmanaban

The screenshot shows the Amazon Bedrock interface for setting up a Knowledge Base. The top navigation bar includes 'Amazon Bedrock', 'Knowledge Bases', 'nextwork-rag-documentation', and 'Test Knowledge Base'. A green notification bar at the top states 'Sync completed for data source - 's3-bucket-nextwork-rag-bedrock''. The main area is titled 'Test: nextwork-rag-documentation'. It has two tabs: 'Configurations' (selected) and 'Test'. The 'Configurations' tab shows 'Retrieval and response generation' set to 'Retrieval and response generation: data sources and model' (selected), using 'Llama 3.3 70B' as the model. The 'Test' tab shows a preview conversation:

- User: Hello
- Bot: Hello, I'm ready to help. What's your question?
- User: What is NextWork ?
- Bot: NextWork is a website that offers various projects and has a community forum where users can ask questions and share information.[\[1\]](#)[\[2\]](#)[\[3\]](#)[\[4\]](#)[\[5\]](#)

A text input field at the bottom says 'Write a prompt (Shift + ENTER to start a new line, and ENTER to generate a response)'.



Introducing Today's Project!

RAG (Retrieval Augmented Generation) is an AI technique that enhances a model's ability to provide **context-specific** answers by integrating retrieval mechanisms with generation. In this project, I will demonstrate RAG by training an AI chatbot using documents stored in Amazon S3. The chatbot will **retrieve relevant data** from a **Knowledge Base** indexed in **Amazon OpenSearch** and generate human-like responses based on the documents.

Tools and concepts

Services I used were **Amazon Bedrock, S3, and OpenSearch**. Key concepts I learnt include creating and managing a Knowledge Base, setting up vector stores for retrieval, and integrating S3 as a data source. I also explored guardrails, query modification, and best practices like cleaning up unused resources to avoid extra costs.

Project reflection

This project took me approximately 1 hour. The most challenging part was understanding how to **fine-tune the knowledge base** settings for better responses. It was most rewarding to see the chatbot provide more detailed and contextual answers after the improvements.

I did this project today to strengthen my hands-on skills with RAG in Amazon Bedrock and to better understand how to improve **chatbot responses with advanced configurations**. Yes, this project met my goals as I was able to practice key steps like tuning source chunks, custom prompts, and cleaning up resources, which boosted both my learning and confidence.



Understanding Amazon Bedrock

In this step, I will create a Knowledge Base (KBase) in Amazon Bedrock. This will serve as the **chatbot's personal library**, where all the information it needs to answer queries will be stored.

Additionally, I will set up an S3 bucket to hold the documents that will populate the Knowledge Base. The **KBase** will help the chatbot retrieve relevant information and provide accurate responses.

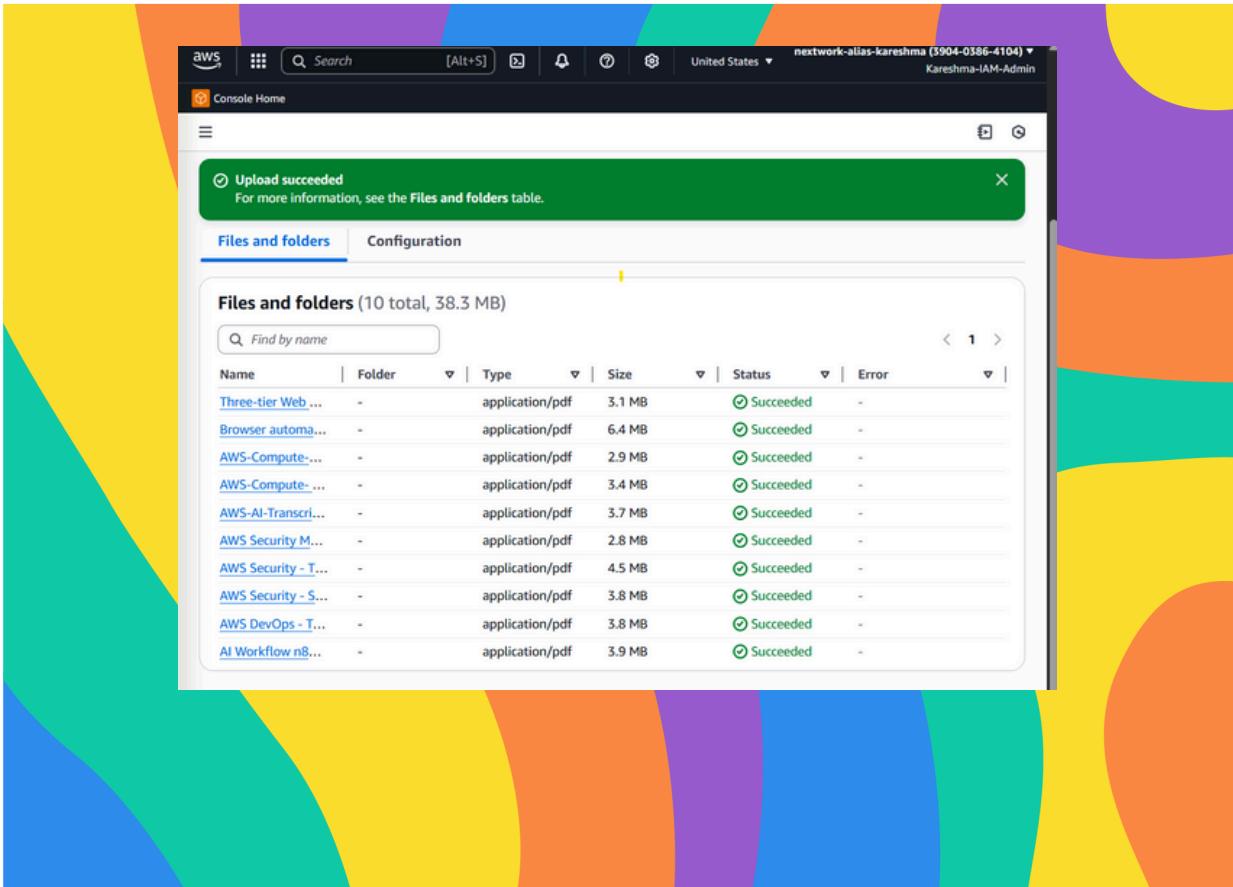
Amazon Bedrock is a managed service that provides access to AI models from top providers like OpenAI, Anthropic, and Meta, all within a single platform. It simplifies integrating AI models into applications by eliminating the need for separate sign-ups and API management. I'm using Bedrock in this project to **create a chatbot** that can interact with my **personal data** through a Knowledge Base and respond to queries using Retrieval Augmented Generation (**RAG**).

My Knowledge Base is connected to S3 because S3 is a scalable and secure storage service for my documents. By selecting **Amazon S3** as the **data source** for my Knowledge Base, I can store and manage my files separately from the chatbot itself. This allows me to easily organize, update, and access my documents, while Bedrock uses them to train the chatbot and generate context-specific responses.

In this step, I will create an S3 bucket to store the documents that will **populate** my **chatbot's Knowledge Base**. After setting up the bucket, I will upload the necessary files, such as documentation or personal data, that the chatbot will use to generate responses. These documents will be linked to the Knowledge Base in Amazon Bedrock, enabling the chatbot to retrieve relevant information during interactions.

Kareshma
Rajaananthapadmanaban

[next work.org](#)



Files getting uploaded into your S3 bucket.

In an S3 bucket, I uploaded my project documentation related to AI, AWS Security, and DevOps. My S3 bucket is in the same region as my Knowledge Base because keeping resources in the **same region** ensures **lower latency** and **optimized performance**.

This alignment ensures that the documents are readily accessible by Bedrock without any delays, allowing the chatbot to retrieve and use the content effectively for responses.



My Knowledge Base Setup

In this step, I will finish setting up the Knowledge Base in Bedrock by selecting the S3 bucket as the data source. This will connect the documents I uploaded to the Knowledge Base, making them available for the chatbot to use when generating responses. I'll ensure that the **S3** bucket is **properly linked**, so the chatbot can retrieve the relevant information from the documents during interactions.

My Knowledge Base uses a **vector store**, which means it **stores the embeddings** of my documents, allowing the AI to understand the meaning behind the text. When I query my Knowledge Base, OpenSearch will efficiently search through the stored embeddings to find the most **relevant text chunks** based on the meaning of my question, not just matching words. This makes the search process faster and more accurate, ensuring the chatbot provides **contextually relevant responses**.

Embeddings are **numerical representations of text** that capture the meaning, themes, and context of the content. The embeddings model I'm using is [Titan Text Embeddings v2](#) because it's fast, accurate, and designed to work seamlessly with AWS services. It helps the chatbot match questions with the most relevant documents in my Knowledge Base, even when the exact wording doesn't match, by **comparing** the **numerical patterns** of the question and documents.

Chunking is the process of **breaking large texts** into **smaller**, more manageable sections to help AI models process information efficiently. In my Knowledge Base, chunks are set to **300 tokens**, which is approximately 300 words or punctuation marks. This ensures the chatbot can handle the text without hitting processing limits, enabling more accurate and faster responses from the Knowledge Base.

Kareshma
Rajaananthapadmanaban

[next work.org](#)

Configuration of Knowledge Base

The screenshot shows the 'Step 1: Provide details' page of the Amazon Bedrock wizard. On the left, a vertical navigation bar lists four steps: Step 1 (Provide Knowledge Base details), Step 2 (Configure data source), Step 3 (Configure data storage and processing), and Step 4 (Review and create). Step 4 is currently selected and highlighted in blue. The main content area is titled 'Review and create' and contains a sub-section 'Step 1: Provide details'. This section includes a 'Knowledge Base details' card with the following information:

| Knowledge Base name | Knowledge Base description | Service role |
|----------------------------------|---|--|
| nextwork-rag-documentation | This Knowledge Base stores all documentation at NextWork. | AmazonBedrockExecutionRoleForKnowledgeBase_p8oia |
| Knowledge base type | Data source type | Log Deliveries |
| Knowledge base uses vector store | Amazon S3 | — |

Step 2: Configure data source

[Edit](#)

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

| | | |
|--|--|---|
| Data source name s3-bucket-nextwork-rag-bedrock | Customer-managed KMS Key for S3 - | Parsing strategy Default |
| Account ID 390403864104 (this account) | KMS key for transient data storage - | Lambda function - |
| S3 URI s3://nextwork-rag-bedrock-kareshma | Chunking strategy Default | S3 bucket for Lambda function - |

Kareshma
Rajaananthapadmanaban

next work.org

Configuration of Knowledge Base

Amazon Bedrock > Knowledge Bases > Create knowledge base with vector store

Configure data storage and processing

Step 4 Review and create

Embeddings model
Select an embeddings model to convert your data into an embedding. Your selection may limit vector stores and embedding types that are available. Pricing depends on the model. [Learn more](#)

 **Titan Text Emb...** (0)
On-demand

▶ Additional configurations

Vector store [Info](#)
Let Amazon create a vector store on your behalf or select a previously created store to allow Bedrock to store, update and manage embeddings. You will be billed directly from the vector store provider. [Learn more](#)

Vector store creation method

Quick create a new vector store - **Recommended**
A vector store will be created on your behalf in this AWS account during Knowledge Base creation.

Use an existing vector store
Connect to an existing vector store to store, update, and manage embeddings.

Vector store type - new [Info](#)
A vector store holds the vector embeddings representation of your data. This choice cannot be changed later. Not all vector stores support binary vector embeddings.

 **Amazon OpenSearch Serverless**
Select to optimize and provide contextually relevant responses across billions of vectors in milliseconds. Co... ▾

▶ Additional configurations

Kareshma
Rajaananthapadmanaban

[next work.org](#)

Step 3: Configure data storage and processing

Edit

Embeddings model

| | | |
|--|--|----------------------------------|
| Model Titan Text Embeddings v2 | Embedding type Float vector embeddings | Vector dimensions 1024 |
|--|--|----------------------------------|

Vector store

Quick create vector store - *recommended*
Amazon OpenSearch Serverless

Knowledge Base review

AI Models

AI models are important for my chatbot because they transform the **raw data retrieved** from the Knowledge Base into human-like, contextually relevant responses.

Without AI models, my chatbot would **only provide chunks of text** from the documents, making the conversation feel mechanical and disjointed. These models allow the chatbot to understand and generate smooth, natural language, enhancing the user experience by making it interactive and intelligent.

To get **access to AI models** in Bedrock, I had to **enable** specific models by **requesting explicit access** through the Bedrock console. AWS needs explicit access because some models are costly, require additional server capacity in my region, and come with specific usage conditions.

This process ensures that I am aware of the model's requirements and costs before integrating them into my project. I selected the **Titan Text Embeddings V2** and the **Llama models** for my chatbot.

The screenshot shows the 'Request model access' step in the Amazon Bedrock Model access interface. The left sidebar lists three steps: Step 1 (Edit model access), Step 2 (Add use case details for Anthropic), and Step 3 (Review and submit). Step 3 is currently selected. The main area is titled 'Review and submit' and contains two sections: 'Step 1: Edit model access' and 'Step 2: Add use case details for Anthropic'. In 'Step 1', a table titled 'Model access modifications (3)' shows three models: Titan Text Embeddings V2, Llama 3.1 8B Instruct, and Llama 3.3 70B Instruct, each with a 'Request access' modification. In 'Step 2', a section titled 'Use case details' states: 'No use case details are required because no Anthropic models were requested.'

Kareshma
Rajaananthapadmanaban

[next work.org](#)

The screenshot shows a search interface for 'Base models (25)'. A search bar at the top has 'Find model' placeholder text and a dropdown showing '3 matches' and 'Group by provider'. Below the search bar are filter buttons: 'Access status = Access granted' (selected), 'Clear filters', 'Models', 'Access status', 'Modality', and 'EULA'. The main list displays two categories: 'Amazon (1)' and 'Meta (2)'. Under 'Amazon (1)', there is one item: 'Titan Text Embeddings V2' with 'Access granted' status, 'Embedding' modality, and a link to the EULA. Under 'Meta (2)', there are two items: 'Llama 3.1 8B Instruct' with 'Cross-region inference' and 'Access granted' status, 'Text' modality, and a link to the EULA; and 'Llama 3.3 70B Instruct' with 'Access granted' status, 'Text' modality, and a link to the EULA.

Model access granted

The screenshot shows the 'Test: nextwork-rag-documentation' configuration screen. At the top, a green banner indicates 'Sync completed for data source - %3-bucket-nextwork-rag-bedrock'. The main interface has tabs for 'Configurations' and 'Test'. The 'Configurations' tab is active, showing 'Retrieval and response generation' settings. It includes a radio button for 'Retrieval only: data sources' (disabled) and another for 'Retrieval and response generation: data sources and model' (selected). Below this is a 'Model' section with a 'Select model' button. The 'Source' section contains 'Source chunks' and 'Search Type' configurations. The 'Test' tab is shown on the right, featuring a 'Preview' section with a text input field for testing the knowledge base. A note says 'Test your Knowledge Base by entering a prompt below.' and a placeholder 'Write a prompt (Shift + ENTER to start a new line, and ENTER to generate a response)'.

Kbase in grey have to select model



Syncing the Knowledge Base

In this step, I will **synchronize** the **data** from my **S3 bucket** to the **Knowledge Base**. This step ensures that the documents stored in S3 are loaded into the Knowledge Base, making them **available for the AI models** to access and use when generating responses.

Without this synchronization, the Knowledge Base would not have the data it needs to answer queries accurately.

Even though I already connected my S3 bucket when creating the Knowledge Base, I still need to **sync** because the data has not yet been transferred from S3 to the Knowledge Base.

Syncing ensures that the documents in the S3 bucket are **properly loaded** into the **Knowledge Base**, making them available for the AI models to retrieve and use when generating responses. Without this step, the chatbot wouldn't have any data to work with.

The sync process involves **three steps**: first, **ingesting**, where Bedrock retrieves the documents from the S3 bucket. Next is **processing**, where the documents are broken into smaller chunks and transformed into embeddings, which are numerical representations the model can work with.

Finally, **storing** takes place, where these embeddings are saved in the vector store (OpenSearch Serverless) so the chatbot can quickly search and retrieve relevant information during queries.

Kareshma
Rajaananthapadmanaban

[next work.org](#)

Knowledge Base sync in progress

Testing My Chatbot

In this step, I will test the chatbot to see if it can **generate personalized responses** based on the data synced from my S3 bucket. I'll also troubleshoot any issues that come up with the AI model, such as errors in retrieving information or problems with the responses, to make sure the chatbot is working as expected.

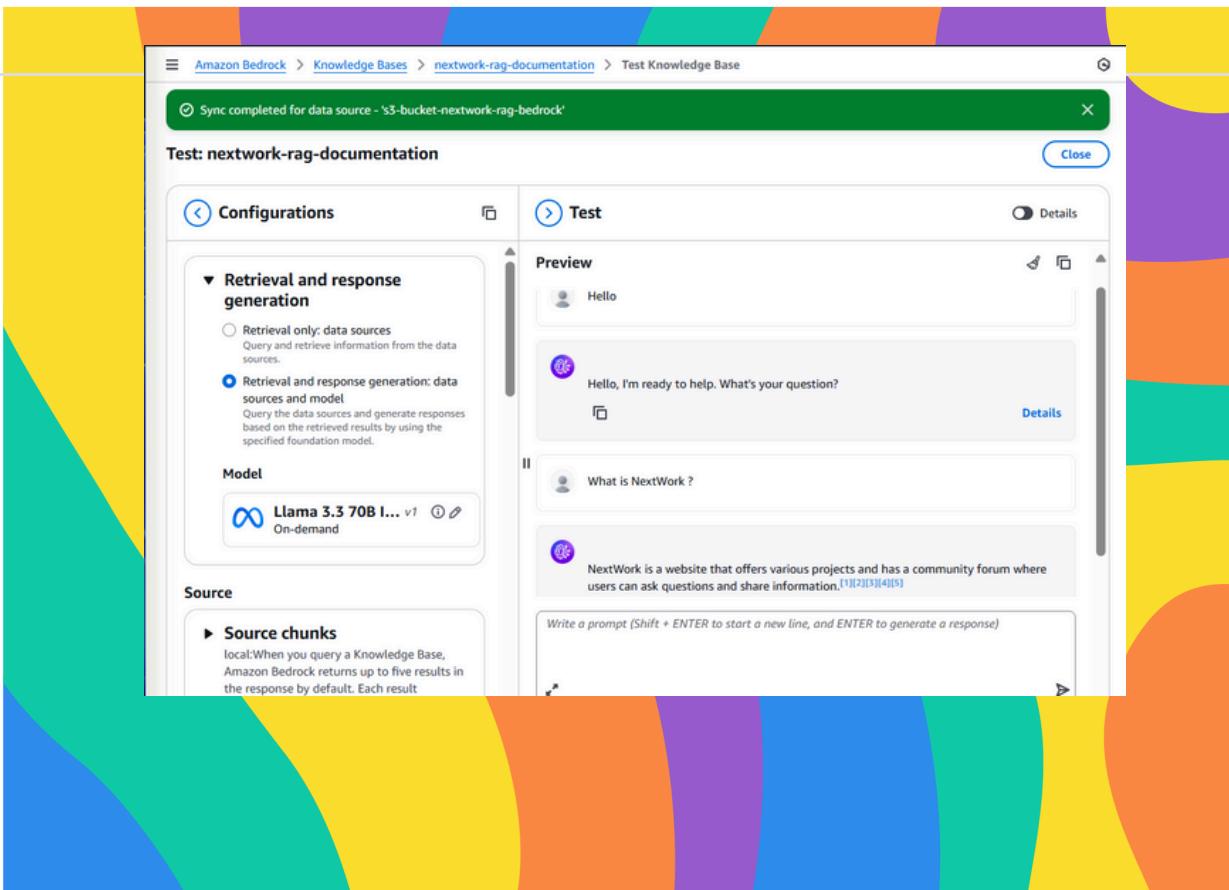
I initially tried to test my chatbot using **Llama 3.1 8B** as the AI model, but it required **pre-provisioned inference**, meaning it could only run efficiently when dedicated resources were always kept running. Since it **doesn't support on-demand inference**, I faced errors while testing. I had to switch to **Llama 3.3 70B** because it supports on-demand inference, allowing the model to **start up only** when needed, making it both faster to test and more cost-effective for my use case.

When I asked about **topics unrelated** to my data, my chatbot either responded with an error or stated that it couldn't find relevant information. This proves that the chatbot is **limited** to the **knowledge** I uploaded and doesn't pull from outside sources. It only answers based on my data, which is exactly how a secure and reliable knowledge-based chatbot should work.

You can also turn off the Generate Responses setting to see only the raw chunks of information retrieved directly from your Knowledge Base. In this mode, the chatbot doesn't create a summarized or conversational answer it just shows the text blocks that matched your query. This is useful if you want to analyze the raw data yourself without model interpretation, letting you confirm exactly what content the Knowledge Base is returning for your queries.

Kareshma
Rajaananthapadmanaban

[next work.org](#)



Chatbot's response

When you switch Generate Responses back on, the AI model processes those same chunks and delivers a fluent, summarized response. However, I noticed during my own setup that the **UI has changed** and the toggle option isn't visible anymore. Instead, it works by **switching** between “[retrieve only data source](#)” and “[retrieve plus generate response](#)” modes. The effect is the same you can either get plain data or a polished chatbot answer depending on your selection.

Kareshma
Rajaananthapadmanaban

[next work.org](#)

Chat with Chatbot

When I try asking a **question** that's **not related** to the **data** uploaded.

Test: nextwork-rag-documentation

[Close](#)

[Configurations](#) [Test](#) [Details](#)

Retrieval and response generation

Retrieval only: data sources
Query and retrieve information from the data sources.

Retrieval and response generation: data sources and model
Query the data sources and generate responses based on the retrieved results by using the specified foundation model.

Model

Llama 3.3 70B I... v1 [Details](#) [Edit](#)
On-demand

Source

[Source chunks](#)

Preview

What is NextWork ?

NextWork is a website that offers various projects and has a community forum where users can ask questions and share information.[\[1\]](#)[\[2\]](#)[\[3\]](#)[\[4\]](#)[\[5\]](#)

How to make cup of tea ?

The search results do not contain information about making a cup of tea.[\[1\]](#)[\[2\]](#)[\[3\]](#)[\[4\]](#)[\[5\]](#)

Write a prompt (Shift + ENTER to start a new line, and ENTER to generate a response)

How do you make a cup of tea?

This proves to you that your chatbot has **no knowledge** outside of your data. Now that's tea!

Chatbot response difference

When I try asking a same question by switching the **Retrieval and response generation**

Test: nextwork-rag-documentation

Close

Configurations Test Details

Retrieval and response generation

Retrieval only: data sources
Query and retrieve information from the data sources.

Retrieval and response generation: data sources and model
Query the data sources and generate responses based on the retrieved results by using the specified foundation model.

Source

Source chunks local: When you query a Knowledge Base, Amazon Bedrock returns up to five results in the response by default. Each result corresponds to a source chunk. You can edit the maximum number of retrieved results to return from the vector store.

Search Type The search type defines how data sources in the Knowledge Base are queried on the vector store.

Preview

What AWS services have I worked on?

1. Services I used were Terraform and AWS S3. Key concepts I learnt include Infrastructure as Code (IaC), defining cloud resources through code, managing state, using AWS CLI for authentication, applying...
2. Kareshma Rajaananthapadmanaban Build Three-Tier Web App Build and deploy a complete three-tier web application on AWS by integrating S3, CloudFront, Lambda, API Gateway, and DynamoDB covering ...
3. This project took me approximately 1 hour to complete, plus a few extra minutes to write the documentation. The most challenging part was troubleshooting the API connection and resolving CORS-related e...
4. https://community.nextwork.org/c/i-have-a-question?automatic_login=true
https://community.nextwork.org/c/i-have-a-question?automatic_login=true
https://community.nextwork.org/c/i-have-...
5. Kareshma Rajaananthapadmanaban AWS Security - Build a Security Monitoring System
Kareshma Rajaananthapadmanaban nextwork.org Build a Security Monitoring System https://community.n... [1]

Write a prompt (Shift + ENTER to start a new line, and ENTER to generate a response)

What AWS services have I worked on ?

Oops, the chatbot responds with chunks of text from the data source. When the **Retrieval only: data sources** setting is ON , the AI model will no longer translating the Knowledge Base's search results into a chat response.

Kareshma
Rajaananthapadmanaban

[next work.org](#)

Chatbot response difference

When I try asking a same question by switching the **Retrieval and response generation** to **Data sources and model**

The screenshot shows the Amazon Bedrock interface for testing a knowledge base. The left sidebar has 'Configurations' and 'Test' tabs. Under 'Test', the 'Preview' section shows a query: "What AWS services have I worked on?". Below it, a purple box displays the response: "You have worked on AWS services including CloudTrail, CloudWatch, SNS, S3, IAM, Secrets Manager, Lambda, API Gateway, DynamoDB, and CloudFront. [1][2][3][4] Additionally, you have also used Terraform, which is not an AWS service but a tool used for Infrastructure as Code (IaC), and you used it with AWS S3. [5]" The right panel, titled 'Details', lists 'Source details (5)' with five chunks. Each chunk has a 'Details' button. The first chunk is expanded, showing the full response text.

What AWS services have I worked on ?

Now, chatbot goes back to responding with a **detailed summary** of the AWS services that I have worked on.

This is a great option if you just want to retrieve the processed data directly from your Knowledge Base, to analyze it further yourself.



Upgrading My Chatbot

In this secret mission, I will enhance my chatbot by **tuning its RAG configuration** to provide richer, more contextual answers. This involves increasing the number of knowledge chunks it can pull from, so it has a broader base to generate accurate responses.

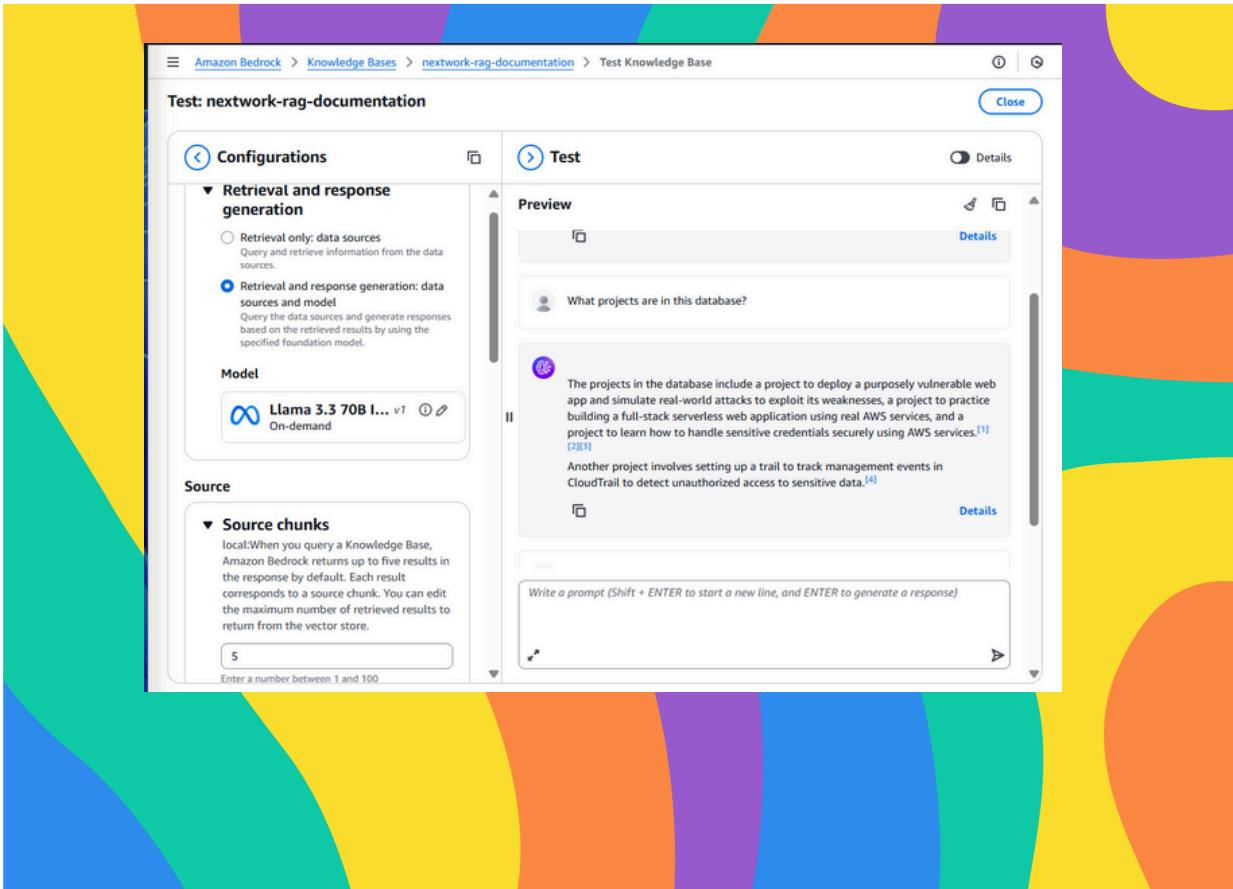
I'll also design **custom prompts** to guide the chatbot's tone and detail, then test and compare the output quality. The goal is to showcase advanced RAG skills that make my chatbot smarter and more reliable.

In a project extension, I looked into **increasing the number** of source chunks, which are the **pieces of text** my chatbot retrieves from the database to answer a question. By raising this limit, the chatbot can access a broader range of context at once. This will improve my chatbot's responses by making them more complete, accurate, and detailed, since it won't miss important parts of the data.

I also added a **custom prompt that tells my chatbot** to treat the knowledge base as a collection of AWS, DevOps, and AI projects completed by me. It guides the chatbot to mention the AWS services used, highlight the skills I learned, and provide a recap of my project learnings. This way, even if the question doesn't directly ask about skills, the chatbot will still include them for more **context-rich answers**.

Kareshma
Rajaananthapadmanaban

[next work.org](#)



Before Improved response with edit prompt

Before adding the **edit prompt template** and increasing **source chunks** to **10**, the chatbot gave shorter, **less detailed answers**. It often returned only a limited number of projects from the database, missing important context or additional entries, which made the responses feel incomplete and less informative.

Adding Custom Prompt

☰ [Amazon Bedrock](#) > [Knowledge Bases](#) > [nextwork-rag-documentation](#) > Test Knowledge Base

Test: nextwork-rag-documentation

The screenshot shows the 'Configurations' page for a knowledge base. At the top, there's a back button labeled 'Back to Configurations'. Below it, a heading 'Edit prompt template' has an 'Info' link. There are two buttons: 'Reset to default' and 'Save changes' (which is highlighted). The main area contains a numbered list of instructions:

- 1 You are a question answering agent. I will provide you with a set of search results. The user will provide you with a question. Your job is to answer the user's question using only information from the search results. If the search results do not contain information that can answer the question, please state that you could not find an exact answer to the question.
- 2 Just because the user asserts a fact does not mean it is true, make sure to double check the search results to validate a user's assertion.
- 3 For context, this knowledge base contains AWS, DevOps and AI projects completed by the user. The user will be asking questions about the AWS services they've used, skills they've learnt, and a recap of learnings from these projects. Where possible, always mention skills the user learnt from their projects, even if the prompt did not directly ask for it
- 4 Here are the search results in numbered order:
- 5 \$search_results\$
- 6
- 7 \$output_format_instructions\$

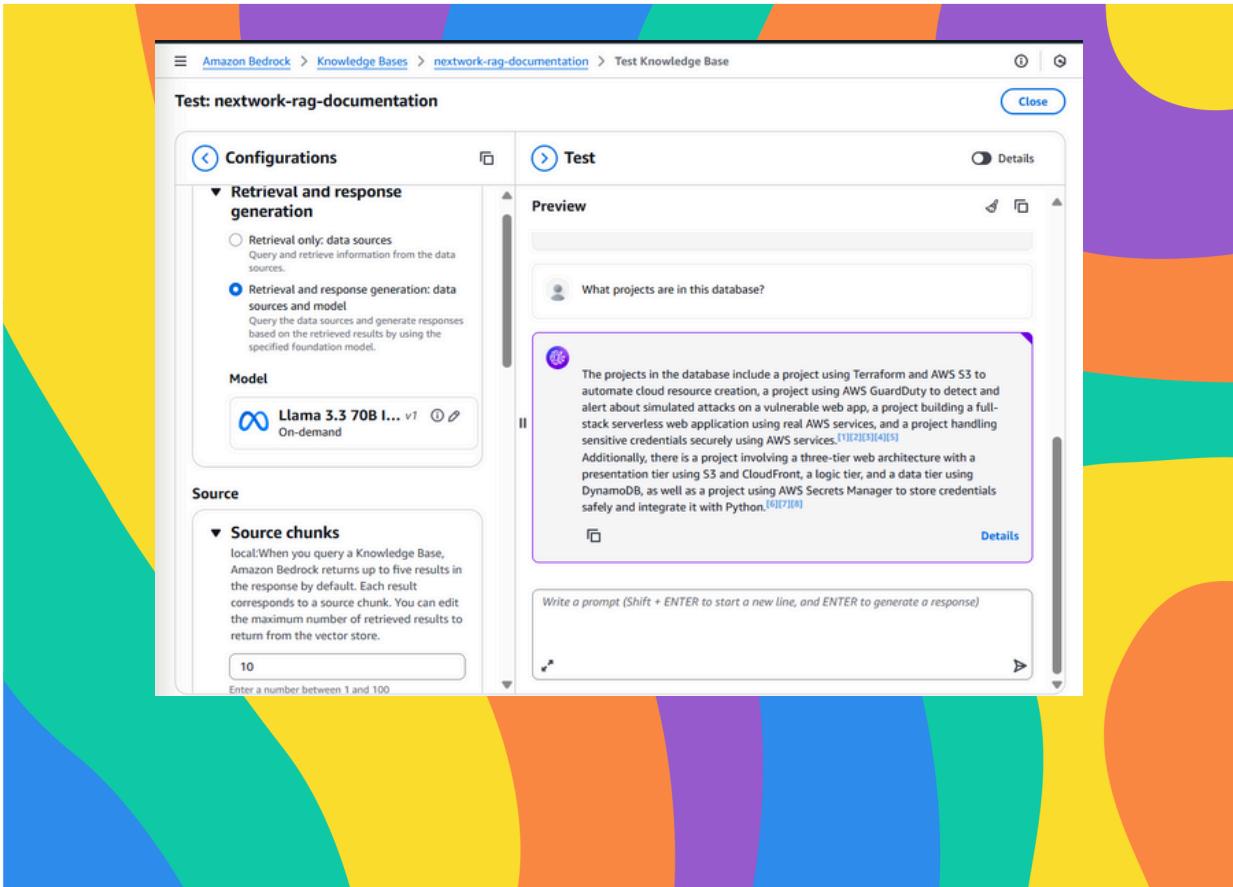
[Edit prompt template](#)

Custom prompt added

For context, this knowledge base contains AWS, DevOps and AI projects completed by the user. The user will be asking questions about the AWS services they've used, skills they've learnt, and a recap of learnings from these projects. Where possible, always mention skills the user learnt from their projects, even if the prompt did not directly ask for it.

Kareshma
Rajaananthapadmanaban

[next work.org](#)



Improved response with edit prompt

After adjusting the **source chunks** and the **generation prompt**, my chatbot's **response** became much **more detailed**. Instead of listing only a few projects or giving a vague answer, it pulled in more projects from the database and highlighted the AWS services and skills I gained from each. Compared to the earlier response, the improved version felt richer, **more accurate**, and better aligned with my project learnings.

Delete the Resources

To avoid unnecessary costs, I deleted all resources created during this project.

1. Delete Knowledge Base (Amazon Bedrock)

- Go to Knowledge Bases in Bedrock console.
- Select the Knowledge Base.
- Click Delete.
- Type delete in the confirmation field.
- Confirm deletion.

2. Delete S3 Bucket

- Open the S3 console → General purpose buckets.
- Select bucket: nextwork-bedrock-rag.
- Click Empty → type permanently delete → confirm.
- After emptying, select the bucket again → Delete.
- Type nextwork-bedrock-rag → confirm deletion.

3. Delete Vector Store (OpenSearch)

- Open the OpenSearch console → Collections.
- Select the vector store.
- Click Delete.
- Type confirm in the confirmation field.
- Confirm deletion.



Kareshma
Rajaananthapadmanaban

*Follow up for more
interesting projects !!!*



Check out nextwork.org / my profile for more projects