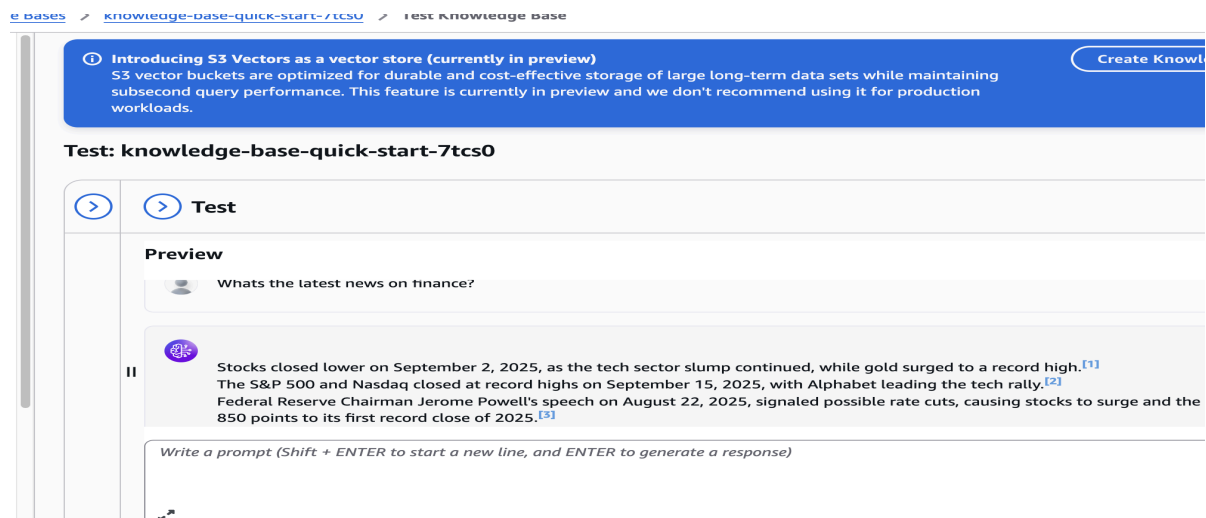
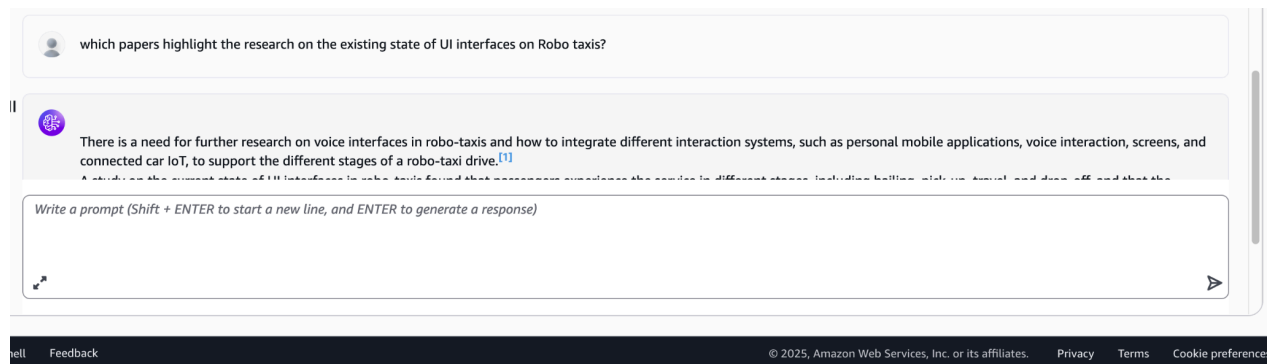


# Set Up a RAG Chatbot in AWS Bedrock



*Jerin George*

# Table of Contents

<b>Introduction .....</b>	<b>3</b>
<b>Used Tools &amp; Concepts .....</b>	<b>3</b>
<b>Objectives .....</b>	<b>3</b>
<b>Project Challenges .....</b>	<b>4</b>
<b>Amazon Bedrock .....</b>	<b>4</b>
<b>Data Source .....</b>	<b>4</b>
<b>Amazon S3 .....</b>	<b>4</b>
<b>Web Crawler .....</b>	<b>5</b>
<b>Knowledge Base Configuration .....</b>	<b>7</b>
<b>AI Models .....</b>	<b>8</b>
<b>Knowledge Base Synchronization .....</b>	<b>9</b>
<b>Testing My Chatbot .....</b>	<b>10</b>
<b>Demonstration .....</b>	<b>11</b>
<b>References .....</b>	<b>13</b>

**Introduction:**


**RAG (Retrieval Augmented Generation) chatbot** is an AI-powered chatbot that enhances AI language models by incorporating documentation as a data source. This technique enables the chatbot to provide accurate and up-to-date responses, even if the required information is the latest or unavailable online. While non-RAG chatbots rely on pre-trained data and may struggle to offer current or specialized knowledge, the RAG chatbots can retrieve information from a specific knowledge base or external sources. This makes RAG chatbots particularly useful for applications, where up-to-date information is critical, such as legal or government updates, product releases, or personalized content.


## **Used Tools and Concepts:**

The AWS services I used in this project were **Amazon Bedrock**, **S3**, and **OpenSearch Serverless**. Key concepts include storing data in S3, creating a Knowledge Base, requesting access to AI models, how chatbot generates responses using AI models and Knowledge Base, and utilizing vector stores for efficient retrieval.

## **Objectives:**

## Use Case 1 – Thesis Research Assistant

 which papers highlight the research on the existing state of UI interfaces on Robo taxis?




There is a need for further research on voice interfaces in robo-taxis and how to integrate different interaction systems, such as personal mobile applications, voice interaction, screens, and connected car IoT, to support the different stages of a robo-taxi drive.<sup>[1]</sup>


A study on the current state of UI interfaces in robo-taxis found that passengers experience the service in different stages, including hailing, pick-up, travel, and drop-off, and that the design of the interface should take into account these stages to provide a better passenger experience.<sup>[2][3]</sup>

Another study used a user-design activity to develop an experience-design framework for autonomous taxi services, which included the development of a Human-Machine Interface (HMI) design for future autonomous taxis.<sup>[4]</sup>

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




hell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preference




There is a need for further research on voice interfaces in robo-taxis and how to integrate different interaction systems, such as personal mobile applications, voice interaction, screens, and connected car IoT, to support the different stages of a robo-taxi drive.<sup>[1]</sup>

A study on the current state of UI interfaces in robo-taxis found that passengers experience the service in different stages, including hailing, pick-up, travel, and drop-off, and that the design of the interface should take into account these stages to provide a better passenger experience.<sup>[2][3]</sup>

Another study used a user-design activity to develop an experience-design framework for autonomous taxi services, which included the development of a Human-Machine Interface (HMI) design for future autonomous taxis.<sup>[4]</sup>

 Details

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



hell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preference

## Use Case 2 – Financial News Updates

The most challenging part of this project was creating the Knowledge Base since I encountered several errors related to permission policies. I struggled to determine which permission policies to add and how to configure them in the IAM console to grant the necessary service access in the IAM account. Additionally, setting up the data source for the Web Crawler was difficult – my initial attempts to sync the data either failed or never completed successfully. By researching solutions through the AWS blog, I was able to resolve these issues. Overcoming these challenges provided valuable insights into AWS service integration and permission management.

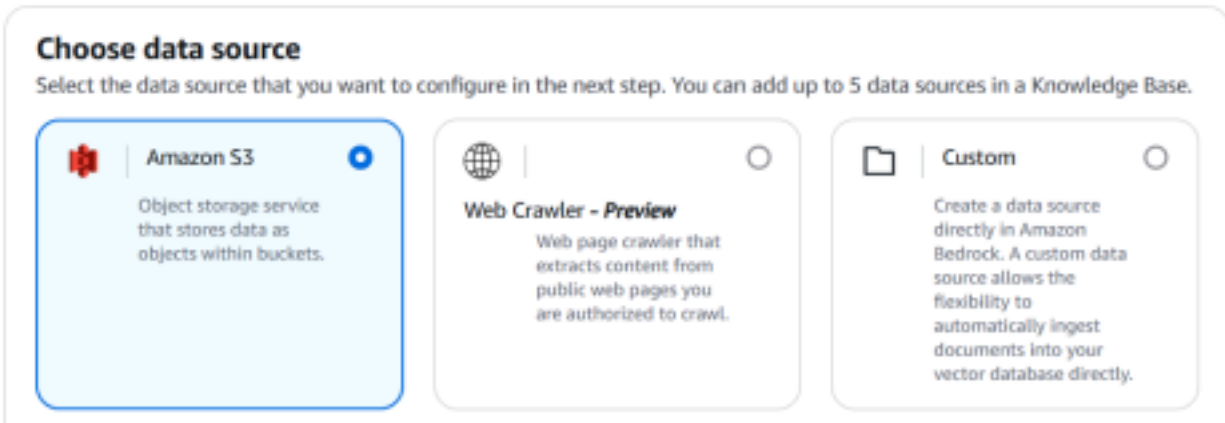
Amazon Bedrock is one of the AWS services that simplifies the development of generative AI applications. It acts as an AI model marketplace, where developers and engineers can find, use, and test AI models from different providers like Meta. In this project, I used Amazon Bedrock to create the Knowledge Base, enabling the chatbot to retrieve relevant information efficiently.

### Data Source (Amazon S3 and Web Crawler):

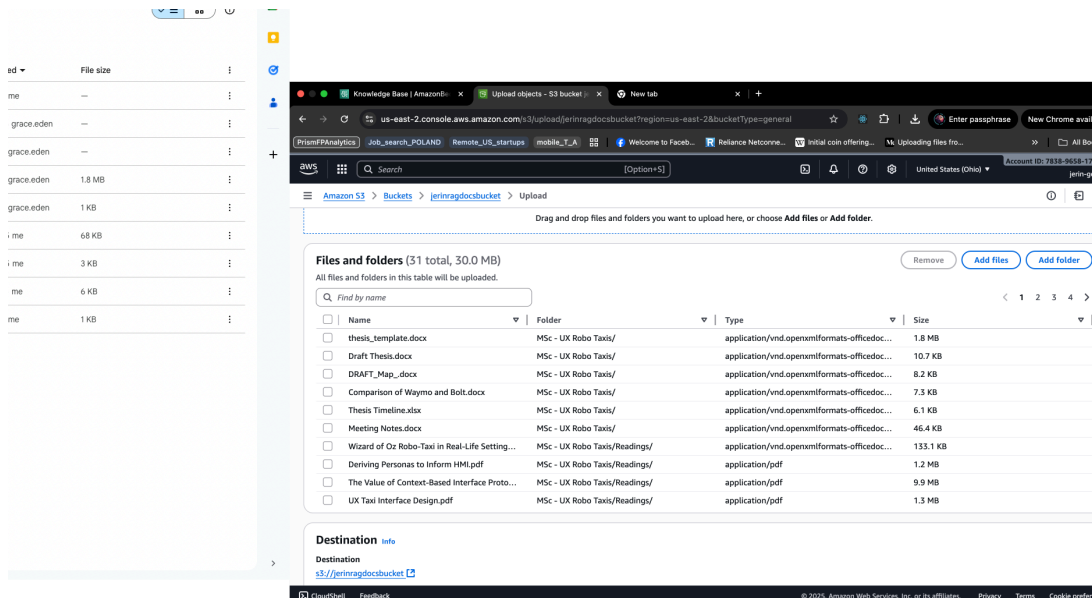
## Use Case 1 – Store Documentation in S3

The Knowledge Base is linked to Amazon S3, which serves as the storage source for the Knowledge Base's row documentation. Amazon S3 is AWS's scalable storage service, where developers and engineers can store various types of objects such as documents, videos, and audio within the same bucket.

### Data source details



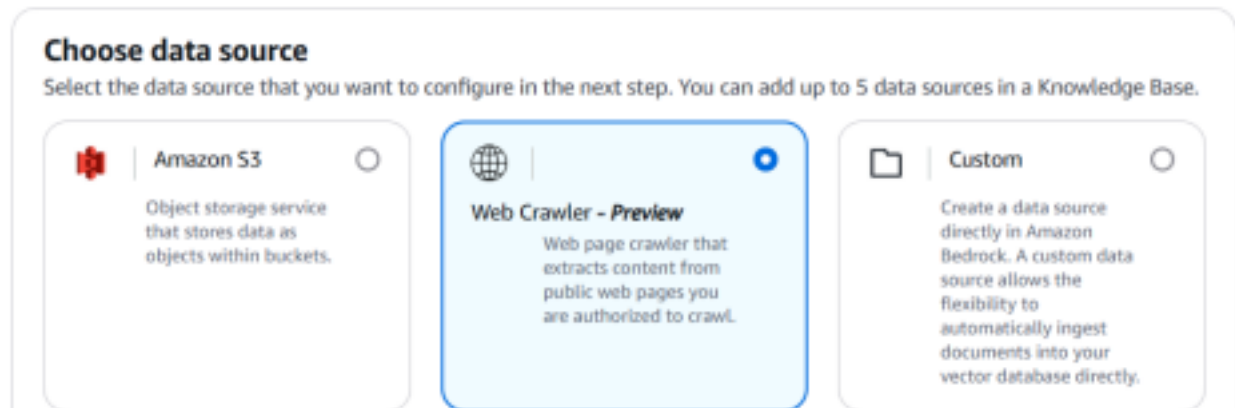
In the S3 bucket, I uploaded the documentation that makes up the AI chatbot's knowledge. Initially, the setup did not work, but I figured out that Amazon Bedrock is a regional service soon after. This means that both the S3 bucket and the Knowledge Base must be in the same AWS region for proper integration. After ensuring they were aligned, the connection functioned as expected.



## Use Case 2 – Add Source URLs

Instead of using S3, I selected **Web Crawler** as a data source and added source URLs. This is because the Web Crawler extracts content from the specified public web pages, enabling the chatbot to provide the latest information as long as those pages are updated. This ensures that responses remain current and reliable without requiring manual data uploads.

### Data source details



Under **Sync Scope**, I initially selected “Default”, causing errors or preventing the data source from successfully synchronizing with the Knowledge Base. To resolve this issue, I switched to “Host only” and added the following URL regex filters. This adjustment finally allowed me to complete the synchronization successfully.

## Knowledge Base Configuration:

***Note:** From this step, the workflows for Use Case 1 and Use Case 2 are the same.*

My Knowledge Base uses a vector store, which is a search engine/database that stores data based on their semantic meaning. When users query the Knowledge Base, OpenSearch will find the relevant chunks of data and pass them to Bedrock for response generation.

Embeddings are vector representations of the semantic meaning of a text chunk. I selected **Titan Text Embeddings v2** as the embedding model because it is fast, accurate, and affordable.

**Chunking** is the process of splitting up text into smaller pieces or chunks. This helps improve data search efficiency within the vector store. In the Knowledge Base, chunks are configured to be approximately 300 tokens in size each.

## AI Models:

AI models play a crucial role in the chatbot, as they translate the search results from the Knowledge Base into human-like text. Without AI models, the chatbot would only respond with raw chunks of text from the documents, which would not provide the best user experience.

To get access to AI models in Bedrock, I had to visit the “Model Access” page and request access explicitly. AWS requires this explicit access since some AI model providers have extra forms or rules for usage, and AWS needs to verify the availability of these models.

Models	Access status	Modality
▼ Amazon (4)	1/4 access granted	
Titan Text Embeddings V2	✓ Access granted	Embedding
Nova Pro <a href="#">Cross-region inference</a>	⊖ Available to request	Text & Vision
Nova Lite <a href="#">Cross-region inference</a>	⊖ Available to request	Text & Vision
Nova Micro <a href="#">Cross-region inference</a>	⊖ Available to request	Text
▼ Anthropic (5)	0/5 access granted	
Claude 3.7 Sonnet <a href="#">Cross-region inference</a>	⊖ Available to request	Text & Vision
Claude 3.5 Haiku <a href="#">Cross-region inference</a>	⊖ Available to request	Text
Claude 3.5 Sonnet v2 <a href="#">Cross-region inference</a>	⊖ Available to request	Text & Vision
Claude 3.5 Sonnet <a href="#">Cross-region inference</a>	⊖ Available to request	Text & Vision
Claude 3 Haiku <a href="#">Cross-region inference</a>	⊖ Available to request	Text & Vision
▼ Meta (8)	2/8 access granted	
Llama 3.5 70B Instruct	✓ Access granted	Text
Llama 3.2 1B Instruct <a href="#">Cross-region inference</a>	⊖ Available to request	Text
Llama 3.2 3B Instruct <a href="#">Cross-region inference</a>	⊖ Available to request	Text
Llama 3.2 11B Vision Instruct <a href="#">Cross-region inference</a>	⊖ Available to request	Text & Vision
Llama 3.2 90B Vision Instruct <a href="#">Cross-region inference</a>	⊖ Available to request	Text & Vision
Llama 3.1 405B Instruct <a href="#">Cross-region inference</a>	⊖ Available to request	Text
Llama 3.1 70B Instruct <a href="#">Cross-region inference</a>	⊖ Available to request	Text
Llama 3.1 8B Instruct <a href="#">Cross-region inference</a>	✓ Access granted	Text

## Knowledge Base Synchronization:

Although the S3 bucket has been connected to the Knowledge Base since its creation, synchronizing the data is still necessary. This is because synchronization is the process that actually transfers the data from S3 into the Knowledge Base and OpenSearch Serverless.

The synchronization process involves three key steps as below:

1. **Ingesting** – Bedrock takes the data from S3.
2. **Processing** – Bedrock chunks and embeds the data.
3. **Storing** – Bedrock stores the processed data in the vector store, OpenSearch Serverless.



## Testing My Chatbot:

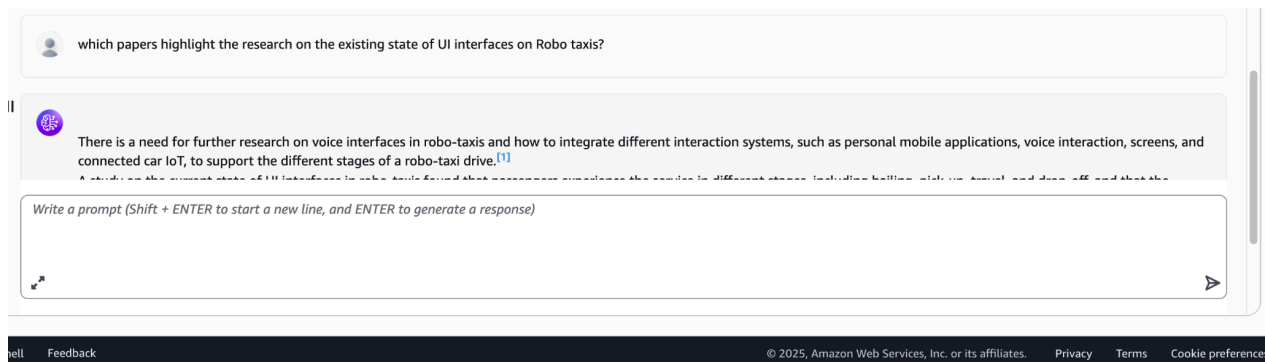
My initial attempt to test my chatbot with **Llama 3.1 8B** as the AI model caused an error because it was unavailable on-demand. To fix this, I switched to **Llama 3.3 70B**, AWS offered on demand as a newer and more efficient model.


When I asked about topics unrelated to the Knowledge Base's data, the chatbot was unable to respond. This means that the chatbot only knows the information I stored in S3. It does not have access to any information outside of the Knowledge Base.

Additionally, you can disable the **Generate Response** setting to see the raw chunks of data directly from the Knowledge Base. During testing, the chatbot just threw a list of paragraphs to answer a question. In contrast, the AI model transforms these chunks into a coherent sentence when generating responses.

## Demonstration:

### Use Case 1 – Thesis Research Assistant








There is a need for further research on voice interfaces in robo-taxis and how to integrate different interaction systems, such as personal mobile applications, voice interaction, screens, and connected car IoT, to support the different stages of a robo-taxi drive.<sup>[1]</sup>

A study on the current state of UI interfaces in robo-taxis found that passengers experience the service in different stages, including hailing, pick-up, travel, and drop-off, and that the design of the interface should take into account these stages to provide a better passenger experience.<sup>[2][3]</sup>

Another study used a user-design activity to develop an experience-design framework for autonomous taxi services, which included the development of a Human-Machine Interface (HMI) design for future autonomous taxis.<sup>[4]</sup>

[Details](#)

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



hell    [Feedback](#)

© 2025, Amazon Web Services, Inc. or its affiliates.    [Privacy](#)    [Terms](#)    [Cookie preferences](#)

## Use Case 2 – Financial News Updates

Introducing S3 Vectors as a vector store (currently in preview)

S3 vector buckets are optimized for durable and cost-effective storage of large long-term data sets while maintaining subsecond query performance. This feature is currently in preview and we don't recommend using it for production workloads.

Create Knowledge Base

## Test: knowledge-base-quick-start-7tcs0



Test

### Preview



Whats the latest news on finance?



II

Stocks closed lower on September 2, 2025, as the tech sector slump continued, while gold surged to a record high.<sup>[1]</sup>  
The S&P 500 and Nasdaq closed at record highs on September 15, 2025, with Alphabet leading the tech rally.<sup>[2]</sup>  
Federal Reserve Chairman Jerome Powell's speech on August 22, 2025, signaled possible rate cuts, causing stocks to surge and the 850 points to its first record close of 2025.<sup>[3]</sup>

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



## References:

*How to set up a RAG chatbot in AWS Bedrock: A Step-by-Step Guide.* (n.d.). Community.aws.  
<https://community.aws/content/2tnxBahN1n18B1Aqfsn1yoNQ7cq/how-to-set-up-a-rag-chatbot-in-aws-bedrock-a-step-by-step-guide>

*Implement web crawling in Amazon Bedrock Knowledge Bases | Amazon Web Services.* (2024, August 29). Amazon Web Services. <https://aws.amazon.com/blogs/machine-learning/implement-web-crawling-in-knowledge-bases-for-amazon-bedrock/>

Sonawane, P. (2025, February 22). Building an AI-Powered RAG chatbot on Bedrock via the AWS console. *Medium*. <https://medium.com/@piyushsonawane10/building-an-ai-powered-rag-chatbot-on-bedrock-via-the-aws-console-8ac62cf3e6b6>