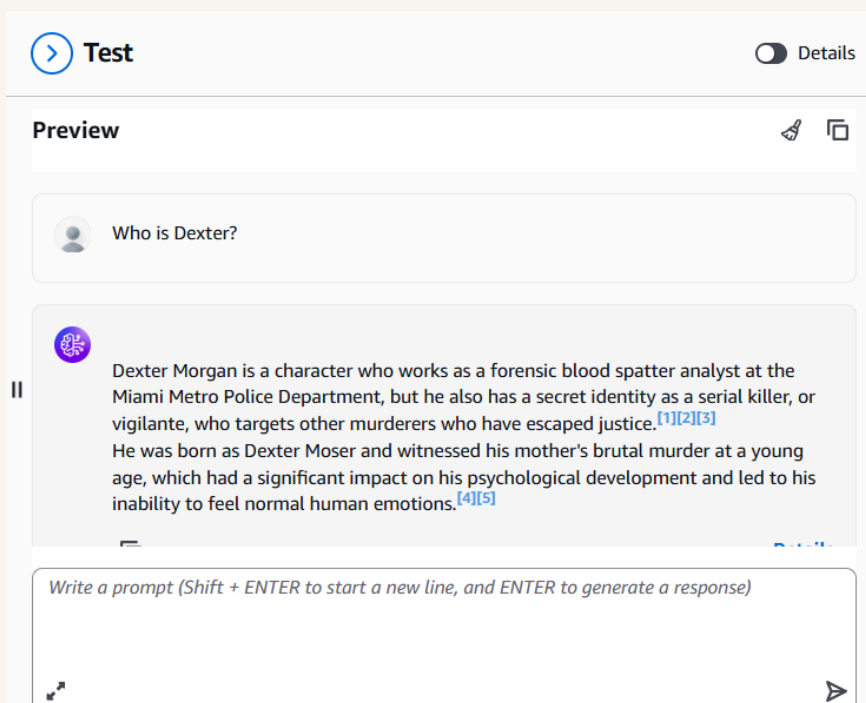# Set Up a RAG Chatbot in Bedrock

Hamza Sajid

# Introducing Today's Project!

RAG is an AI-technique where it recieves data from an external database and uses it to generate an accurate response. In this project, I will demonstrate RAG by using a fake-person info (which I should know) and use it to ask about that person.

## Tools and concepts

Services I used were: 1. Claude - Generating everything about Dexter Morgan and store that Information in the S3 Bucket. 2. S3 Bucket - Storing the Info 3. K-Base - Chunking and get the necessary info

## Project reflection

This project took me approximately took 30 mins to complete.

Enjoyed making a Info about a Psycho-Murder XD.

# Understanding Amazon Bedrock

Amazon Bedrock is an AWS service that amkes it easy for us to build generative AI appliations - it's like a AI Model marketplace that lest us find, use and test moedls from different providers We're using Bedrock to create a Knowledge Base.

My Knowledge Base is connected to S3 becuase I've chosen it as a default data storage to know everything about Dexter Morgan. S3 is the go-to storage for your K-Base to use for retrieving volitial information about Dexter Morgan while being safe.

In an S3 bucket, I uploaded everything about Dexter Morgan. From his Psychological Analysis to his Key Relationships. My S3 bucket is in the same region as my Knowledge Base because it helps reduce latency and most importantly, cost.

⊘ **Upload succeeded**
For more information, see the **Files and folders** table.                          ✕

ⓘ After you navigate away from this page, the following information is no longer available.

**Summary**

| Destination | Succeeded | Failed |
|---|---|---|
| s3://dexter-morgan-rag-bedrock | ⊘ 6 files, 5.3 MB (100.00%) | ⊖ 0 files, 0 B (0%) |

**Files and folders**    Configuration

**Files and folders** (6 total, 5.3 MB)

🔍 Find by name                                                        ‹ **1** ›

| Name | Folder ▽ | Type ▽ | Size ▽ | Status ▽ | Error ▽ |
|---|---|---|---|---|---|
| Dexter Morgan - … | - | application/pdf | 583.2 KB | ⊘ Succeeded | - |
| Dexter Morgan - … | - | application/pdf | 709.1 KB | ⊘ Succeeded | - |
| Dexter Morgan - … | - | application/pdf | 1021.6 KB | ⊘ Succeeded | - |
| Dexter Morgan - … | - | application/pdf | 1.3 MB | ⊘ Succeeded | - |
| Dexter Morgan - … | - | application/pdf | 927.8 KB | ⊘ Succeeded | - |
| The Code of Harr… | - | application/pdf | 908.8 KB | ⊘ Succeeded | - |

**Hamza Sajid**
NextWork Student

nextwork.org

# My Knowledge Base Setup

My Knowledge Base uses a vector store, which means it searchs for the most relevant information based on the User's Input Context. When I query my Knowledge Base, OpenSearch will search, analyze, and visualize large amounts of data quickly.

Embeddings are great to label and oraganise different datas that is stored in a database. The embeddings model I'm using is Titan Text Embeddings V2 because it works very well with AWS related-services.

Chunking is an effecient way to manage large amount text into small chunks for the AI Model to process efficiently. In my K-Base, Chunks are set to 300 tokens per scan (meaning 300 words).

**Review and create**

Step 1: Provide details                                    Edit

**Knowledge Base details**

| Knowledge Base name | Knowledge Base description | Service role |
|---|---|---|
| dexter-morgan-rag-documentation | Everything there is to know about Dexter Morgan | AmazonBedrockExecutionRoleForKnowledgeBase_90rv3 |

| Knowledge base type | Data source type | Log Deliveries |
|---|---|---|
| Knowledge base use vector store | Amazon S3 | — |

Step 2: Configure data source                              Edit

**Data source: d-morgan-rag-bedrock**

| Data source name | Customer-managed KMS Key for S3 | Parsing strategy |
|---|---|---|
| d-morgan-rag-bedrock | - | Default |

| Account ID | KMS key for transient data storage | Lambda function |
|---|---|---|
| 009160060339 (this account) | - | - |

| S3 URI | Chunking strategy | S3 bucket for Lambda function |
|---|---|---|
| s3://dexter-morgan-rag-bedrock | Default | - |

| | | Data deletion policy |
|---|---|---|
| | | DELETE |

# AI Models

AI models are important for my chatbot because they would convert the KBase output into human like response. Without AI models, my chatbot would only output raw results.

To get access to AI models in Bedrock, I had to go to Model Access and select the specific models required for my KBase. AWS needs explicit access so it can know how to use the AI Model

# Syncing the Knowledge Base

Even though I already connected my S3 bucket when creating the Knowledge Base, I still need to sync because the data hasn't actually moved from S3 into your Knowledge Base yet.

The sync process involves three steps: 1. Ingesting 2. Processing 3. Storing
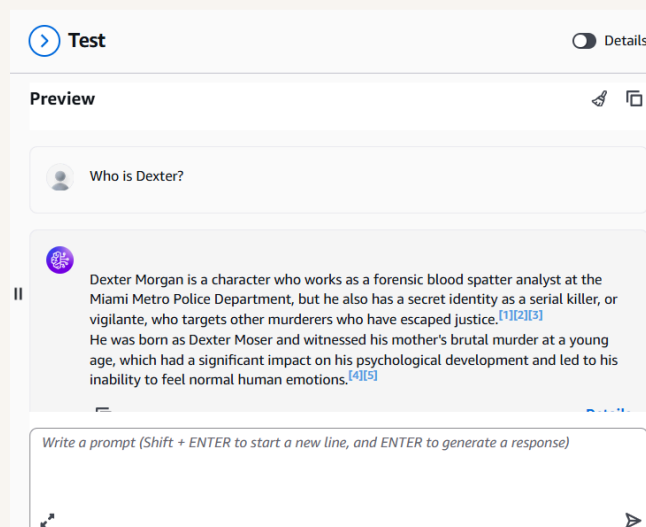
# Testing My Chatbot

'I initially tried to test my chatbot using Llama 3.1 8B as the AI model, but I had to switch to Llama 3.3 70B because it supports on-demand inference!

When I asked about topics unrelated to my data (like Batman), my chatbot replied 'There is no information about Batman in the provided search results.' This proves that the model focus only on what it has been given and had been taught.

You can also turn off the Generate Responses setting to recieve RAW results from your Data.

# The place to learn & showcase your skills

Check out nextwork.org for more projects