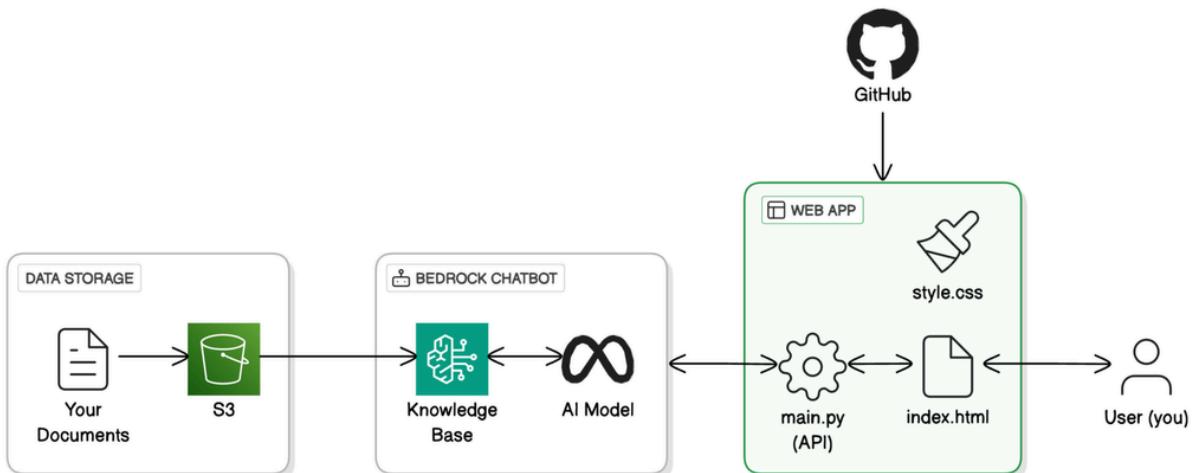




Kareshma
Rajaananthapadmanaban

Build a Web App for RAG Chatbot

Let's build a web app that makes it easy to connect and chat
with my RAG chatbot!

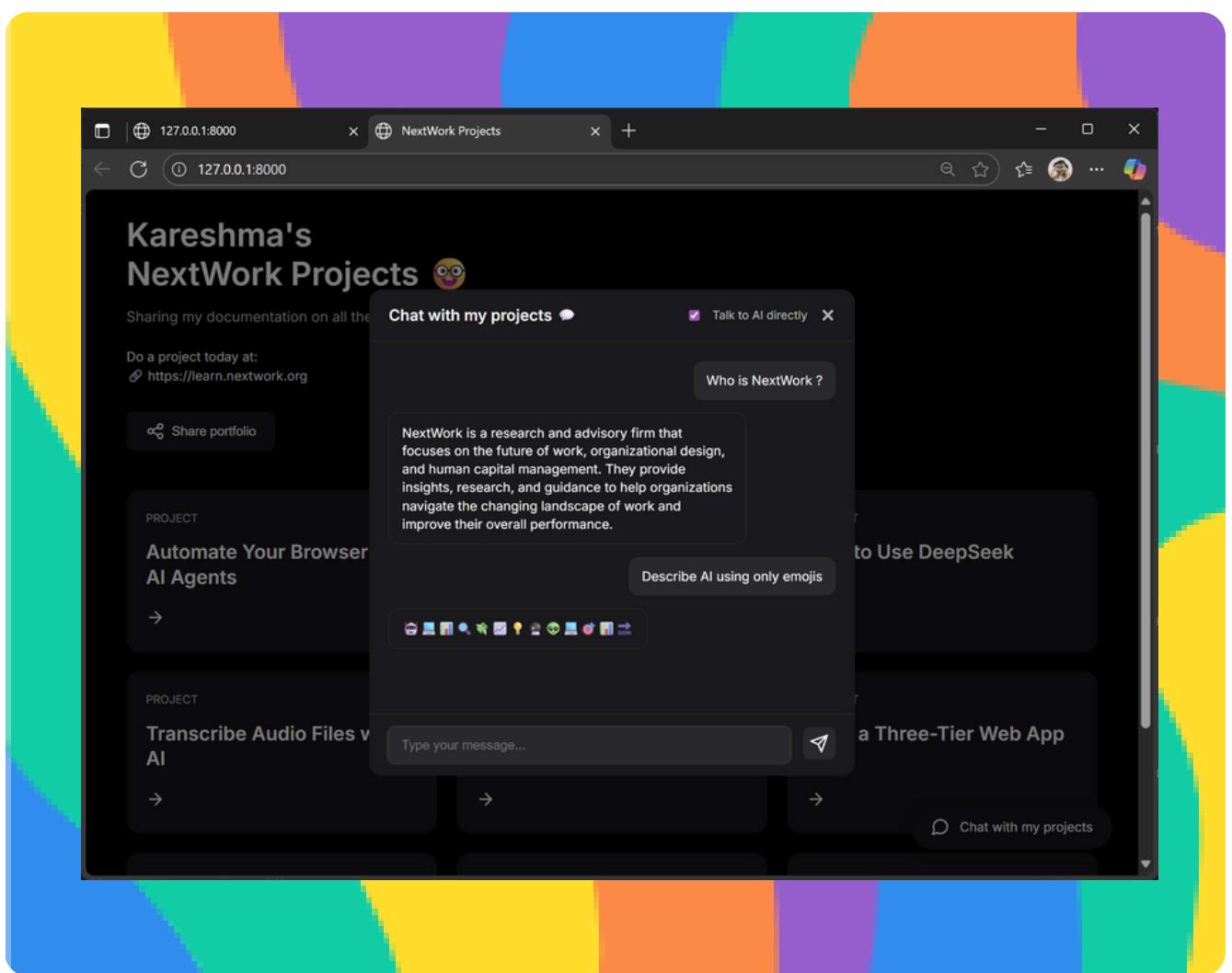




Building a RAG Chatbot with a Web Interface



Kareshma Rajaananthapadmanaban





Introducing Today's Project!

In this project, I will demonstrate how to build a **fully functional web app** for my **RAG chatbot**. I'm doing this project to learn how to connect a **frontend chat interface** with an **API-powered backend** that uses **Amazon Bedrock** and a **knowledge base** in **S3**. This final step is about making the chatbot accessible, scalable, and user-friendly, while gaining hands-on skills in web development and AI integration.

Tools and concepts

Services I used were Amazon Bedrock for model inference, **AWS IAM** for secure access, and **FastAPI** with **Uvicorn** to build and run the chatbot API. I also worked with **dotenv** to manage environment variables and integrated a web frontend to interact with the backend. Key concepts I learnt include Retrieval-Augmented Generation (RAG), separating frontend and backend layers, secure credential handling, and how **enabling or disabling** RAG changes chatbot behavior.

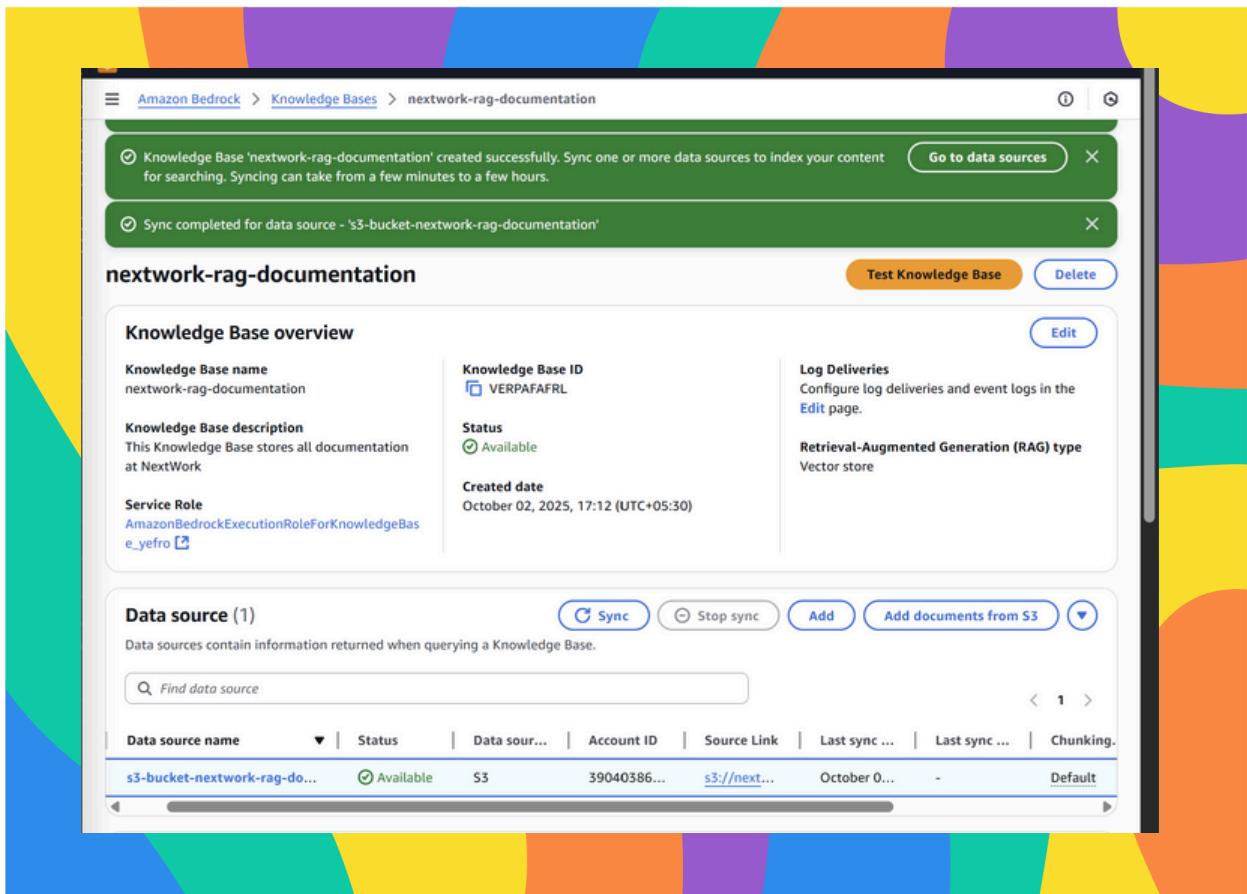
Project reflection

This project took me approximately 2 hours and 30 minutes. The most challenging part was **troubleshooting API integration** and ensuring the chatbot responded correctly with and without RAG enabled. It was most rewarding to see the chatbot pull accurate answers from my knowledge base and experience the difference in responses when switching between **general AI knowledge and RAG-based retrieval**.

I did this project today to practice building with Amazon Bedrock and get hands-on with **Retrieval-Augmented Generation concepts**. My goal was to understand how different AWS services connect to create a working chatbot with a **custom knowledge base**. The project met my goals because I learned how to set up, test, and clean up the entire workflow while strengthening both my **cloud and AI integration skills**.

Chatbot setup

In this step, I will create an **S3 bucket** to store the **documents** my chatbot will use as its knowledge source. I'm doing this to give my RAG chatbot a foundation of information it can reference when generating answers. By uploading my documents to S3 and linking them to a **Bedrock Knowledge Base**, I ensure the chatbot can retrieve accurate, context-aware responses based on my own data.



I created a **Knowledge Base** to give my **RAG chatbot** a structured way to store, organize, and retrieve information from my documents in S3. This is important because without a Knowledge Base, the chatbot wouldn't know how to process or search through the data. By using **embeddings** and a **vector store** in Bedrock, my chatbot can understand context, match questions to relevant content, and deliver accurate, meaningful answers.



Setting Up the API

In this step, I will **clone** a **pre-built web app** from GitHub to my local machine. I'm doing this to quickly set up both the frontend and backend code needed for my RAG chatbot without writing everything from scratch. This starter code includes the **chat interface**, **API handling**, and **integration points** with AWS, so I can focus on connecting it to my Knowledge Base and customizing it.

A **virtual environment** helps me by keeping all the **Python dependencies** for this project **isolated** from my **global setup**. This means any libraries I install for the web app won't affect other projects on my computer.

It also makes the project more portable and reproducible, since I can recreate the same environment anywhere. For this RAG chatbot web app, it ensures clean, conflict-free package management.

I will **test** and **run** the **backend code** of my **RAG web app locally**. I'm doing this to make sure the backend can actually connect to my Bedrock Knowledge Base and handle requests before I combine it with the frontend. Verifying the backend works first helps me catch issues early and ensures my app has a solid foundation before moving forward.

To run the API, I had to install requirements like **fastapi**, **unicorn**, **boto3**, and **python-dotenv**. These packages are important because FastAPI builds the API itself, Unicorn runs the server so my API is accessible, Boto3 connects my app to AWS Bedrock, and python-dotenv loads environment variables securely. Together, they make my chatbot's backend run smoothly and talk to AWS.

```
1.36.20 botocore-1.36.20 click-8.1.8 colorama-0.4.6 fastapi-0.115.8 h11-0.14.0 idna-3.10 jmespath-1.0.1 pydantic-2.10.6 pydantic_core-2.27.2 python-dateutil-2.9.0.post0 python-dotenv-1.0.s3transfer-0.11.2 six-1.17.0 sniffio-1.3.1 starlette-0.45.3 typing_extensions-4.12.2 urllib3-2.3.0 uvicorn-0.34.0
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> pip3 list
Package           Version
-----
annotated-types  0.7.0
anyio             4.8.0
boto3              1.36.20
botocore           1.36.20
click              8.1.8
colorama           0.4.6
fastapi             0.115.8
h11                0.14.0
idna               3.10
Jinja2              3.1.5
jmespath            1.0.1
MarkupSafe          3.0.2
pip                 25.2
pydantic            2.10.6
pydantic_core       2.27.2
python-dateutil     2.9.0.post0
python-dotenv        1.0.1
s3transfer           0.11.2
six                 1.17.0
sniffio              1.3.1
starlette            0.45.3
typing_extensions    4.12.2
urllib3              2.3.0
rn                  0.34.0
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp>
```

Successful Installation of Required Packages



Breaking down the API

In our web app, we need an **API** because it acts as the **bridge** between the **frontend chat interface** and the **backend logic** that talks to Amazon Bedrock. The frontend alone can't process or fetch answers it just displays them. The API handles requests from the user, connects to Bedrock, processes responses, and then sends them back to the frontend for display.

Digging deeper into the API code, the main tools we need to import are **FastAPI**, **boto3**, **os**, and **load_dotenv**. These tools help us by letting us build the API quickly (FastAPI), connect to AWS services like Bedrock (boto3), securely access environment variables (os), and load them from a file (load_dotenv) without hardcoding sensitive information.

The environment variables we need are

`AWS_REGION`, `KNOWLEDGE_BASE_ID`, and `MODEL_ARN`. We store them separately because they contain sensitive information like credentials and IDs. Keeping them in a **separate .env file protects** this information from being exposed in our code or accidentally pushed to GitHub.

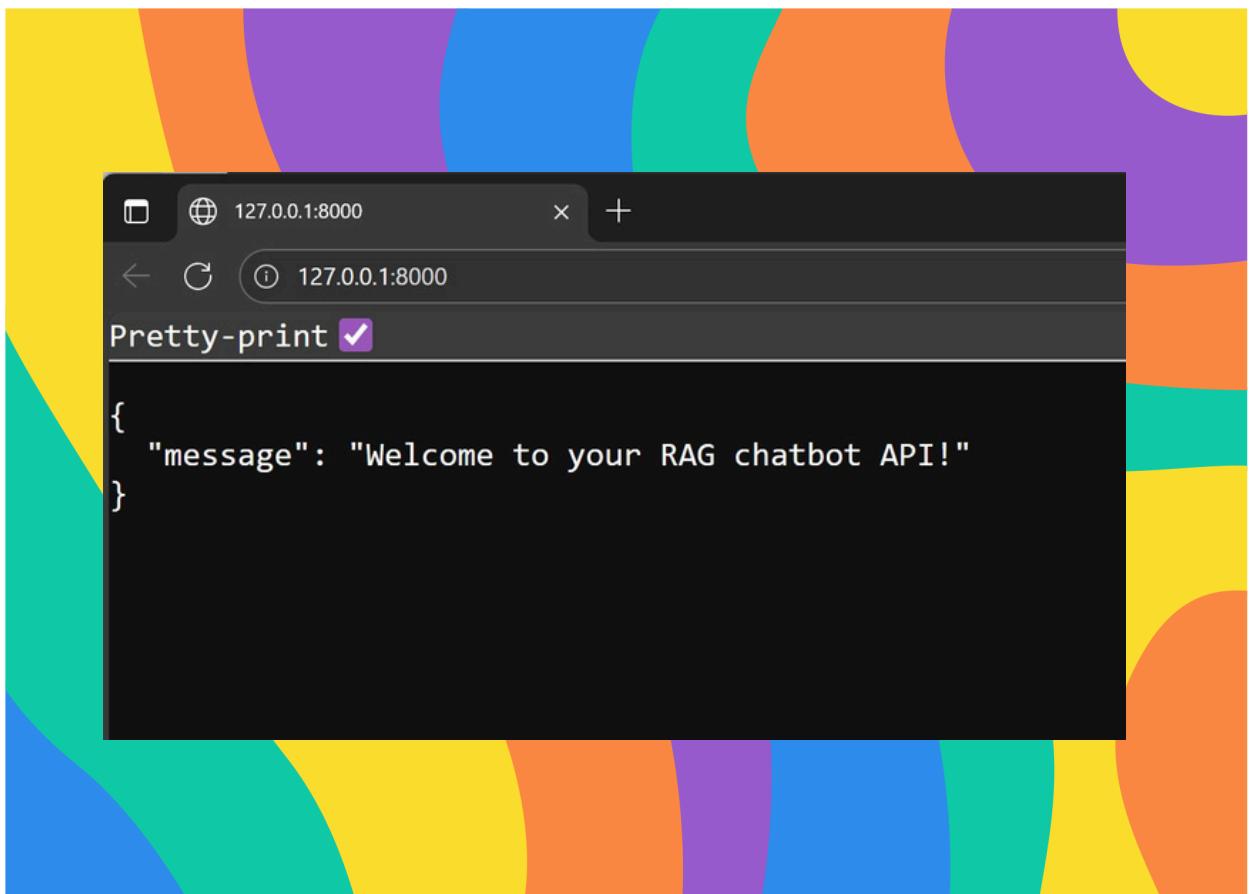
My API uses the Python **AWS SDK** to connect to **Bedrock** from within the code. The difference between the AWS CLI and an SDK is that the CLI is a command-line tool for manual tasks, while an SDK is a **library** for building applications that programmatically interact with AWS services. SDKs allow seamless integration into code, unlike the CLI.

Our API has two main routes: `/` is the **root route** that returns a welcome message, and `/bedrock/query` is where **users send questions** to the chatbot. In general, routes help us define paths in the API that perform specific functions, making it possible for different endpoints to handle different types of requests.



Testing the API

When I visited the root endpoint, I saw `{"message": "Welcome to your RAG chatbot API"}`. This confirms that my **API server** is **running correctly** and responding to requests from the browser, which means the backend setup is working as expected.





Troubleshooting the API

In this step, I will **test the chatbot** by sending a **question to the API** using the **/bedrock/query endpoint**. This is important because it lets me verify that the API can successfully retrieve answers from our Bedrock Knowledge Base and ensures the backend is fully functional before connecting it to the frontend chat interface.

My second endpoint had errors!

To fix the **Parameter Failed error**, I had to create a **.env** file in my project directory and add the correct environment variables: **AWS_REGION**, **KNOWLEDGE_BASE_ID**, and **MODEL_ARN**. Doing this helps because it ensures the API knows which AWS region, Knowledge Base, and AI model to use, preventing invalid parameter errors when querying Bedrock.

To fix the **credentials error**, I ran `aws configure` in my terminal and entered my **IAM user's Access Key ID**, **Secret Access Key**, and the correct **AWS region** (us-east-2). This was necessary because the AWS CLI needs these credentials to authenticate my terminal commands and allow my API to communicate securely with Bedrock.

My chatbot initially failed to generate responses, so I made sure I **had requested access** to the **Titan Text Embeddings V2** and **Llama 3.3 70B Instruct models** in Bedrock, then **synced** my Knowledge Base with the S3 bucket containing my documents. This should help the AI properly process my data and generate accurate, meaningful answers.

A error message for **Invalid type for parameter** is related to **KnowledgeBaseID** or **modelArn** means API isn't correctly configured with the Knowledge Base ID and Model ARN

```
{  
    "detail": "Parameter validation failed:\nInvalid type for parameter  
retrieveAndGenerateConfiguration.knowledgeBaseConfiguration.knowledgeBaseId,  
value: None, type: 'NoneType'\u003E, valid types: 'str'\u003E  
'str'\u003E\nInvalid type for parameter  
retrieveAndGenerateConfiguration.knowledgeBaseConfiguration.modelArn, value:  
None, type: 'NoneType'\u003E, valid types: 'str'\u003E"  
}
```

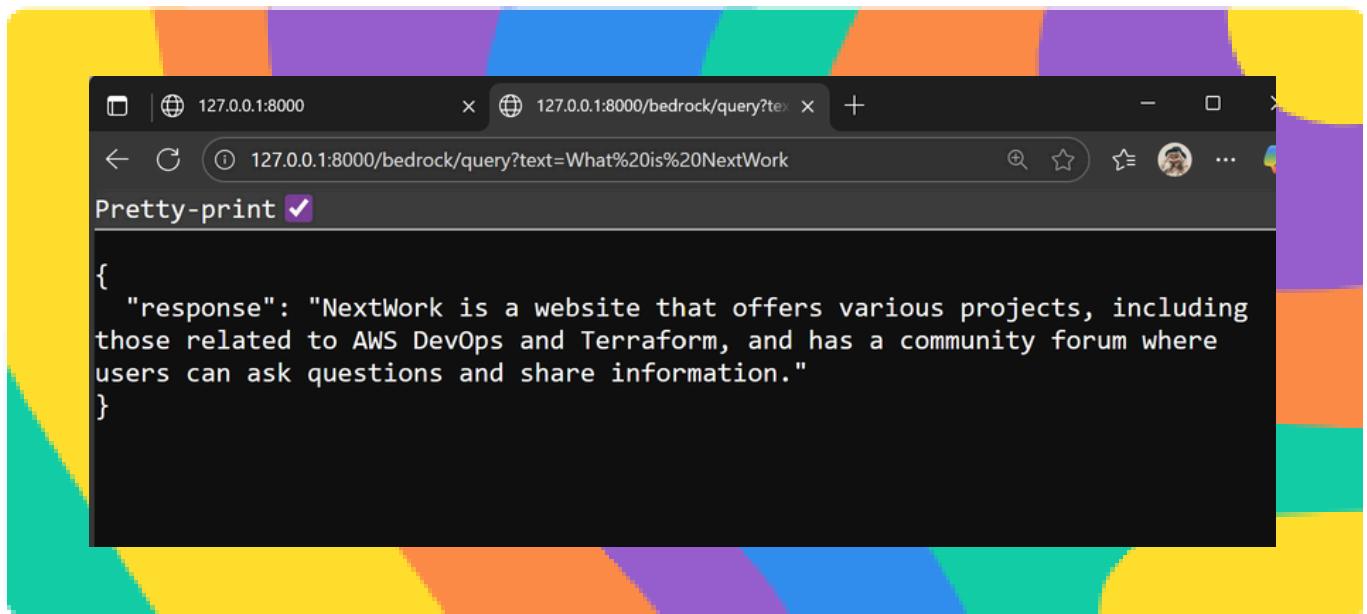
Parameter failed error

```
PowerShell 7 (x64)  
INFO: 127.0.0.1:52633 - "GET /favicon.ico HTTP/1.1" 404 Not Found  
INFO: 127.0.0.1:53473 - "GET / HTTP/1.1" 200 OK  
INFO: 127.0.0.1:64002 - "GET /bedrock/query?text=What%20is%20NextWork HTTP/1.1" 500 Internal Server Error  
INFO: 127.0.0.1:64002 - "GET /favicon.ico HTTP/1.1" 404 Not Found  
INFO: 127.0.0.1:60018 - "GET /bedrock/query?text=What%20is%20NextWork HTTP/1.1" 500 Internal Server Error  
INFO: Shutting down  
INFO: Waiting for application shutdown.  
INFO: Application shutdown complete.  
INFO: Finished server process [21720]  
INFO: Stopping reloader process [10296]  
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> aws configure  
AWS Access Key ID [*****GEWQ]: AKIAVVZPCGYUK6HYBFH3  
AWS Secret Access Key [*****aH43]: WzSgb9hb5gmoFqRwDAo1fvARU5fzm1VuTmfikh  
Default region name [us-east-2]:  
Default output format [json]:  
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> New-Item .env  
Directory: D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp  


| Mode | LastWriteTime     | Length | Name |
|------|-------------------|--------|------|
| -a-- | 10/2/2025 5:20 PM | 0      | .env |

  
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> notepad .env  
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> cat .env  
AWS_REGION = us-east-2  
KNOWLEDGE_BASE_ID = VERPAFAFRL  
MODEL_ARN= arn:aws:bedrock:us-east-2::foundation-model/meta.llama3-3-70b-instruct-v1:0  
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> python -m uvicorn main:app --reload
```

Parameter failed error set up



API's successful response to user's question



Running the Web App

In this step, I will **run** the web app using **web_app.py**, set the **MODEL_ID** environment variable in my **.env file**, and **interact** with the **chatbot** through the browser. I'm doing this to see how the full web interface works, **test responses with and without RAG**, and understand how the frontend, backend, and data layer work together in real time.

The difference between the API and the web app for the user is how they interact with the chatbot. The **API** gives **raw JSON responses** when you hit endpoints like [`/bedrock/query`](#), which is useful for testing and backend integration.

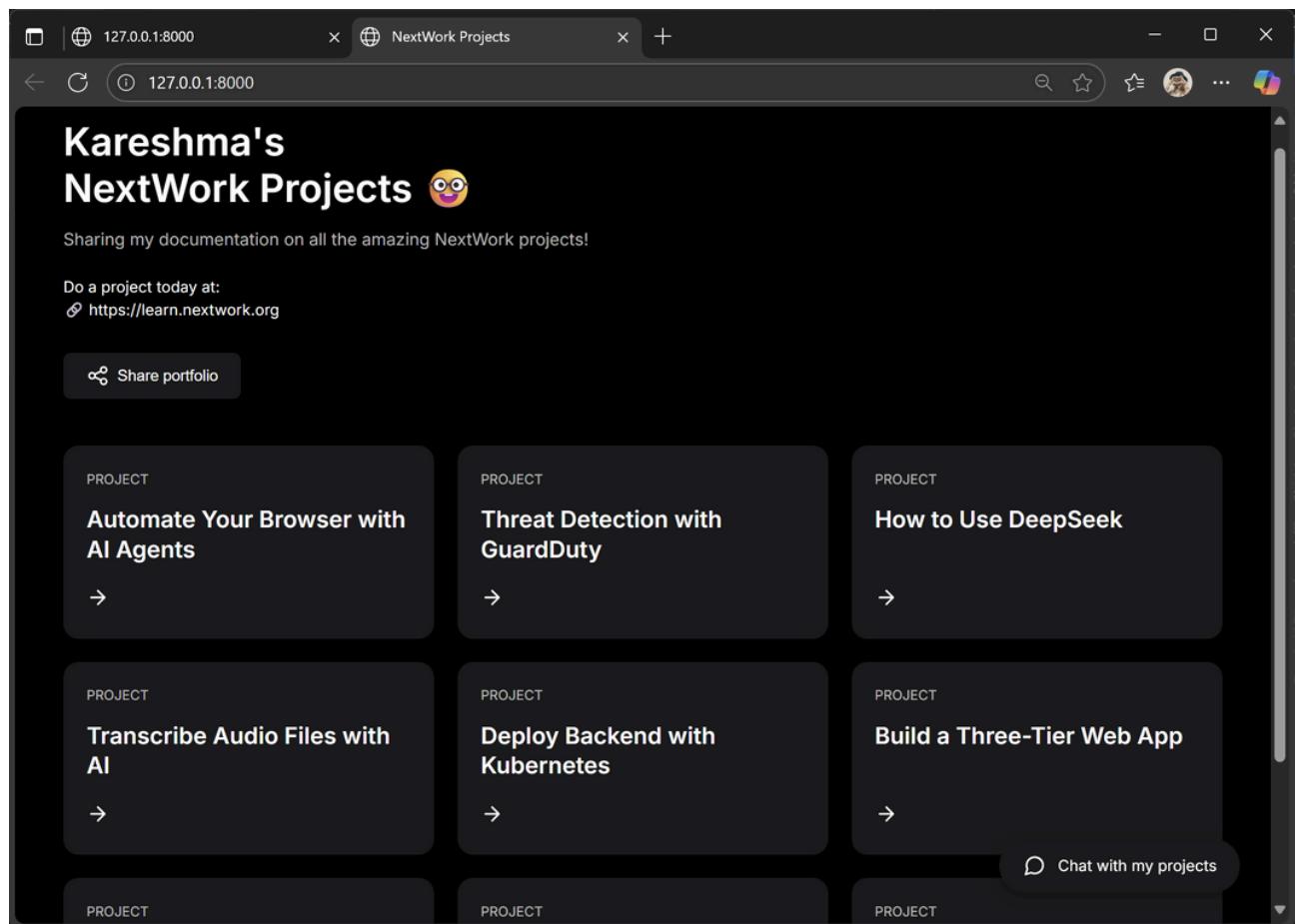
The **web app**, on the other hand, provides a **full chat interface in the browser** where users can type questions and see responses in a **conversational format**.

The **web app** has the **ability to switch** between using **RAG with my Knowledge Base** and **talking directly to the AI model**.

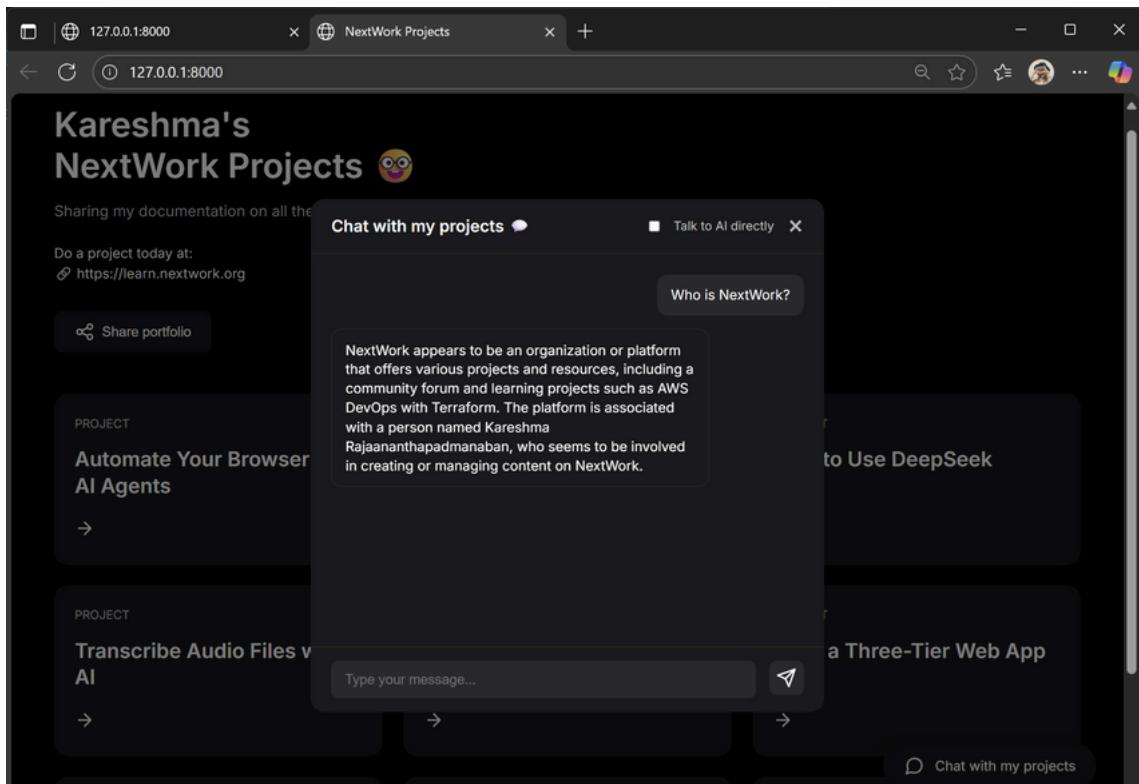
When I asked the same question, “Who is NextWork?”, with **RAG enabled** it gave me a **grounded answer** based on my documents. When I checked “**Talk to AI directly**,” the response came only from the **model’s general knowledge**, which was less accurate.

```
INFO: Started server process [19408]
INFO: Waiting for application startup.
INFO: Application startup complete.
INFO: 127.0.0.1:60355 - "GET /bedrock/query?text=What%20is%20NextWork%20HTTP/1.1" 200 OK
INFO: Shutting down
INFO: Waiting for application shutdown.
INFO: Application shutdown complete.
INFO: Finished server process [19408]
INFO: Stopping reloader process [3556]
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> python -m unicorn web_app:app
--reload
INFO: Will watch for changes in these directories: ['D:\\\\Projects\\\\Kareshma\\\\Nxtwrkcode\\\\nextwork-rag-webapp']
INFO: Uvicorn running on http://127.0.0.1:8000 (Press CTRL+C to quit)
INFO: Started reloader process [16392] using StatReload
INFO:botocore.credentials:Found credentials in shared credentials file: ~/.aws/credentials
INFO: Started server process [1460]
INFO: Waiting for application startup.
INFO: Application startup complete.
|
```

Web app running on Terminal



Web app running on browser



Web app chat with projects RAG Enabled

```
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> python -m uvicorn web_app:app --reload
INFO:     Will watch for changes in these directories: ['D:\\\\Projects\\\\Kareshma\\\\Nxtwrkcode\\\\nextwork-rag-webapp']
INFO:     Uvicorn running on http://127.0.0.1:8000 (Press CTRL+C to quit)
INFO:     Started reloader process [16392] using StatReload
INFO:botocore.credentials:Found credentials in shared credentials file: ~/.aws/credentials
INFO:     Started server process [1460]
INFO:     Waiting for application startup.
INFO:     Application startup complete.
INFO:     127.0.0.1:61717 - "GET / HTTP/1.1" 200 OK
INFO:     127.0.0.1:61717 - "GET /static/style.css HTTP/1.1" 200 OK
INFO:     127.0.0.1:61718 - "GET / HTTP/1.1" 200 OK
INFO:     127.0.0.1:58472 - "GET /bedrock/query?text=Who%20is%20NextWork%20%3F HTTP/1.1" 200 OK
INFO:     127.0.0.1:58477 - "GET /bedrock/invoke?text=Who%20is%20NextWork%20%3F HTTP/1.1" 500 Internal Server Error
```

500 Internal server error when check Talk to AI directly

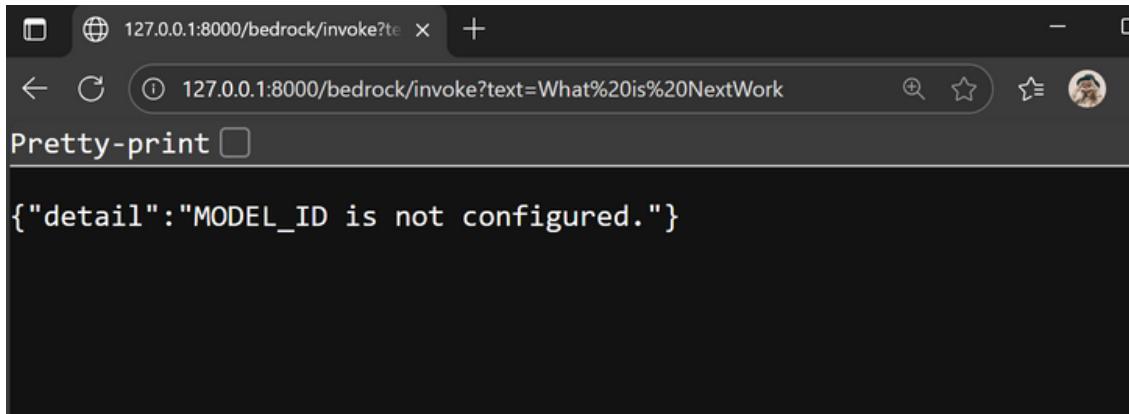


The screenshot shows a web browser window titled "NextWork Projects" at the URL "127.0.0.1:8000". The page displays "Kareshma's NextWork Projects" with a subtitle "Who is NextWork?". It features a sidebar with projects like "Automate Your Browser AI Agents" and "Transcribe Audio Files v AI". A central chat interface has a message from the AI stating: "NextWork appears to be an organization or platform that offers various projects and resources, including a community forum and learning projects such as AWS DevOps with Terraform. It is associated with a person named Kareshma Rajaananthapadmanaban." Below this, there is a message input field with placeholder "Type your message..." and a send button. A checkbox labeled "Talk to AI directly" is checked. The status bar at the bottom right says "Chat with my projects".

No response from chatbot when Checked Talk to AI directly

The terminal says **500 Internal Server Error**. Looks like there's an error in the API...

The **Talk to AI directly** checkbox controls whether the chatbot uses Retrieval-Augmented Generation (RAG) or not. When **checked**, the chatbot will **bypass the Knowledge Base** and talk directly to the AI model (Llama 3.3 70B Instruct). In this mode, the chatbot will answer questions based on its general knowledge, without using your specific documents.



Model ID error

```
INFO:    127.0.0.1:55010 - "GET /bedrock/invoke?text=What%20is%20NextWork HTTP/1.1" 500
Internal Server Error
INFO: Shutting down
INFO: Waiting for application shutdown.
INFO: Application shutdown complete.
INFO: Finished server process [1460]
INFO: Stopping reloader process [16392]
(vENV) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> notepad .env
(vENV) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> cat .env
AWS_REGION = us-east-2
KNOWLEDGE_BASE_ID = VERPAFAFRL
MODEL_ARN= arn:aws:bedrock:us-east-2::foundation-model/meta.llama3-3-70b-instruct-v1:0
MODEL_ID = meta.llama3-3-70b-instruct-v1:0
(vENV) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> python -m uvicorn web_app
:app --reload
```

Adding Model ID to .env file

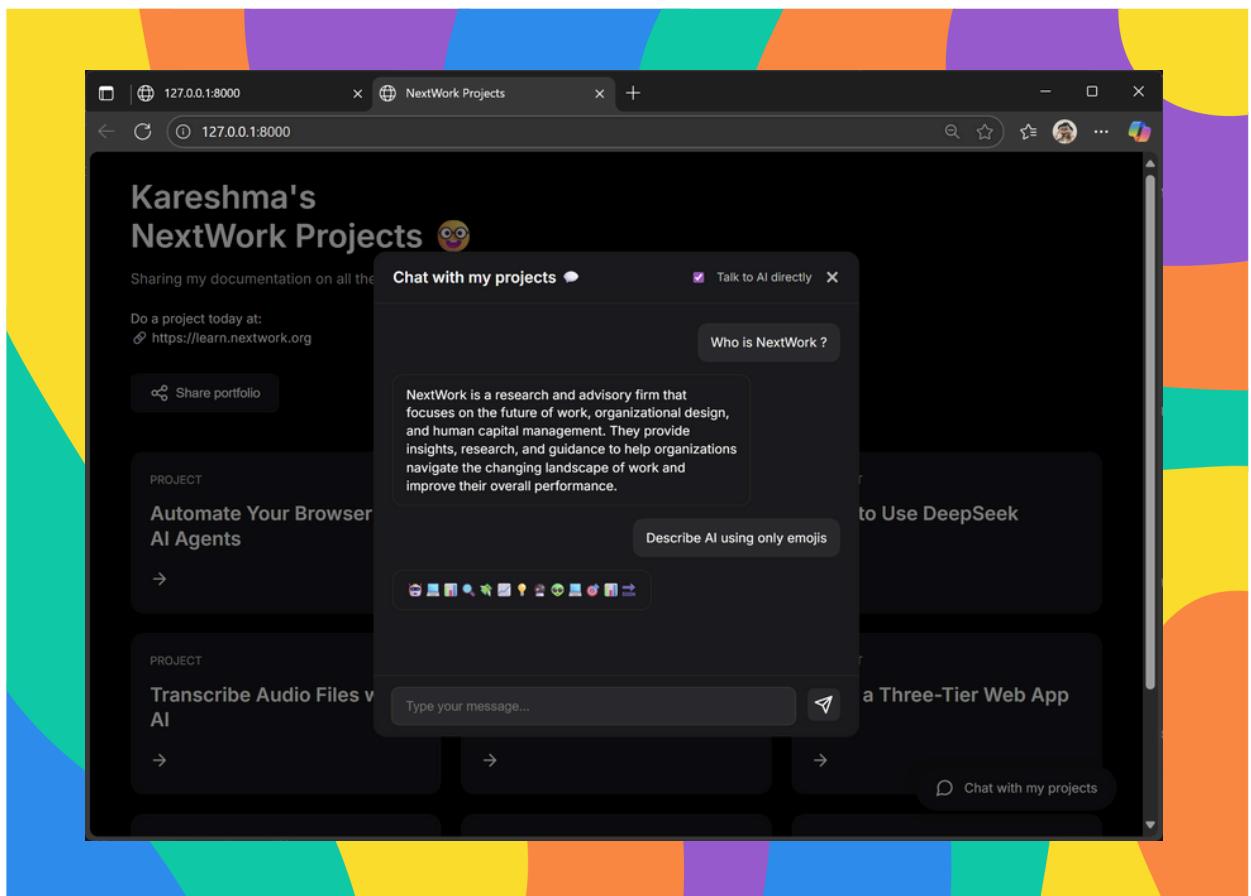
Error {"detail": "Model ID is missing"}

This is because **MODEL_ID haven't configure** in environment variable.

MODEL_ID is the **ID of the AI model** to call it directly. It's a simpler ID that just names the model, without any region or account information. That's why it's called an ID, not a Model ARN!

**Kareshma
Rajaananthapadmanaban**

nextwork.org



Web app showing response AI directly

Breaking down the Web App

In this secret mission, I'm going to investigate the full web app's code, understand how its frontend, backend, and data layers work together, and explore how chat interactions are handled. I'll also compare **running the API alone (main.py)** versus **running the complete web app (web_app.py)** to see how everything integrates into a functional chatbot interface.

When a user sends a chat message, the frontend JavaScript captures the text and sends it as a **GET request** to the [/bedrock/query endpoint](#). The backend (web_app.py) receives the request, queries the Knowledge Base via Bedrock, and generates a response. That response is sent back to the frontend, where the JavaScript displays both the **user's message** and the **chatbot's reply** in the **chat window**.

web_app.py extends the API by adding a **full web interface** on top of the core backend in main.py. It serves **HTML pages, CSS, and JavaScript** via **FastAPI**, uses **Jinja2 templates** to dynamically render pages, configures logging for debugging, and introduces a new endpoint ([/bedrock/invoke](#)) that lets the **AI model respond directly**, bypassing the Knowledge Base.

📁	.git	10/2/2025 12:00 PM	File folder
📁	__pycache__	10/2/2025 5:33 PM	File folder
📁	static	10/2/2025 12:00 PM	File folder
📁	templates	10/2/2025 12:00 PM	File folder
📁	venv	10/2/2025 12:03 PM	File folder
📄	.env	10/2/2025 6:40 PM	ENV File
📄	main.py	10/2/2025 12:00 PM	Python.File
📄	requirements	10/2/2025 12:00 PM	Text Document
📄	web_app.py	10/2/2025 12:00 PM	Python.File

All files to run a full web app for RAG chatbot

The screenshot shows a code editor window with the file 'index.html' open. The code is written in HTML and includes CSS and SVG elements. The editor has a dark theme with syntax highlighting. The status bar at the bottom shows 'Ln 1, Col 1' and other file-related information.

```
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>NextWork Projects</title>
    <link rel="stylesheet" href="/static/style.css">
    <link href="https://fonts.googleapis.com/css2?family=Inter:wght@400;500;600&display=swap" rel="stylesheet">
</head>
<body>
    <div class="container">
        <header>
            <h1>Kareshma's <br> NextWork Projects <img alt="NextWork logo" style="vertical-align: middle;"></h1>
            <p class="description">Sharing my documentation on all the amazing NextWork projects!</p>
            <p class="cta">
                Do a project today at:<br>
                <span class="link"><a href="https://learn.nextwork.org">https://learn.nextwork.org</a></span>
            </p>
            <button class="share-btn">
                <svg width="24" height="24" viewBox="0 0 24 24" fill="none" xmlns="http://www.w3.org/2000/svg">
                    <path d="M8.68439 10.6578L15.3124 7.34375M15.3156 16.6578L8.69379 13.3469M21 6C21 7.65685 19.6569 9.65685 19.6569 16C19.65685 13.3469 16 10.6578 12 7.34375Z" stroke="currentColor" stroke-width="2" stroke-linecap="round" stroke-linejoin="round"></path>
                </svg>
                Share portfolio
            </button>
        </header>
        <div class="projects-grid">
            <div class="project-card">
                <div class="project-label">PROJECT</div>
                <h2>Automate Your Browser with AI Agents</h2>
                <div class="arrow-icon">
                    <svg width="24" height="24" viewBox="0 0 24 24" fill="none" xmlns="http://www.w3.org/2000/svg">
                        <path d="M5 12H19M19 12L12 5M19 12L12 19" stroke="currentColor" stroke-width="2" stroke-linecap="round" stroke-linejoin="round"></path>
                    </svg>
                </div>
            </div>
        </div>
    </div>
</body>
```

Index.html file of Web app

Customizing the Frontend

In this secret mission, I'm going to give my chatbot web app a fresh look by **customizing the frontend**. The goal is to make the interface feel more personal and engaging, instead of just sticking with the default design. I'll either generate a **brand-new HTML** file or tweak the existing index.html, then integrate it into my project and test the updated design in action.

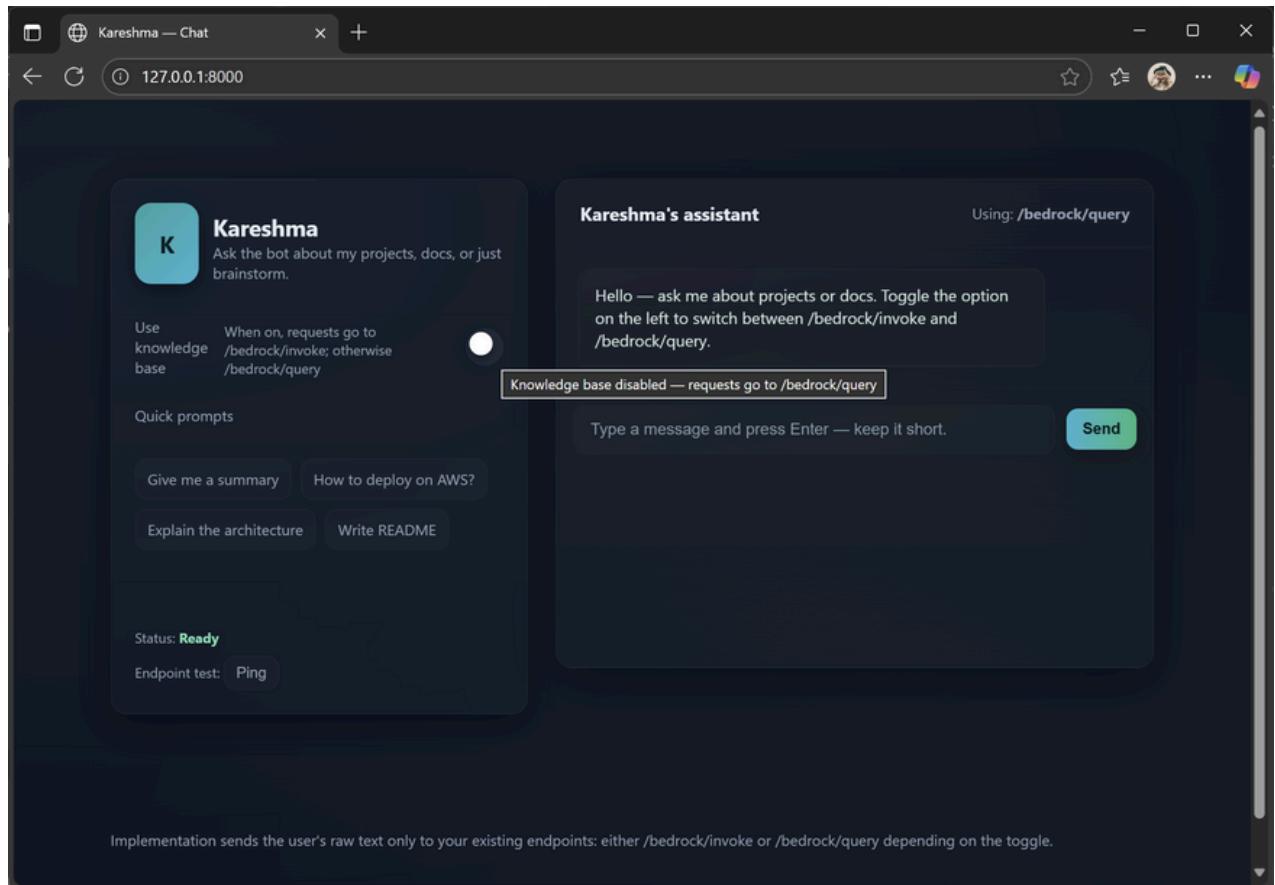
I want to experiment with customizing the frontend because it gives me the freedom to design a chatbot interface that feels more engaging and unique, instead of being stuck with a default look. This is possible because as long as the **API endpoints stay the same**, the **backend doesn't care** how the **frontend looks** it just sends and receives data. That separation lets me play with design while keeping the functionality intact.

My customized interface now features a clean, modern dark theme with two distinct panels. On the left, there's a profile-style card showing my name, a **toggle** to switch between `/bedrock/invite` and `/bedrock/query`, quick prompt buttons, and a status section. On the right, the chat assistant panel feels more polished, with a neatly styled conversation area, an input box, and a bright send button. It looks like a professional dashboard rather than a plain chat box.

```
(venv) PS D:\Projects\Kareshma\Nxtwrkcode\nextwork-rag-webapp> python -m uvicorn web_app:app --reload
INFO: Will watch for changes in these directories: ['D:\\Projects\\Kareshma\\Nxtwrkcode\\nextwork-rag-webapp']
INFO: Uvicorn running on http://127.0.0.1:8000 (Press CTRL+C to quit)
INFO: Started reloader process [6084] using StatReload
INFO:botocore.credentials:Found credentials in shared credentials file: ~/.aws/credentials
INFO: Started server process [1764]
INFO: Waiting for application startup.
INFO: Application startup complete.
INFO: 127.0.0.1:62236 - "GET / HTTP/1.1" 200 OK
INFO: 127.0.0.1:62237 - "GET /bedrock/invite?text=Who%20is%20NextWork%20%3F HTTP/1.1" 200 OK
INFO: 127.0.0.1:61233 - "GET /bedrock/invite?text=Describe%20AI%20using%20only%20emojis HTTP/1.1" 200 OK
INFO: 127.0.0.1:53597 - "GET /bedrock/invite?text=What%20would%20Shakespeare%20say%20about%20cloud%20outages HTTP/1.1" 200 OK
INFO: 127.0.0.1:55884 - "GET /bedrock/invite?text=What%20is%20NextWork HTTP/1.1" 200 OK
INFO: 127.0.0.1:55885 - "GET / HTTP/1.1" 200 OK
INFO: 127.0.0.1:57922 - "GET /bedrock/query?text=ping HTTP/1.1" 200 OK
INFO: 127.0.0.1:57922 - "GET /bedrock/query?text=ping HTTP/1.1" 200 OK
INFO: 127.0.0.1:61944 - "GET /bedrock/query?text=What%20is%20the%20service%20used%20in%20Threat%20detection%20on%20GuardDuty%20%3F HTTP/1.1" 200 OK
INFO: 127.0.0.1:64437 - "GET /bedrock/query?text=What%20are%20the%20projects%20I%20have%20done%20%3F HTTP/1.1" 200 OK
INFO: 127.0.0.1:64438 - "GET / HTTP/1.1" 200 OK
INFO: 127.0.0.1:52609 - "GET /bedrock/invite?text=What%20are%20the%20key%20differences%20between%20AWS%2C%20Azure%2C%20and%20Google%20Cloud%3F HTTP/1.1" 200 OK
INFO: 127.0.0.1:54258 - "GET / HTTP/1.1" 200 OK
```

Kareshma
Rajaananthapadmanaban

nextwork.org



Customize Web app Query option - Enabled RAG



Kareshma — Chat

127.0.0.1:8000

Kareshma

Ask the bot about my projects, docs, or just brainstorm.

Use knowledge base When on, requests go to /bedrock/invoke; otherwise /bedrock/query

Quick prompts

Give me a summary How to deploy on AWS?

Explain the architecture Write README

Status: Ready

Endpoint test: Ping

Kareshma's assistant

Using: /bedrock/invoke

Hello — ask me about projects or docs. Toggle the option on the left to switch between /bedrock/invoke and /bedrock/query.

Knowledge base disabled — requests go to /bedrock/query

Type a message and press Enter — keep it short.

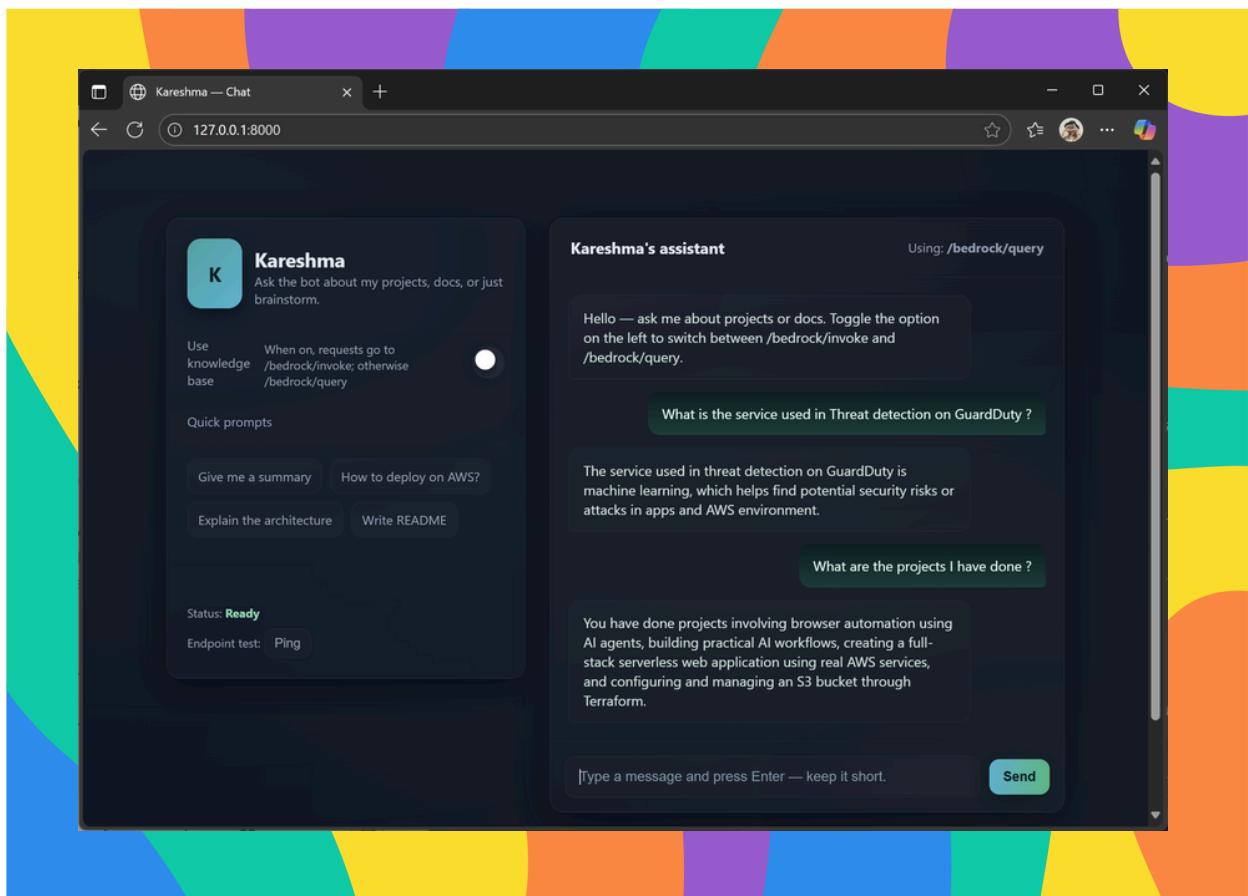
Send

Implementation sends the user's raw text only to your existing endpoints: either /bedrock/invoke or /bedrock/query depending on the toggle.

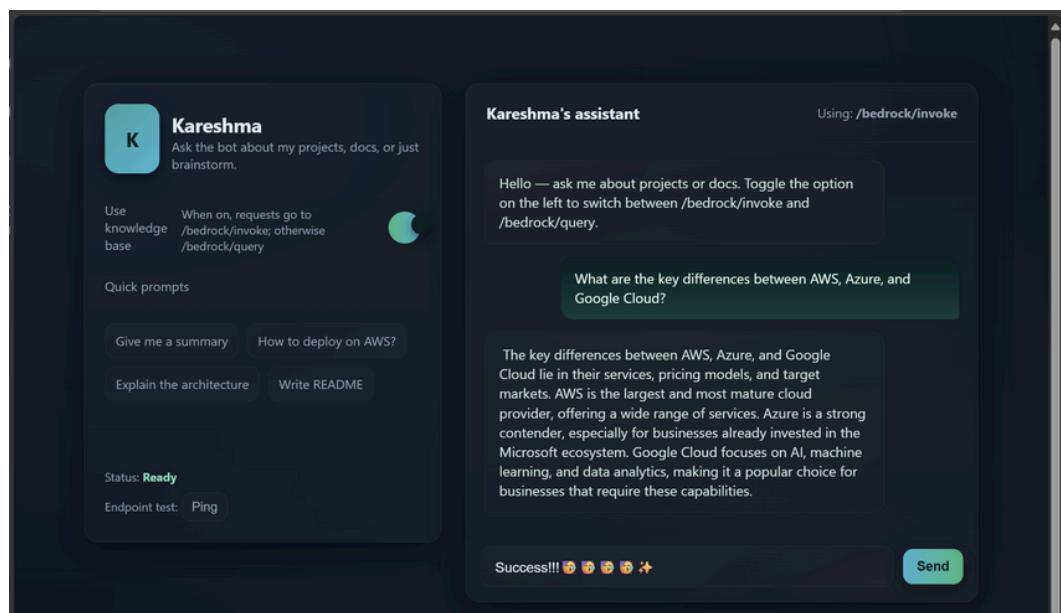
Customize Web app Invoke - Talk to AI directly

Kareshma
Rajaananthapadmanaban

nextwork.org



Query response



Invoke response

Delete the Resources

Now that we've built and tested our chatbot, it's time to clean up the AWS resources we created. Cleaning up ensures you don't get charged for unused resources and keeps your account tidy.

Resources to delete:

- Bedrock Knowledge Base
- OpenSearch Serverless Collection (vector store)
- IAM Access Key
- S3 Bucket
- Virtual Environment and cloned repository

1) Delete the Knowledge Base

1. Open the Bedrock Console.
2. From the left menu, choose Knowledge Bases.
3. Select your Knowledge Base.
4. Choose Delete.
5. Type delete in the confirmation field and confirm.

2) Delete the Vector Store (OpenSearch)

💡 What's a vector store?

A vector store organizes data by meaning, not just keywords. For example, searching “quick breakfast ideas” might also match “5-minute morning meals.”

Delete the Resources

- Open the OpenSearch Console.
- From the left menu, choose Collections.
- Select your vector store.
- Choose Delete.
- Type confirm in the confirmation field and confirm.

3) Delete IAM Access Key

1. Open the IAM Console.
2. From the left menu, select Users.
3. Select your IAM user → Security credentials tab.
4. In Access keys, find your key.
5. Select Actions → Delete.
6. Deactivate it, then type the access key ID to confirm deletion.

💡 Why delete access keys?

They provide programmatic access to your AWS account. Removing them prevents unauthorized access.



Delete the Resources

4) Delete the S3 Bucket

1. Open the S3 Console.
2. Select your bucket (e.g., bedrock-kb-data-<your_initials>).
3. Choose Empty → type permanently delete → confirm.
4. Once empty, select the bucket → Delete.
5. Type the bucket name (e.g., nextwork-rag-documentation-<your_initials>) to confirm.

5) Deactivate Virtual Environment

1. Go back to your terminal.
2. Stop the server with Ctrl+C.
3. Deactivate your environment:

CMD → deactivate

You should see the (venv) prompt disappear.

💡 Why deactivate?

It ensures clean separation for future projects and avoids accidental dependency usage.



Kareshma
Rajaananthapadmanaban

*Follow up for more
interesting projects !!!*



Check out nextwork.org / my profile for more projects