

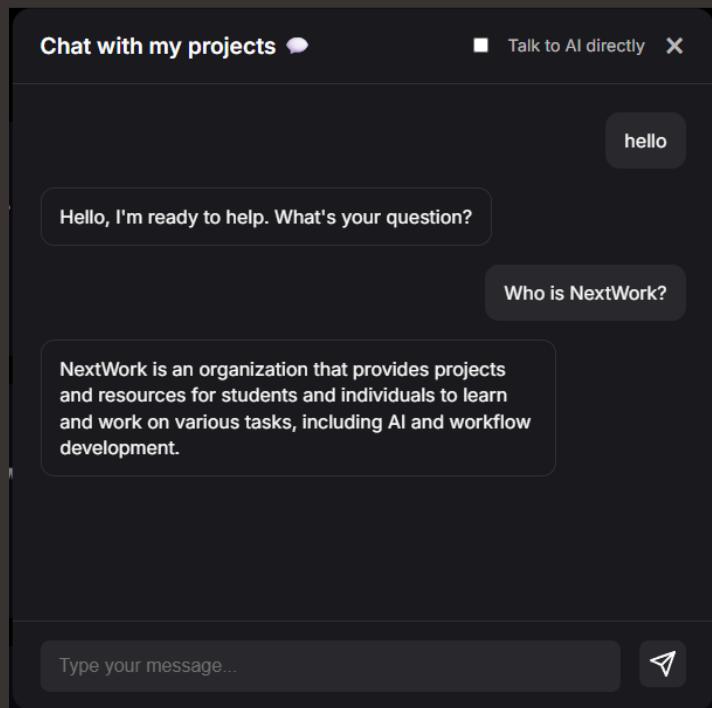


nextwork.org

Building a RAG Chatbot with a Web Interface

ME

Melvin J Bonner



Introducing Today's Project!

In this project, I will build a fully functional web app where users can chat with my chatbot and get answers based on my documents. I'm doing this project to gain hands-on experience building complete web applications that integrate AI capabilities.

Tools and concepts

Services I used were Amazon Bedrock, S3, Amazon OpenSearch Service, and IAM. Key concepts I learned included Knowledge Base, RAG chatbot, Query Endpoint, and vector store.

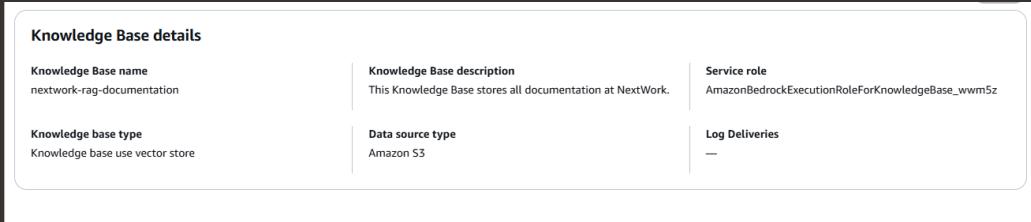
Project reflection

This project took me approximately 3 hours. the most challenging part was troubleshooting the parameter failed error. It was the most rewarding to have a fully functioning chatbot.

I wanted to gain more experience with AI and a better understanding of Amazon Bedrock. This project met my goals.

Chatbot setup

I created a Knowledge Base to create the knowledge management system for my chatbot. This is important because when I ask my chatbot a question, the Knowledge Base is responsible for quickly finding the information that will answer my question.



Setting Up the API

To run the API, I had to install requirements like fastapi, uvicorn, boto3, and python-dotenv. These packages are important because they are the required packages needed for Python to run the API.

A virtual environment helps me by telling my computer to start using the virtual environment I created to allow all the Python packages I installed to only exist in this environment. A virtual environment will not allow my project to interfere with other Python projects on my computer.

```
(venv) PS C:\Users\melvi\nextwork-rag-webapp> pip3 list
Package           Version
------------------ -----
annotated-types   0.7.0
anyio             4.8.0
boto3              1.36.20
botocore           1.36.20
cffi               0.8.9
colorama           0.4.6
fastapi            0.115.8
h11                0.14.0
idna               3.10.0
Jinja2              3.1.5
jmespath            1.0.1
MarkupSafe          3.0.2
pdf                25.0.1
pydantic            2.10.6
pydantic_core       2.27.2
python-dateutil     2.9.0.post0
python-dotenv        1.0.1
s3transfer           0.17.2
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 C:\Users\melvi\nextwork-rag-webapp> |
```

Breaking down the API

In our web app, we need an API because the API will let our web app to talk to the Amazon Bedrock chatbot.

Digging deeper into the API code, the main tools I need to import are FastAPI, AWS SDK, environment variables, and load. These tools help by giving the tools I need for the API to work. Instead of coding everything from scratch, I'm grabbing pre-built tools. The tools make it way faster than making every piece myself.

The environment variables we need are AWS_REGION, KNOWLEDGE_BASE_ID, and MODEL_ARN. I store them separately because putting these variables directly in the code is a security risk. If you accidentally push your code to GitHub with AWS credentials or Knowledge Base ID in it, anyone could use them. Environment variables keep sensitive stuff separate from your code.

My API uses the Python AWS SDK to talk to Bedrock. The difference between the AWS CLI and SDK is you use the CLI when you need to do something quickly yourself in the terminal, and you use an SDK when you're building an application that needs to work with AWS.



ME

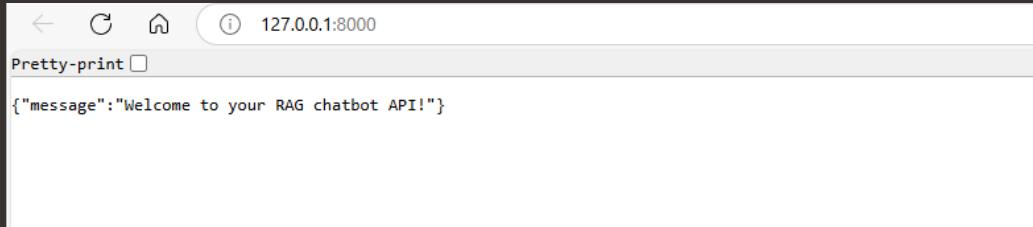
Melvin J Bonner
NextWork Student

nextwork.org

Our API has two main routes: / is the root route, which is like the default route for my API. If the user doesn't specify a route, my API will use this route. /bedrock/query: is a route that lets the user query my chatbot. In general, an API route is like a path to a specific part or function in my API.

Testing the API

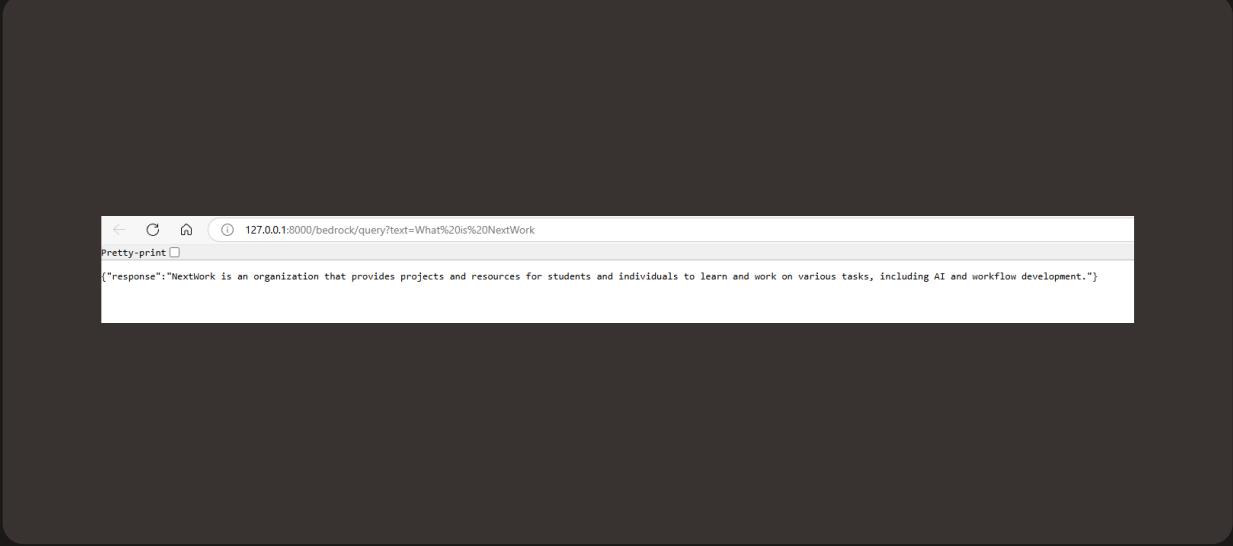
When I visited the root endpoint, I saw the message {"message": "Welcome to your RAG chatbot API"} in the browser. This confirms that the API is running.



Troubleshooting the API

My second endpoint had errors!

To fix the Parameter Failed error, I had to update the environment variables. Doing this helps to fix my error and get the API's query endpoint running.



A screenshot of a web browser window displaying a JSON response from an API. The URL in the address bar is 127.0.0.1:8000/bedrock/query?text=What%20is%20NextWork. Below the address bar, there is a "Pretty-print" checkbox. The main content area shows a single JSON object with a key "response" and a value describing NextWork as an organization that provides projects and resources for students and individuals to learn and work on various tasks, including AI and workflow development.

```
{"response": "NextWork is an organization that provides projects and resources for students and individuals to learn and work on various tasks, including AI and workflow development."}
```

Running the Web App

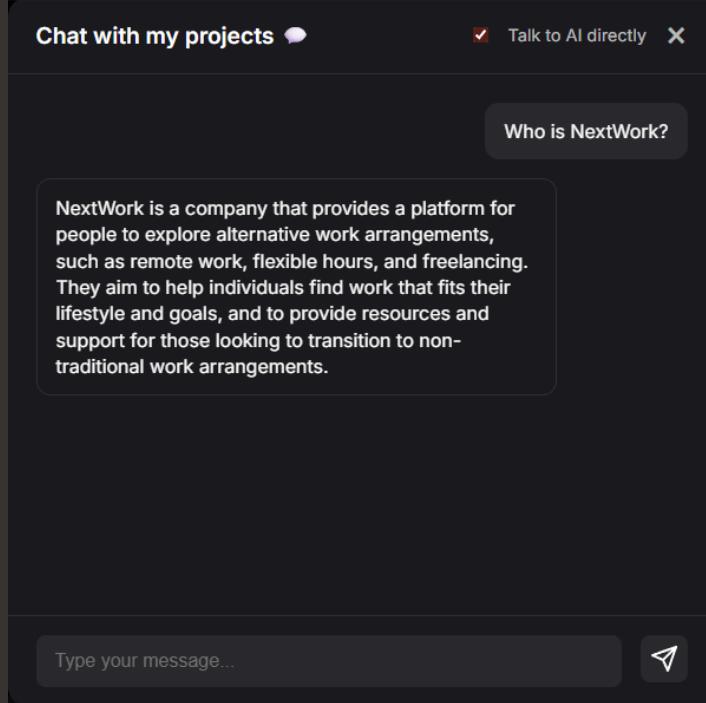
The difference between the API and the web app for the user is how they react inside the app after receiving a question. The web app sends your question to the API endpoint, the API receives the question, and uses the Bedrock command to query my Knowledge Base with the question, the api send Bedrock's response back to the web app, and the web app displays the response.

The web app has the ability to switch between my knowledge base and general knowledge. When RAG is disabled, I ask the same question and get a lot less accurate response, because the chatbot is now using the AI model directly, without using my Knowledge Base.

ME

Melvin J Bonner
NextWork Student

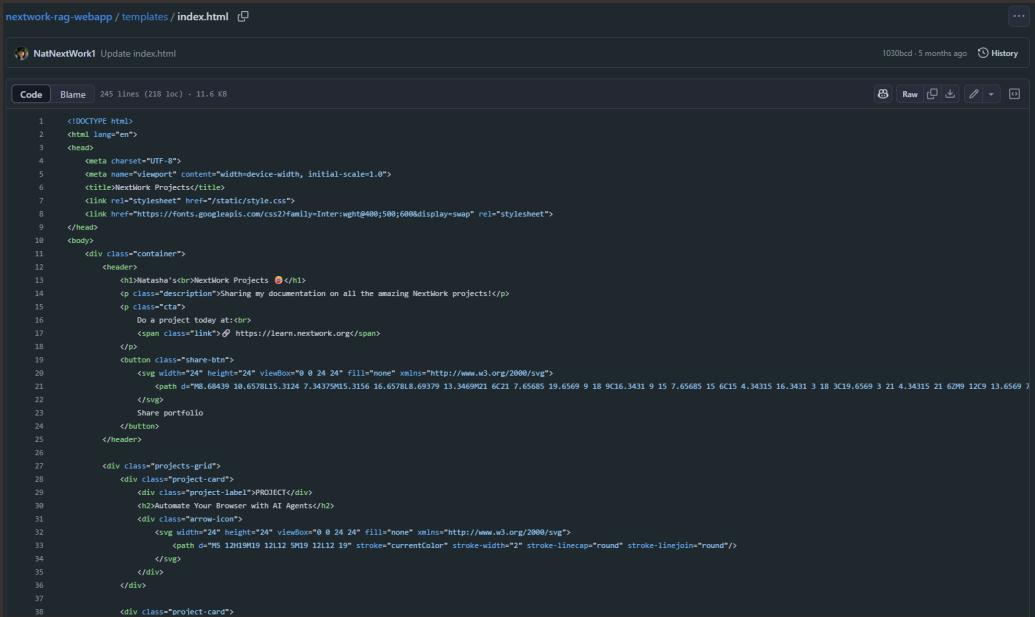
nextwork.org



Breaking down the Web App

When a user sends a chat message JavaScript makes an API call to /bedrock/query, the API (powered by web_app.py) processes your message using Bedrock, the response comes back and JavaScript displays it in the chat window.

web_app.py extends the API by including features needed to serve a web interface.



The screenshot shows a GitHub code editor interface for the file `index.html` in the `nextwork-rag-webapp / templates` directory. The code is as follows:

```
1  <!DOCTYPE html>
2  <html lang="en">
3  <head>
4      <meta charset="UTF-8">
5      <meta name="viewport" content="width=device-width, initial-scale=1.0">
6      <title>NextWork Projects</title>
7      <link rel="stylesheet" href="/static/style.css">
8      <link href="https://fonts.googleapis.com/css?family=Inter:wght@400;500;600&display=swap" rel="stylesheet">
9  </head>
10 <body>
11     <div class="container">
12         <header>
13             <h1>NextWork Projects <a href="#">Sharing my documentation on all the amazing NextWork projects!</a></h1>
14             <p>Sharing my documentation on all the amazing NextWork projects!</p>
15             <p>Do a project today at:<a href="https://learn.nextwork.org">https://learn.nextwork.org</a></p>
16             <button class="share-btn">
17                 <svg width="24" height="24" viewBox="0 0 24 24" fill="none" xmlns="http://www.w3.org/2000/svg">
18                     <path d="M8.68439 10.657815.3124 7.34375M15.3156 16.657818.69379 13.3469M21.621 7.65685 19.6569 9.18 9C16.3431 9.15 7.65685 19.6C15.4.34315 16.3431 3.18 3C19.6569 3.21 4.34315 21.62M9.12C9.13.6569 7.14.6569 5.15.6569 3.16.6569 1.17.6569 0.18.6569 0 0Z" style="stroke:#000; stroke-width:1px; stroke-linecap:round; stroke-linejoin:round;" />
19                 </svg>
20                 Share portfolio
21             </button>
22         </header>
23         <div class="projects-grid">
24             <div class="project-card">
25                 <div class="project-label">PROJECT</div>
26                 <div>Automate Your Browser with AI Agents</div>
27                 <div class="arrow-icon">
28                     <svg width="24" height="24" viewBox="0 0 24 24" fill="none" xmlns="http://www.w3.org/2000/svg">
29                         <path d="M5.129M19.12L12M5.129 19.12L12 19" stroke="currentColor" stroke-width="1px" stroke-linecap="round" stroke-linejoin="round;" />
30                     </svg>
31                 </div>
32             </div>
33         </div>
34     </div>
35 
```

ME

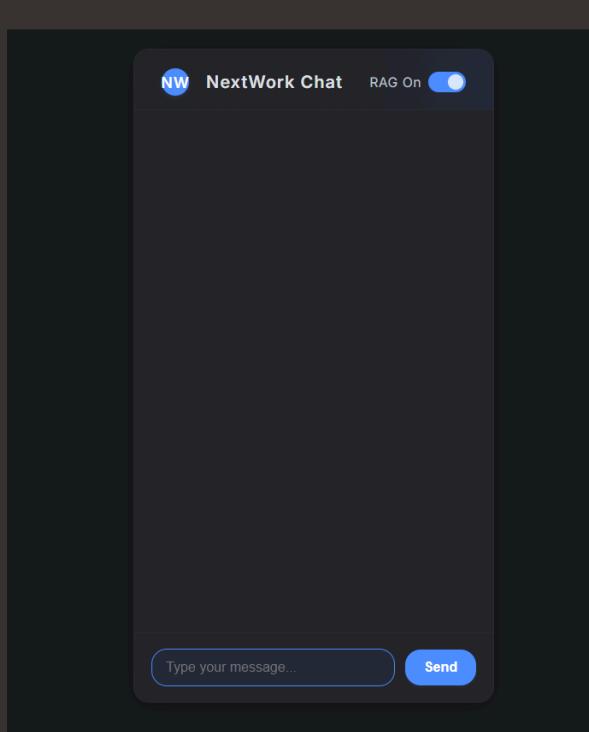
Melvin J Bonner
NextWork Student

nextwork.org

Customizing the Frontend

I want to experiment with customizing the frontend because I want to have an app that is unique and has a more professional look, tailored for sophisticated users. This is possible because as long as the API endpoints stay the same the application will continue to work.

My customized interface now features a standalone chatbot that has an uncluttered UI.





nextwork.org

The place to learn & showcase your skills

Check out nextwork.org for more projects

