**Link to the demo:** https://dc8eru1q29n2p.cloudfront.net

# 1. Domain

The chosen domain for this project is Education and personalised learning, more specifically, it is an AI-Assisted Code Understanding. Sometimes the code files I am looking at for open-source development are very complex to understand because of the use of libraries that I might not have used or seen before, or because of poor code writing practises. The objective of making this tool is to find or analyse the exact location in a file where a certain functionality exists, and answer code related questions.

# 2. Problem Statement

Modern software teams deal with large, complex codebases where understanding legacy code, onboarding new developers, or finding relevant functionality is time-consuming and error-prone. Traditional keyword-based search tools (like grep, or built-in IDE search) fail to capture the semantic meaning of code, leading to poor discovery, redundant work, and slow development cycles.

# 3. Business Case

Organizations waste hundreds of hours per developer annually on tasks like:

- Understanding unfamiliar code sections
- Searching for functions, APIs, or logic
- Asking teammates questions that could be answered from code.

This results in:

- Delayed feature delivery
- Frustrated engineers
- Higher onboarding time
- Increased bugs due to misunderstood logic.

The proposed system enables:

- Uploading a codebase
- Embedding it semantically using **CodeBERT**
- Creating an efficient vector index with **FAISS**
- Letting users ask **natural-language questions**
- Returning accurate, contextually relevant answers via model deployed on **SageMaker.**

# 4. Functional Requirements

- The system should allow users to upload either one zip file or a Github link to a codebase.
- The system should process the files so that when the user asks questions related to the content in the files, the system can answer the questions.
- The system should return answers that are relevant to the content of the uploaded code.
- The system should allow users to ask natural-language questions about their uploaded codebase.
- The received response from the system must be in English and should be visible on screen.
- A user should be able to ask multiple questions after they upload files or the Github link.
- The system should notify users if there is an error during upload, processing, or querying.

# 5. Non-Functional Requirements

**Performance:**

- The system should not take more than 10 minutes to answer user queries.
- The system should not take more than 10 minutes to process the uploaded files.
- The system should support more than one user concurrently.

**Scalability:**

- The system should be able to spawn more compute power (EC2 instances) if the demand spikes so that system doesn't become slow or throw error when it receives multiple requests at the same time.

**Security:**

- The system should restrict access to user-uploaded data such that only the owning user can query their code.
- All data transfers between client and server shall be encrypted.
- Back-end functionalities should only be accessible from the front-end and nowhere else (like terminal or other tools).
- Only the front-end / website / code should be able to access the back-end processing functionalities.

**Data Handling:**

- Data should be deleted after 30 days unless compliance to a government policy requires the storage of data for longer periods, or there is an explicit business requirement.

**Monitoring and Logging:**

- Requests to the Lambda functions and other AWS services must be monitored through CloudWatch.
- The system should expose operational metrics (e.g., job queue length, task duration, query latency) for observability.

**Availability and Reliability:**

- The system should be available most of the time to take user input.
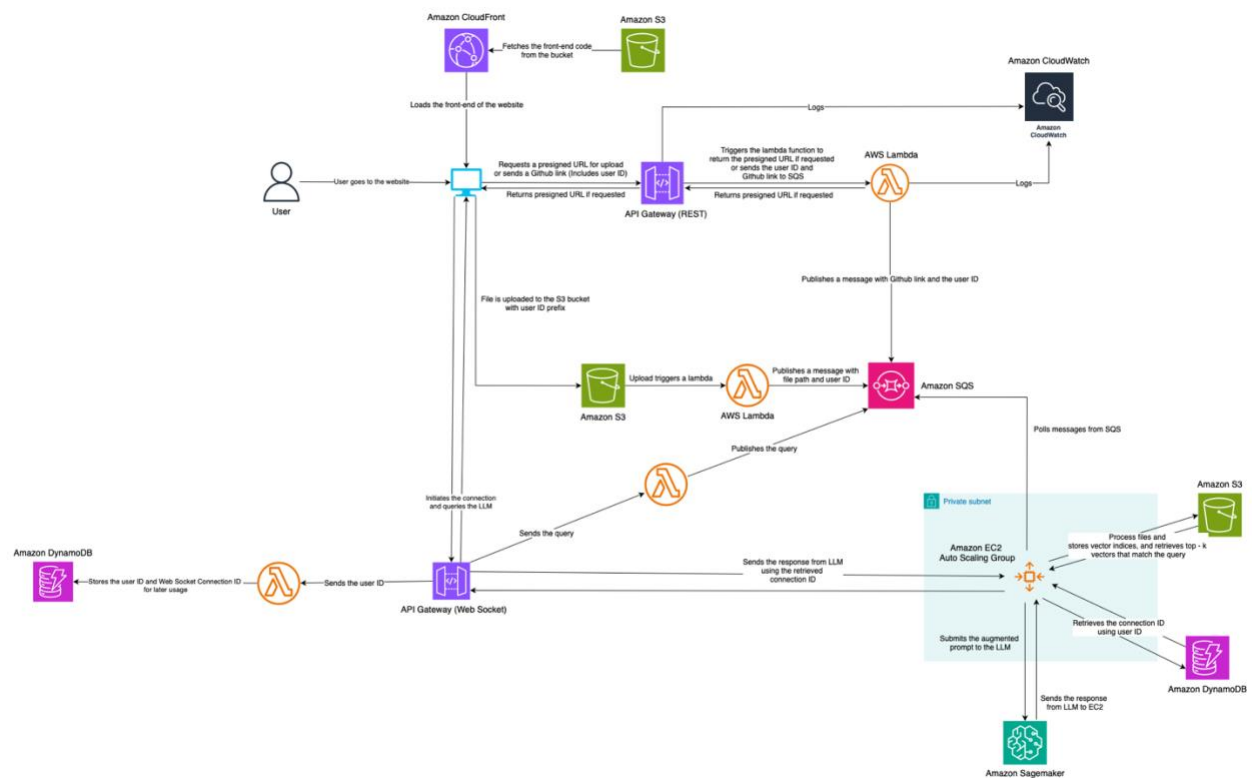
## 6. Architecture Diagram



**Fig. 6.1** Architecture Diagram

I. User first access the website through Amazon CloudFront.

II. When the user lands on the website, the browser initiates the connection to web socket so that the responses could be retrieved later on.

III. Once the connection is open, a unique user ID is generated and wrapped in a message. This message is sent to the web socket.

IV.   This causes a lambda function to trigger. The lambda function will then extract the user ID from the message, and the connection ID from its invocation metadata, and stores them in the DynamoDB.

V.   Now user either uploads the zip file of the code, or enters the corresponding Github URL, and sends a request to API Getway.

VI.   API Gateway triggers a lambda function that will process these inputs. If the user has uploaded zip file of the code, it will generate a presigned URL and send that URL to the front-end. If the user has uploaded Github URL, the URL along with the user ID will be passed to the Amazon SQS. Once the user uploads the file using the presigned URL, the file path will be passed to the same Amazon SQS.

VII.   EC2 instance managed by Elastic Beanstalk will poll the messages from the queue and start processing the uploaded files. It will embed the files in chunks using CodeBERT model and store the vectors in FAISS index, the metadata in another S3 bucket.

VIII.   Once the processing is finished, it will query the DynamoDB for the connection ID using the user ID to notify the front-end using the web socket.

IX.   The front-end will allow user to input query.

X.   Once the query is entered, it will be sent as a message along with the user ID to the web socket, which will trigger a Lambda function.

XI.   The lambda function will put this message in a queue. EC2 will poll the message and starts processing the query.

XII.   The query will be embedded using CodeBERT. Then the program will retrieve vectors stored in FAISS index and search for the top – k similar embeddings / vectors that match the query vector.

XIII.   Once the matches have been found, they will be augmented into the user prompt and passed to the AI model deployed (Code Llama / Meta Llama / CodeBERT) on Amazon Sagemaker.

XIV.   Amazon Sagemaker will make use of augmented context into the prompt and return a response to the EC2 instance.

XV.   EC2 instance will send the response in chunks to the front-end of the website through the web socket.

XVI.   The user will be able to see the responses as they are received from the back-end.

**Fig 6.2** UML Data Sequence diagram

## 7. AWS Services and Tech Stack

**Services:**

Amazon CloudFront
Amazon S3
API Gateway (REST and Web Socket)
AWS Lambda
Amazon SQS
Amazon DynamoDB
Amazon EC2 (Auto Scaling Groups)
Amazon Sagemaker
Amazon CloudWatch

**Programming Languages:**

JavaScript for AWS Lambda
Python for EC2 (back-end) code

**Additional Tools:**

CodeBERT from Hugging Face (for converting text data to vector)
CodeLLama from Meta (to get response for the augmented prompt)

## 8. Potential Architectural Challenges

- Maintaining long-lived web socket connection, especially when there are multiple concurrent connections, could be tricky because of connection timeouts and scaling issues with concurrent users.

- If the model used to create embeddings gets updated, it can break the semantic alignment. Hence, it's important to keep track of specific version of the model used to create embeddings.

- Users with malicious intent can compromise the system while querying the model. Input sanitization is necessary to avoid this.

- If chunks are too large or too many, program might hit the token limit. It is important to embed the code in smaller chunks.

## 9. Initial cost estimation for overall implementation

| Service | Estimate & Assumptions | Cost* |
|---|---|---|
| Amazon S3 (static hosting) | First 50 TB / month + per 1000 requests | $0.023 + $0.0054 |
| CloudFront | First 1 TB monthly transfer + requests | $0 |
| API Gateway (REST + WebSocket) | 1M REST calls & 1M WS messages | $3.50 + $1.00 |
| AWS Lambda x 3 | Price per milli seconds for 128MB memory functions + 1M requests | $0.0000000021 + $0.20 (x 3) |
| Amazon SQS | 1 M requests + First 10 TB / Month data out | $0.50 + $0.09 |
| Amazon DynamoDB | 1M WRU/RRU + per GB per month (First 25 GB is free) | $0.625 + $0.125 + $0.25 |
| Amazon EC2 | C6i.2xlarge on demand hourly with 15 GB storage | $0.34 |
| Amazon SageMaker | ml.g5.2xlarge instance for model deployment | $1.212 |
| Amazon S3 (vector store) | First 50 TB / month + per 1000 requests | $0.023 + $0.0054 |
| Data transfer between services | (Mostly free within same AZ, negligible) | $0 |
| | | |
| **Total** | | **$8.33** |

*Costs have been retrieved from [AWS's official pricing pages](). Free tier for services is ignored.*

## 10. Explanation

- **Front-end hosted on S3 Bucket**

  The bucket that hosts the front-end code is inaccessible to the outside world. Only the CloudFront distribution can access the bucket contents due to the policy set on the S3 Bucket. The S3 bucket has versioning enabled so that one can always go to the previous version in case of an anomaly.

  

- **API Gateway**

  There is a REST API and a Websocket API configured on API Gateway. While it doesn't require authorization to access the API at this point, the URLs for the APIs are stored along with the front-end code in the S3 bucket. This allows the front-end to fetch the API during runtime so the APIs are hidden. Bothe the APIs invoke lambda integrations.

- **Lambda functions**

  There are 4 lambda functions in this architecture. Each has been assigned specific IAM role so that they only have enough access to show logs on CloudWatch, access dynamo DB to store WebSocket connection information, integrate with S3 to generate presigned URLs, and publish messages to the queue (SQS). They have all been assigned environment variables to avoid hardcoding different AWS services URLs. They all run on ARM architecture to save costs.

## Functions (4)

🔍 Filter by attributes or search by keyword

| | Function name | Description | Package type | Runtime | Last modified |
|---|---|---|---|---|---|
| ☐ | sendQuery | - | Zip | Node.js 22.x | 1 day ago |
| ☐ | storeConnection | - | Zip | Node.js 22.x | 1 day ago |
| ☐ | returnPresignedURL_or_uploadLink | - | Zip | Node.js 22.x | 1 day ago |
| ☐ | SendFilePathToQueue | - | Zip | Node.js 22.x | 1 day ago |

## Environment variables (2)

The environment variables below are encrypted at rest with the default Lambda service key.

🔍 Find environment variables

| Key | Value |
|---|---|
| AWS_LAMBDA_EXEC_WRAPPER | /opt/otel-instrument |
| QUEUE_URL | https://sqs.us-east-2.amazonaws.com/257394469414/MessageQueue.fifo |

## Environment variables (2)

The environment variables below are encrypted at rest with the default Lambda service key.

🔍 Find environment variables

| Key | Value |
|---|---|
| AWS_LAMBDA_EXEC_WRAPPER | /opt/otel-instrument |
| TABLE | UserSocketConnections |

## Environment variables (3)

The environment variables below are encrypted at rest with the default Lambda service key.

🔍 Find environment variables

| Key | Value |
|---|---|
| AWS_LAMBDA_EXEC_WRAPPER | /opt/otel-instrument |
| QUEUE_URL | https://sqs.us-east-2.amazonaws.com/257394469414/MessageQueue.fifo |
| S3_BUCKET_NAME | 5411-copilot-project-v2 |

## Environment variables (2)

The environment variables below are encrypted at rest with the default Lambda service key.

🔍 Find environment variables

| Key | Value |
|---|---|
| AWS_LAMBDA_EXEC_WRAPPER | /opt/otel-instrument |
| QUEUE_URL | https://sqs.us-east-2.amazonaws.com/257394469414/MessageQueue.fifo |

- **User upload store bucket**

  This is the bucket where user uploaded files are stored. Each uploaded file is prefixed with a unique userId for efficient retrieval. A lifecycle management rule here deletes all the files after 1 day because user wouldn't need them after a session is over. This bucket is also inaccessible to the outside world and only resources with appropriate role can access this bucket (EC2 instances in this case).

- **Dynamo DB**

UserSocketConnection table stores the user to connection mapping. There is no resource based policy attached to it.



- **SQS**

This queue stores messages which has URL of files in S3 or GitHub link, and user query. It has a resource-based policy attached to it so that only resources (lambda functions) with appropriate role can access it.

- **S3 vector store bucket**

This bucket stores the vectors computed for files uploaded by users. It also has a lifecycle rule that deletes all the files after one day, and it's public access is also blocked by default.



- **Amazon EC2**

The Auto Scaling Group launches instances in a private subnet. Each instance runs a user data script that retrieves the back-end code from an S3 bucket, installs the necessary dependencies, and starts the Python application. These EC2 instances do not have public IPv4 access. The Auto

Scaling Group monitors the SQS queue, and if the number of messages exceeds five, it provisions additional instances. The attached security group restricts traffic to outbound-only access. C6i instances are used to support multiple worker threads and handle transformer models efficiently. However, instance size may be reduced in the future based on CPU and memory usage analysis.

EC2 > Auto Scaling groups > 5411-midterm

## 5411-midterm

### 5411-midterm Capacity overview

Edit

arn:aws:autoscaling:us-east-2:257394469414:autoScalingGroup:2a54e8bf-a910-4508-8aa8-a3a2964f02cd:autoScalingGroupName/5411-midterm

| Desired capacity | Scaling limits (Min - Max) | Desired capacity type | Status |
|---|---|---|---|
| 1 | 1 - 5 | Units (number of instances) | ⊖ Updating capacity |

**Date created**
Fri May 30 2025 16:17:51 GMT-0300 (Atlantic Daylight Time)

**Details** | Integrations - *new* | Automatic scaling | Instance management | Instance refresh | Activity | Monitoring

### Launch template

Edit

| Launch template | AMI ID | Instance type | Owner |
|---|---|---|---|
| lt-0f4e2ff19f7dea8f1 | ami-06c8f2ec674c67112 | c6i.2xlarge | arn:aws:iam::257394469414:root |
| EC2-5411-MIDTERM | | | |

| Version | Security groups | Security group IDs | Create time |
|---|---|---|---|
| Default | - | sg-09ecef8323913c057 | Fri May 30 2025 16:02:41 GMT-0300 (Atlantic Daylight Time) |

| Description | Storage (volumes) | Key pair name | Request Spot Instances |
|---|---|---|---|
| 1.0 | /dev/xvda | first | No |

View details in the launch template console

### Network

Edit

| Availability Zones | Subnet ID | Availability Zone distribution |
|---|---|---|
| us-east-2c | subnet-03e9328257f0ce1b9 | Balanced best effort |

---

## sg-09ecef8323913c057 - allow ec2 access

SecurityGroup

Actions ▼

### Details

| Security group name | Security group ID | Description | VPC ID |
|---|---|---|---|
| allow ec2 access | sg-09ecef8323913c057 | SSH, HTTP and HTTPS | vpc-0c9b2203850ed1067 |

| Owner | Inbound rules count | Outbound rules count | |
|---|---|---|---|
| 257394469414 | 0 Permission entries | 1 Permission entry | |

**Inbound rules** | Outbound rules | Sharing - *new* | VPC associations - *new* | Tags

### Inbound rules

Manage tags    Edit inbound rules

| Name ▽ | Security group rule ID ▽ | IP version ▽ | Type ▽ | Protocol ▽ | Port range ▽ | Source |
|---|---|---|---|---|---|---|

No security group rules found

---

## sg-09ecef8323913c057 - allow ec2 access

Actions ▼

### Details

| Security group name | Security group ID | Description | VPC ID |
|---|---|---|---|
| allow SecurityGroup | sg-09ecef8323913c057 | SSH, HTTP and HTTPS | vpc-0c9b2203850ed1067 |

| Owner | Inbound rules count | Outbound rules count | |
|---|---|---|---|
| 257394469414 | 0 Permission entries | 1 Permission entry | |

Inbound rules | **Outbound rules** | Sharing - *new* | VPC associations - *new* | Tags

### Outbound rules (1)

Manage tags    Edit outbound rules

| Name ▽ | Security group rule ID ▽ | IP version ▽ | Type ▽ | Protocol ▽ | Port range ▽ | Destination |
|---|---|---|---|---|---|---|
| – | sgr-0b535b020bd65fe1a | IPv4 | All traffic | All | All | 0.0.0.0/0 |

- **Amazon SageMaker AI**

The SageMaker domain's network is also deployed in the private subnet where the security rules doesn't allow inbound access from the network. The CodeLlama model is deployed on a default endpoint. The endpoint is hosted on ml.g5.2xlarge instance because of the CodeLlama model's compatibility.

- **Amazon CloudWatch Logs**

All Lambda functions and the API Gateway generate logs for each request. Alarms configured for SQS trigger EC2 auto scaling when thresholds are met. The monitoring dashboard displays key metrics, including the number of messages in the queue, EC2 instance CPU utilization, and the number of requests received by the API Gateway.

If the back end becomes inactive for any reason and messages begin to accumulate, an SNS topic sends email notifications to alert me. Application traces provide visibility into the availability of Lambda functions.

## Panel 1

**ConsumerUnresponsive**
Metric alarm
⊘ OK

**SQS few messages**
Metric alarm
⚠ In alarm

**SQS Overcrowd**
Metric alarm
⊘ OK

0

| 05/26 | 05/26 | 05/27 | 05/27 | 05/28 | 05/28 | 05/29 | 05/29 | 05/30 | 05/30 | 05/31 | 05/31 | 06/01 | 06/01 |

● ApproximateNumberOfMessagesVisible

Click timeline to see the state change at the selected time.

| 5/26 | 5/26 | 5/27 | 5/27 | 5/28 | 5/28 | 5/29 | 5/29 | 5/30 | 5/30 | 5/31 | 6/1 | 6/1 |

■ In alarm   ■ OK   ■ Insufficient data   ■ Disabled actions

[CloudWatch console feature]

Details | Tags | **Actions** | History | Parent alarms

### Actions ⊘ Actions enabled

| Type ▽ | Description ▽ | Config ▽ |
|---|---|---|
| Notification | When in alarm, send message to topic "Consumer_Unresponsive" | - |

## Panel 2

**ConsumerUnresponsive**
Metric alarm
⊘ OK

**SQS few messages**
Metric alarm
⚠ In alarm

**SQS Overcrowd**
Metric alarm
⊘ OK

0.059

| 05/22 | 05/23 | 05/24 | 05/25 | 05/26 | 05/27 | 05/28 | 05/29 | 05/30 | 05/31 | 06/01 |

● ApproximateNumberOfMessagesVisible

Click timeline to see the state change at the selected time.

| 5/22 | 5/23 | 5/24 | 5/25 | 5/26 | 5/27 | 5/28 | 5/29 | 5/30 | 5/31 | 6/1 |

■ In alarm   ■ OK   ■ Insufficient data   ■ Disabled actions

Details | Tags | **Actions** | History | Parent alarms

### Actions ⊘ Actions enabled

| Type ▽ | Description ▽ | Config ▽ |
|---|---|---|
| Auto Scaling | When in alarm, use policy ScaleInInstances (Remove 1 instance) for the group 5411-midterm | - |

[CloudWatch console feature]

## Panel 3

# Services  Info

30m | 1h | **3h** | 12h | Custom ▦        Local timezone ▼    ⟳ ▼

### Services by SLI status

■ Healthy (0)   ■ Unhealthy (0)
■ No SLO (4)   ■ Insufficient data (0)

### Top services by fault rate

| Service | Fault rate |
|---|---|
| No services with faults | |

### Top dependency paths by fault rate

| Remote service | Service | Fault rate |
|---|---|---|
| No dependencies with faults | | |

### Services (4) Info

⟳   Actions ▼   **Enable Application Signals**

🔍 Filter services and resources by text, property or value

‹ 1 › ⚙

| | Name ▽ | SLI status ▽ | Service Availa... ▽ | Application ▽ | Hosted in |
|---|---|---|---|---|---|
| ⦿ | SendFilePathToQueue | Create SLO | 100% | - | Lambda function SendFilePathToQueue |
| ○ | storeConnection | Create SLO | 100% | - | Lambda function storeConnection |
| ○ | returnPresignedURL_o... | Create SLO | 100% | - | Lambda function returnPresignedURL_or_uploadLink |
| ○ | sendQuery | Create SLO | 100% | - | Lambda function sendQuery |

These are the log groups that contain logs from lambda functions, applications signals and API gateways.

## 11. Application working

User can only upload a single file or a Github link. Repository should not exceed the size of 25 MB. Users are only allowed to upload 1 file or link per session.



## 12. Why?

I used SQS to decouple the front end from the back end to simplify debugging. Another reason was security. Since the API Gateway lacks proper authorization, allowing the front end to directly access the back end would expose it to potential threats. Due to decoupling, the back end doesn't need to handle incoming requests immediately, and the requests are guaranteed to be processed. Without this, high CPU utilization upon receiving requests could prevent them from being processed.

I used DynamoDB to store the userId to connectionId mapping because the application requires frequent read and write access. Using RDS would introduce unnecessary overhead, as each query goes through a SQL engine and maintains a TCP connection with each client, which limits scalability due to connection constraints. Since the data is simple and non-relational, a NoSQL key-value store like DynamoDB is a better fit for this use case.

## 13. Is it Well-Architected?

I have used AWS Managed service DynamoDB so that I don't need to manage the database myself. I have also implemented proper monitoring logs and dashboard so that a failure in any system would be visible. My front-end code can be changed independently of the back-end as both are decoupled using a queue service. The auto scaling EC2 instances fetch code from the source code bucket so manual configuration of code and the environment is not required. Front-

end also changes as per the code changes in the S3 bucket. All of these ensures operational excellence. Usage of CloudFront provides performance efficiency.

I have used roles with least privileges which ensures security of the architecture. A back-end failure would notify me via email, and all the components scale automatically, hence reliability is also ensured.

AWS cost optimizer has been enabled which would recommend me settings for cost optimization as well.

## 14. Security Measures

All the services in this architecture uses specific roles that provide least privilege access. Back-end computation is done in the EC2 instances deployed inside the Private subnet with no inbound access from the internet. Also, the front-end and the back-end has been decoupled by using SQS service. I have used NAT gateway so that EC2 can send messages to connected clients and download files from Github.

## 15. Model Selection

**Microsoft's codebert-base** model was used to create embeddings of query and user uploaded data because the model is trained on code and natural language to create better embeddings.

**CodeLlama 7B Instruct** model is used to answer the user queries. This model has been trained on code and natural languages and is fine tuned to support instruction handling. I didn't make any changes in the model. I deployed the model as it was.

## 16. Cost Optimization strategies

Each stored object is deleted after a day from S3 using lifecycle rul. This reduces cost because of less storage usage.

EC2 auto scaling with dynamic policies allows EC2 instances to shut down and start up whenever necessary. Hence, there is no manual reservation of EC2 instances.

## 17. Lesson learned and Future Improvements

Even though the model works reasonably well, if the embeddings aren't good enough and if the code chunks aren't properly embedded, the generated response could become irrelevant. In future, a better embedding model for query, and a better method for chunking user data could be used so that the context is relevant to the query and the LLM can answer the question. Currently, each file has its own chunk, and the embeddings aren't as close to the chunk's embedding as required for better performance. This causes the top_k result to vary, and therefore the results.