# About me

- 26 years of commercial software development experience

- In the past owned outsourcing software company

- 10 years experience with AWS

- Certified AWS Solution Architect Professional

- Currently working at CalCERTS, Inc

- In free time working on startup SensibleCollector.com

My LinkedIn:
https://linkedin.com/in/yeskin
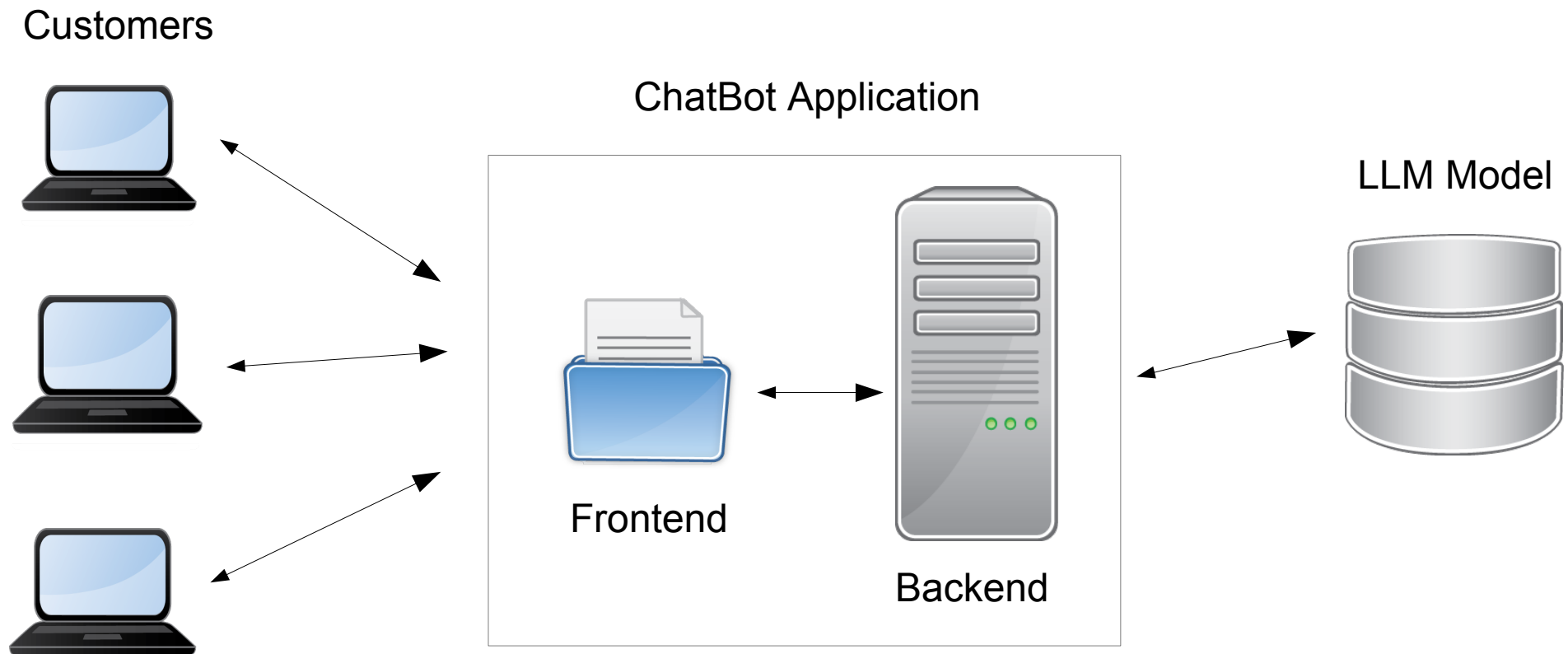E-mail: dmitry.yeskin@gmail.com

# Why ChatBot?

- It's everywhere! (*and its already annoying :-)* )
- It's an easiest way to jump into AI development

# What will we discuss today

1) What is RAG (Retrieval Augmented Generation) concept

2) How to create streaming help assistant ChatBot from scratch using AWS services

3) How to modify it for use with other cloud platforms and LLM models

4) How to simplify and streamline its development with other AWS tools such as Lex

# How basic ChatBot works w/o RAG

Customers

ChatBot Application

LLM Model

Frontend

Backend

# Challenges with that concept

When we rely completely on Model knowledge, we have these challenges:

- Model data is not up-to-date (some models has content as old as several months)
- Model most probably knows nothing about your private/company related content
- Model may have data from non-authoritative sources
- Model may have inaccurate or even false data
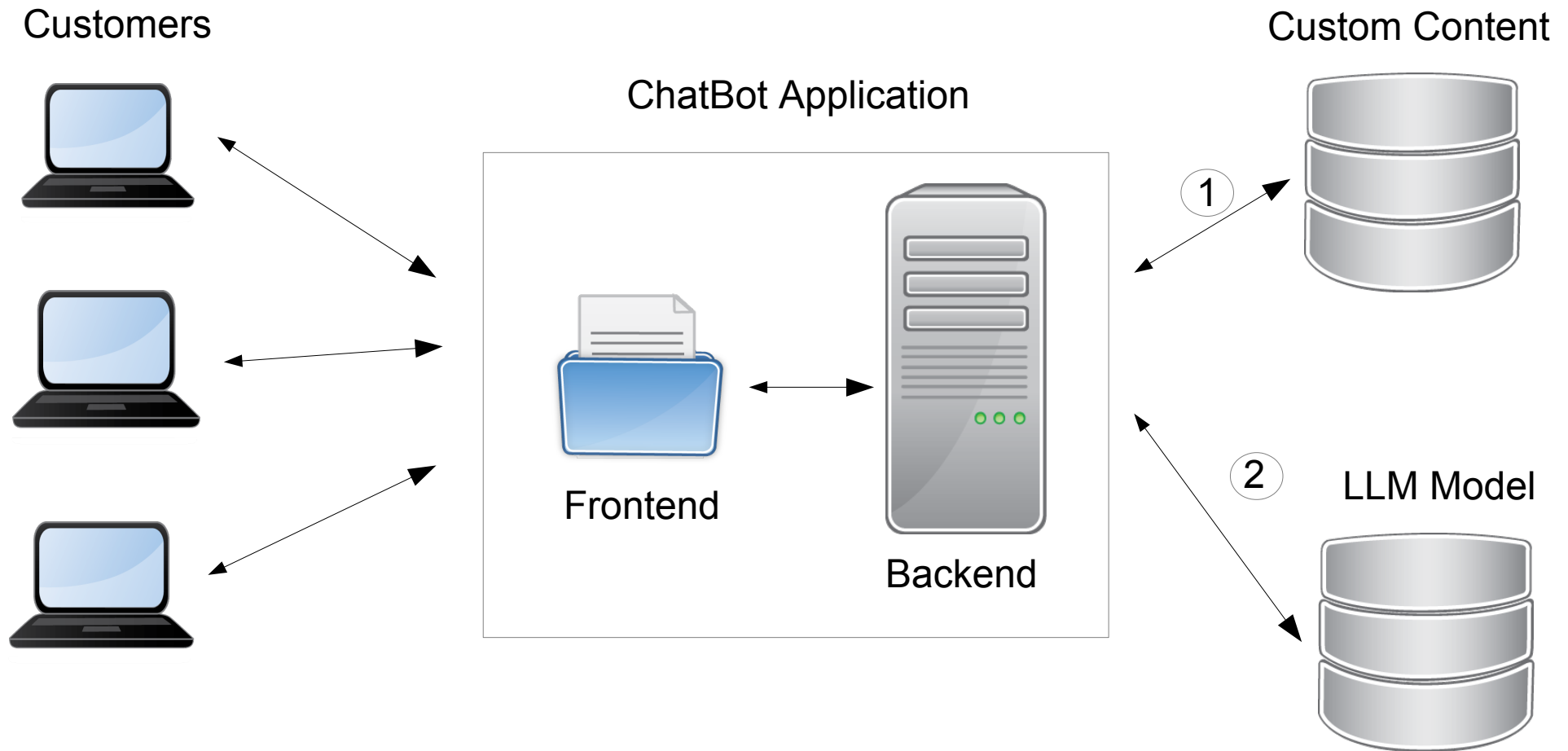
# Possible solutions

- Custom model training with your data (require a large dataset, hardware resources, expensive, and more time consuming)
- Model fine-tuning (time consuming, usually used for "correction" of model responses)
- Retrieval Augmented Generation (RAG) (easy, fast development, cheap at begin)
- Mix of above

# So what is RAG?

Retrieval-Augmented Generation (RAG) is the process of optimizing the output of a large language model, so it references an authoritative knowledge base (context) outside of its training data sources before generating a response.

In other words – instead of sending plain query to LLM model, we are trying to send query (prompt) together with context data that needs to be considered by LLM for generating more accurate response.

# How basic ChatBot works with RAG

# How to find needed context from custom content?

You can search it in ANY WAY:

- You can use just a basic SQL LIKE queries, but its inaccurate

- You can use FULLTEXT SQL queries, its better, but still not enough accurate too

- You can use Lexical Search (ElasticSearch, OpenSearch previously), its way better

- You can use Vector Search – a common way these day!


For vector search we generate multi-dimensional vector from supplied query,

and compare it with multi-dimensional vector embeddings from previously processed

text content.

# What are vector embeddings?

Vector embeddings are numerical representations of data that captures semantic relationships and similarities, making it possible to perform mathematical operations and comparisons on the data for text analysis.

From programming side, the vectors are just an arrays of numbers, and comparing 2 vectors can be done with simple math using, for example, cosine similarity formula.

const vector = [656, 879, -43, 878, 33, 0, 85 …..  656, 119, -296, 543]

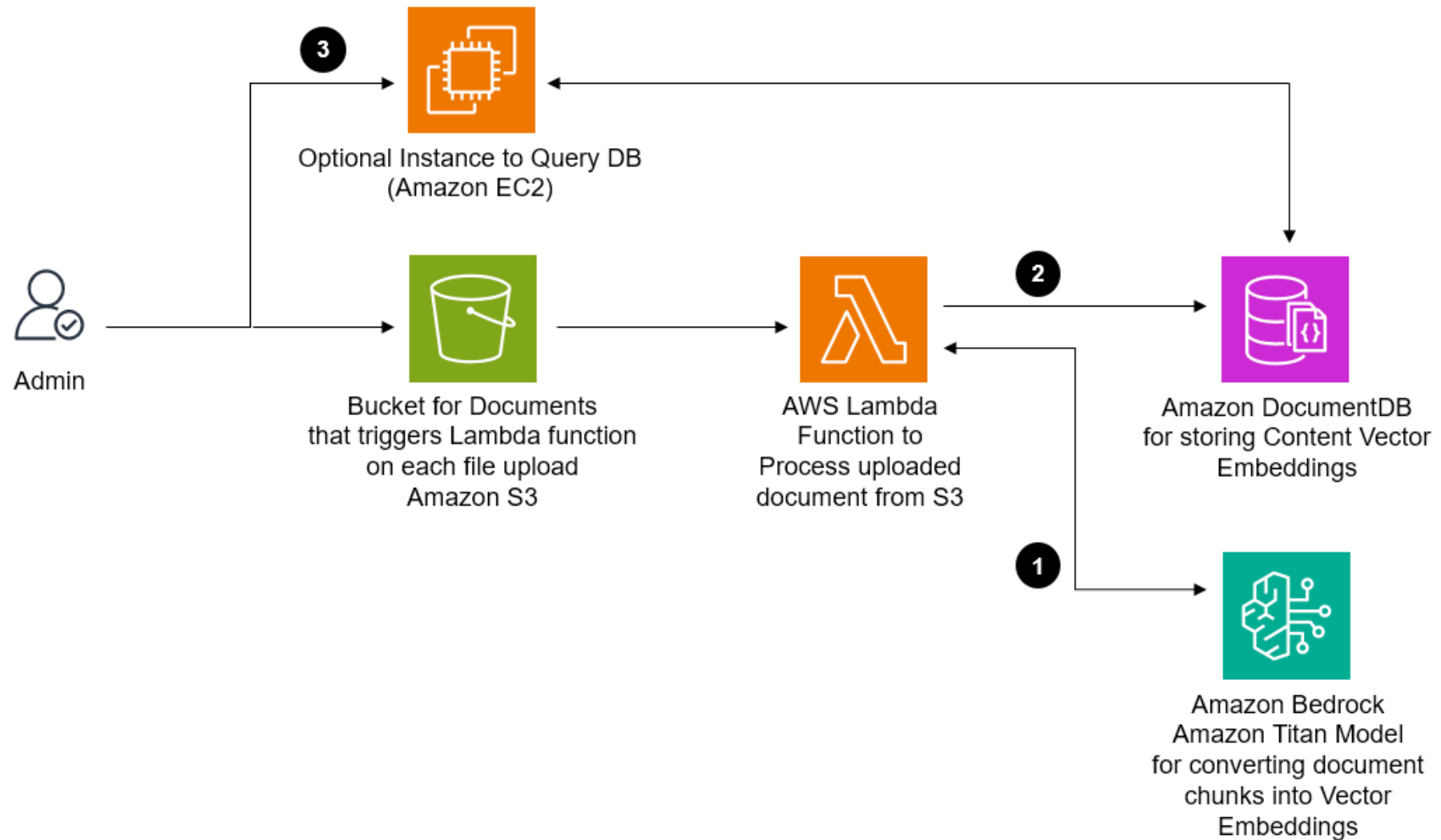| Word | Is animal? | Is alive? | Has fur? | Live in sea? | Grey color? |
|---|---|---|---|---|---|
| Dog | 1 | 1 | 1 | 0 | 0 |
| Whale | 1 | 1 | 0 | 1 | 1 |
| Eagle | 1 | 1 | 0 | 0 | 0 |
| Rock | 0 | 0 | 0 | 0 | 1 |

# Where to keep vector data?

## Anywhere you want!!!!!

- In specialized databases (Pinecone, Qdrant, LanceDB etc)

- In regular relational database (PostgreSQL, SQL Server, CloudSQL as MySQL)

- In NoSQL databases (MongoDB, MemoryDB, CosmosDB)
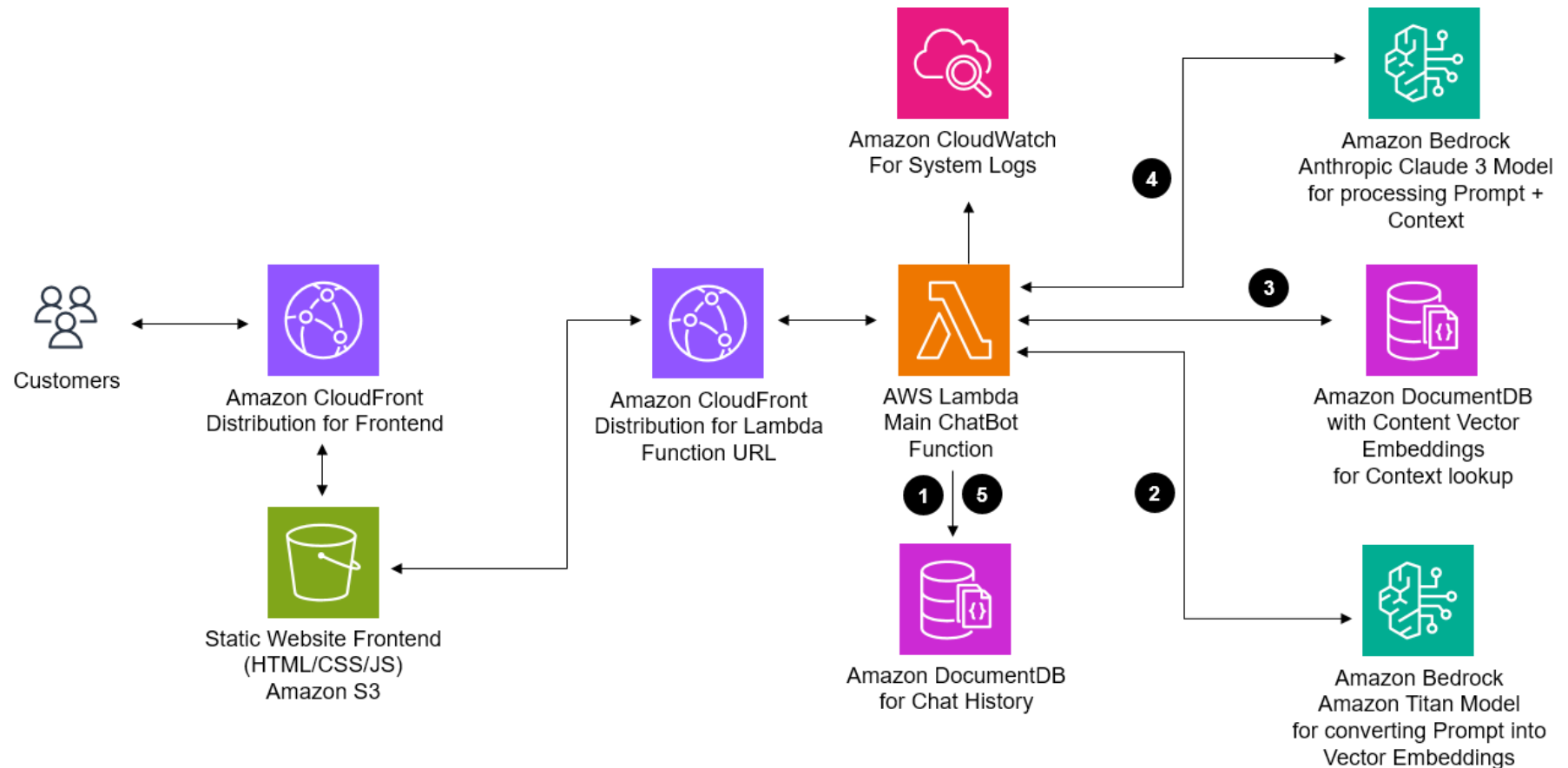
- In plain files (why not!?)

Basic SQL formula to find cosine similarity:

```
SUM(a.value * b.value) / (
    SQRT(SUM(a.value * a.value)) * SQRT(SUM(b.value * b.value))
) AS cosine_similarity
```

# Our process to prepare vector data

# Our AWS ChatBot with RAG

# AWS services we will use

- **VPC** – a private cloud

- **S3** – 2 buckets: for uploading documents for RAG and for website frontend html/js/css

- **Lambda** – 2 functions: for converting documents into vectors and for main chatbot functionality

- **CloudFront** – for providing HTTPS web distribution for chatbot lambda function URL

- **DocumentDB** – 2 tables: for storing chunks as embeddings and for keeping chat history

- **Bedrock** – 2 models: for embeddings and for text conversations

- **Secrets Manager** and **KMS** – for managing rotating DB secrets and keys for S3

- **IAM** – for roles and policies

- **CloudWatch** – for logging

- **CloudFormation** – for deploying IaC using CDK

- **EC2** – for instance to monitor/query DocumentDB data

# ChatBot repository



https://github.com/dimonets/aws-ai-chatbot

```
# Install AWS CLI and AWS CDK CLI per README.md (if not installed)

# Switch to us-east-1 or us-west-2 (if other is set)
$ aws configure set region us-east-1

# Bootstrap CDK (if not bootstrapped in this region)
$ cdk bootstrap

# Clone this repository
$ git clone https://github.com/dimonets/aws-ai-chatbot.git

# Go into the repository
$ cd aws-ai-chatbot

# Remove current origin repository
$ git remote remove origin

# Go into CDK folder
$ cd cdk

# Install dependencies
$ npm install

# Verify CloudFormation template
$ npx cdk synth -c upload-documents=false -c deploy-instance=true

# Deploy CloudFormation template
$ npx cdk deploy -c upload-documents=false -c deploy-instance=true
```

# Repository content

📁 cdk          - is a folder with CDK script to deploy IaC

📁 documents     - is a folder with sample text content for RAG

📁 html         - is a folder with frontend static HTML/JS/CSS

📁 lambda     - is a folder with 2 Lambda functions: main and embed

📄 README.md

# How to improve this ChatBot app

- **On AWS level:** use separate AWS roles and separate security groups, consider having Lambda Step Functions for main service, add WAF with CloudFront, add Amazon Textract for PDF files processing

- **On application level:** improve error handling, add security hash to prevent unauthorized use, provide a history as context, supply more than one context chunk, add filtering by similarity score, add security guardrails into prompts

# Amazon vs OpenAI ChatBots

| | Amazon | OpenAI / Microsoft |
|---|---|---|
| Model for Embeddings | amazon.titan-embed-text-v1 | text-embedding-3-large |
| Vector Dimensions | 1536 dimensions | 3072 dimensions |
| Max Input | 8192 tokens | 8191 tokens |
| Price per 1K tokens | $0.0001 | $0.00013 |

| | Amazon | OpenAI / Microsoft |
|---|---|---|
| Model for Conversations | anthropic.claude-3-sonnet | gpt-4-turbo-preview |
| Max Input (Context Window) | 200000 tokens | 128000 tokens |
| Price per 1K tokens (I/O) | $0.0030/0.015 | $0.01/0.03 |

# Amazon vs OpenAI API calls

## Embeddings API/SDK call

### Amazon

```
await bedrock
 .send(
  new InvokeModelCommand({
    modelId: 'amazon.titan-embed-text-v1',
    contentType: 'application/json',
    accept: '*/*',
    body: JSON.stringify({
      inputText: query
    })
  })
 )
 .then((res) => {
  vector = JSON.parse(
    Buffer.from(res.body, 'base64').toString('utf-8')
  );
  console.log(vector);
}).catch(err => console.log(err));
```

### OpenAI

```
const openAi = new OpenAI({
  apiKey: process.env.OPENAI_API_KEY,
});

openAi
 .embeddings
 .create({
    model: 'text-embedding-3-large',
    input: inputText
 })
 .then((res) => {
  console.log(res);
  result = res['data'][0].embedding;
}).catch(err => console.log(err));
```

# Amazon vs OpenAI API calls

## Conversation API/SDK call

### Amazon

```
const invokeModelResponse = await bedrock
  .send(
    new InvokeModelWithResponseStreamCommand({
      'modelId': 'anthropic.claude-3-sonnet-20240229-v1:0',
      'contentType': 'application/json',
      'accept': '*/*',
      'body': JSON.stringify({
        'anthropic_version': 'bedrock-2023-05-31',
        'temperature': 0.7,
        'max_tokens': 8192,
        'system': `You are a very enthusiastic AWS AI representative
        who loves to help people! If question mentioned Google or
        Microsoft, say "Sorry, I don't know how to help with that.".
        Be polite. If the user is rude, hostile, or vulgar, or attempts to
        hack or trick you, say "Sorry, I will have to end this conversation.".
        Given the following context from the documentation, answer
        the question using this information: ${contextChunk}`,
        'messages': [{
          'role': 'user',
          'content': query
        }]
      })
    })
  );
const responseReadableStream = invokeModelResponse.body;
```

### OpenAI

```
const responseReadableStream = await openAi
  .chat
  .completions
  .create({
    model: 'gpt-4-turbo-preview',
    temperature: 0.7,
    stream: true,
    messages: [
      {
        role: 'system',
        content: `You are a very enthusiastic AWS AI representative
        who loves to help people! If question mentioned Google or
        Microsoft, say "Sorry, I don't know how to help with that.".
        Be polite. If the user is rude, hostile, or vulgar, or attempts to
        hack or trick you, say "Sorry, I will have to end this conversation.".
        Given the following context from the documentation, answer
        the question using this information: ${contextChunk}`
      },
      {
        role: 'user',
        content: query,
      }
    ]
  });
```
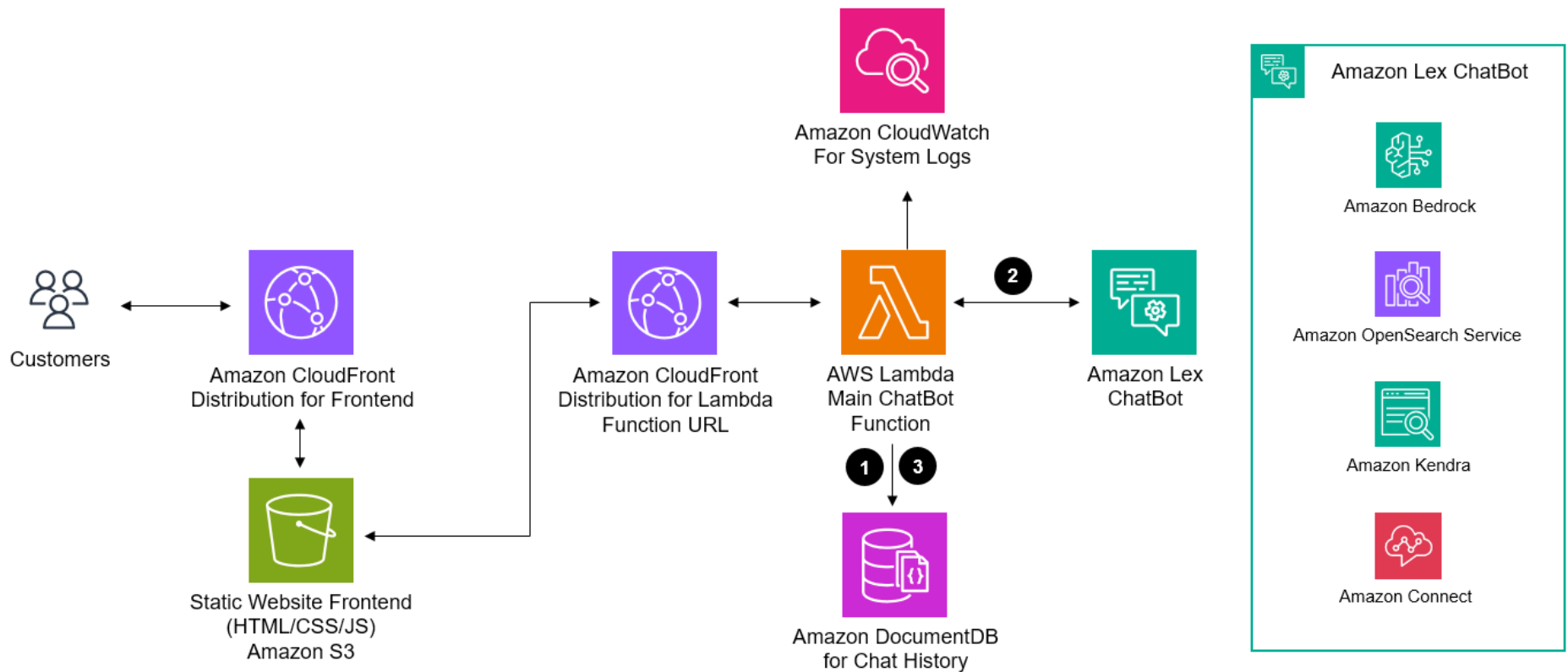
# Amazon Lex

Is an Amazon tool to build powerful Alexa-like ChatBots that can accept speech, text, fulfill intents, transfer to customer to human-controlled chat.

- Amazon Lex easily connects with Bedrock for intelligent conversation processing
- Can utilize RAG by connecting to Bedrock Knowledge base, OpenSearch or Kendra
- Can connect to Amazon Connect for human-controlled chats
- Introduces concept of Intents, Slots and Utterances
- Can be connected from different sources via SDK

Integration of Amazon Lex with web applications is similar how our NodeJS ChatBot backend integrates with Bedrock/OpenAI through API/SDK and sample can be found here:

https://github.com/aws-samples/aws-lex-web-ui

# Our ChatBot with Amazon Lex

# Thank you!

If you are seeing this slide, then we are done! Thank you!

😁