aws
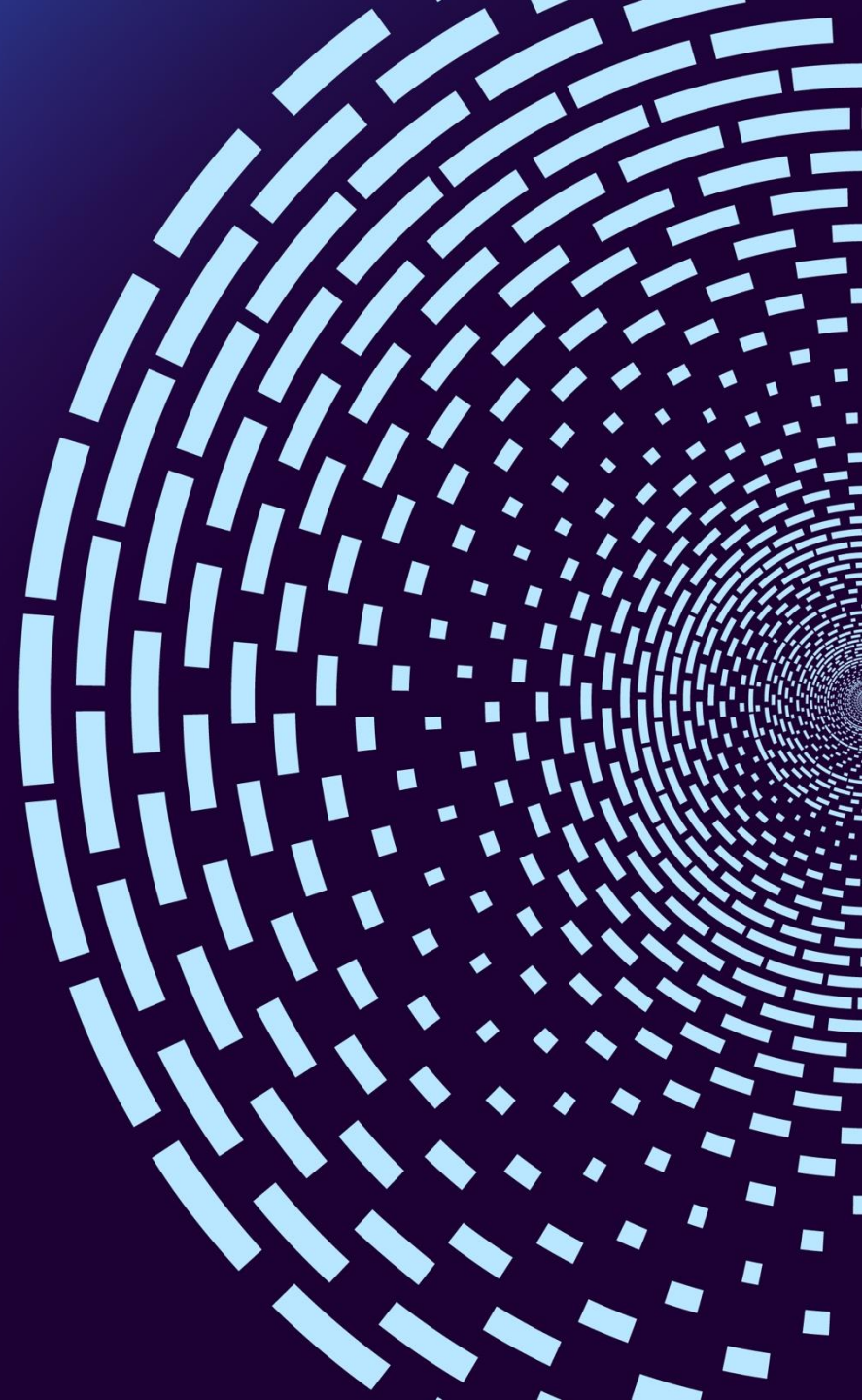
# AI Conclave

Online
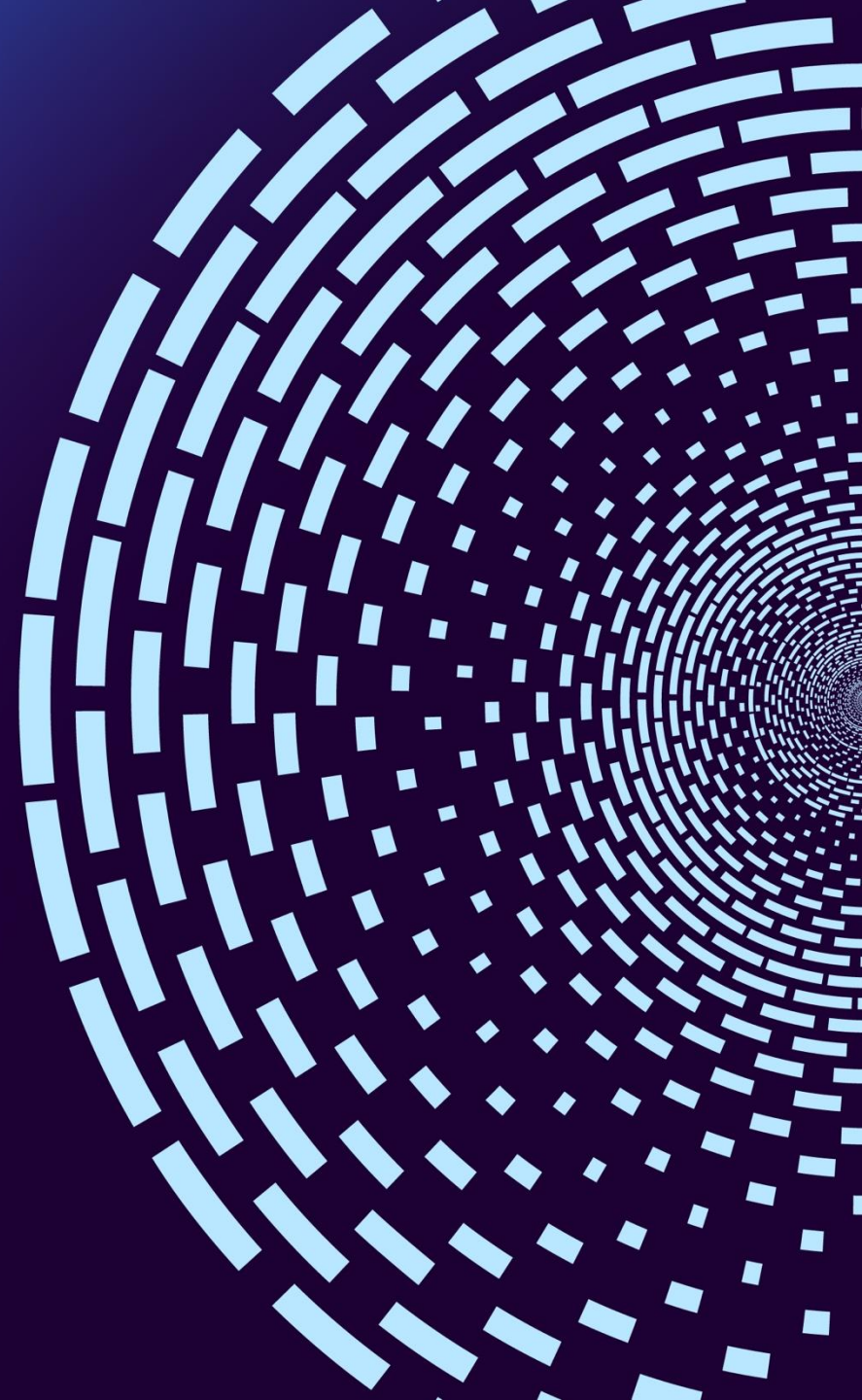
aws

# Build application using advanced RAG methods and validate using different evaluation mechanism

**Shailesh Shivakumar**

Senior Solutions Architect
AWS India

# Agenda

- RAG introduction
- Challenges with Normal RAG
- Advanced RAG – Parsing
- Advanced RAG - Chunking
- Advanced RAG – Query reformulation
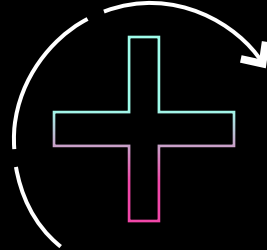- RAG evaluation
- Demo

# RAG introduction

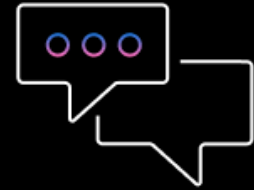# What is Retrieval Augmented Generation?

## Retrieval

Fetches the relevant content from the external knowledge base or data sources based on a user query

## Augmentation

Adding the retrieved relevant context to the user prompt, which goes as an input to the foundation model

## Generation

Response from the foundation model based on the augmented prompt

# RAG use cases

**Improved content quality**

**E.g.**, helps in reducing hallucinations and connecting with recent knowledge including enterprise data
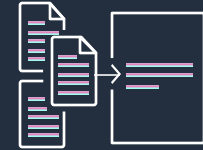
**Contextual chatbots and question answering**

**E.g., e**nhance chatbot capabilities by integrating with real-time data

**Personalized search**

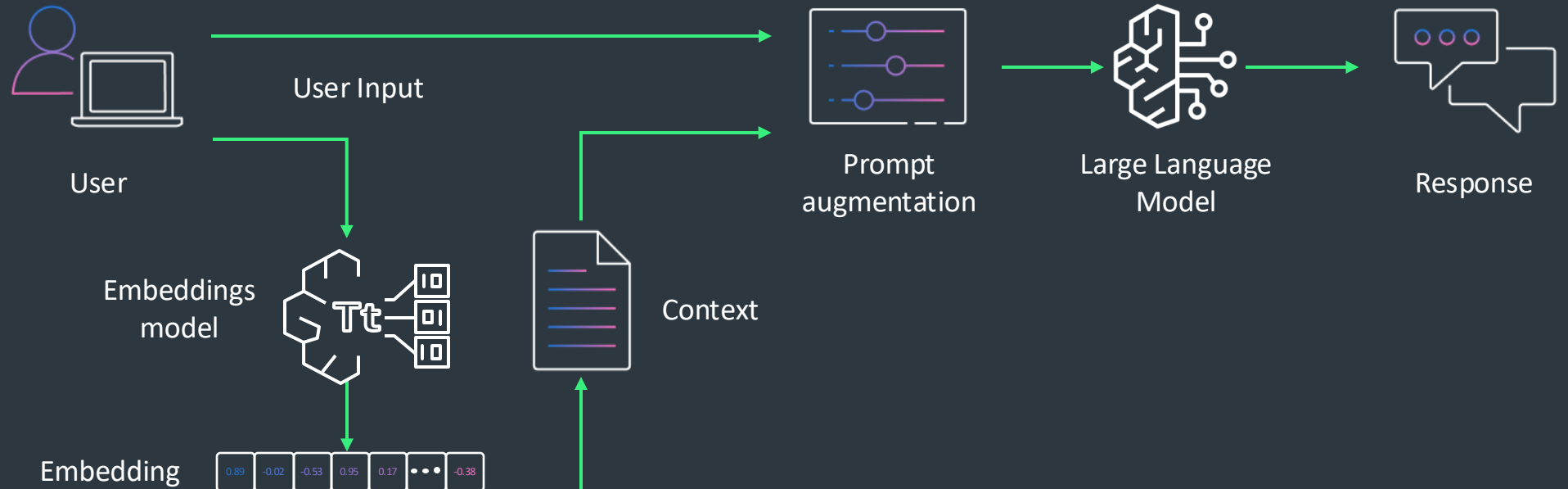**E.g.,** searching based on user previous search history and persona
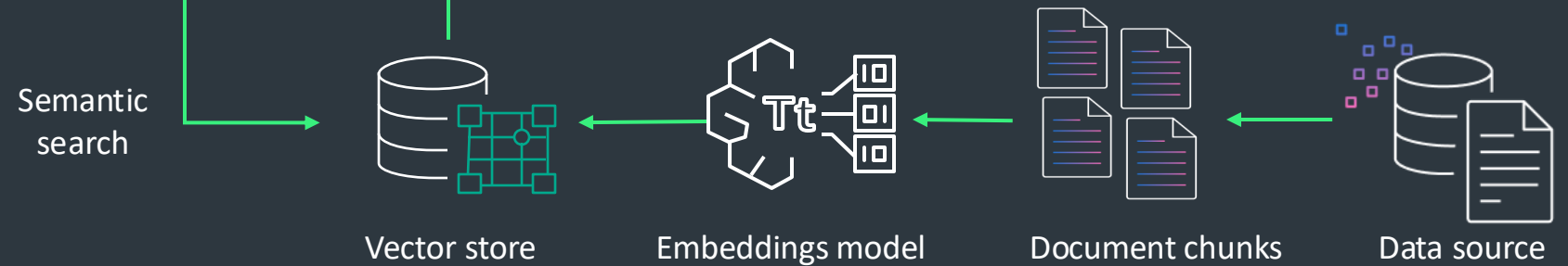
**Real-time data summarization**

**E.g., r**etrieving and summarizing transactional data from databases, or API calls

# RAG in Action



Text Generation Workflow

User Input

User

Embeddings model

Embedding

| 0.89 | -0.02 | -0.53 | 0.95 | 0.17 | ••• | -0.38 |

Context

Prompt augmentation

Large Language Model

Response

Data Ingestion Workflow

Semantic search

Vector store

Embeddings model

Document chunks

Data source

# Knowledge Bases for Amazon Bedrock

Gives FMs and agents contextual information from your private data sources for Retrieval Augmented Generation (RAG) to deliver more relevant, accurate, and customized responses.

**Fully managed support for end-to-end RAG workflow**

**Securely connect FMs and agents to data sources**

**Easily retrieve relevant data and augment prompts**

**Provide source attribution**

# RAG challenges

# Challenges with Regular RAG Approach

- **Parsing complex structures** – Entities nested tables, images, graphs are not parsed accurately.

- **Use relevant data from large dataset** – As the dataset grows, only the configured set of top results are used for context impacting accuracy.

- **Responding to complex questions** – A compound question (such as ones with different sub-questions) pose challenges to the RAG solution.

- **Retrieve right context with complex document structures** – Few documents such as legal documents or technical manuals have semantically related sections in different pages impacting the overall accuracy.

- **Generate complete, correct and grounded response** without any hallucination.

# Workarounds

- **Programatic parsing** – Use different service (such as Amazon Textract) to extract the table content into .md file

- **Programmatic split of compound questions** – Programmatically split the question or reformulate the question and get the responses and programmatically aggregate the response.

- **Programatic chunking** – Programmatically create custom chunks of the documents and mange them in vector store.

# Advanced RAG - Parsing

# Advanced RAG with Amazon Bedrock Knowledge Bases

**Chunking and parsing configurations** Info
Choose between default or advanced customization.

○ **Default**
Uses default parsing and chunking strategy.

● **Custom**
Customize the parsing and chunking strategy, including using advanced parsing.

**Parsing strategy**
Parsing analyses and extracts useful information from documents.

☑ Use foundation model for parsing 🔗
Suitable for parsing more than standard text in supported document formats, including tables within PDFs with their structure intact. View pricing 🔗

**Choose foundation model for parsing**

A\ **Claude 3 Sonnet v1** 🔗        ○
By Anthropic

A\ **Claude 3 Haiku v1** 🔗        ○
By Anthropic

▼ **Instructions for the parser - *optional***

```
1  Transcribe the text content from an image page and output in Markdown syntax (not code blocks). Follow these
       steps:
2
3  1. Examine the provided page carefully.
```

# Advanced RAG - Chunking

# Chunking strategy

**What**

**Chunking**: Segmenting text into meaningful, concise chunks
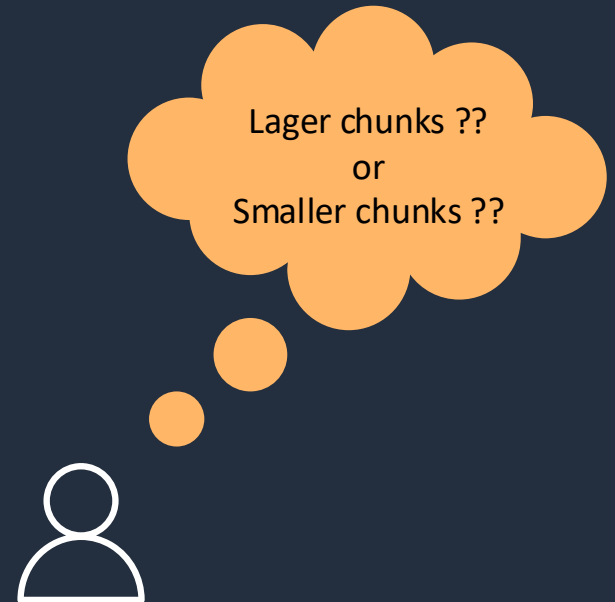
**Why**

**Enables** semantic text embedding

**How**

Effectiveness of chunking strategy largely depends on **quality and structure of the chunks**.

Optimal chunk size balances relevant context and speed

Lager chunks ??
or
Smaller chunks ??

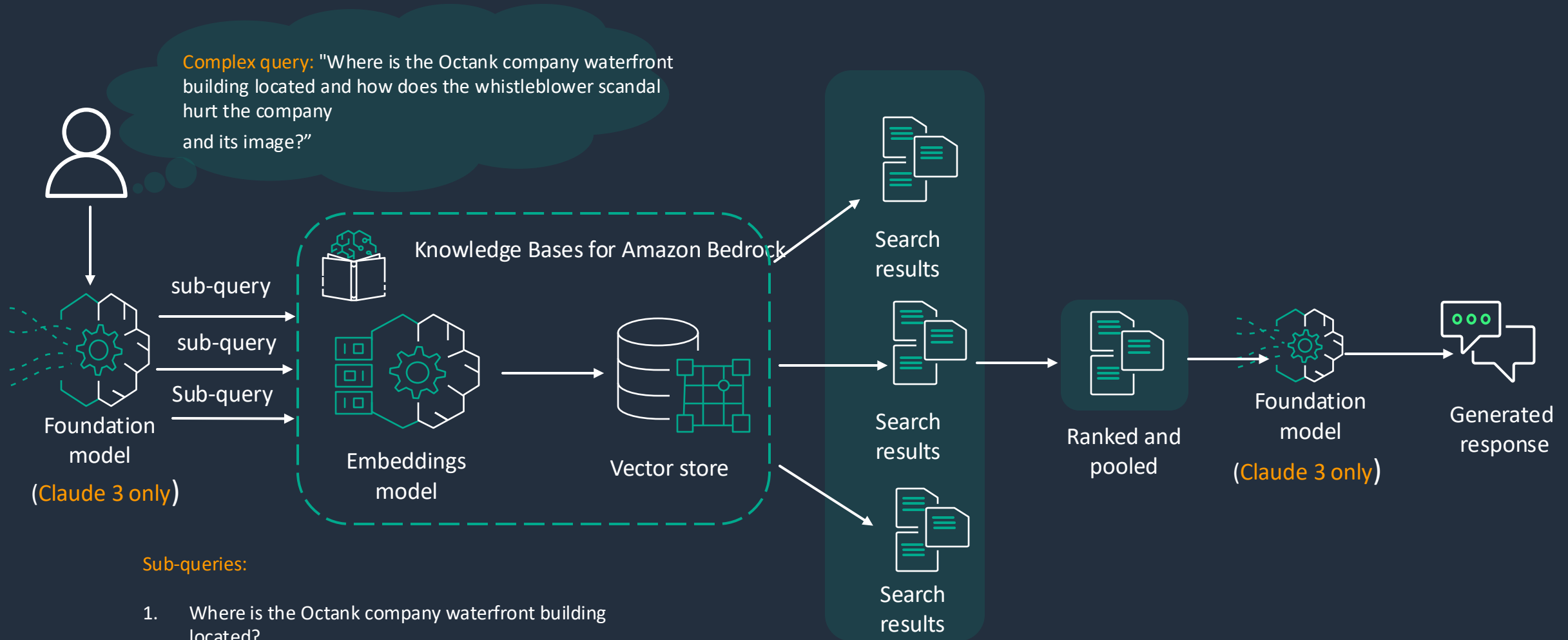# Overview of different chunking approaches

| Type | Description | Pros | Cons | Remarks |
|---|---|---|---|---|
| Fixed chunking | - Fixed character count division<br>- Recursive character text splitting | Quick and efficient | Lacks context awareness | Data is relatively uniform in length and structure. |
| Semantic chunking | - Meaningful and complete chunks based on semantic similarity | Better retrieval quality | | Suitable for lengthy documents that have related sections. |
| Hierarchical chunking | - Organize chunks in a structured manner<br>- Nodes with summaries, parent-child relationships | Improve retrieval efficiency and reliability | Require domain-specific or problem-specific expertise | When required context is split across multiple different documents |

# Advanced RAG – Query reformulation

# RAG API - Query reformulation

New

Complex query: "Where is the Octank company waterfront building located and how does the whistleblower scandal hurt the company

and its image?"

Knowledge Bases for Amazon Bedrock

sub-query

sub-query

Sub-query

Foundation model
(Claude 3 only)

Embeddings model

Vector store

Search results

Search results

Search results

Ranked and pooled

Foundation model
(Claude 3 only)

Generated response

Sub-queries:

1. Where is the Octank company waterfront building located?
2. What is the whistleblower scandal involving Octank company?
3. How did the whistleblower scandal affect Octank company's reputation and public image?

# RAG evaluation

# Special challenges with RAG evaluation

Use relevant data from your knowledge base

Retrieve the right context from documents

Generate a correct, complete, and grounded answer minimizing hallucinations

Iteratively improve your RAG system and compare across changes

Evaluate biases, safety, and trust

# RAG evaluation input data format

Input dataset contains 2 things in JSONL format
1. Prompt
2. Optional golden ground truth

```
{

    "conversationTurns": [{
            "referenceResponses": [{
                    "content": [{
                            "text": "This is a reference response"
                    }]
            }],
            "prompt": {
                    "content": [{
                            "text": "This is a prompt"
                    }]
            }
    }]
}
```

Retrieve jobs have referenceContexts

# Choice of evaluation metrics

**Retrieval**

Context Coverage

Context Relevance

**Retrieval + Generation**

Correctness

Completeness

Helpfulness

Logical coherence

Faithfulness

**Responsible AI**

Harmfulness

Stereotyping

Refusal

# How correctness works

## Example input

prompt: What is the capital of Spain?

referenceResponse: Madrid

Model response: Barcelona

## Judge prompt (simplified)

You are a helpful assistant…

You are given a question, a candidate response from an LLM, and reference response.

Your task is to check if the candidate response is correct compared to the reference response…

Here is the actual task:
Question:  {prompt}
Reference Response: {referenceResponse}
Candidate Response: {Model response}

Explain your response, followed by your evaluation:
2) Correct
1) Partially correct
0) Incorrect

Note: Correctness can be with or without ground truth

# How RAG evaluation works with Knowledge Bases



Choose evaluator model

Choose your Knowledge Base

Choose to evaluate retrieval only or retrieve + generate

Choose your generator model

Choose your metrics

Upload your prompt dataset

Inference and evaluation

View results

# Get results in a few clicks



- Simple to read scores
- See distributions visually
- See ratings explanations

# Demo

# Thank you!

**Shailesh Shivakumar**

Senior Solutions Architect
AWS India