# Leveraging RAG architecture for generative AI applications

**Jatin Pal Singh**

Principal Solutions Architect
Amazon Web Services

**Ashutosh Dubey**

Sr Solutions Architect
Amazon Web Services

aws

# Agenda

Customizing foundation model for your use case

Introduction to retrieval augmented generation (RAG)

Use cases RAG

Demo of RAG pattern using Vector DB and LLM(Amazon Bedrock)

Overview and demo Knowledge bases for Amazon Bedrock

Overview and demo Amazon Q for business

# Why customize?

### Adapt to domain-specific language

E.g., Healthcare – Understand medical terminology and provide accurate responses related to patient's health

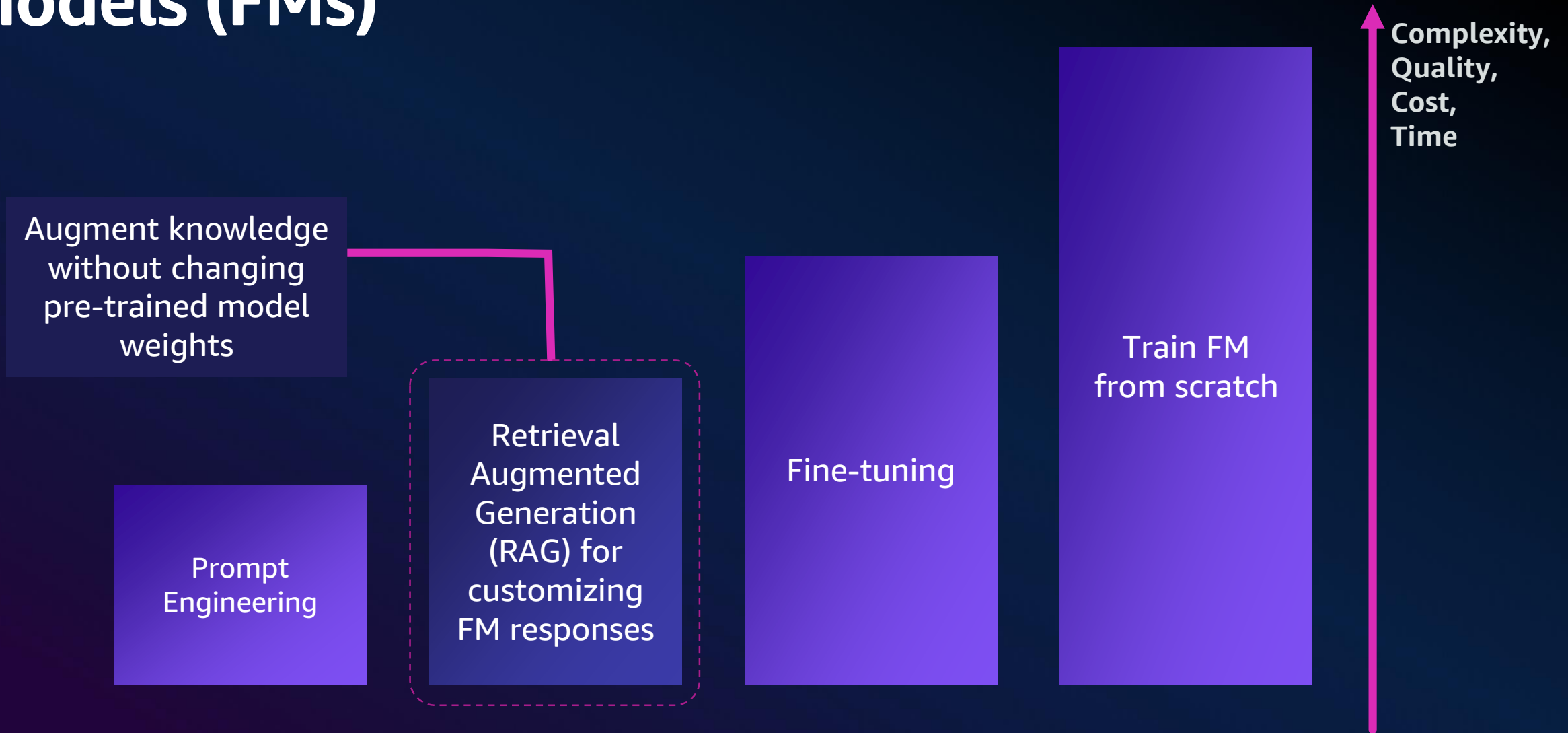### Enhance performance for specific tasks

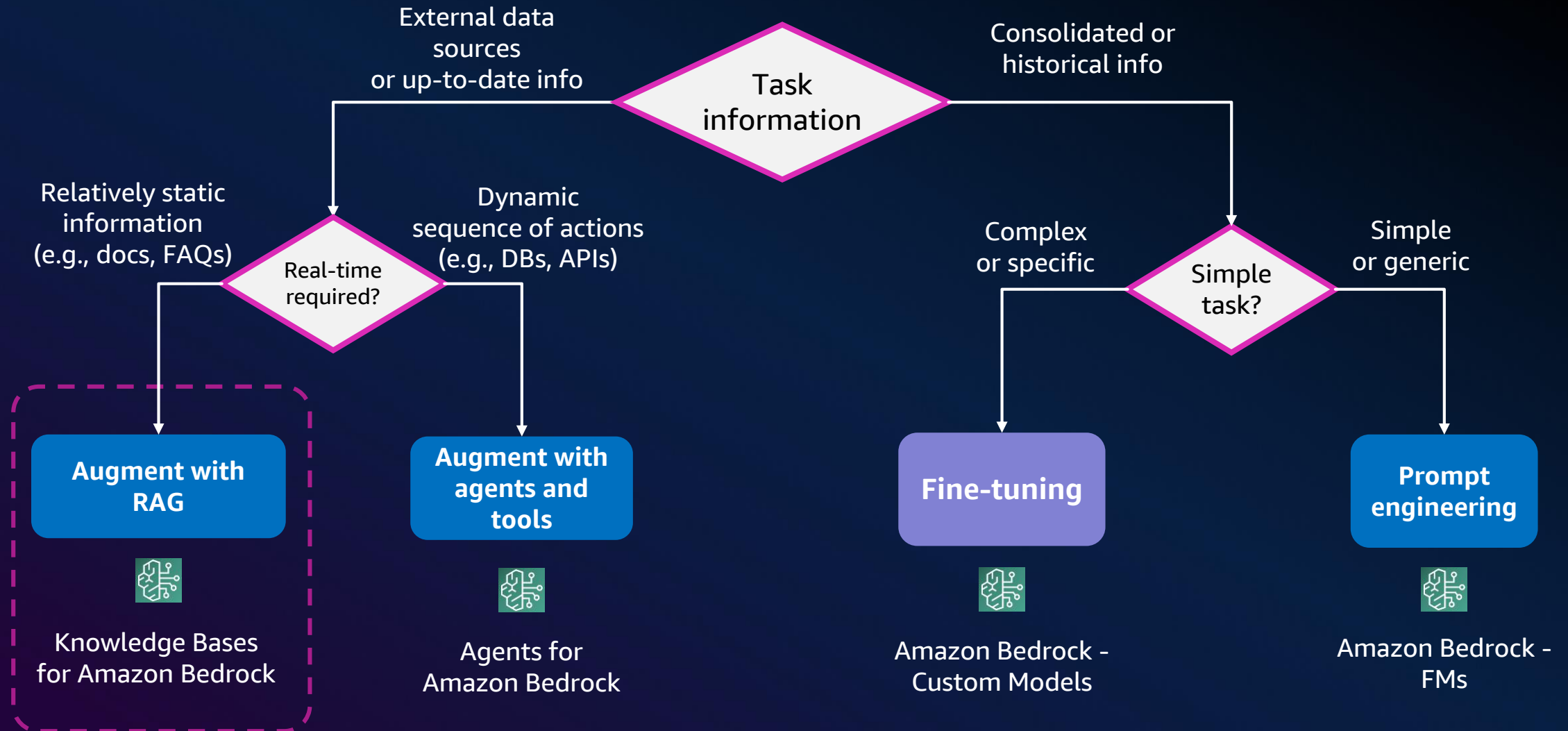E.g., Finance – Teach financial & accounting terms to provide good analysis for earnings reports

### Improve context-awareness in responses

E.g., Customer Service – Improve ability to understand and respond to customer's inquires and complaints

# Common approaches for customizing foundation models (FMs)

Complexity,
Quality,
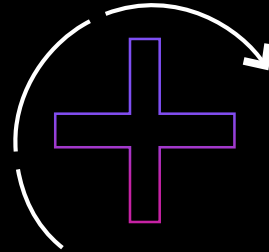Cost,
Time

Augment knowledge without changing pre-trained model weights

Prompt Engineering

Retrieval Augmented Generation (RAG) for customizing FM responses

Fine-tuning

Train FM from scratch

# Customize vs. augment



**Task information** (decision diamond)

- External data sources or up-to-date info →
  - **Real-time required?**
    - Relatively static information (e.g., docs, FAQs) →
      **Augment with RAG**
      Knowledge Bases for Amazon Bedrock
    - Dynamic sequence of actions (e.g., DBs, APIs) →
      **Augment with agents and tools**
      Agents for Amazon Bedrock
- Consolidated or historical info →
  - **Simple task?**
    - Complex or specific →
      **Fine-tuning**
      Amazon Bedrock - Custom Models
    - Simple or generic →
      **Prompt engineering**
      Amazon Bedrock - FMs

# 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.,** enhance 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

# Types of retrieval

## Rule Based

Fetches unstructured data like documents

e.g., Key word searches

## Structured data

Transactional retrieval from database or API

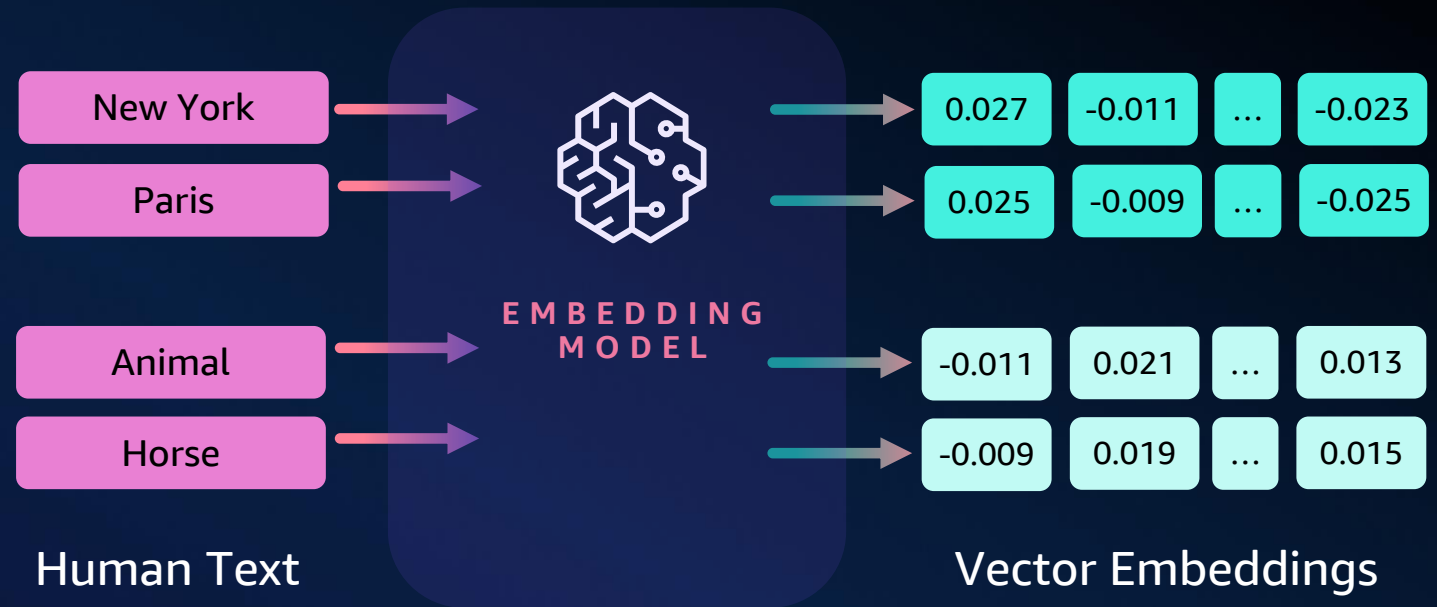e.g., Select customers from All_orders where order == 'XYZ'

## Semantic Search
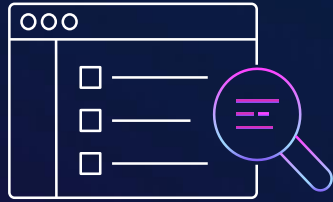
Get relevant documents based on text embeddings

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

New York ⟶ Subway
Statue Liberty
Tall buildings

# What are embeddings?

- Numerical representation of text (vectors) that captures semantics and relationships between words.

- Embedding models capture features and nuances of the text.

- Rich embeddings can be used to compare text similarity.

- Multilingual Text Embeddings can identify meaning in different languages.



Human Text

EMBEDDING MODEL

| 0.027 | -0.011 | ... | -0.023 |
| 0.025 | -0.009 | ... | -0.025 |

| -0.011 | 0.021 | ... | 0.013 |
| -0.009 | 0.019 | ... | 0.015 |

Vector Embeddings

New York
Paris
Animal
Horse

# Why are embeddings important for RAG?

Powers **text retrieval based on semantic meaning.**

Used to augment prompts with **more accurate context** from vector stores using the Retrieval Augmented Generation (RAG).

High-accuracy embeddings leads to improved context **and higher quality LLM-generated** responses to a user query.

# Titan text embeddings model

## Amazon Titan Text Embeddings
V**2.0**

Translates text inputs (words, phrases) into numerical representations (embeddings). Comparing embeddings produces more relevant and contextual responses than word matching.

Max Tokens: **8,000**
Output Vectors**: 1,536**
Language: **Multilingual** (25 languages)

Model ID: *amazon.titan-embed-g1-text-02*

## Highlights

• Titan Text Embeddings offers fast, cost effective, high-performance, accurate embeddings in 25 languages.

• Optimized for text retrieval tasks, semantic similarity and clustering.

• Applications of this model includes semantic search and personalization.

# RAG in Action



**Text Generation Workflow**

User

User Input

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

# However, when it comes to implementing RAG, there are challenges...

**Managing multiple data sources**

**Creating vector embeddings for large volumes of data**

**Incremental updates to vector store**

**Coding effort**

**Scaling retrieval mechanism**

**Orchestration**

# Demo – Vector DB and Bedrock

# 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**

# Data Ingestion Workflow

**Fully managed data ingestion workflow**

New data → Data source → Document chunks → Embeddings model → Vector store

- Choose your data source (Amazon S3)
- Support for incremental updates
- Multiple data file formats supported

- Choose your chunking strategy
  - Fixed chunks
  - No chunking
  - Default (300 tokens)

- Choose your embedding model
  - Amazon Titan
  - Cohere

- Choose your vector store
  - Open search serverless
  - Pinecone
  - Redis
  - Aurora PostgreSQL pgvector

# RetrieveAndGenerate API



Fully managed RAG

User → User query → RetrieveAndGenerate API → Response

RetrieveAndGenerate API → User query → Generated response → RetrieveAndGenerate API

Generate query embedding → Retrieve similar documents from knowledge base → Augment query with retrieved documents → Generate response from LLM

# Demo – Knowledge Bases for Amazon Bedrock

# Amazon Q is your business expert (Preview)

## BOOST YOUR WORKFORCE PRODUCTIVITY WITH GENERATIVE AI

> Delivers quick, accurate, and relevant answers to your business questions, securely and privately and document repositories.

> Provides responses with references and citations for easy fact-checking

> Respects existing access control based on user permissions

> Connects to over 40 popular enterprise applications and document repositories

> Enables administrators to easily apply guardrails to customize and control responses

# Key Features

Trusted answers generated from enterprise data

Citations and source attribution

Conversation history and context

Upload files and analyze content

Execute Actions across multiple Enterprise Apps

# Key Features

Use pre-built guardrails for toxicity

Restrict responses to enterprise content only

Specify blocked words or phrases that never appear in responses

Define special topics and configure guardrails for such topics

---

**Enterprise controls** Info

**Application guardrails** Info
Application guardrails will apply to all messages returned by Enterprise Q.

**Response settings** Info
You can limit Enterprise Q from using its own knowledge to generate answers when it cannot find relevant content in your enterprise corpus.

☐ Only produce responses from retrieval augmented generation (RAG)
   Responses will be limited to ingested documents in your enterprise corpus.

**Blocked words** Info
Define blocked words for the application. The application will not respond to questions that contain these words or mention them in any responses.

*Enter blocked words*    ( Add )

You can block 20 more words.

**Messaging shown for blocked words**

I am not allowed to talk about this topic. Please contact your Admin for more details.

Valid characters are a-z, A-Z, 0-9, _ (underscore) and - (hyphen). This response can have up to 150 characters.

**Feature settings** Info
Configure features the end users have access to in the web experience.

☑ Allow end users to upload files in chat context
   This feature enables end users to upload files directly to chat in order to ask questions specific to the document

# Key Features

## FASTER TIME TO MARKET

In-built vector index with managed ingestion

In-built application with SSO

3 click setup:
Settings, retriever, and data sources

Accuracy of retriever-augmented generation (RAG)

# Amazon Q for business Demo

# Thank you!

Jatin Pal Singh

sinjatin@amazon.com

Ashutosh Dubey

ashdbey@amazon.com