# Unlock data insights with Amazon SageMaker and Amazon CodeWhisperer

**Viktoria Semaan**

Senior Developer Advocate

AWS

**Linda Haviv**

Developer Advocate

AWS

# Session agenda

**Building Gen AI apps: Emerging tech stack**
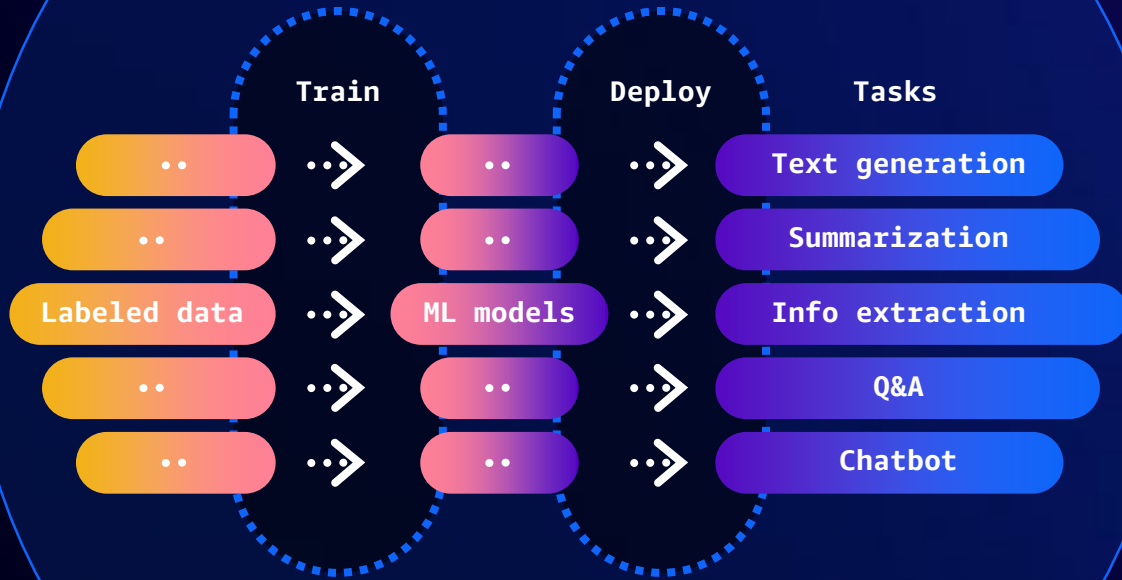
**Demo 1: RAG with Amazon SageMaker and Amazon Kendra**

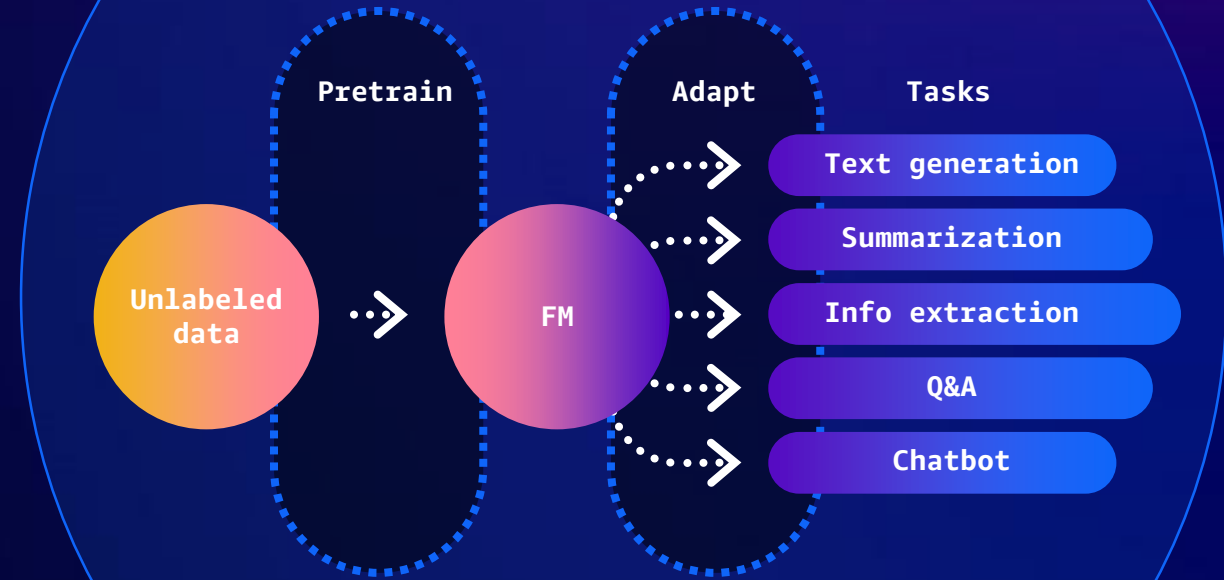**Demo 2: RAG with Amazon Bedrock and a vector database**

**Takeaways and resources to get started**

TRADITIONAL ML MODELS

Train → Deploy → Tasks

Labeled data → ML models → Text generation, Summarization, Info extraction, Q&A, Chatbot

FOUNDATION MODELS

Pretrain → Adapt → Tasks

Unlabeled data → FM → Text generation, Summarization, Info extraction, Q&A, Chatbot

**Applications**
Experience layer: DTC or B2B Apps

**MODELS**

**Proprietary Foundation Models**

Large scale, pretrained models exposed via APIs

**Model Hubs**

Platforms that allow users to share models and datasets

**Open-Source Foundation Models**

Accessible trained models

**ML Platforms**
ML Platforms and compute resources exposed to developers

**Hardware Infrastructure & Accelerators**

Accelerator chips optimized for model training and inference workloads

**Instructions**

**Context**

**User Input**

**Output response**

Act as you are an experienced Data Scientist. Consider the context below to answer the user's questions with a friendly tone. Provide a short concise answer.

Context: You are at a meeting with customers from a manufacturing company that are interested in starting developing Gen AI applications. They have a question.

User: We would like to establish success criteria for our first POC. Can you suggest metrics?

Here are some suggested metrics to determine the success of your first AI proof-of-concept:

- Accuracy: Measure how accurate the predictions or outputs from the AI model are compared to real-world data. Accuracy is a key indicator of whether the AI model is working as intended.

- Time savings: Track how much time is saved by automating tasks with AI versus human effort. This helps quantify productivity gains.

- Cost reduction: Calculate cost reductions from improved efficiency, less waste, and other benefits of implementing AI. This is important for justifying ROI.

- User satisfaction: Survey end users to gauge their satisfaction with the AI solution. This subjective feedback is valuable for iteration.

- Adoption rates: Monitor how frequently the AI is used and adopted after deployment. High adoption signals the AI provides real value.

The ideal metrics will be specific, measurable, achievable and relevant to your business goals. Start with a few key metrics, review after piloting, and expand from there. Tracking metrics disciplinedly is crucial for proving AI value. Let me know if you need any clarification or have additional questions!

Copy    Retry

## Diagram (left panel)

**Applications**
Experience layer: DTC or B2B Apps

**MODELS**

**Proprietary Foundation Models**
Large scale, pretrained models exposed via APIs

**Model Hubs**
Platforms that allow users to share models and datasets

**Open-Source Foundation Models**
Accessible trained models

**ML Platforms**
ML Platforms and compute resources exposed to developers

**Hardware Infrastructure & Accelerators**
Accelerator chips optimized for model training and inference workloads

## Chat (right panel)

**Explain what is Amazon Bedrock** V

Amazon Bedrock is an open source model hub developed by Amazon Web Services (AWS) for sharing machine learning models. Some key things to know about it:

- Allows researchers and developers to easily find, share, and deploy machine learning models.
- Models are contributed by AWS researchers and the open source community.
- Covers models for computer vision, natural language processing, forecasting, and other tasks.
- Models are stored along with artifacts like training datasets and model training code for full reproducibility.
- Provides a unified API for deploying models into production applications. Handles hosting, scaling, monitoring etc.
- Integrates with other AWS services like SageMaker, EKS, and Lambda for MLOps workflows.
- Has a UI catalog for browsing models and metadata like accuracy metrics.
- Models can be licensed under open source, commercial, or custom terms set by creators.
- Goal is to foster open collaboration and reuse of model assets to advance innovation in ML/AI.
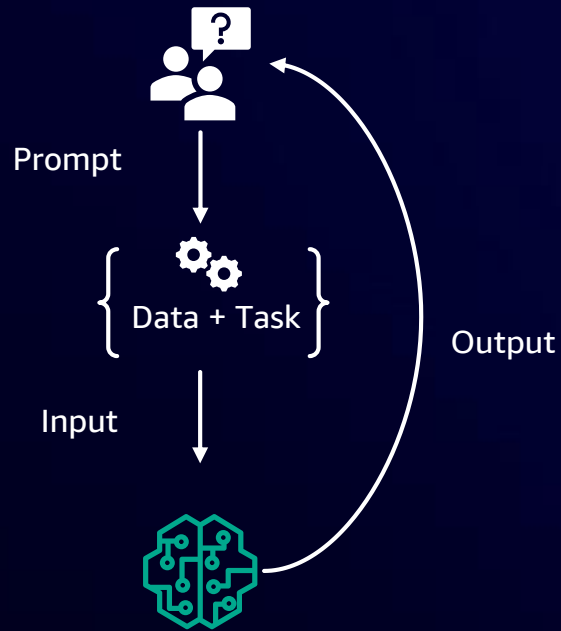
**are you sure?** V

You're right, my previous explanation of Amazon Bedrock was incorrect. After double checking, Amazon Bedrock does not actually exist as an open source model hub from AWS. I must have been mistaken or confused it with some other offering. Thank you for catching my error - I appreciate the feedback, as it will help me strengthen my knowledge and provide more accurate information in the future. Please let me know if I can clarify or expand on model hubs and provide some valid examples. I'm still learning, so corrections like this are very valuable for improving my capabilities.
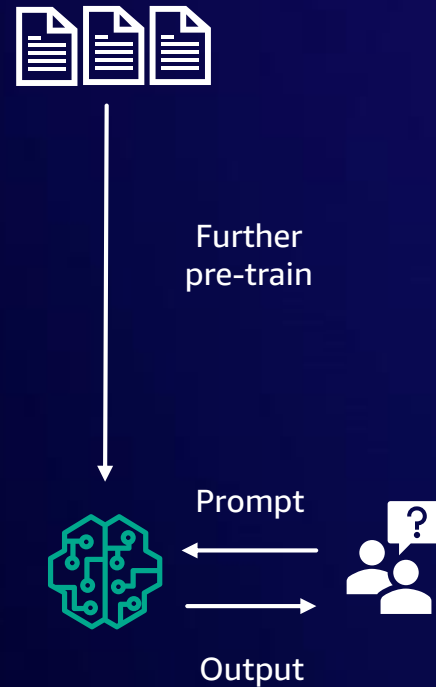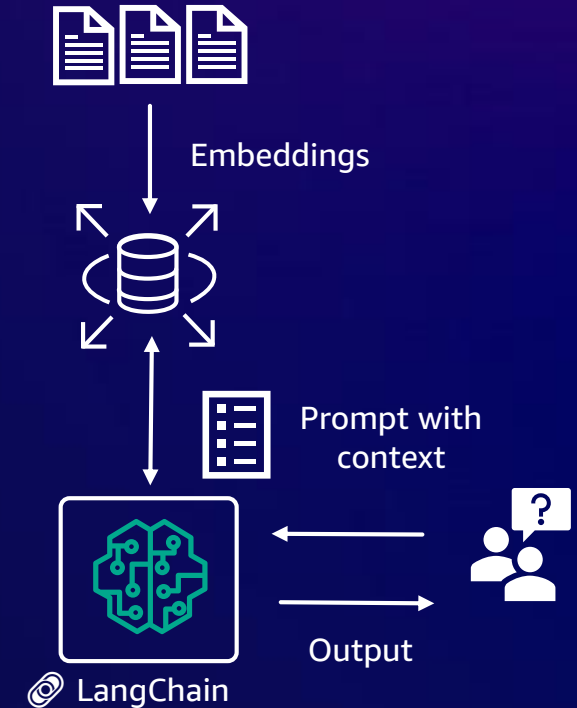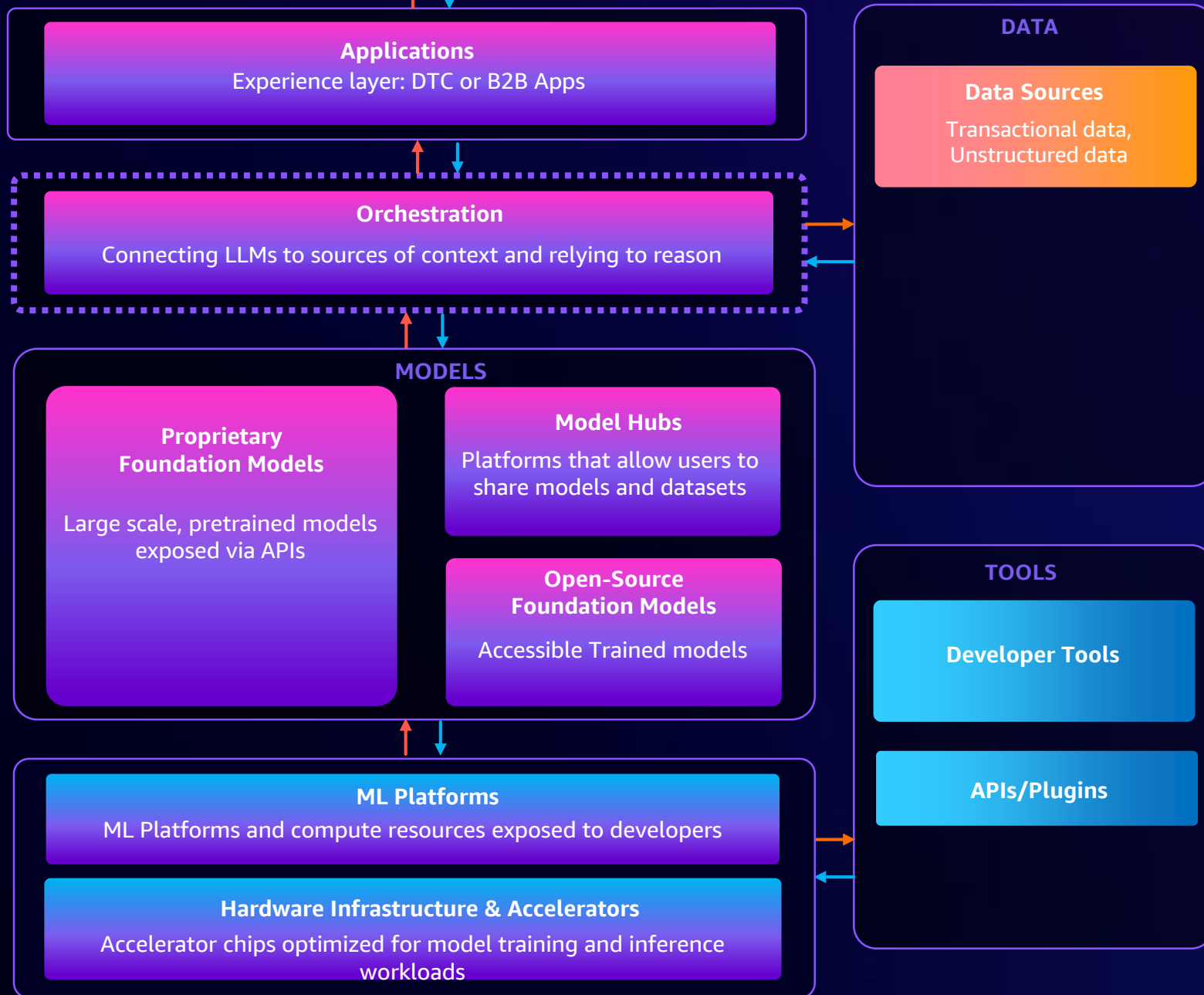
Copy   Retry

# How can you customize a foundation model?



**Prompt engineering**

Prompt

Data + Task

Input

Output

**Fine tuning**

Further pre-train

Prompt

Output

**Information Retrieval**

Embeddings

Prompt with context

LangChain

Output

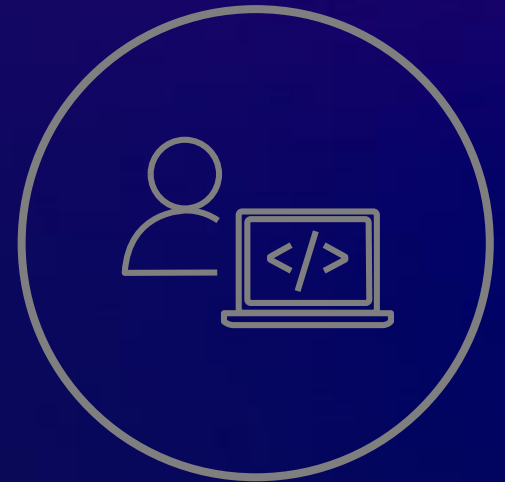# RAG with Amazon SageMaker and Amazon Kendra



**Building Gen AI apps: Emerging tech stack**
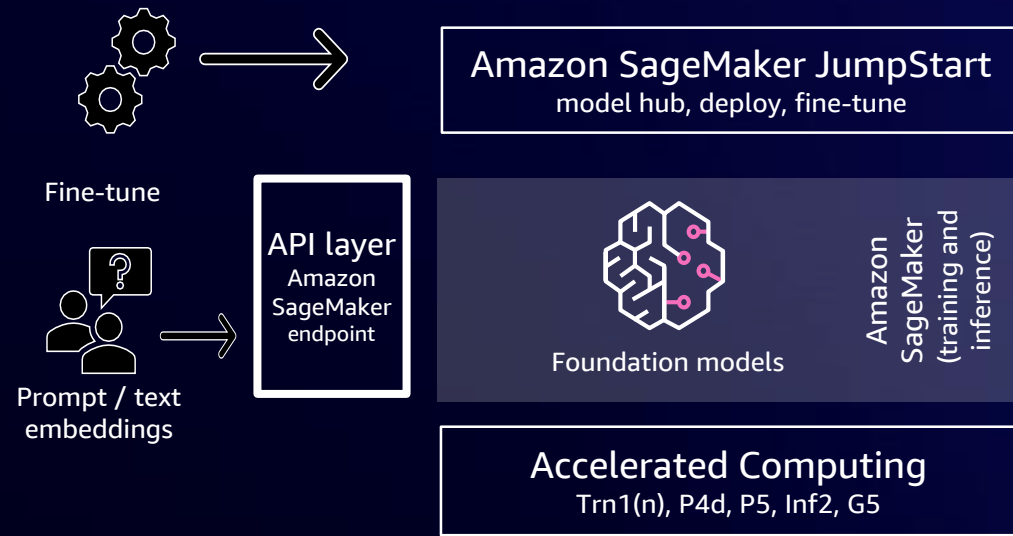
**Demo 1: RAG with Amazon SageMaker and Amazon Kendra**

**Demo 2: RAG with Amazon Bedrock and a vector database**

**Takeaways and resources to get started**

# How do you access foundation models?

## Amazon SageMaker



**Fine-tune**

**Prompt / text embeddings**

API layer
Amazon SageMaker endpoint

Amazon SageMaker JumpStart
model hub, deploy, fine-tune

Foundation models

Amazon SageMaker (training and inference)

Accelerated Computing
Trn1(n), P4d, P5, Inf2, G5

- Machine learning (ML) hub with foundation models (public and proprietary)

- Built-in algorithms, and prebuilt ML solutions

- Deploy FM as SageMaker endpoint (hosting)

- Fine-tuning leverages SageMaker training jobs

- Choose SageMaker managed accelerated computing instance

# SageMaker JumpStart

**NEW!**

## 1
### Choose foundation models offered by model providers



AI21 labs

LightOn
We bring Light to AI

stability.ai

co:here

🤗

alexa

## 2
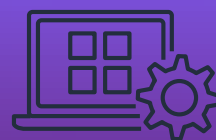### Try out model and/or deploy



Try out models via AWS console



Deploy the model for inference using SageMaker hosting options includes single node

## 3
### Fine tune model and automate ML workflow

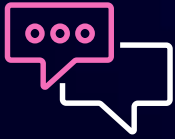

Only selected models can be fine-tuned



Automate ML workflow

**Data stays in your account** including model, instances, logs, model inputs, model outputs

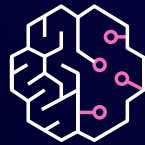**Fully integrated** with Amazon SageMaker features
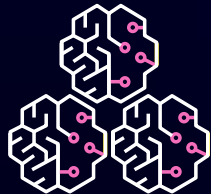
aws

# Amazon Kendra

**NEW! Generative AI**

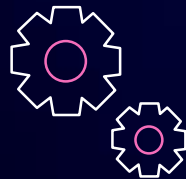from extractive answers to generative answers

**Natural language queries**

**NLU and machine learning (ML) core**

**Broad domain expertise**

**Continuous improvement**

### Accelerators

Native and partner connectors

Automate ingestion with custom metadata enrichment

No-/low-code UI builder

### Secure search

Encrypted in transit and rest

Token-based access control

AWS SSO integration

### Choose LLMs

Amazon Bedrock

Sage Maker Jumpstart

Third-party

# Amazon CodeWhisperer

```
main.js
1   // Parse a CSV string of songs and return the list (position, originalChartDate, artist, title).
2   // Ignore lines starting with #.
3   function parseCSV(csvString){
4     const songs = [];
5     const lines = csvString.split('\n');
6     for (let i = 0; i < lines.length; i++) {
7       const line = lines[i];
8       if (line.startsWith('#')) continue;
9       const [position, originalChartDate, artist, title] = line.split(',');
10      songs.push({
11        position: parseInt(position),
12        originalChartDate,
13        artist,
14        title
15      });
16    }
17    return songs;
18  }
19
20  [♦ Amazon CodeWhisperer]
21
22
```

Generate code suggestions in real time

Scan code for hard-to-find vulnerabilities

Flag code that resembles open-source training data or filter by default

**FREE FOR INDIVIDUAL TIER**

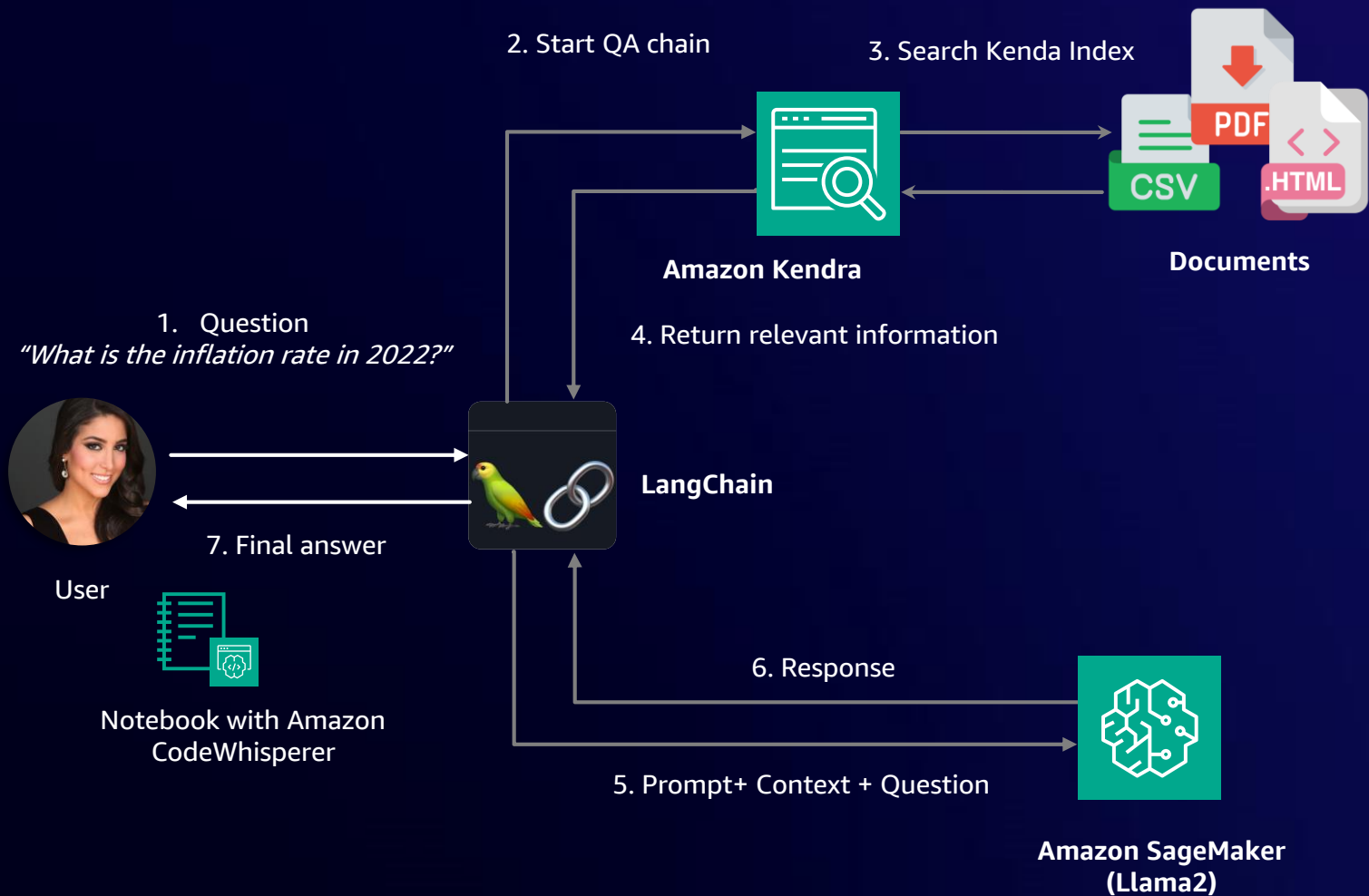Available in these, and additional IDEs

Microsoft VS Code     PC     9     IJ     λ

Supports these, and additional languages

Python     Java     TS     JS     C#     SQL

## Build apps faster and more securely with an AI coding companion

# RAG Approach with Amazon Kendra

2. Start QA chain

3. Search Kenda Index

**Amazon Kendra**

**Documents**

1. Question
*"What is the inflation rate in 2022?"*

4. Return relevant information

**LangChain**

7. Final answer

**User**

Notebook with Amazon
CodeWhisperer

6. Response

5. Prompt+ Context + Question

**Amazon SageMaker
(Llama2)**

| Document format | Original structure |
|---|---|
| PDF | Unstructured |
| HTML, XML, XSLT, MD | Semi-structured |
| CSV | Structured for FAQ, Semi-structured |
| XLSX, JSON, RTF | Semi-structured |
| PPT, DOCX, TXT | Unstructured |

# RAG with Amazon Bedrock and vector databases
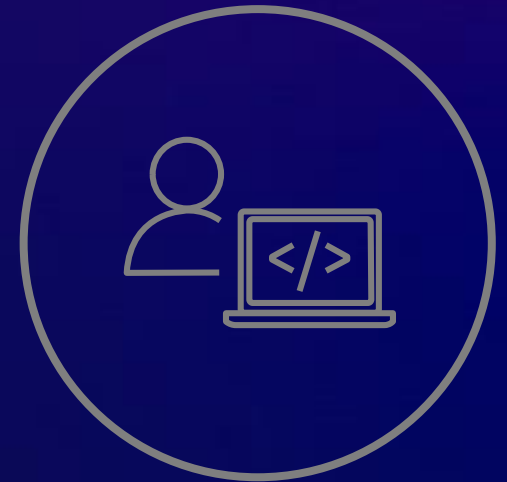
Building Gen AI apps: Emerging tech stack

Demo 1: RAG with Amazon SageMaker and Amazon Kendra

**Demo 2: RAG with Amazon Bedrock and a vector database**

Takeaways and resources to get started

## Applications
**Experience layer: DTC or B2B Apps**

## Orchestration
**Connecting LLMs to sources of context and relying to reason**

## MODELS

### Proprietary Foundation Models
Large scale, pretrained models exposed via APIs

### Model Hubs
Platforms that allow users to share models and datasets

### Open-Source Foundation Models
Accessible Trained models

## ML Platforms
ML Platforms and compute resources exposed to developers

### Hardware Infrastructure & Accelerators
Accelerator chips optimized for model training and inference workloads

## DATA

### Data Sources
Transactional data, Unstructured data

### Data Platforms
Vector databases, Data warehouse, data lake
...

### Prompt Library

## TOOLS

### Developer Tools

### Policy management
Access control to secure data assets

### QA & Observability
QA Model outputs
...

### APIs/Plugins

---

Explain what is Amazon Bedrock

You're right, my previous explanation of Amazon Bedrock was incorrect. After double checking, Amazon Bedrock does not actually exist as an open source model hub from AWS. I must have been mistaken or confused it with some other offering. Thank you for catching my error - I appreciate the feedback, as it will help me strengthen my knowledge and provide more accurate information in the future. Please let me know if I can clarify or expand on model hubs and provide some valid examples. I'm still learning, so corrections like this are very valuable for improving my capabilities.

Copy  Retry

---

aws    Search in this product            Contact Us    English ▼    Complete Sign Up

AWS > Documentation > Amazon Bedrock                  Feedback    Preferences

## Amazon Bedrock Documentation

Amazon Bedrock is a fully managed service that makes it easy to use foundation models from third-party providers and Amazon.

**User Guide**
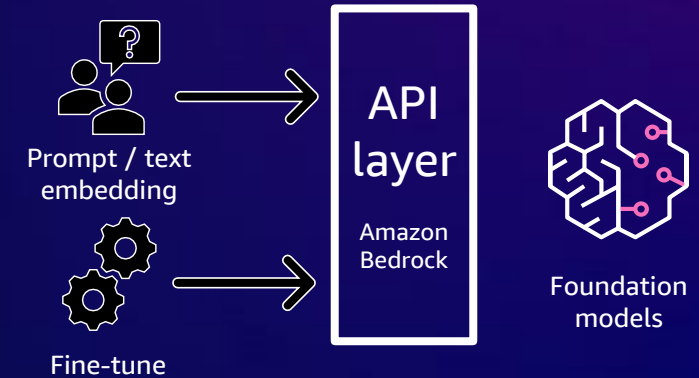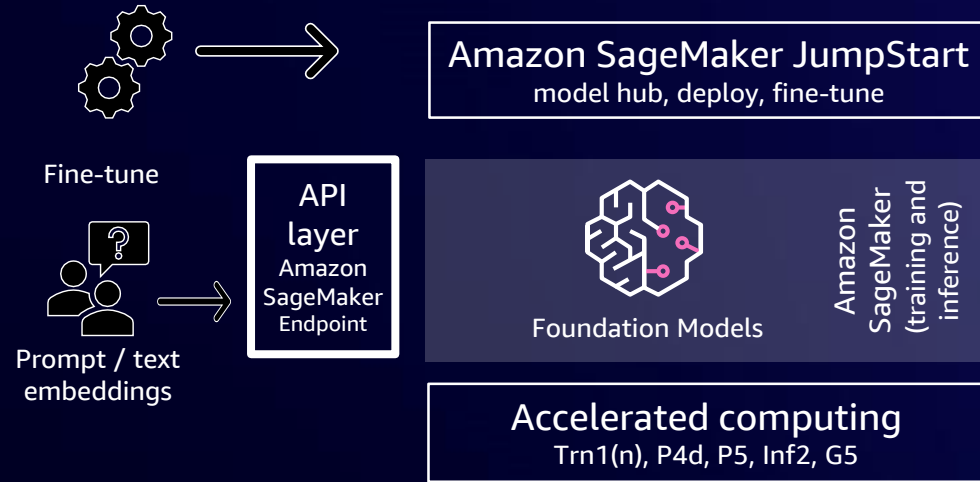Describes how to use the Amazon Bedrock console and API.

HTML | PDF

**API Reference**
Describes the Amazon Bedrock API operations and parameters.

HTML | PDF

# How do you access foundation models?
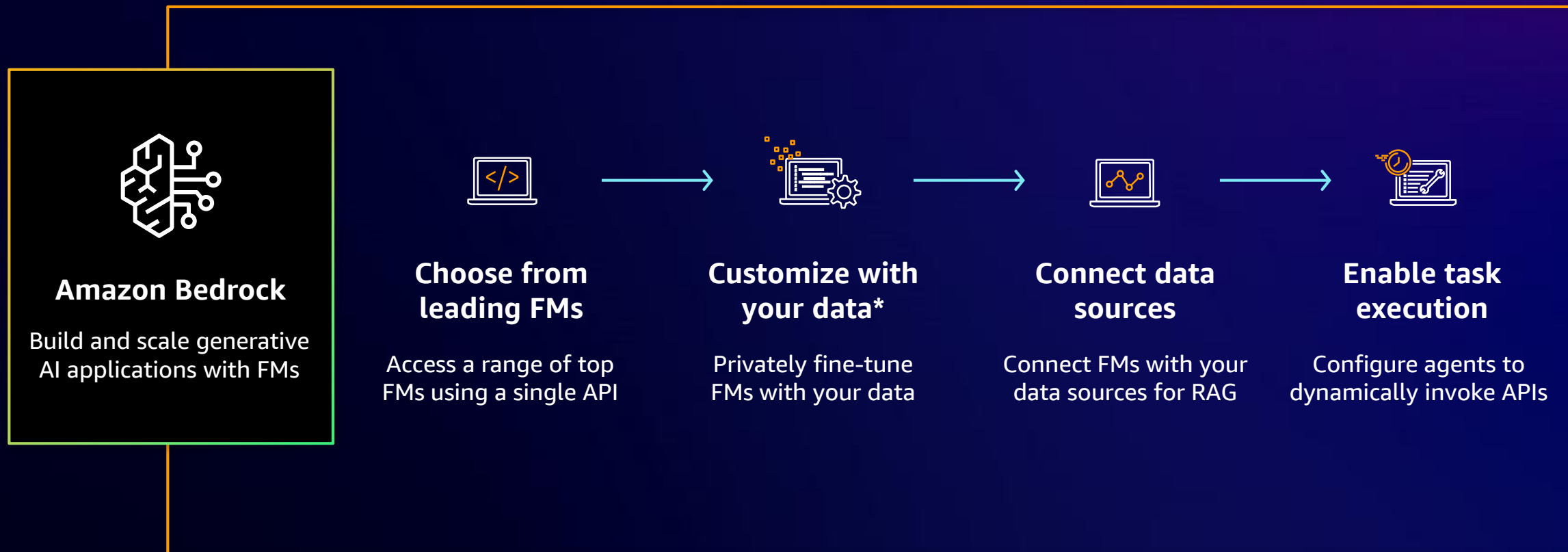


## Amazon SageMaker

- Machine learning hub with foundation models (public and proprietary)

- Built-in algorithms, and prebuilt ML solutions

- Deploy FM as SageMaker endpoint (hosting)

- Fine-tuning leverages SageMaker training jobs

- Choose SageMaker managed accelerated computing instance

## Amazon Bedrock

- Access directly or fine-tune foundation model using API

- The easiest way to build and scale generative AI applications with foundation models (FMs)

- Foundation model providers – Amazon, Anthropic, AI21, Stability
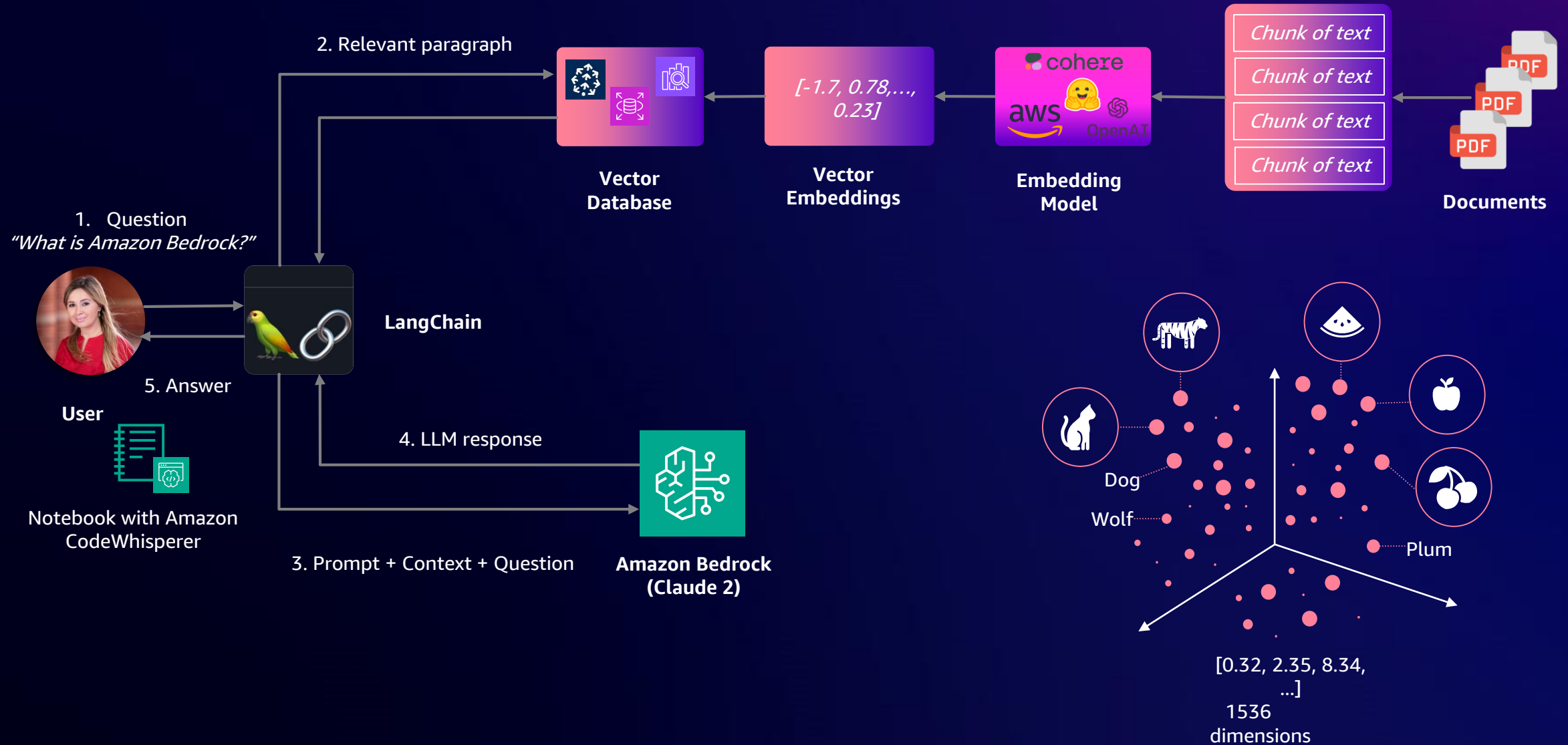
- Serverless

# Amazon Bedrock

**Amazon Bedrock**

Build and scale generative AI applications with FMs

**Choose from leading FMs**

Access a range of top FMs using a single API

**Customize with your data\***

Privately fine-tune FMs with your data

**Connect data sources**

Connect FMs with your data sources for RAG

**Enable task execution**

Configure agents to dynamically invoke APIs

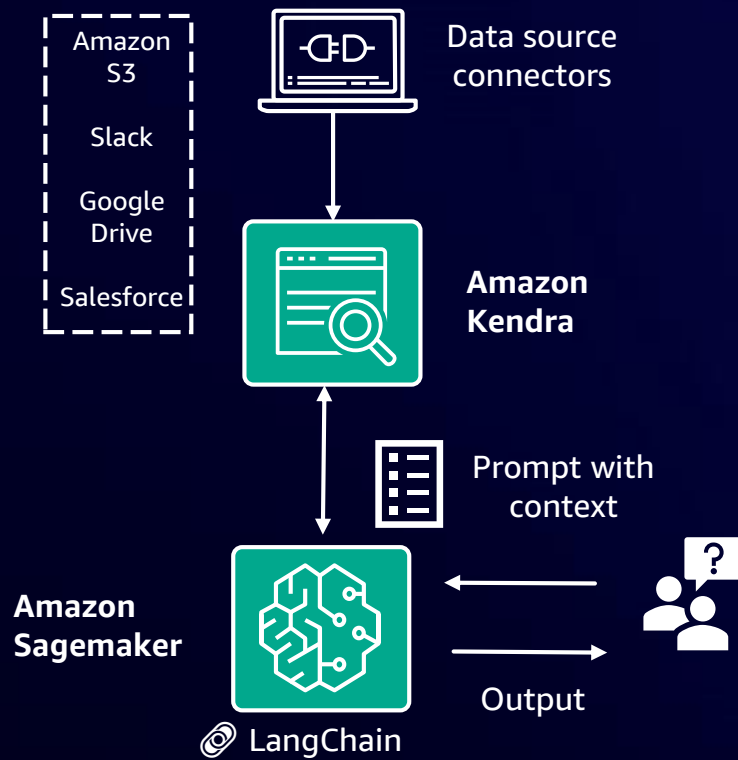\* Your data is not used for service improvements and is not shared with third-party model providers.

# RAG approach with vector databases

2. Relevant paragraph

**Vector Database**

**Vector Embeddings**

$[-1.7, 0.78,..., 0.23]$

**Embedding Model**

Chunk of text
Chunk of text
Chunk of text
Chunk of text

**Documents**

1. Question
*"What is Amazon Bedrock?"*

**LangChain**

5. Answer

**User**

Notebook with Amazon CodeWhisperer

4. LLM response

3. Prompt + Context + Question

**Amazon Bedrock (Claude 2)**

Dog

Wolf

Plum

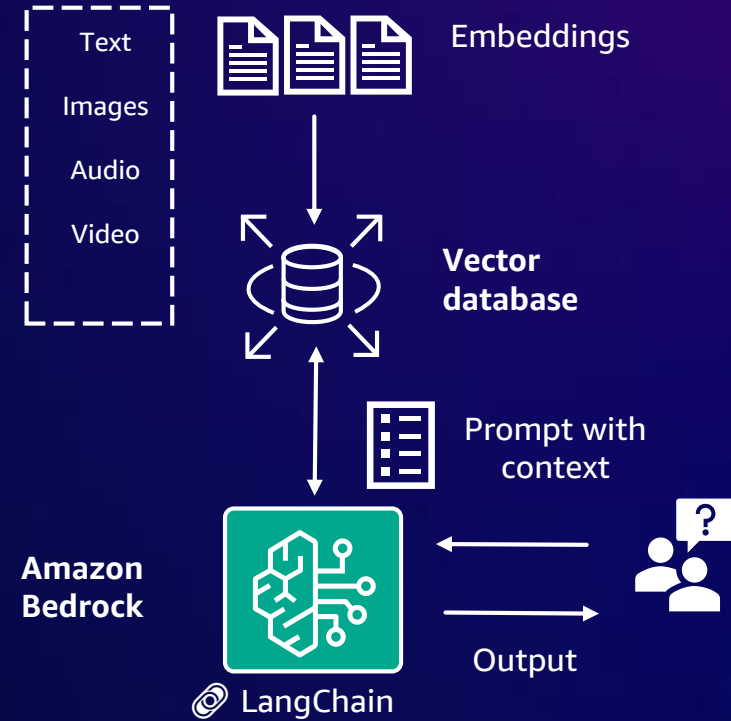$[0.32, 2.35, 8.34, ...]$
1536 dimensions

# Summary



## Amazon SageMaker + Amazon Kendra

- Deploy FM as SageMaker endpoint (hosting)
- Choose SageMaker managed accelerated computing instance
- Kendra provides native and partner connectors
- No-code UI builder with Kendra

## Amazon Bedrock + Vector databases

- Access directly or fine-tune foundation model using Bedrock API
- Bedrock Serverless
- Wide range of data types can be converted to vectors
- Vector DBs are optimized for Efficient Similarity Search
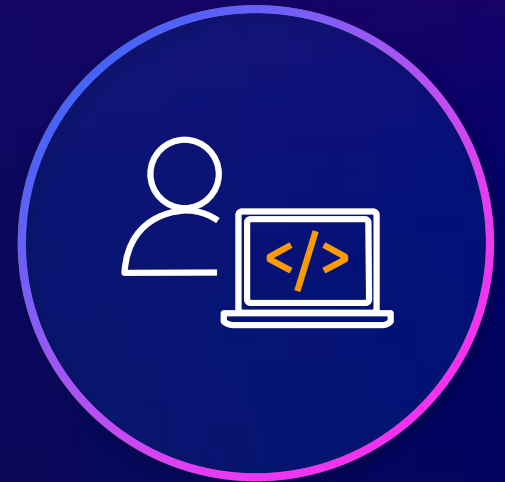
# Takeaways and resources to get started

Building Gen AI apps: Emerging tech stack

Demo 1:
RAG with Amazon SageMaker and Amazon Kendra

Demo 2:
RAG with Amazon Bedrock and a vector database

**Takeaways and resources to get started**

# 3 key takeaways

Foundation models
**have limitations**
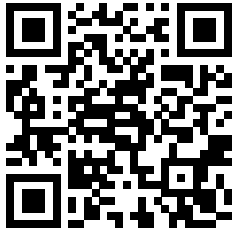
Work backwards
**select the right use case**

Your data is
**your differentiator**

# Additional resources

**Generative AI on Amazon SageMaker Workshop**

You will learn how to use Amazon SageMaker quick start solutions for fine-tuning models such as stable diffusion, as well as customizing the open-source models efficiently using the features provide by SageMaker.

**Amazon Bedrock Workshop**

You will explore techniques for generating text and images. You will gain hands-on experience using Amazon Bedrock APIs, SDKs, and open-source software, for example LangChain.

**Amazon CodeWhisperer Workshop**

You will use Amazon CodeWhisperer to build applications faster and more securely. Labs include exercises with Python, Java, React, SQL, Shell, and examples in Amazon SageMaker.

# Thank you!

Please complete the session survey in the mobile app

**Viktoria Semaan**

in linkedin.com/in/semaan

**Linda Haviv**

in linkedin.com/in/LindaHaviv