



ACOT201

SalesGenie

Leveraging generative AI for assisting life insurance sales agents on the field

Sanjay Thawakar

Senior Vice President & Head of AI Works & BPMA
Max Life Insurance



Senthil Jeyachandran

AI & ML Specialist, Solutions Architect
AWS India

Agenda



- Generative AI landscape within life insurance industry
- Maxlife SalesGenie – Solving sales distribution challenges
- Solution development, key elements and learnings
- Solution architecture
- Achievements and way-forward
- Generative AI overview
- Amazon Bedrock
- Customization, Knowledge base, Agents for Bedrock, Guardrails
- Demo



Generative AI landscape within Insurance

Function/ Category*	Sales	Marketing	HR	Operations	Customer Service	Digital, Ecom, AIW	Finance, Legal, & Others
Content Translation & Classification	<ul style="list-style-type: none"> Translating sales script & material to language of choice 	<ul style="list-style-type: none"> Vernacular Campaign design 	<ul style="list-style-type: none"> Resume Screening 		<ul style="list-style-type: none"> Script translation 	<ul style="list-style-type: none"> Code Refactoring 	
Content Creation	<ul style="list-style-type: none"> Generating Sales Script 	<ul style="list-style-type: none"> KFD video creation Seller activation and buzz Press drafts Celebrity image generation CEO Webcast script generation 	<ul style="list-style-type: none"> Interview Questions Creating Job Descriptions 	<ul style="list-style-type: none"> Renewal KFD Renewal Pitch Generation 	<ul style="list-style-type: none"> Co-pilots for agents to generate responses on customer queries 	<ul style="list-style-type: none"> Code Generator Creating scripts to test code SEO Meta Description Creation 	<ul style="list-style-type: none"> Boiler plate contract generator
Content Clarification	<ul style="list-style-type: none"> Calls transcripts & summarization 	<ul style="list-style-type: none"> Analysis for Annual Reports 	<ul style="list-style-type: none"> Policy Summary 	<ul style="list-style-type: none"> Condense Lengthy Document into Brief Summaries 	<ul style="list-style-type: none"> Calls transcripts & summarization 	<ul style="list-style-type: none"> Annotations to Code Bug Detection 	
Content Analysis, Design & Summarization	<ul style="list-style-type: none"> Product Information 	<ul style="list-style-type: none"> Personalised Product EDM 	<ul style="list-style-type: none"> Interview feedback Candidate Ranking Employee Engagement Survey Analysis 	<ul style="list-style-type: none"> Search Internal Documents, Optimize, & Recommendations 	<ul style="list-style-type: none"> Sentiment Analysis & response generator 	<ul style="list-style-type: none"> Code Optimization 	<ul style="list-style-type: none"> Contract Summarization Financial Report Summarization
Act basis Instruction	<ul style="list-style-type: none"> Agent Training Bots Virtual Sales Agent Virtual Supervisor 	<ul style="list-style-type: none"> Automate ad targeting 	<ul style="list-style-type: none"> Virtual Recruiter Interview Scheduling, processing & summarizing 	<ul style="list-style-type: none"> UW Twin Renewal Tele calling assistant Perform Data Entry 	<ul style="list-style-type: none"> Email & Chat Bot Product feature summarization for CS Support Virtual CS Agent 	<ul style="list-style-type: none"> Coding co-pilot 	<ul style="list-style-type: none"> Investment Guide

Low

Business Team Enablement

Complexity

Enterprise Grade Solution

High

SalesGenie: Solving sales distribution challenges



Low trainer to
FLS ratio

Limited
Supervisors

Retention of
knowledge

Staying
updated with
current affairs

Training new force

Disparate skill of
sellers

1 Lac field sales force. High %ge new force
every year



Sales Pitch Assistance

In-depth MLI product information and comparison



Competitor product differentiation



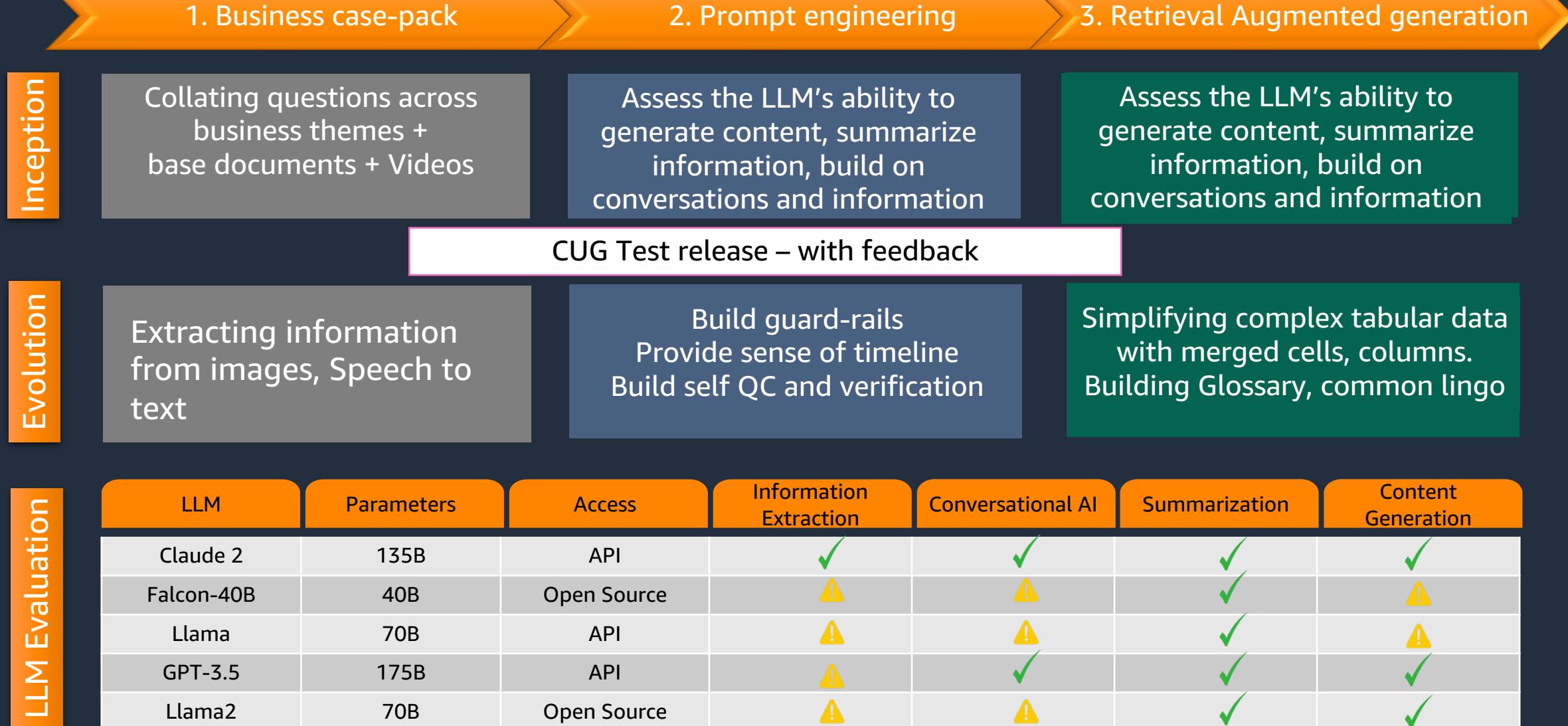
Objection Handling



**Awareness of Financial instrument
and Current Affairs**



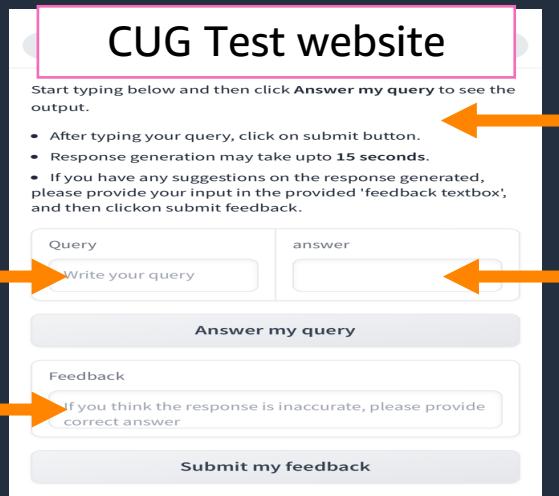
Evolution of solution development



Key elements and learnings

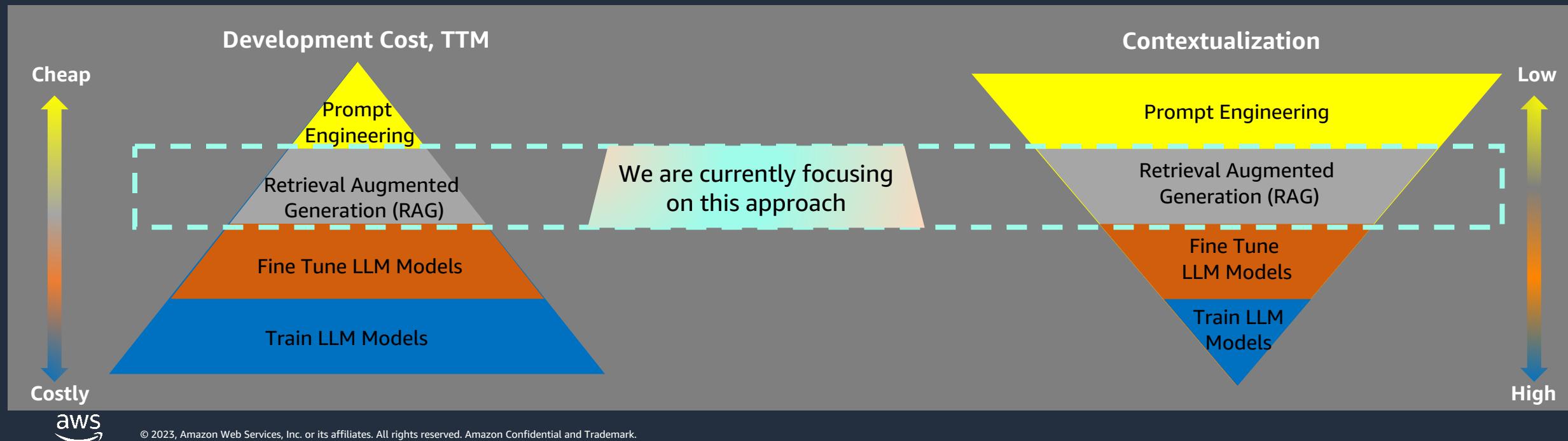
Placeholder to ask query

Gather feedback on generated response



Instructions for new users – How to use and What to expect

Response generated from intelligence engine



Considerations for developing LLM apps

Fraud - jailbreak

- ✓ Strict persona
- ✓ Verify content match to source data
- ✓ Condition user query

Bias - Hallucination

- ✓ Lowest temperature for generation
- ✓ Guard-rails for fact check
- ✓ Verify calculations / inference

Sensitive data leaks

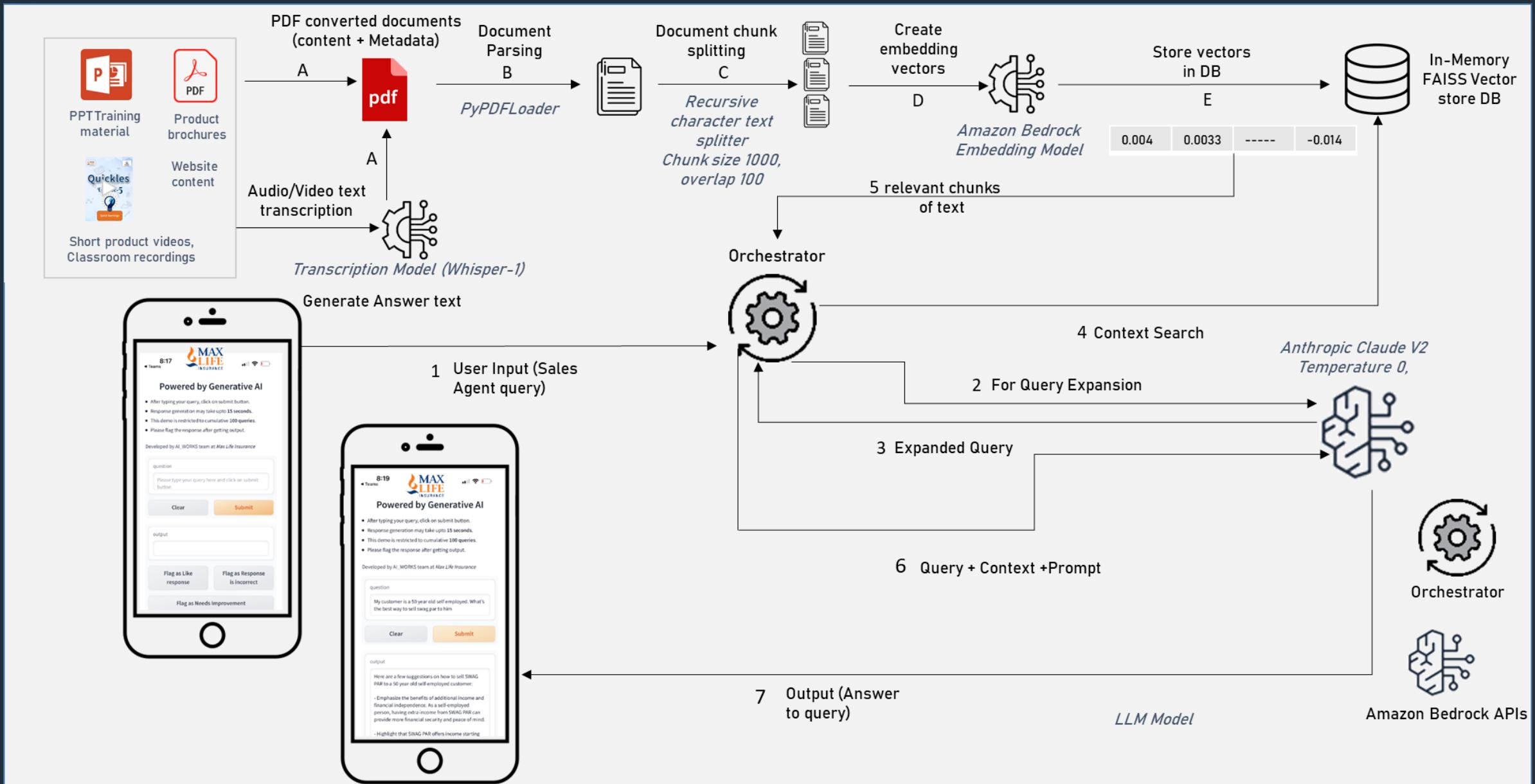
- ✓ Using private models
- ✓ Using private embedding model

Stability of response

- ✓ Using private instances
- ✓ Model version tracking



Solution Architecture

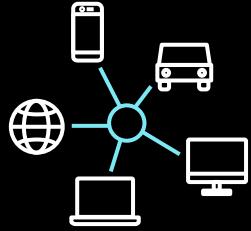


Test results and way forward

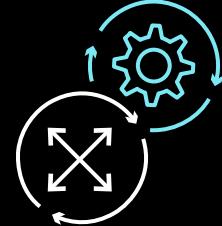
- Amazon bedrock titan text embedding + Anthropic's Claude-V2 **outperformed** GPT3.5-turbo + Ada text embedding
- **Simplifying** complex data tables key to improved response generation
- Maintaining knowledgebase and updating knowledge retrieval systems to fetch latest data is key challenge to future **sustenance**.
- Ability to handle new products, lingo and dialects remain a **key challenge** when they are presented to LLMs. Ex-SWAG product

Question Category	# Questions	Claude V2		GPT 3.5 Turbo	
		Correct Answers as expected	% Accuracy	Correct Answers as Expected	% Accuracy
Sales and Pitching	216	210	97%	210	97%
Products Related (Recommendations and variant information)	150	108	72%	91	60%
Competition comparisons	126	114	90%	54	43%
Financial knowledge	270	264	98%	264	98%
Sales pitching, Conversation Starters and Objection Handling	174	168	97%	168	97%
Total	936	864	92%	787	84%

What generative AI customers are asking for



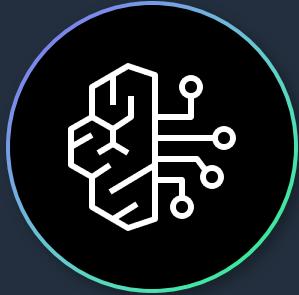
**Which model
should I use?**



**How can I
move quickly?**



**How can I keep
my data secure
and private?**



Amazon Bedrock

The easiest way to build and scale generative AI applications with foundation models

Choice of leading FMs via single API

Model customization

Retrieval Augmented Generation (RAG)

Agents that execute multistep tasks

Security, privacy, and safety

Amazon **Bedrock**

Broad choice of models

AI21labs

amazon

ANTHROPIC

cohere

Meta

stability.ai

Jurassic-2 Ultra

Jurassic-2 Mid

Titan Text Embeddings

Titan Multimodal Embeddings

Titan Text Lite

Titan Text Express

Titan Image Generator

Claude 2

Claude 2.1

Claude Instant

Command + Embed

Cohere Command Light

Cohere Embed English

Cohere Embed Multilingual

Llama 2

Llama 2 13B

Llama 2 70B

Stable Diffusion XL1.0



Why customize?



Adapt to domain-specific language

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



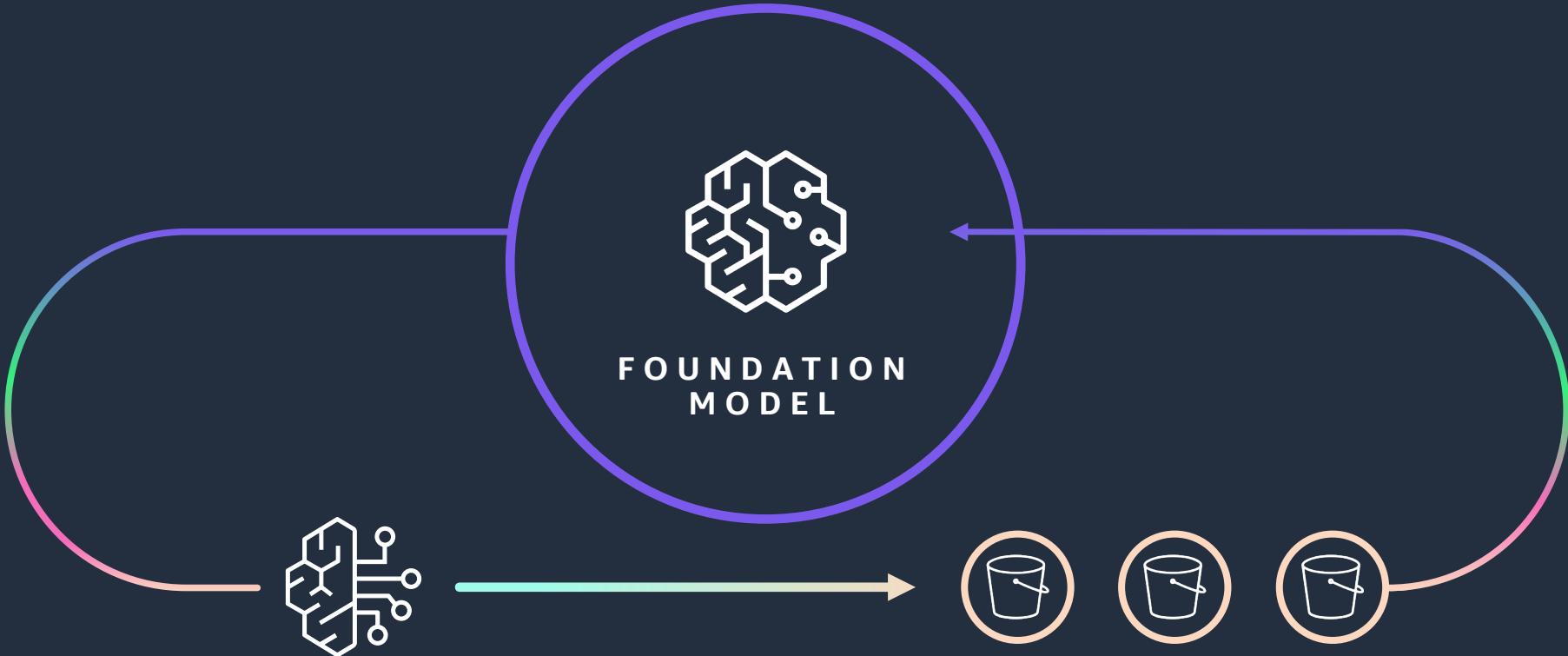
Enhance performance for specific tasks

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



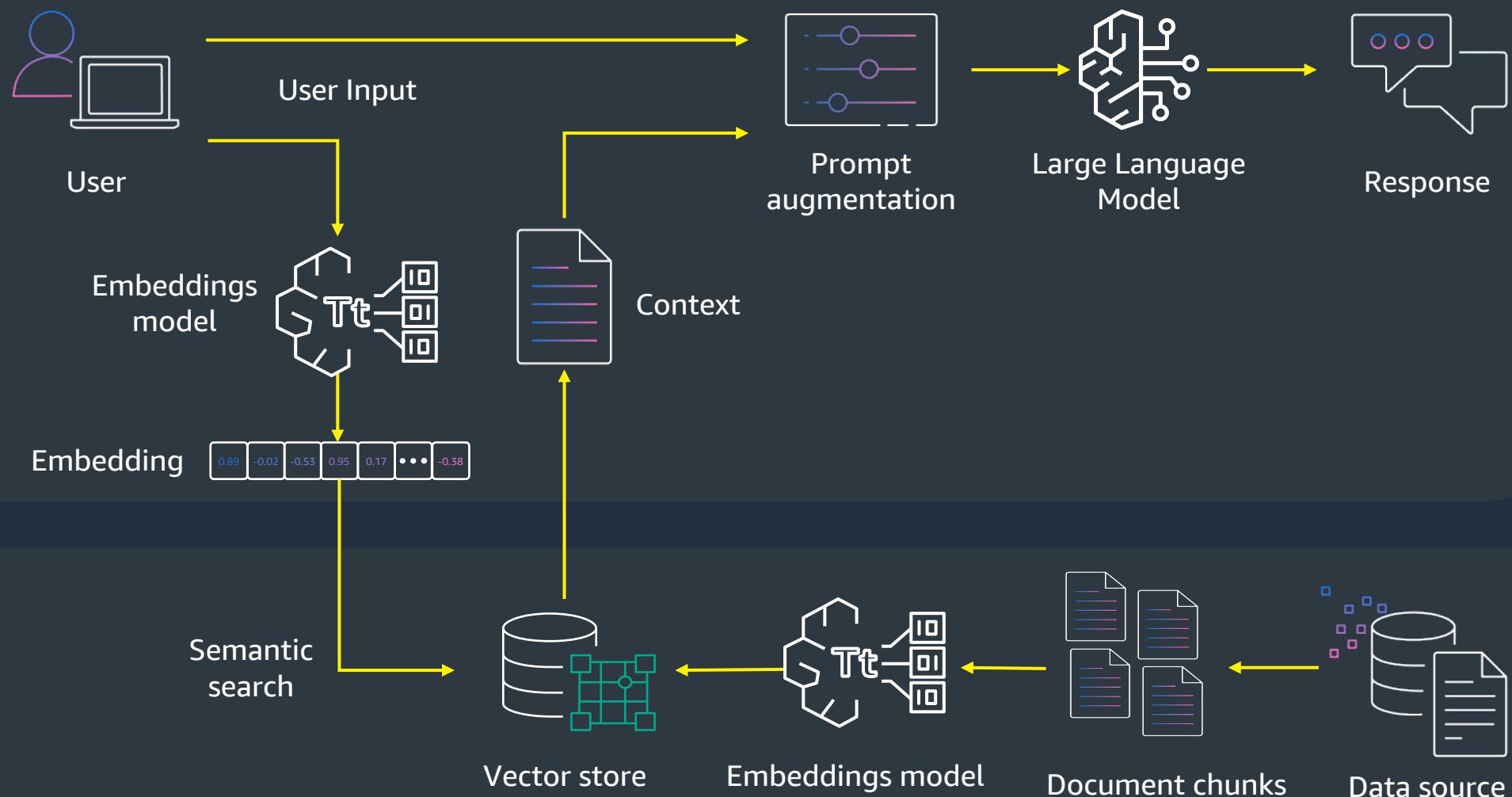
Fine tuning

Retrieval Augmented
Generation (RAG)

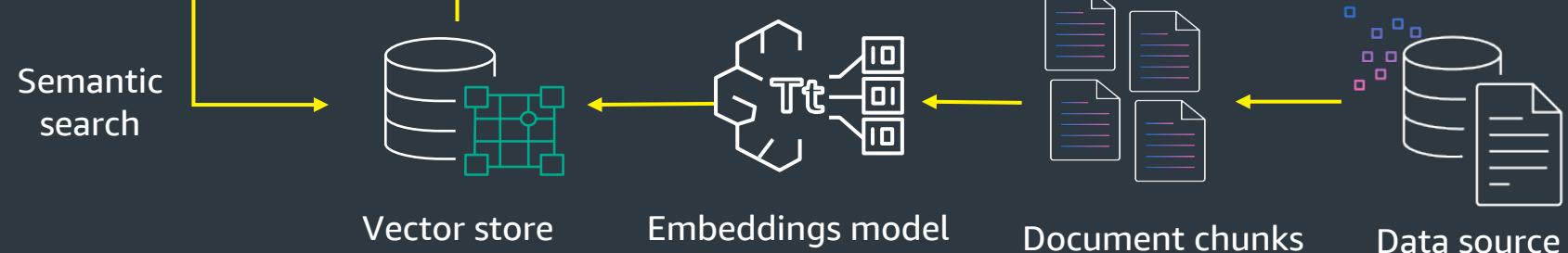
Continued
Pre-training

RAG in Action

Text Generation Workflow



Data Ingestion Workflow



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

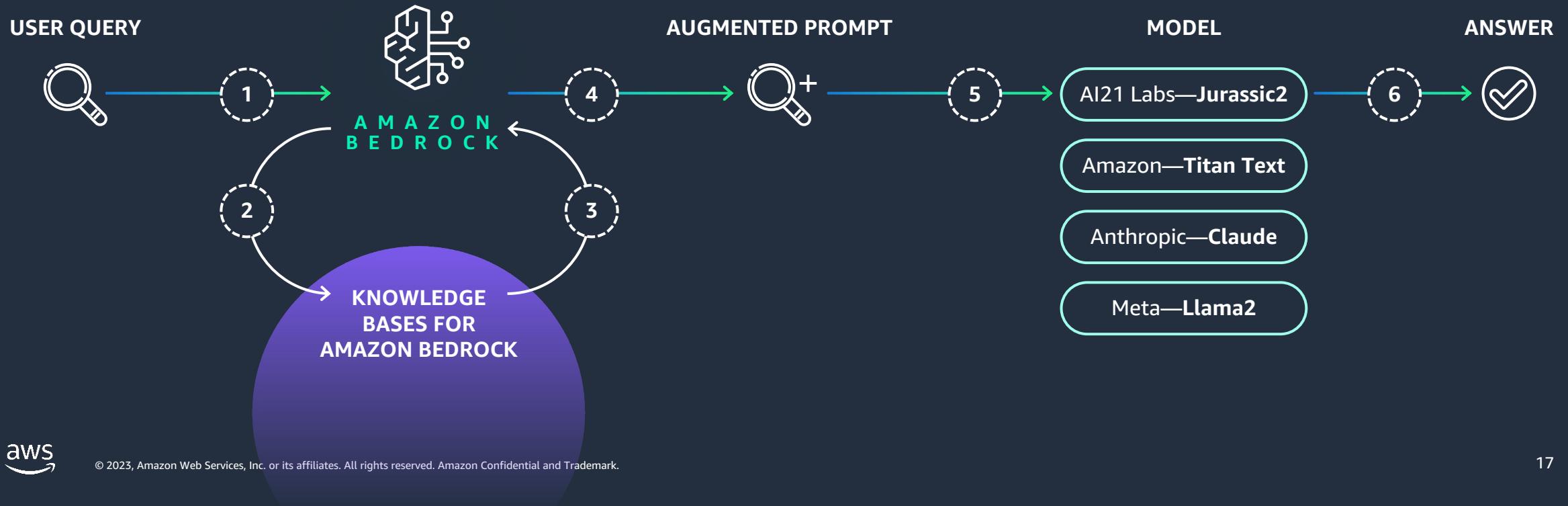


Orchestration

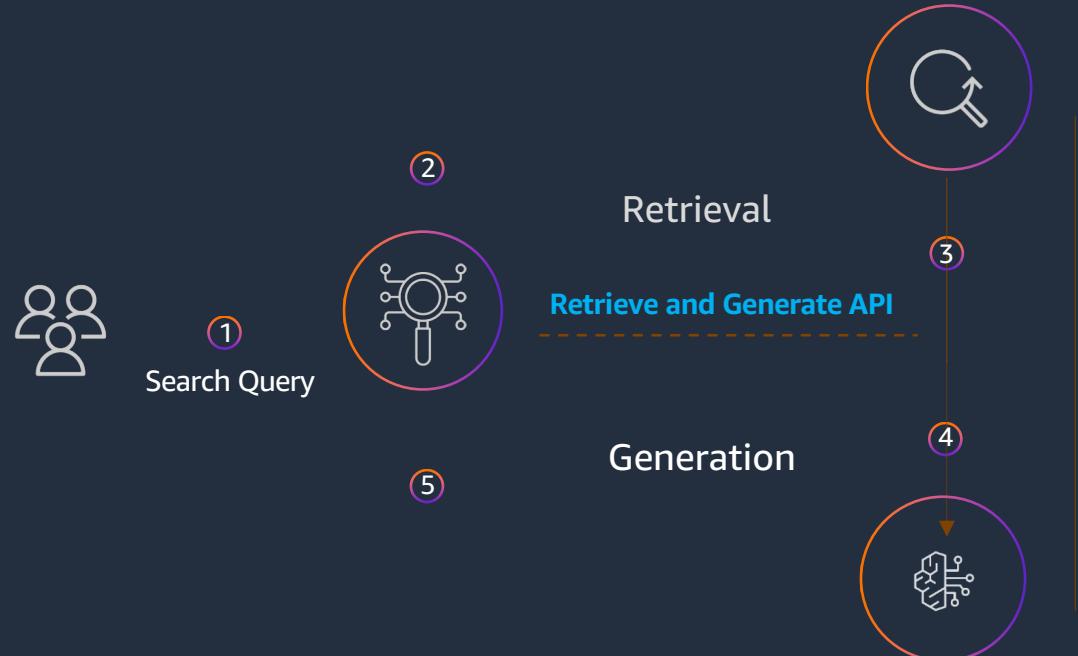
Knowledge bases for Amazon Bedrock

NATIVE SUPPORT FOR RETRIEVAL AUGMENTED GENERATION (RAG)

- Securely connect FMs to data sources for RAG to deliver more relevant responses
- Fully managed RAG workflow including ingestion, retrieval, and augmentation
- Built-in session context management for multi-turn conversations
- Automatic citations with retrievals to improve transparency



Retrieve and Generate API



Single API	One API for RAG
Comprehensive	Multiple knowledge-bases
Session	Built-in session management

Retrieve and generate API will enable a simplified RAG solution

Vector databases for **Amazon Bedrock**



Vector Engine For
Amazon OpenSearch
Serverless



Redis Enterprise
Cloud



Pinecone

COMING SOON



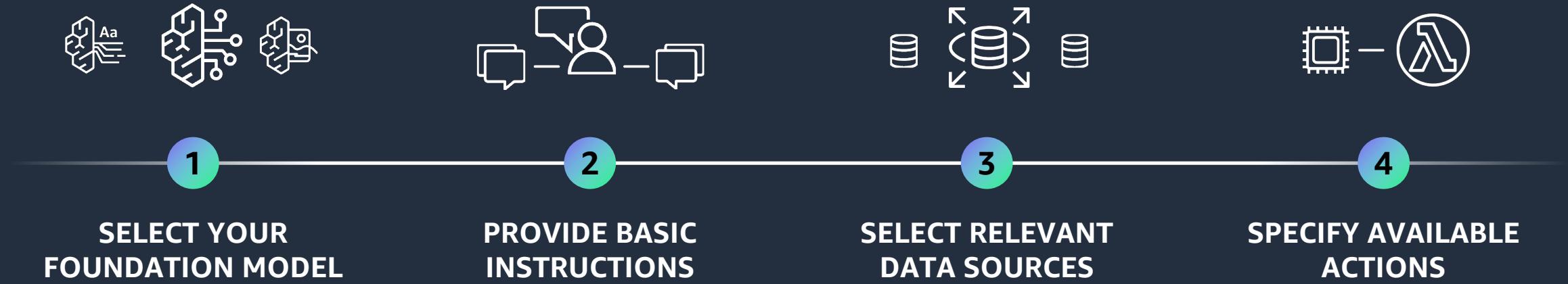
Amazon
Aurora



MongoDB

Agents for Amazon Bedrock

ENABLE GENERATIVE AI APPLICATIONS TO EXECUTE MULTISTEP TASKS USING COMPANY SYSTEMS AND DATA SOURCES



| Breaks down and orchestrates tasks |

| Securely accesses and retrieves company data for RAG |

| Takes action by invoking API calls on your behalf |

| Chain-of-thought trace and ability to modify agent prompts |

Guardrails for Amazon Bedrock

IMPLEMENT SAFEGUARDS TAILORED TO YOUR APPLICATION REQUIREMENTS AND RESPONSIBLE AI POLICIES

Preview

Apply guardrails consistently across FMs including fine-tuned models and agents

Configure filtering of harmful content and topics to avoid based on your responsible AI policies

Redact personally identifiable information (coming soon)

The screenshot shows the Amazon Bedrock Guardrails configuration interface for a 'Working draft' named 'antje-banking-assistant'. The interface is divided into several sections:

- Denied topics (1):** A table with one entry: 'Investment advice' (highlighted with a red box). The instructions state: 'Investment advice refers to guidance or recommendations provided by a financial professional, adv...'.
- Content moderation: filter strengths:** A table comparing prompt and response filter settings.

Prompt filters	Response filters
ON	ON
Toxicity filter strength for prompts	Toxicity filter strength for responses
High	High
Insults filter strength for prompts	Insults filter strength for responses
High	High
Sexual filter strength for prompts	Sexual filter strength for responses
High	High
Violence filter strength for prompts	Violence filter strength for responses
High	High
- Default responses:** A table showing blocked prompts and responses.

Blocked prompts	Blocked responses
Sorry, I can't comment on that.	Sorry, I can't comment on that.
- Test:** Shows a test session with 'Claude Instant v1.2' (ODT) and a prompt: 'Should I open a credit card account?'. The response is: 'Here are a few things to consider when deciding whether to open a credit card account:
 - Having a credit card and using it responsibly can help you establish credit history. This is important for things like qualifying for loans in the future. However, be sure you can pay the bill in full each month to avoid interest charges.'.
- Model response:** Shows the same response as the test session.
- Final response:** Shows the same response as the test session.
- Guardrail check:** Shows a green 'Passed' status with a 'View trace >' button and an orange 'Run' button. A red arrow points from the 'Passed' status to the 'View trace >' button.

Provisioned throughput

- Reserve throughput (input/output tokens per minute)
- Ensure consistent user experience during traffic spikes
- Purchase with commitment term of one month or six months
- Pay hourly rate, discounted for extended commitment



Amazon Bedrock > Provisioned throughput > Purchase provisioned throughput

Purchase provisioned throughput [Info](#)

Provisioned throughput details [Info](#)

Provisioned throughput name

Name can have up to 40 characters, and it must be unique. Valid characters A-Z, a-z, 0-9, and - (hyphen).

Select model

Tags - optional

Model units & commitment term [Info](#)
Select model units & commitment term to purchase Provisioned throughput. To estimate cost use [MU Estimator](#).

Model units
Please request the model units here before purchasing provisioned throughput. [AWS support center](#)

Select commitment term
Commitment terms locks the purchase for the selected duration.

Estimated purchase summary
To view the provisioned throughput pricing please visit [Pricing information](#)

Estimated hourly cost	Estimated daily cost	Estimated monthly cost
-	-	-

Edits to model and model units will be restricted
Once provisioned throughput is purchased, model units cannot be updated and the model can only be updated to another model with the same lineage.
[Learn more](#)

[Cancel](#) [Purchase Provisioned throughput](#)

Demo



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Start typing below and then click **Answer my query** to see the output.



- After typing your query, click on submit button.
- Response generation may take upto **15 seconds**.
- This demo is restricted to cumulative **100 queries**.
- If you have any suggestions on the response generated, please provide your input in the provided 'feedback textbox', and then click on submit feedback.

Query

answer

|Write your query

Answer my query

Feedback

If you think the response is inaccurate, please provide correct answer

Submit my feedback



Thank you!

Sanjay Thawakar

Senior Vice President & Head – AI
Works & BPMA

Max Life Insurance

Senthil Jeyachandran

AI & ML Specialist,
Solutions Architect
AWS India