

WORKSHOP

GenAI Bootcamp



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Agenda

Day 1 Feb 04th

09:30	Intro GenAI & Bedrock
11:30	Bedrock API + Multimodality
13:30	RAG Knowledge Bases
15:00	Bedrock Data Automation + QuickSuite
18:00	Afterwork

Day 2 Feb 05th

09:30	Agents Tools, MCP
11:30	Strands Agents + Multi agents
13:30	AI Security Secure chatbots
15:00	AgentCore Introduction
16:30	Workshop sessions Detailed scoping & Architecture

Day 3 Feb 06th

09:30	AgentCore + Eval Deep dive
11:30	AI coding Kiro, Spec Driven Dev
13:30	Closing session Tech FR lead keynote
13:30	Workshop sessions + Next steps

Build Phase
Feb 09th – Feb 20th

Intro to Generative AI (GenAI) and Amazon Bedrock

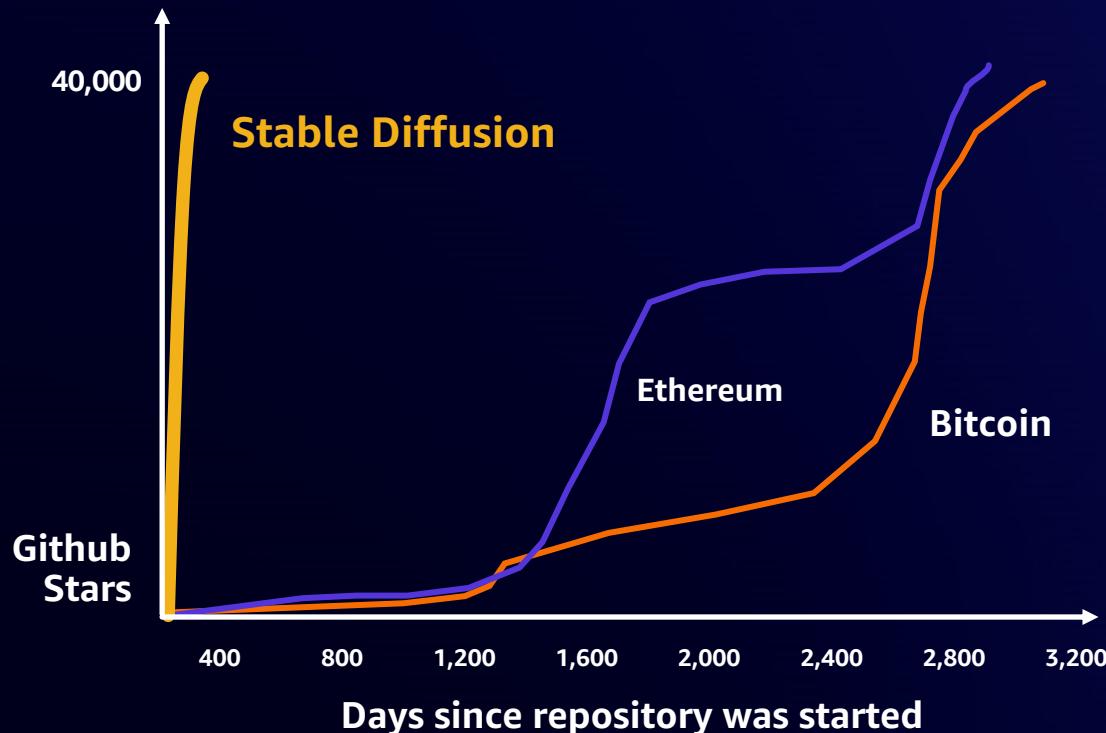


© 2026, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Generative AI is the fastest growing trend in AI

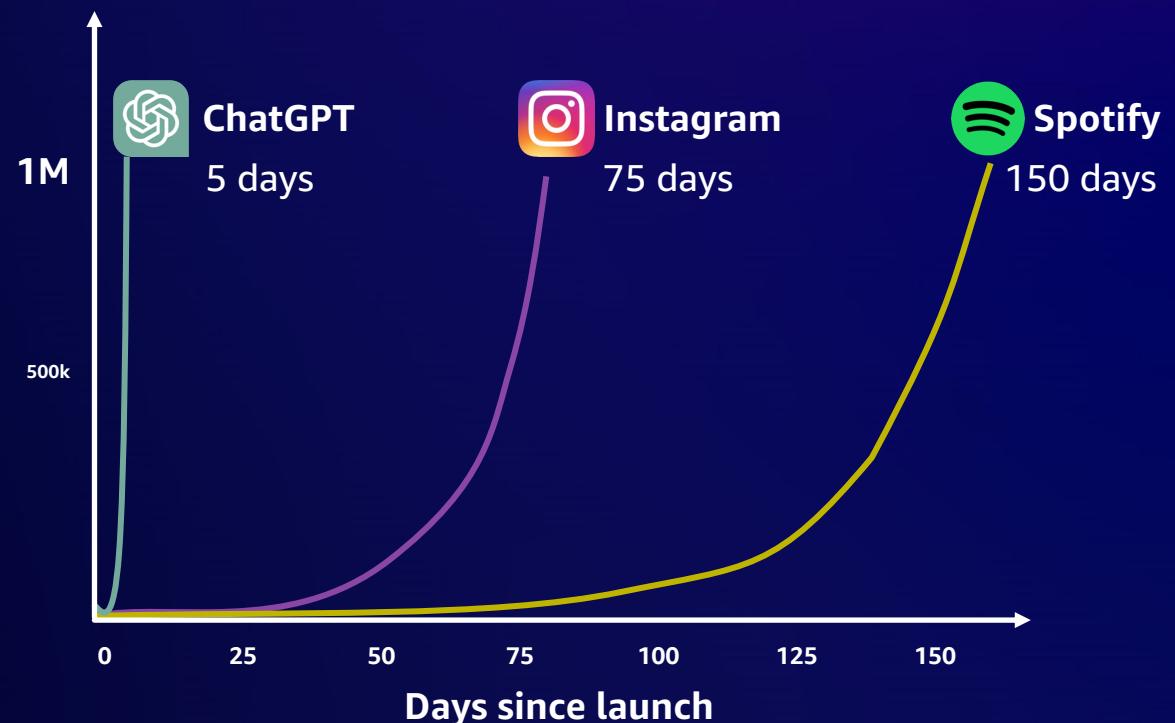
Developer adoption

Stable Diffusion accumulated 40k stars on GitHub in its first 90 days

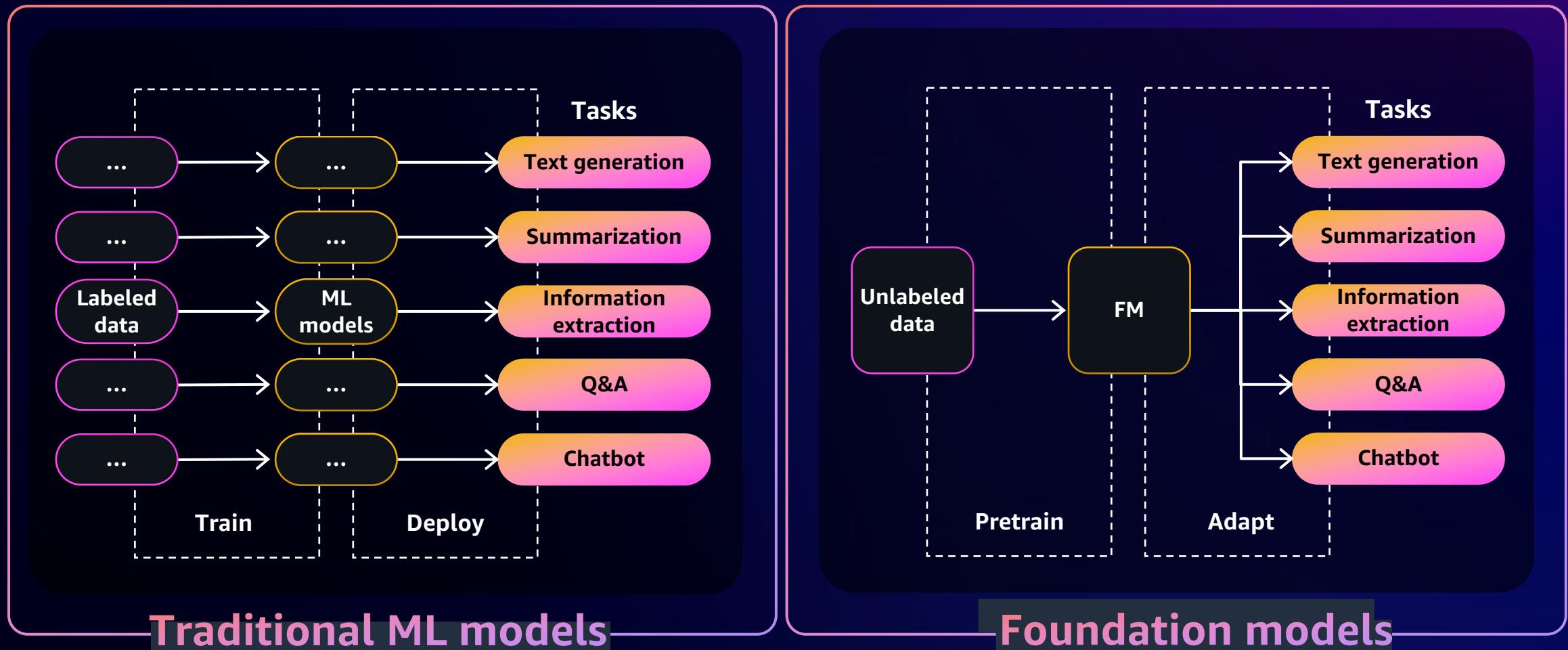


Consumer adoption

ChatGPT reached the 1 million users mark in just 5 days



Why foundation models?



What is Generative AI?



AI that can produce original content close enough to human generated content for real-world tasks



Powered by foundation models pre-trained on large sets of data



Can complete wide range of tasks without additional training

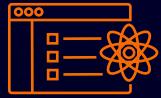


Applicable to many use cases like text summarization, question answering, digital art creation, code generation, etc.



Reduces time and cost to develop ML systems and innovate faster

Generative AI helps across a range of domains ...



Enhance customer experience

CHATBOTS

VIRTUAL ASSISTANTS

AI-POWERED CONTACT CENTER

PERSONALIZATION



Boost employee productivity

CONVERSATIONAL SEARCH

SUMMARIZATION

ENTITY EXTRACTION

CODE GENERATION

DATA TO INSIGHTS



Creativity and content creation

WRITING

MEDIA

DESIGN

MODELING



Improve business operations

DOCUMENT PROCESSING

PROCESS OPTIMIZATION

CYBERSECURITY

DATA AUGMENTATION

What is an LLM?

A Large Language Model is a massive, deep learning AI model pretrained on vast amounts of text.

It understands text; extracts meanings; and learns grammar, languages, and knowledge.

Building generative AI applications can be challenging



Accessing
multiple FMs
and newer
versions



Customizing
FMs is not easy



Data privacy
and security



Getting FMs
to execute tasks



Connecting to
data sources



Difficult
to manage
infrastructure

NEW

AWS AI Portfolio

INTERFACES & PROTOCOLS (MCP/A2A)

SECURITY AND POLICIES

APPLICATIONS AND AGENTS

Kiro

Amazon Quick Suite

AWS Transform

Amazon Connect

AWS Marketplace

FRONTIER AGENTS

Kiro autonomous agent

AWS DevOps Agent

AWS Security Agent

AI & AGENT DEVELOPMENT SOFTWARE & SERVICES

Vertically Integrated

Amazon Nova Act

AGENT BUILDERS

Flexible/OSS

Strands Agents

AMAZON BEDROCK

Models

Amazon Nova

3P Models

Capabilities

Optimization

Customization

Guardrails

Knowledge Bases

AgentCore

Runtime

1P Tools

Gateway

Identity

Memory

Observability

Policy

Evaluations

INFRASTRUCTURE

AMAZON SAGEMAKER

Custom Model Building

Model building

Model training

Deployment

MLOps

Hyperpod

Data Foundation for AI

Data Processing

Governance

Storage and Databases

Vector Search

AI COMPUTE

Trainium

Inferentia

GPUs



Amazon Bedrock



Amazon Bedrock

Broad selection of fully managed models from leading AI companies

AI21labs



ANTHROPIC

cohere

deepseek

Google

Luma

Meta

MINIMAX

MISTRAL
AI

Moonshot AI

NVIDIA

OpenAI

Qwen

stability.ai

TwelveLabs

WRITER

Amazon Bedrock Marketplace

100+ PUBLICLY AVAILABLE AND PROPRIETARY EMERGING, POPULAR, AND SPECIALIZED MODELS THROUGH A SERVERFUL OFFERING



IBM



NVIDIA



HUGGING FACE



MISTRAL AI



STABILITY AI



LG AI RESEARCH



SNOWFLAKE



UPSTAGE



PREFERRED
NETWORKS



GRETTEL



WIDN



EVOLUTIONARY
SCALE



ARCEE AI



WRITER



CAMB.AI



DATABRICKS



STOCKMARK



KARAKURI



LIQUID



JOHN SNOW
LABS



CYBERAGENT



NCSOFT



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Prompt engineering



© 2026, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Prompt Engineering

- What is a **Prompt**?
 - Text input provided to an AI system to elicit a response
- What is **Prompt Engineering**?
 - Using techniques to craft prompts that steer FMs/LLMs towards desired responses
- Why is this important?
 - Enables control over models' behavior
 - Targets desired capabilities



Zero/Few-Shot Prompting

Zero Shot

Input

Tell me sentiment of this statement:
I loved the pizza at that Italian pizzeria

Output

The statement expresses a positive sentiment towards a pizza restaurant.

Few Shot

Input

Best Pakistani restaurant in Zurich: Positive
New York stinks, don't go there: Negative
The talk was on Generative AI: Neutral
This is a rip-off, store not recommended:

Output

Negative

Prompt Engineering - Methodology

1. Collect a set of test user queries to use as test cases
2. Write a prompt [template]
3. Evaluate against test cases
4. Refine prompt
5. Repeat steps 3 & 4 until results are satisfactory
6. Start using the prompt
7. Repeat steps 1 & 5 as requirements change

Prompt Engineering - Tips

Be clear and concise; use simple language

- Include context if needed
- Describe task, specify rules and directions about output formatting
- Provide diverse examples including edge cases (few-shot prompting)
- Instruct the model to say “I don’t know” to prevent hallucinations
- Follow the guidelines for each model (e.g. Reasoning models)
- Metaprompts – use an LLM to generate your prompt
- Leverage automatic prompt tuning/optimization frameworks (e.g. [DSPy](#))

Prompt Engineering – Model adaptation

Different models expect prompt in different formats, depending on how they were trained or fine-tuned

- Model creators provide prompt guides for their model
- e.g. [Amazon Nova](#):
 - Have clear sectional delimitation using **markdown**, **XML**, or other structure.
 - To define the name of the section, use **##Section Name##**, then refer to that section in your prompt with **##Section Name##**.
- e. g. [Claude](#):
 - Any input and output elements should **use XML tags**, e.g. put context inside **<context></context>** tags
- e.g. [Mistral](#)
 - Add instructions like "You will only respond with a **JSON object**"
 - Use delimiters like **###, <<< >>>**

Amazon Bedrock offers several features to balance cost, latency and accuracy

Amazon Bedrock support for prompt caching

Cache repetitive context in prompts across multiple API calls

Securely cache entire prompts

Enhance accuracy through longer, detailed prompts

Reduce costs by **up to 90%** and latency **by up to 85%** for supported models

Amazon Bedrock Intelligent Prompt Routing

Automatically route prompts to different foundation models to optimize response quality and lower costs

Provides a single endpoint to efficiently route prompts

Uses advanced prompt matching to meet cost and latency thresholds

Reduce costs by **up to 30%** without compromising on accuracy

Amazon Bedrock Model Distillation

Utilize smaller, more cost-effective models

Maximize distilled model performance with proprietary data synthesis

Distilled models are up to 500% faster and up to **75% less expensive** than original models, with **less than 2% accuracy loss** for use cases like RAG

Model Parameters



© 2026, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Model parameters - Variety

Temperature

Adjusts randomness of the output

Lower values make output more predictable

Higher values make output more creative

DECREASE

TEMPERATURE

INCREASE



LESS RANDOM

OUTPUT

MORE RANDOM

Model parameters - Variety

Temperature

Adjusts randomness of the output

Lower values make output more predictable

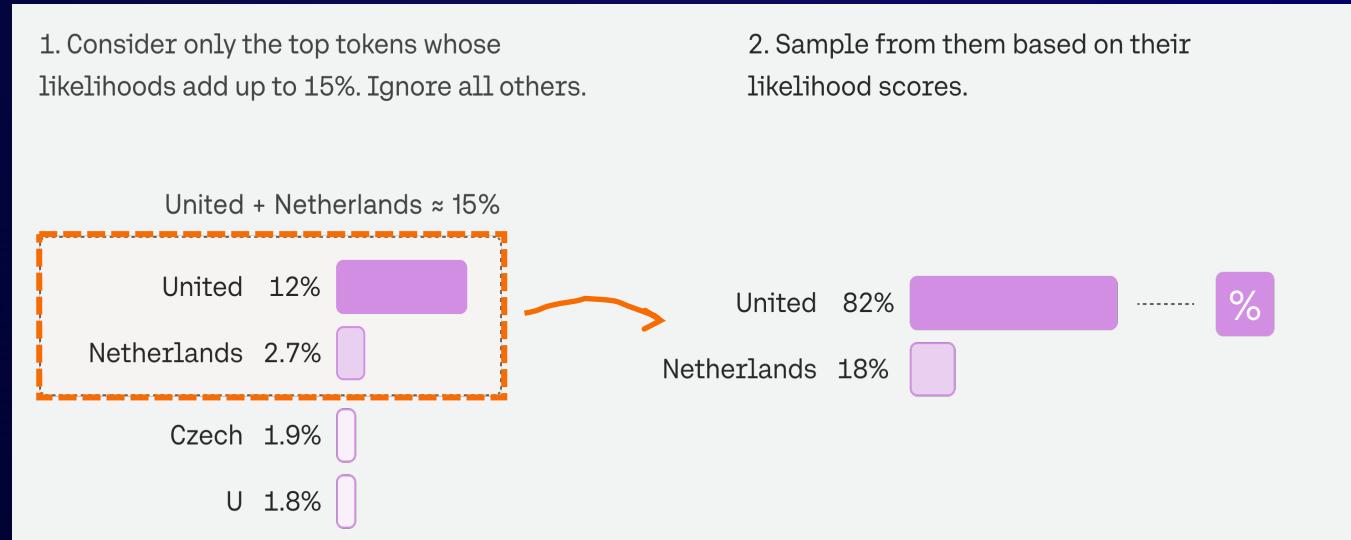
Higher values make output more creative

Top P/Top K

Only include the most likely tokens when choosing the next value

Limits the choices the model can pick from

Prevents choosing highly unlikely options



Model parameters - Length

Stop Words

Stop generation if one of the words is encountered.

Limit model output based on context

Allow model to generate different length of sequences based on input

Model parameters - Length

Stop Words

Stop generation if one of the words is encountered.

Limit model output based on context

Allow model to generate different length of sequences based on input

Max Token Count / Response Length

Stop generation once certain number of tokens is produced

Prevent high usage for unusual prompts

Limit cost of model errors

But wait, how can I use code?

THROUGH AMAZON BEDROCK 2 MAIN APIs

InvokeModel

Control model-specific parameters

Bedrock model identifier

e.g anthropic.claude-sonnet-4-20250514-v1:0

Inference configuration

temperature, top_p , max_tokens etc.

Messages (actual prompt)

The prompt to sent to the model

**The *body* argument
structure will change based
on the model used**

```
1 client = boto3.client("bedrock-runtime", region_name="us-east-1")
2 client.invoke_model(
3     modelId=model_id,
4     body=json.dumps({
5         "anthropic_version": "bedrock-2023-05-31",
6         "max_tokens": 512,
7         "temperature": 0.5,
8         "messages": [
9             {
10                 "role": "user",
11                 "content": [{"type": "text", "text": prompt}],
12             },
13         ]),
14     )
15 )
```

But wait, how can I use code?

THROUGH AMAZON BEDROCK 2 MAIN APIs

Converse

Unified Interface Across Models

Bedrock model identifier

e.g anthropic.claude-sonnet-4-20250514-v1:0

Inference configuration

temperature, top_p , max_tokens etc.

Messages (actual prompt)

The prompt to sent to the model

```
1 client = boto3.client("bedrock-runtime", region_name="us-east-1")
2 client.converse(
3     modelId=model_id,
4     inferenceConfig={"maxTokens": 512, "temperature": 0.5},
5     messages=[{
6         "role": "user",
7         "content": [{"text": prompt}],
8     }],
9 )
```

Recommended API

provides a consistent interface across all Amazon Bedrock models



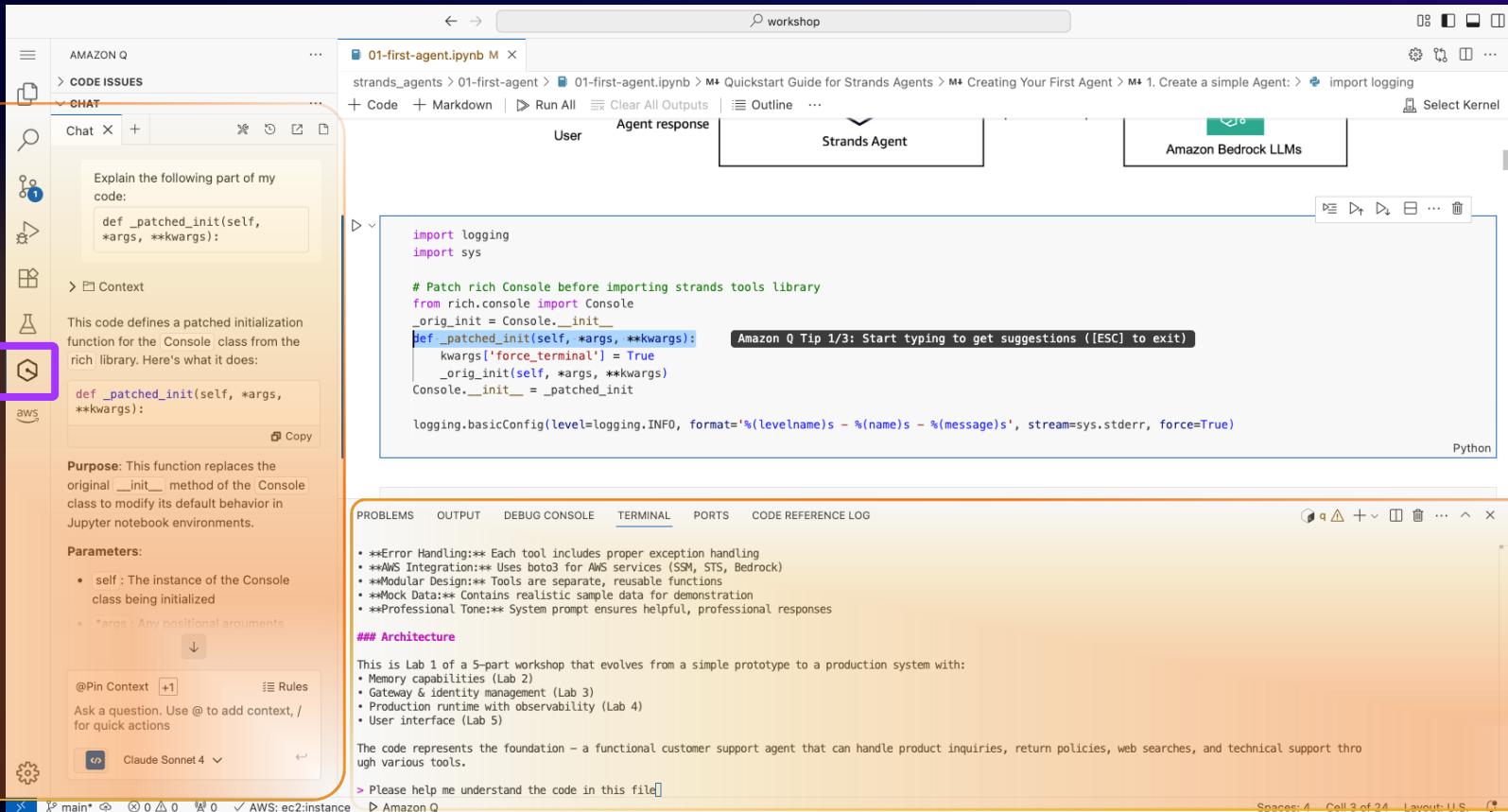
Workshop time!

Ask us,... ...or ask Amazon Q & Kiro

Amazon Q
Developer
Extension



Click the Q extension
icon to begin



Kiro CLI

Type "kiro-cli"
to begin



We've installed the Amazon Q Extension and Kiro CLI for you.
Make use of it during the workshop!



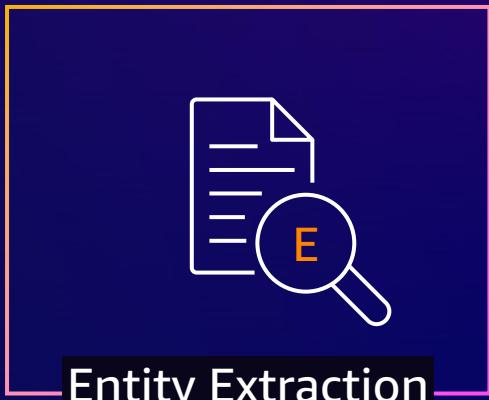
Let's get hands on - Architecture patterns



Text generation



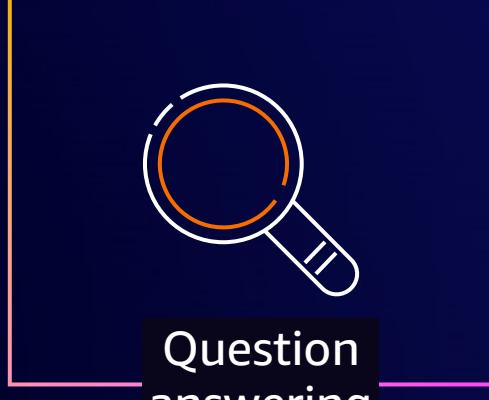
Summarization



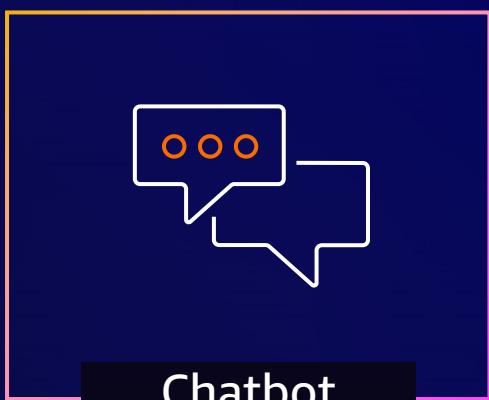
Entity Extraction



Code generation



Question
answering



Chatbot

Multimodality with Amazon Nova



© 2026, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Amazon Nova Foundation Models

State-of-the-art foundation models that deliver frontier intelligence and industry-leading price performance

UNDERSTANDING MODELS

Amazon Nova

Micro

Amazon Nova

Pro

Amazon Nova

Lite

Amazon Nova

Premier

CREATIVE CONTENT GENERATION MODELS

Amazon Nova

Canvas

Amazon Nova

Reel

SPEECH-TO-SPEECH MODEL

Amazon Nova

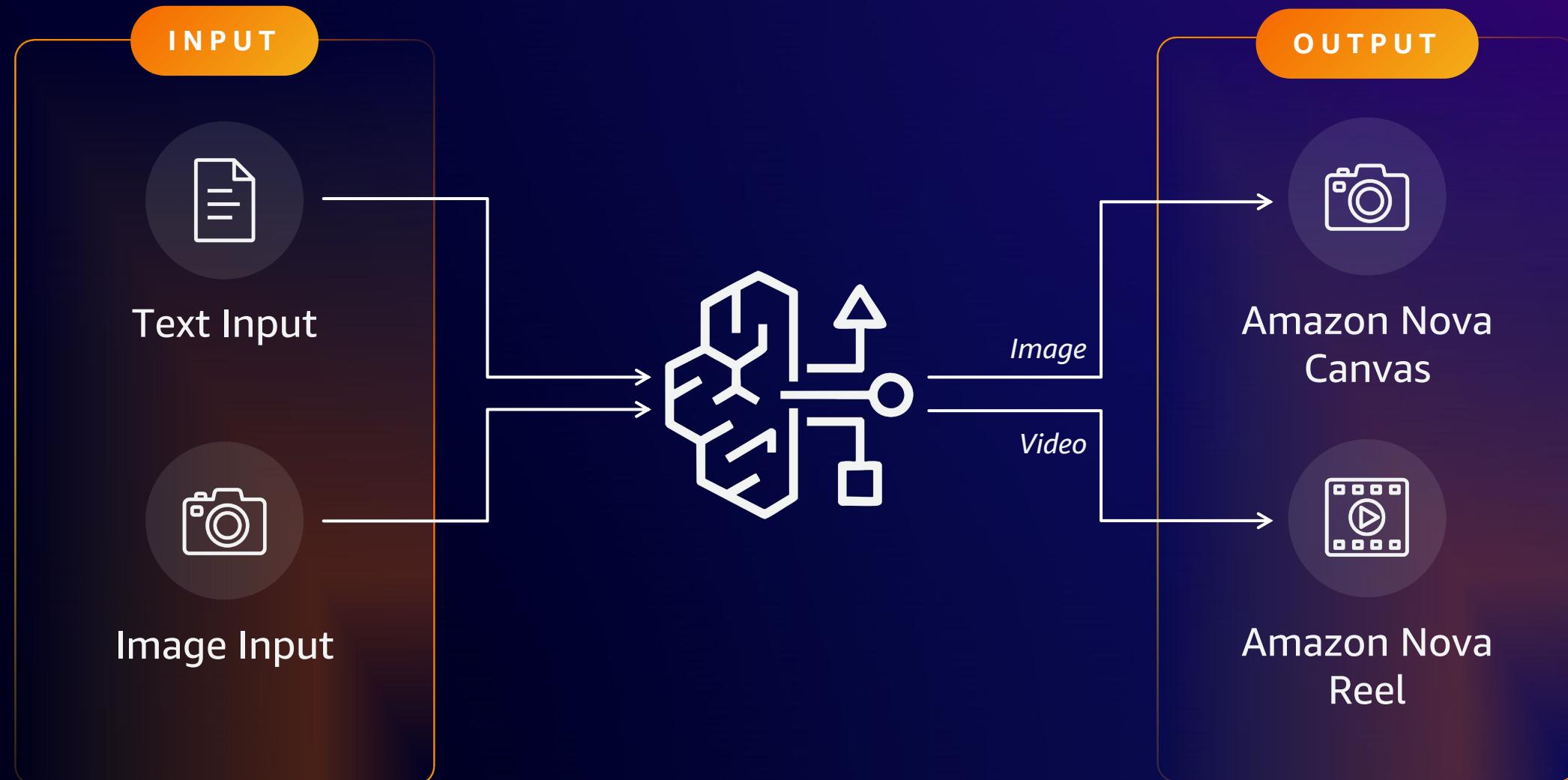
Sonic



Amazon Nova Understanding Models

	Amazon Nova Micro	Amazon Nova Lite	Amazon Nova Pro	Amazon Nova Premier
Availability	GA	GA	GA	GA
Context window	128K	300K	300K (5M coming soon)	1M
Languages	200+ languages	200+ languages	200+ languages	200 + languages
Modalities supported	Text input; text output	Text, image, video input; text output	Text, image, video input; text output	Text, image, video input; text output
Fine-tuning	Yes	Yes	Yes	Yes

Amazon Nova Creative-Content Models



NEW

Amazon Nova 2 Models



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Nova 2 foundation models

GENERALLY AVAILABLE

Amazon
Nova 2 Lite

PREVIEW

Amazon
Nova 2 Pro

PREVIEW

Amazon
Nova 2 Omni

GENERALLY AVAILABLE

Amazon
Nova 2 Sonic

Fast, cost-effective reasoning
model for everyday workloads

Our most intelligent reasoning
model for highly complex tasks

Unified model for multimodal
reasoning and image generation

Speech-to-speech
foundation model for real-
time, human-like
conversational AI



Nova 2 Omni extracts key information from documents

Contract Agreement Between:

KOAT
3801 Carlisle Blvd. NE
Albuquerque, NM 87107
(505)884-7777
www.koat.com

SRCP Media Inc
Attention: Betsy Vonderheid
201 North Union Street
Ste 200
Alexandria, VA 22314

And:

CONTRACT

Print Date: 10/28/16 Page: 1 of 4

Contract / Revision	All Order #
1535975 /	08415555
Product Candidate	
Contract Dates	Estimate #
11/02/16 - 11/08/16	1186 ESPINOZA 4 SEC
Advertiser	
N. Espinoza/R/Secretary of State	Original Date / Revision
10/28/16 / 10/28/16	
Billing Cycle	
EOM/ECC	Broadcast
Property	
Account Executive	Sales Office
KOAT	Mary Tricoli
Special Handling	
Demographic	
Adults 35+	
Adv Code	Advertiser Code
156	409
Agency Ref	
Advertiser Ref	

*Line Ch Start Date End Date Description Start/End Time Days Spots/Length Week Rate PCode/Rtn Type Spots Amount

N 12 KOAT 11/07/16 11/07/16 Jimmy Kimmel Jimmy Kimmel :30 P-6 NM 1 \$45.00

Class of Time - Immediately Pre-emptible without notice

Start Date	End Date	Weekdays	Spots/Week	Rate
Week: 11/07/16	11/13/16	1-----	1	\$50.00

Spot Ch Date Range Description Start/End Time Weekdays Length Rate Type

1 KOAT 11/07/16-11/13/16 Jimmy Kimmel M----- :30 \$50.00 NM

See MG 12.2

2 KOAT 11/02/16-11/04/16 Nightline Nightline ---WTF--- :30 \$45.00 NM

MG for 12.11/07

N 15 KOAT 11/02/16 11/04/16 ELLEN EF 4-5p :30 P-6 NM 3 \$300.00

Class of Time - Immediately Pre-emptible without notice

Start Date	End Date	Weekdays	Spots/Week	Rate
Week: 10/31/16	11/06/16	---WTF---	3	\$100.00

N 21 KOAT 11/02/16 11/04/16 Live with Kelly & Michael 9-10a :30 P-6 NM 3 \$300.00

Class of Time - Immediately Pre-emptible without notice

Start Date	End Date	Weekdays	Spots/Week	Rate
Week: 10/31/16	11/06/16	--WTF--	3	\$100.00

N 30 KOAT 11/07/16 11/07/16 ELLEN EF 4-5p :30 P-6 NM 1 \$105.00

Class of Time - Immediately Pre-emptible without notice

Start Date	End Date	Weekdays	Spots/Week	Rate
Week: 11/07/16	11/13/16	M-----	1	\$105.00

N 31 KOAT 11/02/16 11/04/16 General Hospital 2-3p/1-2p :30 P-6 NM 1 \$95.00

Class of Time - Immediately Pre-emptible without notice

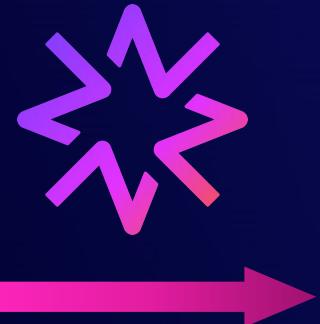
Start Date	End Date	Weekdays	Spots/Week	Rate
Week: 10/31/16	11/06/16	--WTF--	1	\$95.00

N 32 KOAT 11/06/16 11/06/16 ET Weekend Su 4-5p 4-5p :30 P-6 NM 1 \$75.00

Class of Time - Immediately Pre-emptible without notice

Start Date	End Date	Weekdays	Spots/Week	Rate
Week: 10/31/16	11/06/16	-----S	1	\$75.00

Totals 0.00 86 \$31,070.00



Between KOAT & SRCP Media Inc

Date: 10/28/16

Spot Schedule

Line	Program	Date Range	Time	Length	Rate	Spots
12	Jimmy Kimmel	11/07/16	--	:30	\$50.00	1
15	ELLEN	11/02-04/16	4-5p	:30	\$100.00	3
21	Live with Kelly	11/02-04/16	9-10a	:30	\$100.00	3
30	ELLEN	11/07/16	4-5p	:30	\$105.00	1
31	General Hospital	11/02-04/16	2-3p	:30	\$95.00	1
32	ET Weekend	11/06/16	4-5p	:30	\$75.00	1

Financial Terms

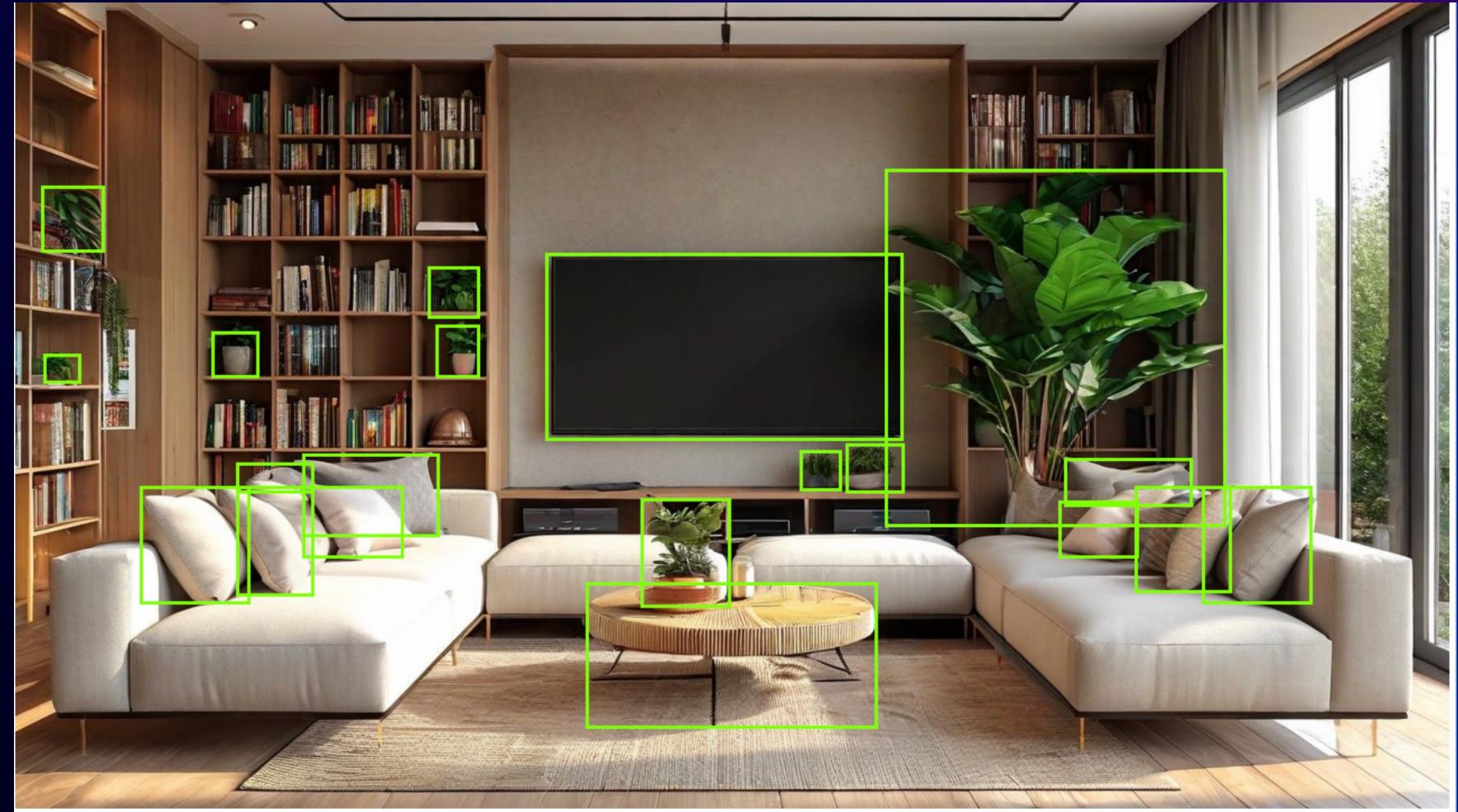
- Gross: \$31,070.00
- Agency Comm: (\$4,660.50)
- Net: \$26,409.50
- Tax (7.313%): \$1,931.33
- Total: \$28,340.83



And **understands** what objects are in your images

Prompt:

Detect plants,
cushions,
tables and
the TV



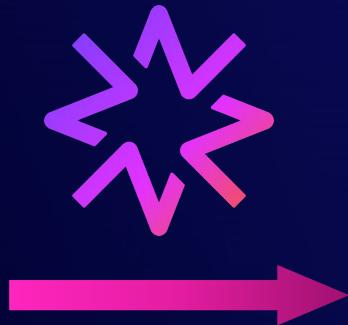
Nova 2 Omni can also **generate** hyper-realistic imagery

Prompt:

Solo female traveler
with a backpack,
standing on mountain
summit with arms
raised in triumph



And **edit** your existing images



Prompt: Change the snowy environment surrounding the house to a beach setting with sand and palm trees.

Nova 2 Lite

A fast cost-effective reasoning model

Input: Text, image, video

Output: Text

Use cases: Everyday AI tasks such as chatbots, code generation

Context window: 1M

Developer controls for thinking effort: Yes

Nova 2 Pro (preview)

Our most intelligent reasoning model

Input: Text, image, video, speech

Output: Text

Use cases: Complex AI tasks such as agentic coding

Context window: 1M

Developer controls for thinking effort: Yes

Nova 2 Omni (preview)

All-in-one model for multimodal reasoning and image generation

Input: Text, image, video, speech

Output: Text, image

Use cases: Image generation, multimodal understanding

Context window: 1M

Developer controls for thinking effort: Yes

Nova 2 Sonic

A speech-to-speech model for conversational AI

Input: Text, speech

Output: Text, speech

Use cases: Customer support, voice-enabled personal assistant

Context window: 1M

Nova Multimodal Embeddings

A state-of-the-art multimodal embedding model

Input: Text, image, video, audio

Output: Embedding

Use cases: Semantic search, agentic RAG

Context window: 8K

G E N E R A L L Y A V A I L A B L E

Amazon Nova Forge

Build your own frontier models using
Amazon Nova

1

Access checkpoints across **all phases of model development**, and leverage new Nova models before they are widely available

2

Blend **your proprietary data** with Amazon Nova-curated training data

3

Perform **reinforcement learning** with reward functions in your environment

4

Use **push-button recipes** that are optimized to build with Nova through visual workflows or a command line interface

5

Use the built-in responsible AI toolkit to implement **custom safety guardrails**

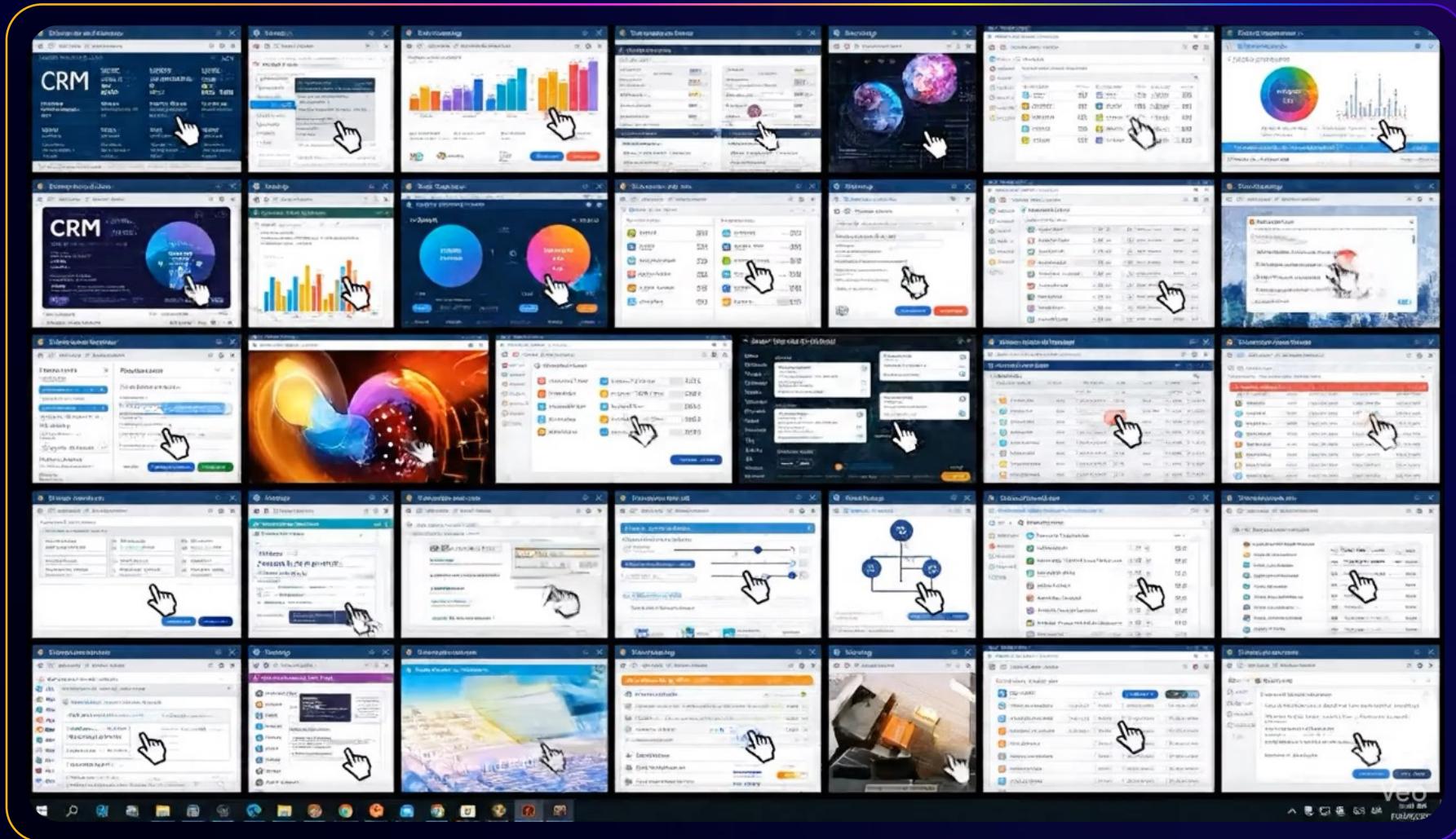
GENERALLY AVAILABLE

Amazon Nova Act

Our AWS service for building, deploying
and running reliable AI agents

- Browser-user model
- New AWS Service & Console
- Nova Act Playground
- SDK
- IDE Extension
- CLI to deploy to AWS
- Human-in-the-loop
- Tool Use

Trained using real web interactions



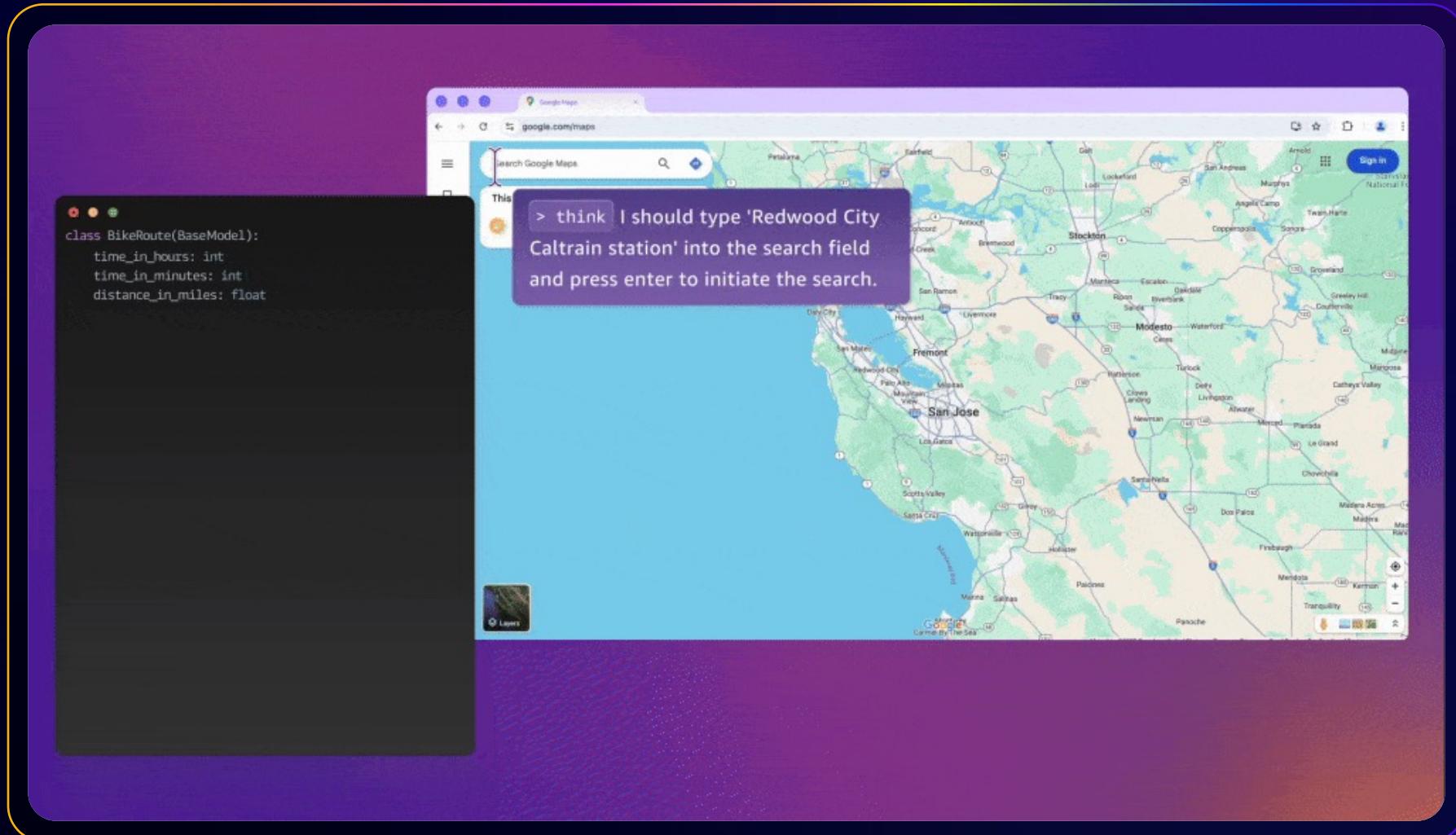
Act can interact with key browser elements

Date Pickers

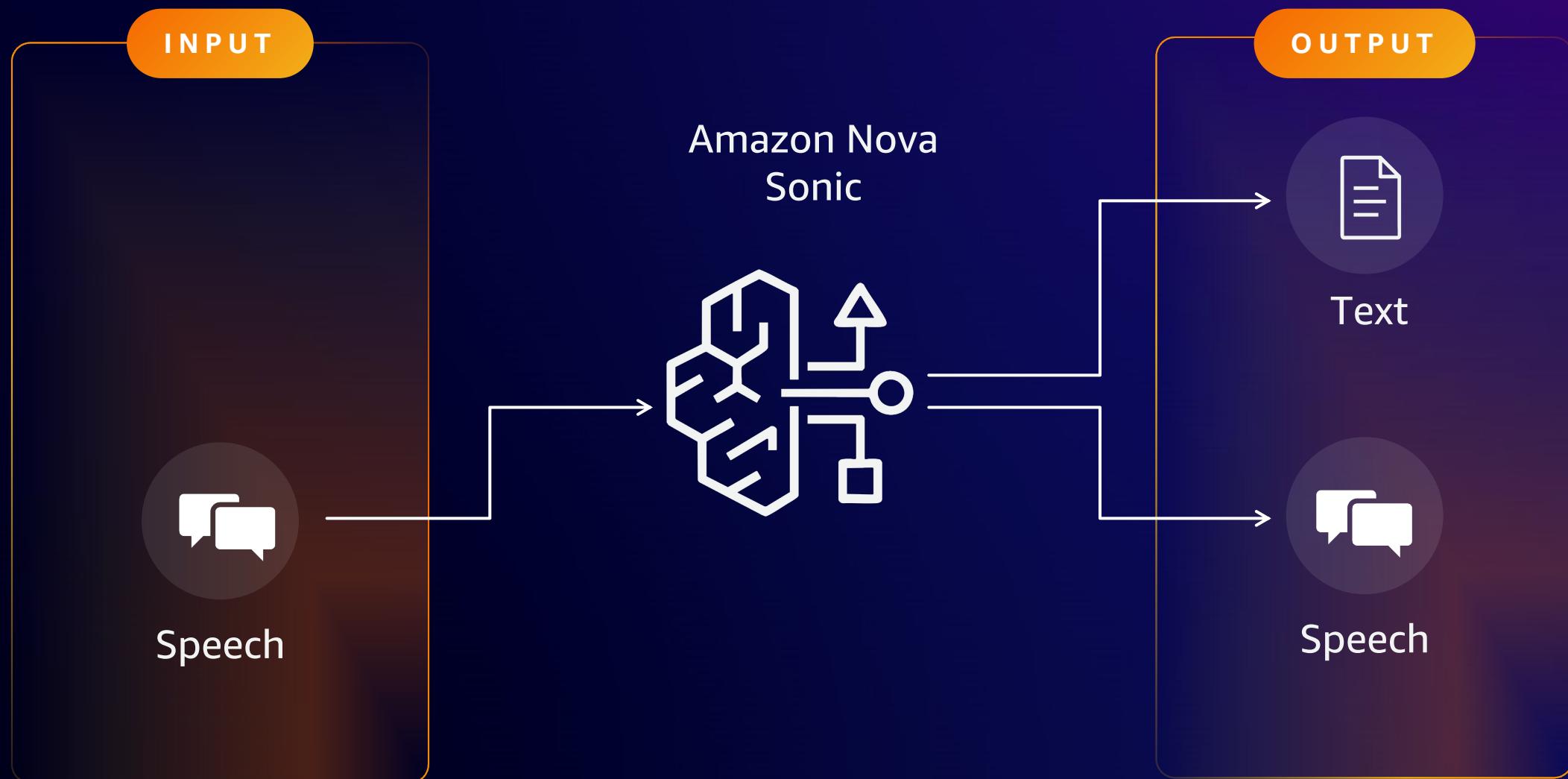
Search/Filter

Dropdowns

Mimicking real user actions

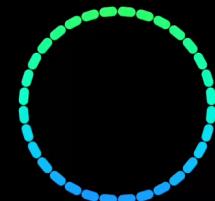


Amazon Nova Speech-to-Speech



Amazon Nova Sonic In Action

Customer service call automation



Recapping Amazon Nova

Comprehensive suite of
models and services
with industry leading
price performance

AGENT SERVICES

CUSTOMIZATION

Amazon
Nova Act

Amazon
Nova Forge

SECOND GENERATION MODELS

Amazon
Nova 2 Lite

Amazon
Nova 2 Sonic

PREVIEW MODELS

Amazon
Nova 2 Pro

Amazon
Nova 2 Omni





GENAI BOOTCAMP

Notebook Exercises

Retrieval Augmented Generation (RAG)



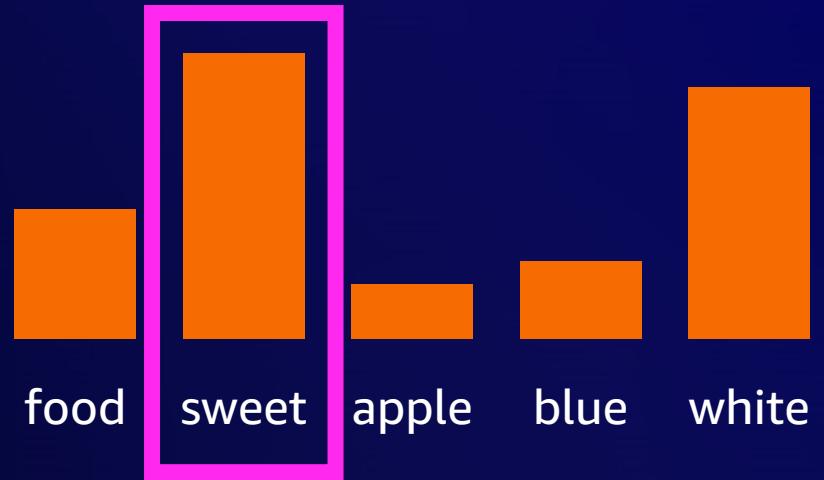
© 2026, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Large Language Models Recap

Roses are red, violets are blue, sugar is _____

Large Language Models Recap

Roses are red, violets are blue, sugar is sweet



Private Data

What is AnyCompany's vacation policy?

**How do we “teach” the model
about AnyCompany?**

Use Cases



Improved content quality

E.g., helps in reducing hallucinations and connecting with recent knowledge including enterprise data



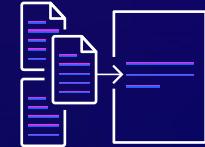
Contextual chatbots and question answering

E.g., enhance chatbot capabilities by integrating with real-time data



Personalized search

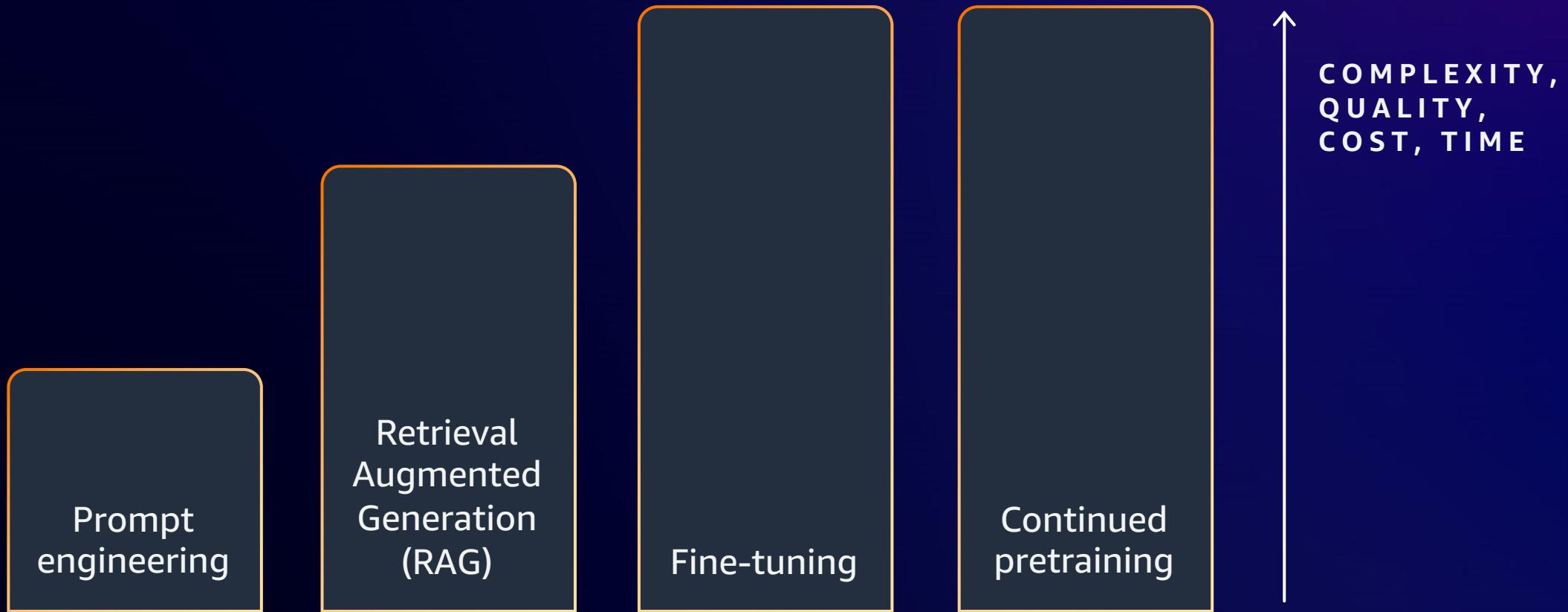
E.g., searching based on user previous search history and persona



Real-time data summarization

E.g., retrieving and summarizing transactional data from databases, or API calls

Customizing LLMs



File Attachment



Based on the following file, what is the vacation policy?



File Attachment



Based on the following file, what is the vacation policy?

<file_content>

Employees are entitled to 10 days (80 hours) of paid vacation annually after completing 6 months of employment. Vacation time must be requested at least 2 weeks in advance and approved by the immediate supervisor.

</file_content>



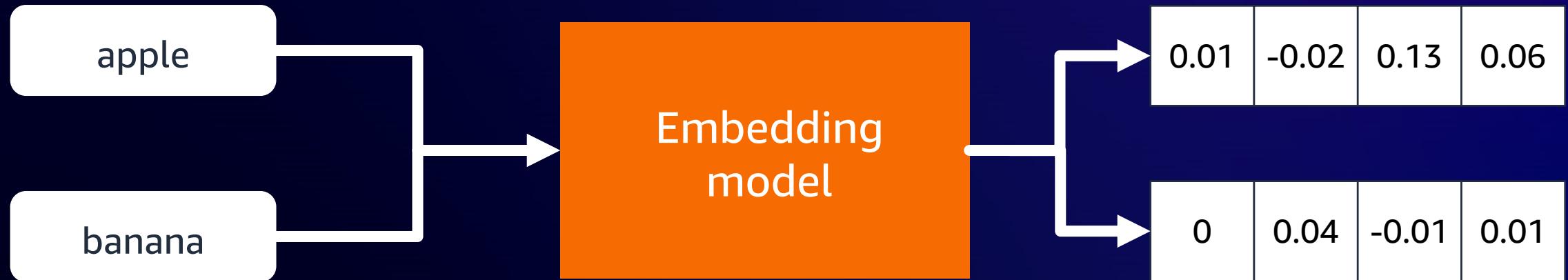
Based on our company policy, you get 10 paid vacation days (or 80 hours) each year once you've been with us for 6 months. Just make sure to submit your vacation requests at least 2 weeks ahead of time so your supervisor can approve them. Let me know if you have any other questions!

We know how to use private data!

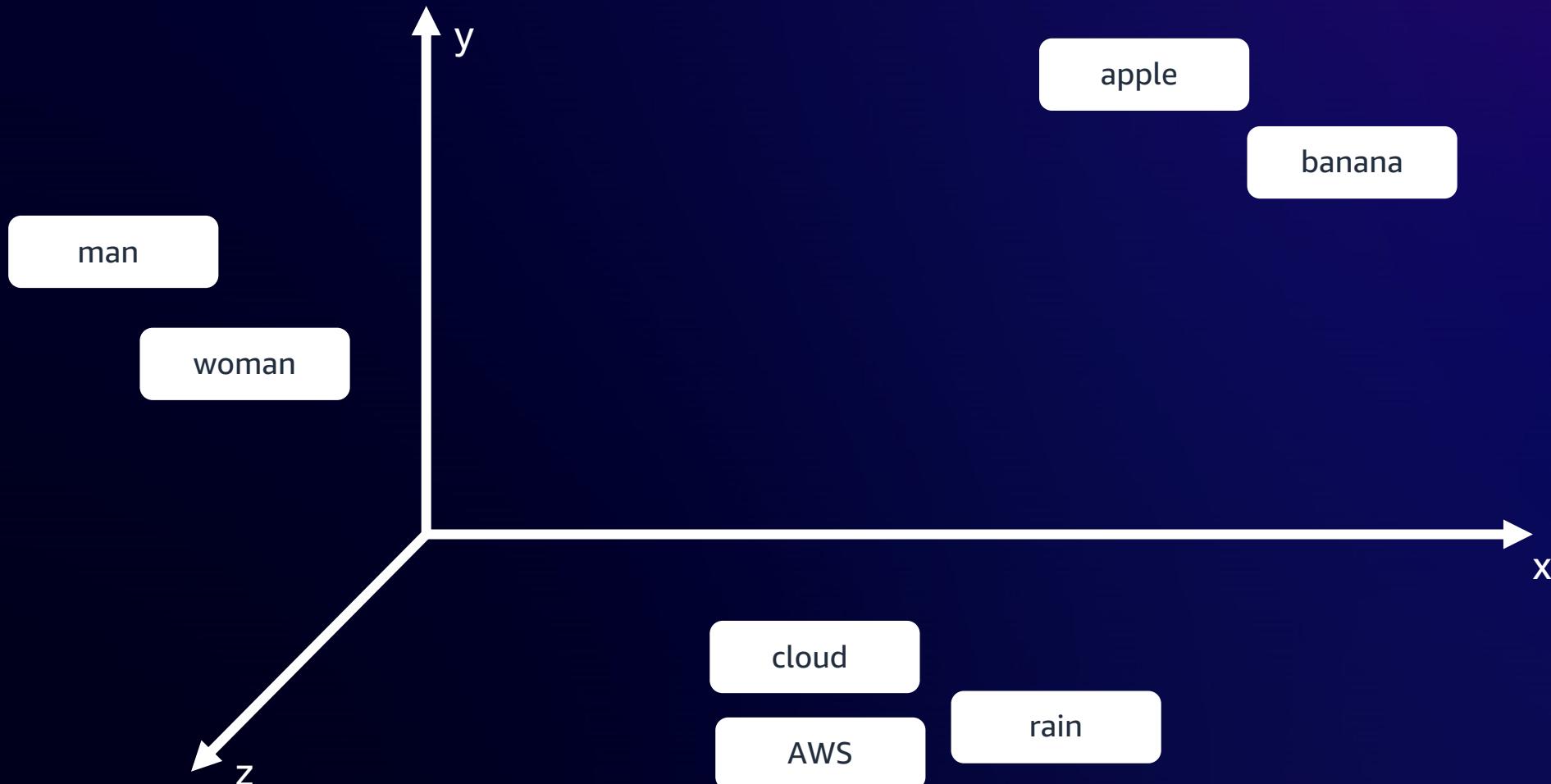
But how do we find the right file?



Embeddings



Embeddings



Vector Databases

“Retrieve the top 2 closest vectors to [0.04, -0.02, 0.1, 0.3]”



0.01	-0.02	0.13	0.06
-0.11	-0.02	0.49	0

0.01	-0.02	0.13	0.06	apple
-0.11	-0.02	0.49	0	banana
0	0.13	-0.4	0.01	AWS



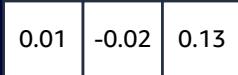
Putting It All Together



AnyCompany's
data sources



Embedding model
converts files into vectors



Store vectors
in database

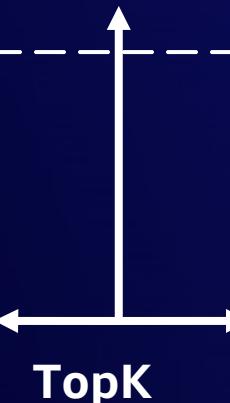
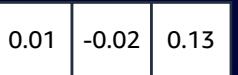
indexing



What is
AnyCompany's
vacation policy?



Embedding model converts
prompt into vectors



Q: Based on the following, what is
AnyCompany's vacation policy?
<file_content>

A: AnyCompany's vacation policy is...

Prompt is augmented and
a response is generated

retrieval

Context Window



Chunking

SECTION 3: COMPENSATION AND SALARY STRUCTURE

3(a) Base Salary Qualification Requirements

- i. Veterinarians holding a Doctor of Veterinary Medicine (DVM) degree from an AVMA-accredited institution shall receive a minimum base salary of \$85,000 per annum.
- ii. Veterinarians holding additional board certifications shall receive a premium of \$15,000 per certification above the base salary.
- iii. Non-licensed veterinary professionals or those with pending licensure shall receive 70% of the base salary until full licensure is obtained.

3(b) Experience-Based Compensation

- i. For each year of post-DVM clinical experience, employees shall receive an annual increment of \$2,500, capped at 15 years.
- ii. Specialized experience in emergency medicine, surgery, or exotic animal care shall warrant an additional \$7,500 per annum.
- iii. Service as a lead veterinarian or department head shall result in a 12% increase to the adjusted base salary.

3(c) Performance and Additional Compensation

- i. Annual performance reviews resulting in "Exceeds Expectations" ratings shall merit a 5% bonus of the total annual salary.
- ii. On-call duties shall be compensated at a rate of \$200 per weekday night and \$300 per weekend day/night.
- iii. Continuing education completion beyond the minimum required hours shall be rewarded at \$500 per 10 additional hours, with an annual cap of \$2,500.

Chunking

SECTION 3: COMPENSATION AND SALARY STRUCTURE

3(a) Base Salary Qualification Requirements

- i. Veterinarians holding a Doctor of Veterinary Medicine (DVM) degree from an AVMA-accredited institution shall receive a minimum base salary of \$85,000 per annum.
- ii. Veterinarians holding additional board certifications shall receive a premium of \$15,000 per certification above the base salary.
- iii. Non-licensed veterinary professionals or those with pending licensure shall receive 70% of the base salary until full licensure is obtained.

3(b) Experience-Based Compensation

- i. For each year of post-DVM clinical experience, employees shall receive an annual increment of \$2,500, capped at 15 years.
- ii. Specialized experience in emergency medicine, surgery, or exotic animal care shall warrant an additional \$7,500 per annum.
- iii. Service as a lead veterinarian or department head shall result in a 12% increase to the adjusted base salary.

3(c) Performance and Additional Compensation

- i. Annual performance reviews resulting in "Exceeds Expectations" ratings shall merit a 5% bonus of the total annual salary.
- ii. On-call duties shall be compensated at a rate of \$200 per weekday night and \$300 per weekend day/night.
- iii. Continuing education completion beyond the minimum required hours shall be rewarded at \$500 per 10 additional hours, with an annual cap of \$2,500.

Chunking

SECTION 3: COMPENSATION AND SALARY STRUCTURE

3(a) Base Salary Qualification Requirements

- i. Veterinarians holding a Doctor of Veterinary Medicine (DVM) degree from an AVMA-accredited institution shall receive a minimum base salary of \$85,000 per annum.
- ii. Veterinarians holding additional board certifications shall receive a premium of \$15,000 per certification above the base salary.
- iii. Non-licensed veterinary professionals or those with pending licensure shall receive 70% of the base salary until full licensure is obtained.

3(b) Experience-Based Compensation

- i. For each year of post-DVM clinical experience, employees shall receive an annual increment of \$2,500, capped at 15 years.
- ii. Specialized experience in emergency medicine, surgery, or exotic animal care shall warrant an additional \$7,500 per annum.
- iii. Service as a lead veterinarian or department head shall result in a 12% increase to the adjusted base salary.

3(c) Performance and Additional Compensation

- i. Annual performance reviews resulting in "Exceeds Expectations" ratings shall merit a 5% bonus of the total annual salary.
- ii. On-call duties shall be compensated at a rate of \$200 per weekday night and \$300 per weekend day/night.
- iii. Continuing education completion beyond the minimum required hours shall be rewarded at \$500 per 10 additional hours, with an annual cap of \$2,500.

Chunking Strategy

	Fixed size	Semantic	Hierarchical	Default	No-chunking
What?	Configure chunks by specifying the number of tokens per chunk and an overlap percentage	Intelligently splits content based on meaning and context, so that related information stays together while maintaining semantic coherence	Creates a two-tiered structure of parent and child chunks, where larger sections contain smaller, related segments while preserving document organization and contextual relationships	Splits content into chunks of approximately 300 tokens, preserving sentence boundaries	Treats each document as a single text chunk
Pros	<ul style="list-style-type: none"> Easy to implement Quick processing 	<ul style="list-style-type: none"> Preserves meaning and context Handle variable content 	<ul style="list-style-type: none"> Maintains structure Rich context preservation to improve search relevance 	<ul style="list-style-type: none"> No configuration needed 	<ul style="list-style-type: none"> Complete control with custom implementation
Cons	<ul style="list-style-type: none"> Might break semantic meaning Less context aware 	<ul style="list-style-type: none"> Unpredictable chunk sizes Variable processing time 	<ul style="list-style-type: none"> Requires a clear structure 	<ul style="list-style-type: none"> Limited control 	<ul style="list-style-type: none"> More development time Higher maintenance
When?	<ul style="list-style-type: none"> Uniform structure Need predictable chunking 	<ul style="list-style-type: none"> Variable content structure Context preservation is crucial 	<ul style="list-style-type: none"> Clear parent-child relationship 	<ul style="list-style-type: none"> Quick implementation No special requirements 	<ul style="list-style-type: none"> Need custom processing Special data format
Example	Logs API doc with standardized format	News/Blog article Customer support conversations	Software doc Technical Manuals Legal Documents Academic Papers	Meeting notes	Multiformat doc requiring custom processing Employee data in CSV

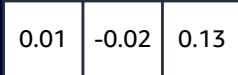
Putting It All Together



AnyCompany's
data sources



Embedding model
converts files into vectors



Store vectors
in database

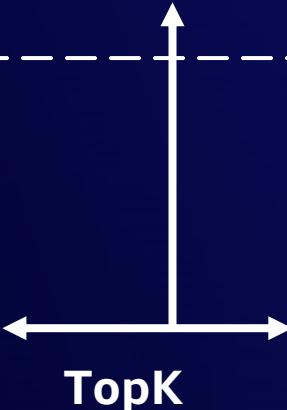
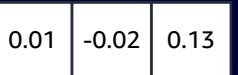
indexing



What is
AnyCompany's
vacation policy?



Embedding model converts
prompt into vectors



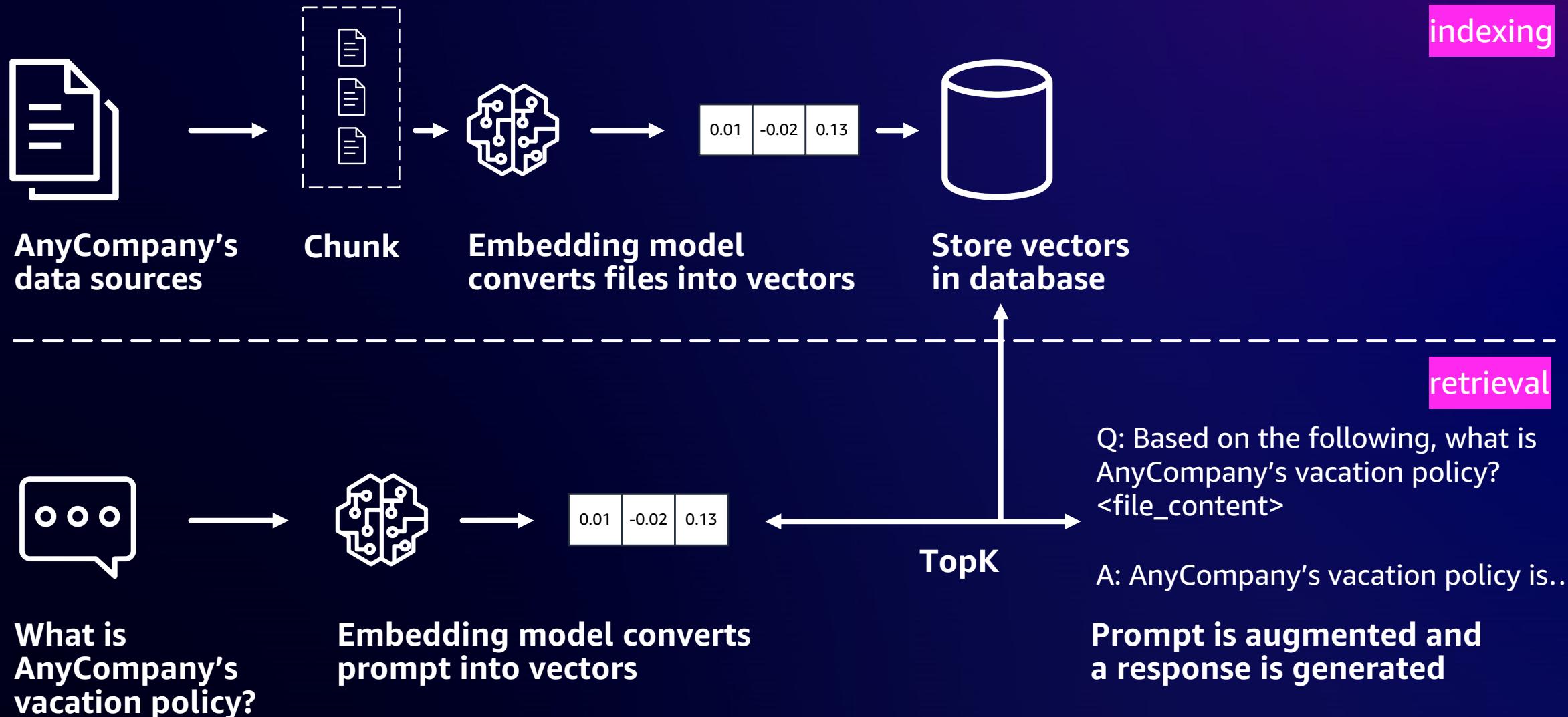
retrieval

Q: Based on the following, what is
AnyCompany's vacation policy?
<file_content>

A: AnyCompany's vacation policy is...

Prompt is augmented and
a response is generated

Putting It All Together



Knowledge Bases for Amazon Bedrock

Fully-managed native support for retrieval augmented generation



Fully managed support for end-to-end RAG workflow



Securely connect FMs and agents to data sources



Automatically converts text documents into embeddings



Stores embeddings in your vector database



Retrieves embeddings and augments prompts



Provide source attribution

Building RAG-based Apps on Amazon Bedrock

Choose your favorite knowledge base or vector search enabled service

NATIVELY AVAILABLE BEDROCK KNOWLEDGE BASES:



Vector engine
for Amazon
OpenSearch Serverless



Amazon Aurora
PostgreSQL
with pgvector



Pinecone



Redis
Enterprise
Cloud



MongoDB

AMAZON DATABASES ENABLED WITH VECTOR SEARCH:



Amazon
OpenSearch
Serverless



Amazon
OpenSearch
Service



Amazon
Aurora
PostgreSQL



Amazon RDS
PostgreSQL



Amazon
DynamoDB
via zero-ETL



Amazon
MemoryDB



Amazon
Neptune



Amazon
DocumentDB



Amazon S3 Vectors

First cloud object store with native support
to store and query vectors

Sub-second

vector query
performance at scale

up to 90%

lower overall costs
for uploading, storing,
and querying vectors

Billions

of vectors stored
at low cost

Using Code



```
1 import boto3
2
3 client = boto3.client('bedrock-runtime')
4
5 response = client.retrieve(
6     knowledgeBaseId='string',
7     retrievalQuery={
8         'text': 'string'
9     }
10 )
11
12 response = client.retrieve_and_generate(
13     input={
14         'text': 'string'
15 },
16     retrieveAndGenerateConfiguration={
17         'type': 'KNOWLEDGE_BASE',
18         'knowledgeBaseConfiguration': {
19             'knowledgeBaseId': 'string',
20             'modelArn': 'string',
21         }
22     }
23 )
```



Using Code

```
● ● ●  
1 import boto3  
2  
3 client = boto3.client('bedrock-runtime')  
4  
5 response = client.retrieve(  
6     knowledgeBaseId='string',  
7     retrievalQuery={  
8         'text': 'string'  
9     }  
10 )  
11  
12 response = client.retrieve_and_generate(  
13     input={  
14         'text': 'string'  
15     },  
16     retrieveAndGenerateConfiguration={  
17         'type': 'KNOWLEDGE_BASE',  
18         'knowledgeBaseConfiguration': {  
19             'knowledgeBaseId': 'string',  
20             'modelArn': 'string',  
21         }  
22     }  
23 )
```



Advanced Techniques



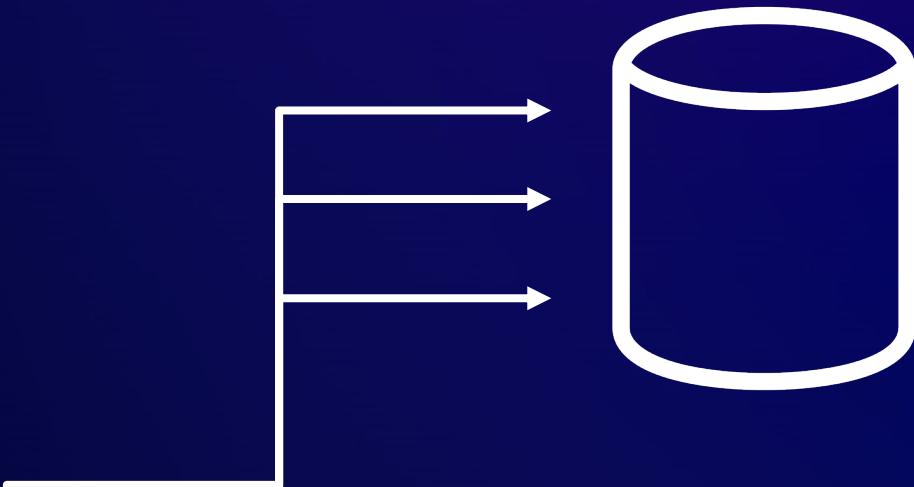
© 2026, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Query Reformulation

"Where is the Octank company waterfront building located and how does the whistleblower scandal hurt the company and its image?"



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



Metadata Filtering

0.01	-0.02	0.13	0.06	apple
-0.11	-0.02	0.49	0	banana
0	0.13	-0.4	0.01	AWS



Metadata Filtering

0.01	-0.02	0.13	0.06	apple	type=food
-0.11	-0.02	0.49	0	banana	type=food
0	0.15	0.1	0.01	AWS	type=cloud



"Retrieve the top 2 closest vectors to [0.04, -0.02, 0.1, 0.3] only if type=food"

HyDE/Hypothetical Questions

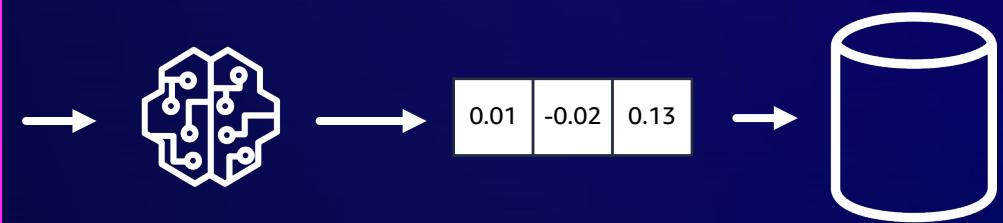
HyDE

What is the vacation policy?



Vacation Day Entitlement at XYZ Corp:
At XYZ Corp vacation day entitlement is structured as follows:
1. New employees: 10 days per year
2. After 2 years of service: 15 days per year
3. After 5 years of service: 20 days per year
4. Managers and above: Additional 5 days per year
Factors affecting vacation entitlement:
- Length of service
- Position level
- Full-time vs. part-time status
- Any negotiated terms in individual contracts
Vacation days are accrued monthly and can be used after a 3-month probationary period. Unused vacation days can be carried over to the next year, up to a maximum of 5 days.

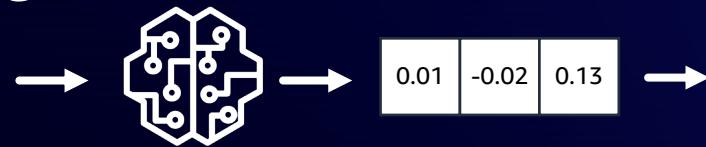
Hypothetical Document Embedding (HyDE)



HyDE/Hypothetical Questions

Hypothetical Questions

What is the vacation policy?

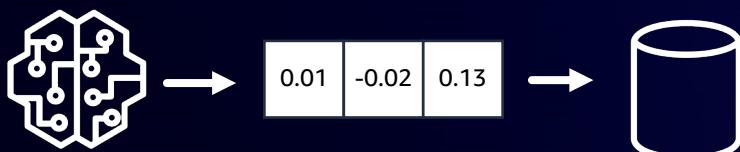
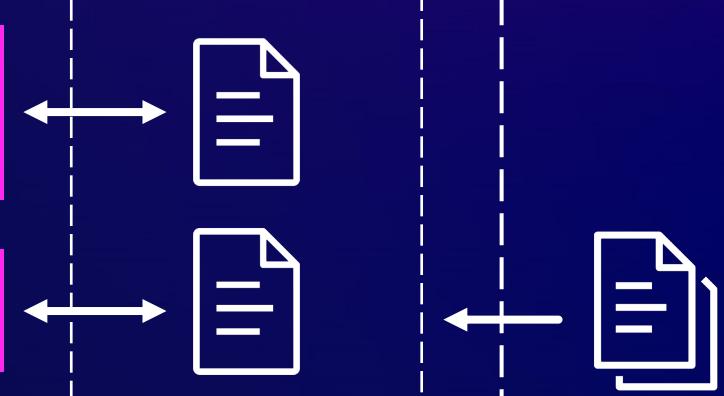


Hypothetical Questions

"What is the vacation policy at our company?"

"How much annual leave do I get?"

"How many days off do I get per year?"





GENAI BOOTCAMP

Notebook Demo

Unstructured data with Amazon Bedrock Data Automation



The multi-modal challenge

80% of data is unstructured in document, image, video, and audio formats, while...

<20% of organizations are able to take advantage of unstructured data at scale



Diverse formats and variations in asset types



Hard to get the desired accuracy



Lack of standardized tooling



ML is promising but requires specialized expertise



Complex postprocessing to adapt and integrate ML output with downstream systems

Amazon Bedrock Data Automation



A generative AI-powered capability to transform multi-modal unstructured content from documents, images, video, and audio into structured data easily, accurately, and at scale.

Amazon Bedrock Data Automation

Reduces development time by up to **80%** with industry leading accuracy at up to **35%** lower cost compared to Off the shelf FM Alternatives



Accurate and trustworthy output helps reduce reliance on manual reviews with confidence



Simpler development effort allows you to focus on building differentiated applications at a lower cost, no ML expertise necessary

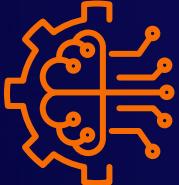


Easily tailor output requirements to diverse business-specific needs and application types

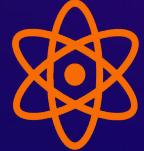
Key Features



Simple and intuitive interface to define output schemas and fine-grained business rules



Orchestration across state-of-the-art task-specific models and foundation models to generate highly accurate, consistent output



Built-in Responsible AI with visual grounding, confidence scores, and toxic content detection



Integration with Amazon Bedrock features and Knowledge base



Single inference API to handle production scale

How it works

TYPES OF OUTPUT THAT BDA CAN RETURN

Standard Output

- ✓ Linearized text representation of asset
- ✓ Gen-AI optimized output: reading / viewing order, semantically related output groupings, etc.
- ✓ Controls to optimize output based on downstream systems with simple selection knobs
- ✓ Automatic modality routing based on semantic modality, *not just file type*

Supported Modalities for Standard Output



Custom Output

- ✓ Developer supplied schema based on your downstream systems (Blueprint)
- ✓ Supports tasks such as extraction, key and value normalization, transformations, reasoning, splitting and classification
- ✓ Simple interface to define business rules and task logic for each field
- ✓ Console-based assistant to create bespoke blueprints in minutes with sample and desired output description

Supported Modalities for Custom Output



Getting Started with BDA is Easy

INPUTS / OUTPUTS

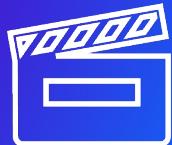


Key usecases

What are the types of applications where multi-modal content processing is needed?



Intelligent
Document Search



Media Asset Analysis
& Monetization



Intelligent Speech
Analytics



Multimodal Intelligent search



Agentic AI Workflow Automation

BDA for Intelligent Document Processing



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Document Processing across Industries

Financial Services



Loan applications, Tax processing, Accounts Payable automation, Financial statement analysis, Identity verification (KYC), etc.

Healthcare & Life Sciences



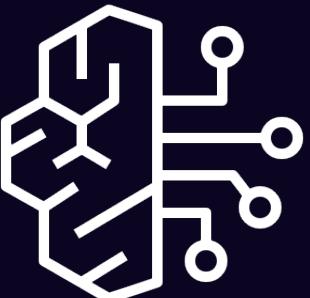
Claims Processing, Patient Enrollment, Medical Billing, Grievances and Appeals, etc.

Public Sector



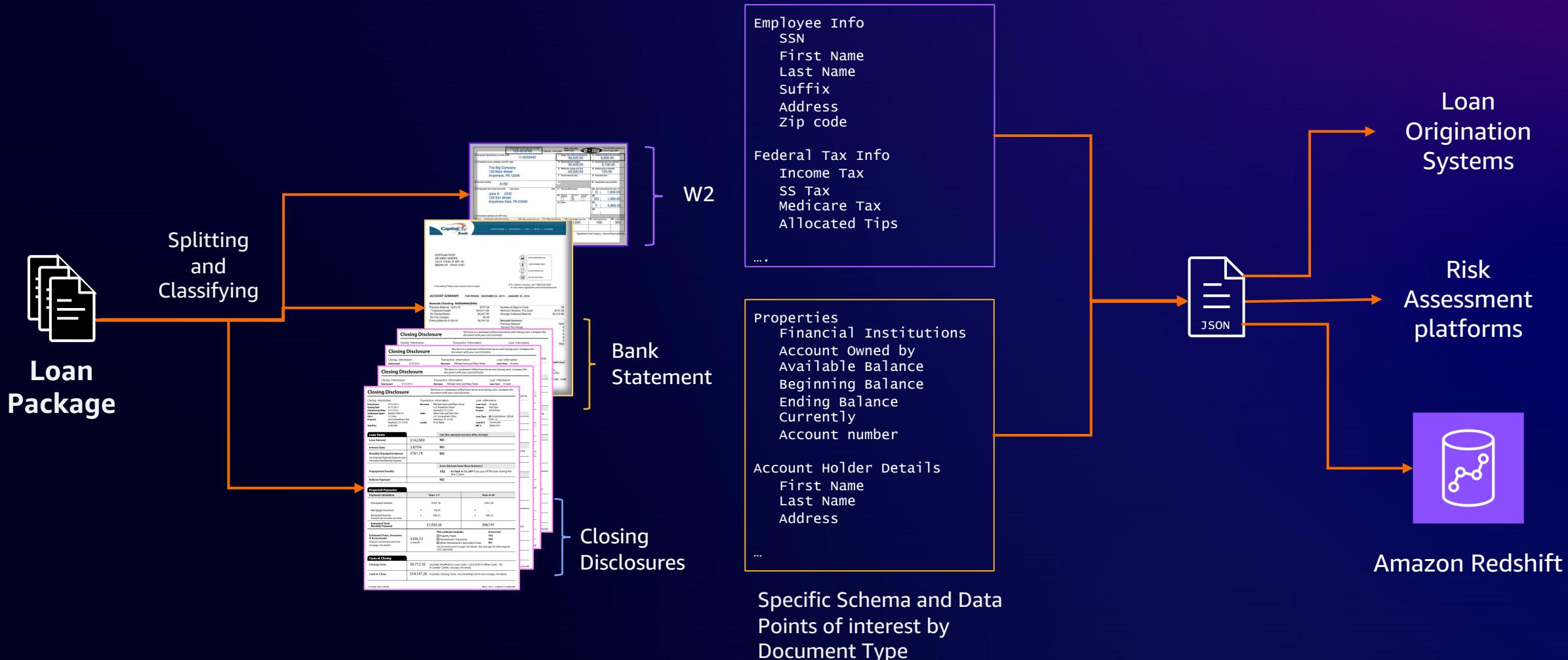
Small Business loans, Benefits processing, Regulatory approvals etc.

The Typical IDP Workflow



**Amazon Bedrock
Data Automation**

Use Case: IDP Loan Processing



Putting it together for IDP



Demo – BDA for Media Analysis



Demo – Media Analysis

Bedrock Media Demo X Bedrock Media Demo X Amazon Bedrock | us-west-2 X + https://d1xvhy22zwmw77y.cloudfront.net/index.html?demo=ca&id=a5214d51-6d89-4e93-9ee5-4d5819abaed5

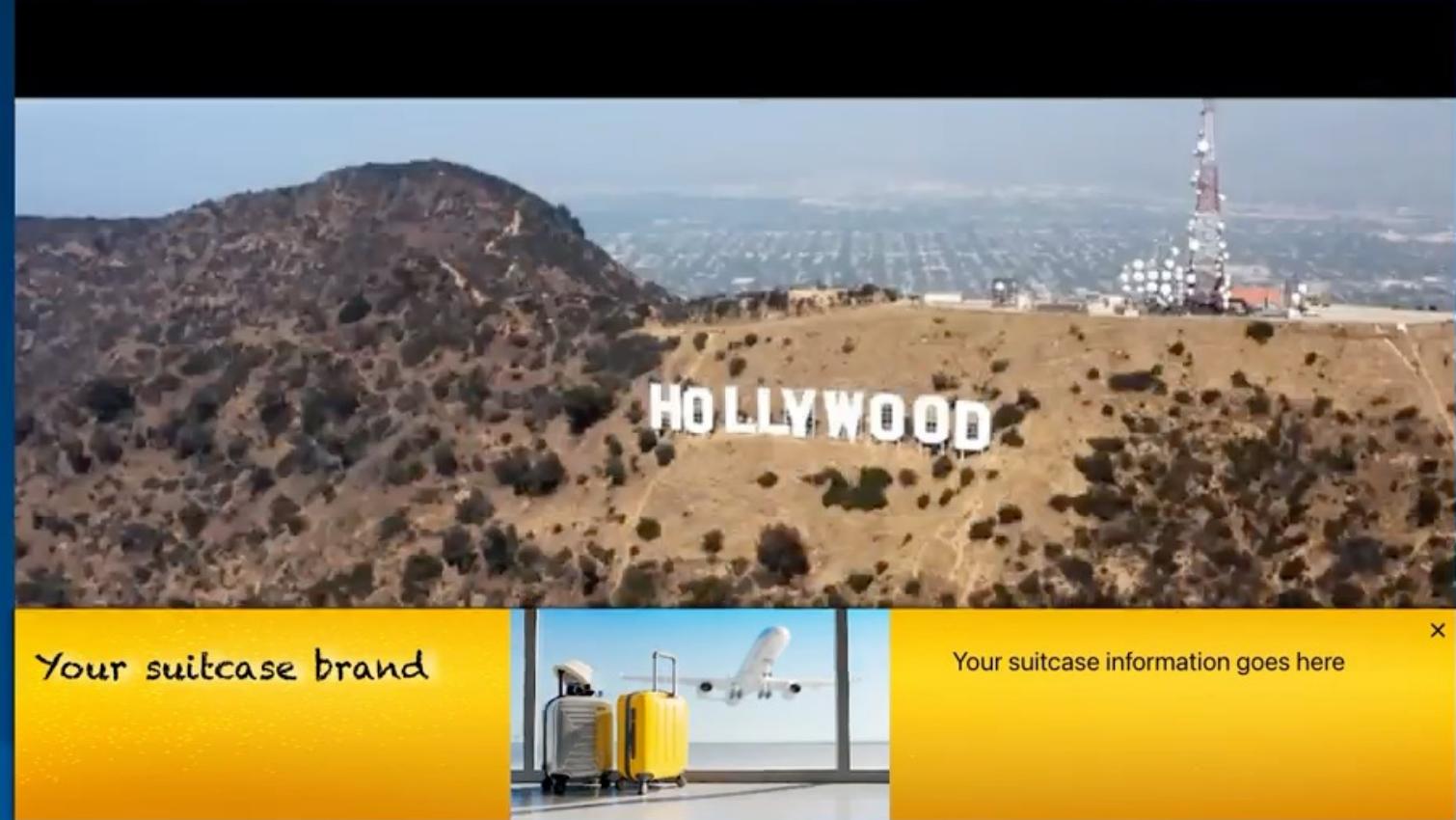
Contextual Advertising on CTV

Info
This demo illustrates how viewers see ads while streaming videos on their CTV devices. It leverages Amazon Bedrock Data Automation to analyze the video, breaking it into scenes with IAB advertising categories. Using the identified IAB and video context, Amazon Nova generates relevant ad content.

Current cuepoint
Scene 0 [00:00:00 - 00:00:00]
Travel

The video showcases a nighttime cityscape of Los Angeles, California. It features a panoramic view of the city skyline, with numerous illuminated skyscrapers and high-rise buildings dominating the landscape. The scene is dominated by the iconic Dodger Stadium, a large baseball stadium located in the heart of the city. The city's infrastructure, including roads, bridges, and street signs, is also visible, providing a sense of the urban environment. The overall atmosphere is vibrant and energetic, capturing the essence of a bustling metropolitan city.

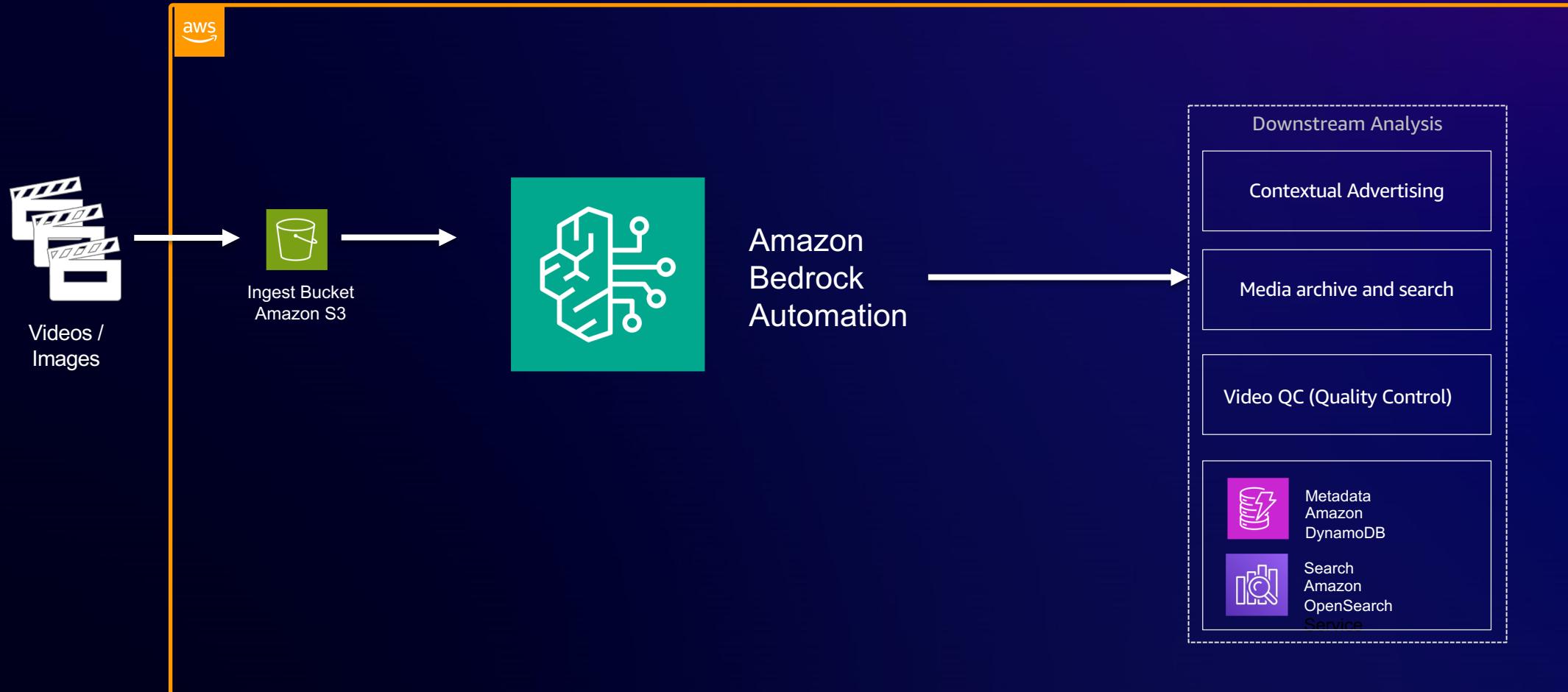
Next cuepoint
Scene 2 [00:00:01 - 00:00:05]
Travel



Your suitcase brand

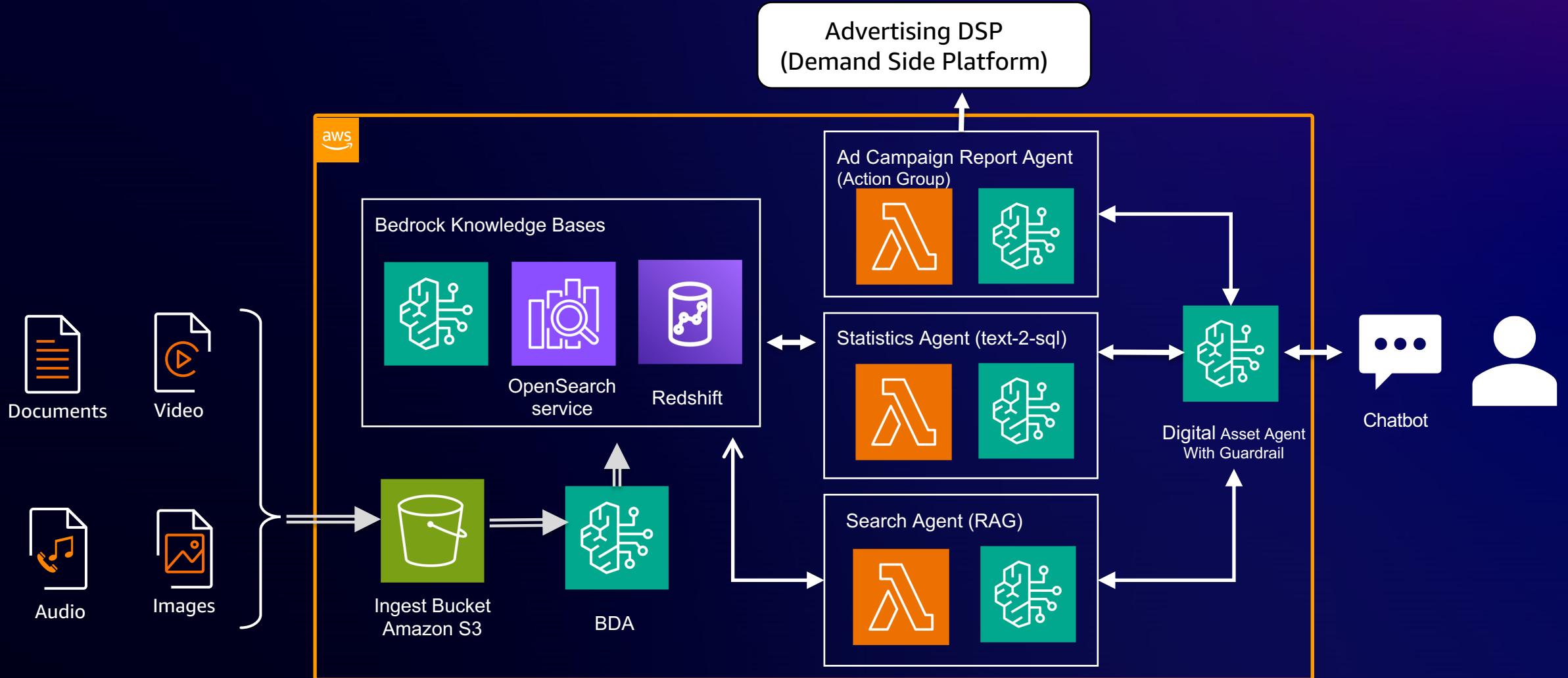
Your suitcase information goes here

Media Analysis with BDA



From BDA to Agentic workflows:

WITH BEDROCK KNOWLEDGE BASES, MULTI-AGENTS, GUARDRAIL, AND AMAZON NOVA





GENAI BOOTCAMP

Notebook Demo



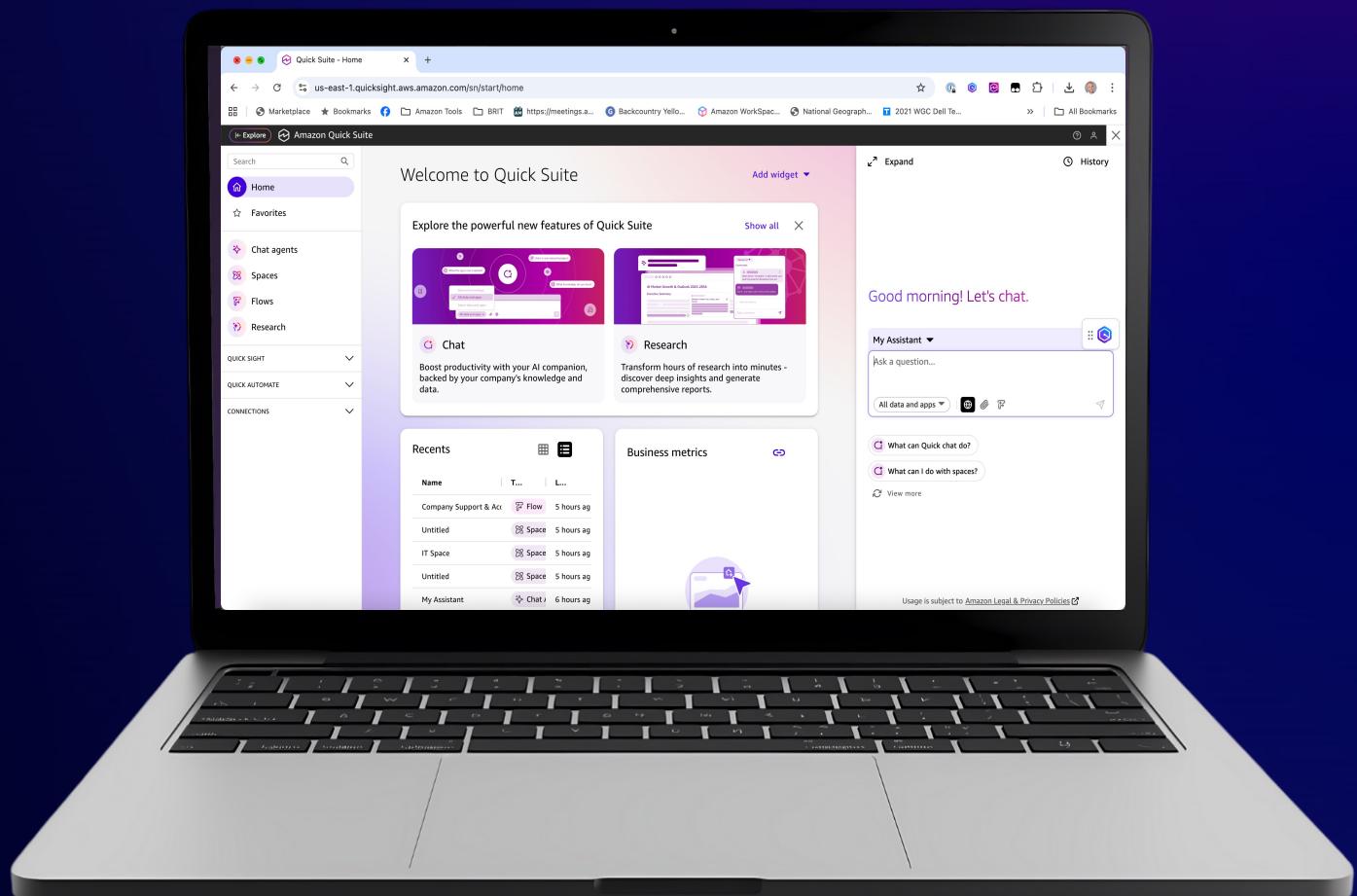
Amazon Quick Suite

Introducing Amazon Quick Suite

Unify your data sources to get clear, context rich answers

Take action alongside agentic AI teammates

Experience an intuitive AI experience your company will love



Amazon Quick Suite



Spaces

- Organize files, dashboards, and data sources
- Collaborate workspaces for your projects



Chat Agents

- Custom AI assistants with your business knowledge
- Share and collaborate across teams



Research

- Deep-dive analysis and comprehensive reports
- Professional, exportable documentation



Quick Sight

- Business Intelligent and data visualization
- Interactive dashboards and analytics



Flows

- Automate repetitive tasks with pre-defined steps
- No-code workflow automation



Automate

- Complex, multi-step workflow automation
- Transform entire business processes

Governance,
data security

Access controls

Guardrails

Responsible AI

Regulatory compliance

Company data: Spaces and datasets



40+ Data connectors



User uploaded files



Quick sight data

World Knowledge



Bedrock models



Web search



Actions in 3P apps

Save time searching thanks to **Indexed Data**



Quick Suite can connect to information across internal repositories, documents, applications, and services

Quick Suite unites teams and agents

Create a custom **chat agent** that understands your unique role and your team's specific needs

 **Describe your team's needs in natural language**

 **Upload your team's relevant data sets to receive accurate, context-rich answers**

 **Get recommendations on how to improve your agent's actions and responses**

Training Guide

"Create an agent that helps new team members learn about our processes and tools. Should answer common questio..."

Project Coordinator

"Build an agent that helps track project status, find relevant documentation, and assist with common project tasks. Sho..."

Use AI to collaborate with your team with **spaces** in Quick Suite



Spaces allow you to create a custom knowledge hub for your team

Add knowledge ▾

Upload your team's specialized data, documents, and resources



You can focus agents on the data contained within the space

Dive deep on your business with **Quick Sight**

- Perform advanced data analysis
- Build dashboards
- Generate rich data stories



Become an expert on a new topic with Quick Research

- Tackle complex and time-intensive research
- Leverage all your data, private third-party data, and data from the internet

Every user can save time on routine tasks

Quick Flows

- Save time and energy by automating repetitive tasks
- Create your own flows using natural language

Automate complex workflows with Quick Automate

- Bring multiple agents together to streamline complex business processes

Welcome to the 2nd day

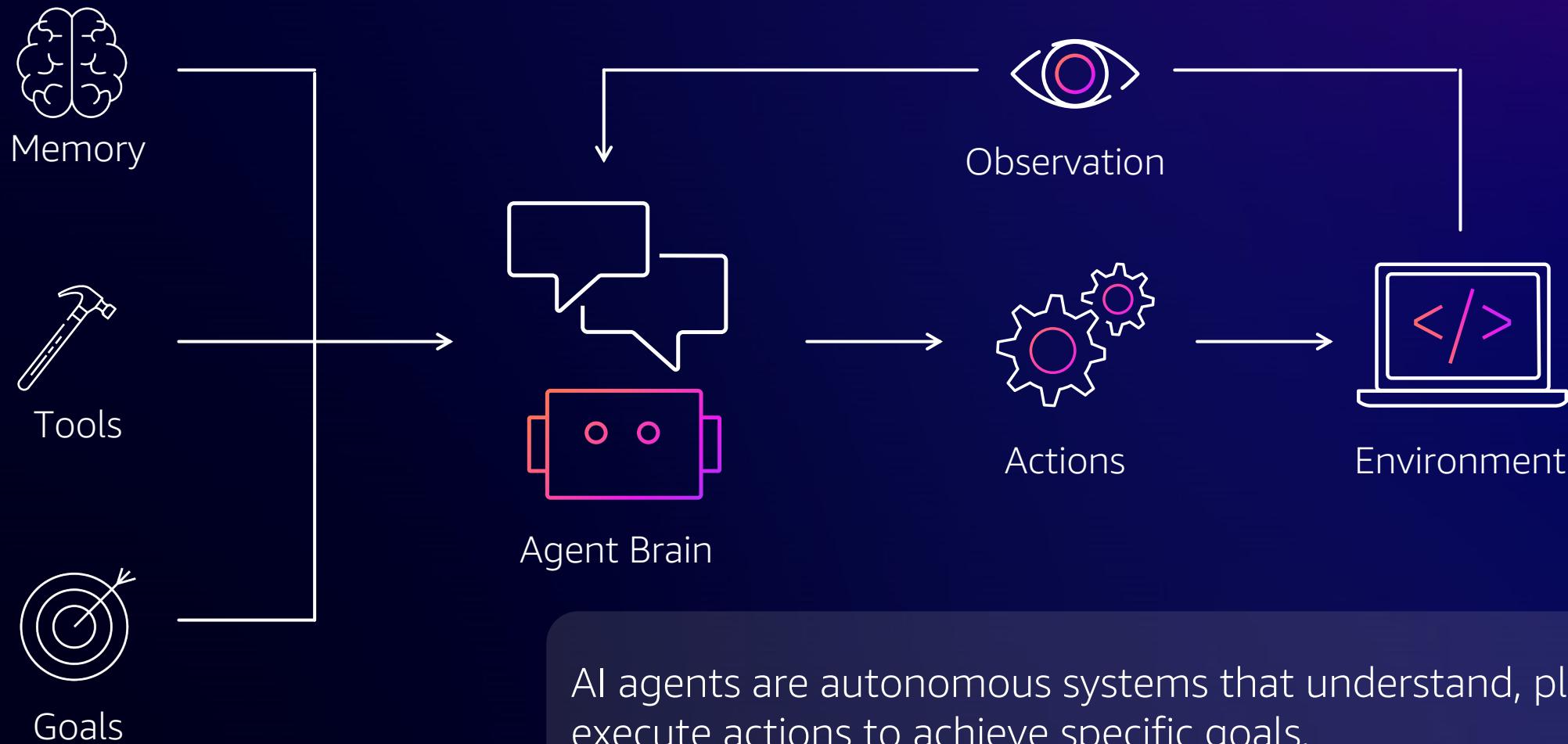


Intro to Agentic AI



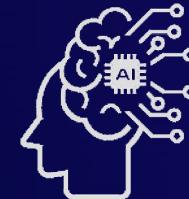
© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

What is agentic AI?



“By 2028, at least
15% of daily work decisions will be
made by
AI agents”

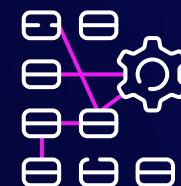
Gartner, “Top Strategic Technology Trends: agentic AI – The evolution of Experience” February 2025



Model Reasoning & Tool execution



Price-performance



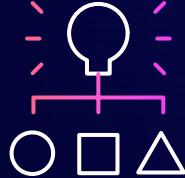
Adaptive Automation

Spectrum of Agentic AI



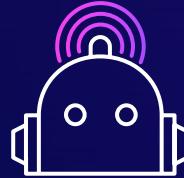
SIMPLE ASSISTANTS

- Basic query responses & actions
- Single step tasks
- Model as Tool caller



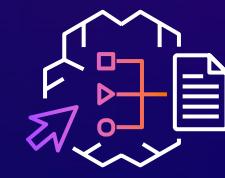
DETERMINISTIC AGENTS

- Multistep tasks
- Plan is predetermined
- Model as Flow Router



AUTONOMOUS AGENTS

- Self directed Agents
- **Expanded toolset – Computer/Browser**
- **Mimic human actions**



AGENTIC VIRTUAL WORKERS

- Multi-agent coordination
- Cross Org boundaries
- **Long running systems**
- **Mimic human teams**

INCREASING AUTONOMY AND BUSINESS IMPACT

Agentic AI use cases

ACTION ENABLED CHAT AGENTS



Real-time chat interactions with API capabilities

WORKFLOW AUTOMATION AGENTS



Event-triggered task automation

AGENTIC PLATFORMS



Create and manage reusable agents with enterprise-wide controls and governance

When to use (and not use) agentic AI

Use agentic AI for:

- Complex workflows with variable paths and decisions
- Coordination across multiple tools and data sources
- Tasks requiring understanding context and adapting responses
- Long-running processes needing memory across sessions



Traditional solutions are Better for:

- Simple, single-step automation tasks
- Ultra-low latency operations (microsecond responses)
- Mission-critical systems requiring 100% deterministic results

Agentic use cases across industries

AUTOMOTIVE

- AI-powered development and validation
- Intelligent contact center
- Supply chain optimization

ENERGY & UTILITIES

- Energy trading
- Renewable energy planning
- Upstream production optimization

FINANCIAL SERVICES

- Insurance underwriting optimization
- Intelligent customer experience
- Investment research assistant

GAMES

- AI for game production
- Code assistance for game development
- In-game AI

MANUFACTURING

- Research intelligence and discovery
- Maintenance management
- Supply chain optimization

TELECOM

- Customer experience
- Enterprise productivity
- Network automation

MEDIA & ENTERTAINMENT

- Multi-agent collaboration for video understanding
- Media planning and campaign optimization
- Multi-agent collaboration for advertising sales data

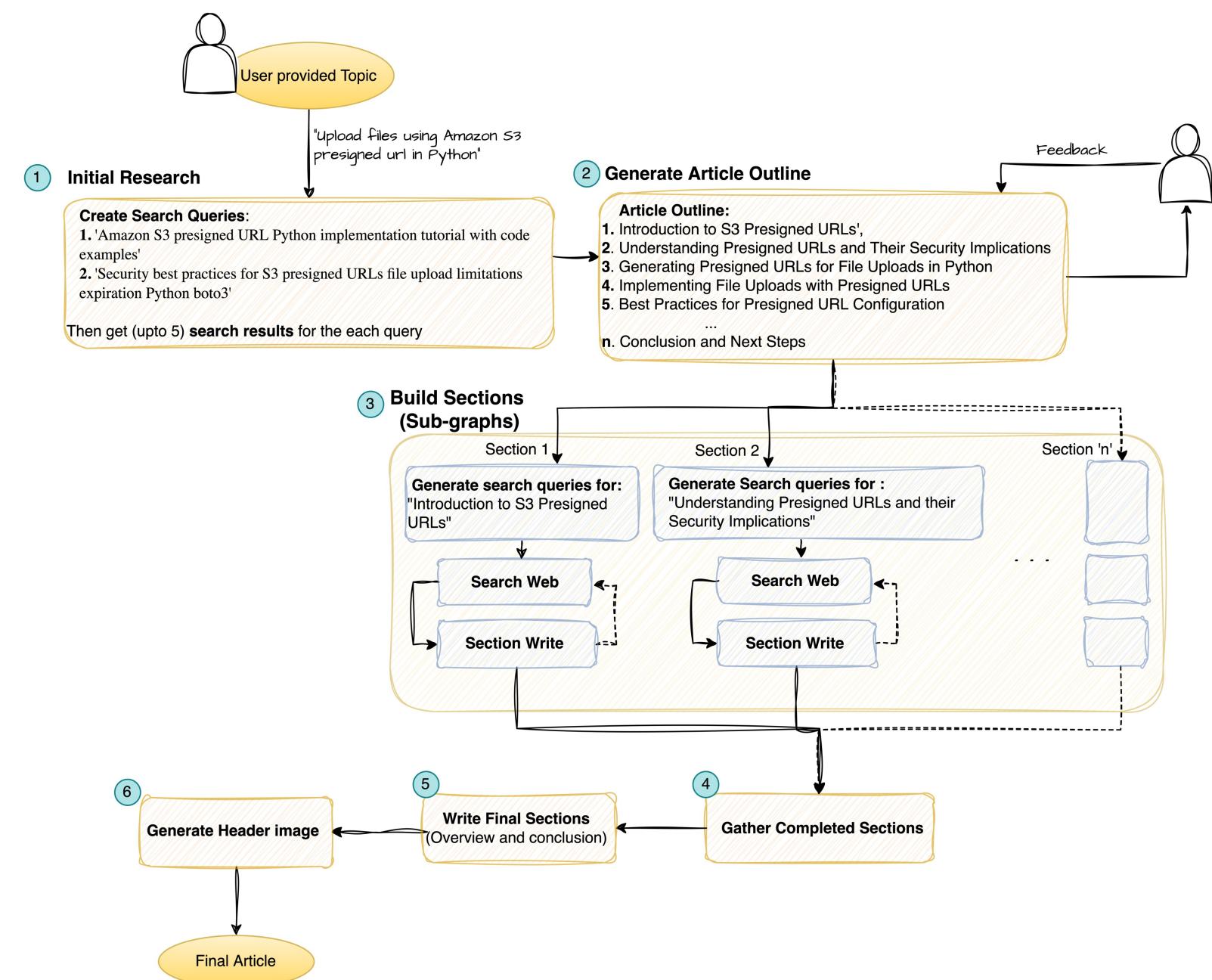
RETAIL & CONSUMER GOODS

- Agentic shopping assistants
- Supply chain reallocation
- Pricing agents

HEALTHCARE & LIFE SCIENCES

- Clinical workflow automation
- Prior authorization automation
- Patient and member engagement

Common Deep Research Agent



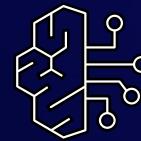
The broadest choice for building and deploying agents

SPECIALIZED



Amazon Q
Pre-built products that
use agents for enhanced
productivity

FULLY-MANAGED



Amazon Bedrock Agents
with built-in FM-powered
orchestration

DIY



Strands Agents
Simple, flexible, and
lightweight open source
SDK for building agents

OUT - O F - T H E - B O X A G E N T S

TOOLS FOR BUILDING AGENTS

Example : Octank Ecommerce – Orders Assistant

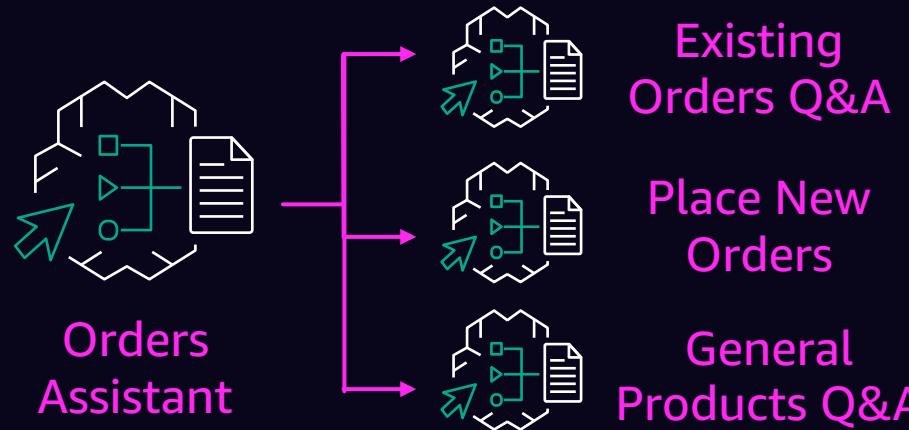
Before



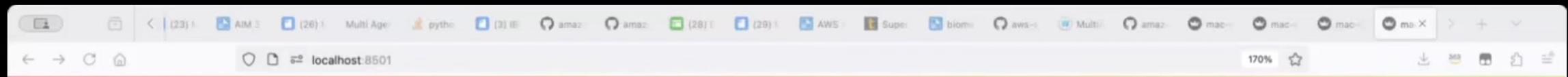
- Voice menus, long waits
- Heavy call center volume
- No chat experience
- Simple web site
- Limited data, static views



With multi-agent collaboration



- Answers in seconds
- Flexible and personalized
- Call center relief
- Extensible



Deploy

Mortgages Assistant

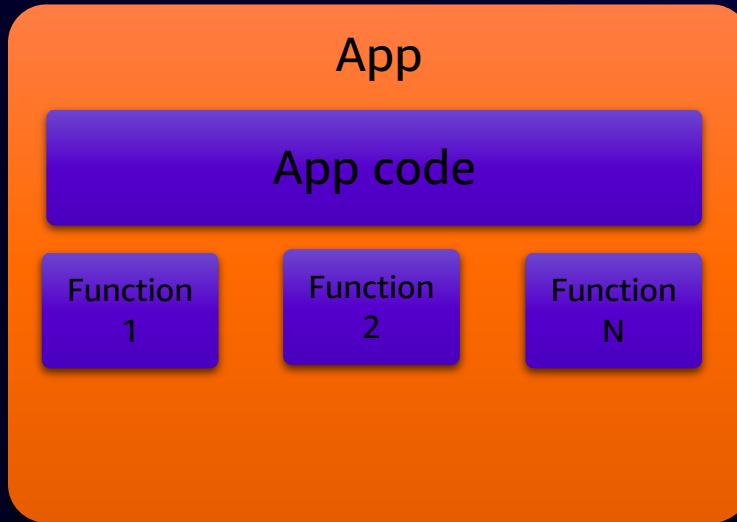
I'm your mortgages assistant. How can I help today? >

Understanding Agents



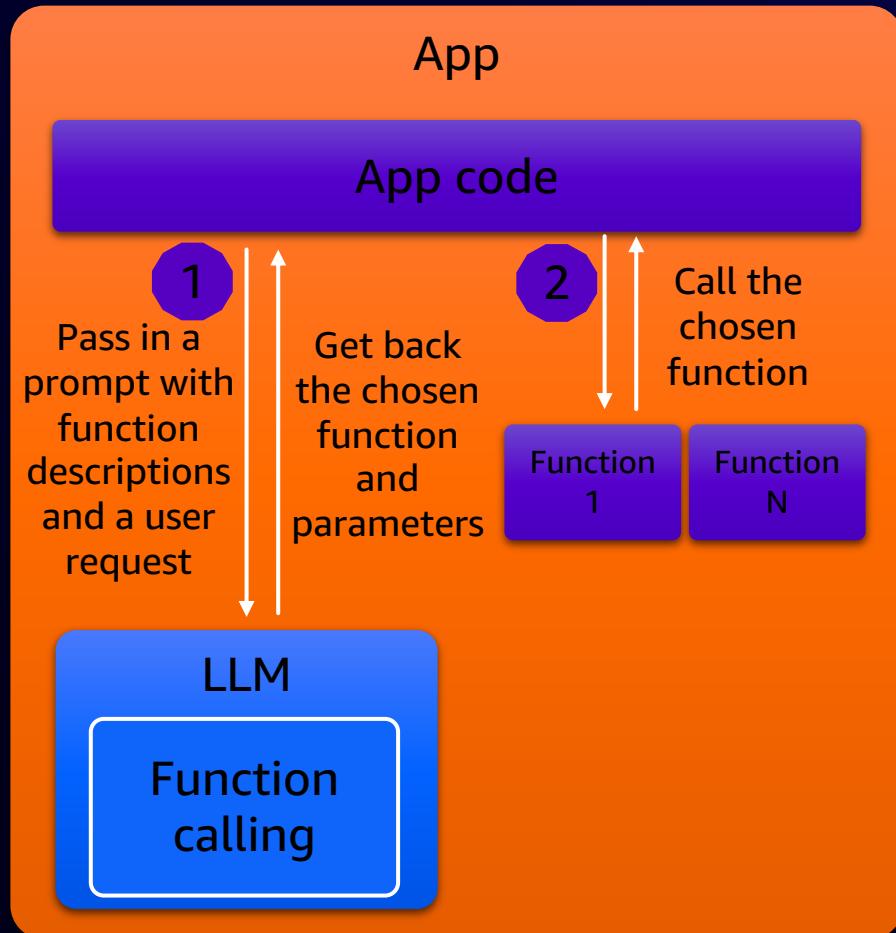
© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

From Apps to Function calling to Agents...



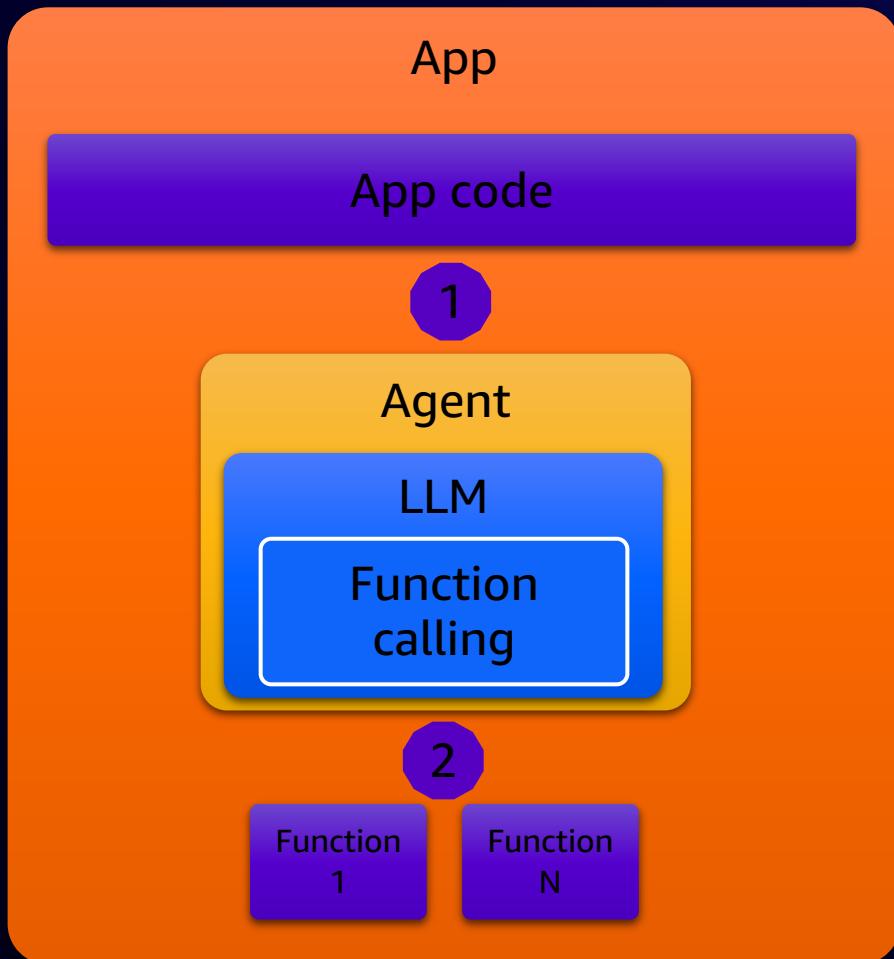
- Has been possible for decades
- Tried and true
- App is coded to accomplish specific tasks and knows what functions to call and how to call them

From Apps to Function calling to Agents...



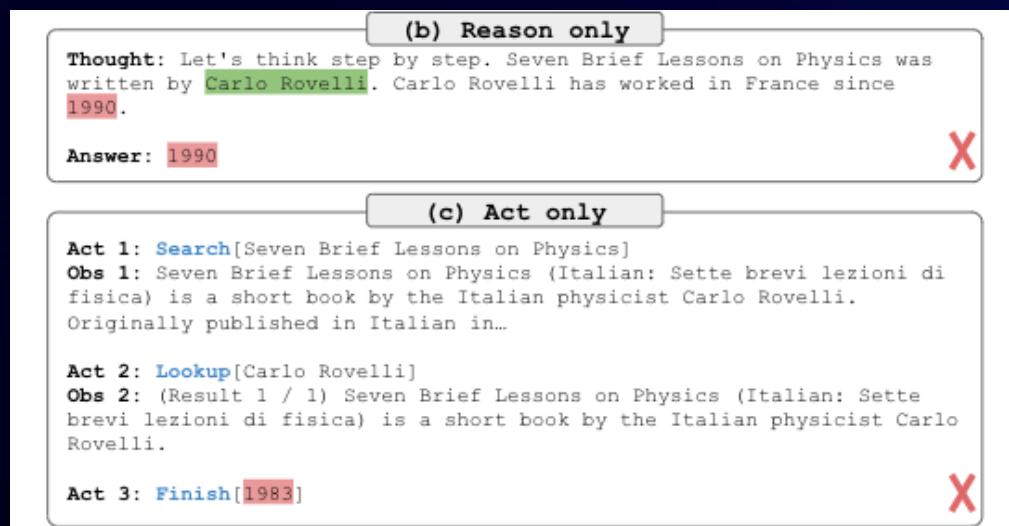
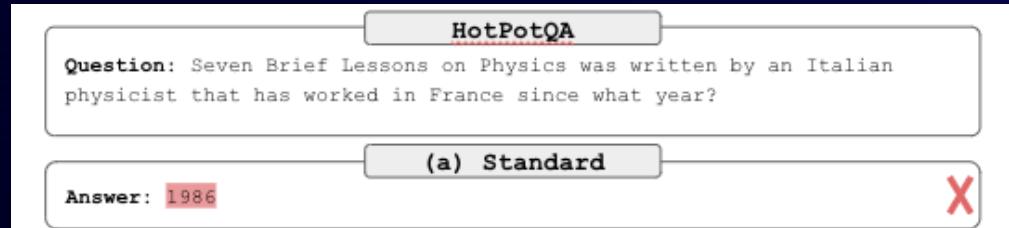
- Now we give the LLM rich descriptions of our available functions and let the LLM tell the app what to call next and with what parameters
- Now a single “app” can handle a full range of tasks without one-off coding

From Apps to Function calling to Agents...

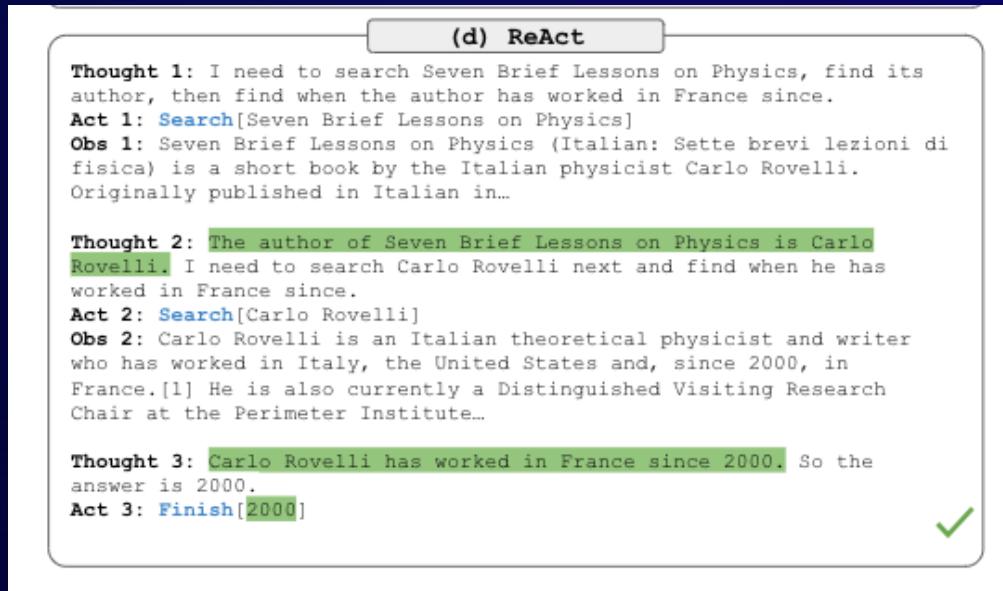


- Now I can create Agents, give them instructions and available actions, and voila, automation with even less coding and more flexibility
- Agent decides what actions, in what order, and even executes those actions
- Many possible implementations - LangChain Agents, Agents for Bedrock

ReAct (Reasoning + Action) Prompting



Combine text reasoning and actions in a single model



ReAct

ReAct is a general paradigm that combines reasoning and acting with LLMs. ReAct prompts LLMs to generate verbal reasoning traces and actions for a task.

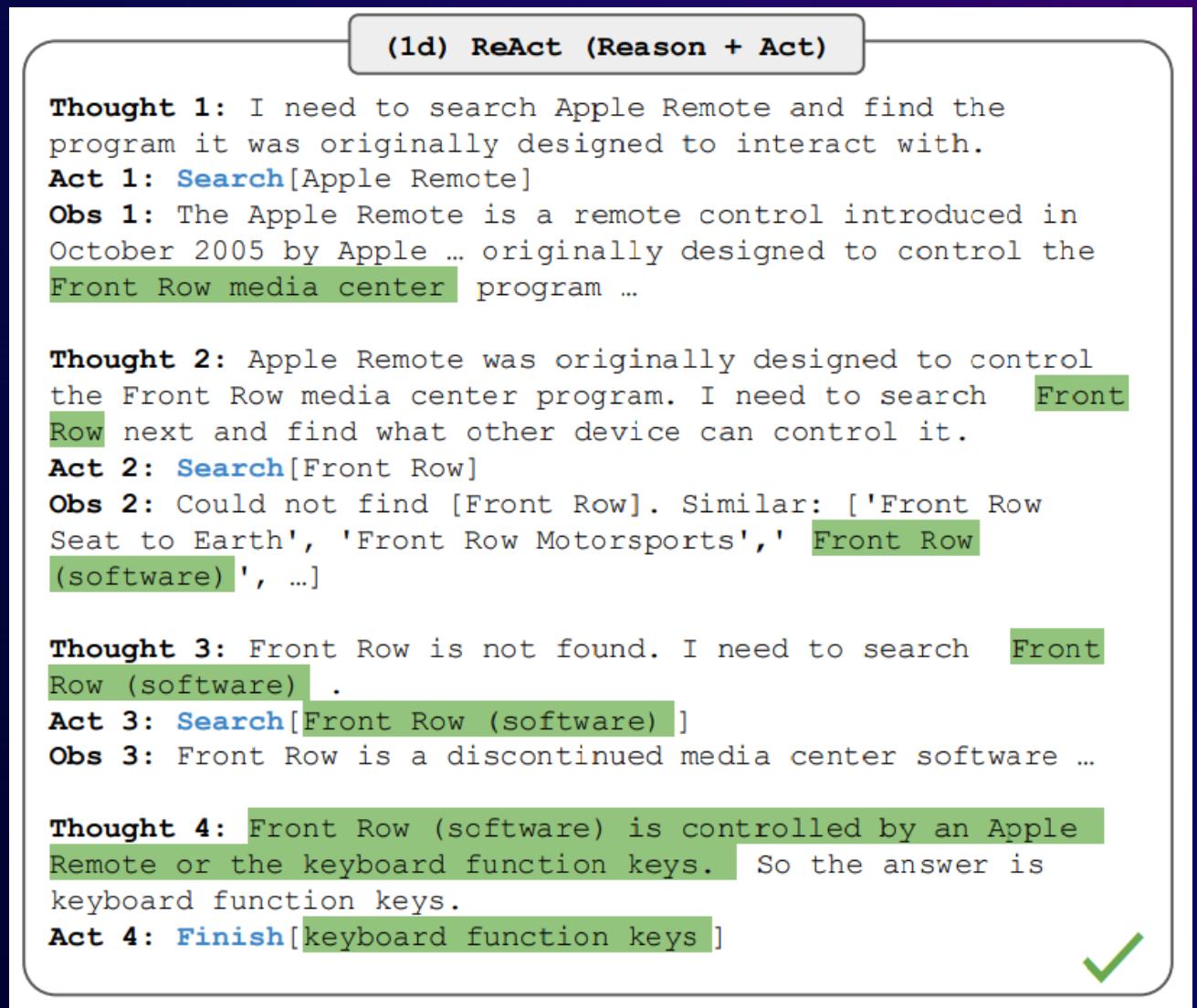


Image Source: [Yao et al., 2025](#)

ReAct Agents as a loop of LLM calls

You can think of LLM agents as a simple while loop, where an LLM decides which tools to be used, is presented with the results of that tool execution, until it finds the final answer or reaches a limit of number of iterations.

ReAct

```
In [3]:  
import json  
import sys  
  
folder = './prompts/'  
prompt_file = 'fever.json'  
with open(folder + prompt_file, 'r') as f:  
    prompt_dict = json.load(f)  
  
webthink_prompt = prompt_dict['webthink_simple3']  
  
def webthink(idx=None, prompt=webthink_prompt, to_print=True):  
    question = env.reset(idx=idx)  
    if to_print:  
        print(idx, question)  
    prompt += question + "\n"  
    n_calls, n_badcalls = 0, 0  
    for i in range(1, 8):  
        n_calls += 1  
        thought_action = llm(prompt + f"Thought {i}:", stop=[f"\nObservation {i}:"])  
        try:  
            thought, action = thought_action.strip().split(f"\nAction {i}: ")  
        except:  
            print('ohh...', thought_action)  
            n_badcalls += 1  
            n_calls += 1  
            thought = thought_action.strip().split('\n')[0]  
            action = llm(prompt + f'Thought {i}: {thought}\nAction {i}:', stop=[f'\n']).strip()  
        obs, r, done, info = step(env, action[0].lower() + action[1:])  
        obs = obs.replace('\n', '')  
        step_str = f"Thought {i}: {thought}\nAction {i}: {action}\nObservation {i}: {obs}\n"  
        prompt += step_str  
        if to_print:  
            print(step_str)  
        if done:  
            break  
    if not done:  
        obs, r, done, info = step(env, "finish[]")  
    if to_print:  
        print(info, '\n')  
    info.update({'n_calls': n_calls, 'n_badcalls': n_badcalls, 'traj': prompt})  
    return r, info
```

Source: [Yao et al., 2025](#)

ReAct Agents as a loop of LLM calls

```
REACT_PROMPT_TEMPLATE = """  
Answer the following question by thinking step by step.  
For each step, follow this format:
```

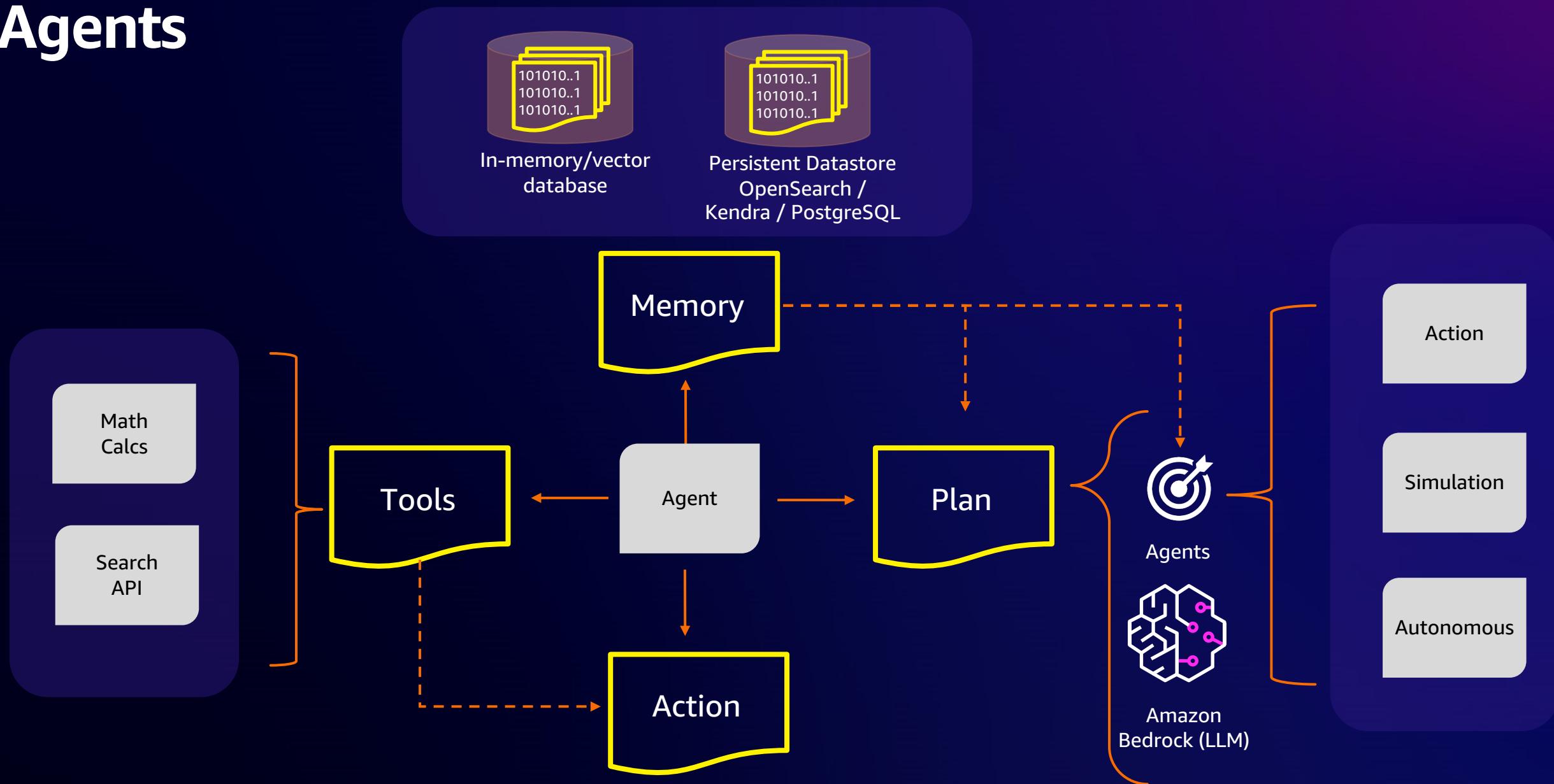
```
Thought [step number]: Reason about the current state and what you need to know.  
Action [step number]: Choose from: Search[query], Lookup[term], Calculate[expression], or Finish[answer]  
Observation [step number]: (This will be filled with the result of your action)
```

```
Question: {question}  
"""
```

You can think of LLM agents as a simple while loop, where an LLM decides which tools to be used, is presented with the results of that tool execution, until it finds the final answer or reaches a limit of number of iterations.

```
def simple_react_agent(question, prompt_template, max_steps=3):  
    """  
    A simplified ReAct agent implementation  
  
    Args:  
        question: The question to answer  
        prompt_template: Template with instructions for the LLM  
        max_steps: Maximum number of reasoning steps  
    """  
    # Initialize the conversation with the formatted prompt  
    history = prompt_template.format(question=question)  
  
    for step in range(1, max_steps + 1):  
        # Get the next thought and action from LLM  
        response = llm(history + f"\nThought {step}:")  
  
        # Parse thought and action from response  
        thought_action = response.strip()  
        thought = thought_action.split("Action [step]:")[0].strip()  
        action = thought_action.split("Action [step]:")[1].strip() if "Action [step]:" in thought_action else ""  
  
        # Execute the chosen tool based on the action  
        if "Search[" in action:  
            tool_input = action.split("Search[")[-1].split("]")[0]  
            observation = search_tool(tool_input)  
        elif "Finish[" in action:  
            return action.split("Finish[")[-1].split("]")[0]  
        else:  
            observation = "Error: Invalid action format"  
  
        # Update history with the complete step  
        history += f"\nThought {step}: {thought}\n"  
        history += f"Action {step}: {action}\n"  
        history += f"Observation {step}: {observation}\n"  
  
        # Check if we've reached a conclusion  
        if "Finish[" in action:  
            break  
  
    return "Failed to reach conclusion within step limit"
```

Agents



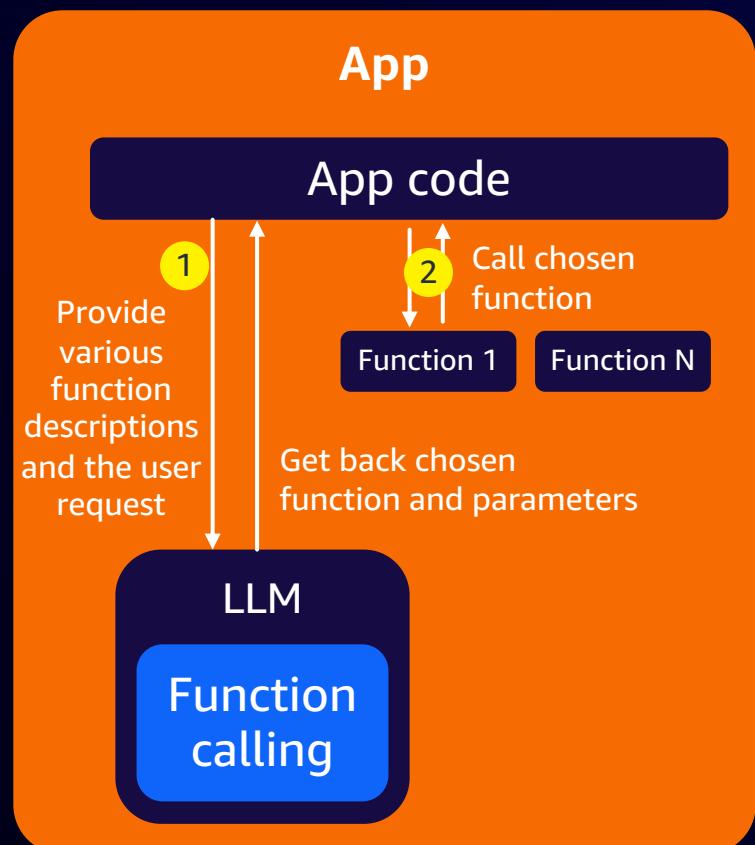
Function calling, Agents, and MCP



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

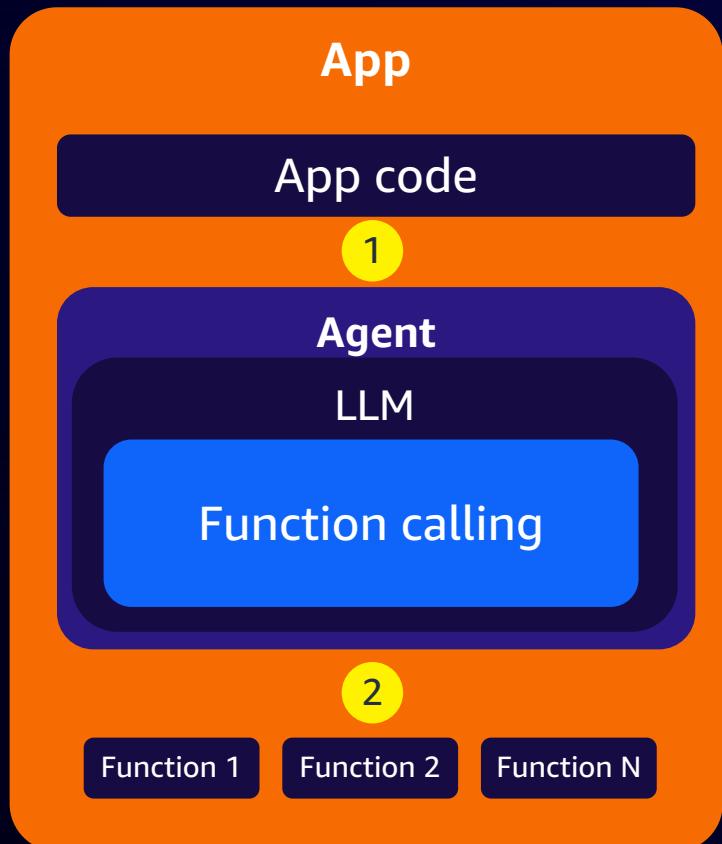
Function calling

LLMs decides which functions to call and how

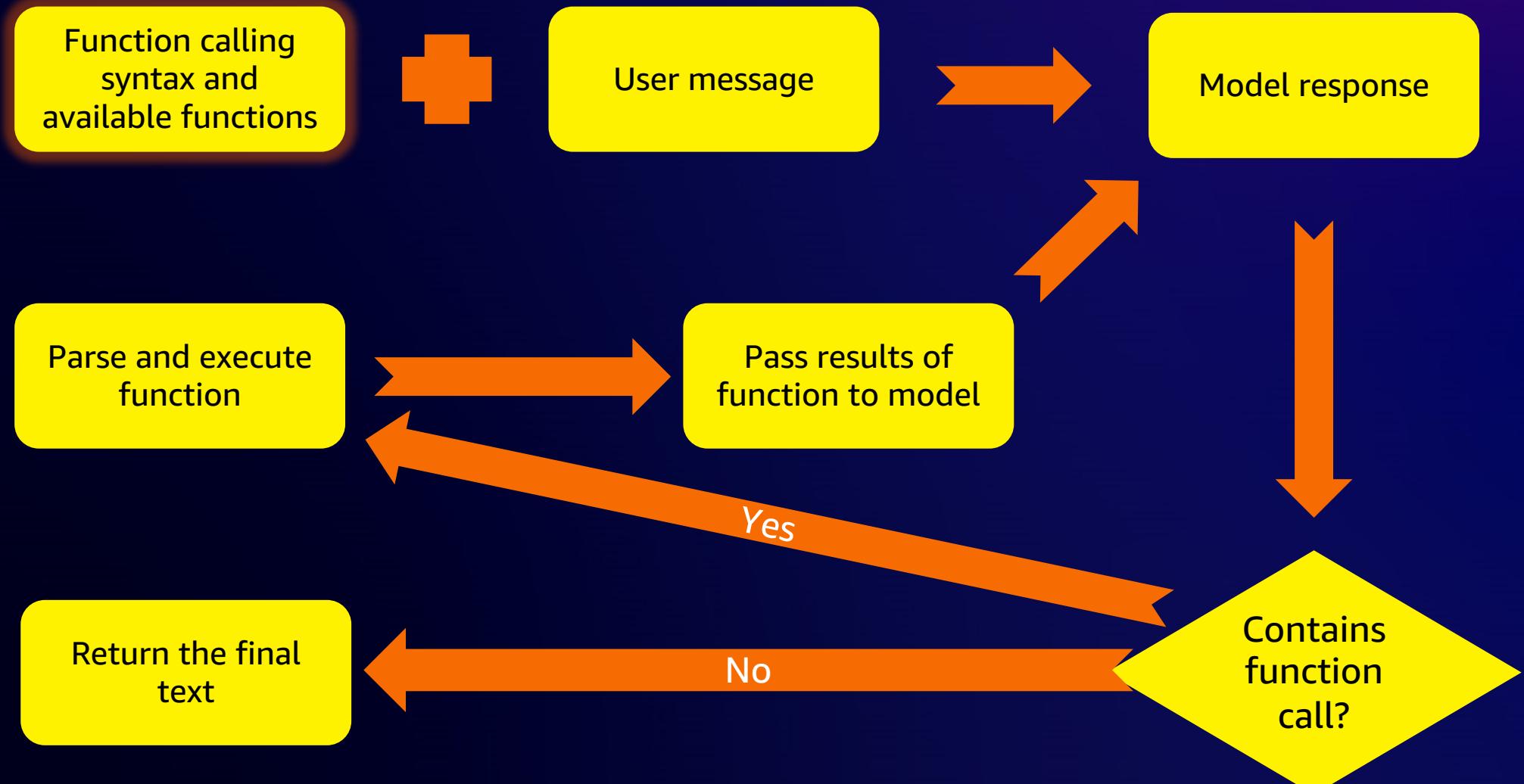


Agents

Let Agent chooses the functions, and also execute them



How function calling works: A closer look



In detail..

WE DEFINE THE FUNCTIONS LLM
CAN SELECT FROM

LLM selects a function

We pass the results back
to LLM

LLM uses the results to
create its answer

The screenshot shows a conversational interface with three message cards:

- User Profile:** You're a helpful assistant. Your training cut off is September 2022. Today's date is Tuesday 8th April, 2024. Be super succinct in your responses. You have access to three tools:
 - get_weather(location)
 - order_pizza(location, toppings)
 - send_email(message, address) Optionally, use these tools to answer the following question. Wait for me to provide results of the tool call after/if you call one:
What's the weather in London right now?
- Tool Call:** Let me check the weather in London.
get_weather("London")
- Result:** It's 21°C and sunny
- Final Response:** In London right now it's 21°C and sunny - a lovely spring day!

MCP



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

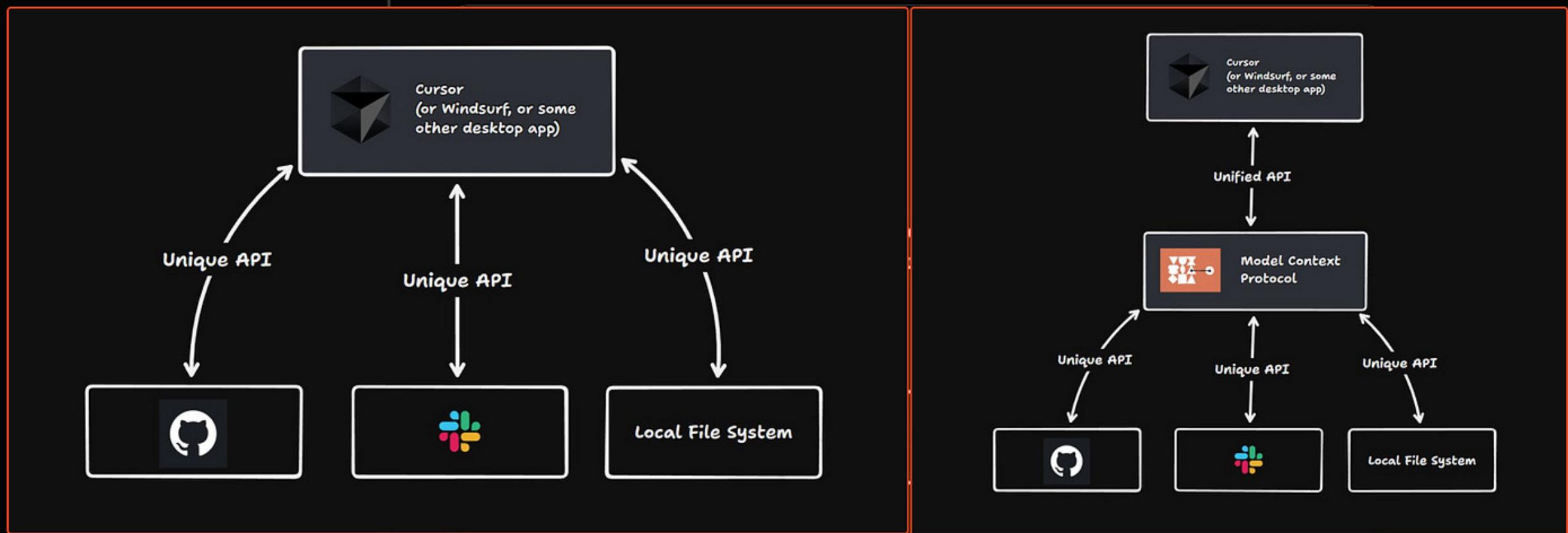


Matt Pocock ✅ @mattpocockuk · Mar 6

🔗 ...

Lots of folks feeling FOMO about MCP.

Here's the problem it solves:



249

695

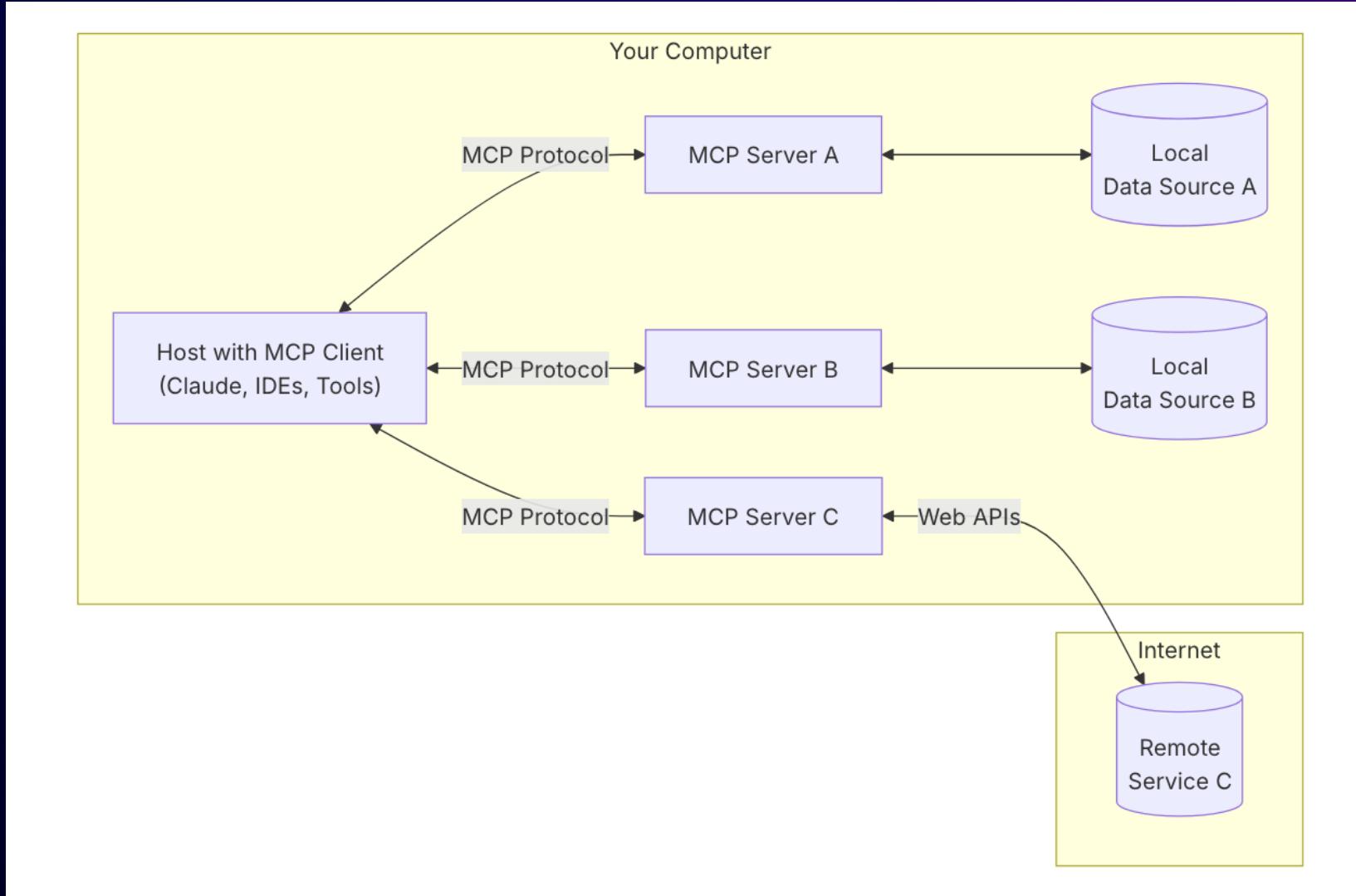
5.6K

656K

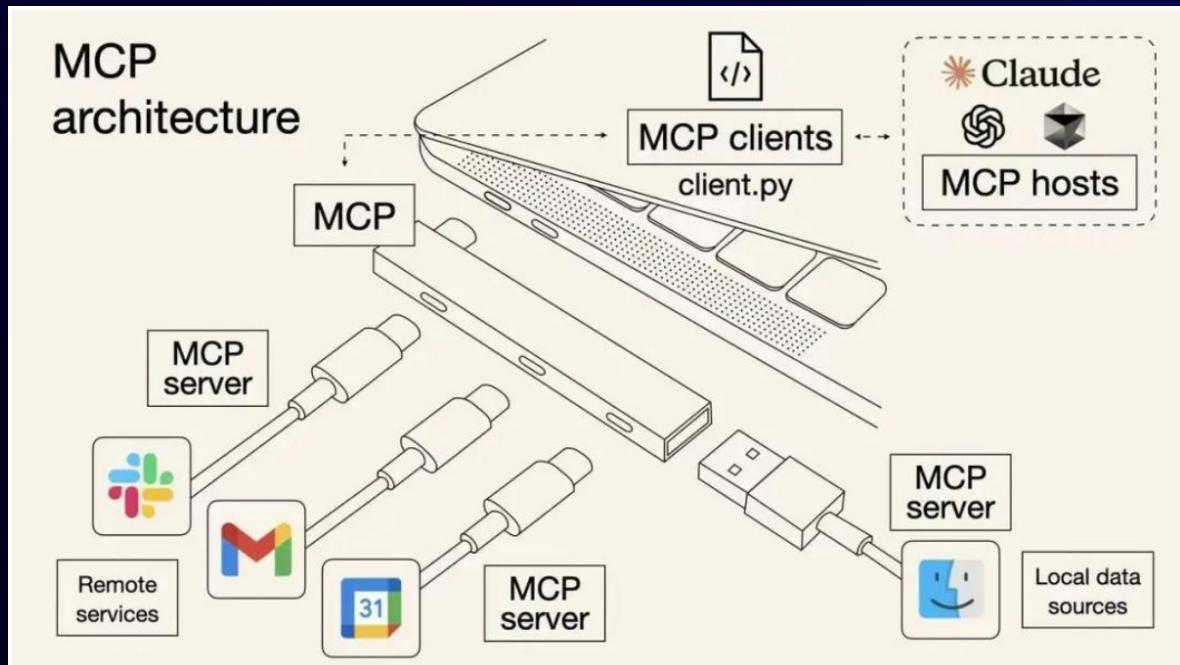
Bookmark Up

Enter MCP

1. Open standard by Anthropic solving AI integration complexity
2. Client-server protocol connecting AI to data and tools.
3. Enables seamless context switching across tools and datasets.
4. Includes SDKs and pre-built integrations (Drive, Slack, GitHub, Postgres).
5. Standardizes AI workflows for simplicity, consistency, and scale.



MCP in simple term



- MCP helps reduced complexity
- MCP decouples AI development from specific backend system
 - MCP client focuses on LLM, agent, application layer focused
 - MCP server focuses on data sources connectivity
- Hence, if you have new data sources, just connect to MCP server to make discoverable by MCP client (hence, AI agent)

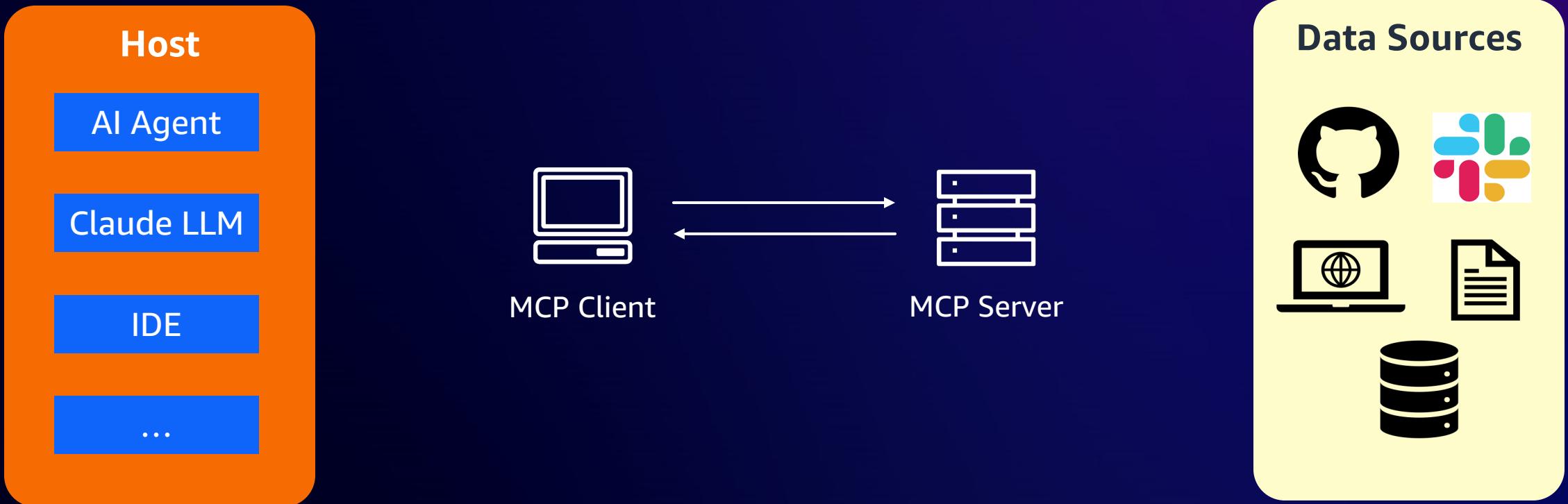
* for business users analogy only, as this is not accurate from technical perspective

Agents with MCP

Connect the host (app) to MCP client, and request execution



MCP main components



1. **Host:** A program, AI tool, or AI agents requiring access to data through MCP
2. **MCP Client:** Protocol clients maintaining one-to-one connections with servers
3. **MCP Server:** Lightweight programs exposing capabilities through standard MCP
4. **Data sources:** both local (databases, files) and remote services (APIs) that MCP server can access

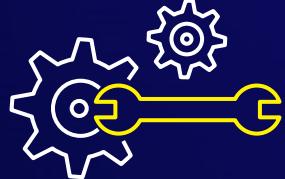
Why is MCP important



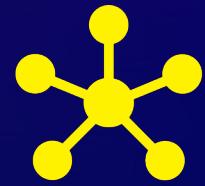
Standardization



**Integration with
External Data Sources**



**Enhanced
Functionality**



**Flexibility and
Scalability**

MCP Concept

Resources



- Allows Servers to expose data
- File contents, database records
- Application-controlled (client decides how and when to use)

Prompts



- Reusable prompt template and workflows
- User-controlled i.e., from servers to clients and then to user

Tools



- Enable LLMs to perform actions through server
- Model-controlled, expose from server to client

In summary ...

Aspect	Before MCP	After MCP
Integration effort	Custom code for each system	Standardized, reusable protocol
Data/tool access	Limited, siloed, often outdated	Real-time, multi-source, live
AI context/memory	Lost between tools and sessions	Persistent, shared across tools/models
User experience	Fragmented, repetitive, impersonal	Seamless, personalized, consistent
Vendor lock-in	High (integration tied to model/tool)	Low (model-agnostic, interoperable)
Developer productivity	Slow, error-prone, high maintenance	Fast, reliable, low maintenance
Ecosystem growth	Fragmented, duplicated effort	Network effects, rapid expansion

Strands Agents



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

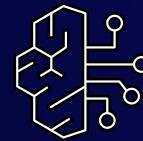
The broadest choice for building and deploying agents

SPECIALIZED



Amazon Q
Pre-built products that
use agents for enhanced
productivity

FULLY-MANAGED



Amazon Bedrock Agents
with built-in FM-powered
orchestration

DIY



Strands Agents
Simple, flexible, and
lightweight open source
SDK for building agents

OUT - O F - T H E - B O X A G E N T S

TOOLS FOR BUILDING AGENTS

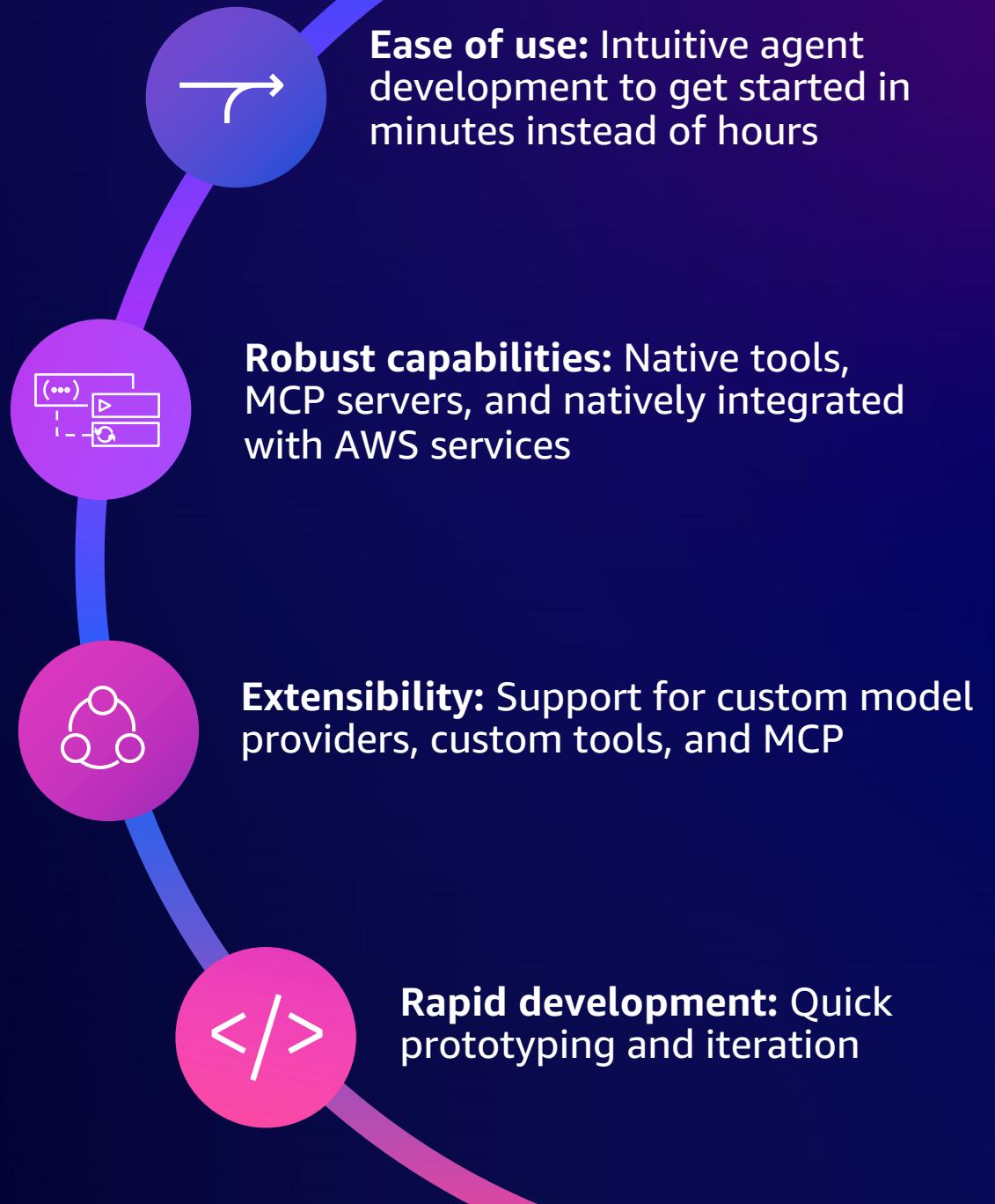
Framework	AWS integration	Autonomous multi-agent support	Autonomous workflow complexity	Multimodal capabilities	Foundation model selection	LLM API integration	Production deployment
AutoGen	Weak	Strong	Strong	Adequate	Adequate	Strong	Do it yourself (DIY)
CrewAI	Weak	Strong	Adequate	Weak	Adequate	Adequate	DIY
LangChain/ LangGraph	Adequate	Strong	Strongest	Strongest	Strongest	Strongest	Platform or DIY
Strands Agents	Strongest	Strong	Strongest	Strong	Strong	Strongest	DIY

<https://docs.aws.amazon.com/prescriptive-guidance/latest/agentic-ai-frameworks/comparing-agentic-ai-frameworks.html>



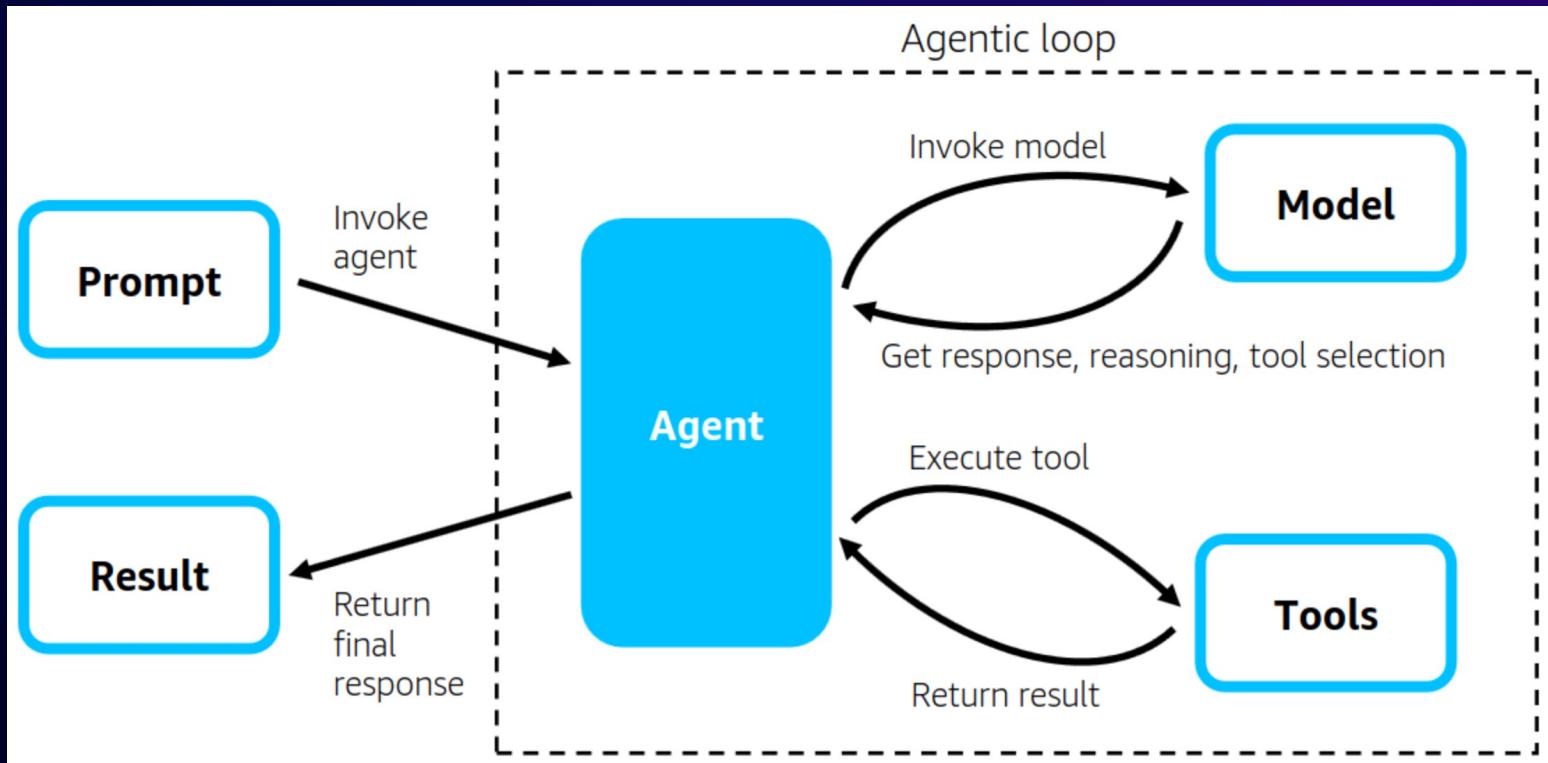
Strands Agents

Strands Agents is an open source SDK for building agents using just a few lines of code

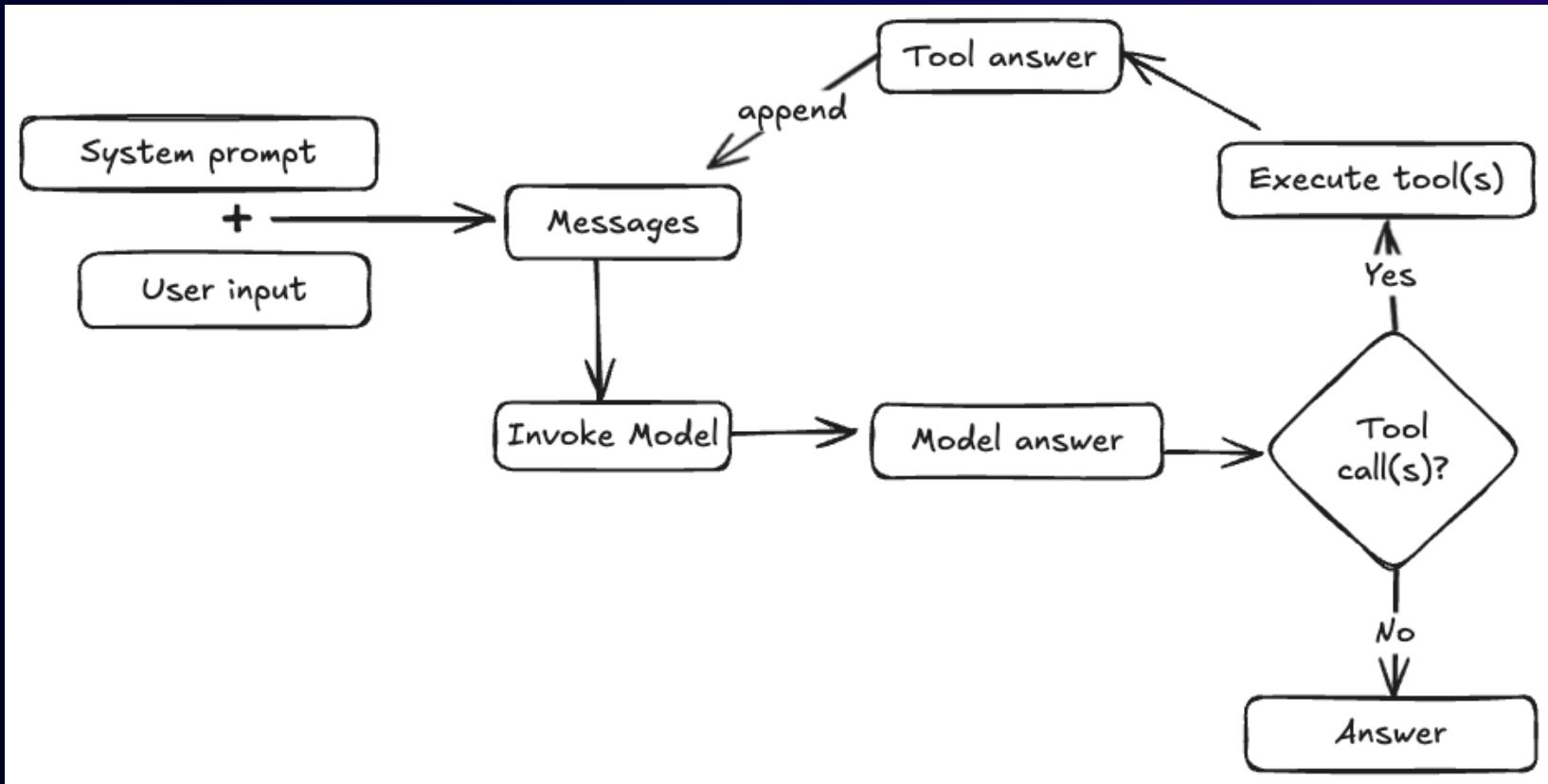


Strands see Agents as “models using tools in cycles”

At the core of Strands Agents is a simple agentic loop which enables models to autonomously plan next steps and execute tools.



Understand the control flow



Built for builders

Who value flexibility, speed, and simplicity

Model & deployment choice

- ✓ Model choice
- ✓ Custom model providers
- ✓ Deploy anywhere

Highly flexible

- ✓ Safeguard with guardrails
- ✓ Native observability
- ✓ Monitoring
- ✓ Evaluation
- ✓ Session Management

Broad selection of tools

- ✓ MCP integration
- ✓ Custom tools
- ✓ Coordinate multiple agents
- ✓ Multi-modal
- ✓ Fetch web data
- ✓ Read and write files
- ✓ Interpret code

Integrations

- ✓ Use AWS services
- ✓ LiteLLM
- ✓ MemO
- ✓ Ragas
- ✓ Tavily
- ✓ Langfuse
- ✓ MCP
- ✓ A2A

Simple hello world Strands Agent

```
from strands import Agent
from strands_tools import calculator

agent = Agent(tools=[calculator])
response = agent("What is 80 / 4?")
```

Model choice – Bedrock

```
from strands import Agent
from strands.models import BedrockModel

bedrock_model1 = BedrockModel(
    model_id= "us.amazon.nova-premier-v1:0",
    params={"max_tokens": 1600, "temperature": 0.7}
)
bedrock_model2 = BedrockModel(
    model="us.anthropic.claude-3-7-sonnet-20250219-v1:0"
)

agent = Agent(model=bedrock_model1)
response = agent("Tell me about Bedrock")
```

Model choice – LiteLLM

```
from strands import Agent
from strands.models.litellm import LiteLLMModel

litellm_model = LiteLLMModel(
    model_id= "azure/gpt-4o-mini",
    params={"max_tokens": 1600,
             "temperature": 0.7}
)
agent = Agent(model=litellm_model)
response = agent("Tell me about LiteLLM")
```

Choice of providers

- Amazon Bedrock
- Anthropic API
- LiteLLM
- Llama API
- MistralAI
- Ollama
- OpenAI
- Writer
- Cohere
- **Custom Providers**

LiteLLM

```
from strands import Agent
from strands.models.litellm import LiteLLMModel
from strands_tools import calculator

model = LiteLLMModel(
    # **model_config
    client_args={
        "api_key": "sk-a[REDACTED]qhk",
        "base_url": "https://api.siliconflow.cn/"
    },
    model_id="openai/Qwen/Qwen3-235B-A22B",
    params={
        "max_tokens": 1000,
        "temperature": 0.7,
    }
)

agent = Agent(model=model, tools=[calculator])
response = agent("What is 2+2")
print(response)
```

OpenAI

```
from strands.models.openai import OpenAIModel
from strands_tools import calculator
import os

model = OpenAIModel(
    # **model_config
    client_args=[
        "api_key": os.environ["OPENAI_API_KEY"],
        "base_url": "https://api.siliconflow.cn/"
    ],
    model_id="Pro/deepseek-ai/DeepSeek-V3",
    max_tokens = 1000,
    temperature = 0.7,
)

agent = Agent(model=model, tools=[calculator])
response = agent("what is 2+2")
```



Tools – repo of pre-built tools to get started quickly

Tool categories

- RAG & Memory
- File Operations
- Shell & System
- Code Interpretation
- Web & Network
- Multi-modal
- AWS Services
- Utilities
- Agents & Workflows

```
from strands import Agent
from strands_tools import http_request

agent = Agent(tools=[http_request])
print(agent("Where is the " +
    "International Space Station?"))
```



Prebuilt Tools Available

-  File Manipulation - Read, write, and edit files with smart modification
-  Shell integration - execute and interact with shell commands
-  Memory tool - stores user and agent memories during operation
-  HTTP client - makes API requests, supports authentication
-  Python execution - Run Python code snippets with state persistence
-  Math tools - Use symbolic math features to perform advanced calculations
-  AWS integration — Common AWS service access tools such as S3, EC2, DDB, etc.
-  Image processing - Generating and processing images for AI applications
-  Environmental management - safe handling of environmental variables
-  Logging - Create and manage structured journals
-  Advanced Reasoning - A tool for complex thinking and reasoning skills
-  Swarm intelligence - Coordinates multiple agents to solve problems in parallel
-  Graph/Workflow Tool — Tool for creating Agent graphs/workflows

Tools – create custom tools

```
from strands import Agent, tool

@tool

def weather_forecast(city: str, days: int = 3) -> str:
    """Get weather forecast for a city.

    Args:
        city: The name of the city
        days: Number of days for the forecast
    """
    return f"Weather forecast for {city} for the next {days} days..."

agent = Agent(tools=[weather_forecast])
print(agent("What's the weather in Seattle tmw?"))
```

Remember MCP (Model Context Protocol)

Is a standard protocol to connect Agents to Tools

Creating agents with MCP tools

Natively integrates with MCP to extend agent capabilities through external tools and services. Strands Agents provides several ways to connect to MCP servers:

Standard I/O (stdio)

For command-line tools and local processes that implement the MCP protocol

Streamable HTTP

For HTTP-based MCP servers that use Streamable-HTTP Events transport

Server-Sent Events (SSE)

For HTTP-based MCP servers that use Server-Sent Events transport (deprecated)

Custom Transport with MCPClient

You can also implement a custom transport mechanism using MCPClient class directly

How to use MCP with Strands Agents? - Stdio

```
from mcp import stdio_client, StdioServerParameters
from strands import Agent
from strands.tools.mcp import MCPClient

stdio_mcp_client = MCPClient(lambda: stdio_client(
    StdioServerParameters(command="uvx",
        args=["awslabs.aws-documentation-mcp-server@latest"])))
)

with stdio_mcp_client:
    tools = stdio_mcp_client.list_tools_sync()
    agent = Agent(tools=tools)
    agent("What does Bedrock InvokeInlineAgent do?")
```

How to use MCP with Strands Agents? - HTTP

```
from mcp.client.streamable_http import streamablehttp_client
from strands import Agent
from strands.tools.mcp.mcp_client import MCPClient

streamable_http_mcp_client = MCPClient(
    lambda: streamablehttp_client("http://localhost:8000/mcp")
)

with streamable_http_mcp_client:
    tools = streamable_http_mcp_client.list_tools_sync()
    agent = Agent(tools=tools)
    agent("What does Bedrock InvokeInlineAgent do?")
```

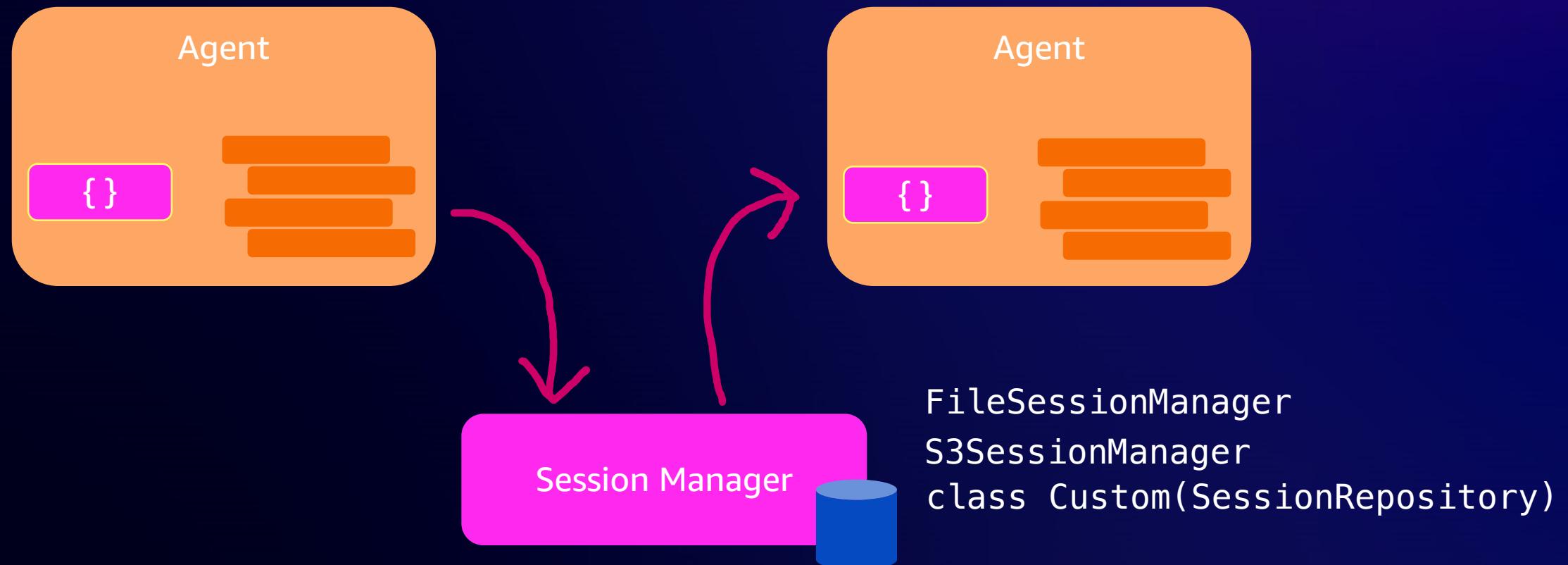
State & Sessions



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Session Management

- Use case: Stop and resume an Agent for a given session Id



State & Sessions

Strands Agents state is maintained in several forms:

- **Conversation History:** The sequence of messages between the user and the agent.
- **Agent State:** Stateful information outside of conversation context, maintained across multiple requests.
- **Request State:** Contextual information maintained within a single request.

Conversation History

```
from strands import Agent  
agent = Agent()  
agent("Hello!")  
print(agent.messages)
```

```
from strands import Agent  
agent = Agent(messages=[  
    {"role": "user", "content": [{"text": "Hello, my name is Strands!"}]},  
    {"role": "assistant", "content": [{"text": "Hi there! How can I help you today?"}]}  
])  
  
agent("What's my name?")
```

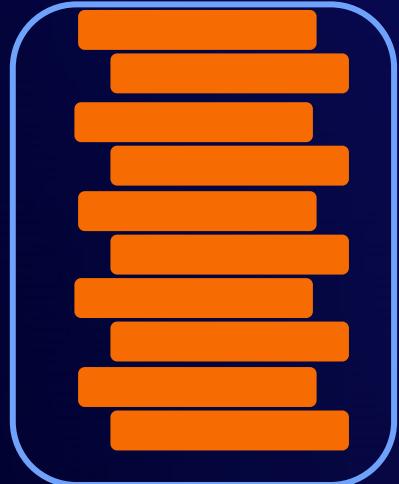
Conversation Management

Use case: Manage context size

No management



Sliding Window (K=40)



Summarization



Conversation Managers

```
from strands import Agent
from strands.agent.conversation_manager import
SlidingWindowConversationManager

# Create a conversation manager with custom window size
conversation_manager = SlidingWindowConversationManager(
    window_size=20, # Maximum number of messages to keep
    should_truncate_results=True, # Enable truncating the result
)
agent = Agent(
    conversation_manager=conversation_manager
)
```

Conversation Managers

```
from strands import Agent

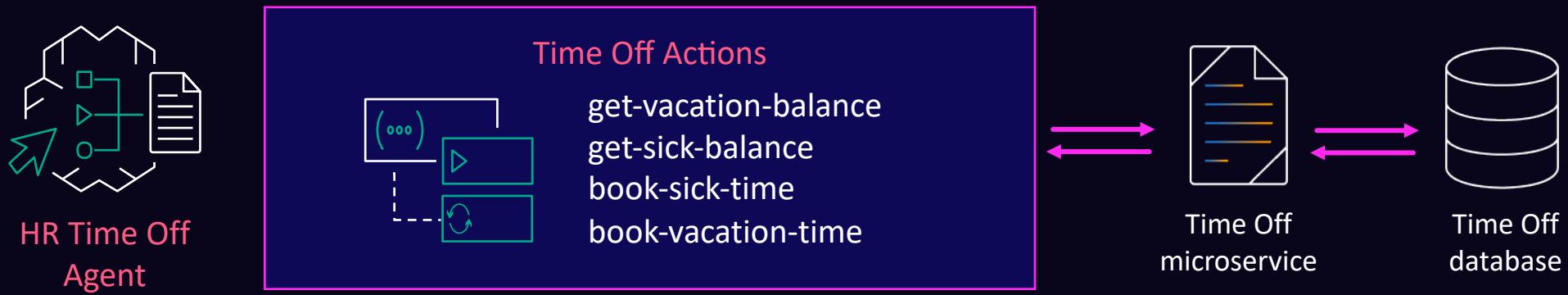
from strands.agent.conversation_manager import SummarizingConversationManager

# Create the summarizing conversation manager with default settings
conversation_manager = SummarizingConversationManager(
    summary_ratio=0.3, # Summarize 30% of messages when context reduction is needed
    preserve_recent_messages=10, # Always keep 10 most recent messages
)

agent = Agent(
    conversation_manager=conversation_manager
)
```

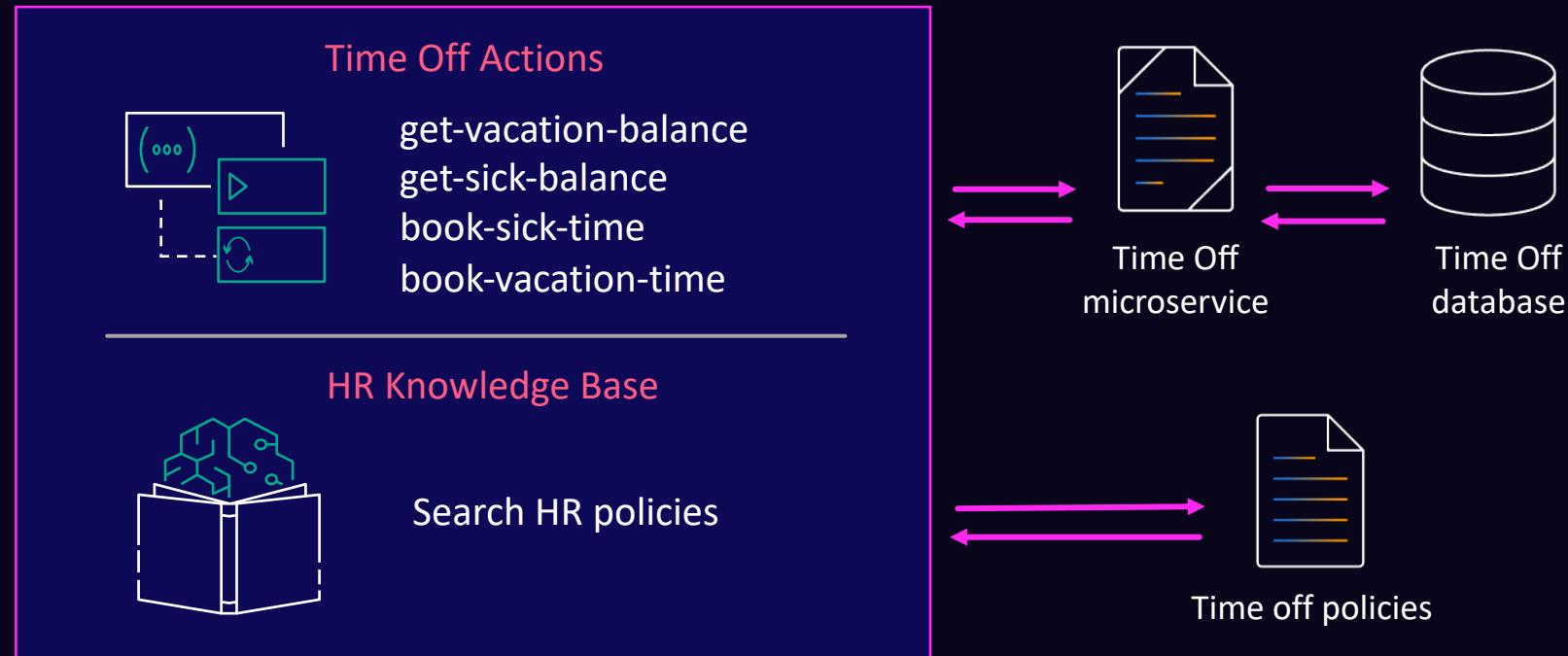
Multi-agent collaboration patterns

Agents start small and focused . . .



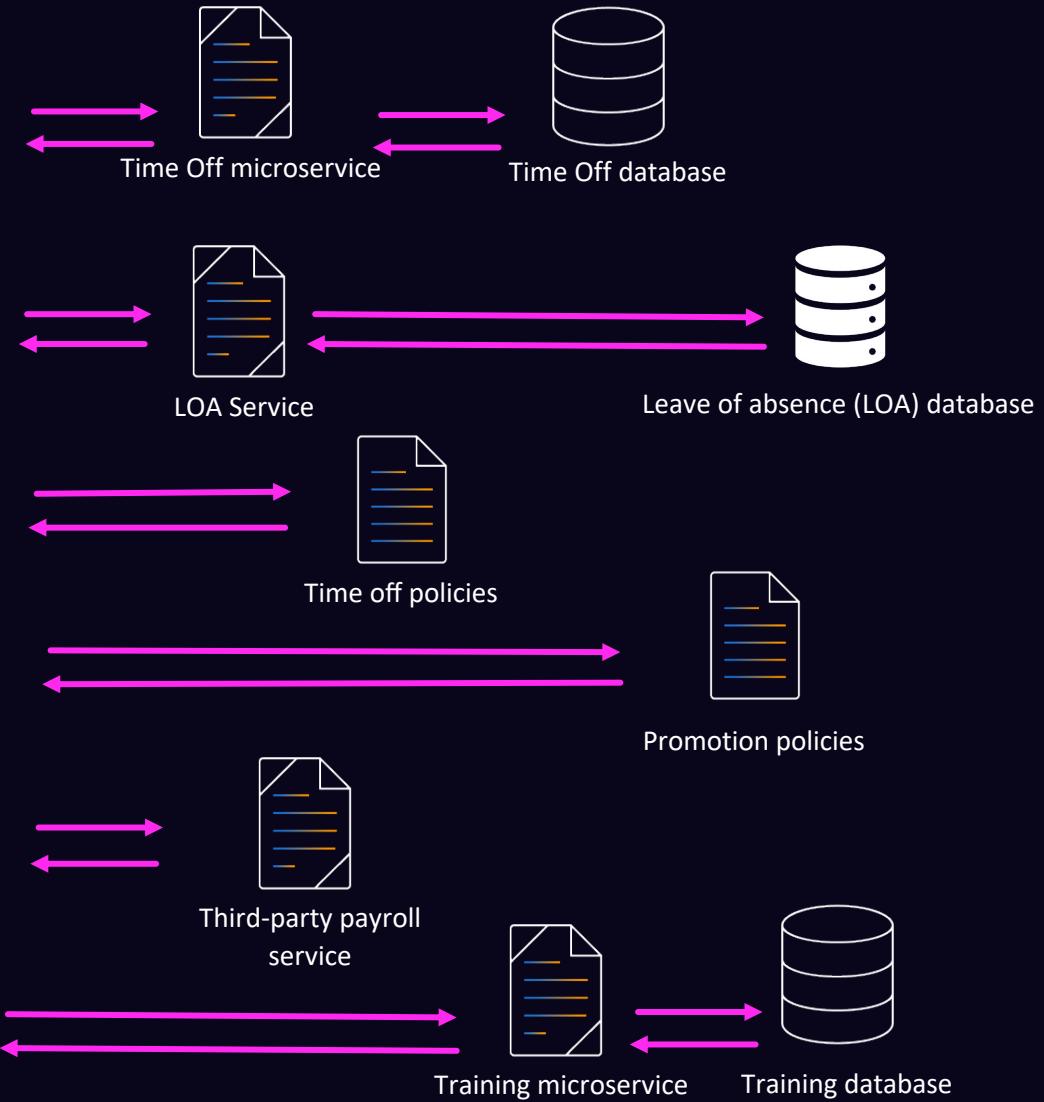
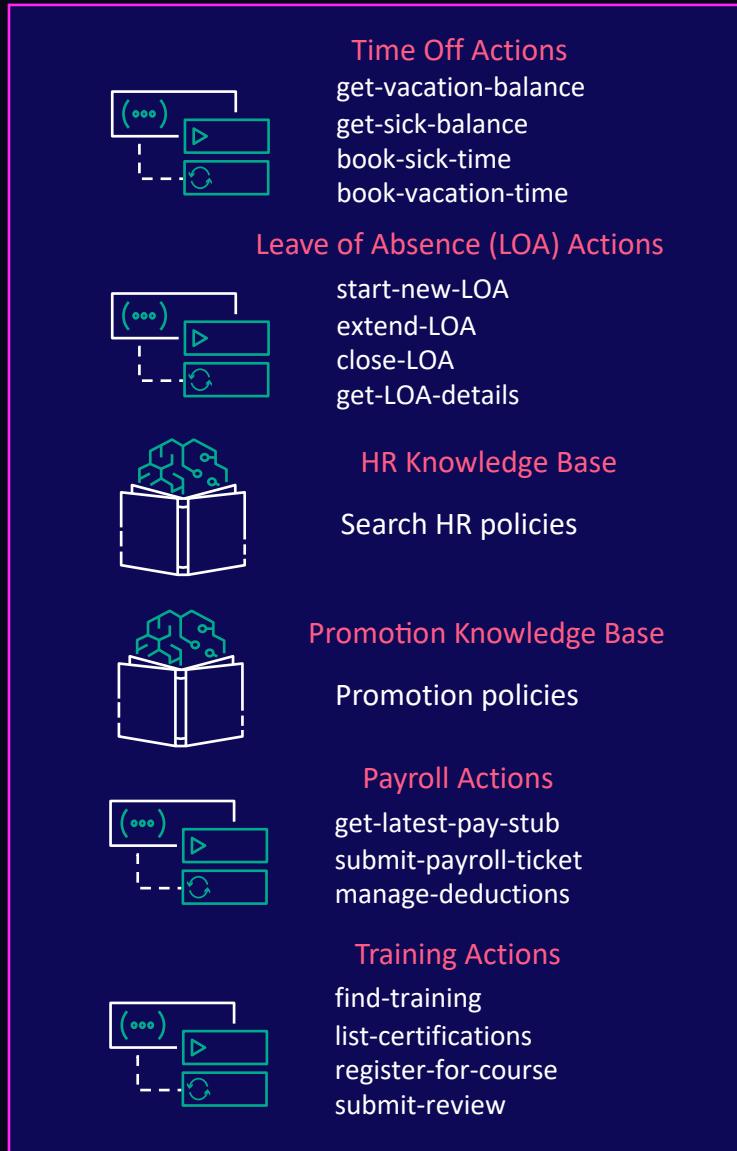
... and can be easily expanded

V2



... but if you take a SINGLE agent too far

V4



... it leads to challenges

Coding gets
complicated



- Complex prompts to limit hallucinations
- Fragile, hard to maintain

Agent gets
confused



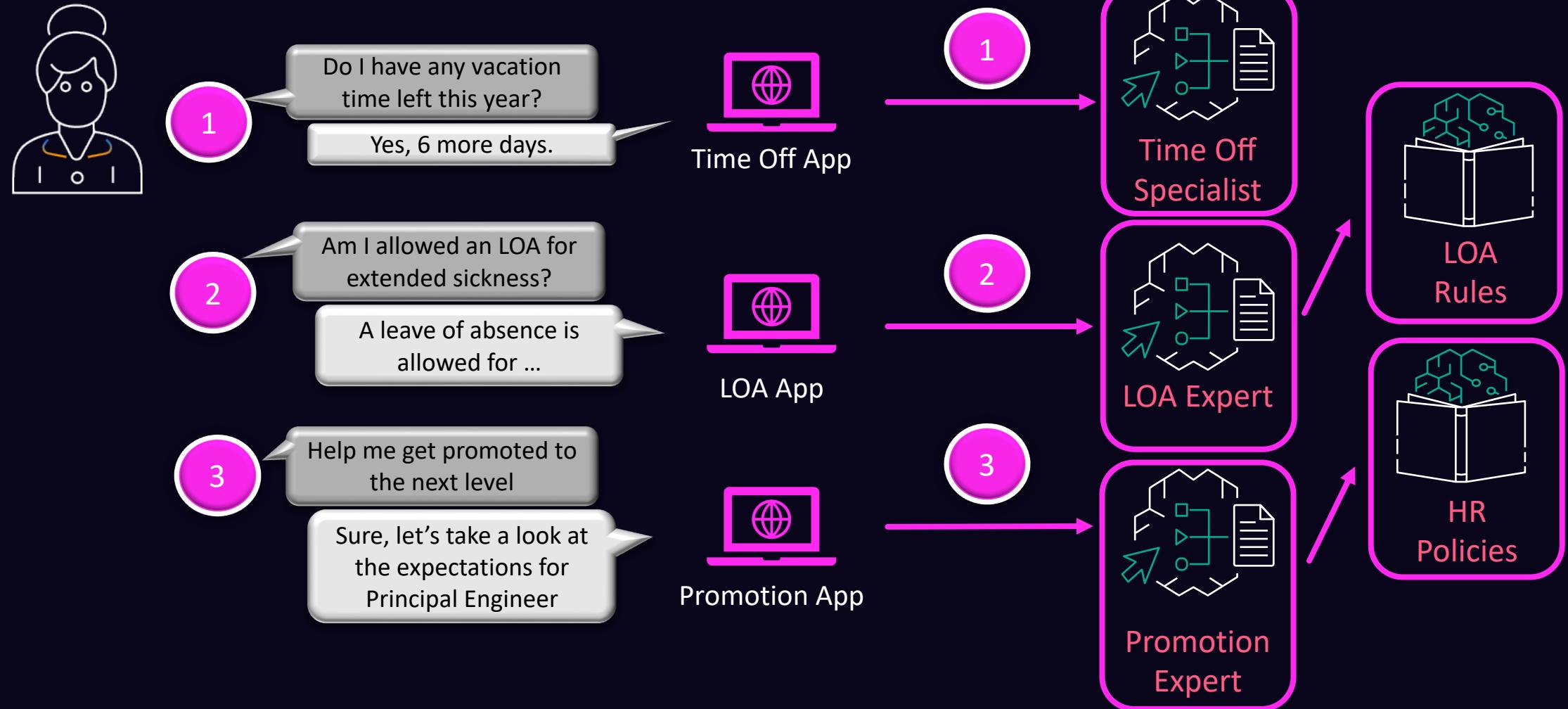
- Calling wrong tools
- Passing wrong arguments
- Inconsistent responses

Agent gets slower and
more expensive

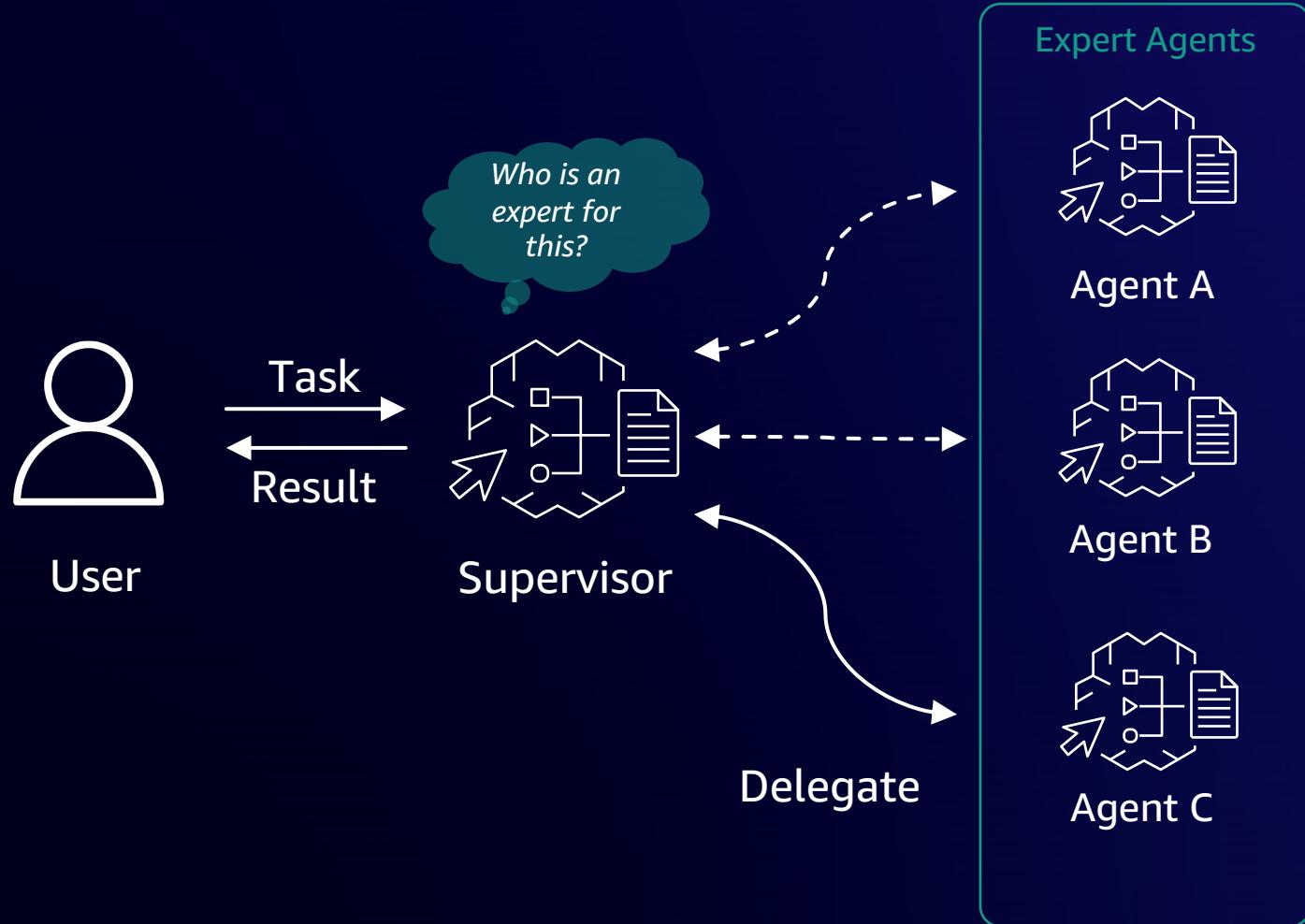


- Frontier models needed
- Prompt sizes grow
- Agents retry steps

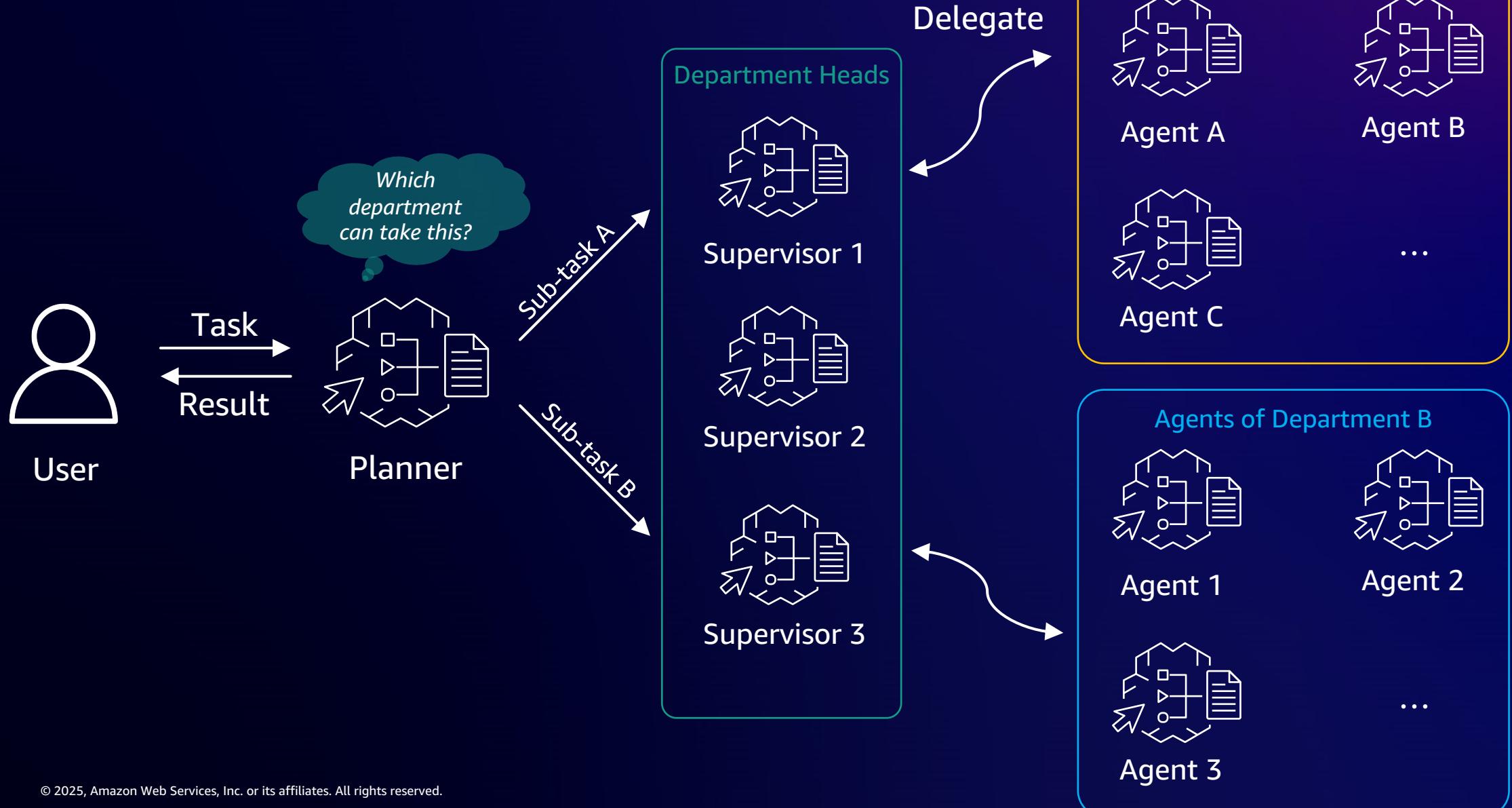
Using multiple agents helps . . .



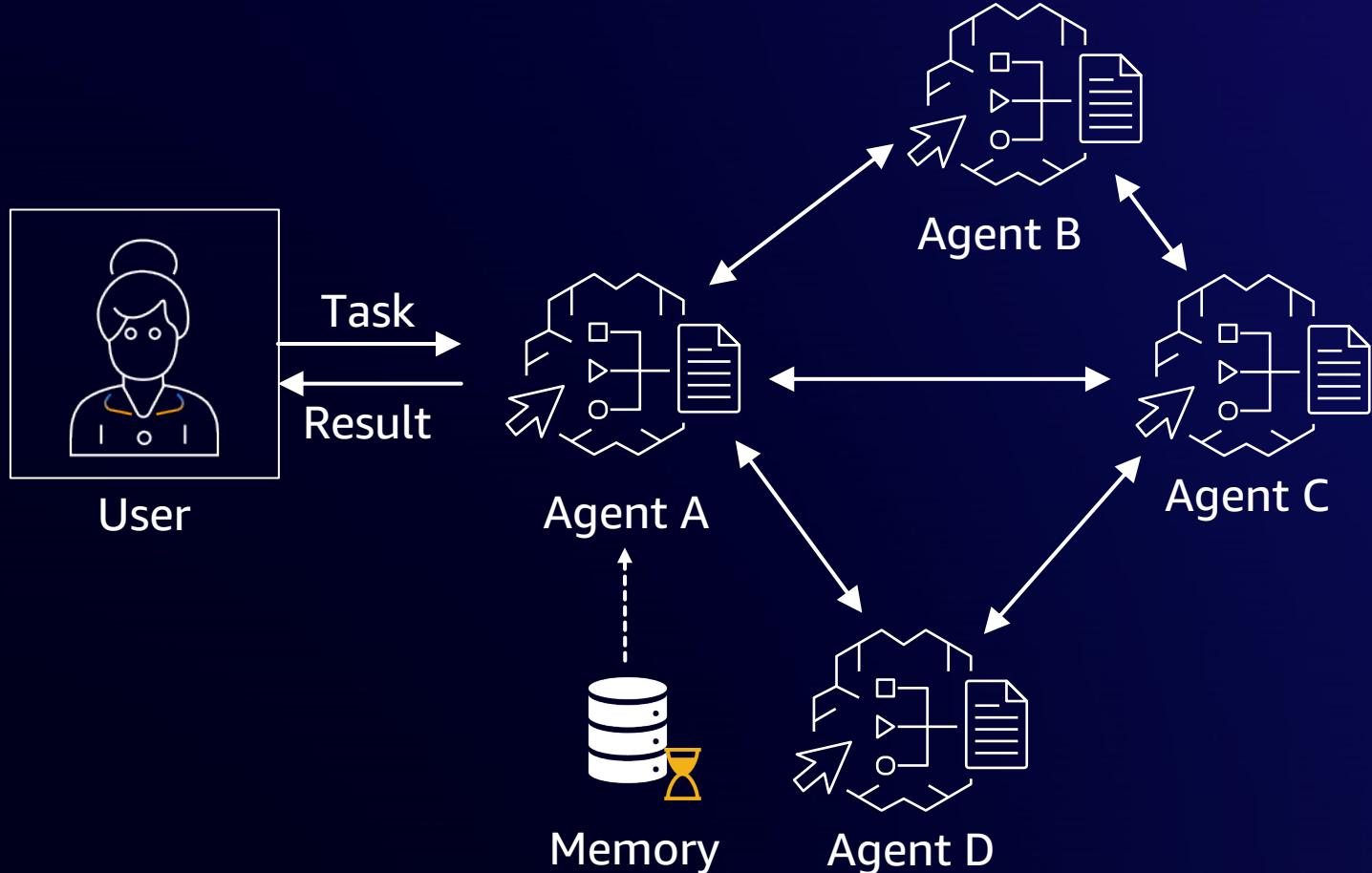
Agents as Tools



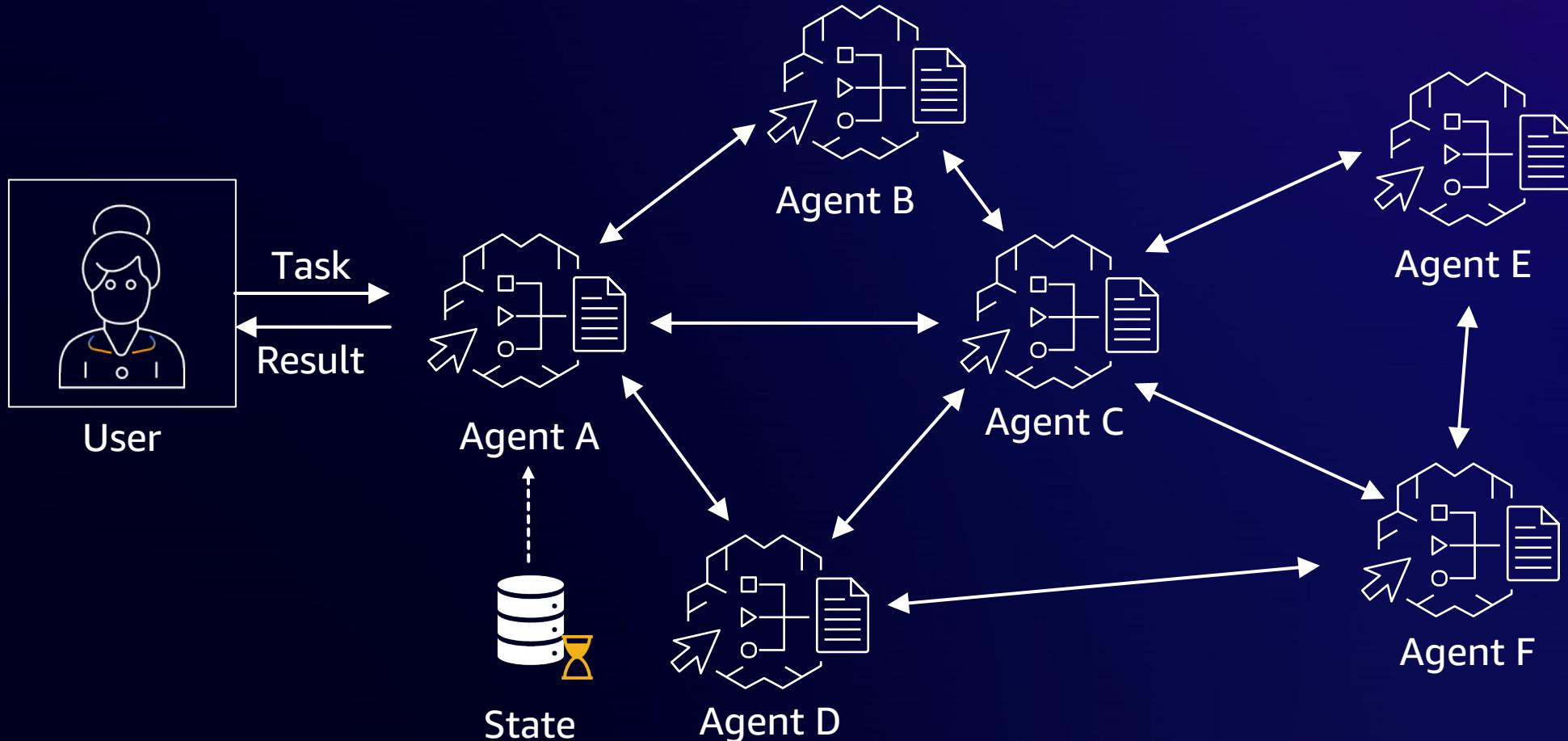
Hierarchical Agents as Tools



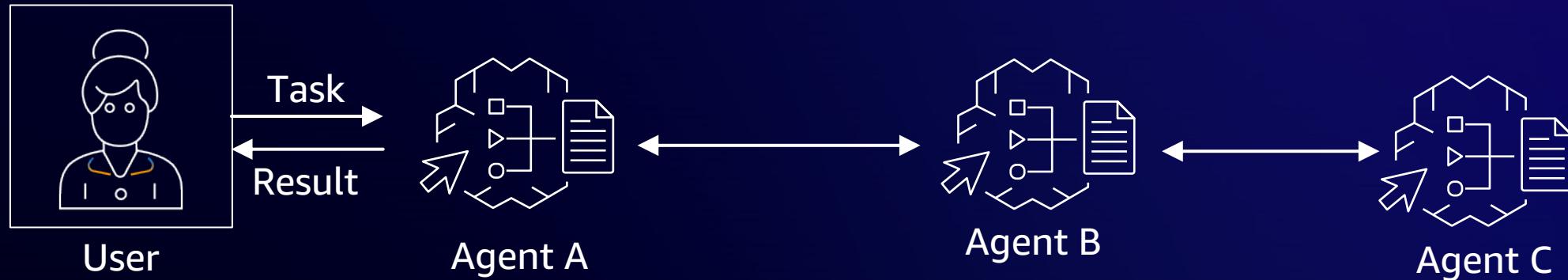
Swarm Agents



Graph Agents



Workflow Agents



Strands Agents multi-agent collaboration patterns

Agent as Tools

Hierarchical structure where:

1. **"orchestrator" agent** handles user interaction and determines which specialized agent to call
2. **Specialized agent** performs domain-specific tasks when called by the orchestrator



Swarm

Multi-agent systems where agents interact within the environment. Architecture consists of 4 key components:

1. Communication Patterns
2. Shared Memory Systems
3. Coordination Mechanisms
4. Task Distribution



Graph

Network of agents where each agent represents a node with specific capabilities, and the connections between agents is defined explicitly.

1. Nodes (agents)
2. Edges (connections)
3. Topology Pattern



Workflow

Structured coordination of tasks across agents. Provides explicit control over execution order, dependencies, and information flow.

1. Task Definition and Distribution
2. Dependency Management
3. Information Flow



Observability and Evaluation

Strands Agents Observability and Evaluation

Strands Agents provide metrics and OTEL traces out-of-the-box

Metrics

- **Token usage:** Input tokens, output tokens, and total tokens consumed
- **Performance metrics:** Latency and execution time measurements
- **Tool usage:** Call counts, success rates, and execution time for each tool
- **Event loop cycles:** Number of reasoning cycles and their durations

Traces

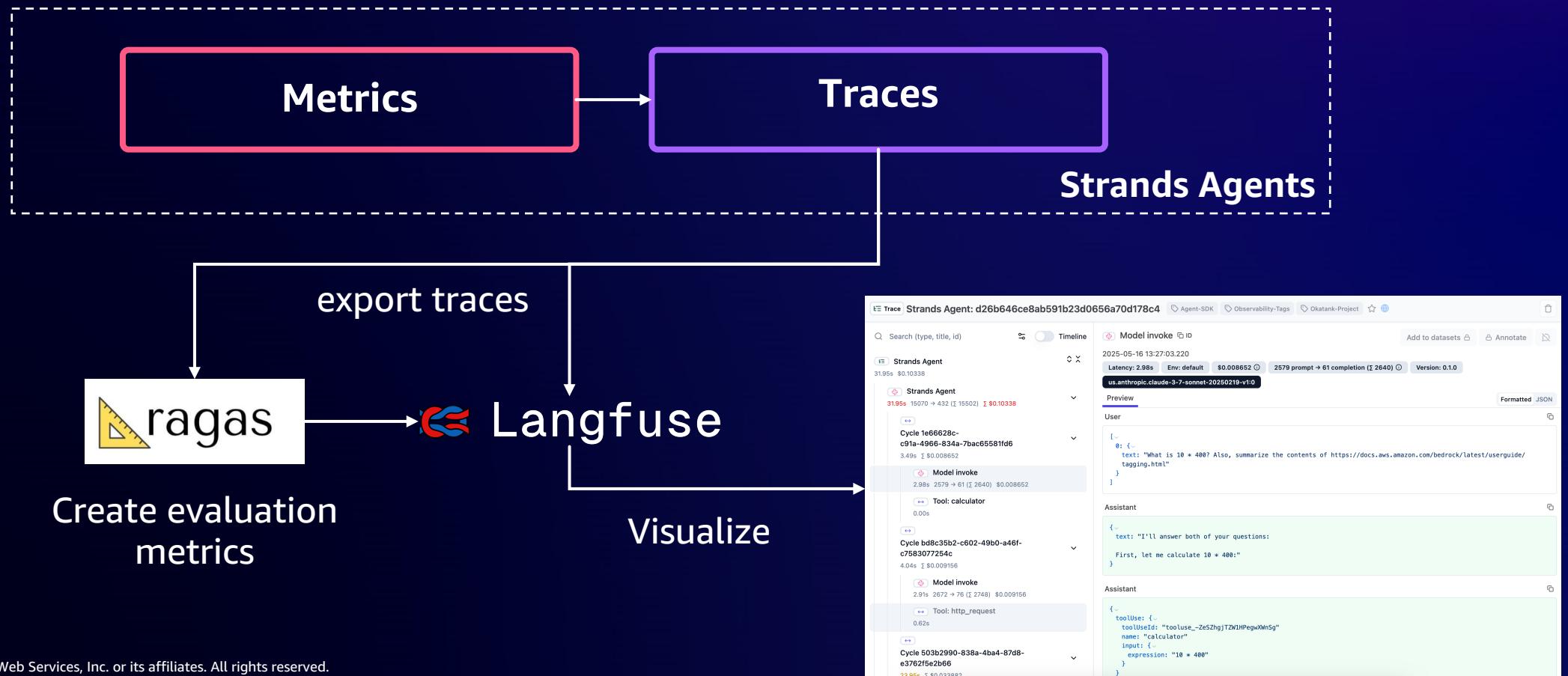
- **Track the entire agent lifecycle:** From initial prompt to final response
- **Monitor individual LLM calls:** Examine prompts, completions, and token usage
- **Analyze tool execution:** Understand which tools were called, with what parameters, and their results
- **Debug complex workflows:** Follow the exact path of execution through multiple cycles

Metrics are also exposed into the traces

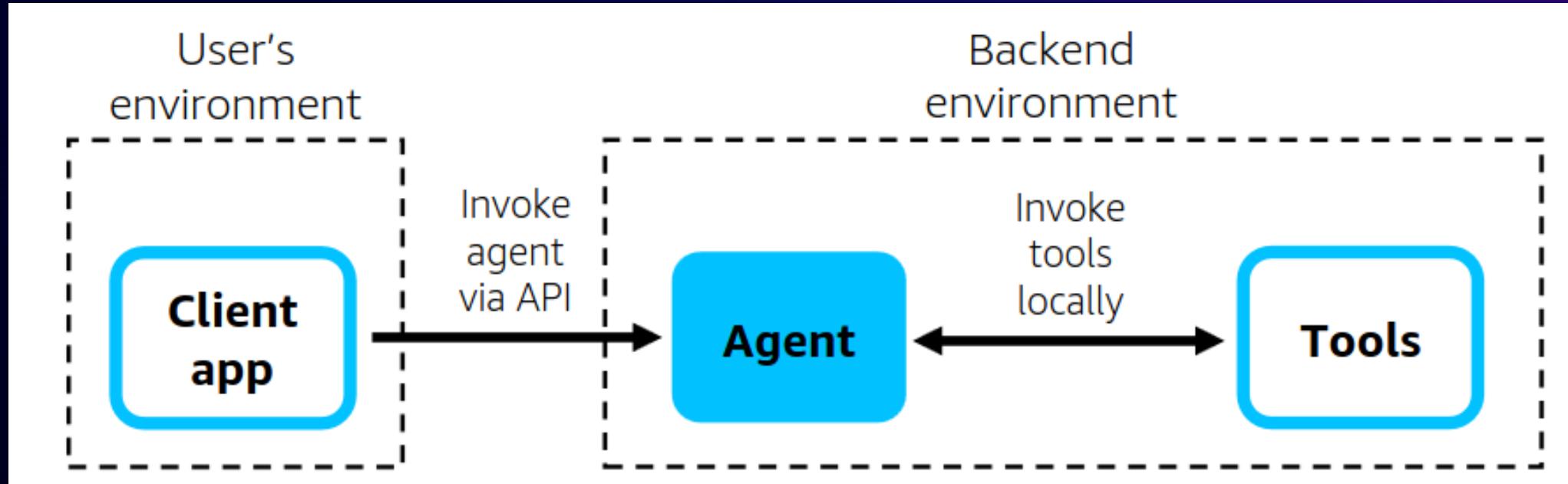


Strands Agents Observability and Evaluation

Strands integrates with different open-source tools to evaluate your agent and visualize metrics and traces on dashboards



Run agent behind an API



Strands Agents comes with examples and documented steps for deploying agents on:
AWS Lambda, Amazon EC2, AWS Fargate, Amazon EKS, Amazon Bedrock AgentCore Runtime

Difference between Bedrock Agents and Strands Agents

Amazon Bedrock Agents

- For developers who want a fully-managed experience for building and deploying agents.
- Fast path to production, using API calls to specify the model, instructions, tools, configuration.
- Avoids overhead of managing hosting and scaling, or stitching together core agent capabilities.
- Comes with multi-agent collaboration, observability, code interpretation, memory management, guardrails and knowledge bases.

Strands Agents

- For developers who want deeper customization of agent behavior, including orchestration, and business logic.
- Works with any model provider, including proprietary LLM gateways.
- Agents can be deployed anywhere you would host a Python application.
- Open source community project maintained by AWS.

Get started with Strands Agents

Start your agentic AI journey today



Docs

Review the user guide, quick start, API documentation for best practices



SDK

Get started building agents and contribute to the open-source project on GitHub



Tools

Access pre-built tools to get started experimenting with agents



Samples

Discover sample agents in our samples repository

Notebook Exercises

GENAI BOOTCAMP

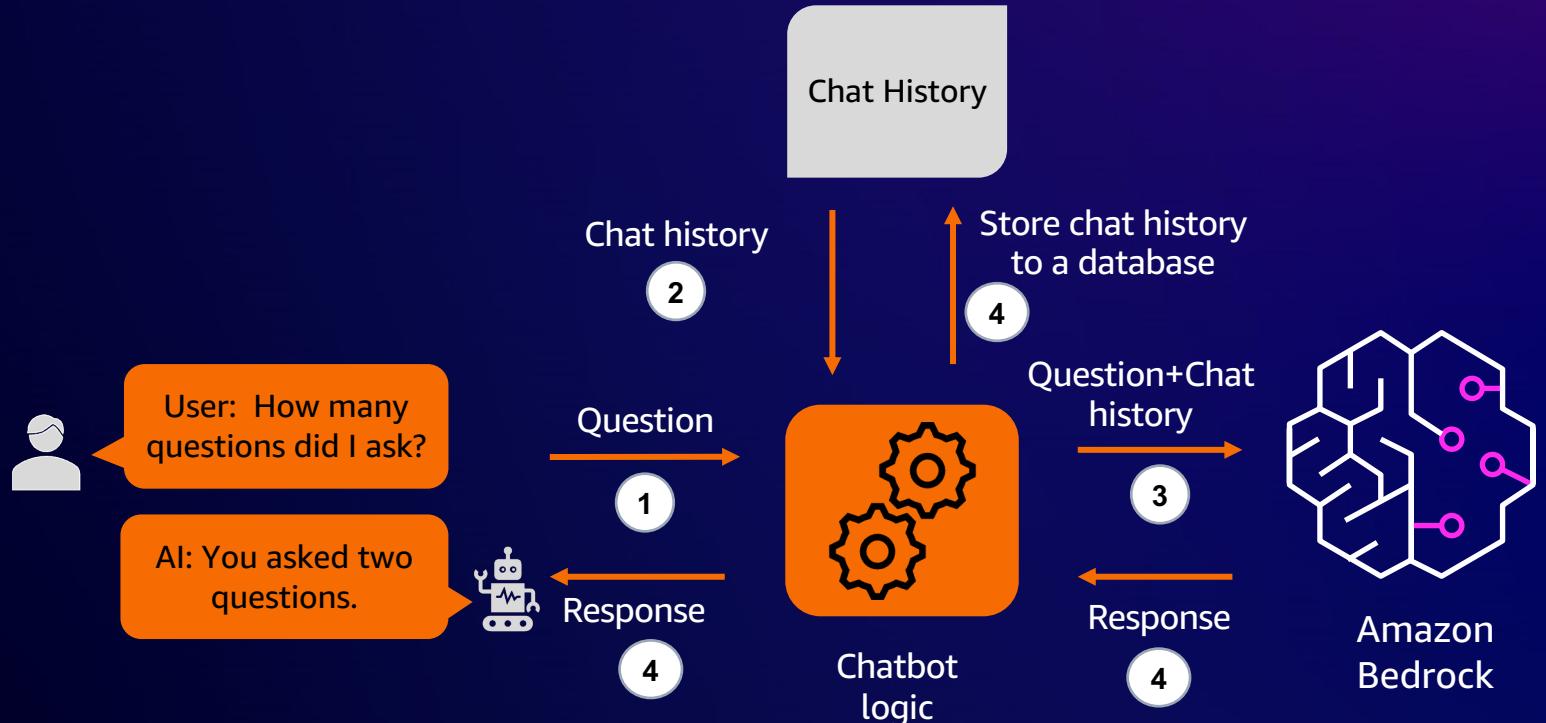


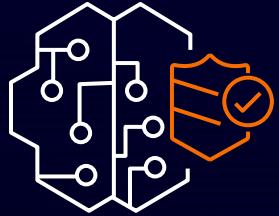
Secure Chatbots with Bedrock Guardrails



Chatbot - Overview

- 1 User asks a question
- 2 Chat history is retrieved
- 3 Answer is generated using **question** and **chat history**
- 4 Answer is returned to the user and saved in chat history





Amazon Bedrock Guardrails

Apply safeguards customized to your generative AI application requirements and responsible AI policies

Why do you need additional controls?

THE COST OF NO GUARDRAILS: HIDDEN RISKS OF UNFILTERED AI INTERACTIONS

Customizations based
on use cases and
organizational policy



Why
guardrails?

The cost
without
guardrails

Safety and privacy controls
for responsible AI



Consistent safeguards
across FMs and applications



Good FM training
isn't enough
Dataset poisoning is persistent



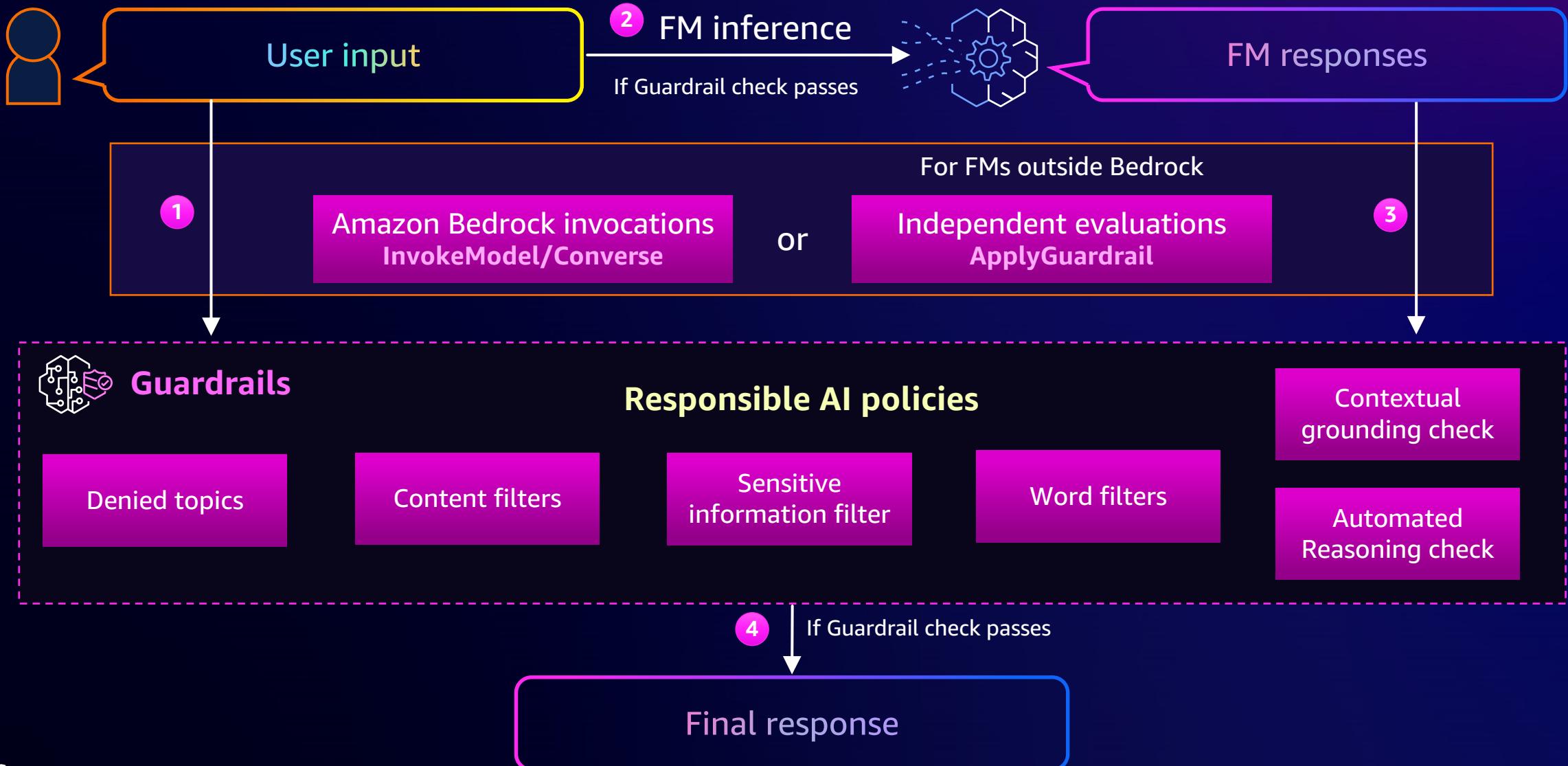
RAG doesn't stop
hallucinations
Risk of exposing
sensitive information



Good FM prompting
isn't enough
Prompt attacks bypass
and FM controls



How Amazon Bedrock Guardrails works



Deep dive into the components of a guardrail

Content filters

FILTER HARMFUL CONTENT ACROSS 6 CATEGORIES

- Hate
- Insults
- Sexual
- Violence
- Misconduct
- Prompt attack

The screenshot shows the 'Edit content filters' page in the Amazon Bedrock Guardrails interface. It displays settings for filtering harmful content across six categories: Hate, Insults, Sexual, Violence, Misconduct, and Prompt attack. Each category has options to enable all harmful categories (Text or Image) and to use the same filters for responses. For each category, there is a dropdown menu for 'Guardrail action' (set to 'Block') and a horizontal slider for 'Set threshold' with four levels: None, Low, Medium, and High.

Filters for prompts

Use the same harmful categories filters for responses

Enable ⓘ	Category	Guardrail action	Set threshold
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Hate	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Insults	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Sexual	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Violence	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Misconduct	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>

Filters for responses

Use the same harmful categories filters for responses

Enable ⓘ	Category	Guardrail action	Set threshold
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Hate	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Insults	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Sexual	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Violence	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Misconduct	Block	None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High <input type="radio"/>

[Reset thresholds](#)

Prompt attacks detection

Similar to harmful categories, prompt attacks are detected based on classification confidence

The screenshot shows the 'Create guardrail' wizard in the Amazon Bedrock console. The current step is 'Step 2 - optional: Configure content filters'. On the left, a sidebar lists steps from 1 to 8. Step 2 is highlighted with a blue circle. The main content area is titled 'Configure content filters - optional' and contains sections for 'Harmful categories' and 'Prompt attacks'. Under 'Harmful categories', there is a note about detecting and blocking harmful inputs and model responses. Under 'Prompt attacks', there is a note about detecting user inputs attempting to override system instructions. Both sections have a 'Configure' button. Below these is a table for configuring filters. The first row has 'Enable' checked, 'Category' set to 'Prompt Attack', 'Guardrail action' set to 'Block', and a 'Set threshold' slider set to 'Medium'. The second row has 'Enable all' checked, 'Category' set to 'Prompt Attack', 'Guardrail action' set to 'Block', and a 'Set threshold' slider set to 'Medium'. Under 'Content filters tier', there are two options: 'Classic' (selected) and 'Standard'. The 'Classic' tier supports English, French, and Spanish languages and does not require cross-Region inference. The 'Standard' tier supports over 50 languages and requires cross-Region inference. At the bottom right are buttons for 'Cancel', 'Skip to Review and create', 'Previous', and 'Next'.

Multi-Modal Toxicity Detection

- ✓ Content filters now supports toxic image detection, in addition to text
- ✓ Enhance user experiences and manage risk by detecting and filtering undesirable content across different media
- ✓ Content moderation for social media, e-commerce and advertising, financial fraud detection, educational and healthcare applications, among other applications
- ✓ Industry-leading, up to 88% harmful multimodal content detection rate

Harmful categories
Enable detection and blocking of harmful user inputs and model responses. Use a higher filter strength to help improve the filtering of harmful content in each category.

Configure harmful categories filters

Filters for prompts
 Use the same harmful categories filters for responses

Enable	Category	Guardrail action	Strength
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Hate	Detect (no action)	None → High
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Insults	Block	None → High
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Sexual	Detect (no action)	None
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Violence	Block	None
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Misconduct	Block	None

Prompt attacks
Enable to detect and block user inputs attempting to override system instructions. To avoid misclassifying filters are selectively applied to user inputs, use input tagging.

Configure prompt attack filter

Enable	Category	Guardrail action	Strength
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Prompt Attack	Block	None

Guardrail trace

Content filters

Category	Test result	Details
Sexual	No action taken	Text Strength: High Confidence: None
Prompt attack	No action taken	Text Strength: High Confidence: None
Hate	No action taken	Text Strength: High Confidence: None
Insults	No action taken	Text Strength: High Confidence: None
Misconduct	No action taken	Text Strength: High Confidence: None

Violence Blocked Text
Strength: High
Confidence: Medium

Guardrail action
干预 (1 例) 干预 View trace

Guardrails Response
Sorry, the model cannot answer this question. -

Run

Guardrail trace

Content filters

Category	Test result	Details
Sexual	No action taken	Text Strength: High Confidence: None
Prompt attack	No action taken	Text Strength: High Confidence: None
Hate	No action taken	Text Strength: High Confidence: None
Insults	No action taken	Text Strength: High Confidence: None
Misconduct	No action taken	Text Strength: High Confidence: None

Violence Blocked Text
Strength: High
Confidence: Medium

Denied topics

Word filters

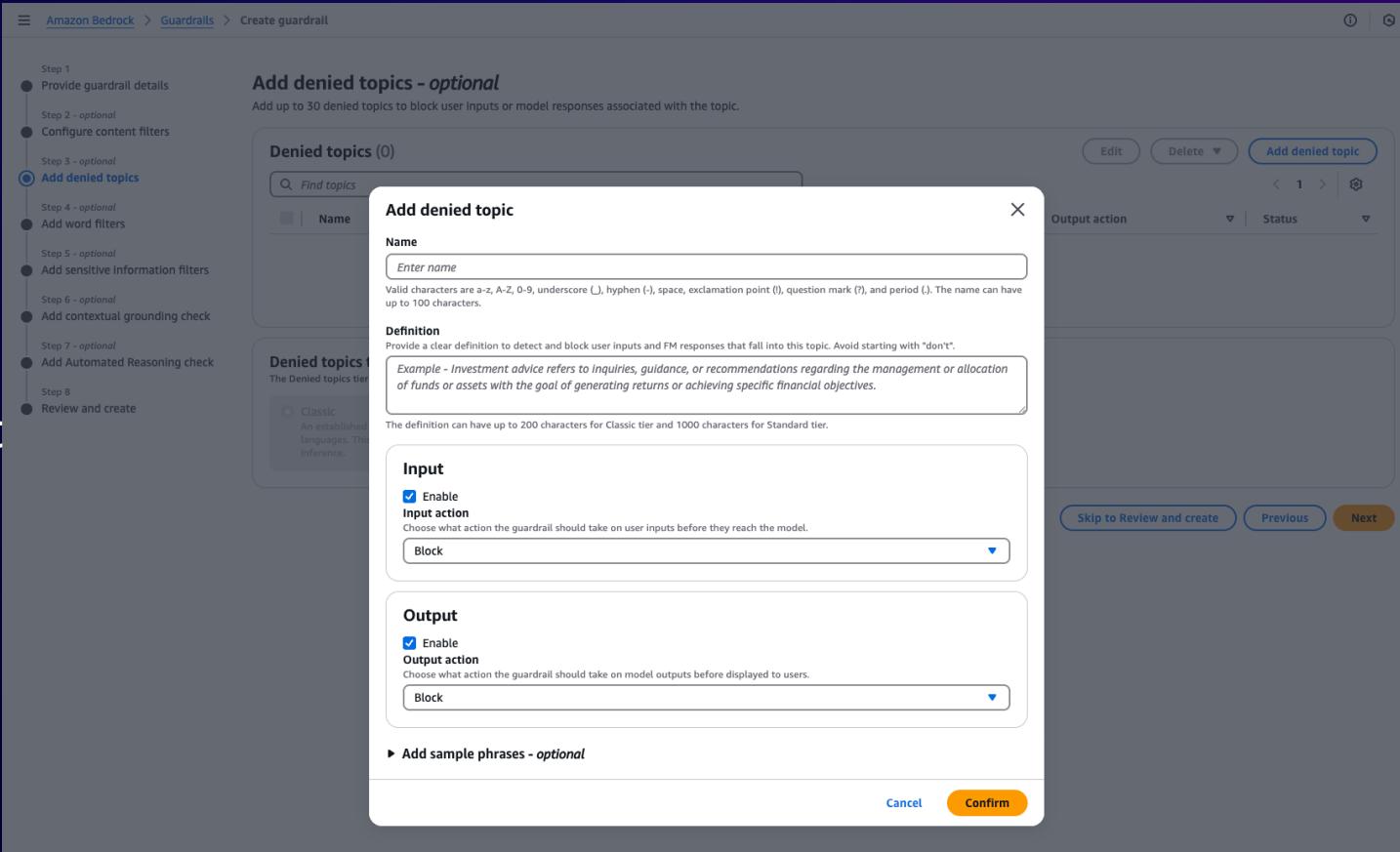
Sensitive information filters

Denied topics

Topics are defined in simple language and compared against user queries/requests to determine similarity

Examples:

- **Substance use history**—past or present use of alcohol, tobacco, drugs, or medications outside the scope defined in the application process
- **Financial information**—income, assets, debts, credit score, or financial details not directly relevant to the insurance product applied for



Sensitive information filter

PROTECT YOUR SENSITIVE DATA

- ✓ Detect, redact or filter PII in FM responses to protect user privacy
- ✓ Detect, redact or filter PII in user inputs
- ✓ Select from a variety of pre-configured PII types, based on application requirements
- ✓ Or define your own sensitive information, by using regular expressions (regex)
- ✓ Apply actions for prompts and/or responses independently

PII types (3)

Choose PII type	Input action	Output action
NAME	Detect (no action)	Mask
EMAIL	Mask	Mask
ADDRESS	Block	Block

Exit Save Save and exit

Add new PII

PII type

Choose PII type

General

- Phone
- Age
- Username
- Password
- Driver ID
- License plate
- Vehicle Identification number (VIN)

Finance

- CVV
- Credit/Debit card expiry
- Credit/Debit card number
- Personal identification number (PIN)
- International Bank Account Number
- SWIFT code

IT

- IPv4 address
- Media access control (MAC) address
- Web address
- AWS access key
- AWS secret key

Confirm

Add regex pattern

Name

Label to identify the pattern. Shown as an identifier if PII is masked, e.g. [BOOKING_ID].

Example - Booking ID

Regex name can have up to 100 characters.

Regex pattern

Example - ^ID\d{3}[A-Z]{3}\$

Input

Enable

Output action

Choose what action the guardrail should take on user inputs before they reach the model

Block

Output

Enable

Output action

Choose what action the guardrail should take on model outputs before displayed to users

Block

▶ Add description - optional

Cancel Confirm

Word filters

BLOCK SPECIFIC WORDS OR PHRASES

- Filter profane words
- Define a set of custom words to block in user input and FM responses

The screenshot shows the AWS Guardrail Word Filters configuration interface. It includes sections for 'Profanity filter' (with 'Filter profanity' checked), 'Input' (with 'Enable' checked and 'Input action' set to 'Block'), 'Output' (with 'Enable' checked and 'Output action' set to 'Block'), 'Add custom words and phrases' (with a note about specifying up to 10,000 words or phrases), and a table for 'View and edit words and phrases'. The table has columns for 'Word or phrase', 'Enable input', 'Input action', 'Enable output', 'Output action', and 'Status'. A search bar at the top of the table allows for filtering by word or phrase.

Profanity filter

Filter profanity
Enable this feature to block profane words in user inputs and model responses. The list of words is based on the global definition of profanity and is subject to change.

Input

Enable
Input action
Choose what action the guardrail should take on user inputs before they reach the model
Block

Output

Enable
Output action
Choose what action the guardrail should take on model outputs before displayed to users
Block

Add custom words and phrases

Specify up to 10,000 words or phrases (max 100 characters per) to be blocked by the guardrail. A blocked message will show if user input or model responses contain these words or phrases.

Add words and phrases manually
Manually add words and phrases to the following table.

Upload from a local file
Populate the following table with words and phrases from a .txt or .csv file from your computer.

Upload from S3 object
Populate the following table with words and phrases from an S3 object.

View and edit words and phrases (0)

Find words and phrases Show all

Word or phrase	Enable input	Input action	Enable output	Output action	Status
No words or phrases added	Upload from file or add manually in the console				
Add a word or phrase					

Contextual Grounding checks & relevance

REDUCE HALLUCINATION BY FILTERING UNGROUNDED AND IRRELEVANT RESPONSES

- Filter hallucinations in RAG and summarization applications
- Check response accuracy based on your enterprise data
- Check if the responses are relevant to user's query or instruction

Grounding threshold : Based on the information in your RAG response and the FM answer, Grounding evaluates the correctness of the answer

Relevance threshold: Evaluate the FM response and the user query. If the confidence that the response answers the query is below the threshold, a blocked message is returned to the user

Add contextual grounding check - optional Info

Use this policy to validate if model responses are grounded in the reference source and relevant to user's query to filter model hallucination.

Grounding
Validate if the model responses are grounded and factually correct based on the information provided in the reference source, and block responses that are below the defined threshold of grounding.

Enable grounding check

Grounding score threshold
Grounding score represents the confidence that the model response is factually correct and grounded in the source. If the model response has a lower score than the defined threshold, the response will be blocked and the configured blocked message will be returned to the user. A higher threshold level blocks more responses. Info



Contextual grounding action
Choose what action the guardrail should take on contextual grounding check.

Relevance
Validate if the model responses are relevant to the user's query and block responses that are below the defined threshold of relevance.

Enable relevance check

Relevance score threshold
Relevance score represents the confidence that the model response is relevant to the user's query. If the model response has a lower score than the defined threshold, the response will be blocked and the configured blocked message will be returned to the user. A higher threshold level blocks more responses. Info



Relevance action
Choose what action the guardrail should take on relevance check.

Configure Guardrails' actions

For each policy, select between:

- **Block** – For preventing content from being delivered when conditions are met

or

- **Detect** – For monitoring only, without intervention in the actual content flow

Available for all configurations, including content filters, denied topics, sensitive information filters, word filters and contextual grounding.

The screenshot shows the AWS Guardrails configuration interface. It includes sections for 'Harmful categories' and 'Prompt attacks'. In the 'Harmful categories' section, there are five categories: Hate, Insults, Sexual, Violence, and Misconduct. Each category has options to enable all filters for Text and Image. Under each category, there is a dropdown for 'Guardrail action' and a slider for 'Strength' (None, Low, Medium, High). A red box highlights the 'Misconduct' category's 'Guardrail action' dropdown, which shows 'Block' (Prevents content from being processed when this category's policy conditions are met) and 'Detect (no action)' (Set up as diagnostic tool without taking action on input or output). The 'Detect (no action)' option is checked. In the 'Prompt attacks' section, there is a dropdown for 'Guardrail action' (Detect (no action)) and a slider for 'Strength' (None, Low, Medium, High). The 'Configure prompt attacks filter' button is also visible.

Harmful categories
Enable detection and blocking of harmful user inputs and model responses. Use a higher filter strength to help improve the filtering of harmful content in each category.
 Configure harmful categories filters

Filters for prompts
 Use the same harmful categories filters for responses

Enable <small>i</small>	Category	Guardrail action	Strength
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Hate	Detect (no action)	<input type="radio"/> None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Insults	Block	<input type="radio"/> None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Sexual	Detect (no action)	<input type="radio"/> None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Violence	Block Prevents content from being processed when this category's policy conditions are met.	<input type="radio"/> None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input checked="" type="checkbox"/> Image	Misconduct	Detect (no action) Set up as diagnostic tool without taking action on input or output.	<input checked="" type="radio"/> Medium <input type="radio"/> High
		Block	<input type="radio"/> None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High

Prompt attacks
Enable to detect and block user inputs attempting to override system instructions. To avoid misclassifying system prompts as a prompt attack and ensure that the filters are selectively applied to user inputs, use input tagging.
 Configure prompt attacks filter

Enable <small>i</small>	Category	Guardrail action	Strength
<input checked="" type="checkbox"/> Enable all <input checked="" type="checkbox"/> Text <input type="checkbox"/> Image	Prompt Attack	Detect (no action)	<input type="radio"/> None <input type="radio"/> Low <input type="radio"/> Medium <input type="radio"/> High

Configure Guardrails' directions

Decide independent configurations and actions for each direction of the flow:

- **Input** – For evaluations on user's inputs
- **Output** – For evaluations on models' outputs

Available for content filters, denied topics, sensitive information filters, and word filters.

Edit denied topic

Name
Investment advice
Valid characters are a-z, A-Z, 0-9, underscore (_), hyphen (-), space, exclamation point (!), question mark (?), and period (.). The name can have up to 100 characters.

Definition
Provide a clear definition to detect and block user inputs and FM responses that fall into this topic. Avoid starting with "don't".
Queries or statements that seek or provide advice about investments in financial assets
The definition can have up to 200 characters.

Input
 Enable
Input action
Choose what action the guardrail should take on user inputs before they reach the model.
Detect (no action)

Output
 Enable
Output action
Choose what action the guardrail should take on model outputs before displayed to users.
Block

▶ Add sample phrases - optional

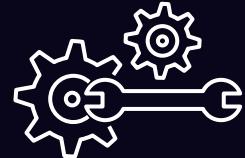
Cancel **Confirm**

Integration with existing systems & applications

SEAMLESS ADOPTION WITHIN YOUR CURRENT INFRASTRUCTURE

Seamless API integration

Easy incorporation into current workflows and applications



separated api to allow use independent of first, third-, or external-party FM use



Guardrails scale to enterprise applications with parallel execution



Tailored guardrails for different use cases and departments



JSON responses allow easy integration with existing observability tools

Strands Agents with Bedrock Guardrails

The screenshot shows the Strands Agents configuration interface. It includes sections for 'Prompt attacks' (Enabled), 'Denied topics' (Fiduciary Advice), and 'Word filters' (Custom words and phrases). The 'Word filters' section lists 10 entries where specific words or phrases are blocked.

No.	Word or phrase	Input action	Output action
1	fiduciary advice	Block	Block
2	investment recommendations	Block	Block
3	stock picks	Block	Block
4	financial planning guidance	Block	Block
5	portfolio allocation advice	Block	Block
6	retirement fund suggestions	Block	Block
7	wealth management tips	Block	Block
8	trust fund setup	Block	Block
9	investment strategy	Block	Block
10	financial advisor recommendations	Block	Block

Code sample

<https://github.com/strands-agents/samples/tree/main/01-getting-started/06-guardrail-integration>

```
# Create a Bedrock model with guardrail configuration
bedrock_model = BedrockModel(
    model_id="us.amazon.nova-premier-v1:0",
    guardrail_id=guardrail_id,
    guardrail_version=guardrail_version,
    boto_session=boto3.Session(region_name="us-west-2"),
    guardrail_trace="enabled"
)

# Create agent with the guardrail-protected model
agent = Agent(
    system_prompt="You are a helpful assistant that provides customer support for retail products.",
    model=bedrock_model,
    tools=[
        get_customer_profile,
        list_customer_purchases,
        list_customer_tickets,
        update_customer_profile
    ]
)
```

User: This is my account number 923987 can you tell me what stocks to buy based on this info

Agent: I apologize, but I am not able to provide fiduciary advice. Additionally, it seems that you may have included some sensitive personal or financial information in your request. For your privacy and security, please modify your input and try again without including any personal, financial, or restricted details. **⚠ GUARDRAIL INTERVENED!**



Notebook Exercises

GENAI BOOTCAMP

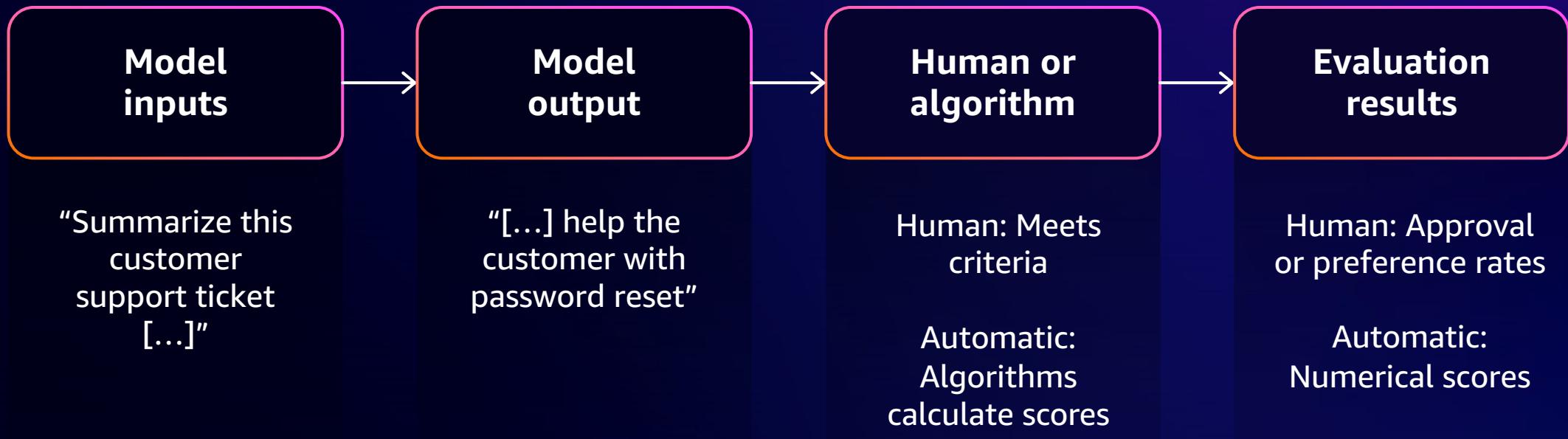


© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

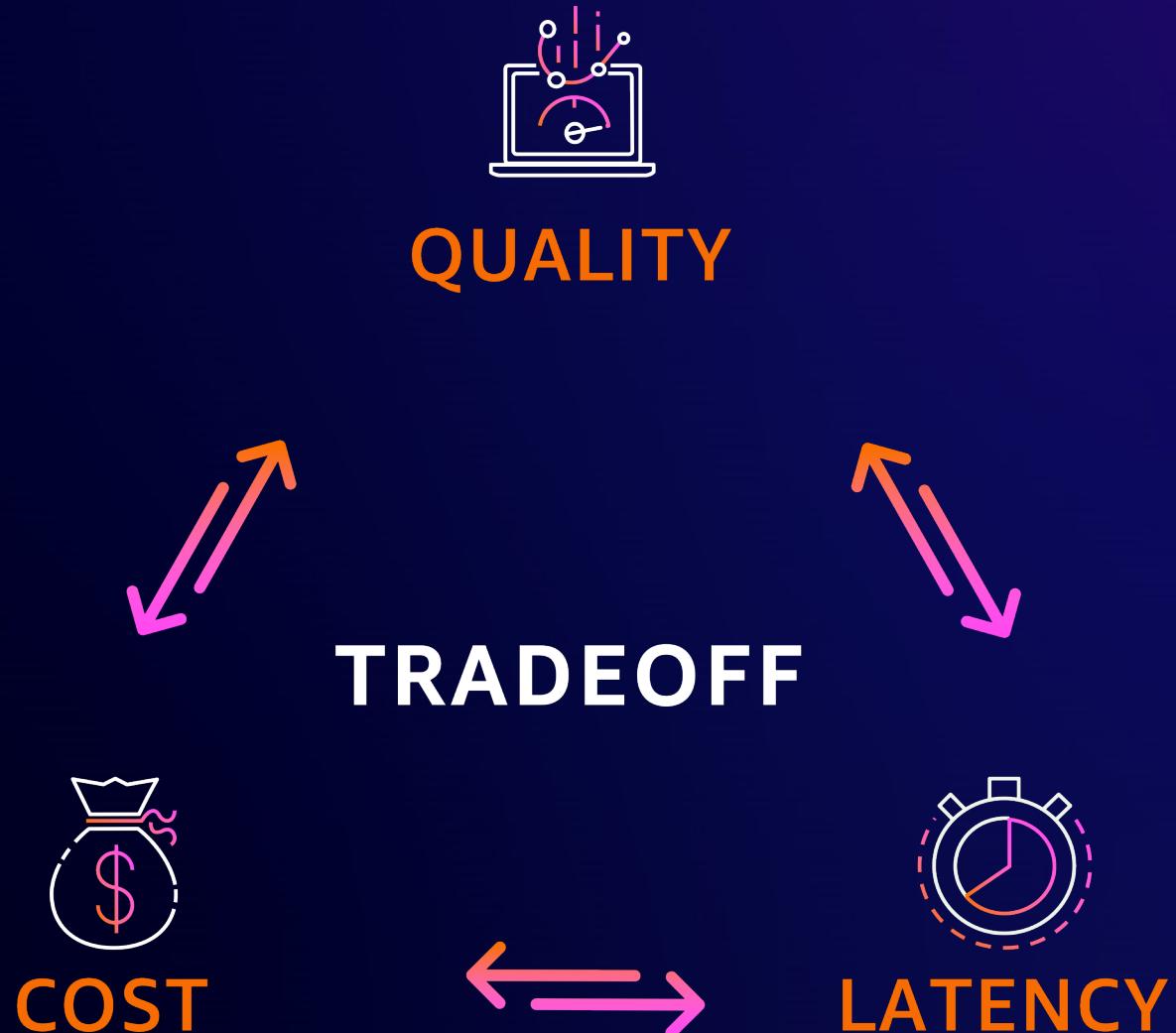
Model + RAG + Agent Evaluation

What is model evaluation?

1. Quality



What is model evaluation?



Selecting an LLM: Multiple priorities

Cost > Precision > Speed

!! Medium priority

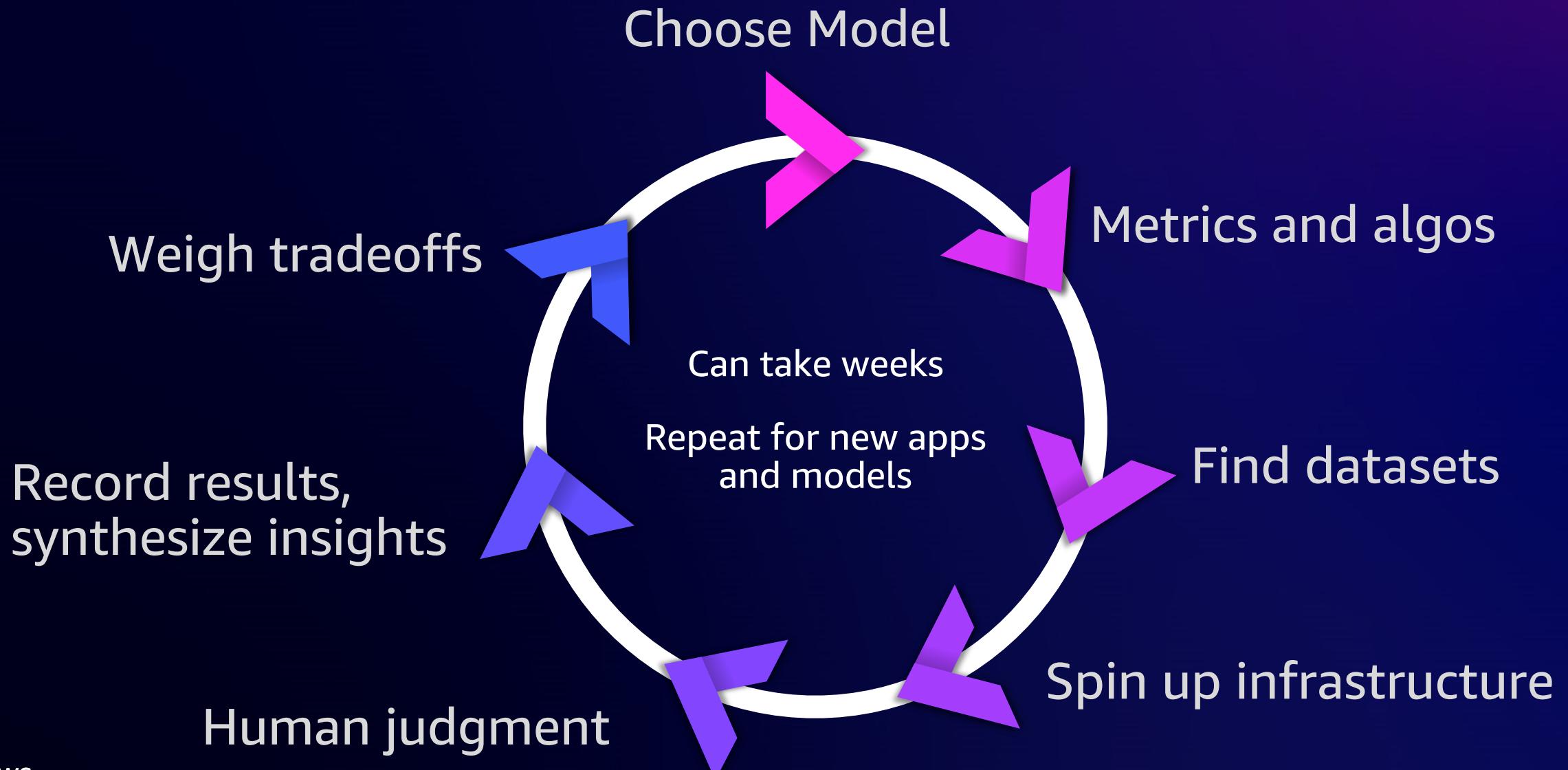
Precision



Why is foundation model evaluation important?

-  Make quality, cost, and latency tradeoffs
-  Align to your company's style and brand voice
-  Evaluate for your specific use cases
-  Evaluate with your company's data
-  Monitor biases, safety, and trust

Model evaluation lifecycle and challenges



Amazon Bedrock model evaluation

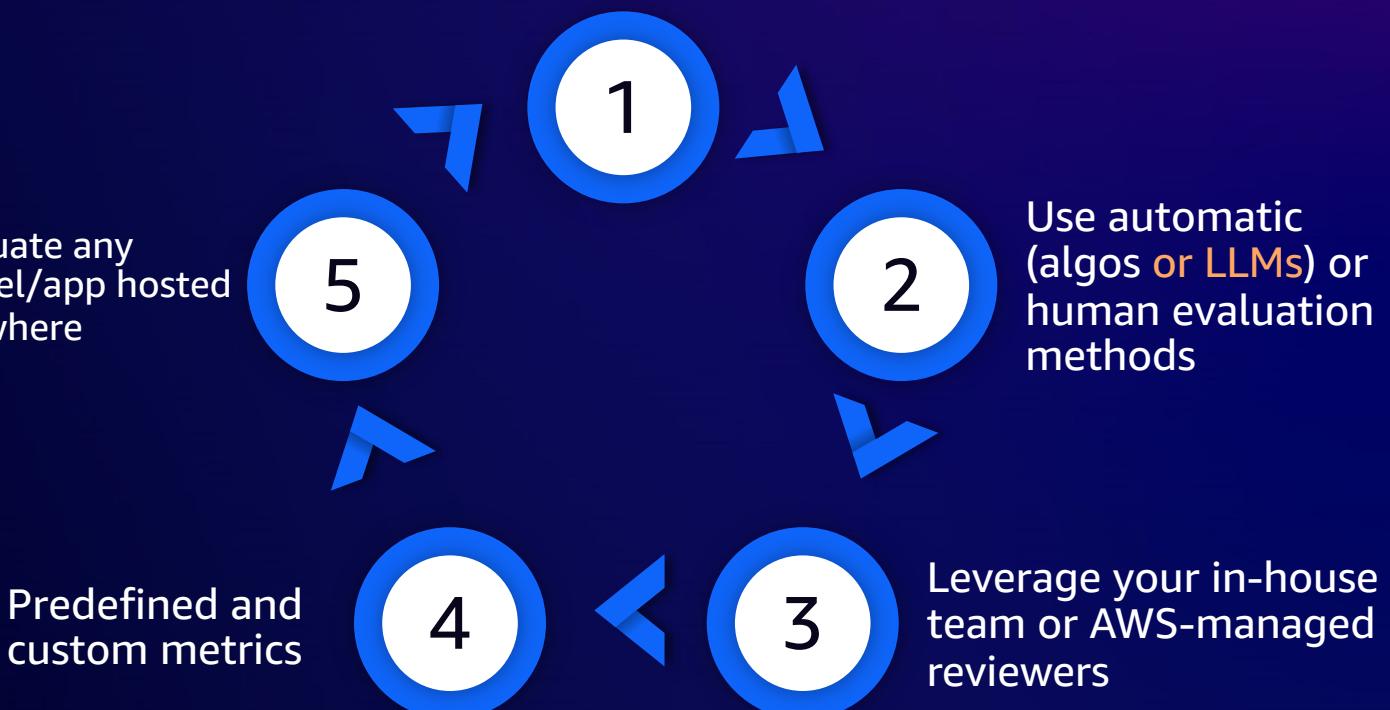
Evaluate, compare, and select the best foundation model for your use case

Public API

- Evaluate custom models
- Evaluate distilled models
- Evaluate imported models
- Evaluate prompt routers
- Any model hosted outside Bedrock
- Use an LLM-as-a-judge**

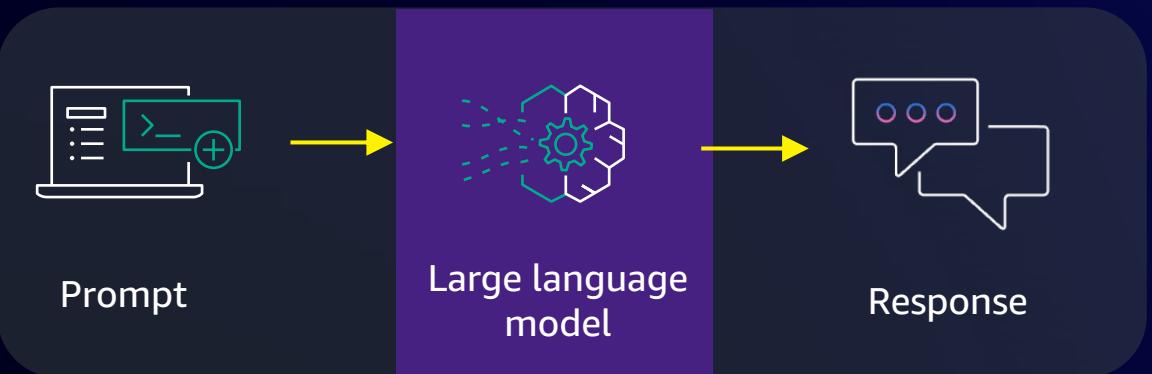
Evaluate any model/app hosted anywhere

Use curated datasets or bring your own for tailored results

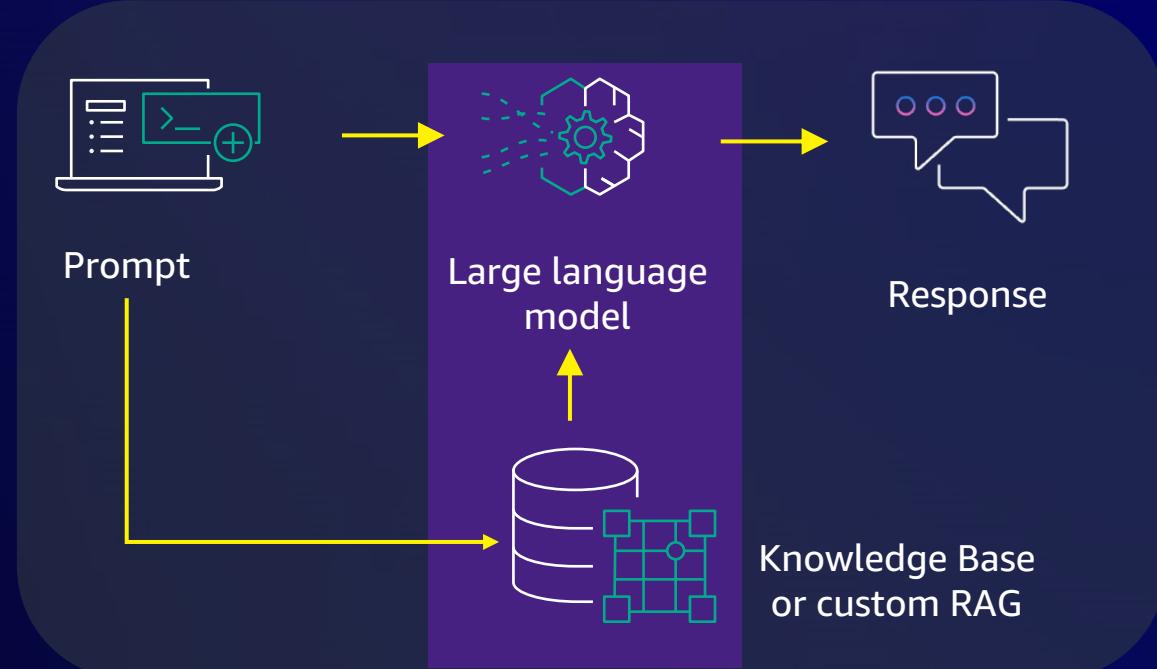


Evaluation Features - Bedrock

GenAI Model Evaluation



RAG Evaluation



Bedrock LLM evaluation options

Programmatic evaluation



Accuracy



Robustness



Toxicity

LLM-as-a-judge



Correctness



Completeness



Helpfulness



Relevance



Coherence



Readability

Human evaluation



Creativity



Style



Tone



Accuracy



Consistency



Brand voice

Algorithms

BERTScore | Classification accuracy
F1 | Real-world knowledge score

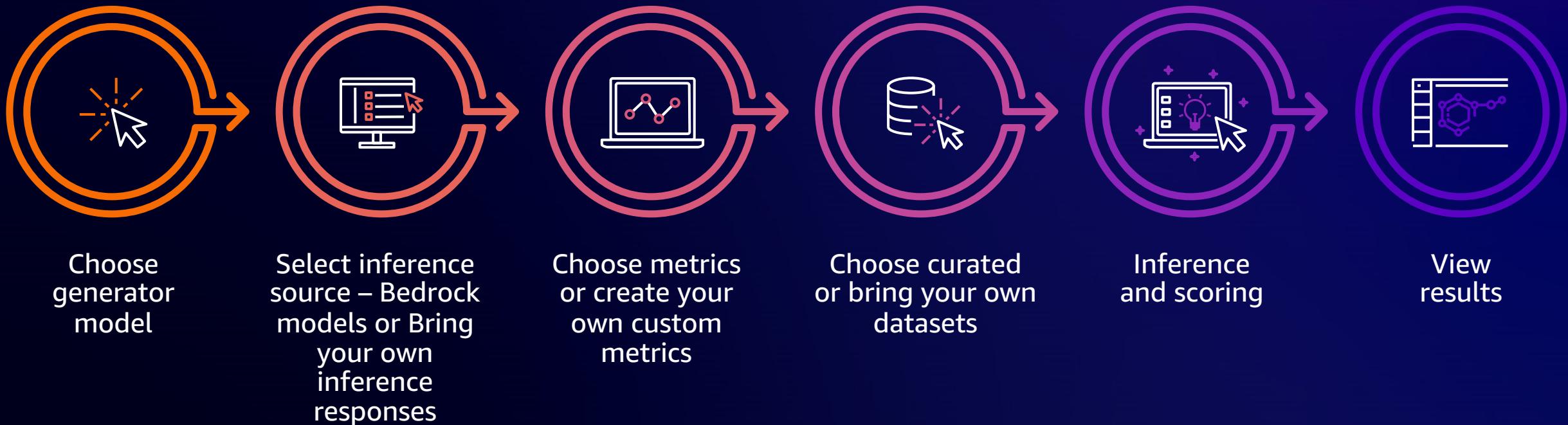
LLM reasoning

Multistep reasoning | Correlation with expert
human evaluators

Rating methods

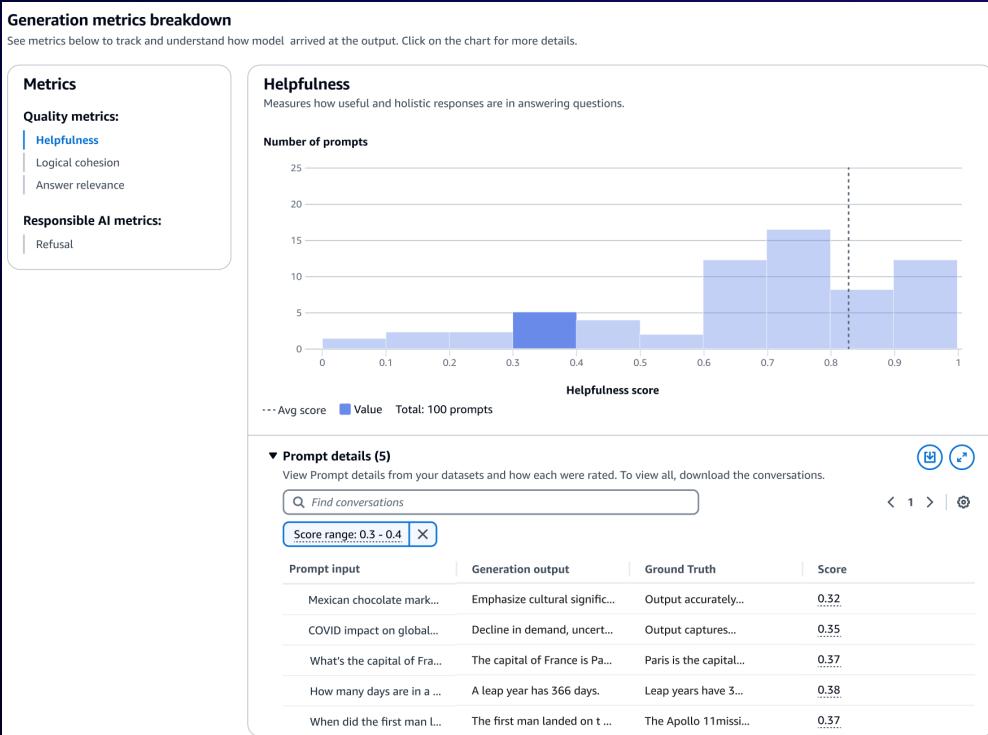
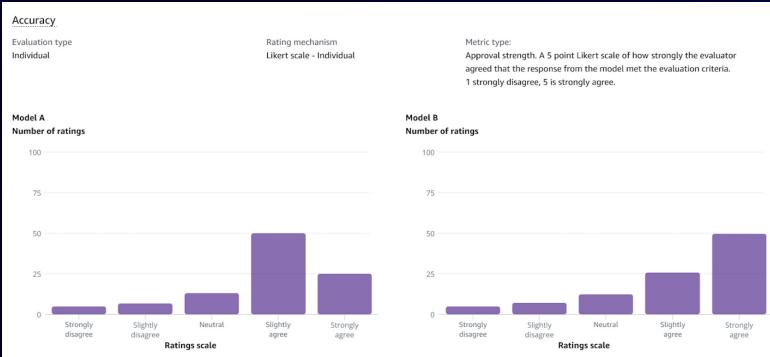
Thumbs up/down | 5-point Likert scales
Binary choice buttons | Ordinal ranking

How to set up LLM-as-a-judge model evaluations



Get results in a few clicks

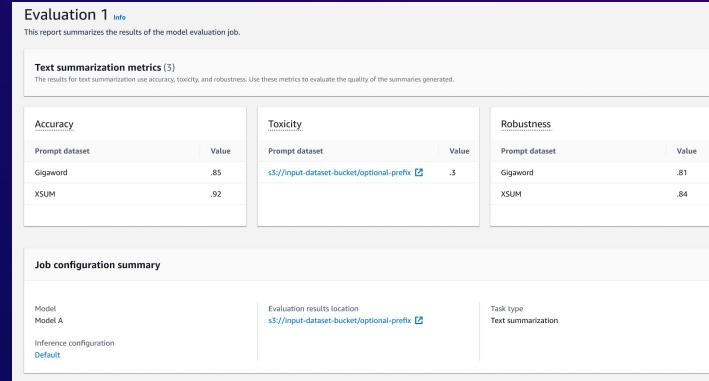
Human evaluation reports



Track ratings from your team
See distributions visually
Simple metric explanations

Simple to read scores
See distributions visually
See ratings explanations

Automatic evaluation reports



Bring your own inference responses – model eval

Bedrock Model

Amazon Bedrock > Evaluations > Create automatic evaluation using model as a judge

Create automatic evaluation Info

Choose the models for your evaluation job. The Generator model is the model that performs inference on the prompt dataset you provide. The Evaluation model is the model that will assess the responses to the prompt based on the metrics you select.

Model evaluation details

Evaluation name
Name must be unique within your account's AWS region.

Valid characters are a-z (lowercase), 0-9, and - (hyphen). The name can have up to 50 characters.

Description - optional

The description can have up to 200 characters.

Evaluator model
Choose a model for automatic evaluation that generates evaluation metrics.
[Select model](#)

Tags - optional

Inference source Info

Select source
Select Bedrock models or provide custom pre-generated inference with dataset.

Bedrock models
Select models provided by the Bedrock service.

Bring your own inference responses
With Bring your own inference responses, you provide both the input prompt and the generated response in your input JSONL file.

Select model
Choose the model you want to evaluate. If you can't find the model you're looking for, check [model access](#).
[Select model](#)

Metrics (2/12)
Select metrics for evaluation. Different metrics have different cost implications.

Quality (1/9)
Evaluate the quality and accuracy of the generated responses in relation to the original input or source.

Helpfulness
Measures how useful and holistic responses are in

Correctness
Measures how correct the responses are in

Any model anywhere (or app)

Amazon Bedrock > Evaluations > Create automatic evaluation using model as a judge

Create automatic evaluation Info

Choose the models for your evaluation job. The Generator model is the model that performs inference on the prompt dataset you provide. The Evaluation model is the model that will assess the responses to the prompt based on the metrics you select.

Model evaluation details

Evaluation name
Name must be unique within your account's AWS region.

Valid characters are a-z (lowercase), 0-9, and - (hyphen). The name can have up to 50 characters.

Description - optional

The description can have up to 200 characters.

Evaluator model
Choose a model for automatic evaluation that generates evaluation metrics.
[Select model](#)

Tags - optional

Inference source Info

Select source
Select Bedrock models or provide custom pre-generated inference with dataset.

Bedrock models
Select models provided by the Bedrock service.

Bring your own inference responses
With Bring your own inference responses, you provide both the input prompt and the generated response in your input JSONL file.

Source name

The name can have up to 256 characters, no spaces. Valid characters are A-Z, a-z, 0-9, - (hyphen), _ (underscore) and . (period).

Metrics (2/12)
Select metrics for evaluation. Different metrics have different cost implications.

Quality (1/9)
Evaluate the quality and accuracy of the generated responses in relation to the original input or source.

Helpfulness
Measures how useful and holistic responses are in

Correctness
Measures how correct the responses are in

- Evaluate any model or app hosted anywhere. No longer limited to Bedrock model calls for model eval
 - Bring your responses in your input dataset (JSONL), skip the Bedrock model call (additional input JSONL fields supported)
- Pro tip:** Not limited to models. If you only want to eval the response, it can be from anything, including whole apps



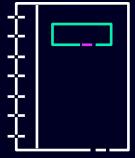
Amazon Bedrock RAG Evaluation

Evaluate your full RAG stack to optimize your application, including any application hosted anywhere

- 1 Bring your own datasets for tailored results
- 2 Evaluate retrieval alone or retrieval + generation with a choice of LLM-as-a-judge
- 3 Built-in metrics for quality and responsible AI, create custom metrics, compatible with Amazon Bedrock Guardrails
- 4 Compare across multiple evaluation jobs
- 5 Evaluate any RAG system hosted anywhere

Choice of RAG evaluation metrics

Retrieval



Context
Coverage



Context
Relevance

Retrieval + Generation



Correctness



Completeness



Helpfulness



Citation
precision



Citation
coverage



Logical coherence



Faithfulness

Responsible AI



Harmfulness



Stereotyping



Refusal

Custom metrics

Type metric rubric (templates available)



Define grading scale



See preview of the full judge prompt



Custom metrics - optional [Info](#)

Custom metrics 1

Metric name
Name must be unique within your account's AWS Region.

Choose metric type [Info](#)

Import JSON file Import predefined metric configuration via JSON. [Learn about file format](#)

Use a template Select from pre-configured metric templates.

Custom Create metric with full configuration control. [View best practices](#)

Instructions [Info](#)
Write a tailored prompt to provide instructions to the model judge to evaluate your Knowledge base. [View best practices](#)

Enter instructions e.g. You are given a task, a candidate answer and a ground truth answer. Based solely on the ground truth answer, assess whether the candidate answer is a correct and accurate response to the task.

Instructions must contain: ✓ Response

Output schema enabled (recommended) [Info](#)
Schema will specify the structure for presenting the custom metric results including visualization, score normalization (if applicable) and model generated explanations.

Scale type

Value
 Description
 [+ Add](#)

Add up to 10 values The description can have up to 63 characters.

Preview
Review the judge prompt that will be provided to the evaluator model based on instructions entered above.

No instructions to display.

[+ Add custom metrics](#)

Add up to 9 more metrics. You can also [download](#) all custom metrics entered above.

Evaluator model for custom metrics [Info](#)
Choose a model to compute evaluation custom metrics.

[Select model](#)

- Define your own metric for LLM-as-a-judge in model and RAG evaluation
- Variables for data injection from dataset and responses
- Save and reuse metrics later



Bring your own inference responses – RAG eval

Bedrock KB

Amazon Bedrock > Evaluations > Create automatic Knowledge Base evaluation

Create automatic evaluation: RAG

Choose a RAG system, select metrics, and relevant datasets to assess its performance.

Evaluation details

Evaluation name
Name must be unique within your account's AWS region.
 rag-evaluation-job-quick-start-20250325033848
Valid characters are a-z (lowercase), 0-9, and - (hyphen). The name can have up to 50 characters.

Description - optional
 Provide a description
The description can have up to 150 characters.

Evaluator model
Choose a model to compute evaluation metrics.
 Select model

Tags - optional

Inference source

Select source
Choose Bedrock Knowledge Bases or provide custom pre-generated RAG inference with dataset.

Bedrock Knowledge Base
Select Knowledge Base created within Bedrock services.

Bring your own inference responses
Provide both prompts and responses you already retrieved in your input dataset.

Choose a Knowledge Base
Select a previously created Knowledge Base or create a new one by navigating to [Knowledge Base page](#).

Evaluation type
Evaluate only the retrieval of text chunks or both retrieval and response generation (different metrics apply).

Retrieval and response generation
Evaluate both information retrieval from data sources and response generation using LLM model.

Retrieval only
Provide both prompts and responses you already retrieved in your input dataset.

Response generator model

Your own RAG system

Amazon Bedrock > Evaluations > Create automatic Knowledge Base evaluation

Create automatic evaluation: RAG

Choose a RAG system, select metrics, and relevant datasets to assess its performance.

Evaluation details

Evaluation name
Name must be unique within your account's AWS region.
 rag-evaluation-job-quick-start-20250325033848
Valid characters are a-z (lowercase), 0-9, and - (hyphen). The name can have up to 50 characters.

Description - optional
 Provide a description
The description can have up to 150 characters.

Evaluator model
Choose a model to compute evaluation metrics.
 Select model

Tags - optional

Inference source

Select source
Choose Bedrock Knowledge Bases or provide custom pre-generated RAG inference with dataset.

Bedrock Knowledge Base
Select Knowledge Base created within Bedrock services.

Bring your own inference responses
Provide both prompts and responses you already retrieved in your input dataset.

Source name

The name can have up to 256 characters and no spaces. Valid characters: a-z, A-Z, 0-9, -(hyphen), _(underscore) and .(period).

Evaluation type
Evaluate only the retrieval of text chunks or both retrieval and response generation (different metrics apply).

Retrieval and response generation
Evaluate both information retrieval from data sources and response generation using LLM model.

Retrieval only
Provide both prompts and responses you already retrieved in your input dataset.

- Evaluate any RAG system hosted anywhere. Not limited to Bedrock Knowledge Bases
- Bring your responses and retrieved passages in your input dataset (JSONL), skip the Bedrock KB call (additional JSONL fields in input dataset)



How RAG evaluation works



RAG Evaluation Reports

aws Mezzanine

[Amazon Bedrock](#) > [Evaluation](#)

Evaluation Info

Models [Knowledge Bases - new](#)

Knowledge Base evaluation Info
Assess the performance or effectiveness of your Knowledge Base using an AI-powered model or a human-powered team. [Learn more](#)

▶ How it works

Knowledge Base evaluations (6)

[Compare](#) [Duplicate](#) [Create](#) [⋮](#)

Find evaluation

Evaluation ...	Status	Knowledge Base name	Evaluation
<input checked="" type="checkbox"/> Travelapp_Eval	Complete	Travelapp_Eval	Clau
<input checked="" type="checkbox"/> Travel&leisure	Complete	knowledge-base...	Clau
<input type="checkbox"/> Financeapp	Complete	knowledge-bases-21	Clau
<input type="checkbox"/> Financeapp_1	Complete	knowledge-bases-21	Clau
<input type="checkbox"/> eval_1	Complete	knowledge-bases-21	Clau
<input type="checkbox"/> eval_1	Complete	knowledge-bases-55	Clau

CloudShell Feedback Language

© 2023, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Compare Knowledge Bases evaluation metrics
Select 2 Knowledge Bases to compare their performance.

▶ Evaluation overview (2)

At a glance Evaluation details

Comparison: Travelapp_Eval and Travel&leisure

[Download](#)

Helpfulness
Harmlessness
Completeness
Logical cohesion
Faithfulness
Stereotyping
Correctness
Refusal

Legend:
Travelapp_Eval (Blue)
knowledge-base-quick-start-fwvzy (Red)

aws Mezzanine

[Amazon Bedrock](#) > [Evaluation](#) > Report: Travelapp_Eval

Report: Travelapp_Eval Info

[Duplicate evaluation](#) [Download](#) [⋮](#)

Evaluation summary

Metrics summary [Info](#)
Evaluate overall performance using metrics (average score across all conversations). Closer to 1 is a higher score, closer to zero is a lower score. For example, closer to 1 for Correctness means more correct answers. You can define custom criteria to highlight any metrics that fall above or below a threshold.

Quality metrics [Info](#)

Metric	Score
Helpfulness	0.32
Faithfulness	0.65
Correctness	0.37

Generation metrics breakdown
See metrics below to track and understand how Knowledge Base arrived at the output. Click on the chart for more details.

Metrics
Quality metrics:
Helpfulness
Faithfulness
Correctness

Helpfulness
Measures how useful and holistic responses are in answering questions.

100 examples

Range	Count
0.0 - 0.1	10 examples
0.1 - 0.2	15 examples
0.2 - 0.3	18 examples
0.3 - 0.4	5 examples
0.4 - 0.5	2 examples
0.5 - 0.6	8 examples
0.6 - 0.7	5 examples
0.7 - 0.8	20 examples
0.8 - 0.9	22 examples
0.9 - 1.0	10 examples

▶ Example conversations

Faithfulness
Measures how aligned the information in the responses are with the information in the original input or source.

Number of example conversations

Score Range	Count
0.0 - 0.1	2
0.1 - 0.2	5
0.2 - 0.3	10
0.3 - 0.4	15
0.4 - 0.5	12
0.5 - 0.6	8
0.6 - 0.7	5
0.7 - 0.8	8
0.8 - 0.9	10
0.9 - 1.0	12

Avg score Value Total: 100 conversations

▶ Example conversations

Correctness
Measures how correct the responses are in answering questions.

Number of example conversations

Score Range	Count
0.0 - 0.1	2
0.1 - 0.2	5
0.2 - 0.3	10
0.3 - 0.4	15
0.4 - 0.5	12
0.5 - 0.6	22
0.6 - 0.7	15
0.7 - 0.8	10
0.8 - 0.9	8
0.9 - 1.0	12

Avg score Value Total: 100 conversations

▶ Example conversations



Agent Evaluation

Why it matters for production?

Agents can take many paths to a correct result

Success isn't just about the outputs

Common LLM evaluation methods are not sufficient

Recommendation 1 Step-by-Step evaluation

Agents components can fail independently (router/skills/tools)

- Assess router performance on tool selection
- Add skill specific evaluation(RAG, Code, API)
- Use both deterministic and LLM-as-a-judge approaches

Recommendation 2 Track path and convergence

Agents can get stuck in loops, or take inefficient paths, this impacts cost, latency and user experience

- Use an iteration counter and track step sequences
- Manually inspect traces during implementation
- Monitor for loops, unnecessary returns to router and redundancy
- Use both deterministic and LLM-as-a-judge approaches

Recommendation 3 Implement test case strategy

Test as many query variations as possible, include edge-cases, off-topic queries and parameter overlaps

- Evolve the cases with production data and discovered failures
- Focus on path coverage and not on volume
- Use both deterministic and LLM-as-a-judge approaches

Evaluation approaches

- Deterministic: anything that can be coded or calculated
- LLM-as-a-Judge(LLMaJ): Qualitative assessment – using stress-tested prompts
- Hybrid: A combination of both for comprehensive coverage
- AgentCore Eval

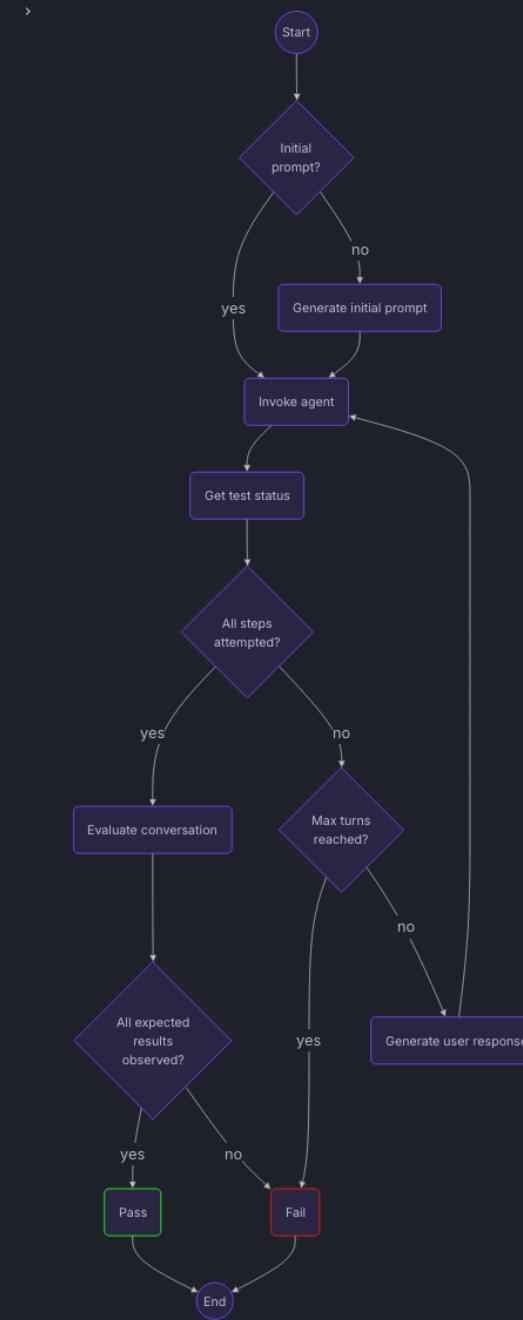
Agent evaluation sample



<https://awslabs.github.io/agent-evaluation/evaluators/>

Evaluation workflow

The diagram below depicts the workflow that is conducted during evaluation.



Open source options for evaluation

- Available templates and guidance to run agent evaluation
- Easy to integrate with Bedrock and AgentCore observability



Welcome to the final day of the Learning Week



Production ready agents with: **Bedrock AgentCore**



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

Implementing an AI Agent is complex

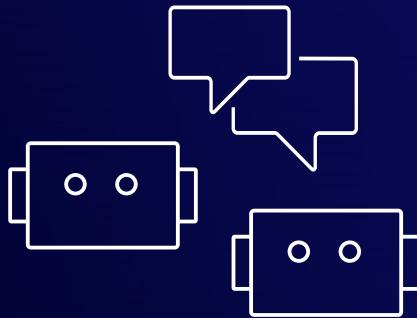
But open source frameworks make it easy to build agents and create proof of concepts (PoCs)



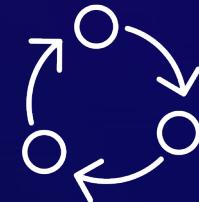
Context



Tools and tools
execution



Multi-agent
collaboration



Orchestrating between
actions



 Strands Agents

 LangChain

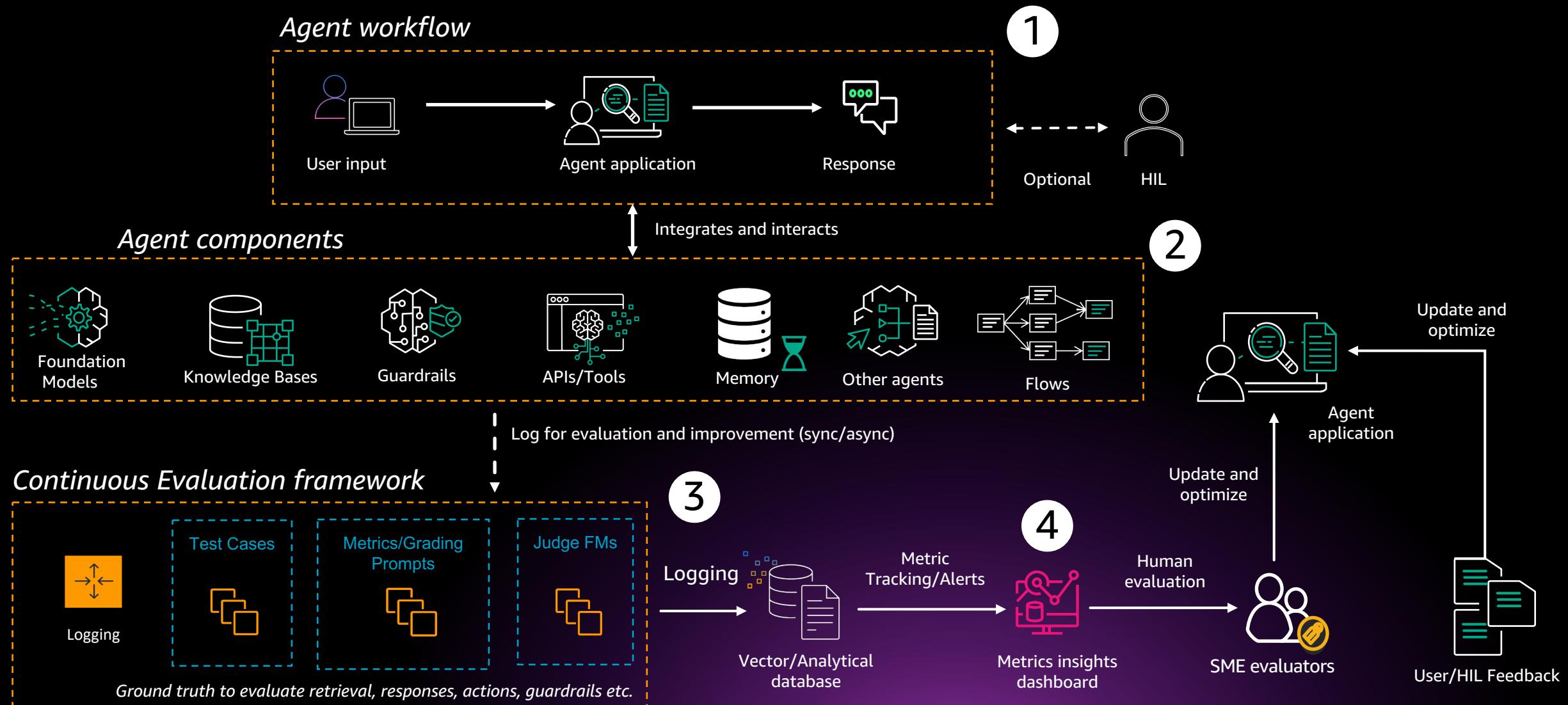
 LangGraph

 crewai

Building an agent is the start, production scale is the goal



What does an agentic system look like?



Scale requires operational excellence



Fundamentals for agentic AI Scale

SECURITY & TRUST

Ensure agents can act safely across system boundaries, protect data, ensure compliance

RELIABILITY & SCALABILITY

Support autonomous actions reliably at enterprise scale

IDENTITY

Authenticate and authorize all actors – agent and users

OBSERVABILITY

Track and monitor agent behaviors and decisions

DATA FOR CONTEXT

Feed agents with rich, contextually aware knowledge

INTEGRATIONS

Connect seamlessly across systems and tools

MEMORY

Agents need to maintain context across interactions

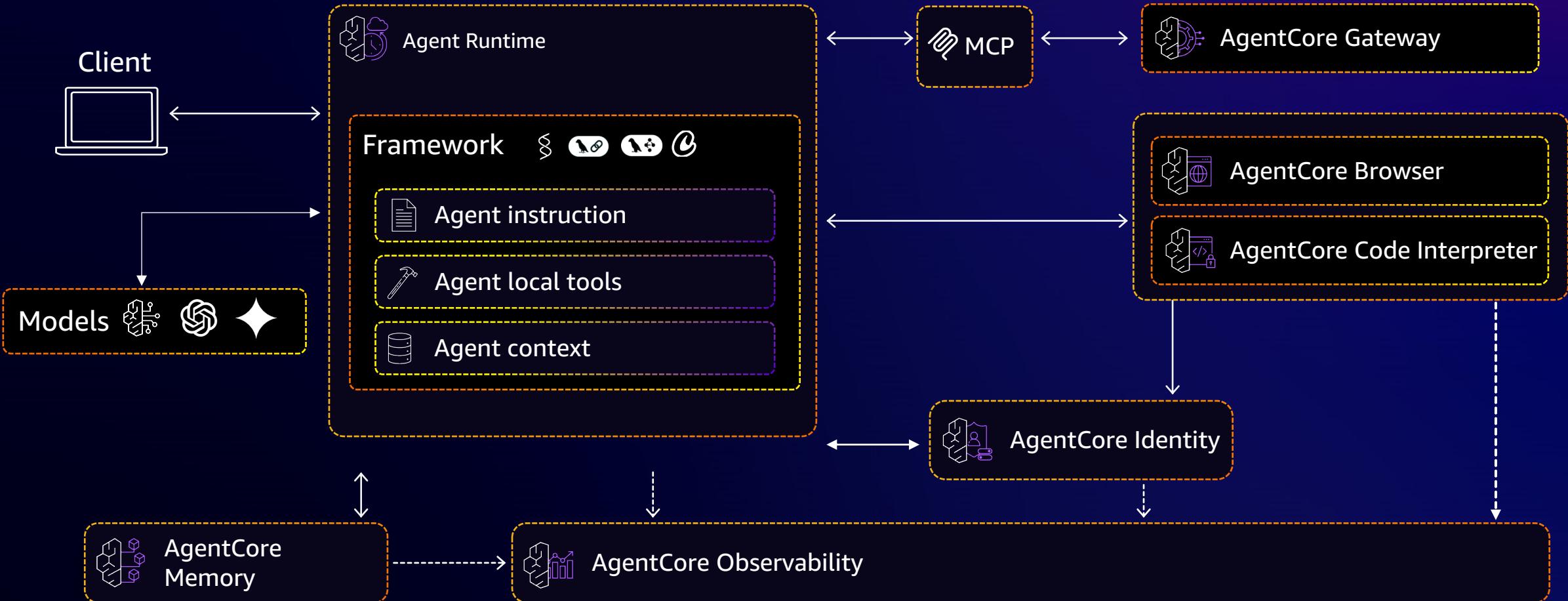
BROWSER ACTIONS

Agents need to navigate websites to perform actions

EXECUTE CODE

Agents need to execute code and run applications on user's behalf

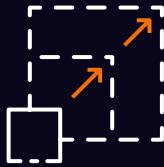
AgentCore capabilities enabling agents at scale





AgentCore Runtime

Scale from real-time to multi-hour workloads



- Scale from real-time to multi-hour workloads with low latency and industry-leading extended runtime (up to 8 hours)
- Supports payloads across modalities

Accelerate time to market



- Deploy any AI agent using common open-source frameworks
- Deploy from POC to production in a few lines of code

Secure workload with enterprise-grade isolation



- True session isolation to protect your data
- Integrates with existing identity providers



AgentCore Runtime

Agent or tool code



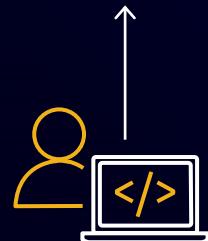
configure



AgentCore Runtime

Amazon ECR Repository

launch



invoke



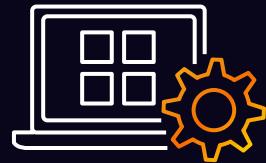
AgentCore Gateway

Simplified tool development & integration



- Turn APIs, Lambda functions, and existing services into MCP-compatible tools
- Access thousands of tools through a standardized interface

Secure and unified access



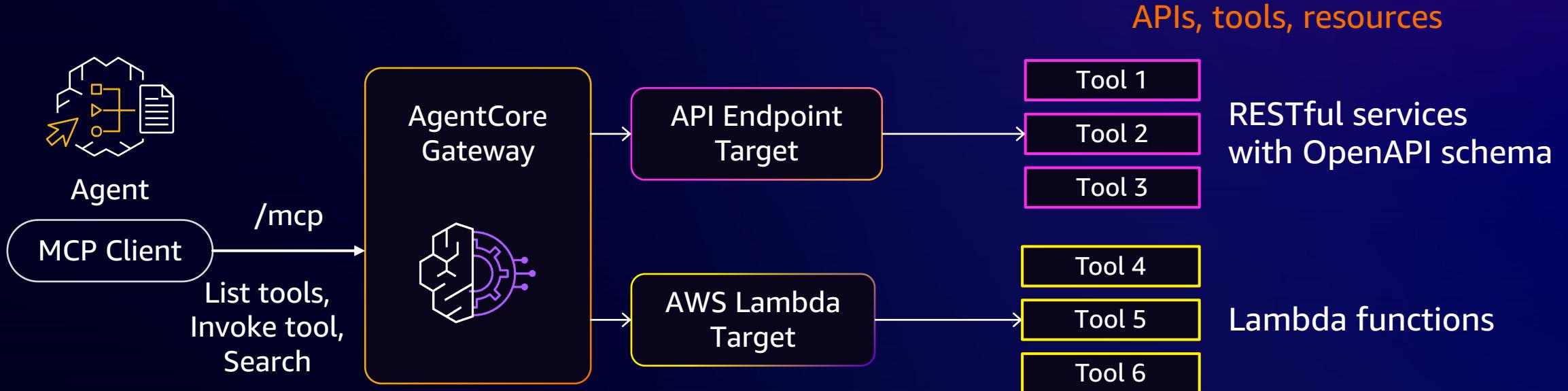
- Discover and use tools through a single, secure endpoint
- Combine multiple tools sources into one unified interface

Intelligent tool discovery



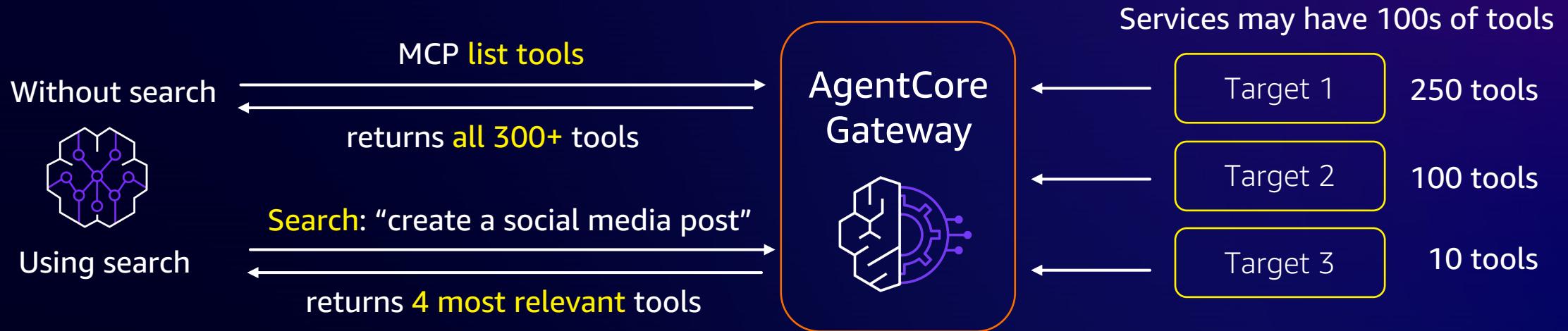
- Enable agents to find the right tools with context aware discovery
- Curated tool collections with granular permissions

AgentCore Gateway





AgentCore Gateway semantic search



Benefits

- AgentCore Gateway automatically indexes tools, and gives serverless semantic search
- Reduces context passed to the agent's LLM, improving accuracy, speed, and cost
- Lets agent focus on tools relevant for a given task

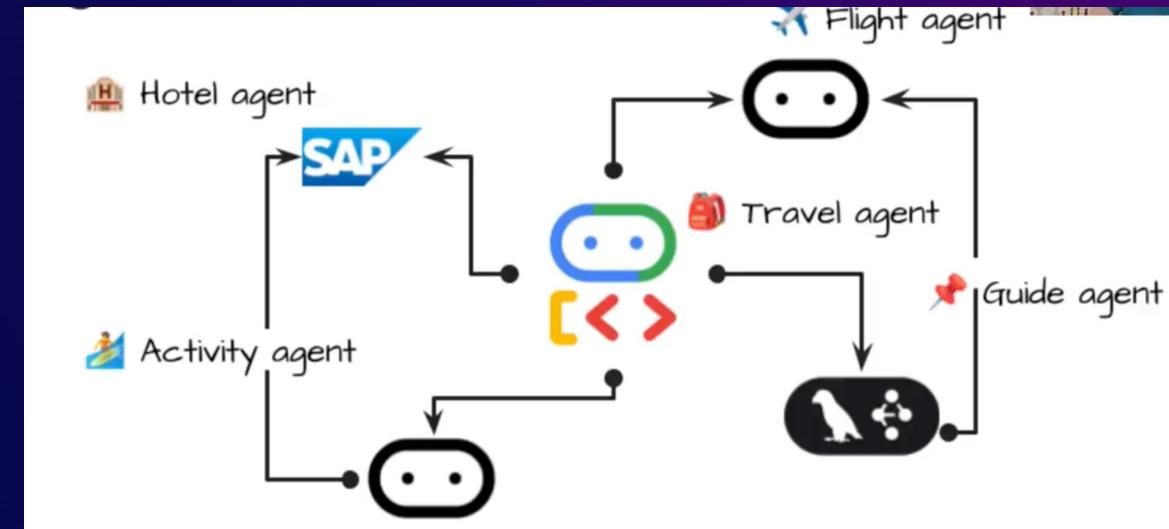
A2A



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

A2A Protocol

A2A is an open protocol that provides a standard way for agents to collaborate with each other, regardless of the underlying framework or vendor.



Key Design Principles

- **Embrace agentic capabilities:** Collaborate in natural, unstructured modalities, even when agents don't share memory, tools and context.
- **Build on existing standards:** Easy integration with existing IT stacks businesses use as the protocol is built on top of existing, popular standards including HTTP, SSE, JSON-RPC
- **Secure by default:** Supports enterprise-grade authentication and authorization, with parity to OpenAPI's authentication schemes at launch.
- **Support for long-running tasks:** Support for long-running tasks like deep research that may take hours and or even days when humans are in the loop. A2A can provide real-time feedback, notifications, and state updates to its users.
- **Modality agnostic:** Support for audio and video streaming

A2A Protocol

HOW IT WORKS

A2A protocol follows a client-server paradigm



Discovery: Agents can advertise their capabilities using an “Agent Card” in JSON format, allowing the client agent to identify the best agent that can perform a task and leverage A2A to communicate with the remote agent.



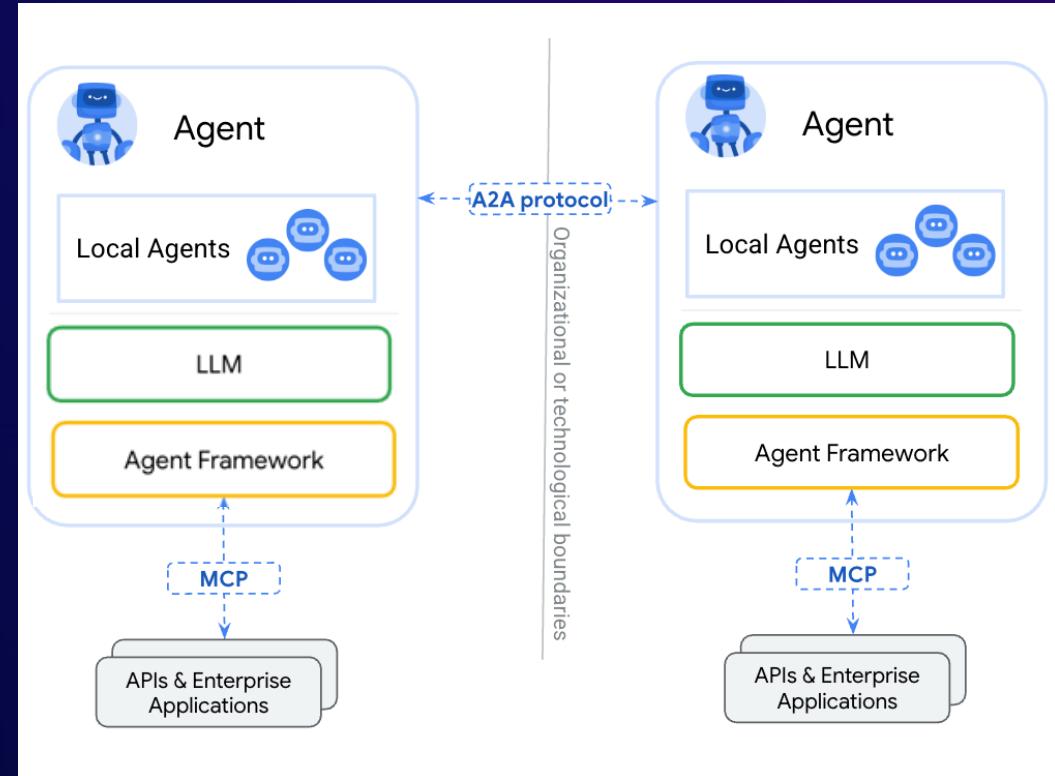
Task management: Client agent is responsible for formulating and communicating “tasks”, while the remote agent is responsible for acting on those tasks.



Collaboration: Agents can send each other messages to communicate context, replies, artifacts, or user instructions (long running tasks).



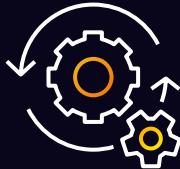
User Ex Negotiation: Each message includes “parts,” which is a fully formed piece of content, like a generated image. Each part has a specified content type, allowing client and remote agents to negotiate the correct format needed and explicitly include negotiations of the user’s UI capabilities—e.g., iframes, video, web forms, and more.





AgentCore Memory

Simplify memory management



- Abstracts memory infrastructure
- Scales automatically with serverless architecture
- Automatically stores and manages agent context across sessions

Enterprise-grade



- Complete data privacy with dedicated storage for each customer
- Enterprise security with encryption and regional data storage options

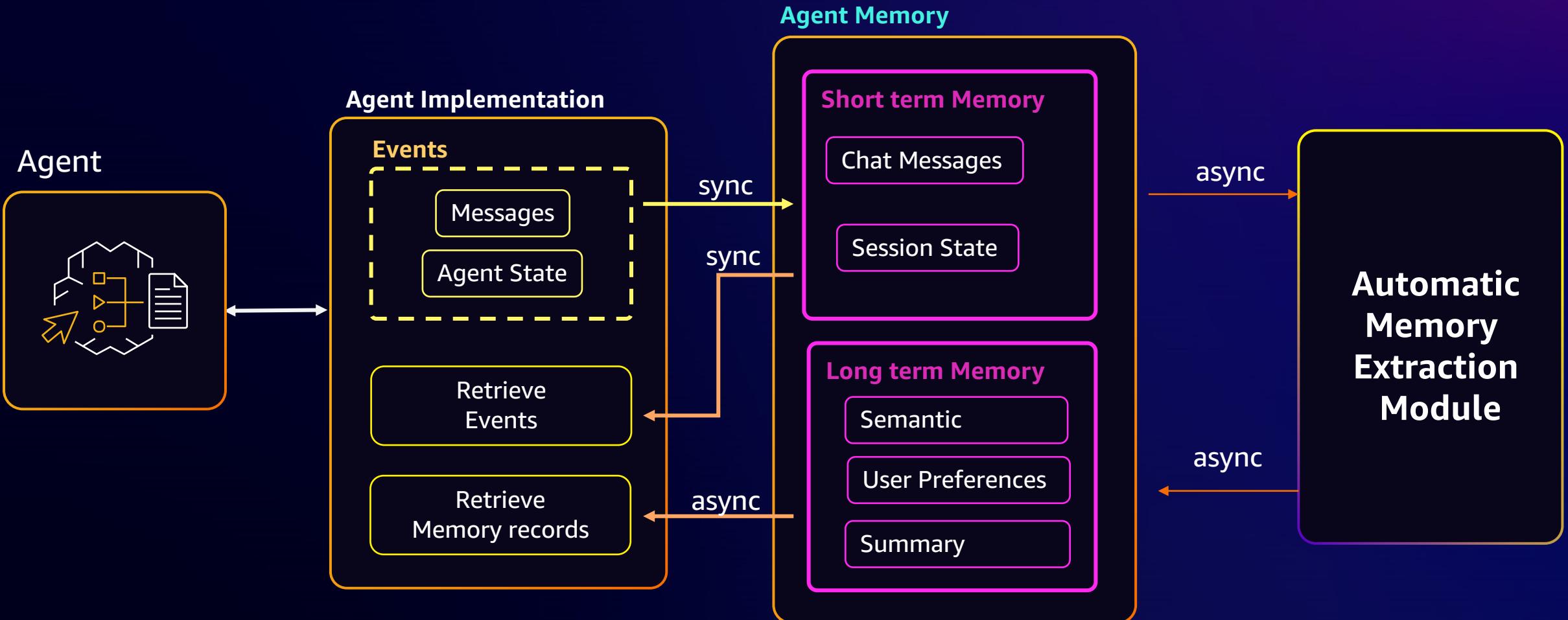
Deep customization



- Define memory patterns based on your use case
- Configure extraction rules
- Choose models and customize prompts for memory extraction



AgentCore Memory





AgentCore Browser

Serverless browser infrastructure



- Low latency browser sessions
- Auto-scales from 0 to hundreds of concurrent sessions

Enterprise-grade security



- Session isolated compute with VM-level isolation per user
- Secure credential handling

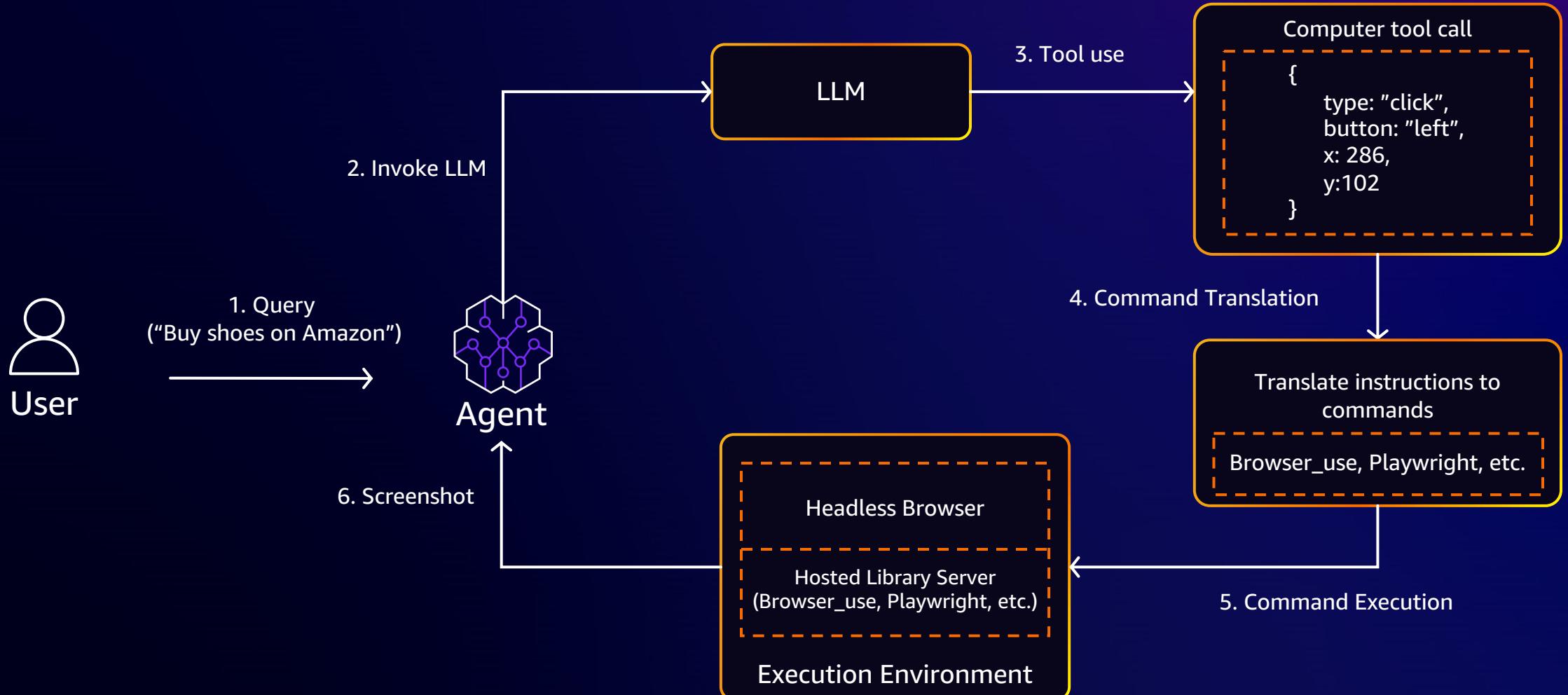
Enterprise observability



- Live streaming URLs for real-time monitoring
- Session replays for debugging
- Extensive logging of all browser commands to CloudTrail



AgentCore Browser





AgentCore Code Interpreter

Execute code securely



- Execute complex workflows and data analysis in isolated sandbox environments
- Access internal data sources securely without exposing sensitive data

Large-scale data processing



- Process gigabyte-scale datasets efficiently using Amazon S3 integration, without API limitations

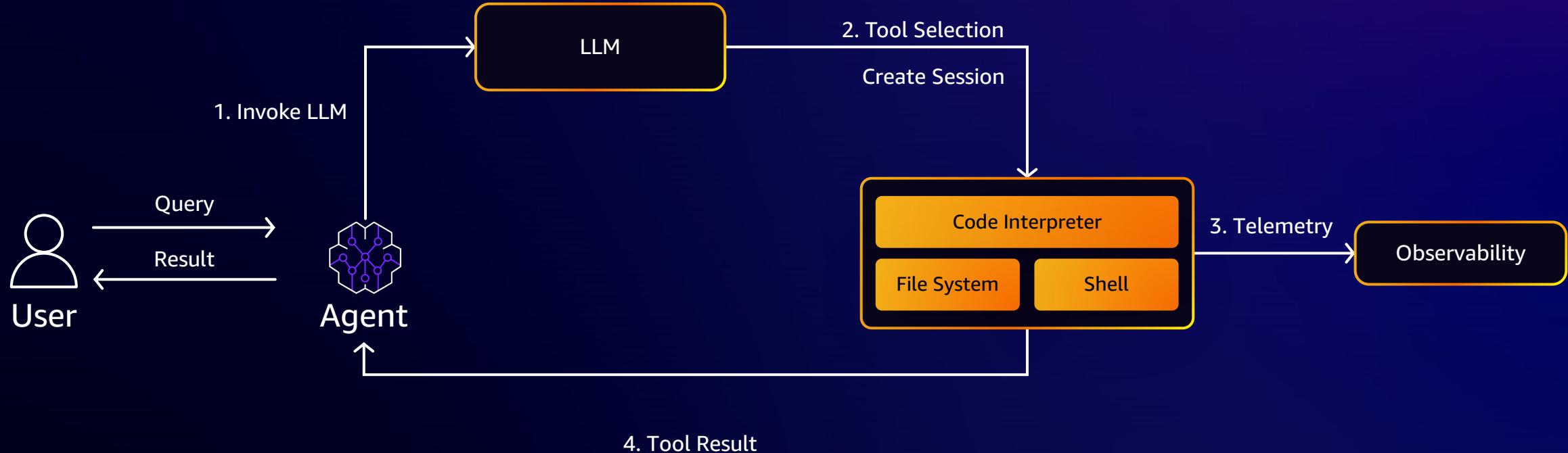
Ease of use



- Quick start with pre-built execution runtimes for JavaScript, TypeScript, and Python with common libraries pre-installed



AgentCore Code Interpreter





AgentCore Identity

Secure, delegated access for AI agents



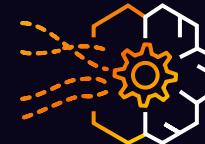
- Enables AI agents to securely access AWS resources and third-party tools such as GitHub, Google, Salesforce and Slack
- Robust access controls with just-enough access and secure permissions delegation

Build streamlined AI agent experiences



- Minimizes consent fatigue with a secure token vault
- Streamlines authentication flows

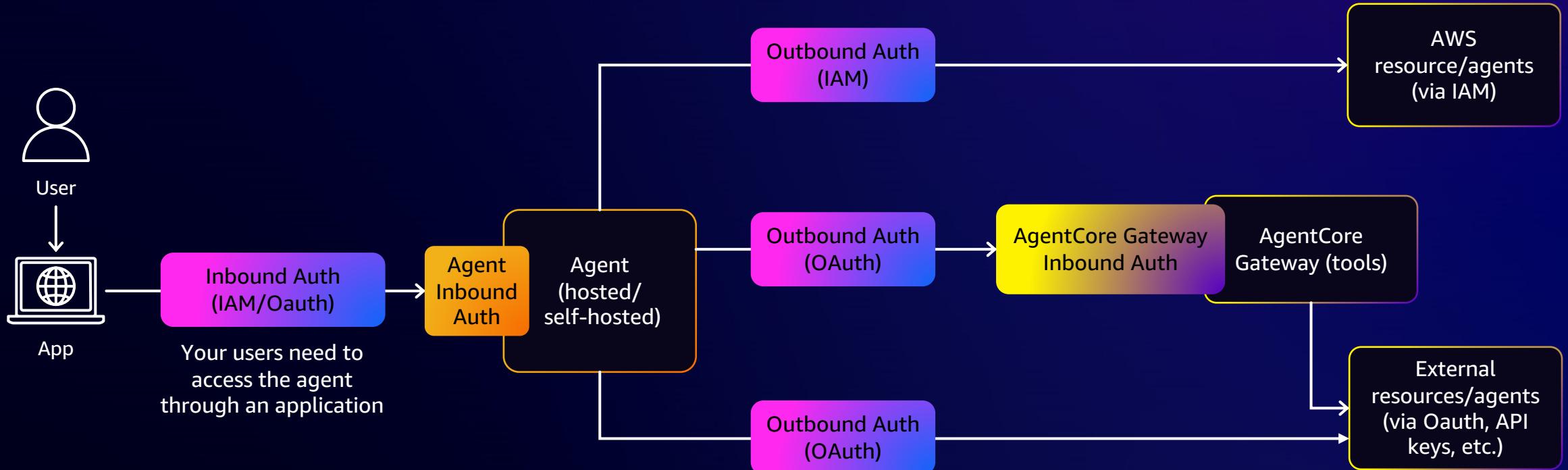
Accelerated AI agent development



- Preserves existing identity systems such as Okta, Azure AD, or Amazon Cognito
- Lowers custom development efforts without need for migrating users or rebuilding authentication flows

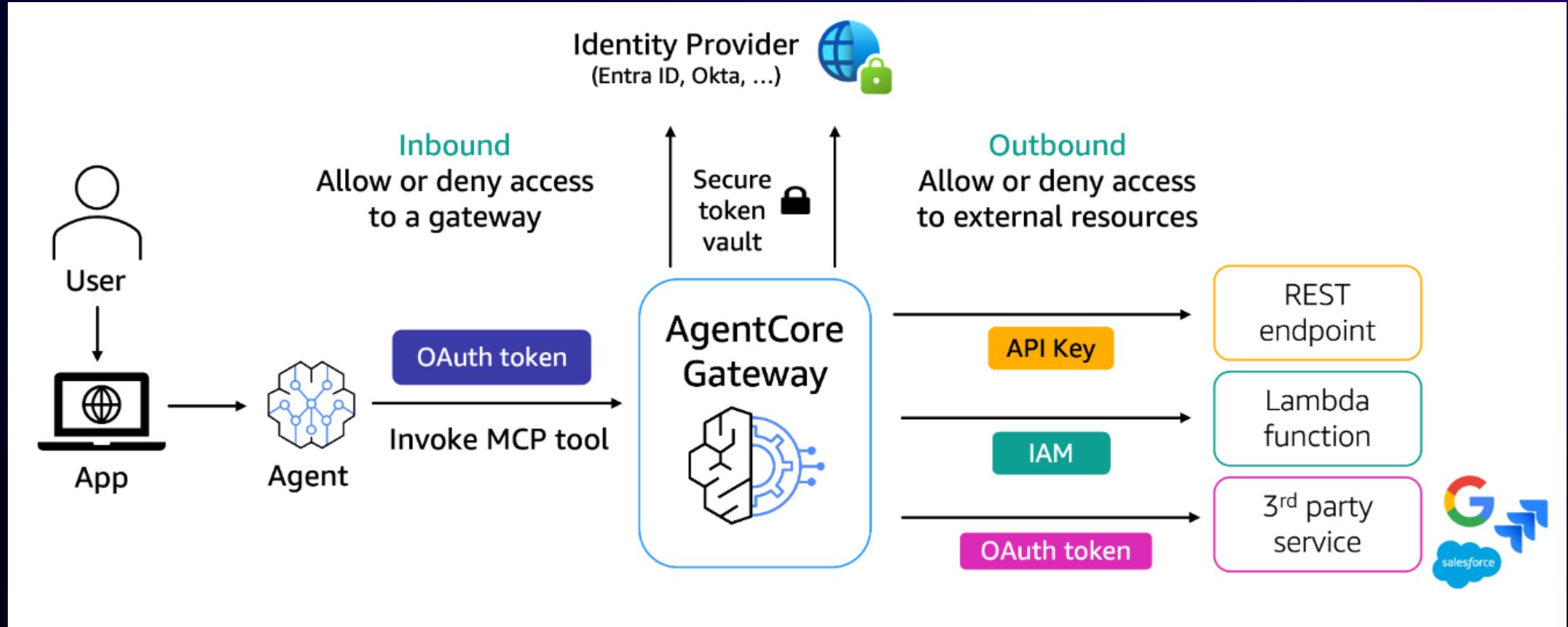


AgentCore Identity





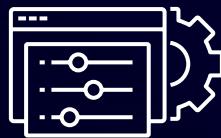
AgentCore Gateway & Identity





Policy in AgentCore: Complete control over agent actions

Control over agent actions



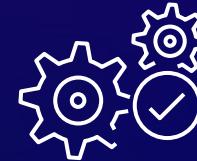
- Intercept every tool call before execution to ensure adherence to defined policy.
- Control agent access, actions across APIs, Lambda functions, and third-party services.

Frictionless policy enforcement

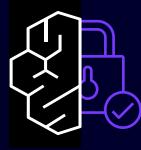


- Process thousands of requests per second with millisecond policy evaluation.
- Support granular access controls with secure VPC and PrivateLink integration.

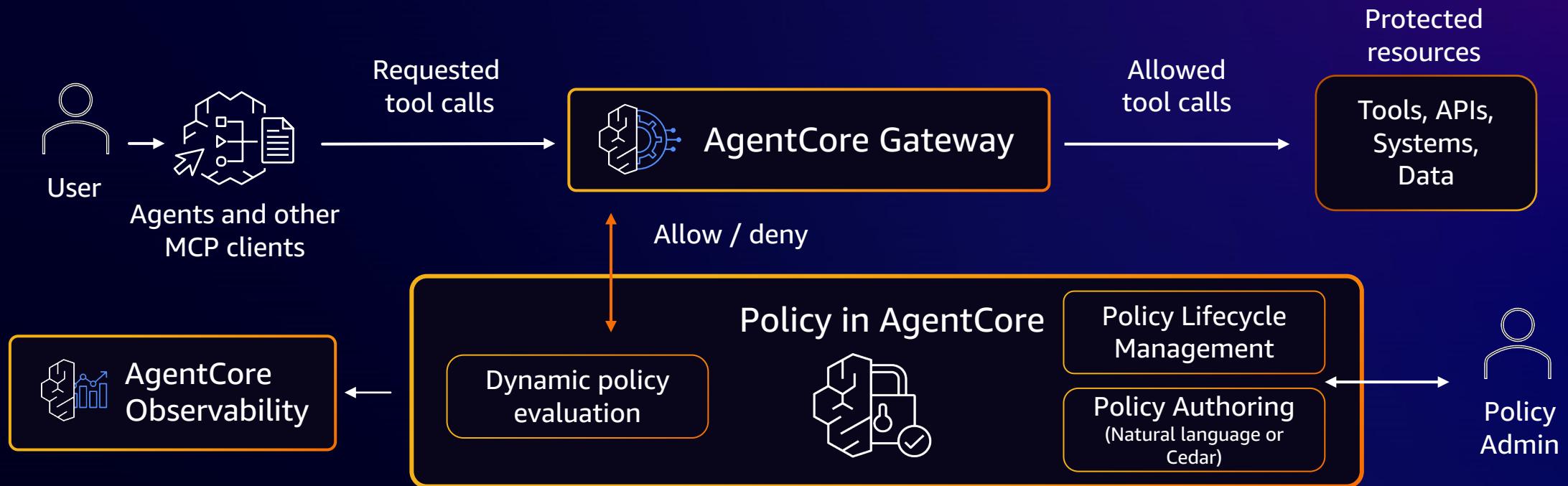
Easy policy creation and management



- Create and manage policies using natural language with automatic Cedar conversion.
- Deploy consistent policies across your organization with centralized control.



Policy in AgentCore



Keep agents in bounds
Act autonomously, but stay within boundaries and compliance

Instant and consistent
Policies evaluated in milliseconds, without slowing agents

Verifiably correct
Built on years of automated reasoning, using Cedar

AgentCore Hands-on



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Observability

Why it matters?

- 1 Agents follow non-deterministic control flows powered by model's reasoning
- 2 Root cause analysis in complex sequence of chained calls
- 3 Assessing system health holistically for operational performance
- 4 Detecting performance degradation

Do I need both?

Short answer: YES!

Observability: What is happening? – Reactive

Evaluation: How well is it happening? - Proactive

Pillars of Observability

Logs



"The input prompt and generated outputs by the user were..."

Metrics



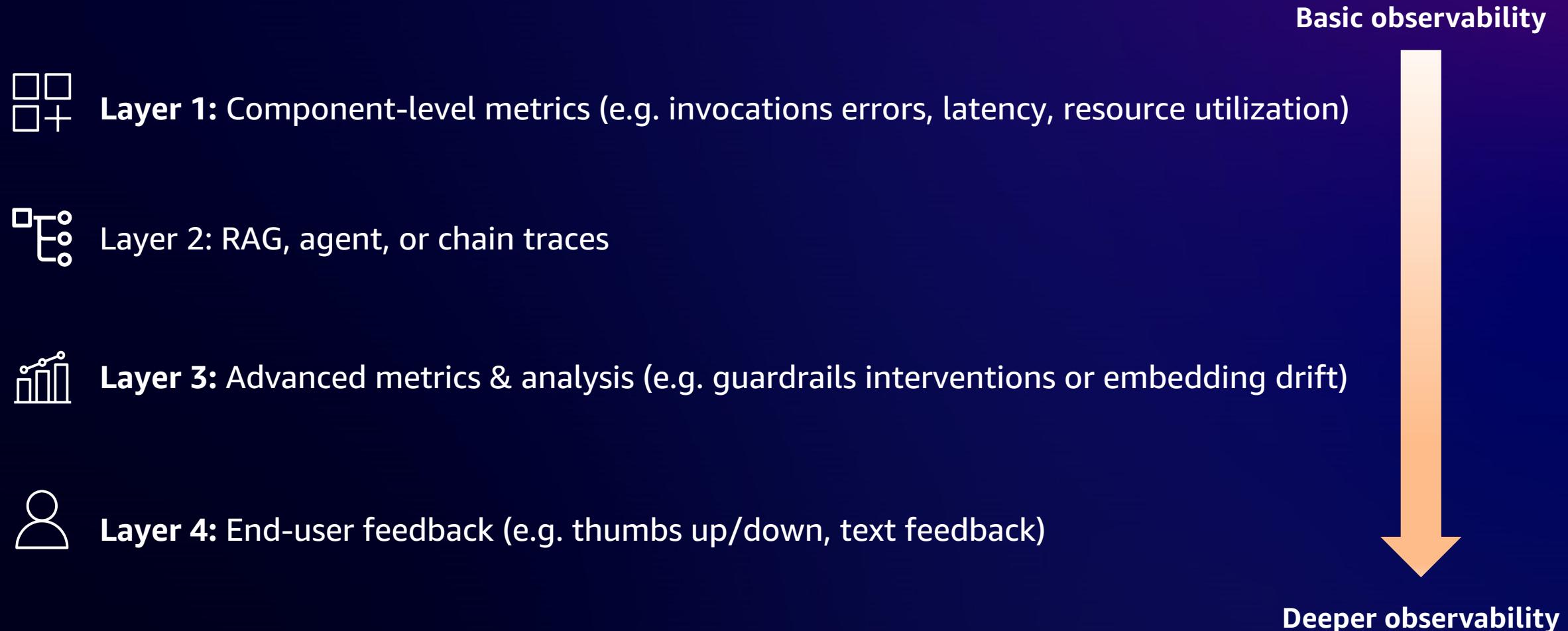
"The prompt to response latency has increased by 35%."

Traces



"What is the distribution of processing time across a RAG architecture?"

A layered approach to instrumenting observability

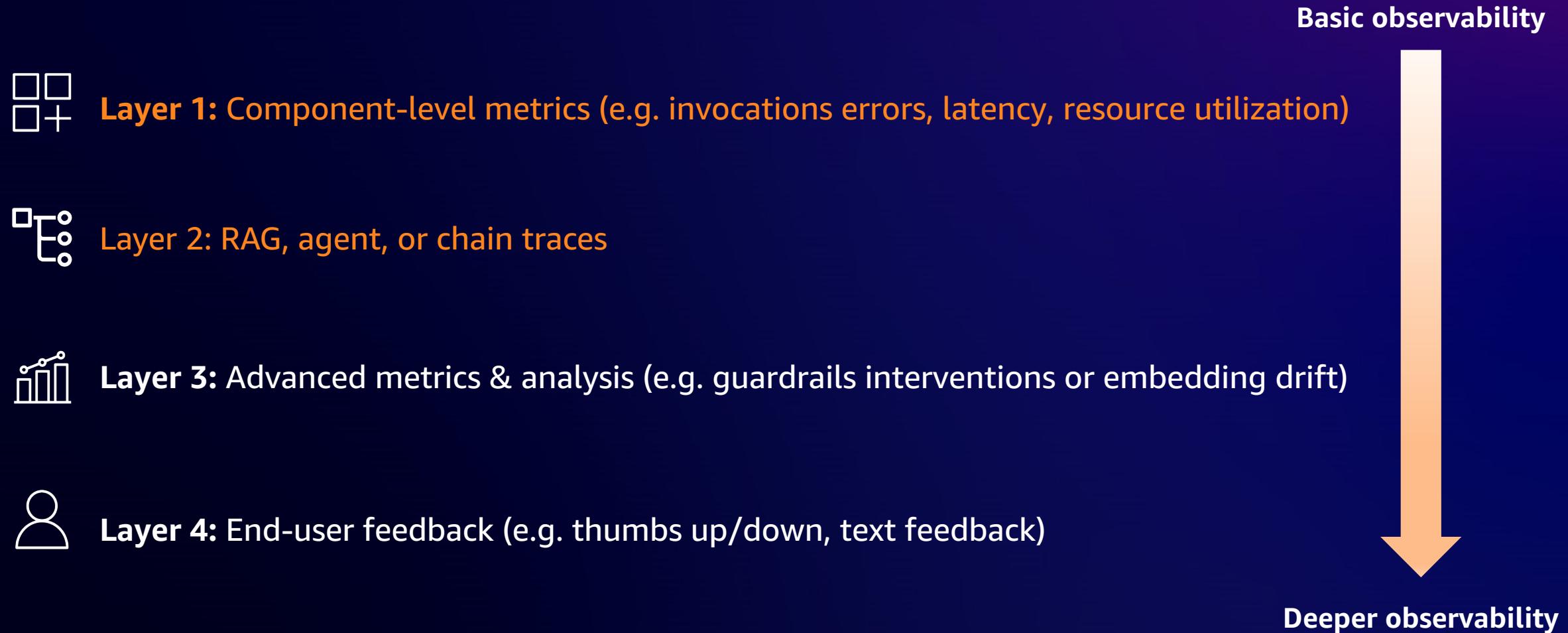


Best practices to consider for this layers

-  1/ Multi-system tracing: across multi-agents, MCPs and tools
-  2/ Session-level tracking: monitor entire interactions and task sequences
-  3/ Path analysis(Obs+Eval): using tracing and step-by-step execution sequences

+Continuous Improvement

A layered approach to instrumenting observability



Observability features in Amazon Bedrock

BUILDING BLOCKS FOR SCALABLE MONITORING



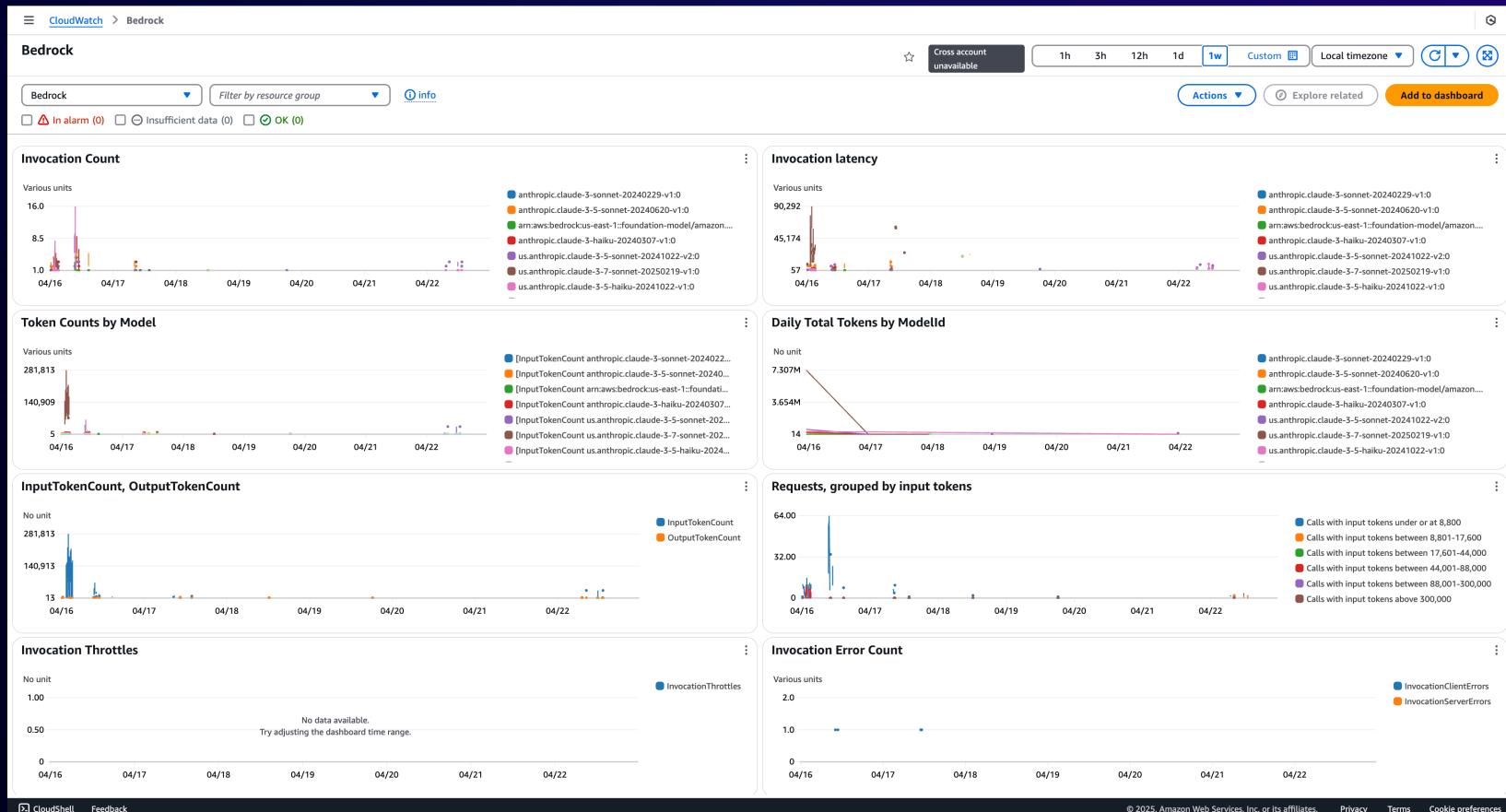
Metrics



Logs



Traces





AgentCore Observability

Maintain quality and trust



- Comprehensive end-to-end visibility into agent behavior
- Accelerated debugging and quality audits
- Quickly detect issues and assess performance trends

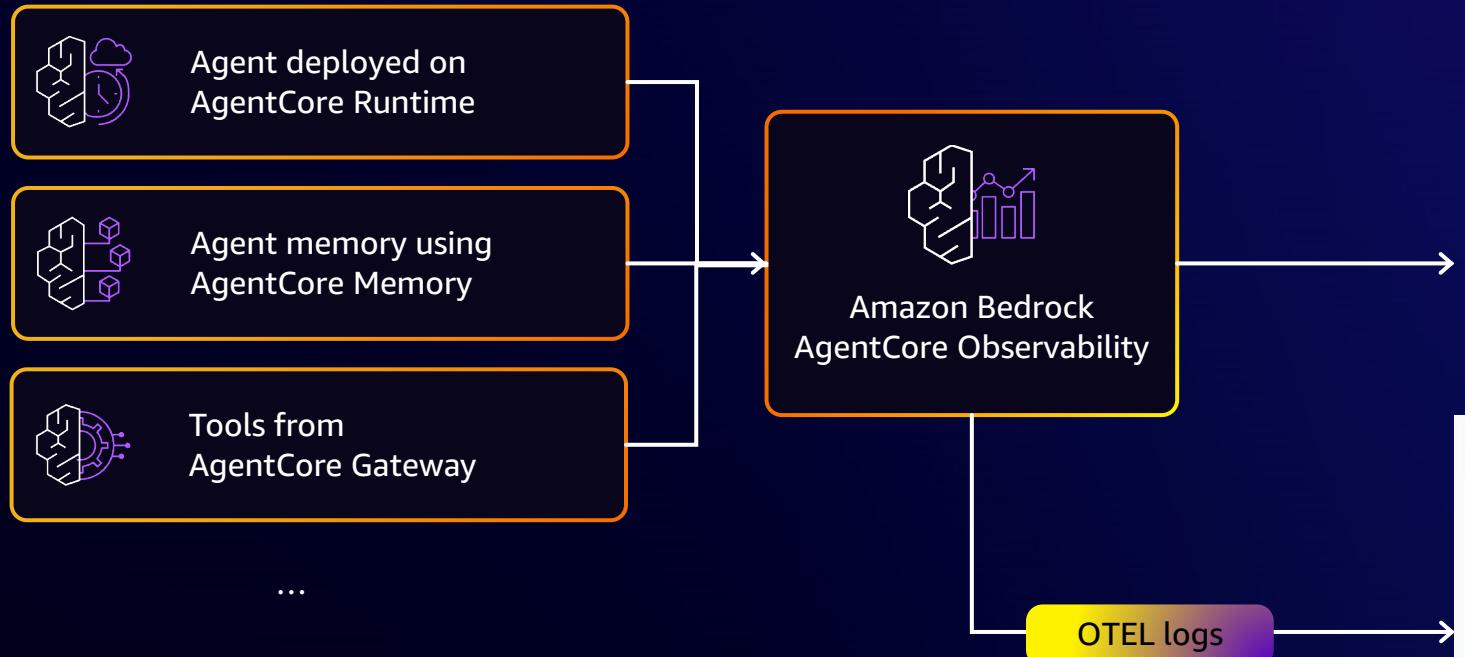
Integrate with 3P observability tools



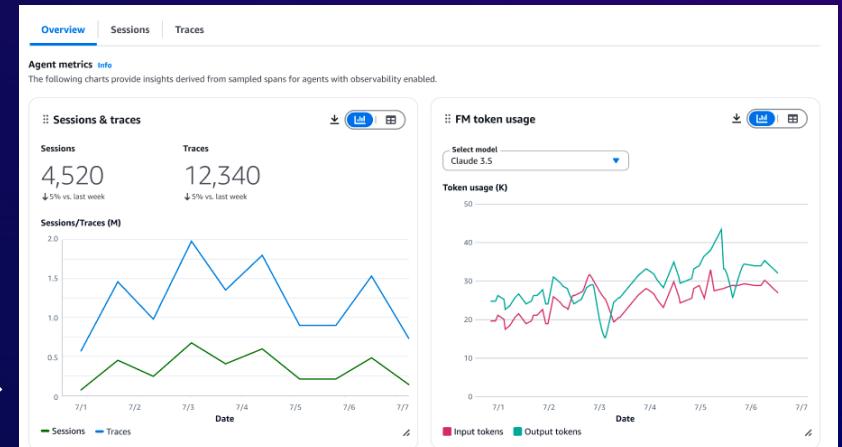
- Integration with a wide range of monitoring and observability tools, including CloudWatch
- Flexibility to leverage your existing observability stack



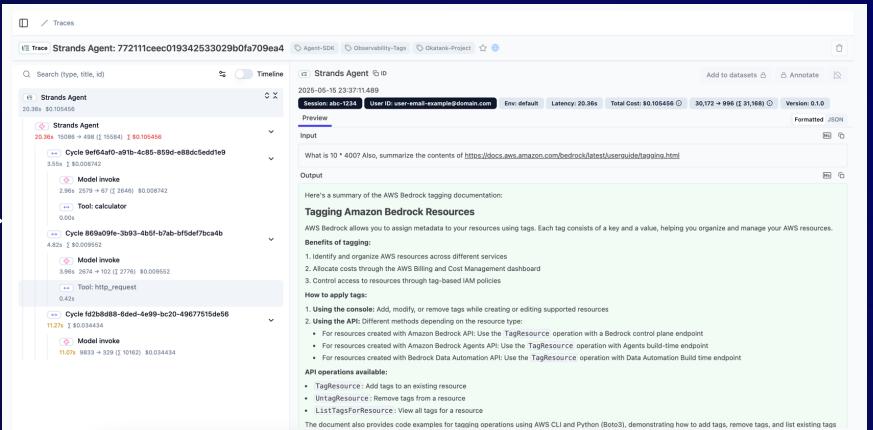
AgentCore Observability



AgentCore Observability dashboards

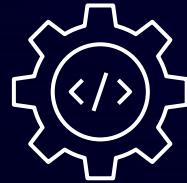


Third-party observability dashboards



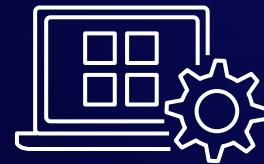
AgentCore Evaluations: Improve agent quality based on real-world performance

Real-time quality intelligence



- Sample and score live interactions with 13 built-in evaluators
- Gain actionable insights from real-world user behavior to continuously improve agent performance.
- On-demand and online evaluations

Custom business scoring



- Build tailored quality assessments for specific use cases using custom evaluators.
- See various levels of granularity (entire trajectory, final response, individual step).

No infrastructure overhead



- Out-of-the-box evaluation from logs available in Amazon CloudWatch.
- Eliminate months of effort required to build infrastructure and manage operational complexities.



AgentCore Evaluations

PREVIEW

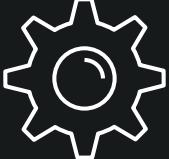


Two evaluation modes for different needs

Agent's evaluation metrics depends on the dimension we want to optimize for

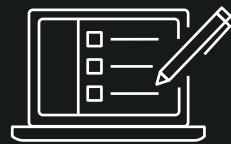
Online Evaluation

- Monitor live agent interactions
- Sample traces continuously
- Detect quality degradation



Use cases

- Catch silent failures
- Track quality trends
- Monitor after deployments
- Ensure ongoing reliability



On-Demand Evaluation

- Test changes before deployment
- Run evaluation tests in CI/CD
- Regression testing for builds
- Gate deployments on quality



Use cases

- Validate prompt changes
- Compare model performance
- Prevent quality regressions
- Automate quality gates



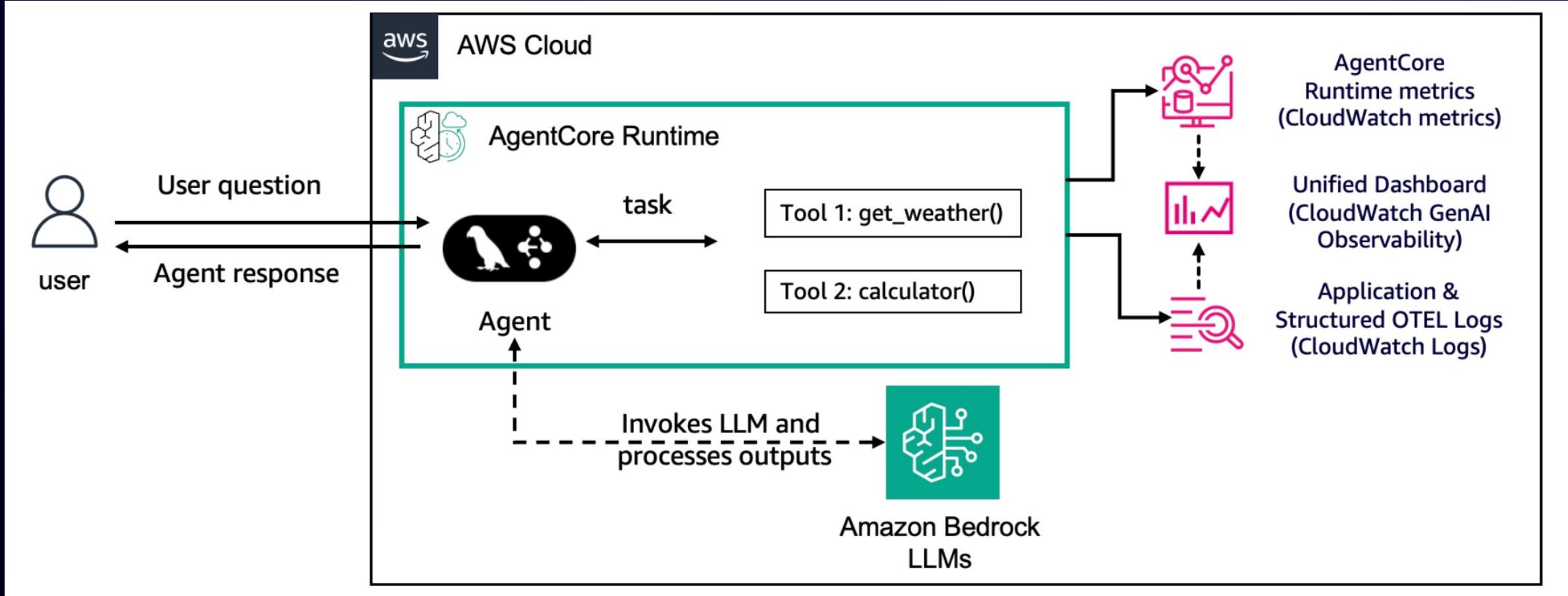
Production Monitoring

Development and CI/CD

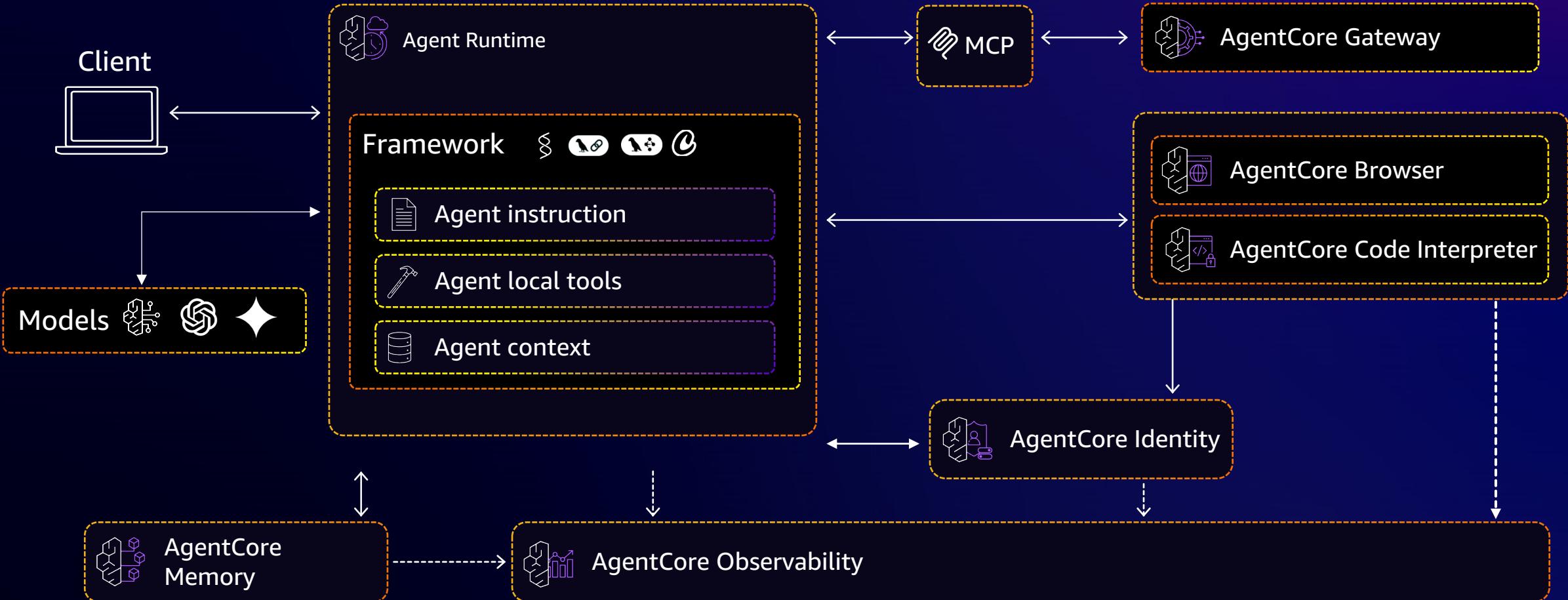


Architectures Diagram

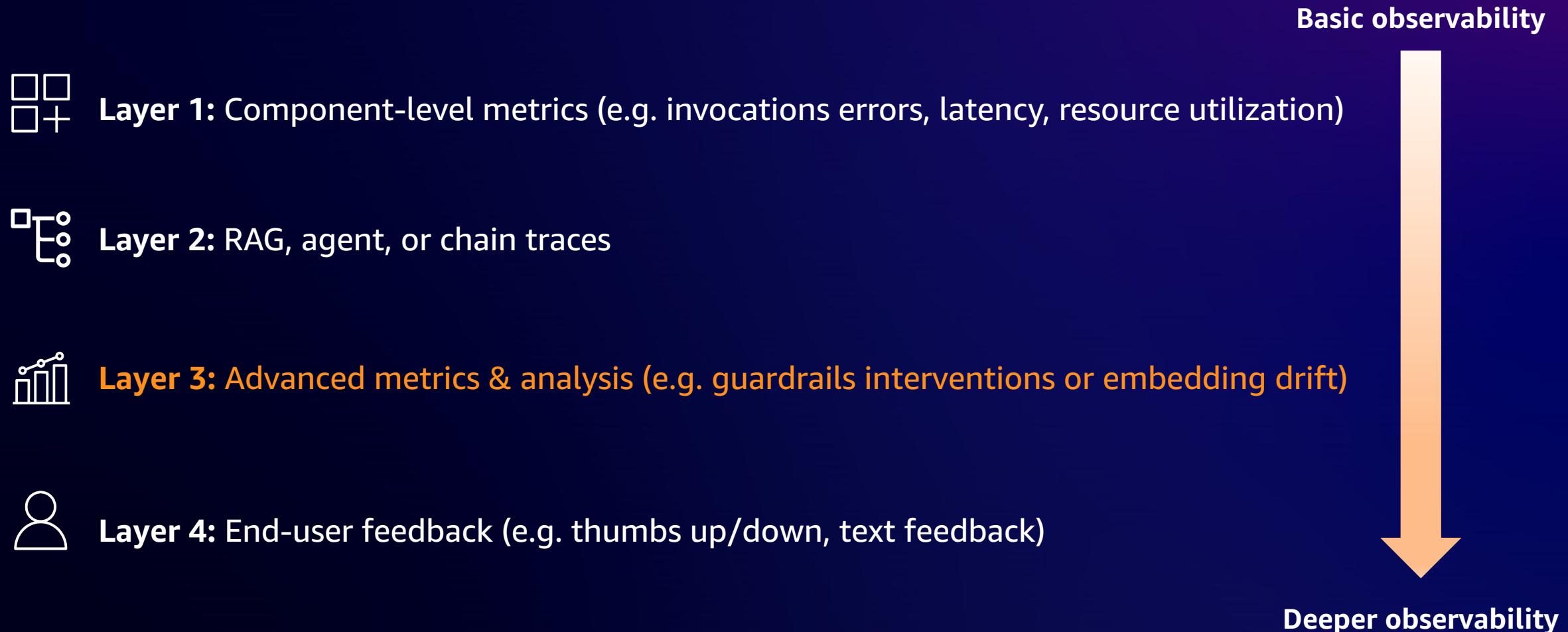
AgentCore Observability with agents deployed on AgentCore Runtime



AgentCore capabilities enabling agents at scale



A layered approach to instrumenting observability



Layer 3: Advanced metrics and analysis

The following table des-

Runtime metrics

Metric name
Invocations
InvocationLatency
InvocationClientErrors
InvocationServerErrors
InvocationThrottles
TextUnitCount
InvocationsIntervene

Dimension name	Dimension values	Available for the following metrics
Operation	ApplyGuardrail	<ul style="list-style-type: none">• Invocations• InvocationLatency• InvocationClientErrors• InvocationServerErrors• InvocationThrottles• InvocationsIntervened• TextUnitCount
GuardrailContentSource	<ul style="list-style-type: none">• Input• Output	<ul style="list-style-type: none">• Invocations• InvocationLatency• InvocationClientErrors• InvocationServerErrors• InvocationThrottles• InvocationsIntervened• TextUnitCount
GuardrailPolicyType	<ul style="list-style-type: none">• ContentPolicy• TopicPolicy• WordPolicy• SensitiveInformationPolicy• ContextualGroundingPolicy	<ul style="list-style-type: none">• InvocationsIntervened• TextUnitCount
GuardrailArn, GuardrailVersion	<ul style="list-style-type: none">• Guardrail Arn• Guardrail Version number or DRAFT	<ul style="list-style-type: none">• Invocations• InvocationLatency• InvocationClientErrors• InvocationServerErrors• InvocationThrottles• InvocationsIntervened• TextUnitCount

CloudWatch Metrics.

operation

errors

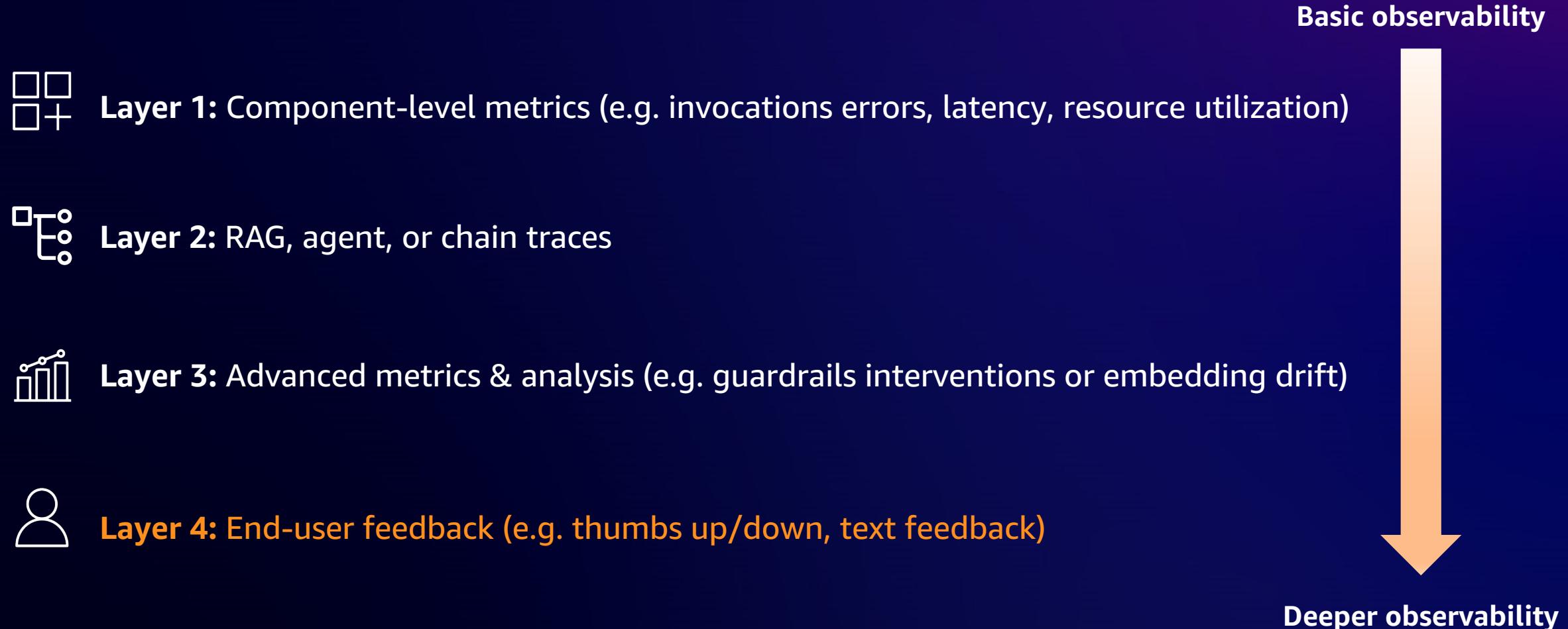
er-side errors

ed

policies

ertained

A layered approach to instrumenting observability

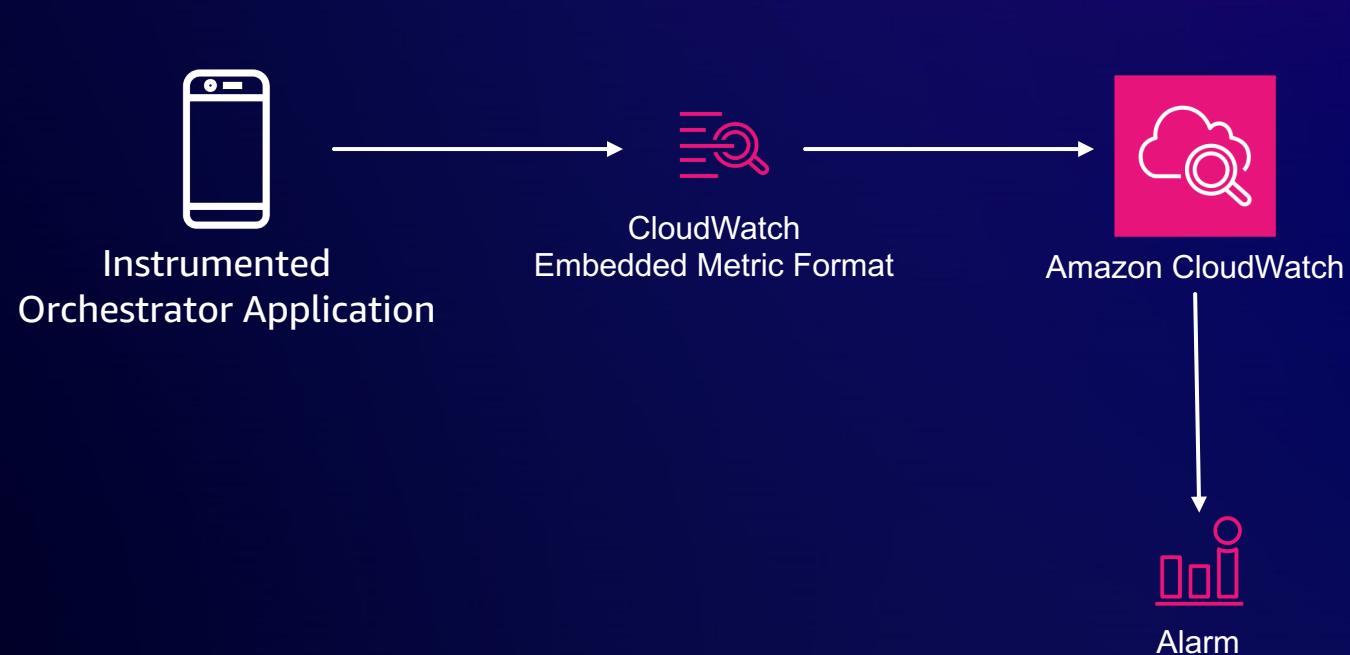


Layer 4: End-user feedback

PROVIDE EARLY INDICATIONS OF POTENTIAL ISSUES

We can leverage CloudWatch Embedded Metric Format (EMF) to instruct CloudWatch Logs to automatically extract metrics values.

This allows us to correlate a thumbs-up/thumbs-down response with written feedback and optionally alarm on negative user experience



CloudWatch Embedded Metric Format (EMF)

CORRELATE YOUR LOGS AND METRICS



CloudWatch Embedded Metric Format (EMF)

```
1 def submit_feedback():
2
3     customMetricJson = {
4         "_aws": {
5             "Timestamp": int(time.time())*1000,
6             "CloudWatchMetrics": [
7                 {
8                     "Namespace": "ChatBot-Sentiment",
9                     "Dimensions": [{"modelId"}],
10                    "Metrics": [
11                        {
12                            "Name": "sentiment",
13                            "Unit": "Count",
14                            "StorageResolution": 60
15                        }
16                    ]
17                }
18            ],
19        },
20        "modelId": model_id,
21        "sentiment": thumbFeedback,
22        "feedback": txtFeedback
23    }
```

```
1 cwLogEvent = json.dumps(customMetricJson)
2
3 logsClient = boto3.client('logs')
4 response = logsClient.put_log_events(
5     logGroupName='/chatbot/feedback',
6     logStreamName='default',
7     logEvents=[
8         {
9             'timestamp': int(time.time())*1000,
10            'message': cwLogEvent
11        }
12    ]
13 )
```



Amazon Bedrock OpenTelemetry Instrumentation



- Decorator implementation
- Can integrate with any observability provider
- 2500+ lines of Python code
- Single click CFN template for Langfuse
- Instrumentation for invoke and Converse API
- Amazon Bedrock Agents
- Amazon Bedrock knowledge Bases

Monitoring Bedrock Agents apps

GETTING VALUE FROM YOUR TRACES AND LOGS USING OPEN-SOURCE SOLUTIONS

Get started now, understanding agent behavior using Phoenix and/or LangFuse:

https://github.com/awslabs/amazon-bedrock-agent-samples/tree/main/examples/agent_observability

The screenshot shows the Langfuse interface for monitoring Amazon Bedrock Agents. At the top, it displays project statistics: Total Traces (6), Total Tokens (41,363), Latency P50 (~9.40s), and Latency P99 (~14.34s). Below this, there are tabs for Traces, Spans, and Sessions, with Traces selected. A sidebar on the left lists various monitoring categories: Tracing (selected), Sessions, Observations, Scores, Evaluation, Users, Prompts, Playground, and Datasets. The main area shows a table of traces with the following columns: ID, Timestamp, Name, User, Session, Latency, Usage, Scores, Total Cost, Observat..., Tags, and Action. The table contains four rows of trace data, each with a star icon and a unique ID.

ID	Timestamp	Name	User	Session	Latency	Usage	Scores	Total Cost	Observat...	Tags	Action
f4e1a7ebb...	2025-02-26 00:44:22	agent_invocation	Somename	default-session1	7.43s	0 → 0 (Σ 0)			10		⋮
0f383ac2e...	2025-02-26 00:43:43	agent_invocation	Somename	default-session1	19.11s	0 → 0 (Σ 0)			14		⋮
9bc20625...	2025-02-26 00:43:11	agent_invocation	Somename	default-session1	12.04s	0 → 0 (Σ 0)			14		⋮
6bd56441...	2025-02-26 00:42:45	agent_invocation	Somename	default-session1	5.71s	0 → 0 (Σ 0)			10		⋮

Integrations



Langfuse



arize



Weights & Biases



DATADOG



elastic

Discovery Questions for Your Teams

Current State

Are we defaulting to agents when a decision tree would suffice?

Scaling Strategy

What is our path from frontier models to cost-optimized inference at scale?

Engineering Investment

How much engineering effort is allocated to evaluation vs. feature development?

Safety Posture

Do we have deterministic guardrails on both input AND output?

Organizational Readiness

Is there executive ownership of AI safety and trust?

Red Flags and Anti-Patterns

Agent-First Thinking

"We will use an agent for everything." This ignores the decision hierarchy.

Eval-Later Mindset

"We will figure out testing once it is in production." Recipe for rework.

Single-Model Dependency

No plan for cost optimization at scale.

Safety as Afterthought

"We will add guardrails when we need them."

Summary

1. Use AgentCore observability as a starting point
2. Introduce advanced metrics such as guardrail interventions and embedding drift detection
3. Finally, incorporate end user feedback for quality assessment and malfunction alerts, completing the observability loop

AgentCore Hands-on



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Cost performance

Leverage the AWS Pricing Calculator

Choose a location type [Info](#)

Region

Choose a Region

US East (N. Virginia)

Amazon Bedrock
Select Amazon Bedrock Options that you want to estimate

Nova Micro Nova Lite Nova Pro
 Nova Canvas Nova Reel Titan Text Embeddings V2

On-Demand [Info](#)

Number of Input tokens Unit

Number of output tokens Unit

[▶ Show calculations](#)

Provisioned Throughput [Info](#)

Is the model customized?

Billing term
Select the billing term commit

Number of model units purchased per month

Amazon Bedrock offers several features to balance cost, latency and accuracy

Amazon Bedrock Model Distillation

Utilize smaller, more cost-effective models

Maximize distilled model performance with proprietary data synthesis

Distilled models are up to 500% faster and up to **75% less expensive** than original models, with **less than 2% accuracy loss** for use cases like RAG

Amazon Bedrock Intelligent Prompt Routing

Automatically route prompts to different foundation models to optimize response quality and lower costs

Provides a single endpoint to efficiently route prompts

Uses advanced prompt matching to meet cost and latency thresholds

Reduce costs by **up to 30%** without compromising on accuracy

Amazon Bedrock support for prompt caching

Cache repetitive context in prompts across multiple API calls

Securely cache entire prompts

Enhance accuracy through longer, detailed prompts

Reduce costs by **up to 90%** and latency **by up to 85%** for supported models

Amazon Bedrock inference consumption options

ON DEMAND

Pay-as-you-go,
no commitment

Pricing based on input and output token
count for LLMs

Great for prototyping, POCs, and small
workloads with more relaxed requirements
for throughput and latency

Requests per minute (RPM) and tokens per
minute (TPM) limits enforced

Includes support for cross-region inference

PROVISIONED THROUGHPUT

Provision sufficient throughput to
meet your application's
performance requirements

Reserve throughput (input/output tokens
per minute) at a fixed cost

Flexible commitment term of 1 month or
6 months

Hourly PT (2MU), no commit, across
different models

Pay hourly rate, discounted for extended
commit

Great for production workloads or
inference on custom models

BATCH INFERENCE

Use Amazon Bedrock to
process prompts in batch to
get responses for model
evaluation, experimentation,
and offline processing.

Amazon Bedrock offers select foundation
models (FMs) from leading AI providers like
Anthropic, Meta, Mistral AI, and Amazon for
batch inference at 50% of on-demand
inference pricing

Efficiently run model inference on large
volumes of data

Avoids throttling when running
large jobs



KIRO



**What would a development
experience look like if it could take
full advantage of working with AI
agents?**

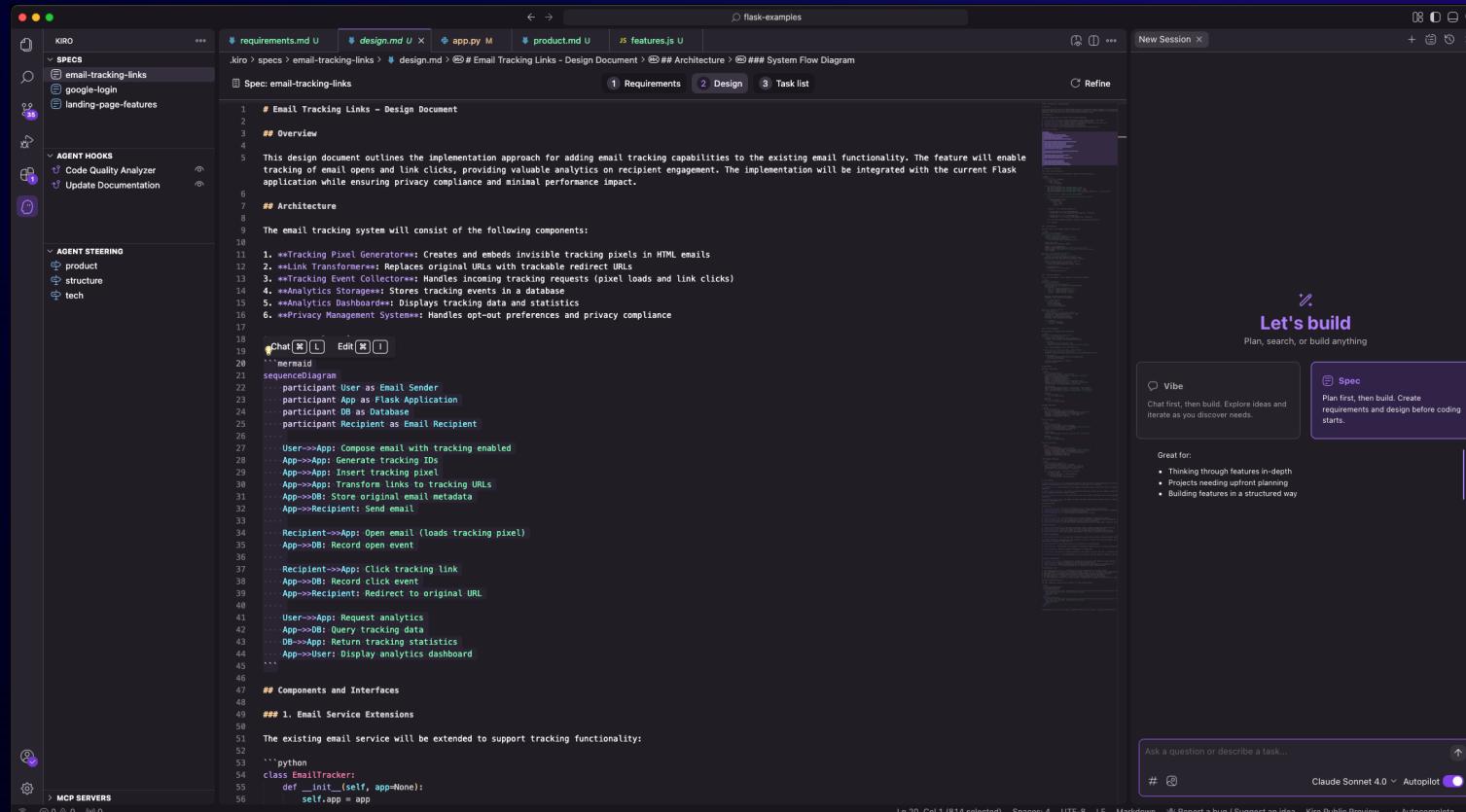
Solution: Spec-driven development

Vibe coding ships the prototype,
traditional SDLC practices ship the product



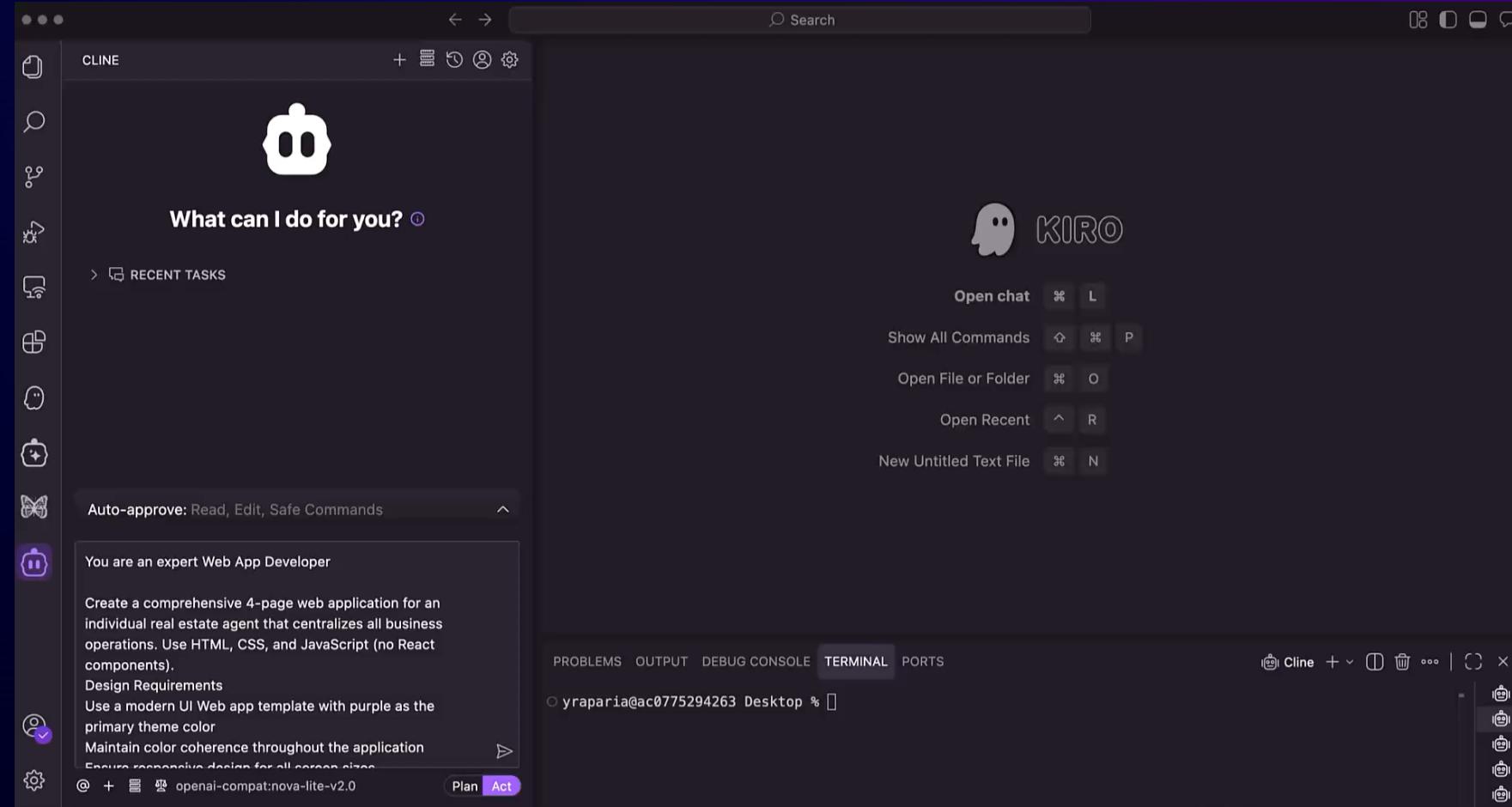
Kiro helps developers and engineering teams
ship high quality software with AI agents

The AI IDE for prototype to production



Kiro helps you do your best work by bringing structure to AI coding with spec-driven development.

Watch Nova 2 create a functioning web app in one shot using Cline extension in Kiro



SageMaker AI - New capabilities, at a glance

Serverless RFT and SFT

*Advanced tuning techniques
without configuring infrastructure*

Flexible deployment

*1-click deployment to SageMaker or Bedrock,
and exportable model weights saved to S3*

Synthetic data generation

*Agent-guided SDG for save
time sourcing training data*

Deep evaluation

*LLM-as-a-Judge, custom scorer, and public
benchmarks, including built-in visualizations*

Serverless MLFlow

*Advanced experiment tracking,
fully integrated and serverless*

Optimized hyperparams

*Built on pre-optimized recipes based on
pairing of model, technique, and instance*

Closing remarks

Thank you!



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.