aws

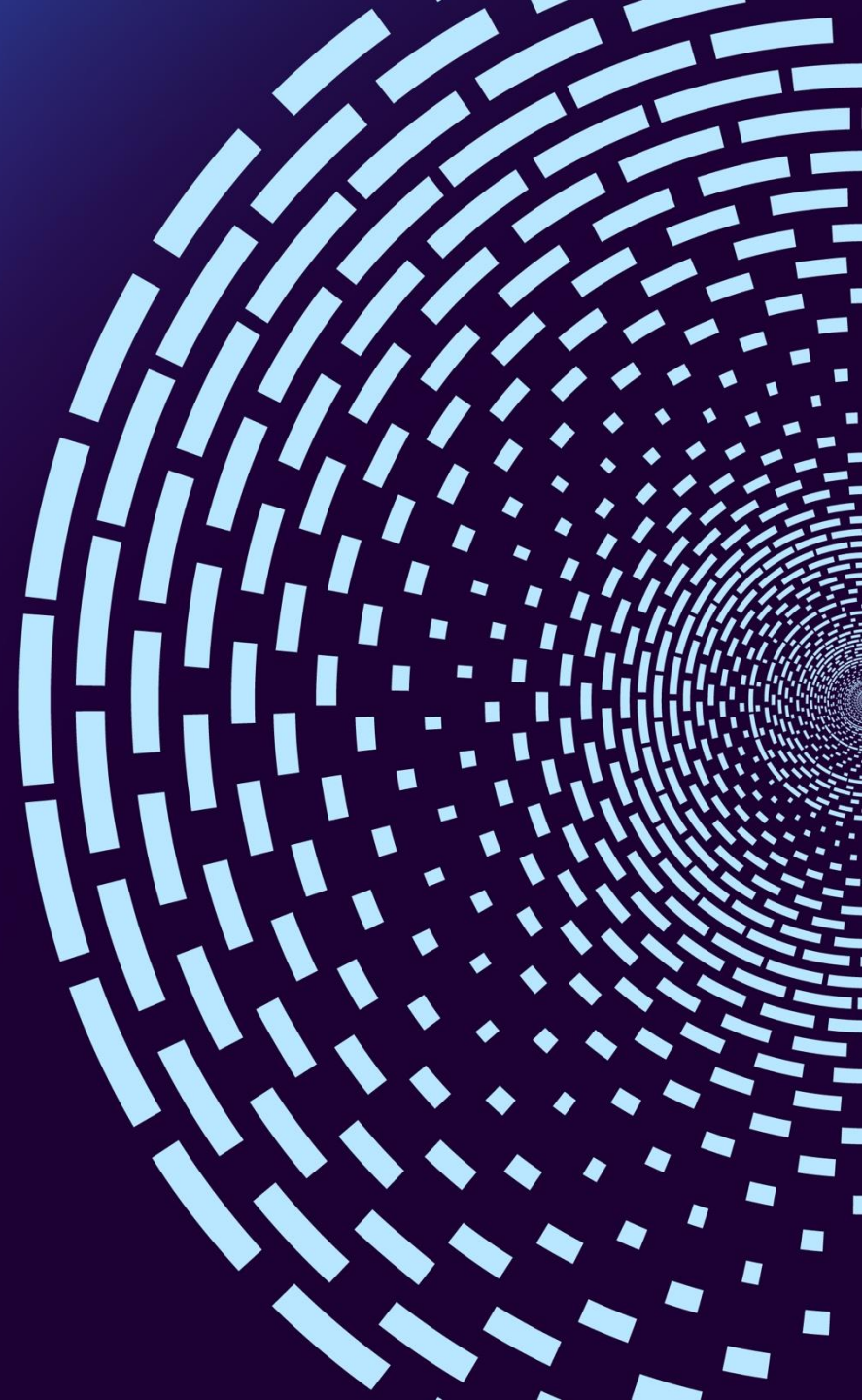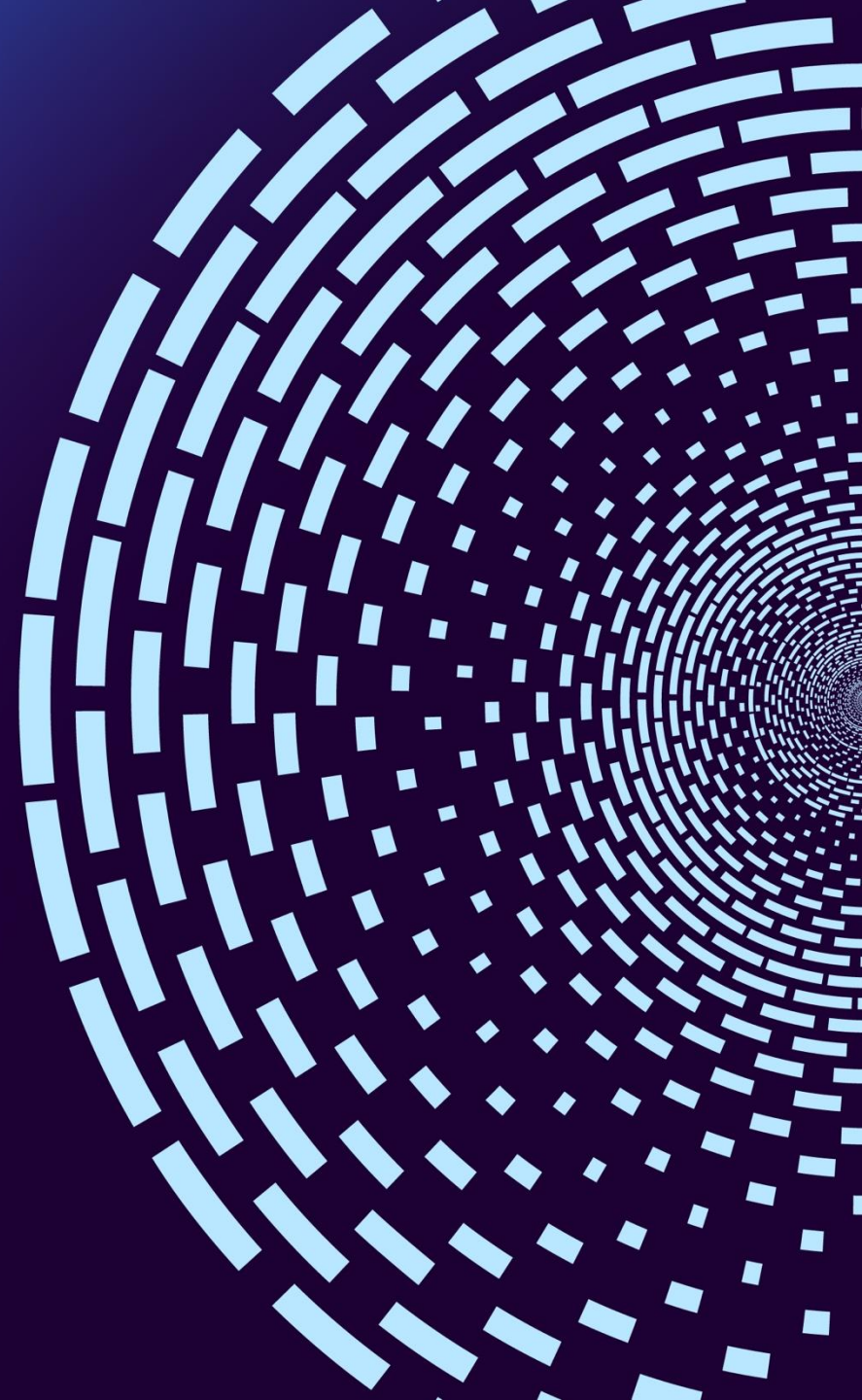# AI Conclave

Online

aws

AIOT101

# Amazon Bedrock state of the union

**Aparajithan Vaidyanathan**

Principal Solutions Architect
AWS India

# Agenda

- **Introduction**

- **Choose the best model**

- **Customize with your data**

- **Apply safety and responsible AI checks**

- **Build and orchestrate agents**

- **Optimize for cost, latency and accuracy**

# Introduction

# Amazon Bedrock

The easiest way to build and scale generative AI applications

Choice of leading FMs through a single API

Model customization

Retrieval Augmented Generation (RAG)

Agents that execute multistep tasks

Security, privacy, and data governance

# Amazon Bedrock | re:Invent launches
## THE EASIEST AND FASTEST WAY TO BUILD AND SCALE GENERATIVE AI APPLICATIONS

**CHOOSE THE BEST MODEL**

amazon

Luma

poolside

Marketplace
(GA)

Model evaluation
LLM-as-a-judge
(preview)

**CUSTOMIZE WITH YOUR DATA**

Data Automation (preview)

Knowledge Bases supports GraphRAG
(preview)

Knowledge Bases supports structured
data retrieval (GA)

Knowledge Bases supports
multimodal data processing (GA)

Knowledge Bases supports
streaming responses (GA)

Knowledge Bases supports real-
time sync from custom data
sources (GA)

Rerank API (GA)

Knowledge Bases supports RAG
evaluation (preview)

Knowledge Bases provides auto-
generated query filters for
retrieval (GA)

**APPLY SAFETY AND
RESPONSBILE AI CHECKS**

Guardrails support
Automated
Reasoning check
(preview)

Guardrails supports
multimodal
toxicity detection
(preview)

**BUILD AND ORCHESTRATE
AGENTS**

Multi-agent collaboration
(preview)

**OPTIMIZE FOR COST, LATENCY,
AND ACCURACY**

Latency-optimized Inference
Options (preview)

Prompt caching
(preview)

Intelligent Prompt Routing
(preview)

Model distillation
(preview)

AWS security, privacy and reliability built-in

# Choose the best model

# Amazon Bedrock

| AI21 labs | amazon | ANTHROP\C | cohere | Luma | Meta | MISTRAL AI_ | poolside | stability.ai |
|---|---|---|---|---|---|---|---|---|
| Effective reasoning & rapid analysis for long context windows | Frontier multimodal intelligence at low-latency, Agent & RAG Applications, high-quality image & video generation | Advanced reasoning & coding capabilities, including computer use skills | Multimodal search & advanced retrieval powering multilingual knowledge agents | High-quality video generation from text & images | Advanced image & language reasoning | Knowledge summarization, expert agents, & code completion | Software engineering AI for large enterprises | High-quality AI image generation, easily deployable at scale |
| JAMBA | AMAZON NOVA<br>**New** | CLAUDE | COMMAND<br>EMBED<br>RERANK **New** | LUMA RAY 2<br>**Coming soon** | LLAMA | MISTRAL<br>MIXTRAL | MALIBU<br>POINT<br>**Coming soon** | STABLE DIFFUSION<br>STABLE IMAGE |

# Amazon Nova foundation models

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

Understanding models

Creative content generation models

## Amazon Nova Micro

Our text only model that delivers the lowest latency responses at very low cost

**GENERALLY AVAILABLE**

## Amazon Nova Lite

Our lowest cost multimodal model that is lightning fast for lightweight tasks

**GENERALLY AVAILABLE**

## Amazon Nova Pro

Our highly capable multimodal model with best combination of accuracy, speed, and cost for a wide range of tasks

**GENERALLY AVAILABLE**

## Amazon Nova Premier

Our most capable multimodal model for complex reasoning tasks and for use as the best teacher for distilling custom models

**COMING SOON**

## Amazon Nova Canvas

State-of-the-art image generation model

**GENERALLY AVAILABLE**

## Amazon Nova Reel

State-of-the-art video generation model

**GENERALLY AVAILABLE**

**Lower cost & latency** ←———————————→ **Increasing intelligence**

# Amazon Nova - Understanding models

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

# Amazon Nova creative content generation models

| | Amazon Nova Reel | Amazon Nova Canvas |
|---|---|---|
| Availability | GA | GA |
| Input characters | 512 | 1024 |
| Languages | EN | EN |
| Modalities supported | Text, image input; video output | Text, image input; image output |
| Duration | 6 seconds (2 minutes coming soon) | N/A |
| Fine-tuning | Coming soon | Coming soon |

# 3rd party penchmarks

## Amazon Nova Lite



Source: artificialanalysis.ai (as of 12/9/24)

# Amazon Bedrock
# Marketplace
## (GA)

Discover and use over 100 popular, emerging and specialized models in Amazon Bedrock

- Streamline development workflows with a unified console experience

- Deploy models on managed endpoints with custom scaling policies

- Leverage Amazon Bedrock's APIs, tools, and security

# Amazon Bedrock
# LLM-as-a-Judge
## (Preview)

Get human-like evaluation quality at a much lower cost than full human-based evaluations, while saving weeks of time

- Ensure you have the right combination of evaluator models and models being evaluated

- Use curated metrics to evaluate objective facts or subjective evaluation of writing style and tone on your dataset

- Compare results across evaluation jobs to make decisions faster

# How does LLM-as-a-judge work?

**Example input**

**Judge prompt (simplified)**

You are a helpful assistant…

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

prompt: What is the capital of Spain?

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

referenceResponse: Madrid

Model response: Barcelona

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

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

Note: Correctness can be with or without ground truth

# Customize with your data

# Amazon Bedrock Knowledge Bases

re:Invent 2024 launches

Structured data retrieval

Auto-generated query filters

Multi-modal data processing

Rerank API

GraphRAG

Real-time sync for custom data sources

Streaming responses

RAG evaluation

# Amazon Bedrock Knowledge Bases structured data retrieval (GA)

Seamlessly integrate structured data for RAG

- Use data stored in Amazon SageMaker Lakehouse, Amazon Redshift and Amazon S3 Tables

- Reduce application development time from months to days

- Improve the accuracy of your queries with customization context

**Amazon Bedrock Knowledge Bases**

# GraphRAG

**(preview)**

Generate more relevant responses for RAG applications using knowledge graphs

- Generate knowledge graphs to link relationships across data sources

- Build more comprehensive, explainable generative AI applications

- Enhance transparency of source information for better fact verification

# What is the name of the capital city of the country where Eiffel Tower is located?

Requires **Multi-step reasoning**

| Identifying the location of the Eiffel Tower | → | Determining the country of that location | → | Finding the capital city of that country |
|---|---|---|---|---|

## Basic RAG

It doesn't inherently understand relationships between entities (Eiffel Tower, Paris, France)

Basic RAG might find information about the Eiffel Tower or about France's capital separately, but may not connect these pieces of information effectively.

## GraphRAG

It understands the relationships: Eiffel Tower is in Paris, Paris is the capital of France

GraphRAG can establish connections between multiple entities: Eiffel Tower, Paris, France

Recognizes the hierarchy: Landmark -> City -> Country -> Capital City

**Amazon Bedrock Knowledge Bases**

**RAG evaluation**

**(preview)**

Evaluate end-to-end RAG workflow

- Get actionable insights to improve your RAG system

- Ensure the generated content is correct, complete, limits hallucinations, and adheres to responsible AI principles

- Accelerate time to value for deploying RAG applications

# RAG evaluation:



Choose evaluator model

Choose your knowledge base

Choose to evaluate retrieval only or retrieve and generate

Choose your generator model

Choose your metrics

Upload your prompt dataset

Inference and evaluation

View results

**Amazon Bedrock Knowledge Bases**

**streaming responses (GA)**

Introduce a new RetrieveAndGenerateStream API to respond fast

- Reduce latency and improve user experience

- Maintain the same level of accuracy

- Adhere to all relevant data privacy and security regulations, including GDPR, HIPAA, and FedRAMP

**Amazon Bedrock Knowledge Bases**

# multimodal data processing (GA)

Analyze and leverage insights from both textual and image data

- Get a fully-managed RAG workflow for visually-rich data

- Retrieve and generate answers to questions derived from text and visual data

- Improve accuracy, relevancy, and depth of responses

# Multimodal data processing



Documents

→ Images

→ Text

→ Tables

**A M A Z O N
B E D R O C K
K N O W L E D G E
B A S E S**

Choose either
Amazon Bedrock Data
Automation or foundation
model for parsing

Ingest documents
with images, tables,
and text

Query your
knowledge base

Get generated
response back with
images as source
attribution

# Amazon Bedrock Knowledge Bases real-time sync from custom data sources (GA)

Custom connector API to directly ingest content and manage individual documents easily

- Enables ingestion of content and metadata from different sources in real-time

- Enables selective updates and deletion

# Amazon Bedrock Knowledge Bases

# auto-generated query filters for retrieval (GA)

Automatically generate the query filter expressions based on the metadata schema

- Improves user experience as the retrieval filter generator automatically generates the appropriate filter expressions based on user inputs

- Support metadata embedding to ensure metadata is considered in RAG process

- Highly configurable, allowing users to select metadata fields used for filtering, and allowing them to customize prompts
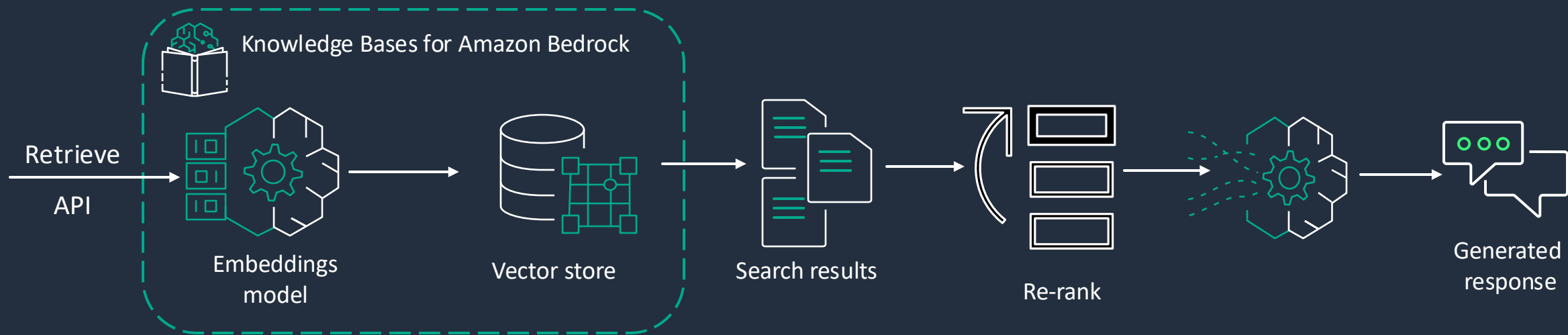
**Amazon Bedrock**

# Rerank API
## (GA)

Rerank the retrieved document chunks based on their relevance to the query

- Improve accuracy of documents by prioritizing the most important content to be passed to generation models

- Invoke the reRank API without any additional model deployment or code
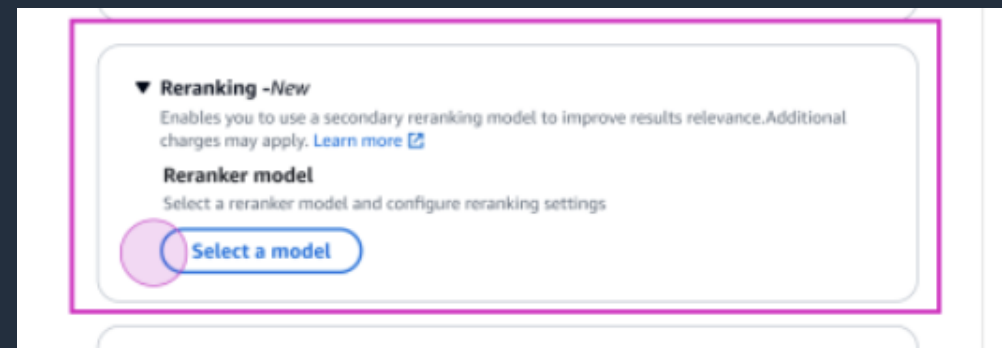
- Provide model flexibility

# Re-ranking

- Model selection – Amazon rerank, Cohere 3.5 rerank

- Accessible independently as well through Rerank API

# Amazon Bedrock
# Data Automation
## (preview)

Transform unstructured multimodal data for generative AI applications and analytics

- Extract, transform, and generate structured data from multimodal content

- Generate customized outputs based on your requirements and business rules

- Streamline application workflows with a fully managed, single API experience

# Apply safety and responsible AI checks

**Amazon Bedrock Guardrails**

# Automated Reasoning checks

## (preview)

Prevent factual errors due to hallucinations

- Verify accuracy of model responses using mathematical proof

- Provides recommendations for correcting factual errors

- Enhance reliability of LLM responses for critical use cases

# Automated reasoning checks – NEW!

## DETECT FACTUAL ERRORS FROM HALLUCINATION WITH LOGICALLY ACCURATE AND VERIFIABLE REASONING

- **Accurate** Identifies and suggests corrections for inaccurate factual claims on supported knowledge

- **Sound** Uses formal logical reasoning to correctly determine accuracy

- **Transparent** Explains why a claim is accurate or not

## Amazon Bedrock Guardrails
# Multimodal toxicity detection
## (preview)

Configurable safeguards for image content

- Enhance security of multimodal generative AI applications

- Available for all foundation models in Amazon Bedrock with image support

- Enable consistent policy control

# Content Filters – NEW: Image support!

## CONFIGURE THRESHOLDS TO FILTER HARMFUL CONTENT

Filter harmful content across categories

- Hate

- Insults

- Sexual

- Violence

- Misconduct *

- Prompt attack *

* Text only

# Build and orchestrate agents

# Amazon Bedrock
# Multi-agent collaboration (preview)

Easily build, deploy and orchestrate teams of agents that work together to handle complex, multi-step tasks

- Accelerate tasks with agents working in parallel

- Effortlessly orchestrate agents without complex coding

# 1. Unify customer experience

1. Unify customer experience

# 2. Automate complex processes with supervisor agents . . .

Give me a complete marketing strategy for my new product XYZ . . .

Here's the XYZ 10-page marketing strategy, including . . .

Orchestrate multi-step plan across agents

**1**

## Generate plan dynamically

1. Conduct thorough market research, include competitors
2. Develop detailed project summary, target persona
3. Formulate comprehensive marketing strategy with goals, tactics, channels, KPIs
4. Create three innovative marketing campaign ideas
5. Develop detailed marketing copy for each campaign, including a video ad script, and a draft video
6. Produce final report and save interim results

## Sub-agents and knowledge bases

Marketing strategist

Content creator

Style guide

Campaign feedback

Market analyst

Copy editor

Project storage agent

Video generation expert

**2** Execute plan

# Inline agents

Configure your agent dynamically at runtime

Quickly experiment with different agent's configurations

Enable subscription-based personalization for different customers

Select persona-based data sources at runtime

Dynamically select actions provided to your agent

Incorporate dynamic identity/behavior to your agent

# Optimize for cost, latency and accuracy

# Internal testing of latency-optimized inference

## 3.5 Haiku
- Output tokens per second, 67 -> 152
- Time to first token, 1.1s -> 0.6s

## Llama 3.1 70b
- Output tokens per second, 32 -> 203
- Time to first token, 0.9s -> 0.4s

| Model | Inference Profile | TTFT P50 | TTFT P90 | OTPS P50 | OTPS P90 |
|---|---|---|---|---|---|
| us.anthropic.claude-3-5-haiku-20241022-v1:0 | Optimized | 0.6 | 1.4 | 85.9 | 152.0 |
| us.anthropic.claude-3-5-haiku-20241022-v1:0 | Standard | 1.1 | 2.9 | 48.4 | 67.4 |
| comparison | | -42.20% | -51.70% | 77.34% | 125.50% |
| us.meta.llama3-1-70b-instruct-v1:0 | Optimized | 0.4 | 1.2 | 137.0 | 203.7 |
| us.meta.llama3-1-70b-instruct-v1:0 | Standard | 0.9 | 42.8 | 30.2 | 32.4 |
| comparison | | -51.65% | -97.10% | 353.84% | 529.33% |

# Amazon Bedrock supports prompt caching (preview)

Cache repetitive context in prompts across multiple API calls

- Securely cache entire prompts

- Enhance accuracy through longer prompts

- Reduce cost by up to 90% and latency by up to 85% for supported models

**Amazon Bedrock**

# Intelligent prompt routing

## (preview)

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

- Provides a single endpoint to efficiently route prompts

- Meets cost and latency thresholds with advanced prompt matching techniques

- Reduces application development costs by up to 30%

# Amazon Bedrock
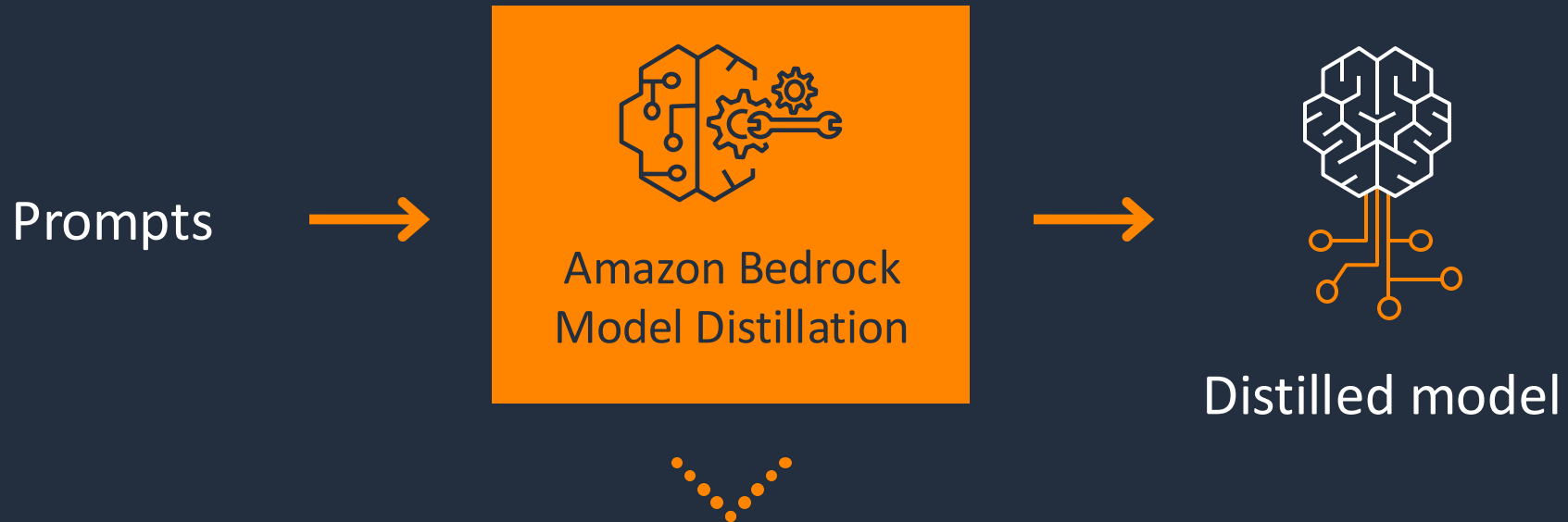# Model Distillation
## (preview)

Create smaller, faster, more cost-effective models

- Easily transfer knowledge from a large, complex model to a smaller one

- Distilled models up to 500% faster and up to 75% less expensive

- Anthropic, Meta, and Amazon models

# Behind the scenes proprietary data synthesis

Prompts →

**Amazon Bedrock Model Distillation**

→ Distilled model

Response generation from teacher model **+** Data synthesis (e.g., generating similar prompts, using golden responses as examples) **+** Training of student model
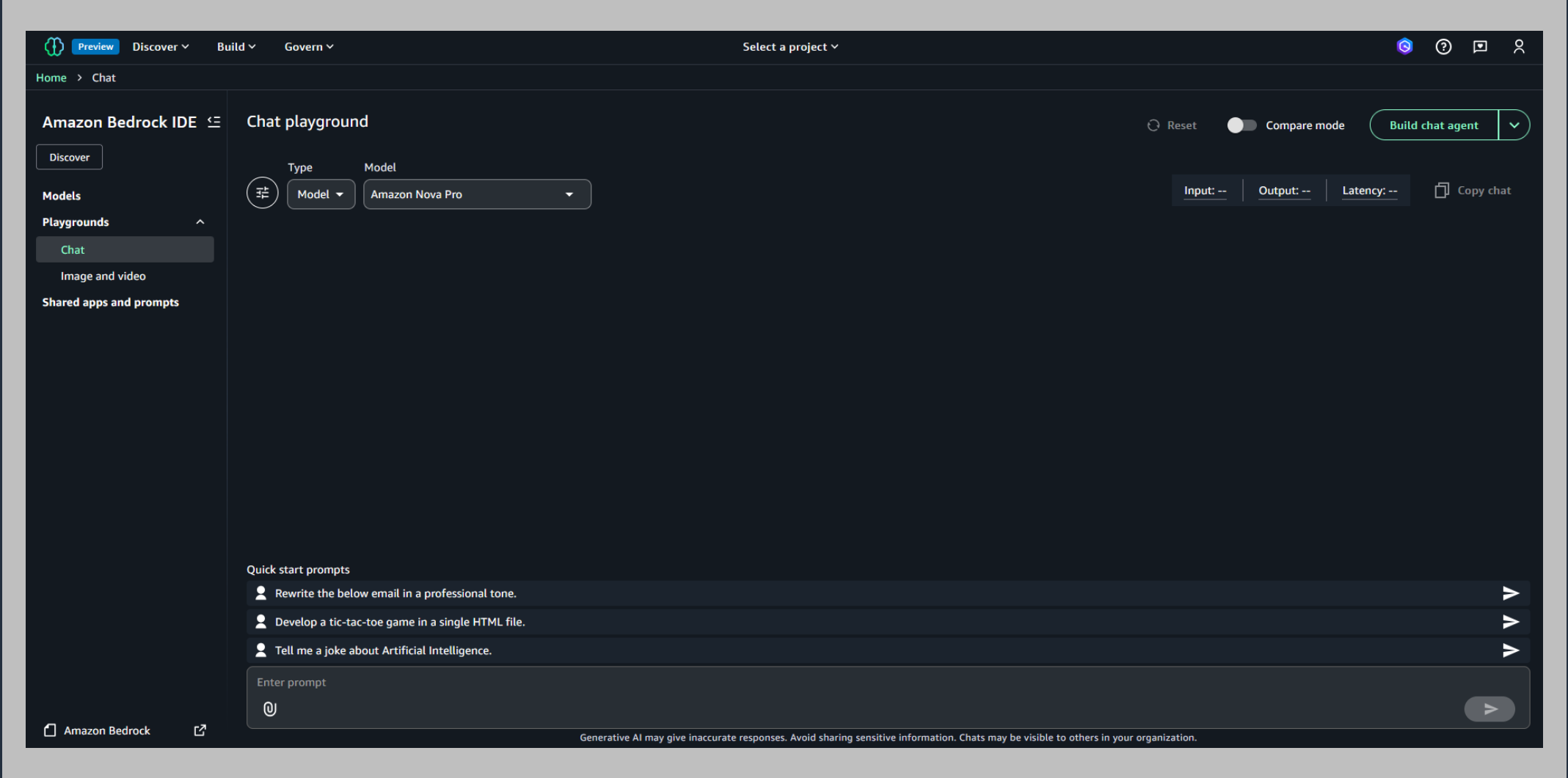
# Amazon Bedrock

# IDE
## (preview)

Enable effortless generative AI development in a governed collaborative environment

- Facilitates custom AI application building with advanced features

- Promotes seamless collaboration among stakeholders

- Simplifies model evaluation and adoption through Playground experience

# Amazon Bedrock IDE

aws

# Thank you!

**Aparajithan Vaidyanathan**

Principal Solutions Architect
AWS India