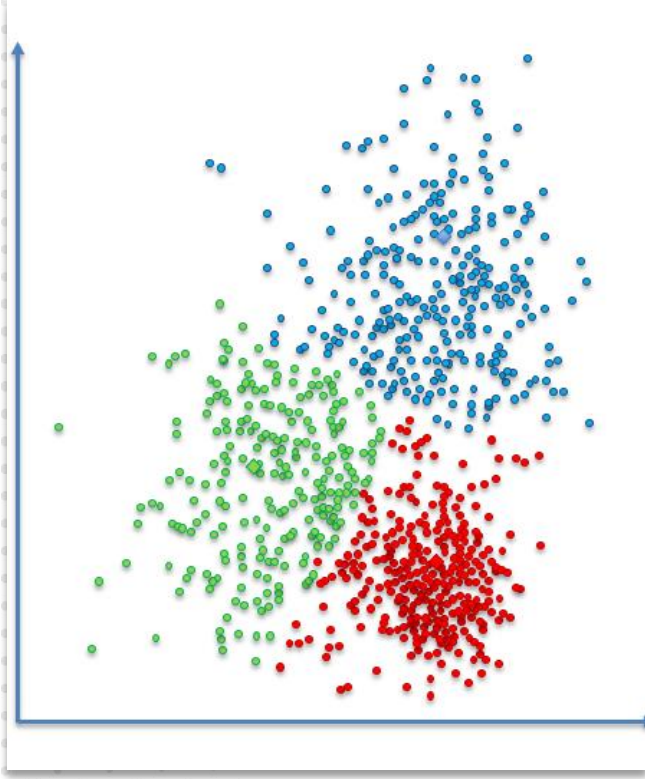


Lesson #5:

Cloud-Based LLM Development Platforms

Objectives:

- Overview of Cloud Platforms
- Enter AWS Bedrock
- Foundation Model Providers and Access
- Anatomy of a Bedrock Call
- Going Deeper



Cloud Development Platforms

Models and Transformers

Llama

OpenAI,
Claude, etc

DeepSeek

Development Environments

Colab

VS Code

Cursor /
Windsurf

Local LLM Management Platforms

Ollama

LM Studio

Jan

RAG & Agentic Frameworks

LlamaIndex

LangGraph

Haystack

Agent Orchestration Platforms

LangGraph

AutoGen

CrewAI

Agentic Protocols

APIs

MCP

A2A

Cloud Development Platforms

AWS Bedrock

Azure
AI Studio

Google
Vertex AI

Monitoring and Observability

LangSmith

Weights &
Biases

Helicone

Cloud LLM Platforms

- AWS, Google, and Azure support LLM model development in a cloud setting.
- They do the following:
 - Provide Access to Foundation Models
 - Offer Serverless Inference APIs
 - Support Model Customization
 - Enable Retrieval-Augmented Generation (RAG)
 - Provide Agentic & Tool-Using Frameworks
 - Support Evaluation, Monitoring & Guardrails

The Big Players

- AWS Bedrock
 - Serverless API access to top foundation models (Claude, Mistral, LLaMA, Titan)
- Azure AI Studio
 - Build AI copilots using GPT-4 etc., with deep integration into Office 365, security, and compliance controls.
- Google Vertex
 - Suited to advanced AI engineers / powerful tools to train/tune/customize LLMs.

Comparison

	AWS Bedrock	Azure AI Studio	Google Vertex
Overview	Fully managed GenAI platform for FMs	Unified AI development platform for Azure OpenAI & more	Full-service MLOps & GenAI platform
Model access	Anthropic (Claude), Meta (LLaMA), Mistral, Cohere, etc.	OpenAI (GPT), Meta (LLaMA), Mistral, etc.	Gemini, PaLM, LLaMA, Mistral, etc.
Serverless inference	Yes	Yes	Yes
Agent framework	Bedrock Agents (API orchestration + tools), Sagemaker integration	Manual tool calling with functions	Agent Builder for tool use, planner, memory
Suited for	Teams needing fast access to multiple hosted FMs	Enterprises with Microsoft stack, Azure OpenAI users	GCP users and advanced AI/ML engineer

Introduction to AWS Bedrock

The screenshot displays the AWS Bedrock console. At the top is the AWS navigation bar with the logo, a search bar, and a language dropdown set to 'United States (Ohio)'. Below this is the 'Amazon Bedrock' header. A left-hand navigation menu lists several categories: 'Getting started' (Overview, Providers), 'Foundation models' (Model catalog, Marketplace deployments, Prompt Routers), 'Playgrounds' (Chat / Text), 'Builder tools' (Agents, Flows, Knowledge Bases, Prompt Management), 'Safeguards' (Guardrails), and 'Inference and Assessment' (Provisioned Throughput, Batch inference, Cross-region inference, Evaluations). The main content area features a blue banner for 'Introducing Prompt routers' with a 'View Prompt routers' button. Below this is the 'Overview' section, which includes a 'Foundation models' card with a 'View Model catalog' button and a 'Discover marketplace models' button. To the right of the Overview section is a 'Model spotlight' for 'Anthropic's Claude' with a link to 'Open in chat playground'. At the bottom, there are sections for 'Chat / Text' (with an 'Open playground' link) and 'Builder tools' (describing the creation of GenAI applications).

aws Search [Option+S] United States (Ohio)

Amazon Bedrock

Amazon Bedrock

- ▼ **Getting started**
 - Overview
 - Providers
- ▼ **Foundation models**
 - Model catalog [New](#)
 - Marketplace deployments [New](#)
 - Prompt Routers
- ▼ **Playgrounds**
 - Chat / Text
- ▼ **Builder tools**
 - Agents
 - Flows
 - Knowledge Bases
 - Prompt Management
- ▼ **Safeguards**
 - Guardrails
- ▼ **Inference and Assessment**
 - Provisioned Throughput
 - Batch inference
 - Cross-region inference
 - Evaluations

Introducing Prompt routers
Route requests between foundational models from the same family, optimizing for response quality and cost. [View Prompt routers](#)

Overview [Info](#)

Foundation models
Amazon Bedrock supports over 100 foundation models from industry-leading providers and emerging leaders. Select a serverless model or Bedrock Marketplace model that is best suited for achieving your unique goals.
[View Model catalog](#) [Discover marketplace models](#)

Model spotlight
AI **Anthropic's Claude**
Choose the exact combination of intelligence, speed, and cost to suit your needs. All of the latest Claude models, including Claude 4, are available in Amazon Bedrock. [Open in chat playground](#)

Chat / Text
Generate text for a vast range of language processing tasks with various pre-trained models. You can use a single prompt or iterate on the result by submitting subsequent prompts that take into account the context of previous prompts and generated responses in a chat format.
[Open playground](#)

Builder tools
Create GenAI applications, augment responses with your proprietary data, and experiment with prompts.

Region Availability

- Primarily us-east-1 and us-west-2 (most Bedrock services are available)
- There is variance by model provider – not all models are available everywhere
- What do I do if I'm in a different region?
 - Some models are accessible through cross-Region inference

- Refer here for updated region support:

<https://docs.aws.amazon.com/bedrock/latest/userguide/models-regions.html>

Who are the Model Providers?

Amazon Bedrock > Providers


Amazon Bedrock <

- ▼ **Getting started**
 - Overview
 - Providers**
- ▼ **Foundation models**
 - Model catalog [New](#)
 - Marketplace deployments [New](#)
 - Prompt Routers
- ▼ **Playgrounds**
 - Chat / Text
- ▼ **Builder tools**
 - Agents
 - Flows
 - Knowledge Bases
 - Prompt Management
- ▼ **Safeguards**
 - Guardrails
- ▼ **Inference and Assessment**
 - Provisioned Throughput
 - Batch inference
 - Cross-region inference
 - Evaluations


Providers (5)
Choose from a range of serverless providers to find the model that is best suited for your use case.

Q Find providers


< 1 > | ⚙️ [Grid View] [List View]




[Amazon](#) | 5 models | Serverless




[Anthropic](#) | 7 models | Serverless



[DeepSeek](#) | 1 models | Serverless



[Meta](#) | 10 models | Serverless



[Mistral AI](#) | 1 models | Serverless

Model Access

aws

Search

[Option+S]

United States (Ohio)

Amazon Bedrock > Model access > Manage model access

Step 1
Edit model access

Step 2
Review and submit

Edit model access

Base models (3/24) Collap

Not seeing a model you're interested in? Check out all supported models by region [here](#).

Group by provider

	Models	Access status	Modality	EULA
<input type="checkbox"/>	▼ Amazon (5)	0/5 access granted		
<input type="checkbox"/>	Titan Text Embeddings V2	Available to request	Embedding	EULA
<input type="checkbox"/>	Nova Pro Cross-region inference	Available to request	Text & Vision	EULA
<input type="checkbox"/>	Nova Lite Cross-region inference	Available to request	Text & Vision	EULA
<input type="checkbox"/>	Nova Micro Cross-region inference	Available to request	Text	EULA
<input type="checkbox"/>	Nova Premier Cross-region inference	Available to request	Text & Vision	EULA
<input type="checkbox"/>	▼ Anthropic (7)	1/7 access granted		
<input type="checkbox"/>	Claude 3 Haiku Cross-region inference	Available to request	Text & Vision	EULA
<input type="checkbox"/>	Claude 3.5 Sonnet Cross-region inference	Available to request	Text & Vision	EULA
<input type="checkbox"/>	Claude 3.5 Haiku Cross-region inference	Available to request	Text	EULA
<input type="checkbox"/>	Claude 3.5 Sonnet v2 Cross-region inference	Available to request	Text & Vision	EULA
<input type="checkbox"/>	Claude 3.7 Sonnet Cross-region inference	Available to request	Text & Vision	EULA
<input checked="" type="checkbox"/>	Claude Sonnet 4 Cross-region inference	Access granted	Text & Vision	EULA
<input type="checkbox"/>	Claude Opus 4 Cross-region inference	Available to request	Text & Vision	EULA
<input type="checkbox"/>	▼ DeepSeek (1)	0/1 access granted		
<input type="checkbox"/>	DeepSeek-R1 Cross-region inference	Available to request	Text	EULA

Playgrounds

The screenshot shows the AWS Bedrock Chat / Text playground interface. At the top, the AWS logo and a search bar are visible. Below the header, the page title is "AWS Console Home" and "Amazon Bedrock" is selected. The main section is titled "Chat / Text playground".

On the left side, there is a "Configurations" panel. It includes a "Mode" dropdown set to "Chat" and a "Compare mode" toggle. The "Configurations" panel lists the model "Llama 3.2 3B In... v1" (US Meta Llama 3.2 3B Instruct). Below this, there are two sections: "Randomness and diversity" and "Length".

The "Randomness and diversity" section has two sliders: "Temperature" (set to 0.5) and "Top P" (set to 0.9). The "Length" section has a "Response length" slider (set to 512). Below these, there is a "Guardrails" section with a dropdown menu and a "Manage guardrails" link. A "Reset all to default" link is also present.

The main area displays a chat conversation. The user prompt is "What is the James Webb Space Telescope?". The model response is a detailed paragraph about the James Webb Space Telescope (JWST), including its launch date, purpose, and key instruments. The response is formatted with a list of four key instruments and a list of four key science objectives.

At the bottom of the chat area, there is a text input field with a placeholder "Write a prompt. Press Shift + Enter to add a new line. Press Enter to generate a response." and a "Run" button.

Model Hyperparameters

- Temperature: a measure of how creative the model is
 - low Temperature = safe, predictable, deterministic answers. High Temperature = bold, imaginative ones with an increasing amount of creativity.

Temp	Sample Model Response
0.2	The sky appears blue because of Rayleigh scattering. This occurs when sunlight passes through the atmosphere and shorter wavelengths of light, such as blue, are scattered more than longer wavelengths like red.
1	The sky looks blue because sunlight interacts with Earth's atmosphere. As sunlight passes through the air, shorter wavelengths — like blue — scatter more easily in all directions. This scattering, known as Rayleigh scattering, causes us to see the sky as blue most of the time. It's the same reason sunsets can appear red or orange — the angle of the sun changes how much the light is scattered.
2	The sky appears blue due to the scattering of solar radiation by molecules in Earth's atmosphere, particularly a phenomenon called Rayleigh scattering. But if you really think about it, maybe the sky is blue because it's Earth's way of dreaming—its way of whispering ocean memories back to the clouds. After all, light is just nature's poetry in motion.
5	Because marshmallow theories of existential yogurt synchronize with nebula-based fashion regulations. Frogs in bowler hats interpret wavelengths through interpretive dance on Tuesdays, especially near cosmic pancakes. Hence, blue. Obviously.
100	Skibble zentropics cloudmonger flibber astro-tuba syntax bananas! Plasmic shoelaces dangle from hypothetical wombats under bureaucratic sky-harmoniums. Reason? Smeep.

Top-K Sampling

- Instead of selecting a token from the entire vocabulary, the model focuses only on the top K most likely tokens and reassigns probabilities within this smaller group. This introduces controlled randomness.

Example: if $K = 5$, the model would narrow its choices to just five tokens:

- "the" (30%)
- "a" (25%)
- "this" (20%)
- "that" (15%)
- "one" (10%)

Top-P (nucleus) Sampling

- Similar to Top-K, but instead of being a fixed number, it dynamically adjusts the size of the candidate pool based on a cumulative probability mass
- Example: if P is set to 0.5, it will only look at the tokens that have a cumulative probability of 50%, and will ignore the rest.
 - E.g. if it takes the top 78 tokens to reach 50%, only these will be examined as a next token option.

Connecting to Bedrock Inference Models

- Need to do a few things to prep:
 - Ensure AWS configuration is set up on your machine
 - Set up a python virtual environment
 - Install the boto3 SDK
 - Ensure you know the right AWS region you will connect to

Setting up AWS Access on your Machine

(may need to install it with brew)

```
[robbarto@ROBBARTO-M-C6YD .aws % aws configure
AWS Access Key ID [None]: 12323
AWS Secret Access Key [None]: 12321341235
Default region name [None]: us-east-1
Default output format [None]: json
```

```
[robbarto@ROBBARTO-M-C6YD .aws % ls -al
total 16
drwxr-xr-x@  4 robbarto  staff   128 Jun 19 12:04 .
drwxr-x---+ 68 robbarto  staff  2176 Jun 19 14:26 ..
-rw-----@  1 robbarto  staff    54 Jun 19 12:04 config
-rw-----@  1 robbarto  staff   117 Jun 19 12:02 credentials
robbarto@ROBBARTO-M-C6YD .aws %
```


Demo Time!

DEMO

1. Confirm you are connecting to AWS Bedrock
2. Perform simple inference to a model
3. Try a Streamlit chatbot with using a Bedrock model