

Unlocking FMOps/LLMOps using Apache Airflow : A guide to operationalizing and managing Large Language Models

Parnab Basak (he/him)

Senior Solution Architect and Apache Airflow Specialist

Amazon Web Services





Agenda

for the next 20 mins

- Machine Learning (ML) Ops overview
- Machine Learning Ops foundation
- Intro to Foundation Model (FM)/Large Language Model (LLM) Ops
- Airflow for FM Ops/LLM Ops
- ML Ops vs FM/LLM Ops Differentiators
- QnA



Machine Learning (ML) Ops

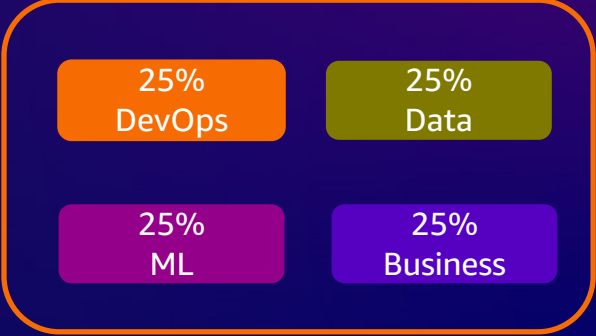
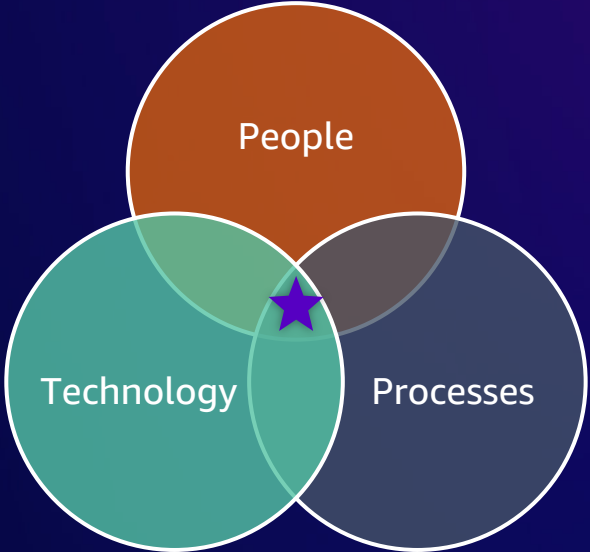


What is MLOps?

MLOps

Machine Learning & Operations

The combination of **people, processes, and technology** to productionize ML solutions efficiently.



Productivity



Reproducibility



Reliability



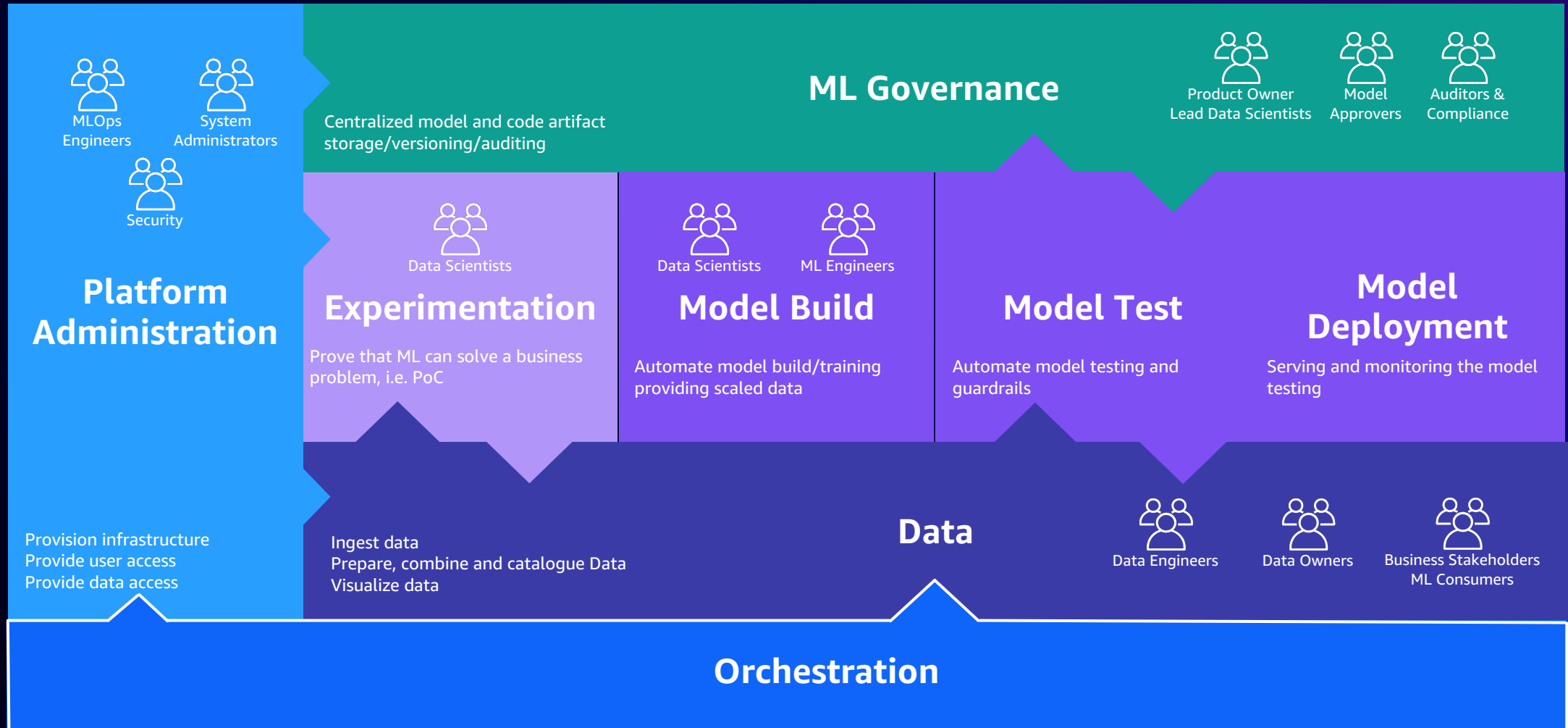
Observability



Lower TCO

MLOps Foundation People & Processes

SEPARATION OF CONCERNS IS KEY FOR SUCCESS



Airflow as the Orchestrator

OUT OF THE BOX FEATURES TO AID MLOPS



Macros & Jinja Templates



Advanced Dataset Scheduling



Dynamic Task Mapping



TaskFlow API



Lineage



Provider Packages



Backfills



Automatic Retries



Task Group



Airflow Plugins



Setup and Teardown



Reruns



Dynamic Compute



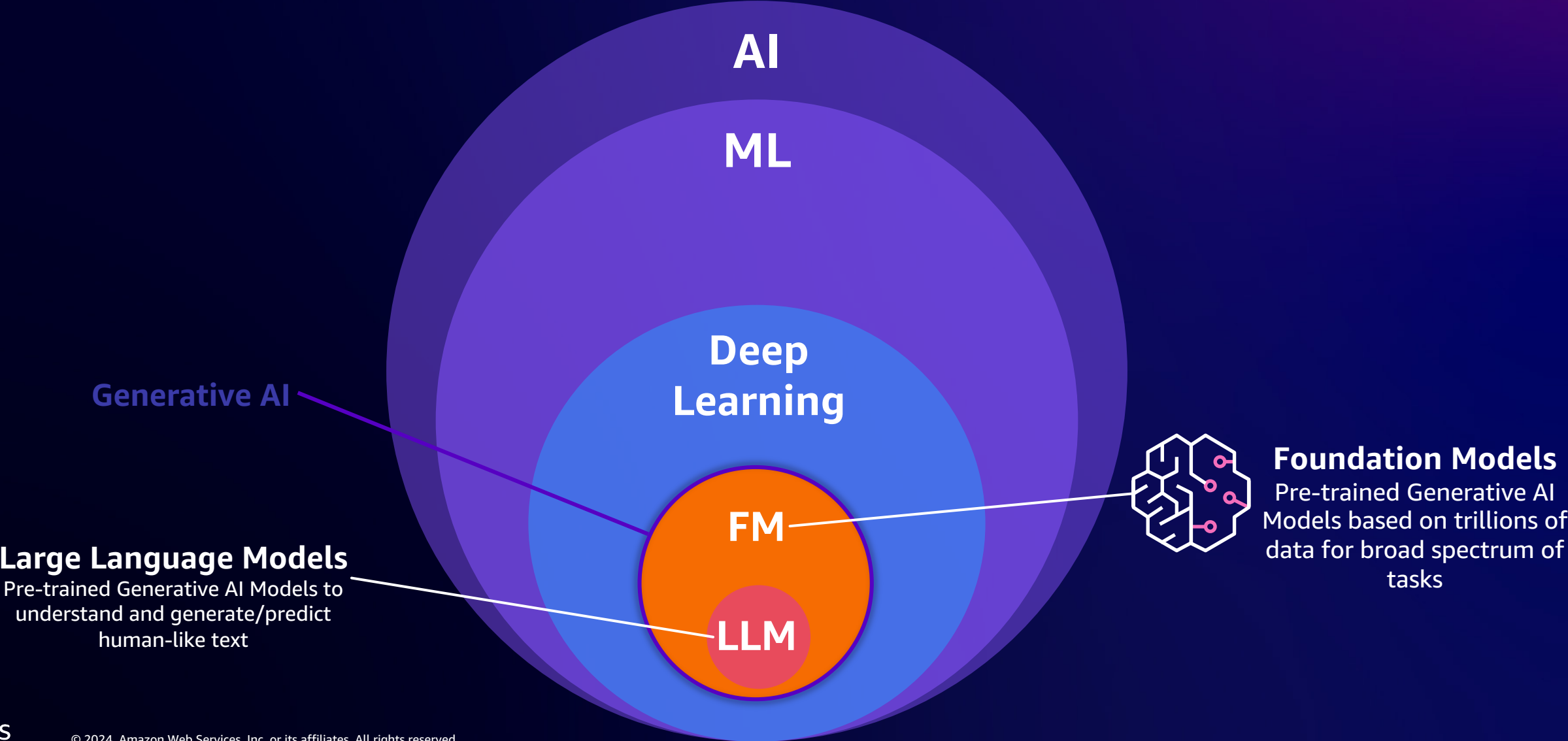
DS Toolkit & IDE Integration



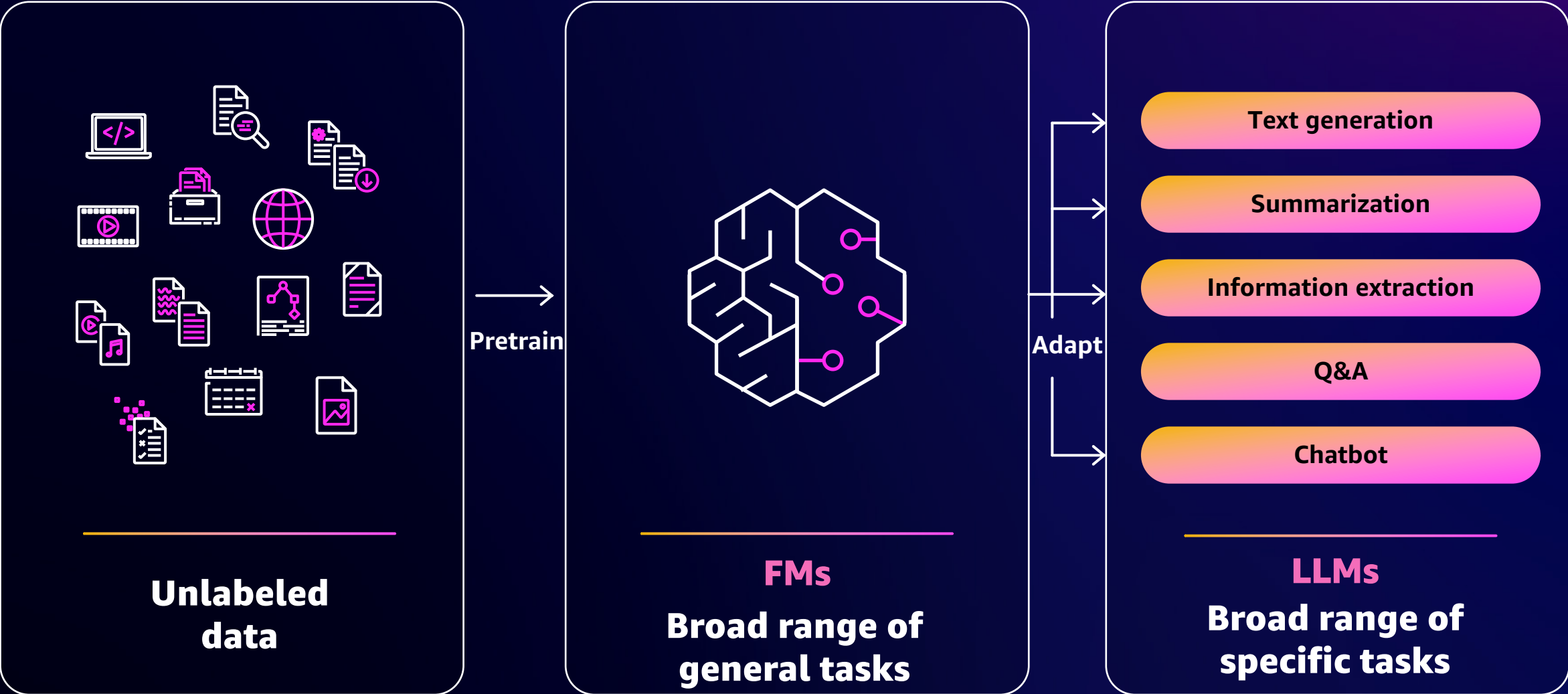
Monitoring & Alerting

Generative AI & Foundation Model Ops/Large Language Model Ops

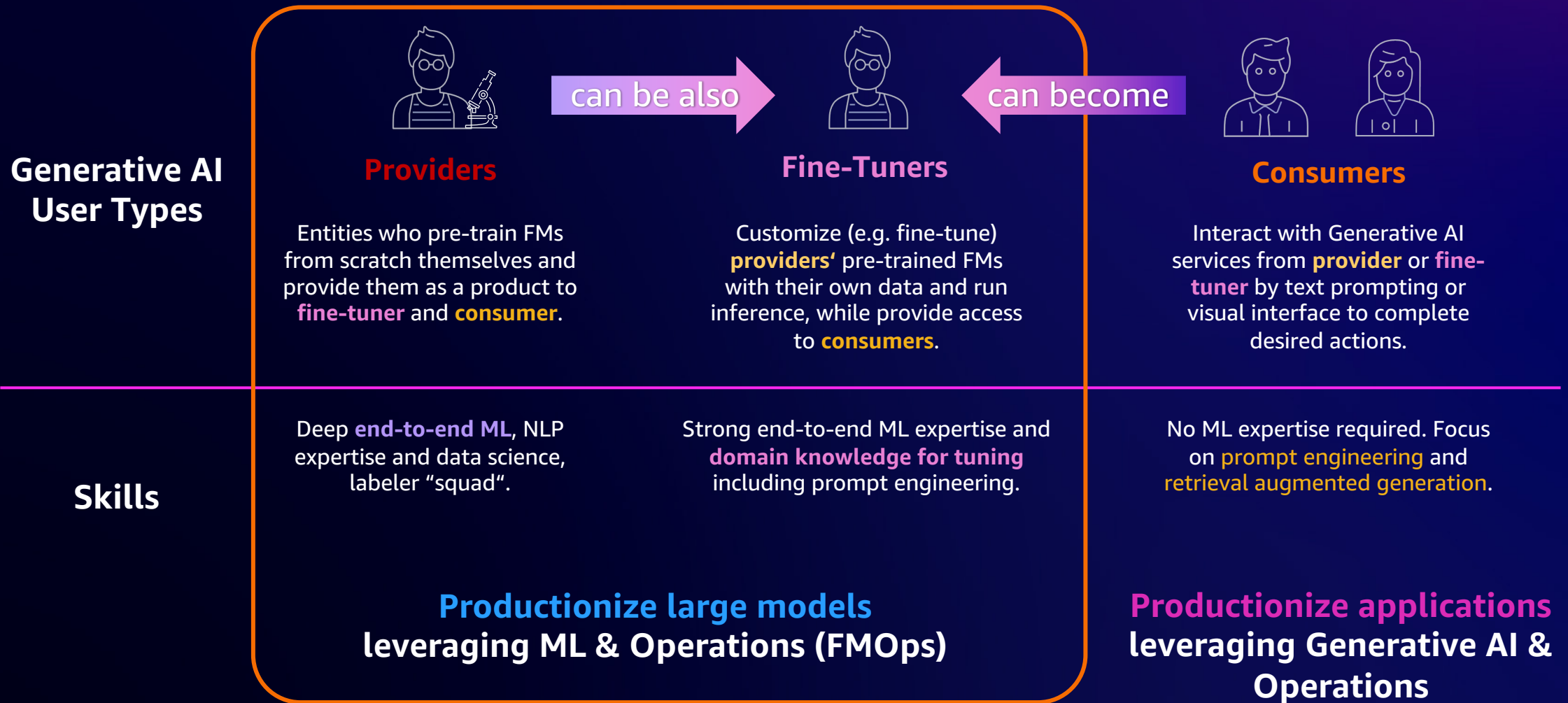
What is Generative AI?



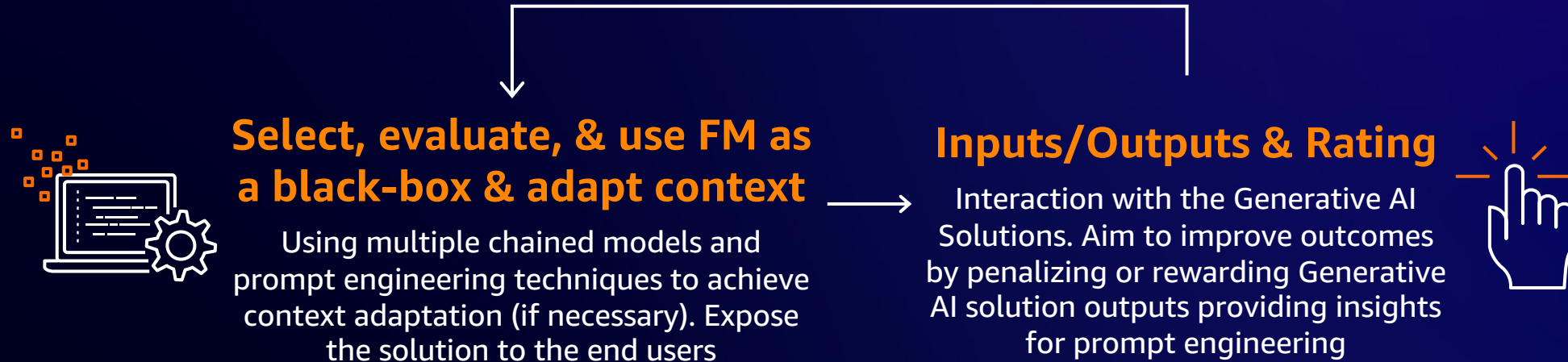
Generative AI overview



Generative AI User Types & Skills

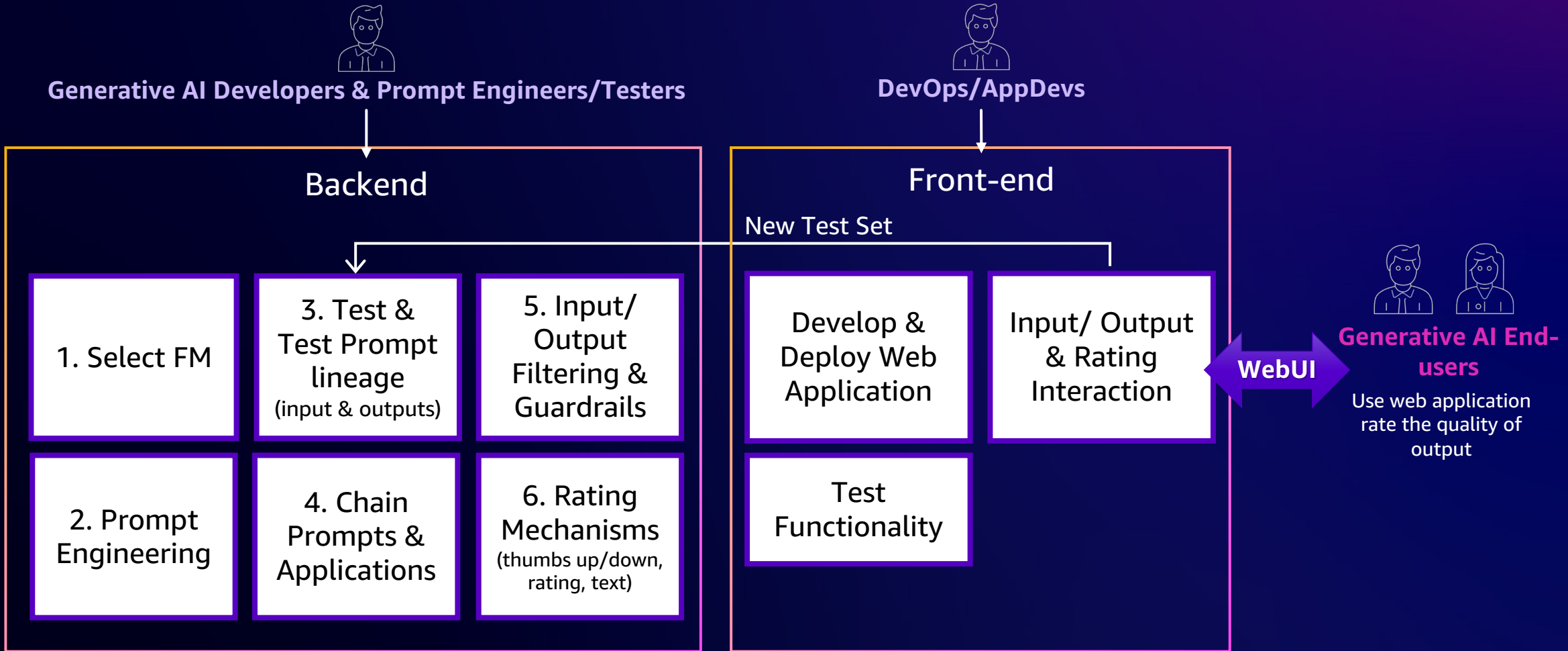


Generative AI Processes – Consumers



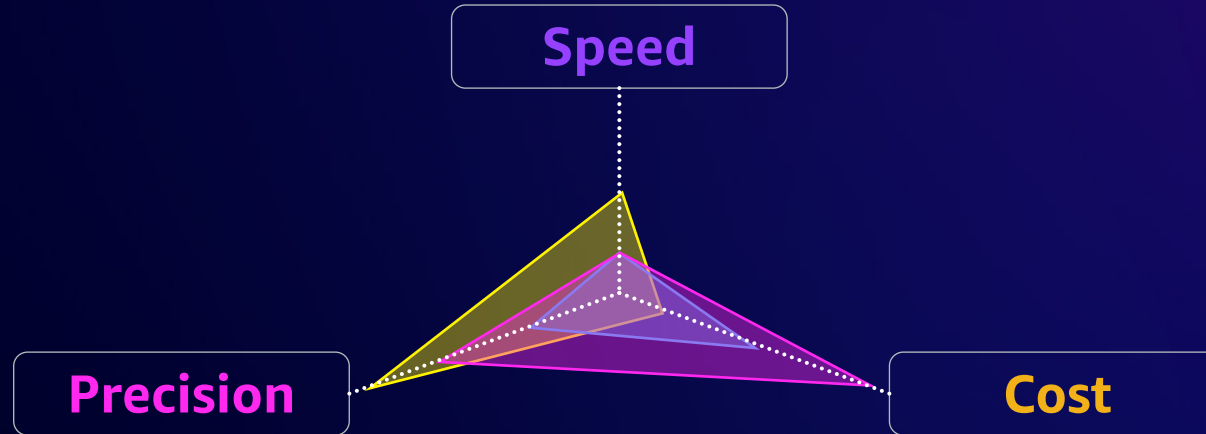
Generative AI Processes for LLM – Consumers

LLM-based Generative AI Solution



Model/Prompt Selection

Model Selection



Automatic evaluation



Accuracy



Robustness



Toxicity

Human evaluation



Creativity



Style



Tone



Relevance

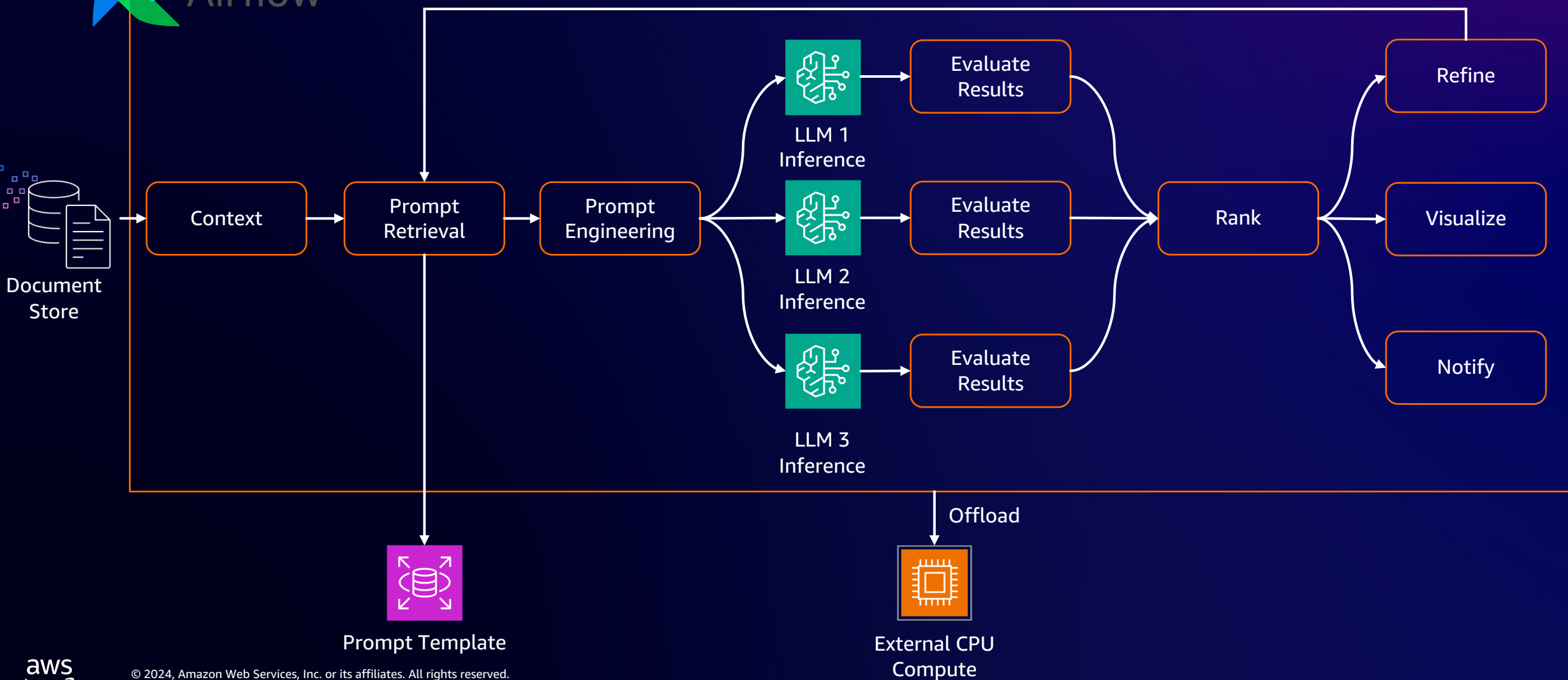


Coherence

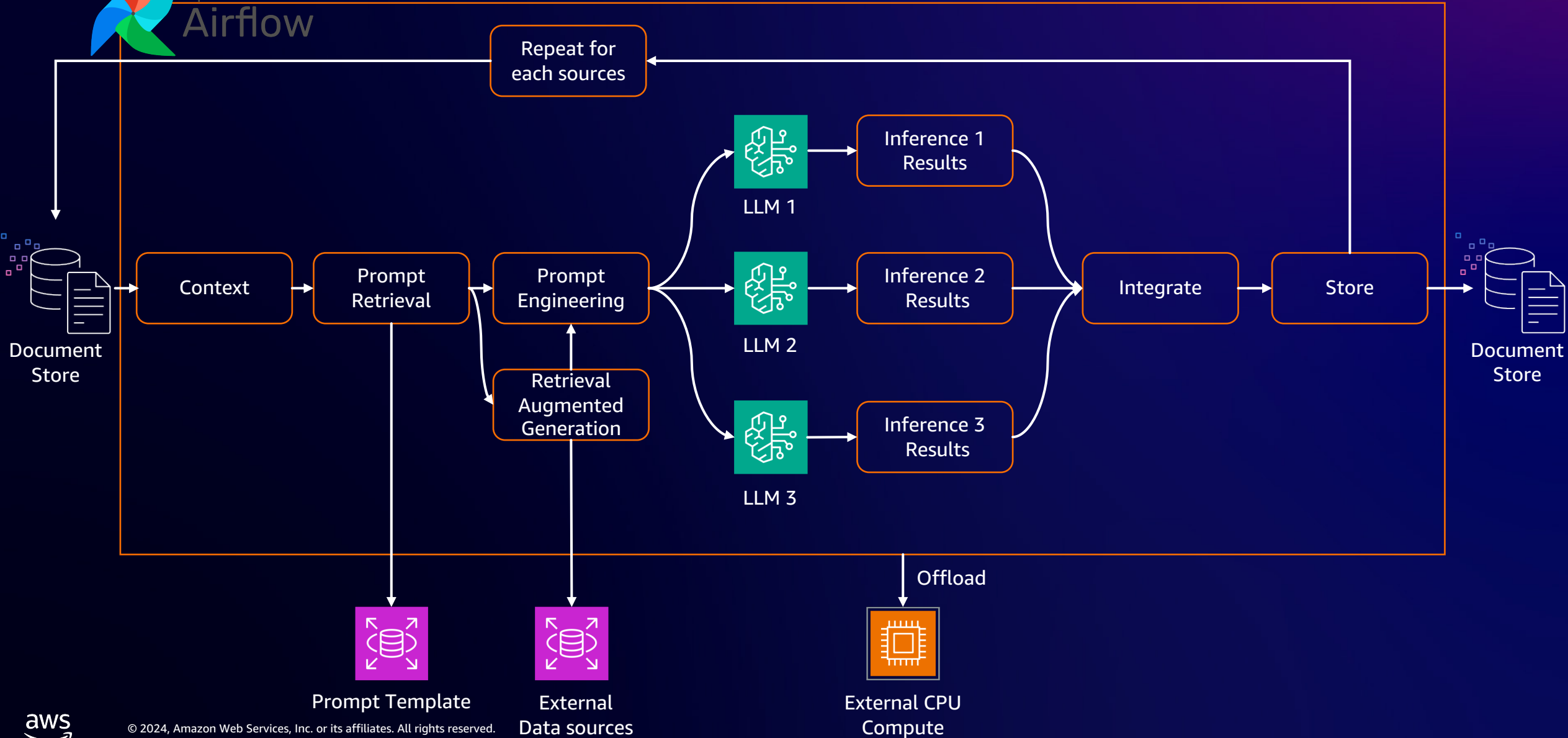


Brand voice

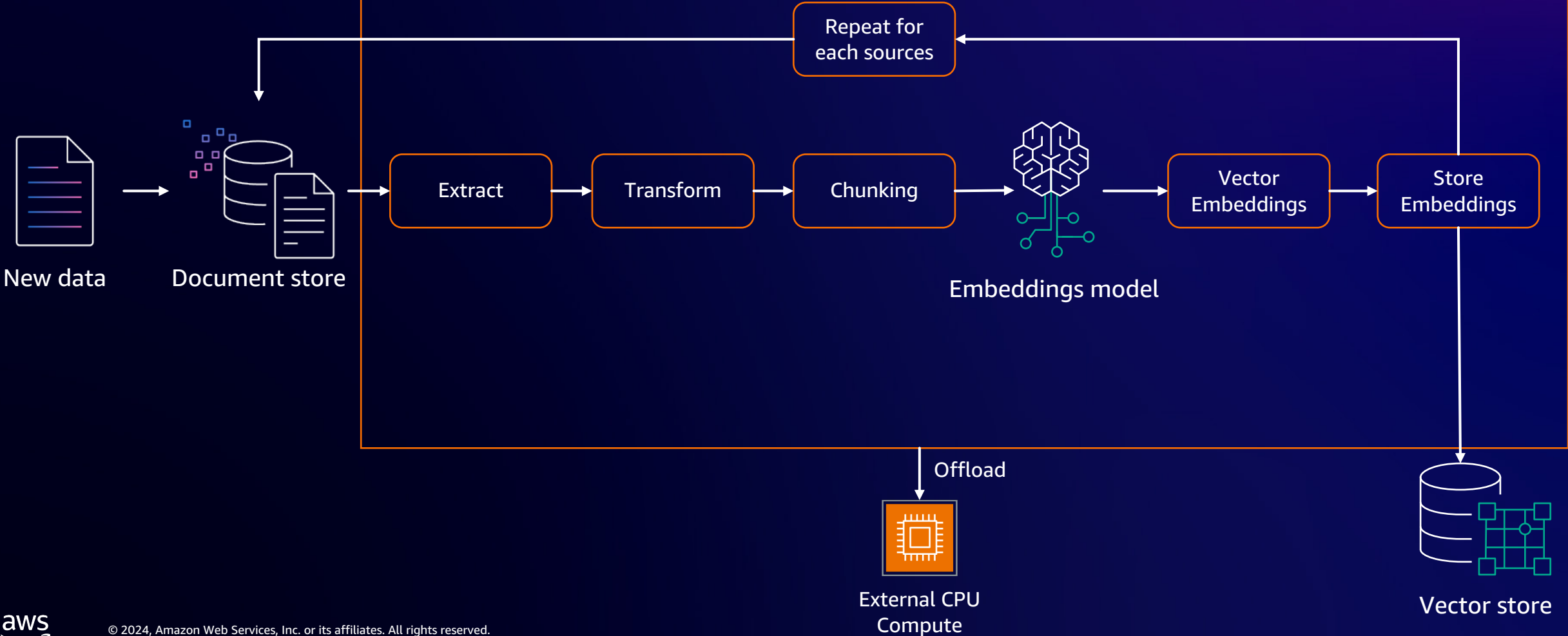
Model/Prompt Evaluation with Airflow



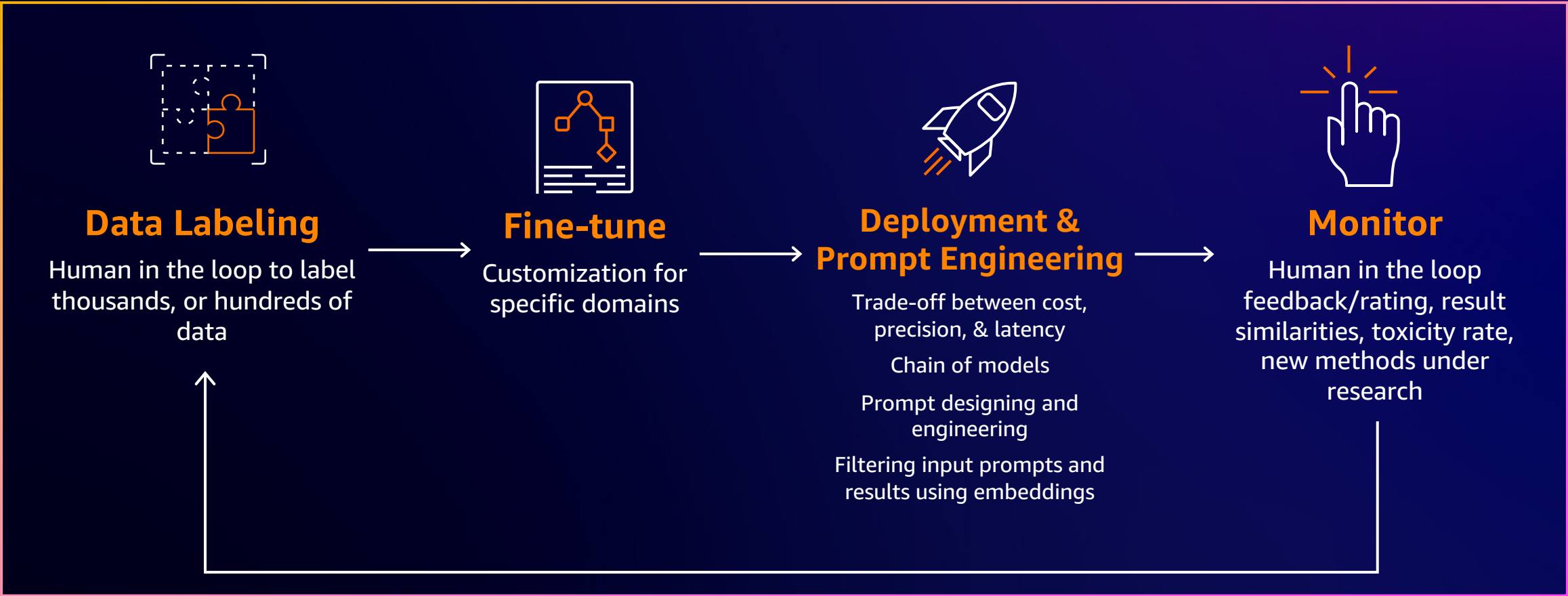
Offline Batch Inferencing with Airflow



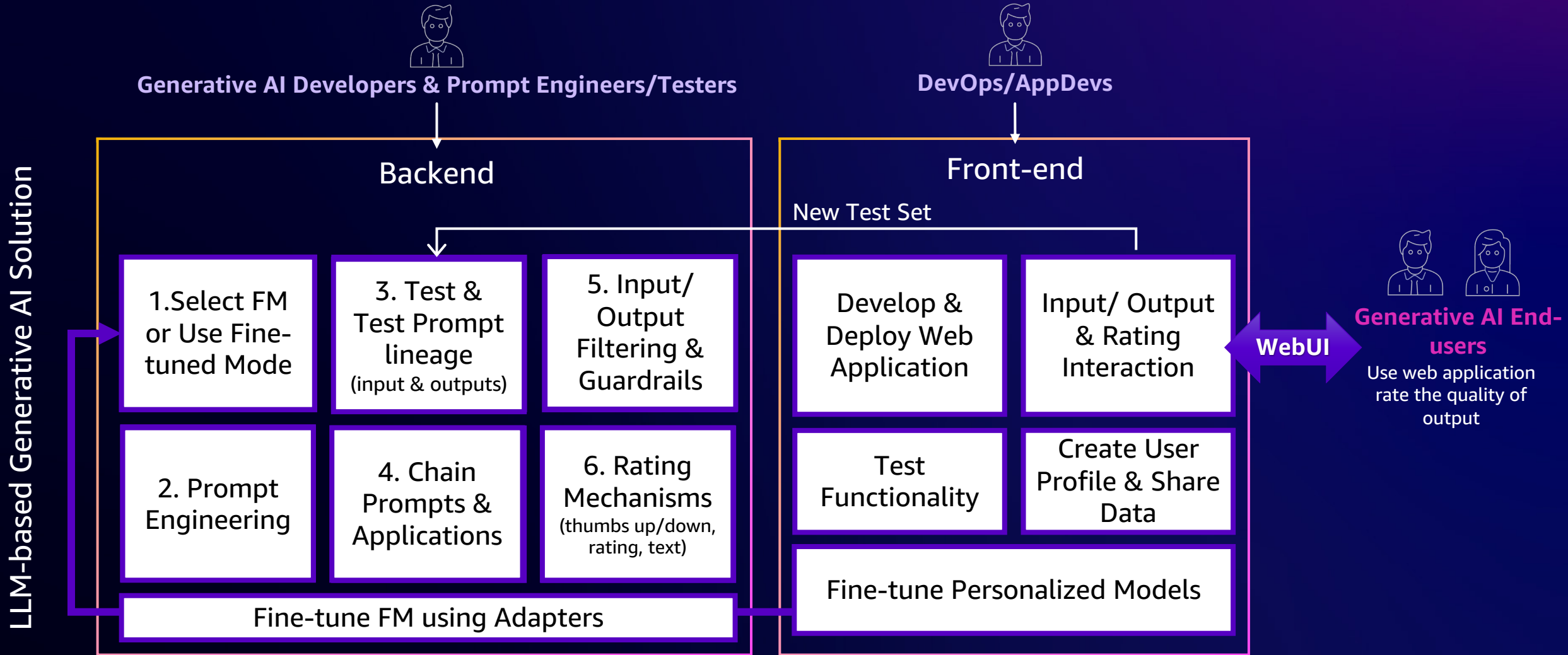
RAG Data Ingestion with Airflow



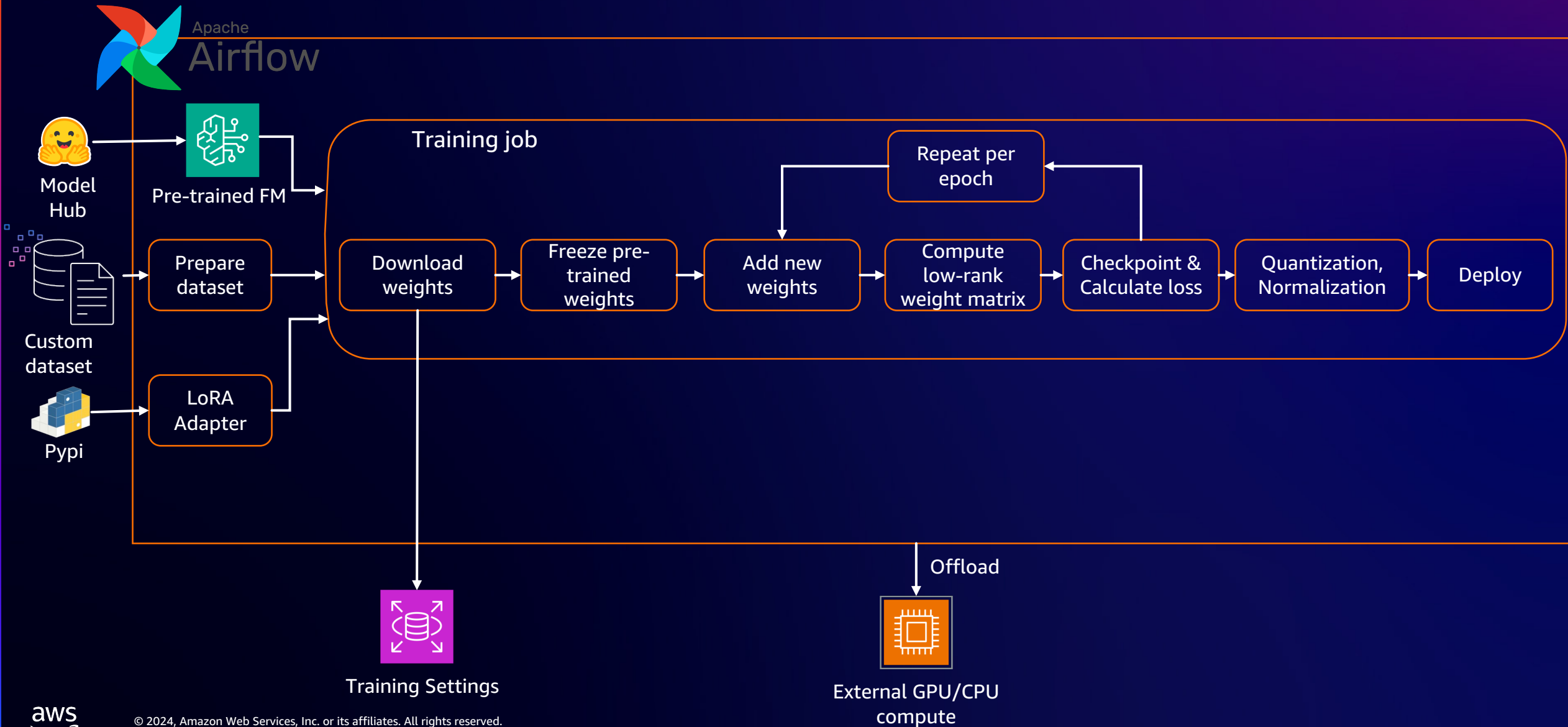
Generative AI Processes - Fine-Tuners



Generative AI Processes for LLM – Fine-Tuners

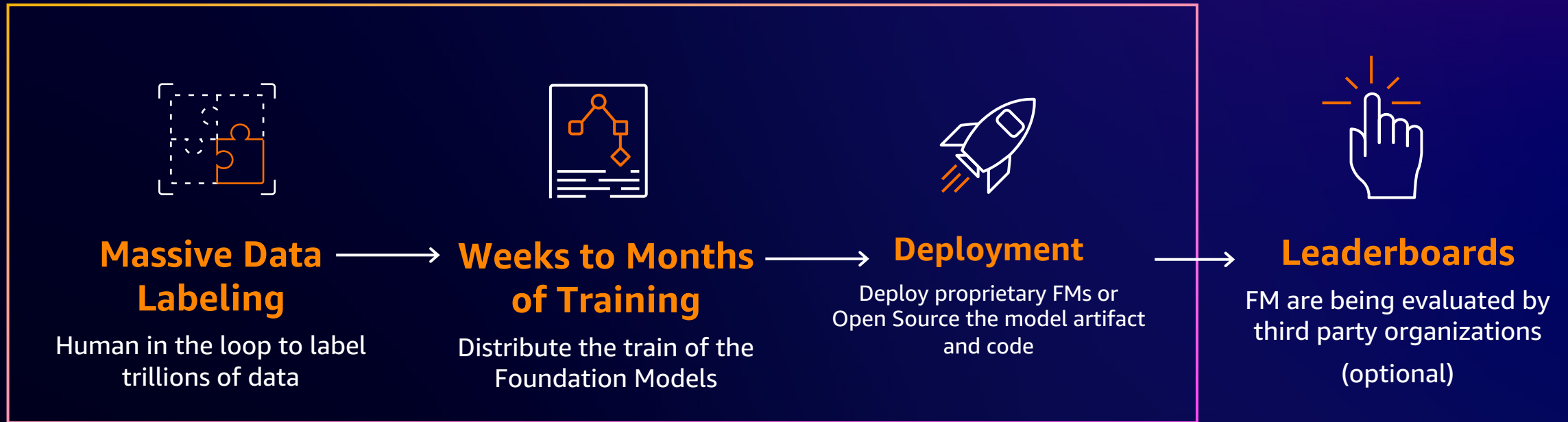


Fine-tuning PEFT (LoRA or QLoRA) with Airflow



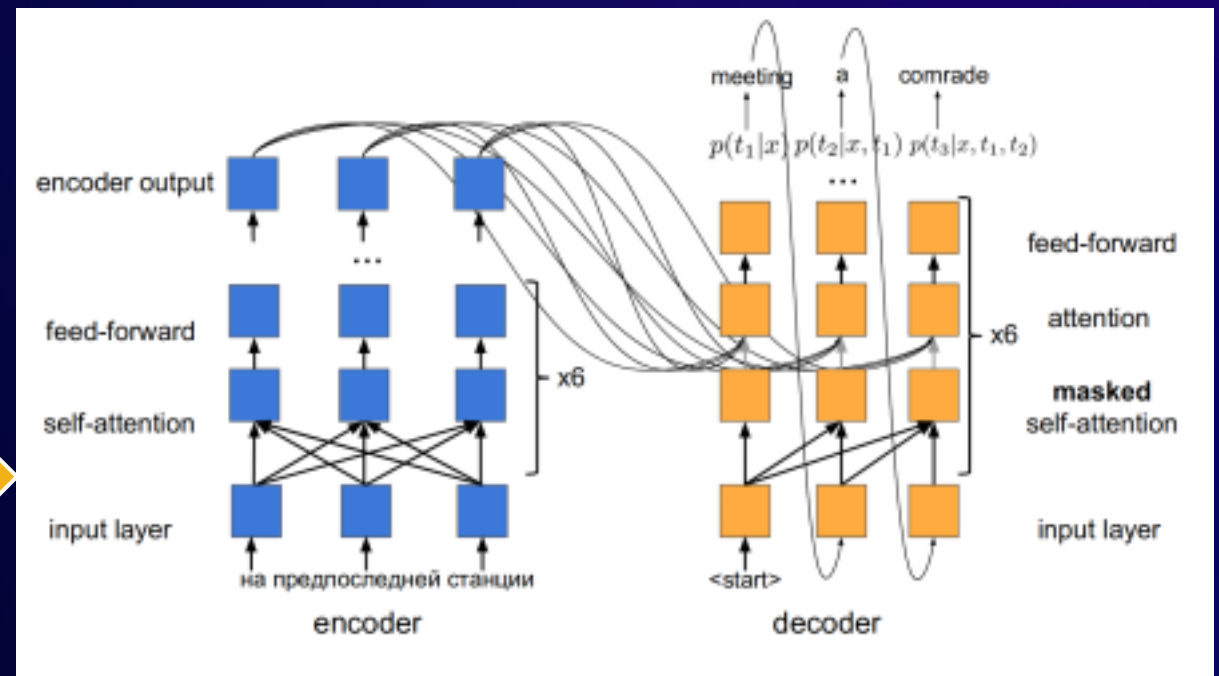
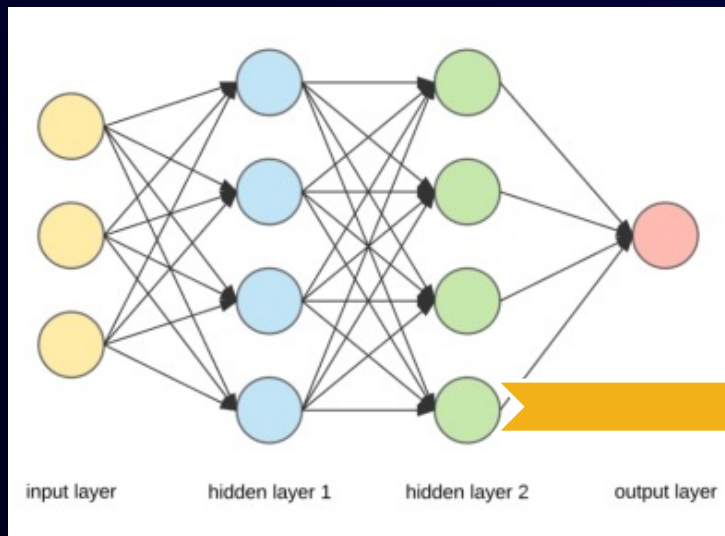
Generative AI Processes - Providers

PRE-TRAINING AN FM IS NOT A TRIVIAL TASK



Encoder/Decoder Transformation Model

HIGH LEVEL IMPLEMENTATION

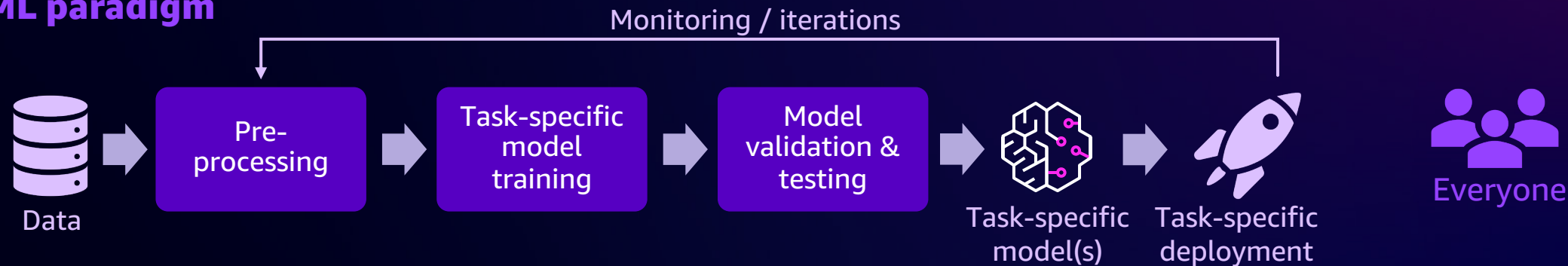


Generative AI & MLOps

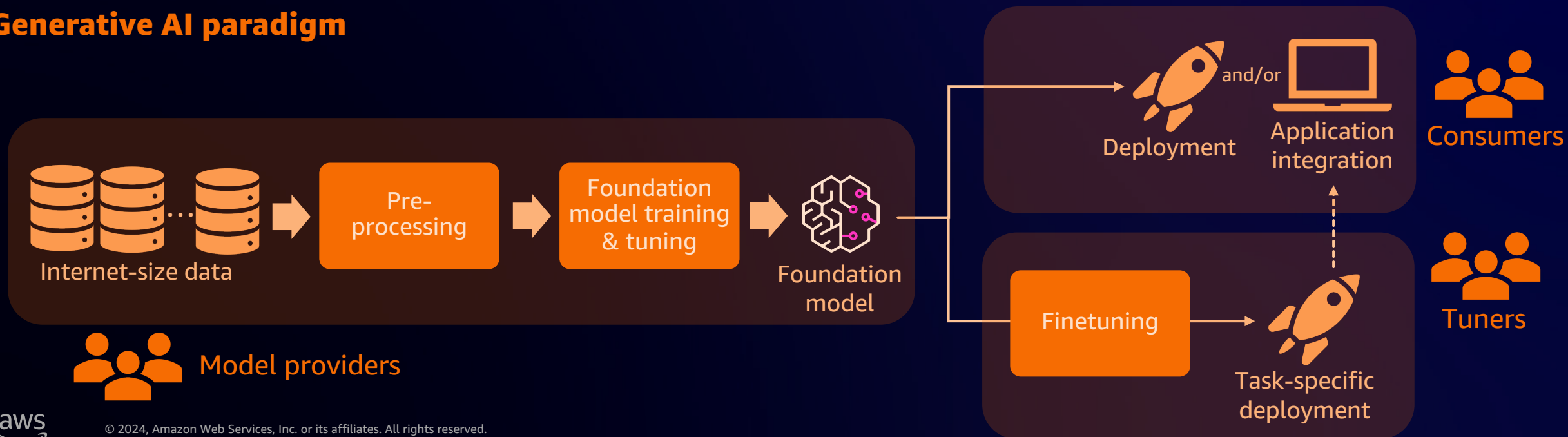
MLOps & FMOPs/LLMOPs Differentiators

Foundation model lifecycle

Conventional ML paradigm



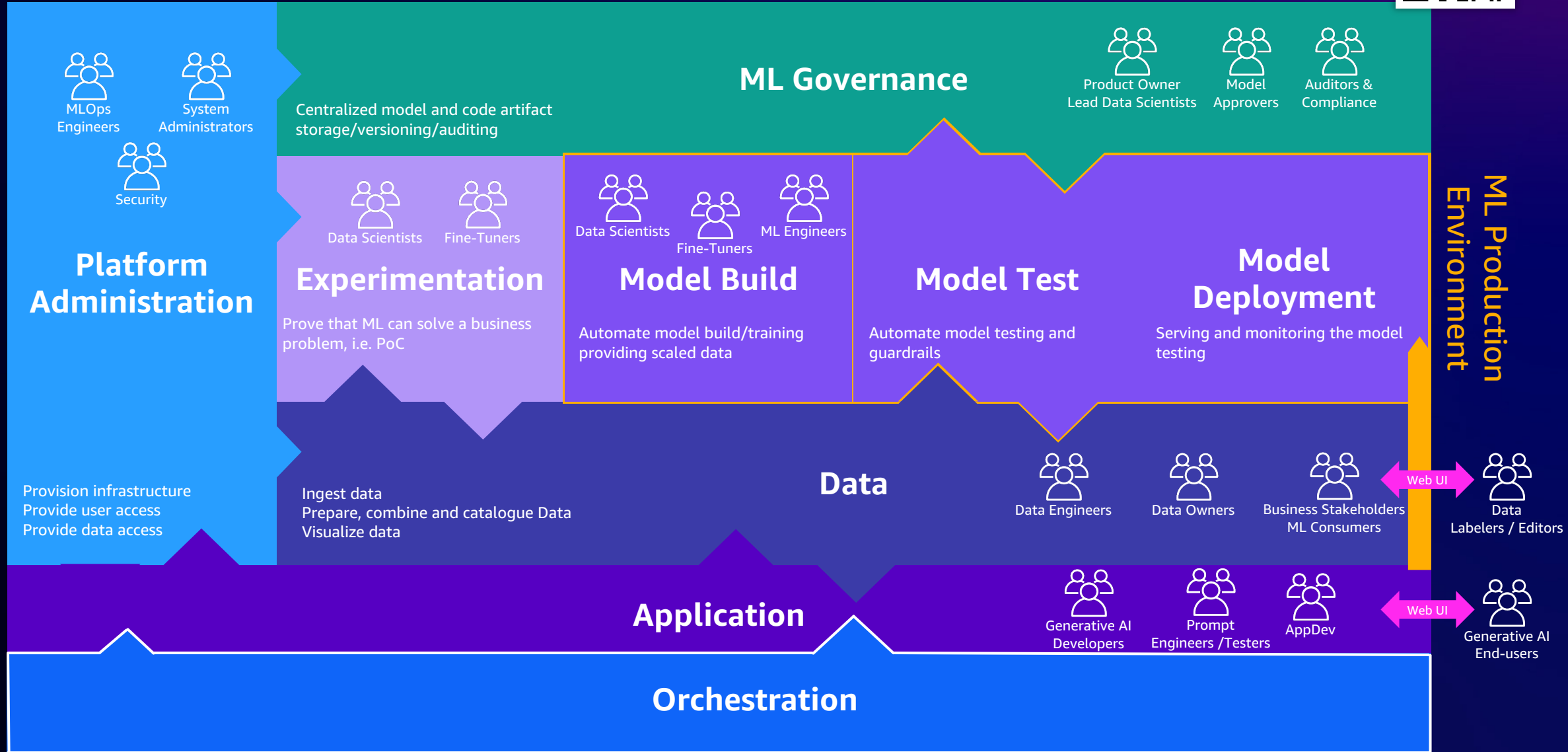
Generative AI paradigm



FMOps Foundation People & Processes



SEPARATION OF CONCERNS IS KEY FOR SUCCESS



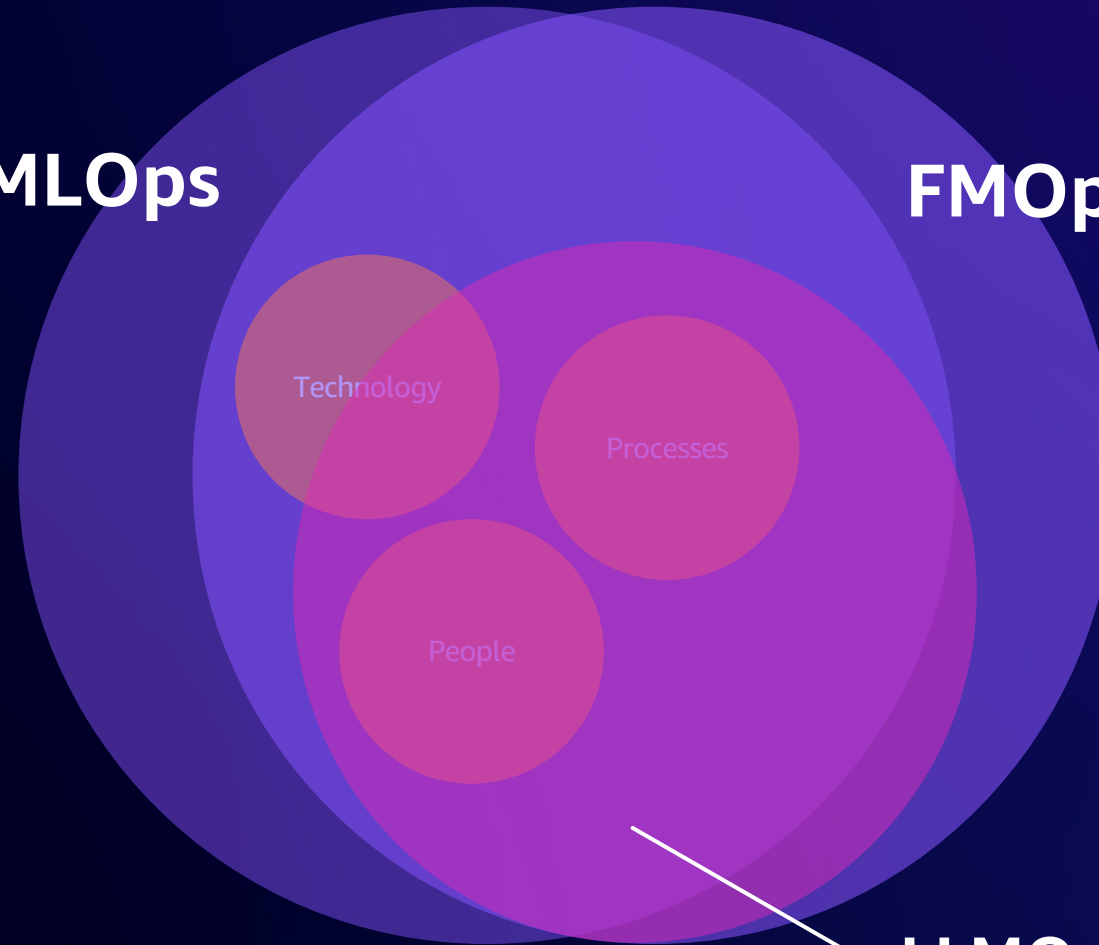
MLOps/FMOps/LLMOps

Machine learning operations
Productionize ML solutions
efficiently

MLOps

FMOps

Foundation model operations
Productionize generative AI solutions
(text-text/image/video/audio/...)



LLMOps

Large language model operations
Productionize large language
model-based solutions

Differences between MLOps and FMOps

MLOps

Technology

People

Processes

FMOps

Providers, fine-tuners, & consumers

Select & Customize the FM on a Specific Context

- Fine-tuning, parameter-efficient fine-tuning, prompt engineering

Proprietary, open source based on the application

Evaluate & Monitor Fine-tuned Models

Human feedback, prompt management, toxicity/bias...

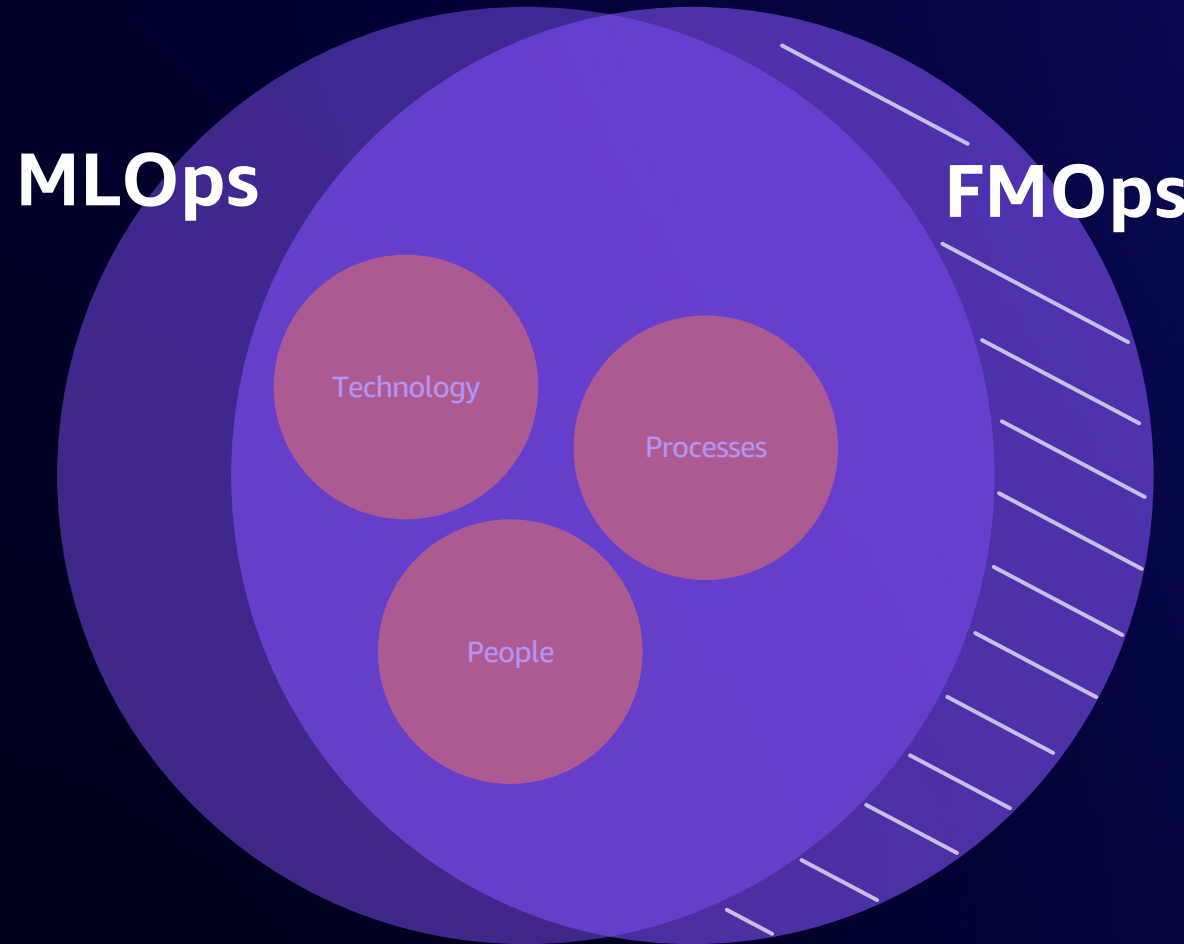
Data & Model Deployment

Data privacy, multi-tenancy, & cost, latency, and precision

Technology

MLOps, data, & application layers

MLOps & FMOps Differentiators



Processes & People

Providers, fine-tuners, & consumers

Select & Customize the FM on a Specific Context

- Fine-tuning, parameter-efficient fine-tuning, prompt engineering
- Proprietary, open source based on the application

Evaluate & Monitor Fine-tuned Models

Human feedback, prompt management, toxicity/bias...

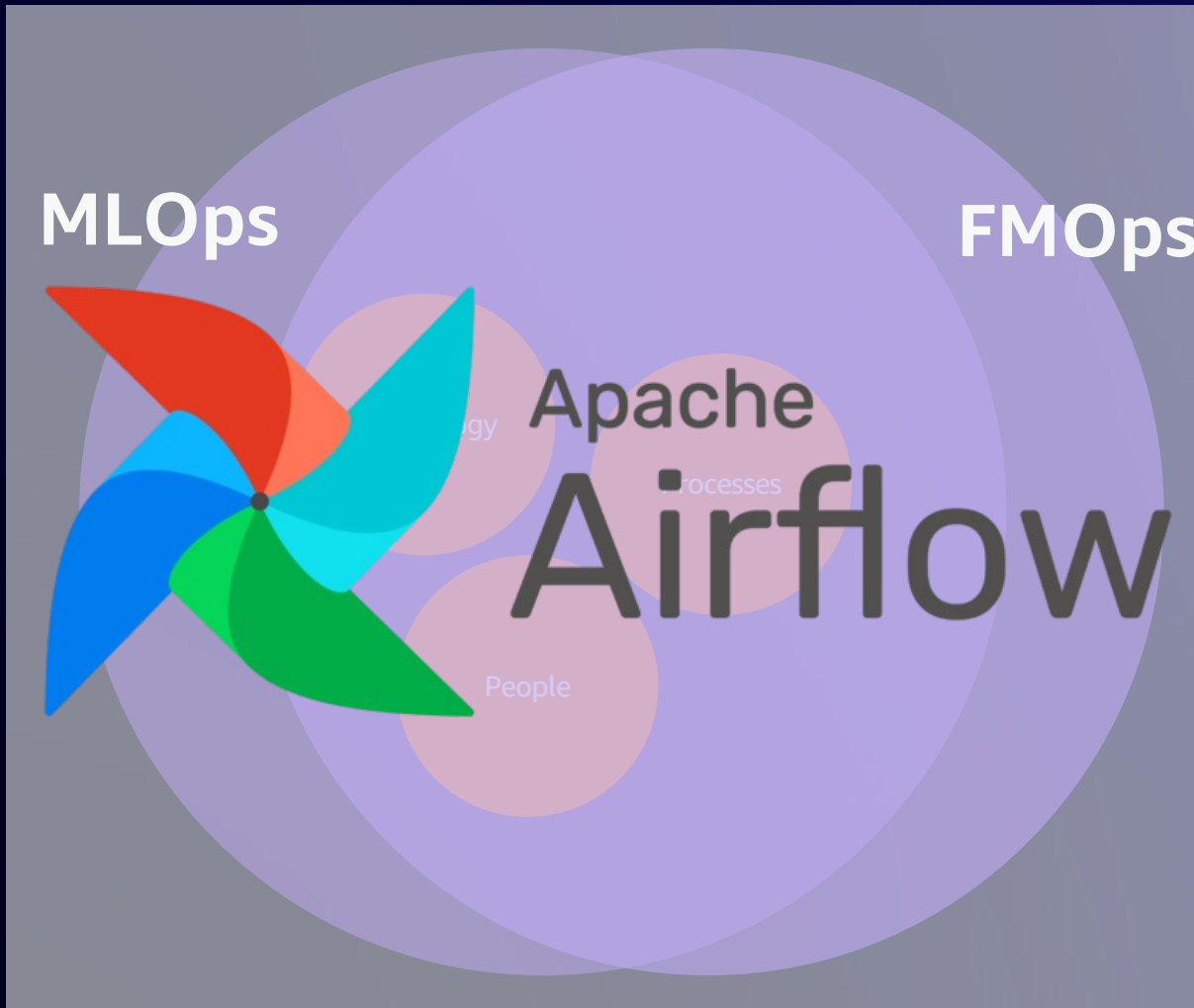
Data & Model Deployment

Data privacy, multi-tenancy, & cost, latency, and precision

Technology

MLOps, data, & application layers

Orchestration with Apache Airflow



Workflows as Code

- Dynamic, Extensible & Flexible
- Data and Compute Agnostic
- Scalability & Reliability
- Rich Ecosystem
- Community Driven
- Continuous Innovation

Generative AI & Operations Resources



Scaling AI Workloads with Apache Airflow

[Today 4:35 PM @Elizabethan A+B](#)



FMOps/LLMOps: Operationalize generative AI and differences with MLOps

<https://aws.amazon.com/blogs/machine-learning/fmops-llmops-operationalize-generative-ai-and-differences-with-mlops>



Operationalize LLM Evaluation at Scale using Amazon SageMaker Clarify and MLOps services

<https://aws.amazon.com/blogs/machine-learning/operationalize-llm-evaluation-at-scale-using-amazon-sagemaker-clarify-and-mlops-services>



Build an internal SaaS service with cost and usage tracking for foundation models on Amazon Bedrock

<https://aws.amazon.com/blogs/machine-learning/build-an-internal-saas-service-with-cost-and-usage-tracking-for-foundation-models-on-amazon-bedrock/>



Session **Survey**



<https://pulse.aws/survey/UUNNEDF8>

Questions ?

Thank you!



Please complete the session survey in the mobile app