

300 LEVEL

FMOps/LLMOps: Operationalise Generative AI using MLOps principles

Dr Sokratis Kartakis (he/him)
Senior MLOps SA Architect, EMEA
Amazon Web Services

Heiko Hotz (he/him)
Senior LLM SA Architect, EMEA
Amazon Web Services



Agenda

MLOps Foundation Overview

- MLOps KPIs, Maturity, People, Processes, Technology

Generative AI (GenAI) & MLOps

- Main Definitions

MLOps & FMOps/LLMOps Differentiators

- Processes & People
- Providers, fine-tuners, & consumers
- Select & Adapt the FM on a Specific Context
- Evaluate & Monitor Fine-tuned Models
- Data & Deployment
- Technology

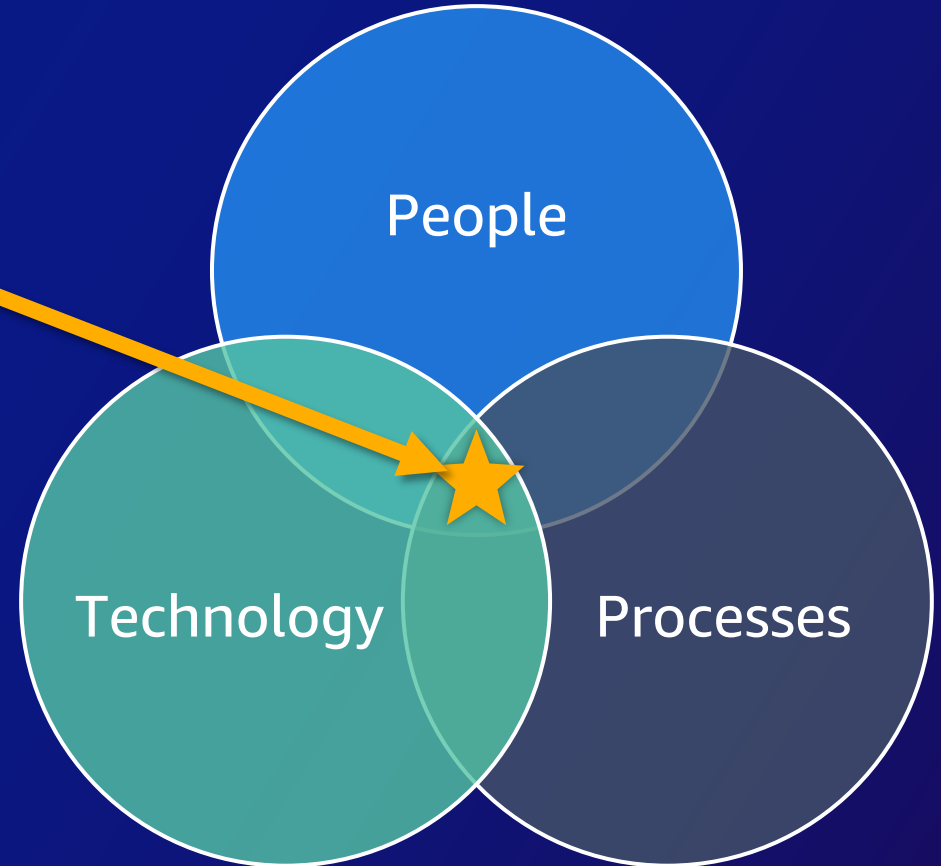
What is MLOps?

MLOps

Machine Learning
& Operations

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

MLOps Definition



MLOps Foundation Expected Outcomes

STANDARDIZE OPERATIONS AND INFRASTRUCTURE FOR YOUR DATA SCIENCE

	Business Goal	Technical Metric	Before MLOps	MLOps Expected Outcomes	Business Value
1	Be more efficient in delivery	Time to value (from idea to production)	up to 12 months	< 3 months	Improve Speed-to-Value by 4x
2	Simplify route-to-live	Time to productionize existing ML use cases	3-6 months	< 2 weeks	Reduce FTE overhead in average 8x
3	Standardize infrastructure, data, & code	% Template driven development	n/a	> 85%	Focus on innovation increasing re-usability by 85%
4	Standardize onboarding of new teams and ML use cases	Time to instantiate a new MLOps infrastructure & ML projects	40 days	< 1 hours	Accelerate ML adoption across all business areas
5	Ensure high security standards	Execute the ML solutions without internet access in a private cloud	n/a	No internet	Your data is safe in your private cloud

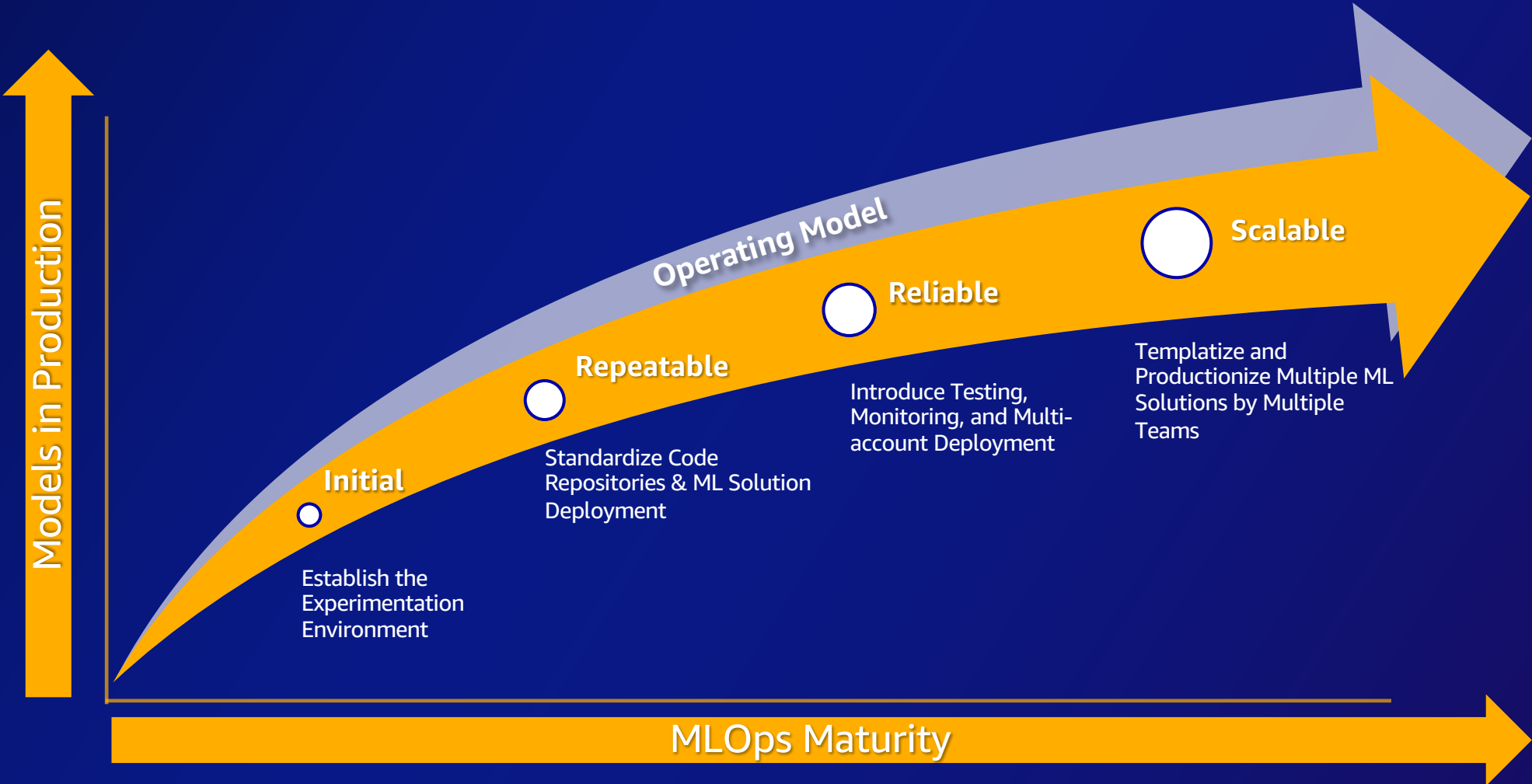
Reduce platform, people and operation costs

Customer references building MLOps foundation and business benefits:

- NatWest: <https://aws.amazon.com/solutions/case-studies/natwest-group-case-study>
- BP: <https://aws.amazon.com/solutions/case-studies/bp-machine-learning-case-study>



MLOps Maturity Model



MLOps Key Personas and Roles

Advance Analytics Team Data Lake



Data Engineer

Prepare & Ingest data building ETL pipelines



Data Owners

Manage data sharing and provide access

Data Science Team Experimentation & MLOps



Data Scientist

Create the best ML models to solve business problems



ML Engineer

Collaborate with DS to productionize ML

Platform Team Secure Cloud/Data/ML Platform



MLOps Engineer/Admin

Standardize CI/CD, user/service role, model consumption, testing and deployment methodology



Security

Assess data, user, and service access creating policies and guardrails



Architects/ SysOps Engineer

Standardize account infrastructure, connectivity, user roles implementation

Business Viz Dashboards, ML Adoption, & ROI



Business Stakeholder Product Owners

Define business problem, business KPIs, and make business decisions



Business Stakeholder Data & ML Consumers

Consumers of ML results from other BUs, driving business decision making

Risk & Compliance Approve & Review Models

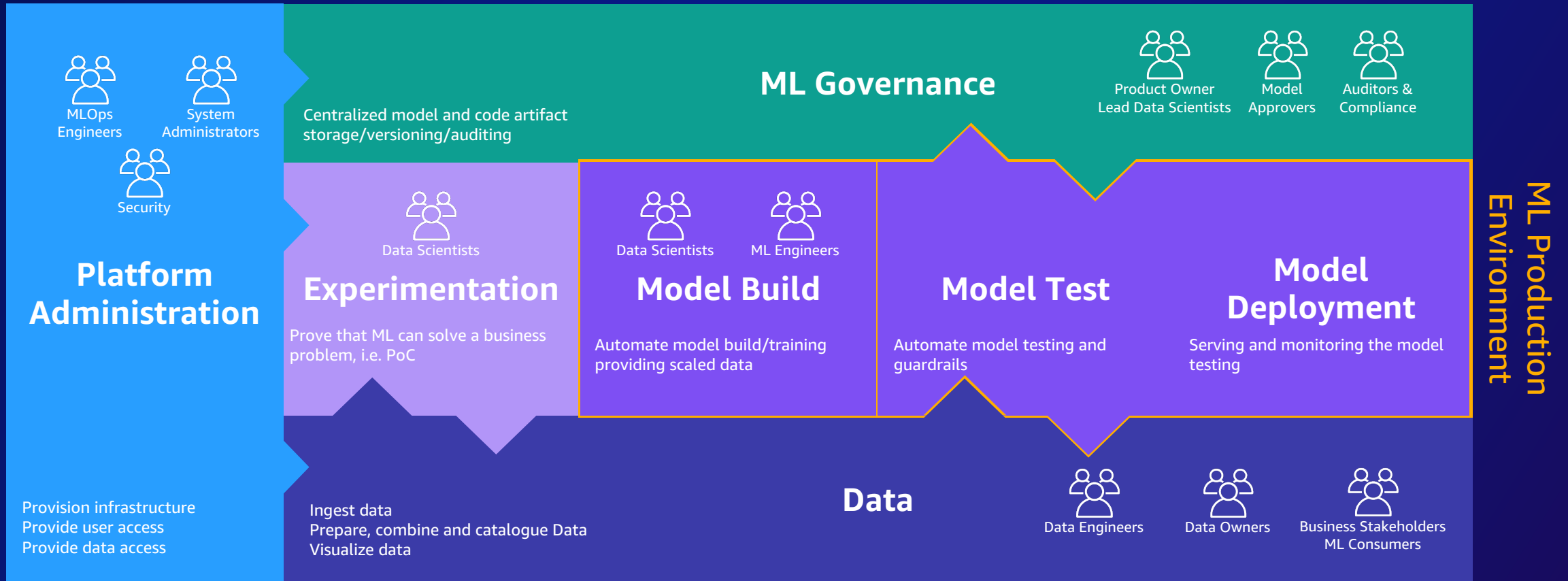


Auditors/Risk & Compliance

Review models, data sources, code artifacts

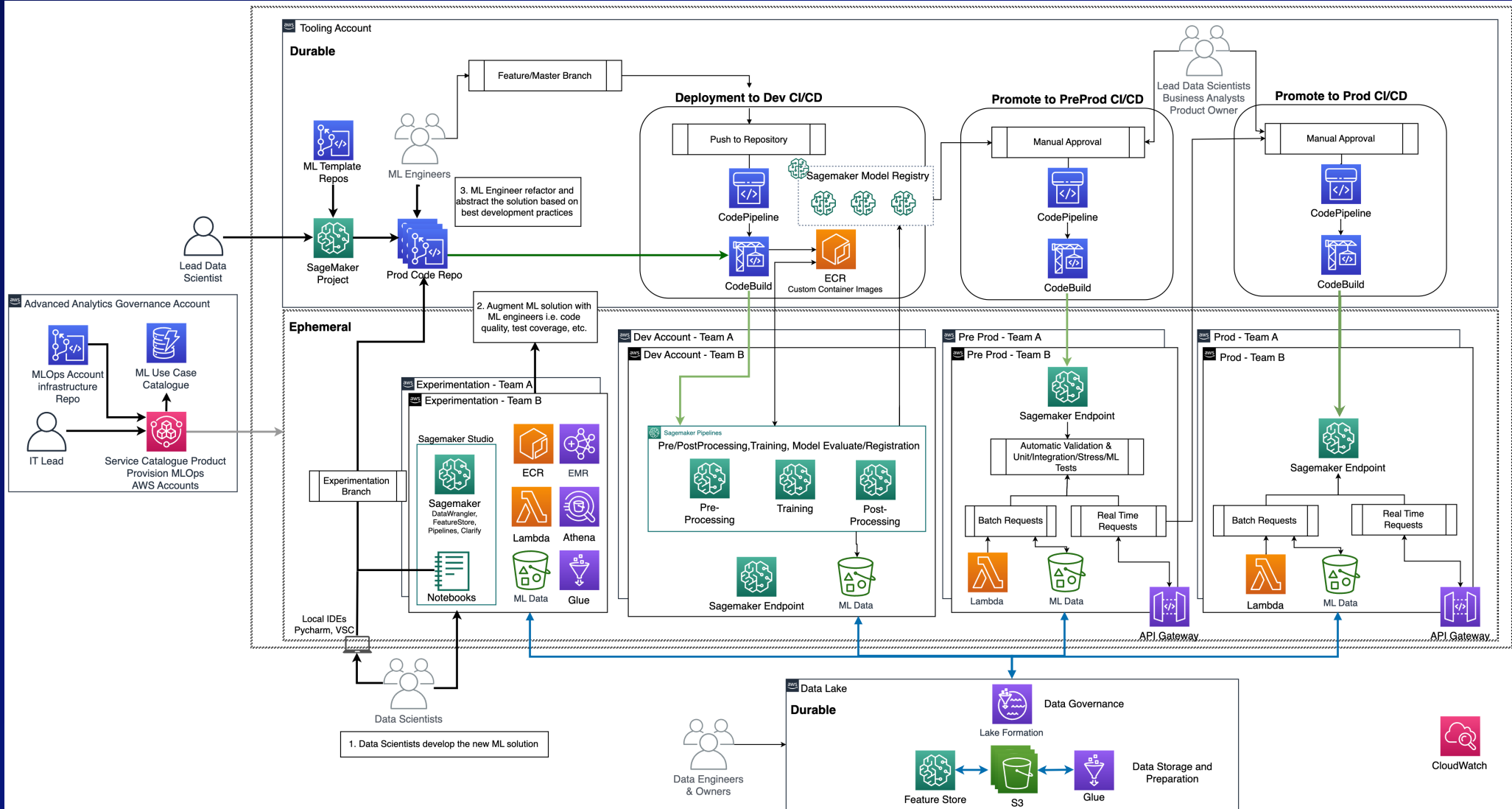
MLOps Foundation **People & Processes**

SEPARATION OF CONCERNS IS KEY FOR SUCCESS



MLOPs Scalable Phase

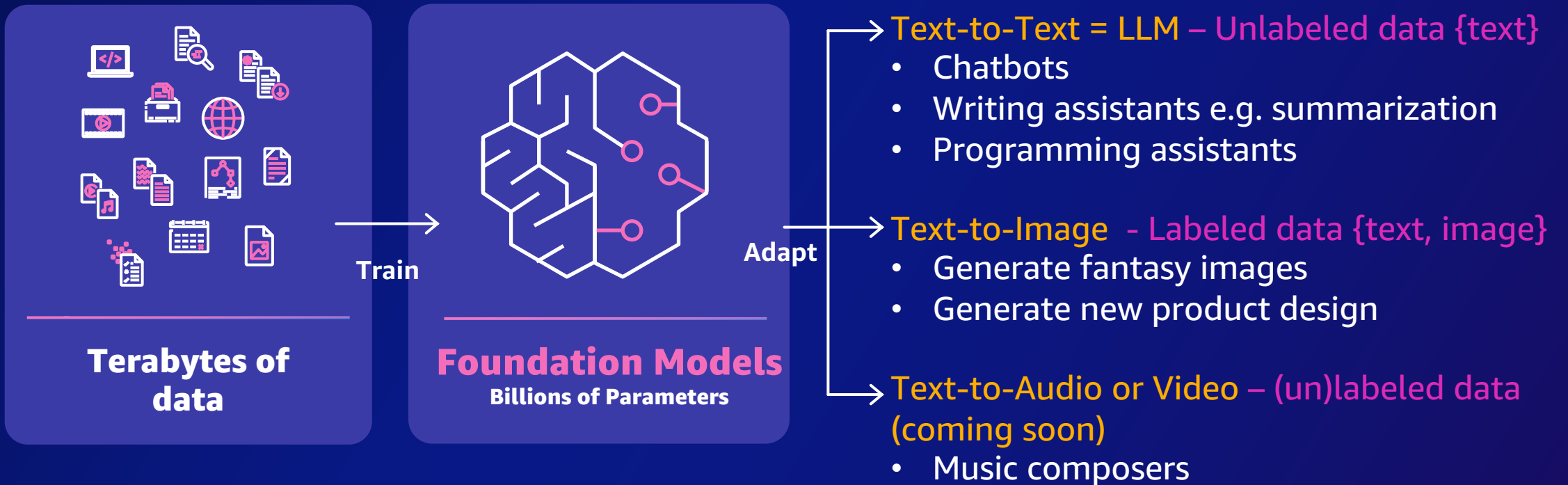
MULTIPLE TEAMS AND ML USE CASES ADOPT MLOPS



Generative AI (GenAI) & MLOps

MLOps & FMOPs/LLMOPs Differentiators

GenAI Use Case Domains



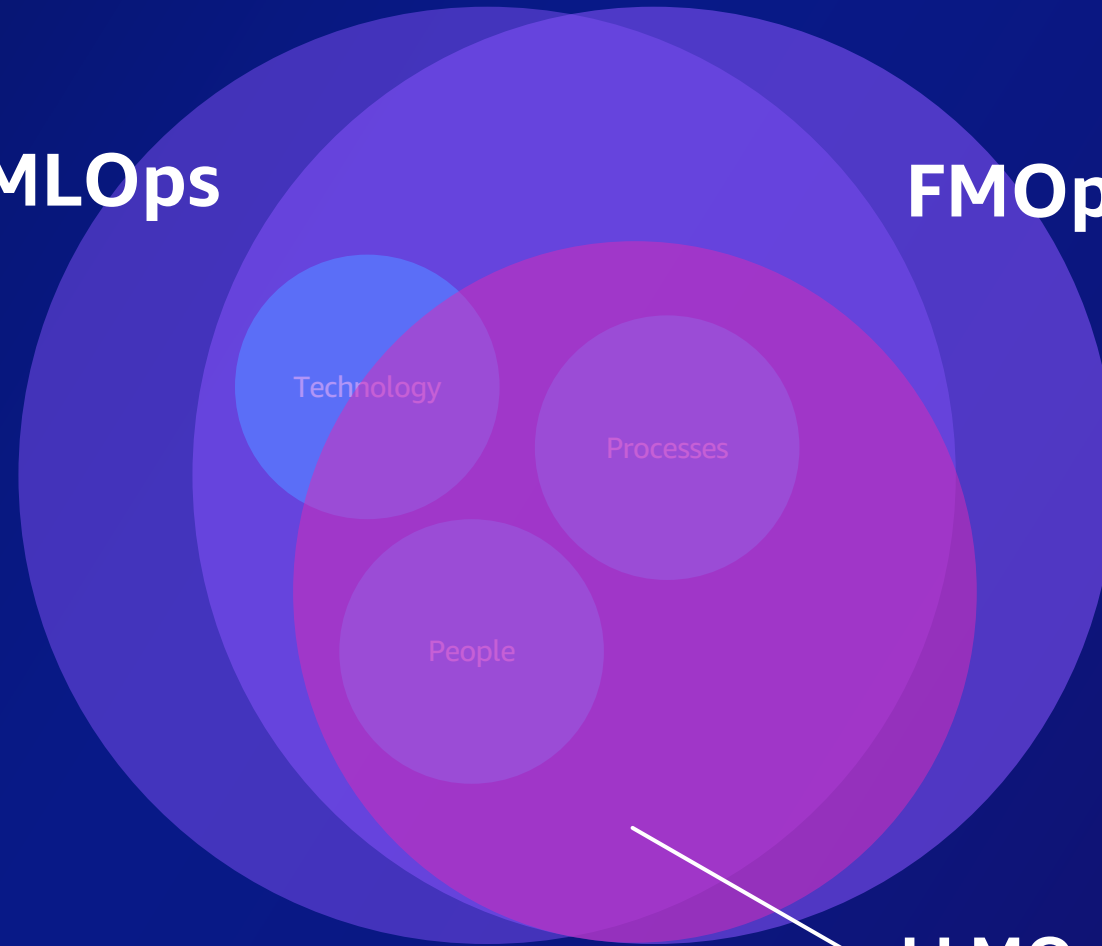
Key Definitions

Machine Learning Operations
Productionize ML solutions
efficiently

MLOps

FMOps

Foundation Model Operations
Productionize GenAI Solutions
(Text-Text/ Image/ Video/ Audio/
...)



LLMOps

Large Language Model Operations
Productionize Large Language
Model-based solutions

MLOps & FMOps Differentiators

MLOps

Technology

People

Processes

FMOps

Processes & People
Providers, fine-tuners, & consumers

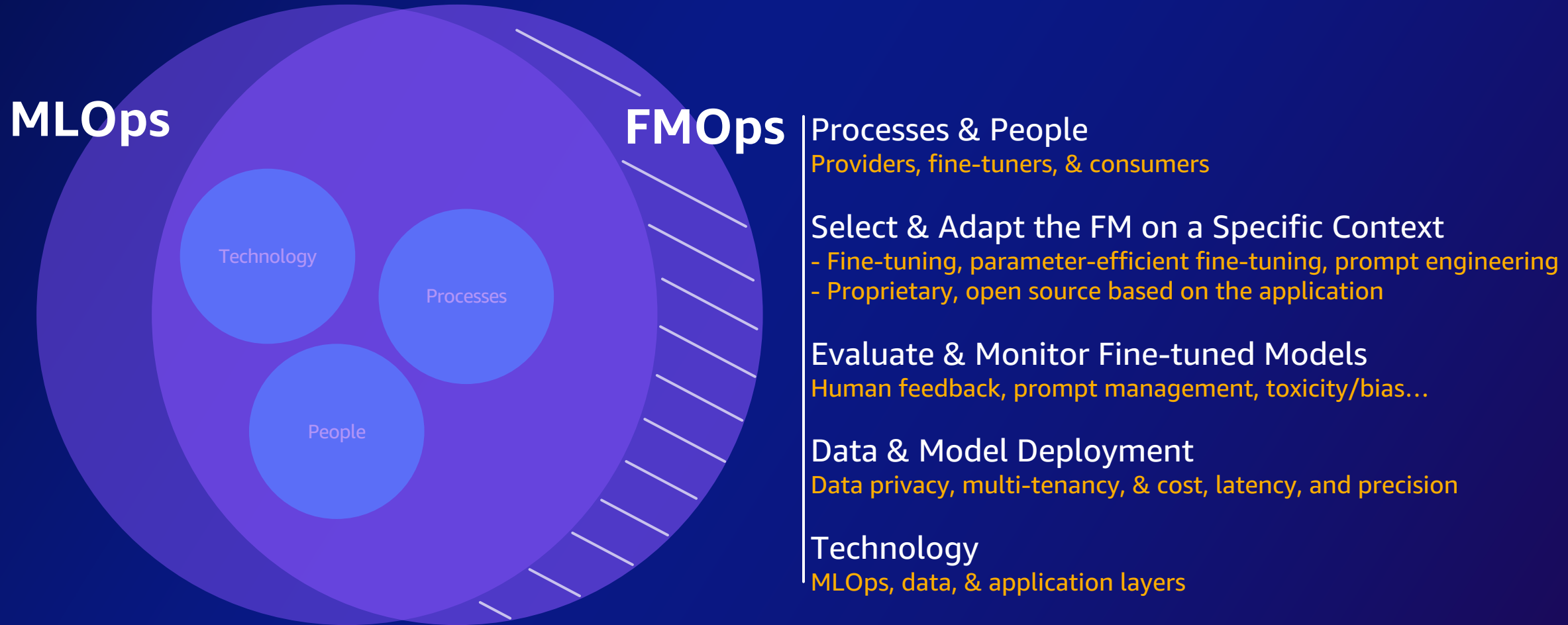
Select & Adapt the FM on a Specific Context
- Fine-tuning, parameter-efficient fine-tuning, prompt engineering
- Proprietary, open source based on the application

Evaluate & Monitor Fine-tuned Models
Human feedback, prompt management, toxicity/bias...

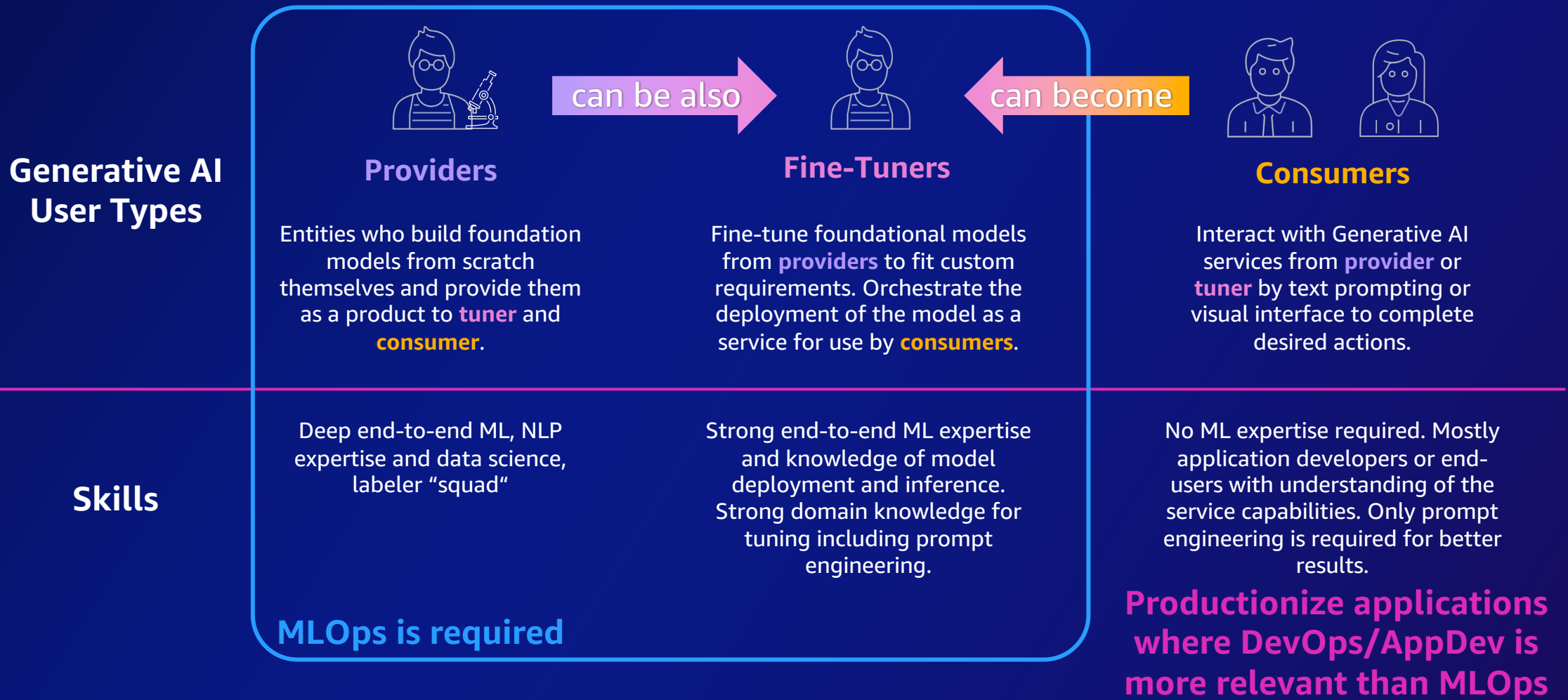
Data & Model Deployment
Data privacy, multi-tenancy, & cost, latency, and precision

Technology
MLOps, data, & application layers

MLOps & FMOps Differentiators

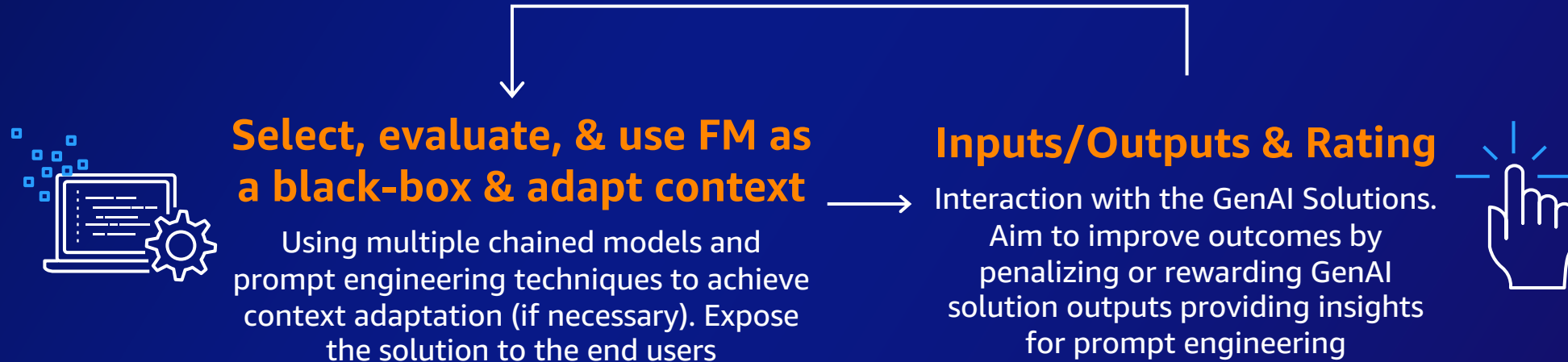


GenAI User Types & Skills

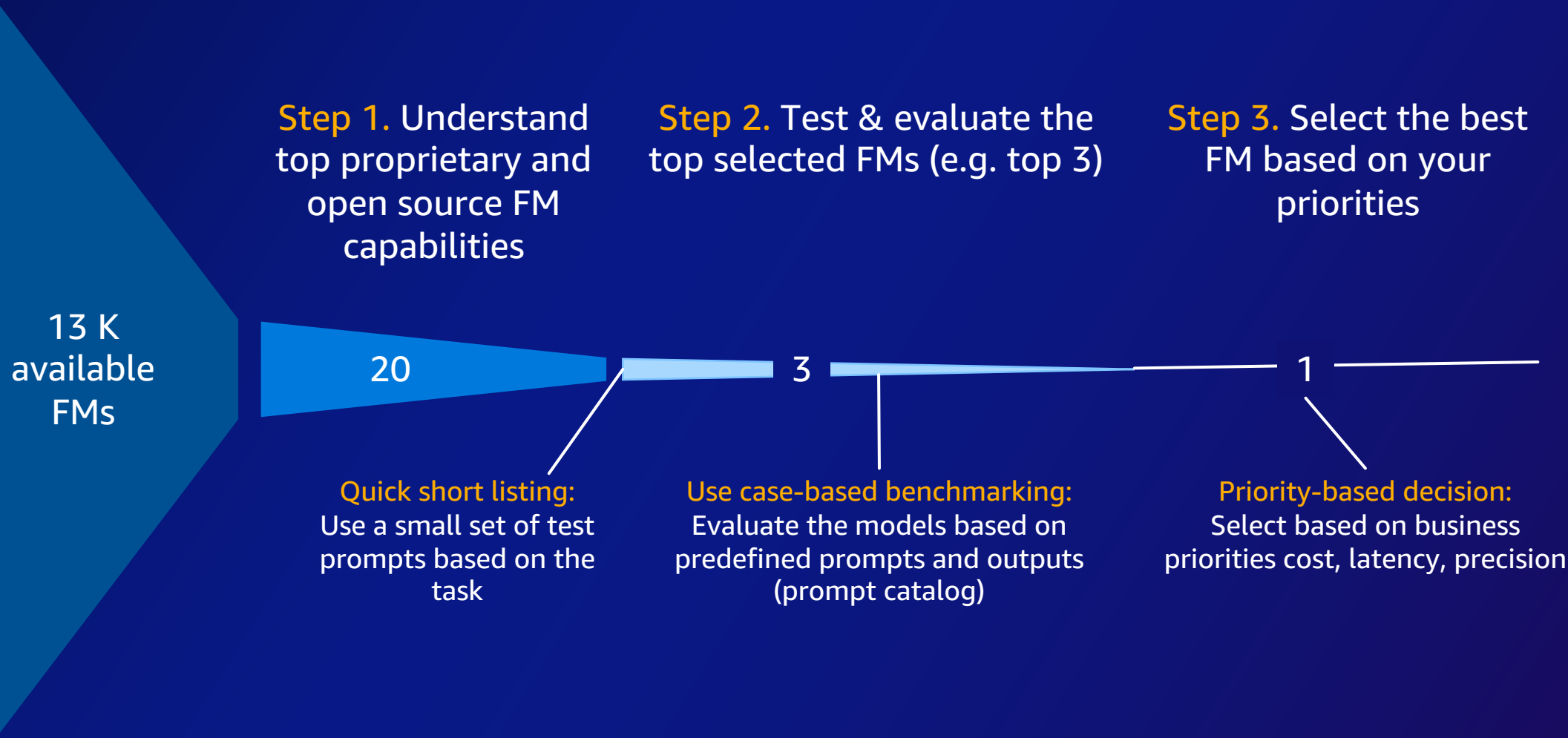


The Journey of Consumers

GenAI Processes - Consumers



Select FM - Consumers



Step 1. Understand top proprietary and open source FM capabilities

Step 2. Test & evaluate the top selected FMs (e.g. top 3)

Step 3. Select the best FM based on your priorities

13 K available FMs

20

3

1

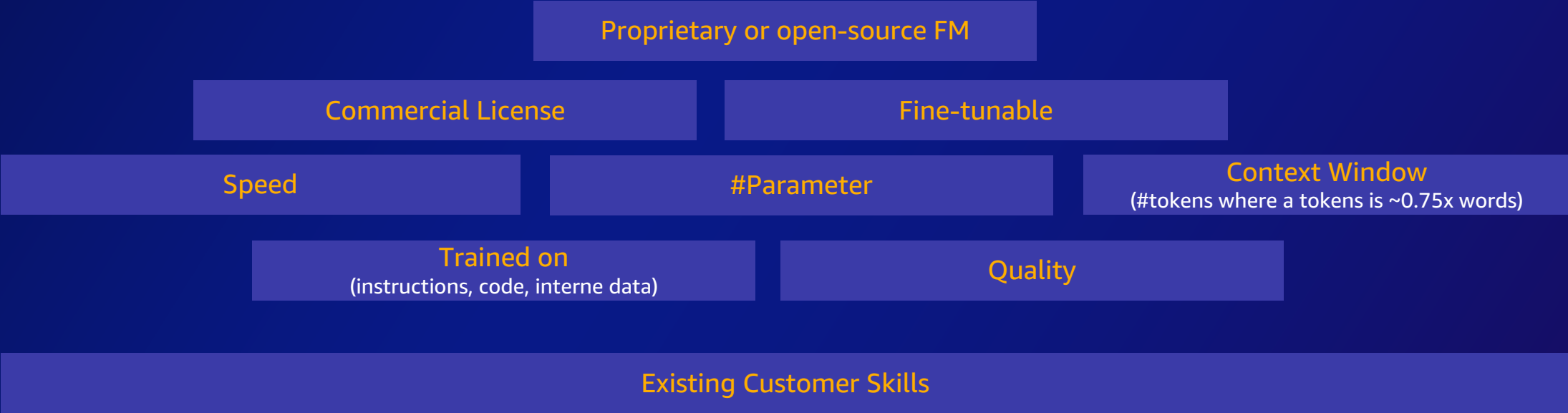
Quick short listing:
Use a small set of test prompts based on the task

Use case-based benchmarking:
Evaluate the models based on predefined prompts and outputs (prompt catalog)

Priority-based decision:
Select based on business priorities cost, latency, precision

Step 1. Understand top FM capabilities

Main FM Capability Matrix



Step 1. **Proprietary** FM Capabilities

Company Name	Model Name	Can be used Commercially	# Params	GPU instance req.	Available on AWS	Speed	Context Window	Trained on	Fine-tunable
AI21	J2 Ultra Instruct	Yes	178 B	p4d.24xl	Bedrock, Jumpstart/SM		8 K	Internet Data, Code, Instructions	No
	J2 Mid Instruct	Yes	17 B	g5.12xl	Bedrock, Jumpstart/SM		8 K	Internet Data, Code, Instructions	No
	AI21 Summarize	Yes		g4dn.12xl	Jumpstart/SM		~13 K	Internet Data, Instructions	No
Amazon	Titan Text Large	Yes	n/a	n/a	Bedrock		4 K	n/a	No
Anthropic	Claude	Yes	n/a	n/a	Bedrock		12 K	Internet Data, Code, Instructions, Human feedback	No
Cohere	Generate Model Command	Yes	n/a (50 B)	n/a	Jumpstart/SM		4 K	Internet Data, Instructions	No
	Generate Model Command-Light	Yes	n/a (6 B)	n/a	Jumpstart/SM		4 K	Internet Data, Instructions	No
LightOn	Lyra-Fr 10B	Yes	10 B	g5.12xl	Jumpstart/SM		?	Internet Data (French)	No
Stability AI	SDXL	Yes	n/a	g5.xl	Bedrock, Jumpstart/SM		-	<Text, Image>	No

Step 1. **Open-source** FM Capabilities

Company Name	Model Name	Can be used Commercially	# Params	GPU instance req.	Available on AWS	Speed	Context Window	Trained on	Fine-tunable
Google	FLAN-UL2	Yes	20 B	g5.12xl	Jumpstart/SM		2 K	Internet Data, Code, Instructions	Yes
	FLAN-T5-XXL	Yes	11 B	g5.xl	Jumpstart/SM		512	Internet Data, Code, Instructions	Yes
Eleuther	GPT-J	Yes	6 B	g5.xl	Jumpstart/SM		512	Internet Data, Code	Yes
TII	Falcon-40B-Instruct	Yes	40 B	g5.12xl	Jumpstart/SM		2 K	Internet Data, Code, Instructions	Yes
	Falcon-7B-Instruct	Yes	7 B	g5.xl	Jumpstart/SM		2 K	Internet Data, Code, Instructions	Yes
BigCode	StarCoder	Yes	15 B	g5.12xl	SM		8 K	Code	Yes
	Santa Coder	Yes	1.1 B	g5.xl	SM		2K	Code	Yes
LMSYS Org	Vicuna-13B	No	13 B	g5.xl	SM		2 K	Internet Data, Code, Instructions	Yes
Meta	Llama-65B	No	65 B	g5.48xl	SM		2 K	Internet Data, Code	Yes
Stability AI	SD 2.1	Yes	-	g5.xl	Jumpstart/SM		-	<Text, Image>	Yes

Step 1. EU AI Act Matters for FM Selection

Grading Foundation Model Providers' Compliance with the Draft EU AI Act

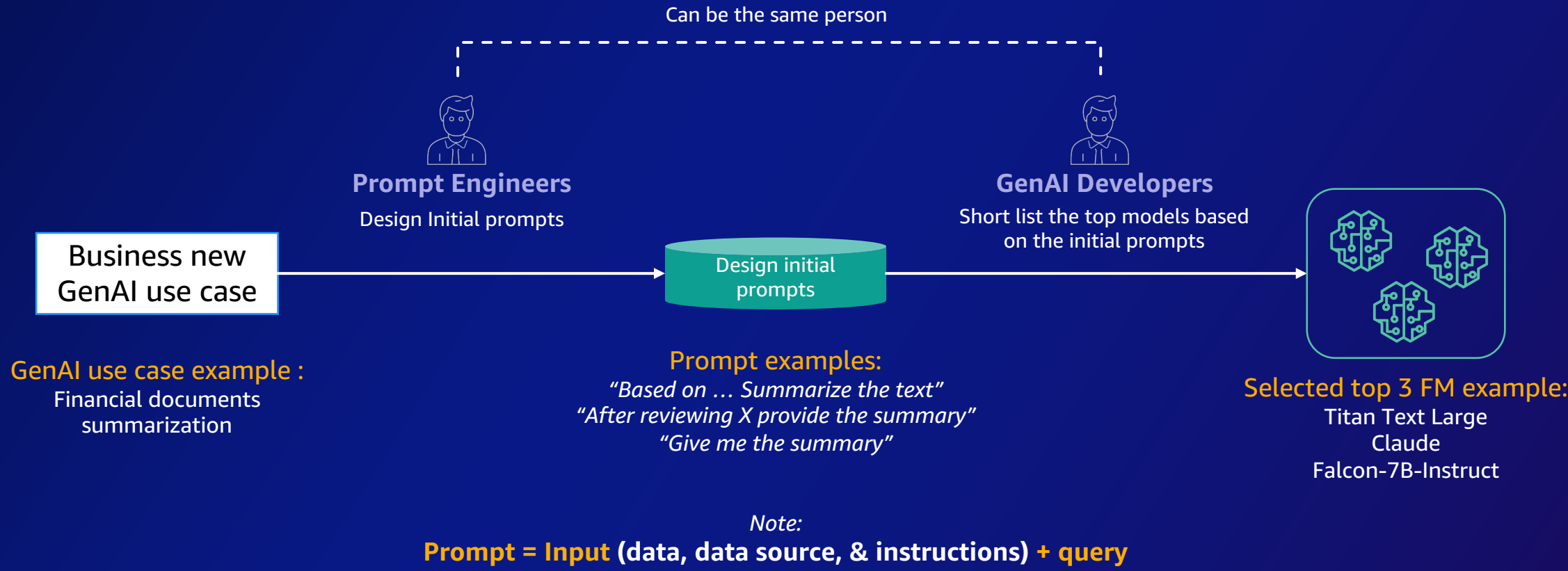
Source: Stanford Research on Foundation Models (CRFM), Institute for Human-Centered Artificial Intelligence (HAI)

	OpenAI	cohere	stability.ai	ANTHROPIC	Google	BigScience	Meta	AI21labs	ALEPH ALPHA	EleutherAI	Totals
Draft AI Act Requirements	GPT-4	Cohere Command	Stable Diffusion v2	Claude	PaLM 2	BLOOM	LLaMA	Jurassic-2	Luminous	GPT-NeoX	
Data sources	● ○ ○ ○ ○	● ● ● ● ○	● ● ● ● ●	○ ○ ○ ○ ○	● ● ● ● ○	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	22
Data governance	● ● ● ○ ○	● ● ● ● ○	● ● ● ○ ○	○ ○ ○ ○ ○	● ● ● ● ○	● ● ● ● ●	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ○	19
Copyrighted data	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	7
Compute	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	● ○ ○ ○ ○	● ● ● ● ●	17
Energy	○ ○ ○ ○ ○	● ○ ○ ○ ○	● ● ● ● ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ●	16
Capabilities & limitations	● ● ● ● ●	● ● ● ● ○	● ● ● ● ●	● ● ● ○ ○	● ● ● ● ●	● ● ● ● ○	● ● ● ○ ○	● ● ● ○ ○	● ● ● ○ ○	● ● ● ● ○	27
Risks & mitigations	● ● ● ● ○	● ● ● ○ ○	● ○ ○ ○ ○	● ○ ○ ○ ○	● ● ● ● ○	● ● ● ● ○	● ● ● ○ ○	● ● ● ○ ○	○ ○ ○ ○ ○	● ○ ○ ○ ○	16
Evaluations	● ● ● ● ●	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ○	● ● ● ● ●	● ● ● ○ ○	○ ○ ○ ○ ○	● ○ ○ ○ ○	● ● ● ○ ○	15
Testing	● ● ● ● ○	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ○	● ● ● ● ○	○ ○ ○ ○ ○	● ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	10
Machine-generated content	● ● ● ● ○	● ● ● ● ○	○ ○ ○ ○ ○	● ● ● ● ●	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	● ● ● ● ○	● ○ ○ ○ ○	● ● ● ○ ○	21
Member states	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ○ ○	● ● ● ● ●	○ ○ ○ ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ○ ○ ○ ○	○ ○ ○ ○ ○	9
Downstream documentation	● ● ● ● ○	● ● ● ● ●	● ● ● ● ●	○ ○ ○ ○ ○	● ● ● ● ●	● ● ● ● ●	● ● ● ○ ○	○ ○ ○ ○ ○	○ ○ ○ ○ ○	● ● ● ● ○	24
Totals	25 / 48	23 / 48	22 / 48	7 / 48	27 / 48	36 / 48	21 / 48	8 / 48	5 / 48	29 / 48	

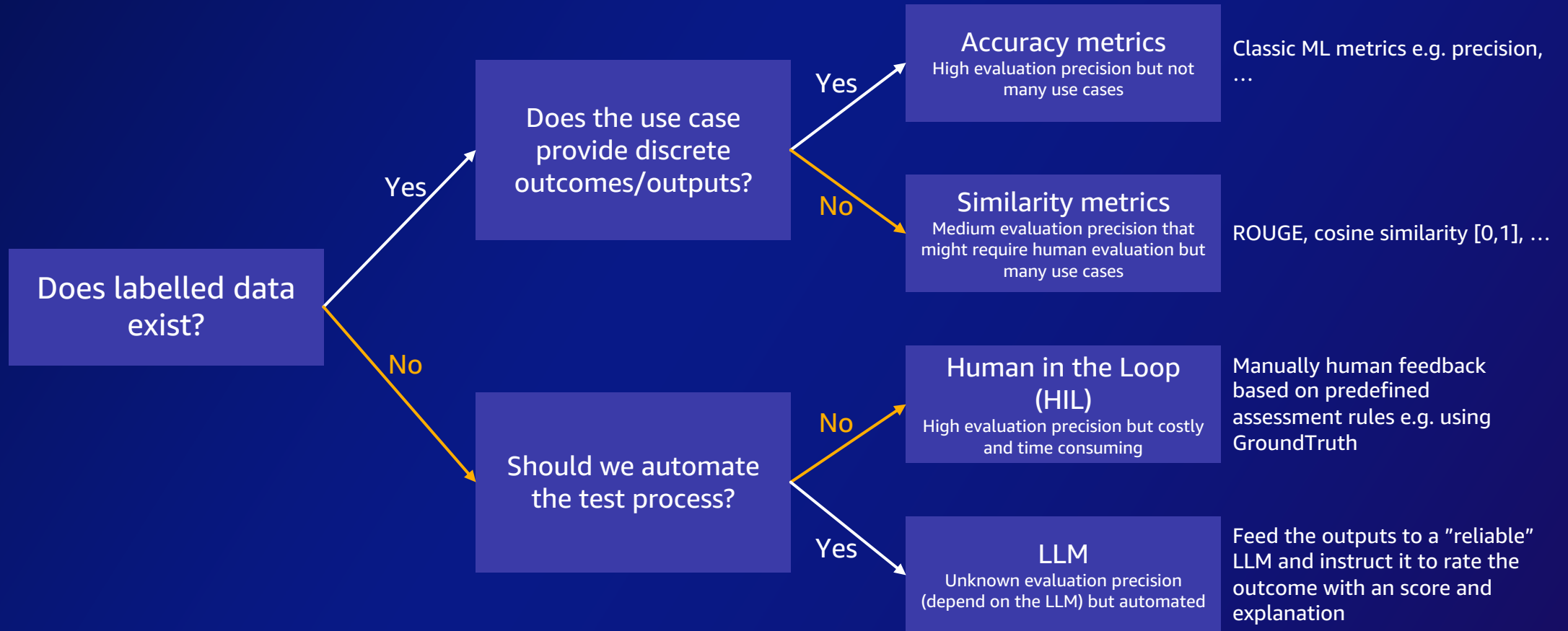
<https://crfm.stanford.edu/2023/06/15/eu-ai-act.html>



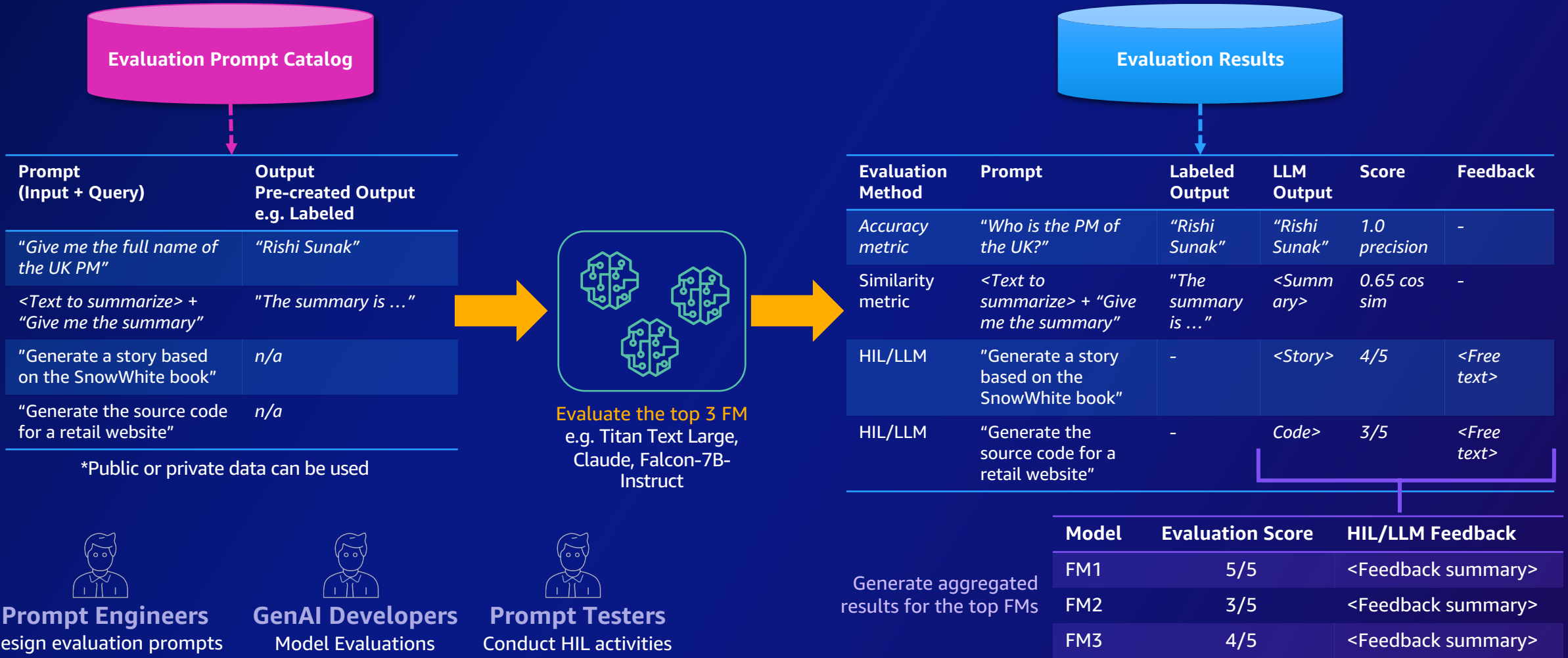
Step 1. Understand FM Capabilities



Step 2. Evaluate the top FMs



Step 2. Evaluate the top FMs - Examples



Step 3. Select the best FM based on priorities

EXAMPLE

Model	Speed
FM1	⚡⚡
FM2	⚡
FM3	⚡

Speed

No priority

High speed, smaller model,
lower precision, smaller cost

Model Selection:
FM2

P1: Precision

Precision

Model	Evaluation Score	HIL/LLM Feedback
FM1	5/5	<Feedback summary>
FM2	4/5	<Feedback summary>
FM3	3/5	<Feedback summary>

Lower speed, larger model,
higher precision, larger cost

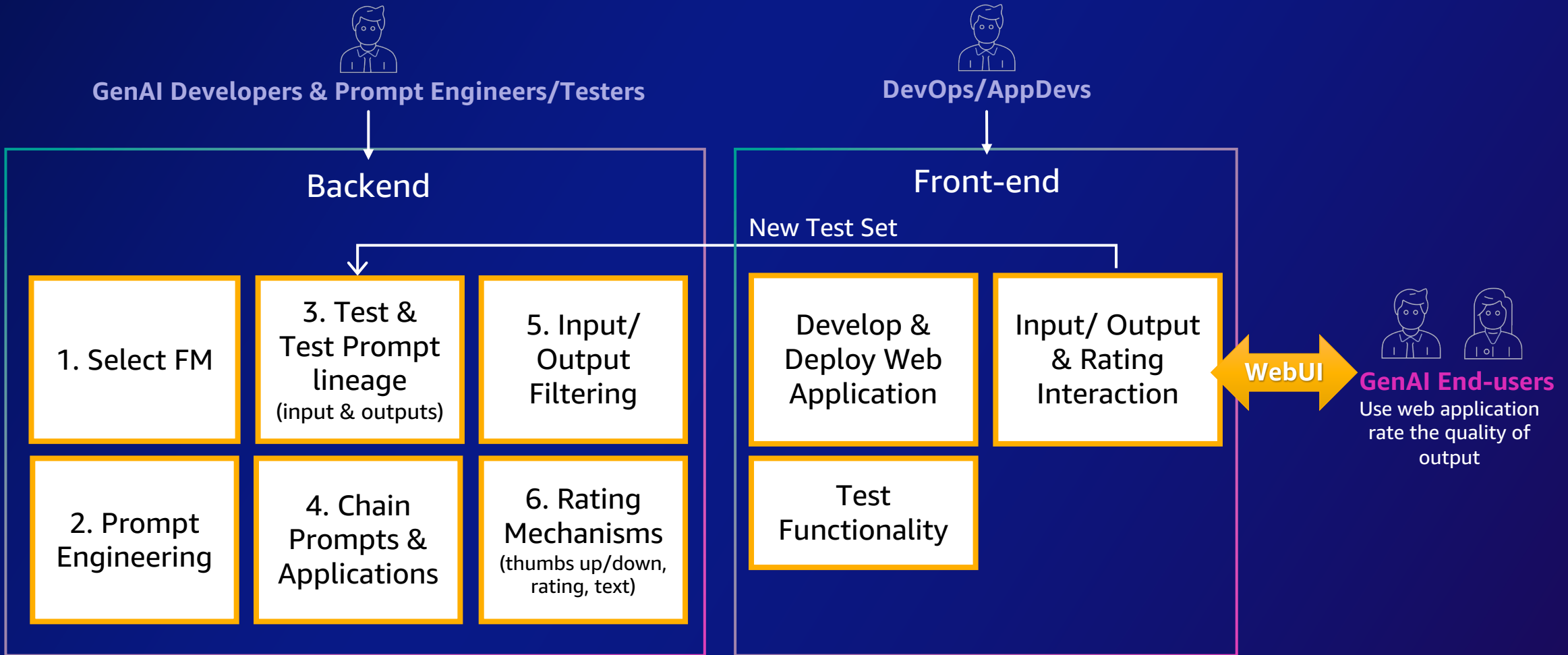
Cost

P0: lower cost

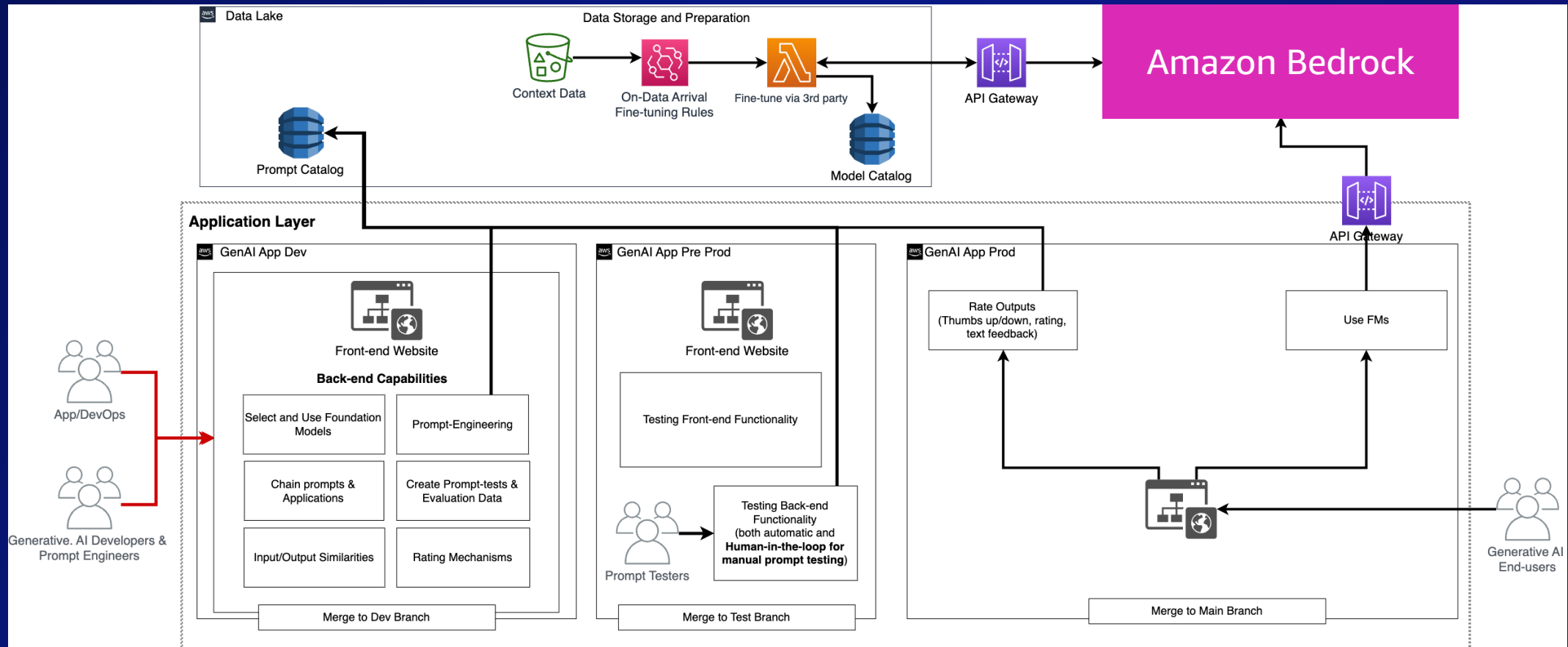
Model	Cost
FM1	\$\$\$\$
FM2	\$
FM3	\$\$\$

GenAI Processes for LLM - Consumers

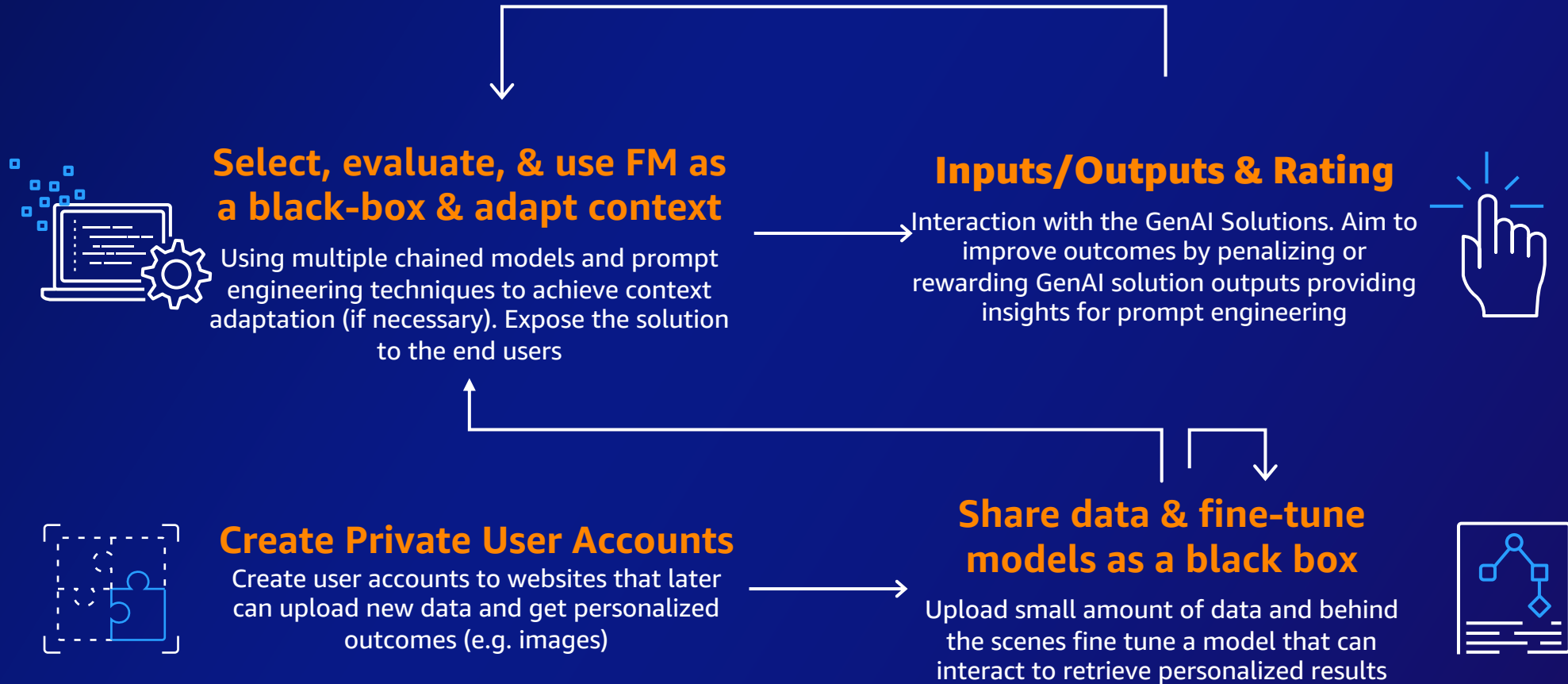
LLM-based GenAI Solution



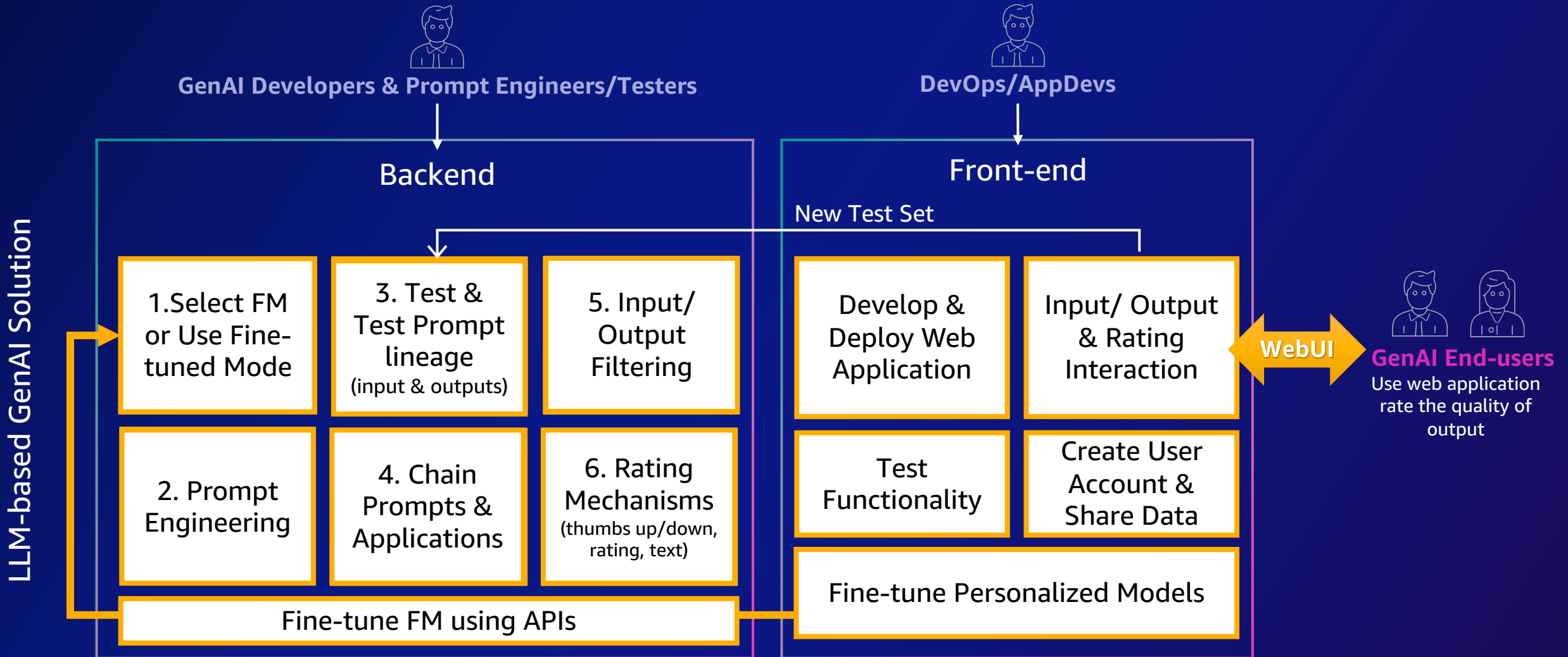
GenAI Technology for LLM - Consumers



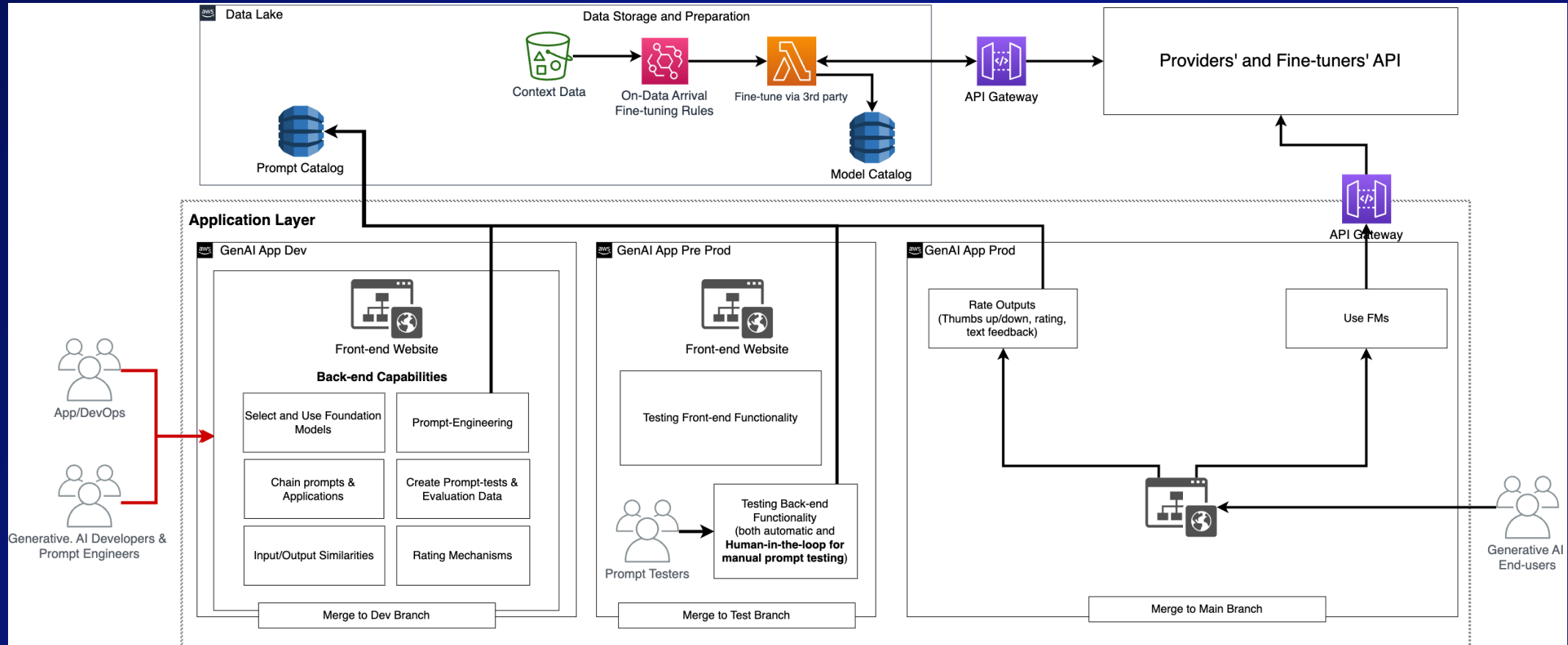
GenAI Processes - Consumers



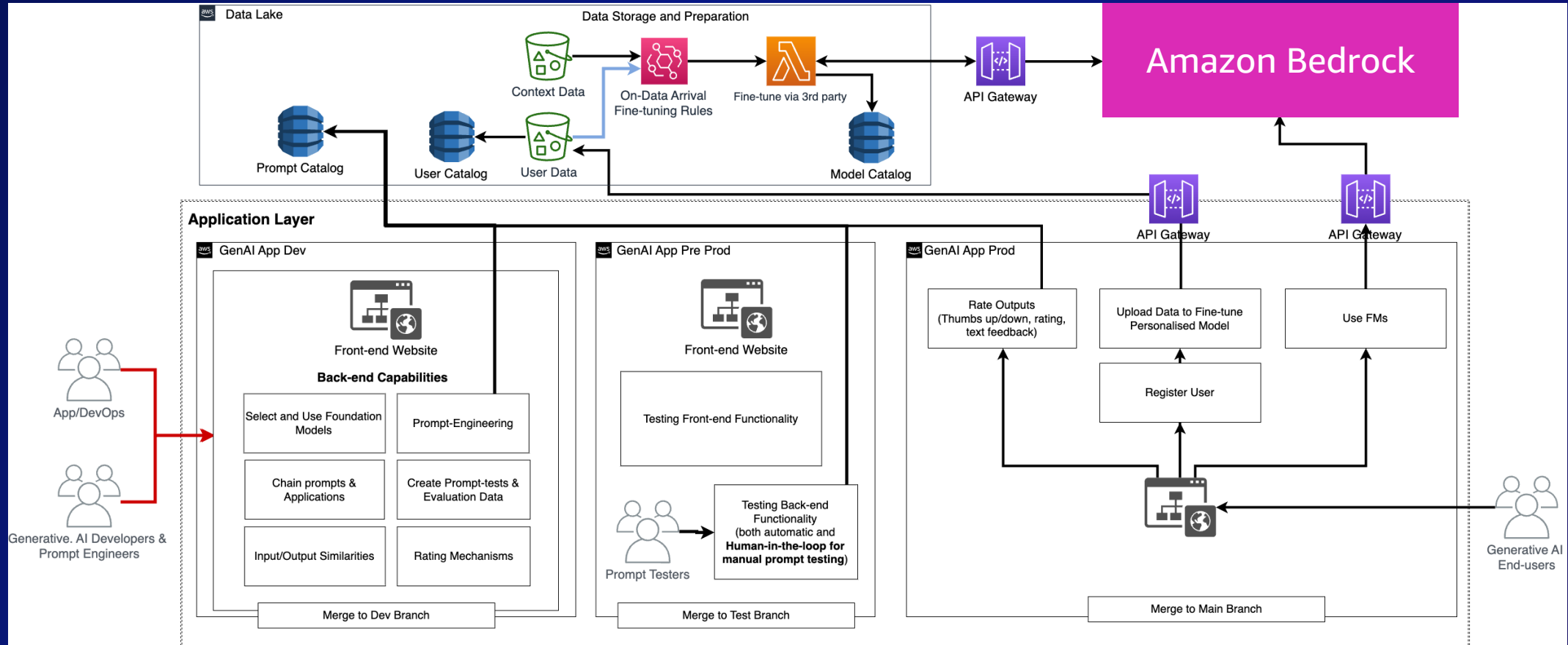
GenAI Processes for LLM - Consumers



GenAI Technology for LLM - Consumers



GenAI Technology for LLM - Consumers

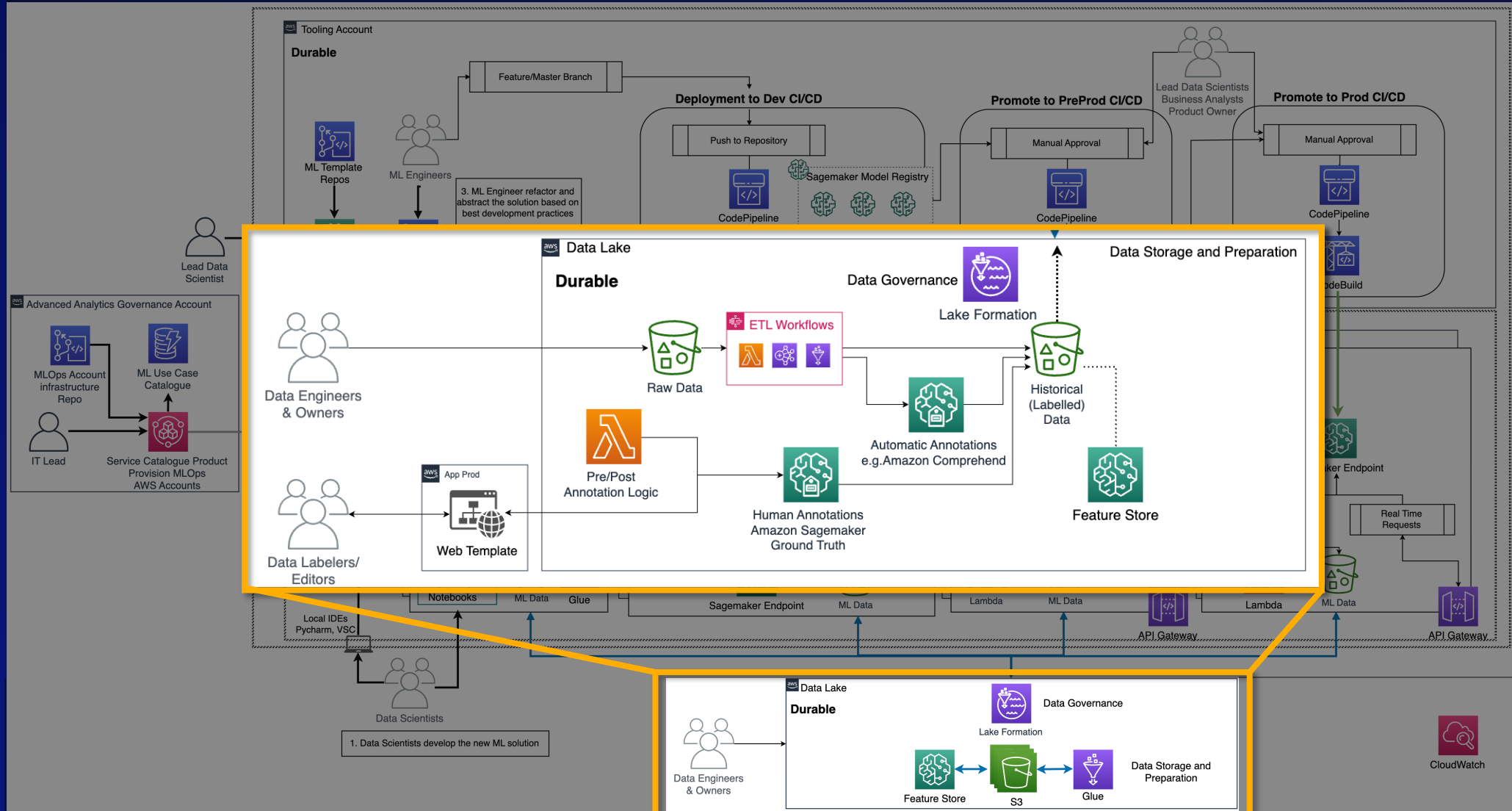


The Journey of Providers



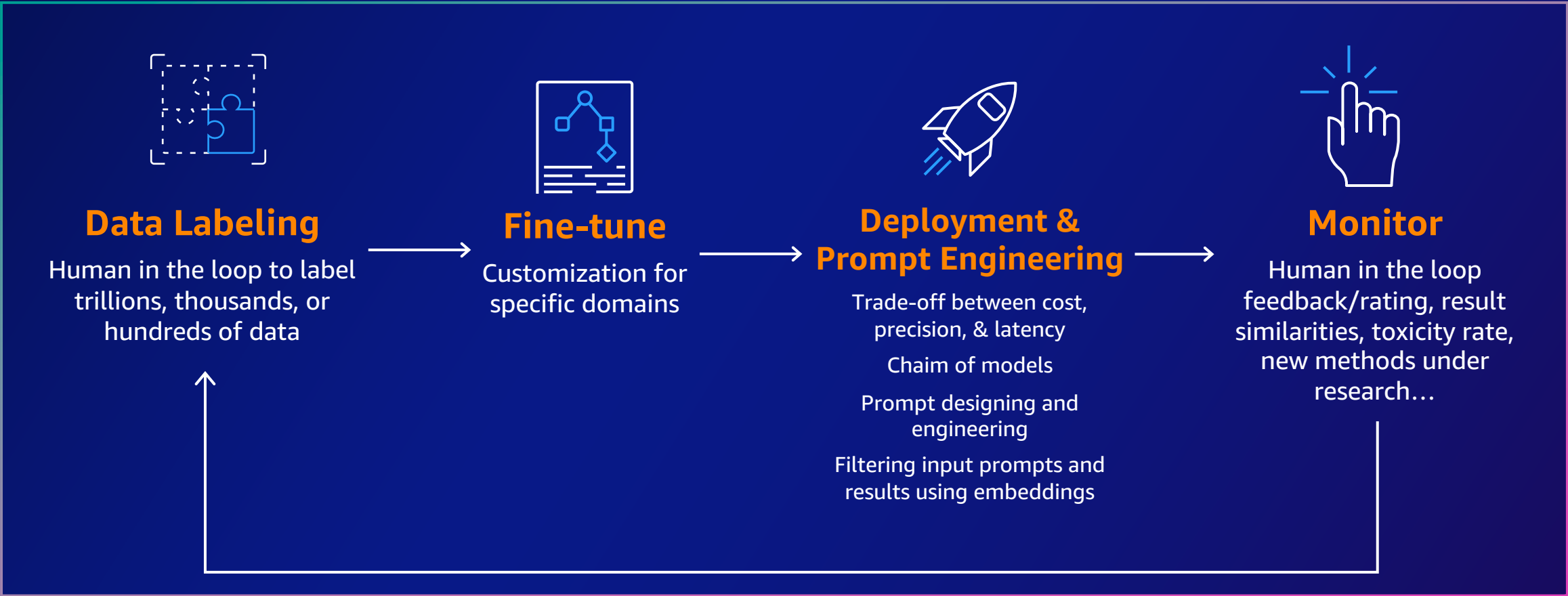
GenAI Providers Productionize FM using MLOps

TRAIN MULTIPLE FOUNDATIONS MODELS



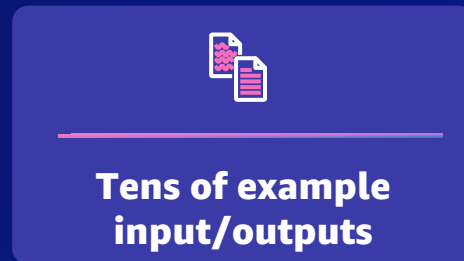
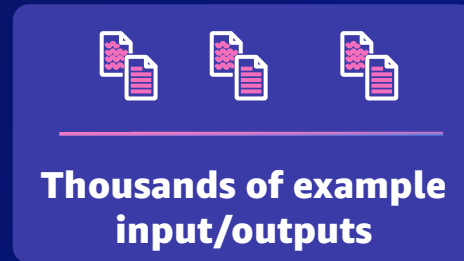
The Journey of **Fine-tuners**

GenAI Processes - Fine-Tuners



Fine-Tuning, PEFT & Training

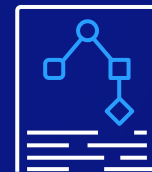
Select the
'right' model



Fine-Tuning (Training Job)

Requires **high computation power (high cost but lower than training)** to calculate all the weights of Large FM (deep learning model)

Higher accuracy

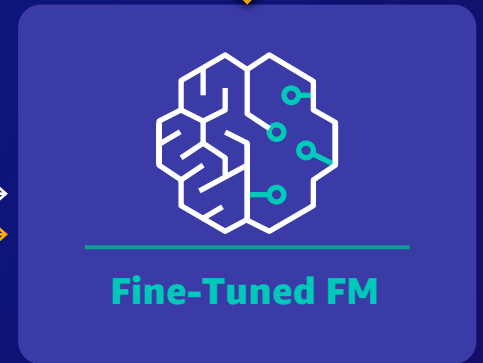


Parameter-Efficient Fine-Tuning (PEFT - Training Job)

Requires **low computation power (reduced cost e.g. 1/10 of fine-tuning)** as it adds small new layers in the Large FM (deep learning model)

Lower accuracy

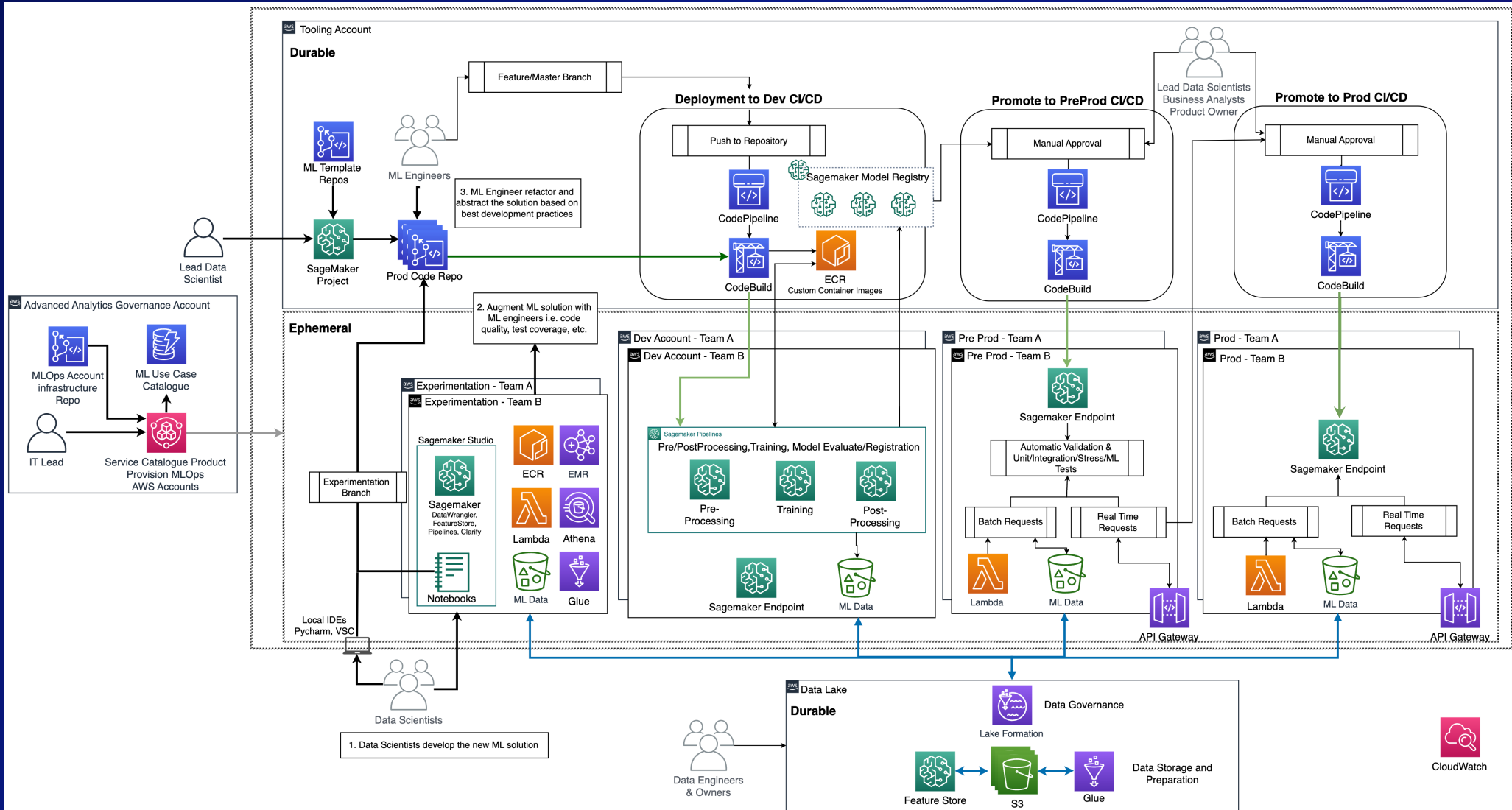
**Prompt
(Inputs)**



**Completion
(Outputs)**

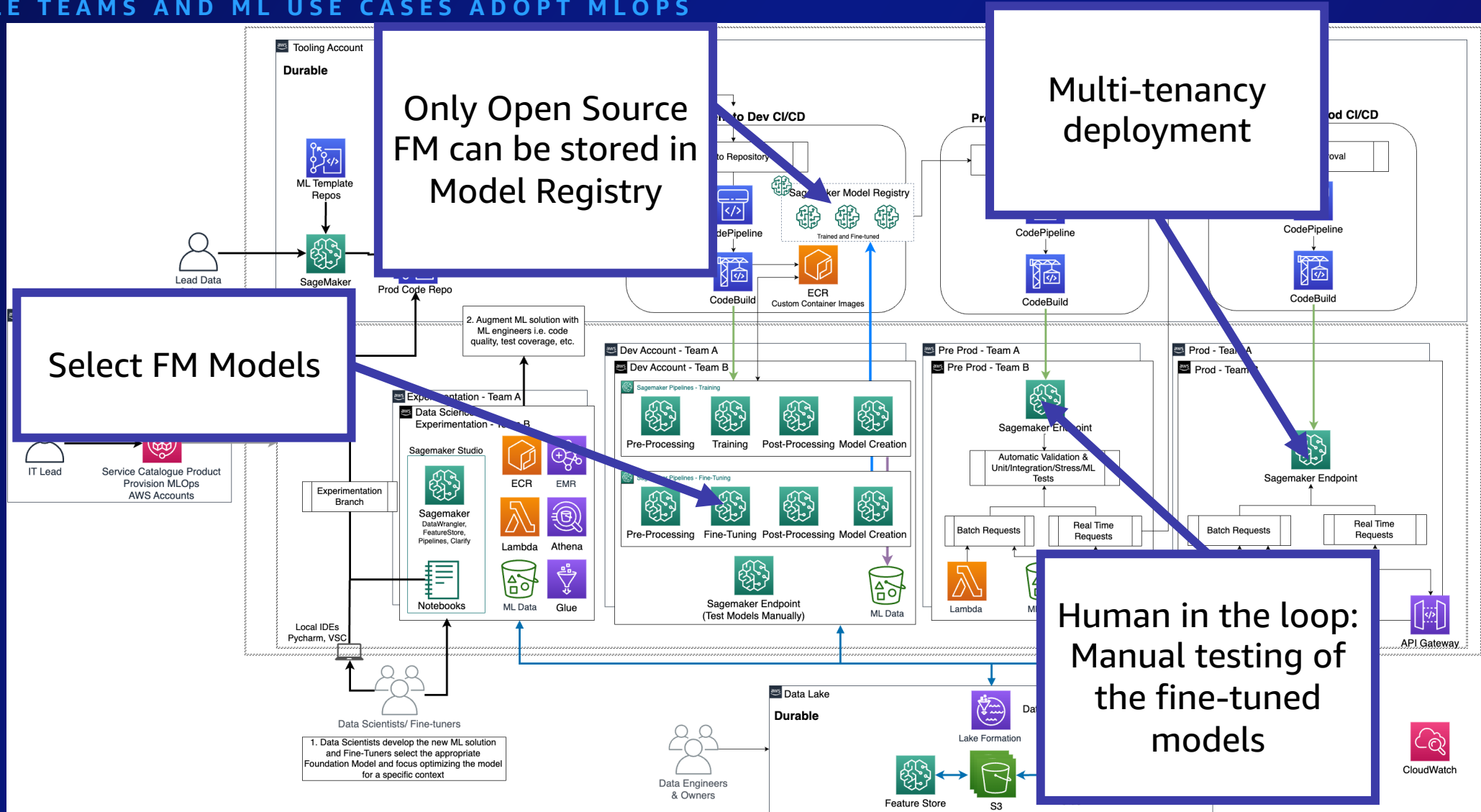
MLOPs Technology

MULTIPLE TEAMS AND ML USE CASES ADOPT MLOPS

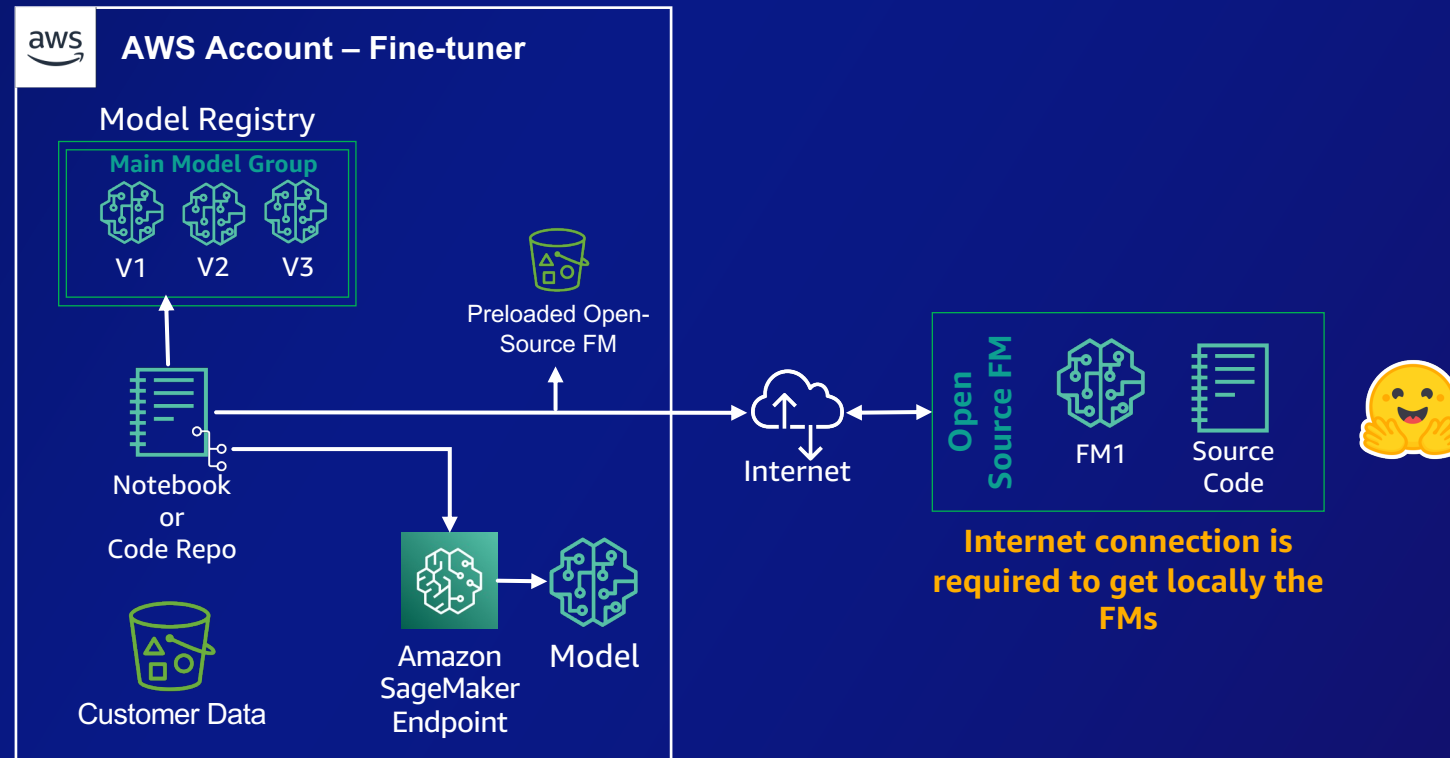


MLOPs & GenAI Technology - Fine-tuner

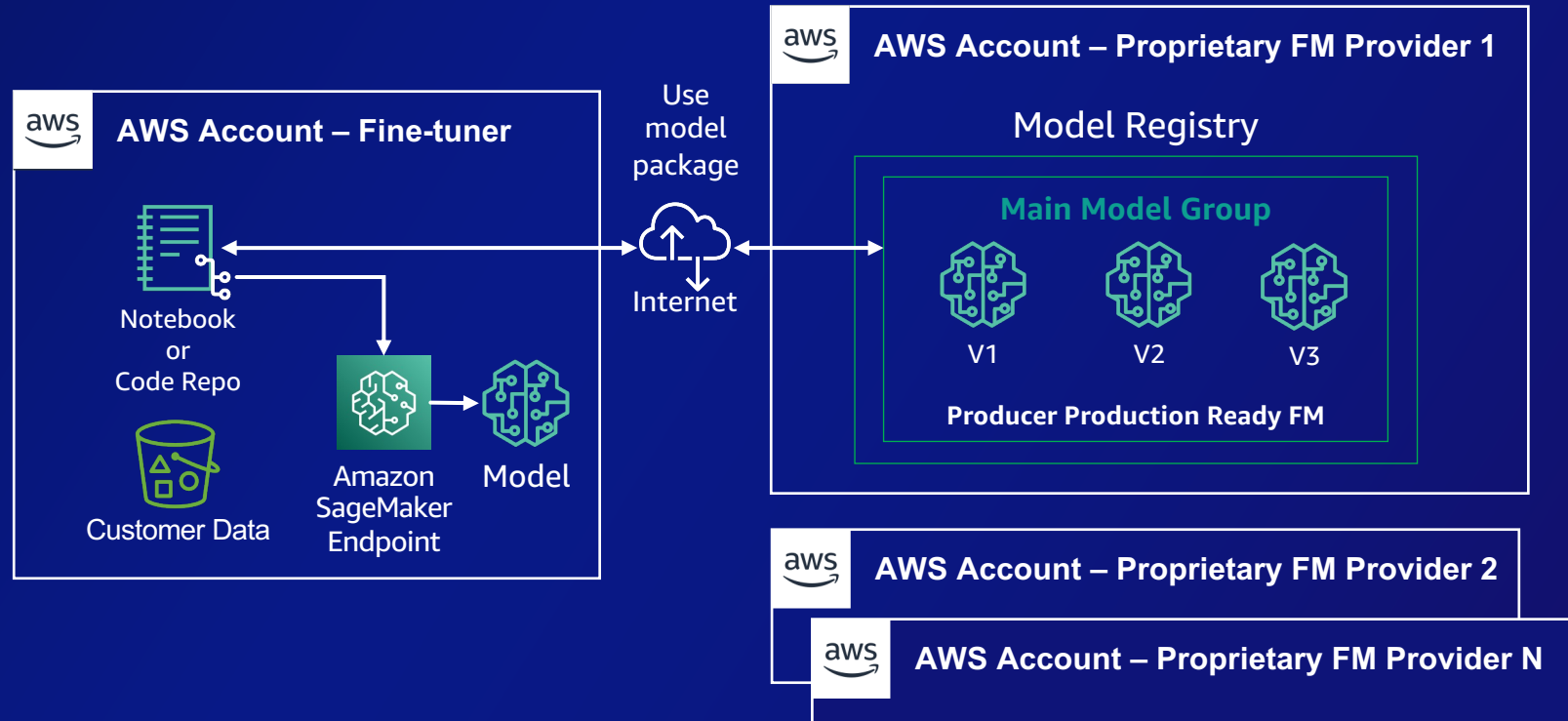
MULTIPLE TEAMS AND ML USE CASES ADOPT MLOPs



Data & Open Source Fine-tuned FM Deployment

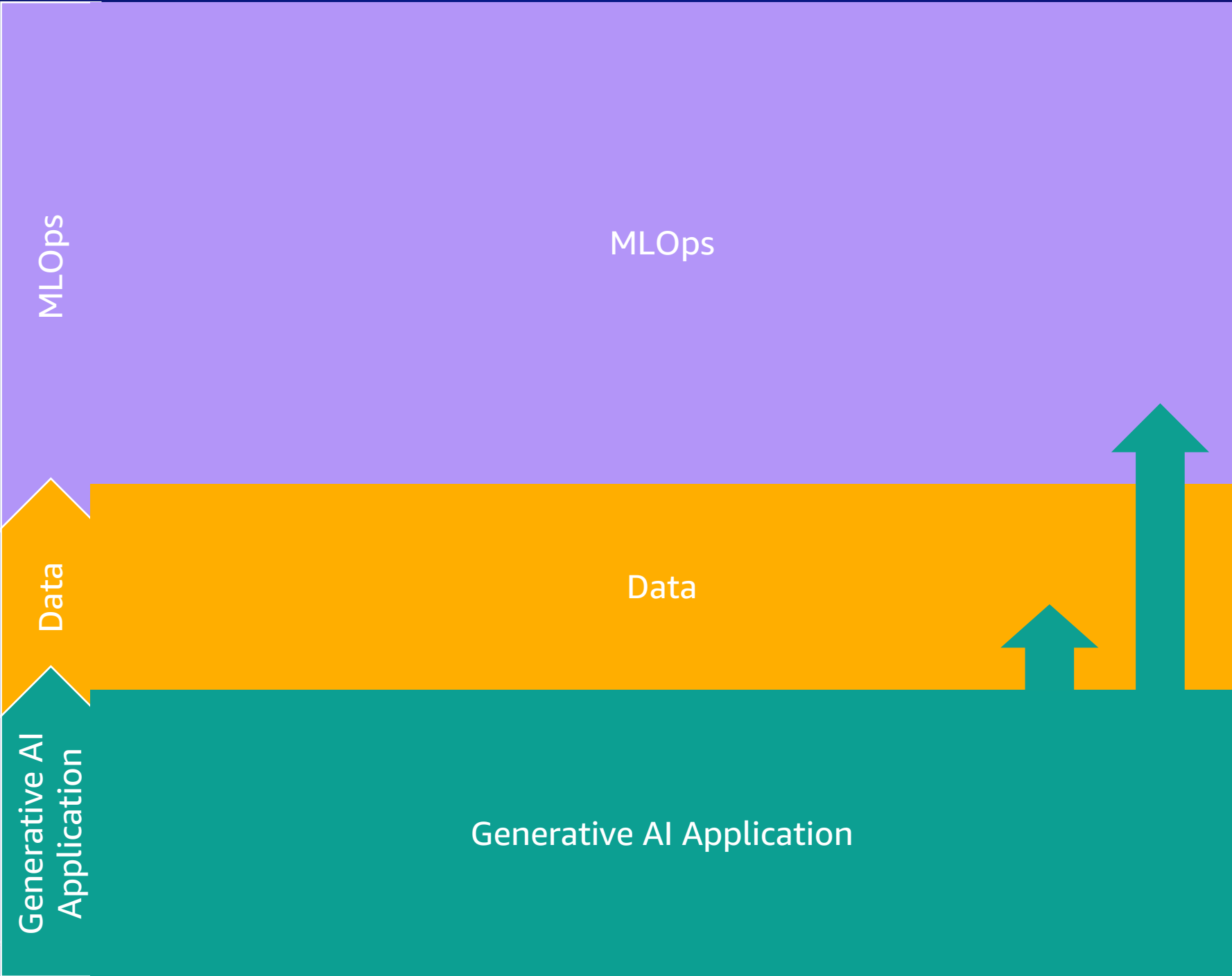


Data & Proprietary Fine-tuned FM Deployment



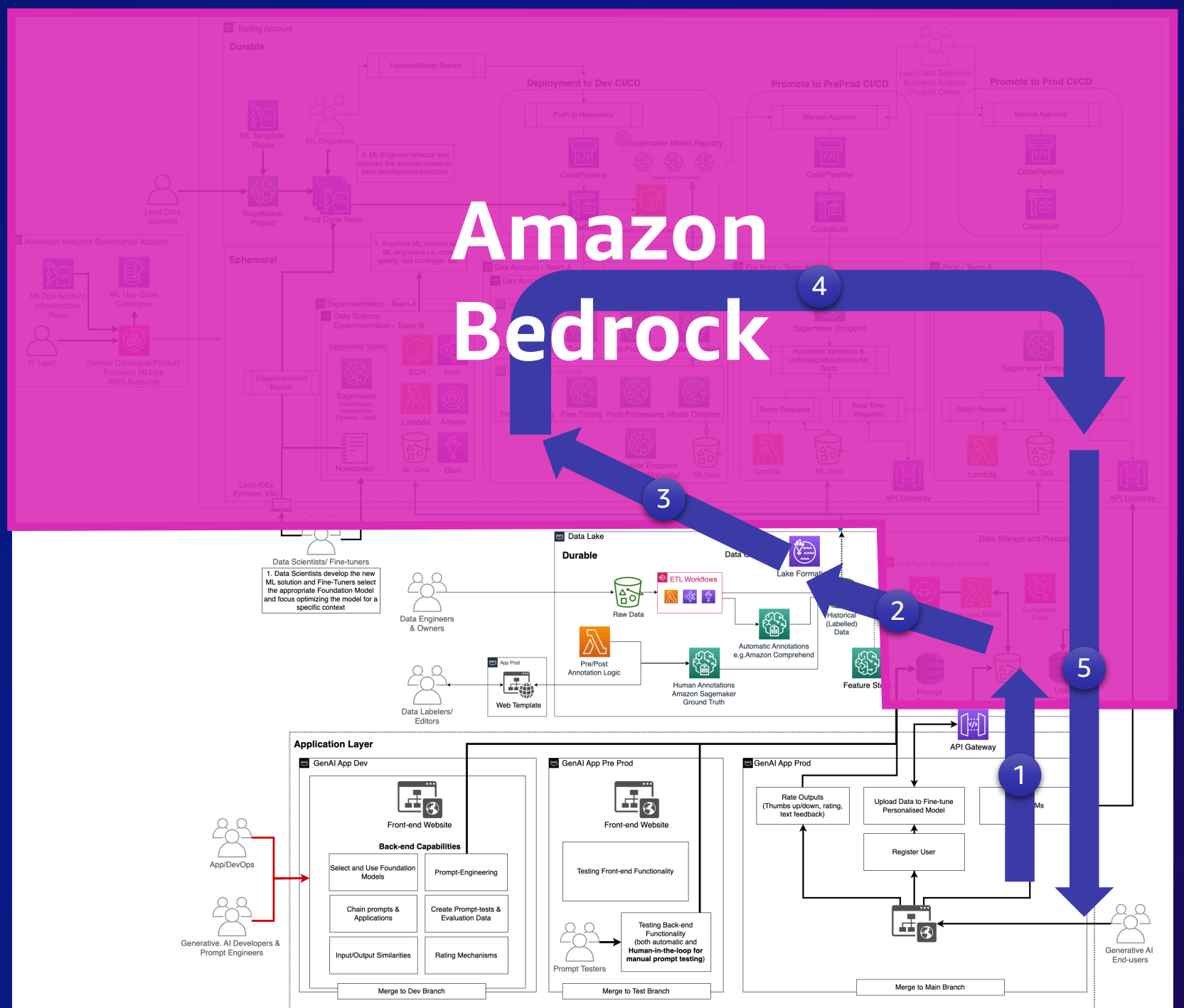
MLOPs & Generative AI Technology – Fine-tuner

THREE MAIN LAYERS ARE INTERCONNECTED



MLOPs & Generative AI Technology – Fine-tuner

THREE MAIN LAYERS ARE INTERCONNECTED



People & Processes



MLOps Key Personas and Roles

Advance Analytics Team Data Lake



Data Engineer

Prepare & Ingest data building ETL pipelines



Data Owners

Manage data sharing and provide access

Data Science Team Experimentation & MLOps



Data Scientist

Create the best ML models to solve business problems



ML Engineer

Collaborate with DS to productionize ML

Platform Team Secure Cloud/Data/ML Platform



MLOps Engineer/Admin

Standardize CI/CD, user/service role, model consumption, testing and deployment methodology



Security

Assess data, user, and service access creating policies and guardrails



Architects/ SysOps Engineer

Standardize account infrastructure, connectivity, user roles implementation

Business Viz Dashboards, ML Adoption, & ROI



Business Stakeholder Product Owners

Define business problem, business KPIs, and make business decisions



Business Stakeholder Data & ML Consumers

Consumers of ML results from other BUs, driving business decision making

Risk & Compliance Approve & Review Models



Auditors/Risk & Compliance

Review models, data sources, code artifacts

MLOps & FMOps Key Personas and Roles

Advance Analytics Team Data Lake



Data Engineer

Prepare & Ingest data building ETL pipelines



Data Owners

Manage data sharing and provide access

Data Science Team Experimentation & MLOps



Data Scientist

Create the best ML models to solve business problems



ML Engineer

Collaborate with DS to productionize ML

Platform Team Secure Cloud/Data/ML Platform



MLOps Engineer/Admin

Standardize CI/CD, user/service role, model consumption, testing and deployment methodology



Security & Architects

Assess data, user, and service access creating policies and infrastructure

Business Viz Dashboards, ML Adoption, & ROI



Business Stakeholder

Product Owners

Define business problem, business KPIs, and make business decisions



Business Stakeholder

Data & ML Consumers

Consumers of ML results from other BUs, driving business decision making

Risk & Compliance Approve & Review Models



Auditors/Risk & Compliance

Review models, data sources, code artifacts

MLOps & FMOps Key Personas and Roles

Advance Analytics Team Data Lake



Data Engineer

Prepare & Ingest data building ETL pipelines



Data Owners

Manage data sharing and provide access

Data Science Team Experimentation & MLOps



Data Scientist

Create the best ML models to solve business problems



ML Engineer

Collaborate with DS to productionize ML

Platform Team Secure Cloud/Data/ML Platform



MLOps Engineer/Admin

Standardize CI/CD, user/service role, model consumption, testing and deployment methodology



Security & Architects

Assess data, user, and service access creating policies and infrastructure

Business

Viz Dashboards, ML Adoption, & ROI



Business Stakeholder

Product Owners

Define business problem, business KPIs, and make business decisions



Business Stakeholder

Data & ML Consumers

Consumers of ML results from other BUs, driving business decision making

Labeler Team Data Preparation at Scale



Data Labelers/Editors

Label or edit billions of Data for FM models and hundreds of data for fine tuning interacting with data lake using a dedicated website

Data Science Team Extension Context Adaptation



Fine Tuners

Select the corresponding FM, evaluate the model & design the deployment method/infrastructure

Application Developer Team Integrate GenAI models in applications



Generative AI Developers, AppDev, & Prompt Engineers/Testers

Design prompt inputs, create examples of prompt input/outputs, and test the engineered prompts, develop the GenAI application and front-end

End-Users Consume Generative AI applications



Generative AI End-users

Consume Generative AI solutions as black box, share data and rate the quality of output



Generative AI Personas

Labeler Team
Data Preparation at Scale



Data Labelers/Editors

Label trillions of Data for FM models and hundreds of data for fine tuning interacting with data lake using a dedicated website

Data Science Team Extension
Context Adaptation



Fine Tuners

Select the corresponding FM, evaluate the model & design the deployment method/infrastructure

Application Developer Team
Integrate GenAI models in applications



Generative AI Developers, AppDev, & Prompt Engineers/Tests

Design prompt inputs, create examples of prompt input/outputs, and test the engineered prompts, develop the GenAI application and front-end

End-Users
Consume Generative AI applications



Generative AI End-users

Consume Generative AI solutions as black box, share data and rate the quality of output

Generative AI Personas

Labeler Team
Data Preparation at Scale



Data Labelers/Editors

Label trillions of Data for FM models and hundreds of data for fine tuning interacting with data lake using a dedicated website

Data Science Team Extension
Context Adaptation



Fine Tuners

Select the corresponding FM, evaluate the model & design the deployment method/infrastructure

Application Developer Team
Integrate GenAI models in applications



Generative AI Developers

Select, test, evaluate the FM, filter inputs/outputs, and develop the GenAI application back-end (e.g. LangChain Experts)



AppDev

Develop the front-end of the GenAI application

End-Users
Consume Generative AI applications



Generative AI End-users

Consume Generative AI solutions as black box, share data and rate the quality of output



Prompt Engineers

Design the input/output prompts to adapt the solution to the context and test the initial version

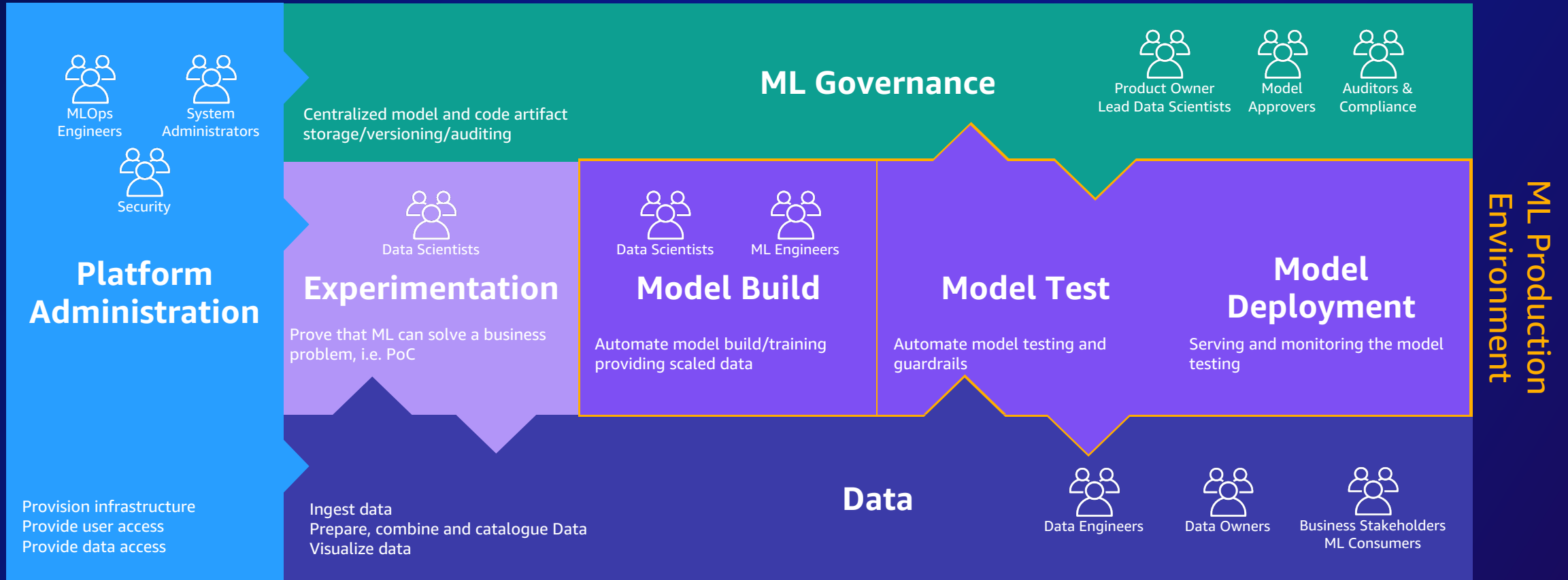


Prompt Testers

Test at scale the Generative AI solution (back-end/front-end) and feed their results to the prompt test repository

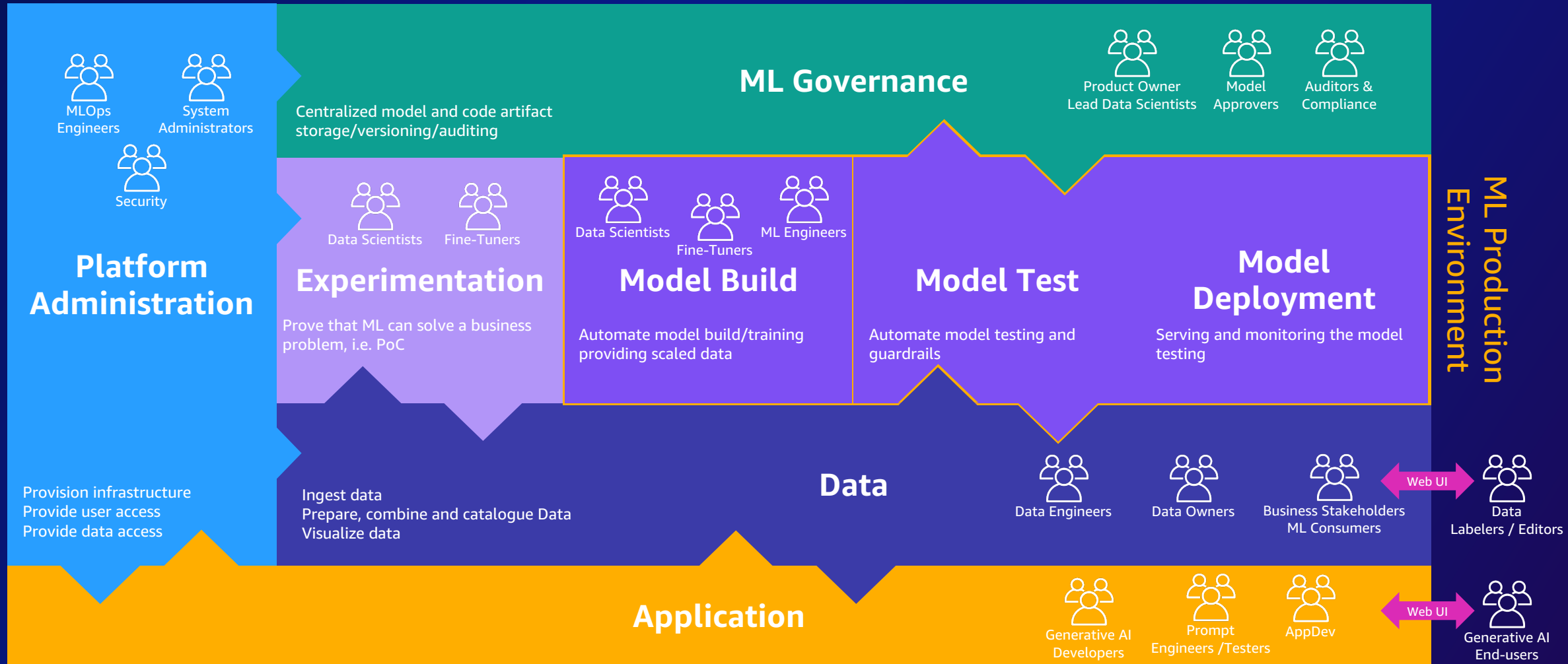
MLOps Foundation **People & Processes**

SEPARATION OF CONCERNS IS KEY FOR SUCCESS



MLOps & GenAI Foundation **People & Processes**

SEPARATION OF CONCERNS IS KEY FOR SUCCESS



Bonus: GenAI/LLM Vulnerabilities

<https://owasp.org/www-project-top-10-for-large-language-model-applications/descriptions/>

GenAI/LLM Vulnerabilities 1/2

Prompt Injections: Bypassing filters or manipulating the LLM using carefully crafted prompts that make the model ignore previous instructions or perform unintended actions.

Data Leakage: Accidentally revealing sensitive information, proprietary algorithms, or other confidential details through the LLM's responses.

Inadequate Sandboxing: Failing to properly isolate LLMs when they have access to external resources or sensitive systems, allowing for potential exploitation and unauthorized access.

Unauthorized Code Execution: Exploiting LLMs to execute malicious code, commands, or actions on the underlying system through natural language prompts.

SSRF Vulnerabilities: Exploiting LLMs to perform unintended requests or access restricted resources, such as internal services, APIs, or data stores.

GenAI/LLM Vulnerabilities 2/2

Overreliance on LLM-generated Content: Excessive dependence on LLM-generated content without human oversight can result in harmful consequences.

Inadequate AI Alignment: Failing to ensure that the LLM's objectives and behavior align with the intended use case, leading to undesired consequences or vulnerabilities.

Insufficient Access Controls: Not properly implementing access controls or authentication, allowing unauthorized users to interact with the LLM and potentially exploit vulnerabilities.

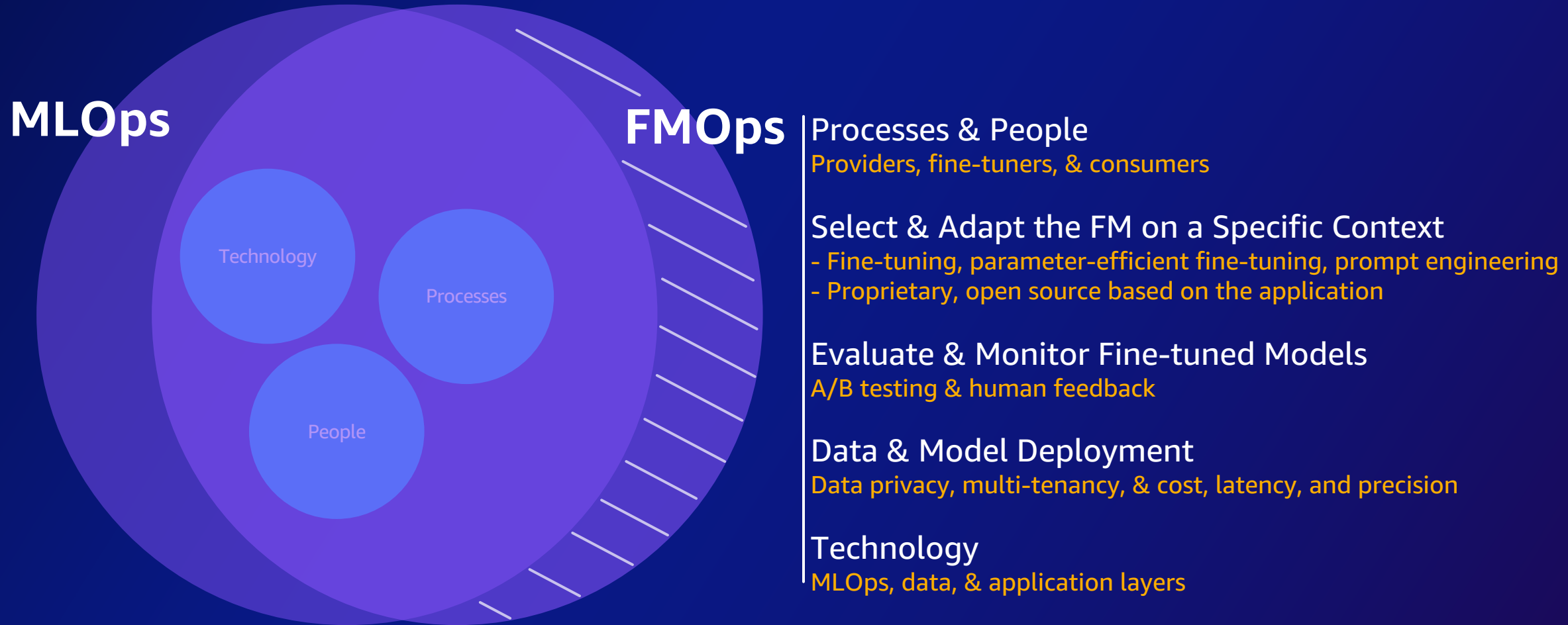
Improper Error Handling: Exposing error messages or debugging information that could reveal sensitive information, system details, or potential attack vectors.

Training Data Poisoning: Maliciously manipulating training data or fine-tuning procedures to introduce vulnerabilities or backdoors into the LLM

Conclusion



MLOps & FMOps Differentiators



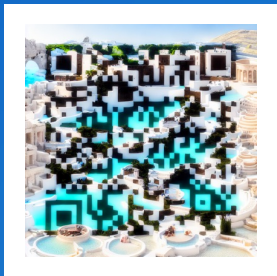
Thank you!



Please complete the session survey

Dr Sokratis Kartakis (he/him)

Senior MLOps SA Architect, EMEA
Amazon Web Services



 @sokratis.kartakis

Heiko Hotz (he/him)

Senior LLM SA Architect, EMEA
Amazon Web Services

