



**Reliably deliver  
ML models at scale with  
Amazon SageMaker MLOps**

# MLOps Overview & Strategy

Shelbee Eigenbrode

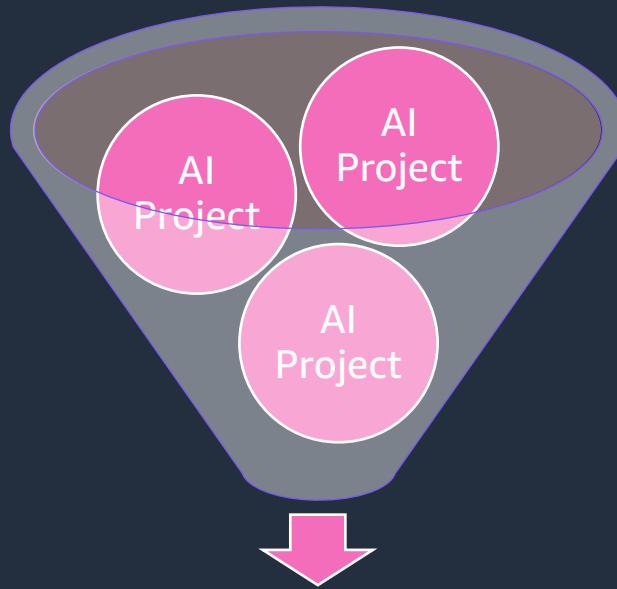
Principal AI/ML Specialist Solutions Architect



# MLOps Drivers

Creating a repeatable, reliable, and scalable path to production

## Creating a repeatable path to Production



54%  
of AI projects make it from  
pilot to production

Source: [Gartner AI in Organizations Survey, 2022](#)

© 2022, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

# MLOps Drivers

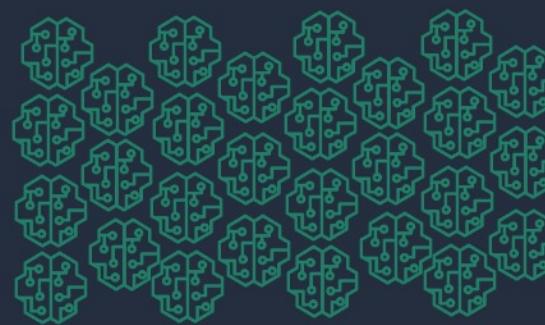
Creating a repeatable, reliable, and scalable path to production

## Reliably manage models at scale

Going From...



To...



MLOps: Recommended

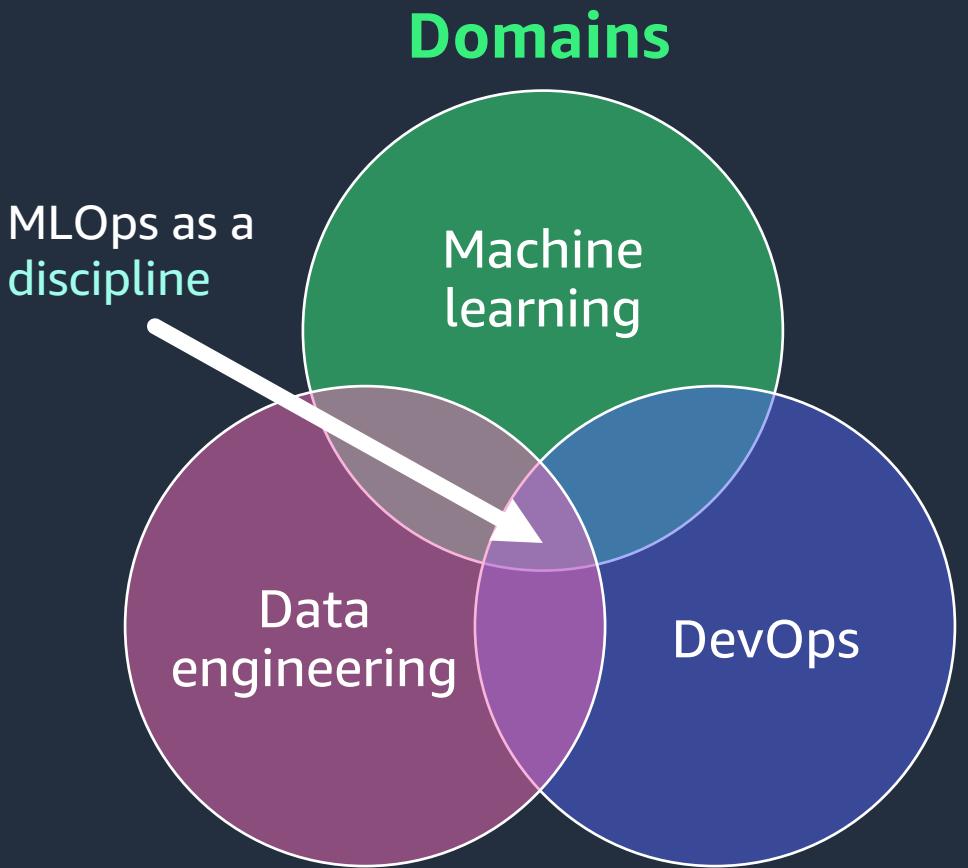
*MLOps is required to scale the adoption of ML*

MLOps: Required

**40%**  
of organizations reported having thousands, or hundreds of thousands of models deployed

# What is MLOps?

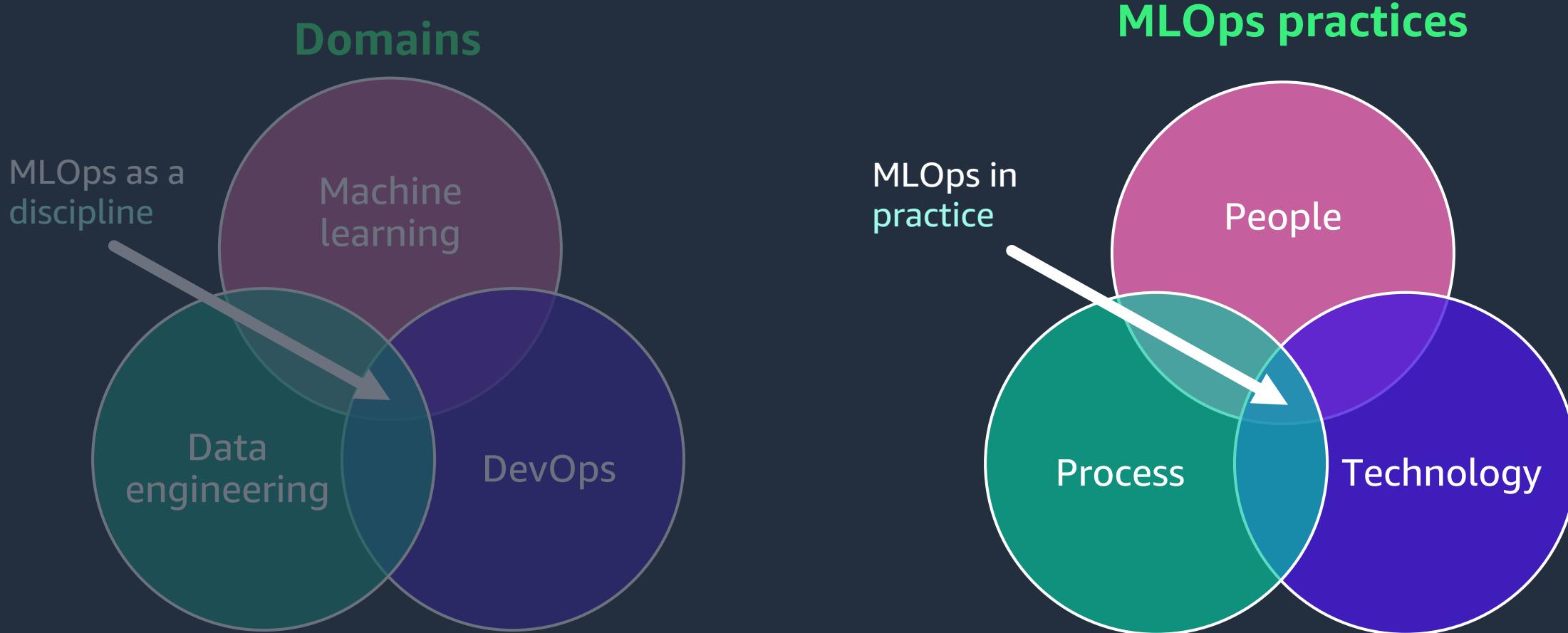
## It's not just technology



*MLOps is a discipline that sits at the intersection of the domains of data science, data engineering, and DevOps.*

# What is MLOps?

## It's not just technology



*MLOps is a discipline that sits at the intersection of the domains of data science, data engineering, and DevOps.*

*MLOps as a discipline is enabled by a core set of practices that span People, Process, and Technology across the Machine Learning Development Lifecycle.*

# MLOps vs DevOps

MLOps builds on DevOps practices with unique considerations for ML

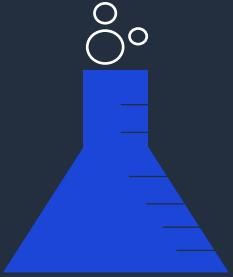
Feature	DevOps	MLOps
Code versioning	✓	✓
Artifact versioning	✓	✓
Continuous integration/continuous delivery (CI/CD)	✓	✓
Continuous Monitoring	✓	✓
Data lineage, governance, & management	✓	
Experiment tracking	✓	
Model lineage, governance, & management	✓	
Model re-training strategies	✓	
Model build & deployment workflows	✓	

# MLOps vs DevOps

## Machine Learning Development Lifecycle != Software Development Lifecycle

1

### Experimentation



*Model Development is unique...*

- *Requires a data strategy that includes ML & often requires access to larger compute/storage needs*
- *Enable rapid experimentation while balancing the needs of IT*
- *Experiment tracking is key for model development activities*

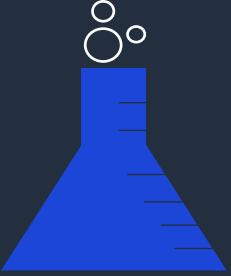


# MLOps vs DevOps

## Machine Learning Development Lifecycle != Software Development Lifecycle

1

### Experimentation

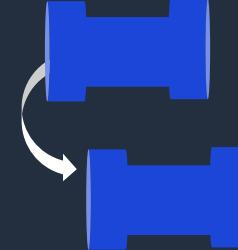


*Model Development is unique...*

- *Requires a data strategy that includes ML & often requires access to larger compute/storage needs*
- *Enable rapid experimentation while balancing the needs of IT*
- *Experiment tracking is key for model development activities*

2

### Pipelines



*Multiple disparate pipelines...*

- *Not always built during model development*
- *Pipelines can have different lifecycles, dependencies, and triggers*
- *End-to-end traceability expands in scope*

# MLOps vs DevOps

## Machine Learning Development Lifecycle != Software Development Lifecycle

1

### Experimentation

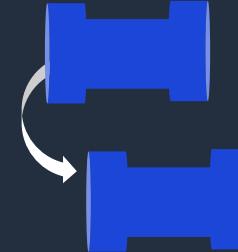


*Model Development is unique...*

- Requires a data strategy that includes ML & often requires access to larger compute/storage needs
- Enable rapid experimentation while balancing the needs of IT
- Experiment tracking is key for model development activities

2

### Pipelines



*Multiple disparate pipelines...*

- Not always built during model development
- Pipelines can have different lifecycles, dependencies, and triggers
- End-to-end traceability expands in scope

3

### Model Lineage & Governance



*Expanded scope for CI/CD...*

- Source control includes multiple code sources
- Version control includes data, features, pipeline outputs and model artifacts
- Governance includes data related to model performance and explainability

# MLOps vs DevOps

## Machine Learning Development Lifecycle != Software Development Lifecycle

1

### Experimentation

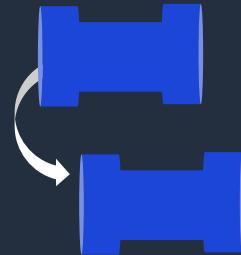


*Model Development is unique...*

- Requires a data strategy that includes ML & often requires access to larger compute/storage needs
- Enable rapid experimentation while balancing the needs of IT
- Experiment tracking is key for model development activities

2

### Pipelines



*Multiple disparate pipelines...*

- Not always built during model development
- Pipelines can have different lifecycles, dependencies, and triggers
- End-to-end traceability expands in scope

3

### Model Lineage & Governance

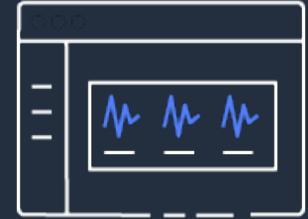


*Expanded scope for CI/CD...*

- Source control includes multiple code sources
- Version control includes data, features, pipeline outputs and model artifacts
- Governance includes data related to model performance and explainability

4

### Monitoring

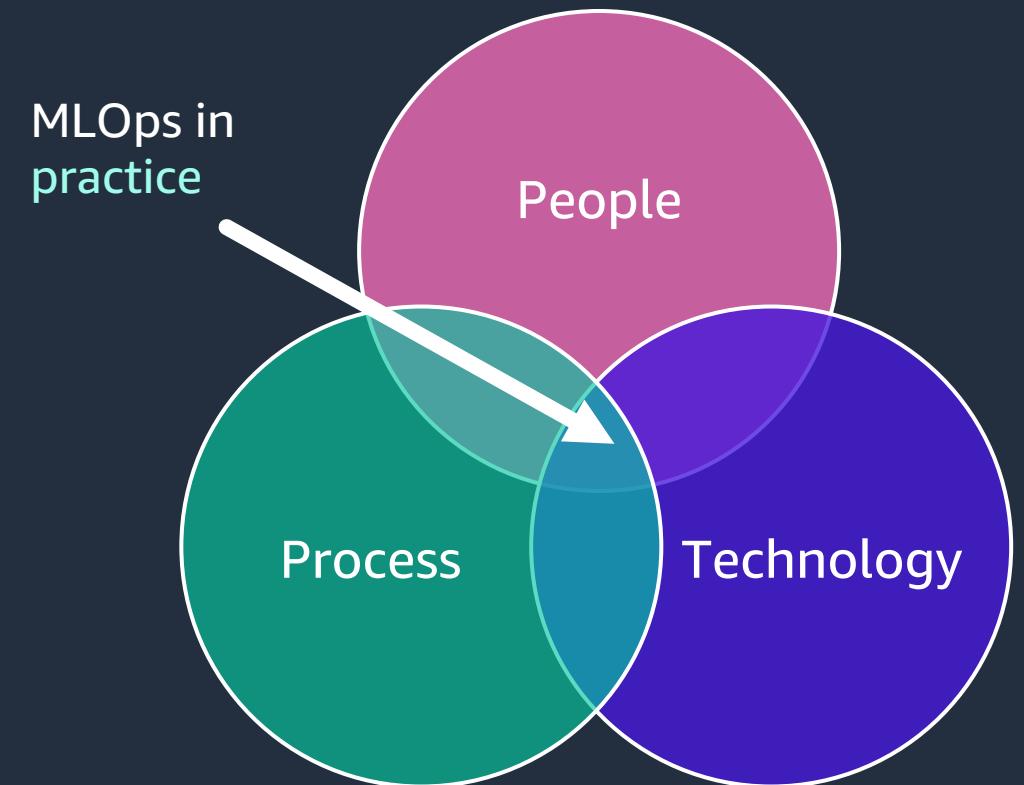


*Monitoring includes models...*

- Continuous monitoring includes system PLUS model specific monitoring (data/concept drift)
- Interpreting/responding to alerts may require multiple personas
- Continuous retraining is not always the goal!

# Why MLOps?

## MLOps drivers



© 2023, Amazon Web Services, Inc. or its affiliates.



### Economy

Improved time-to-market and reusability



### Agility

Faster iterations and increased repeatability



### Quality

Increased quality, reliability, and transparency

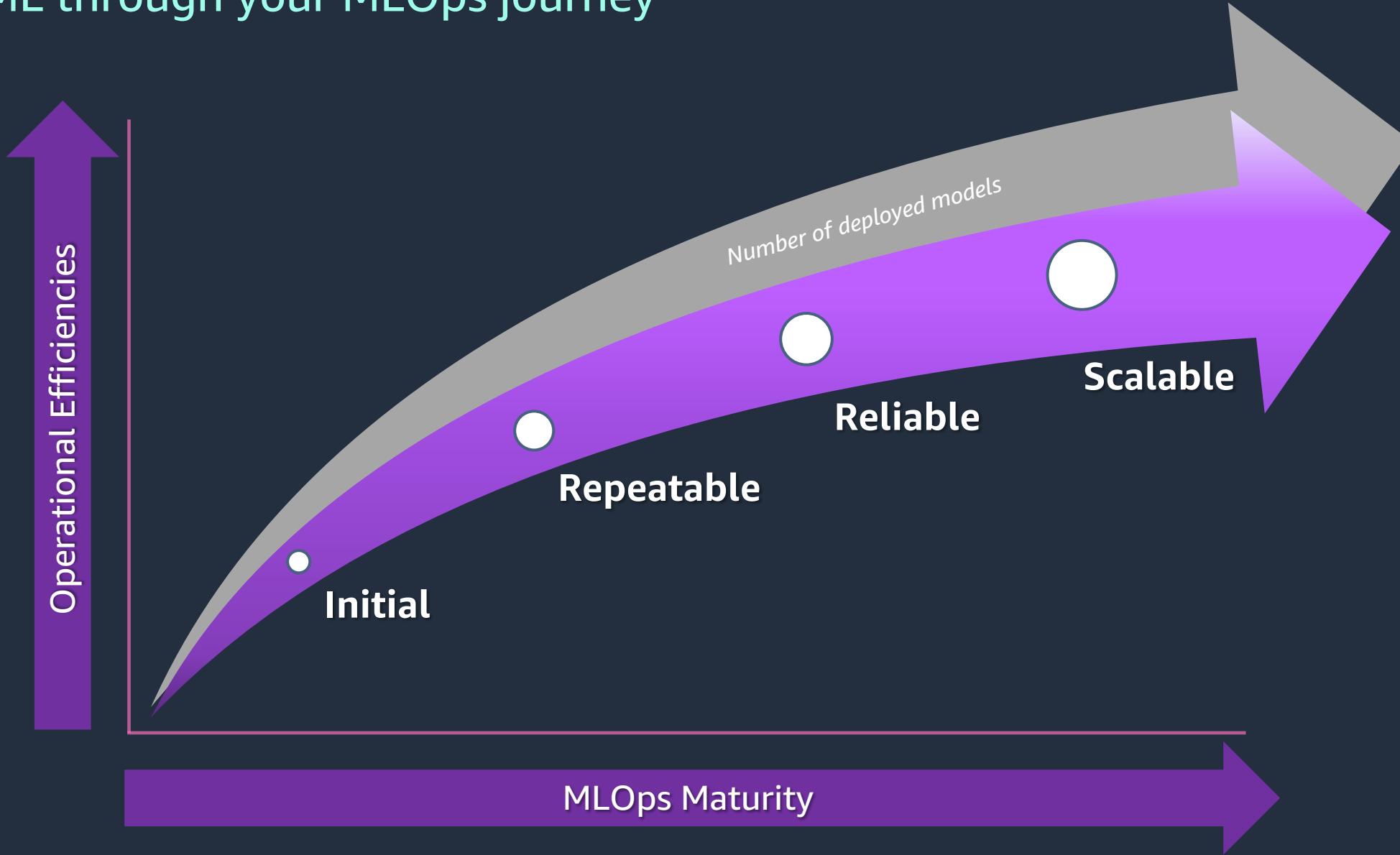


### Scale

Ability to efficiently build, deploy and manage models at scale

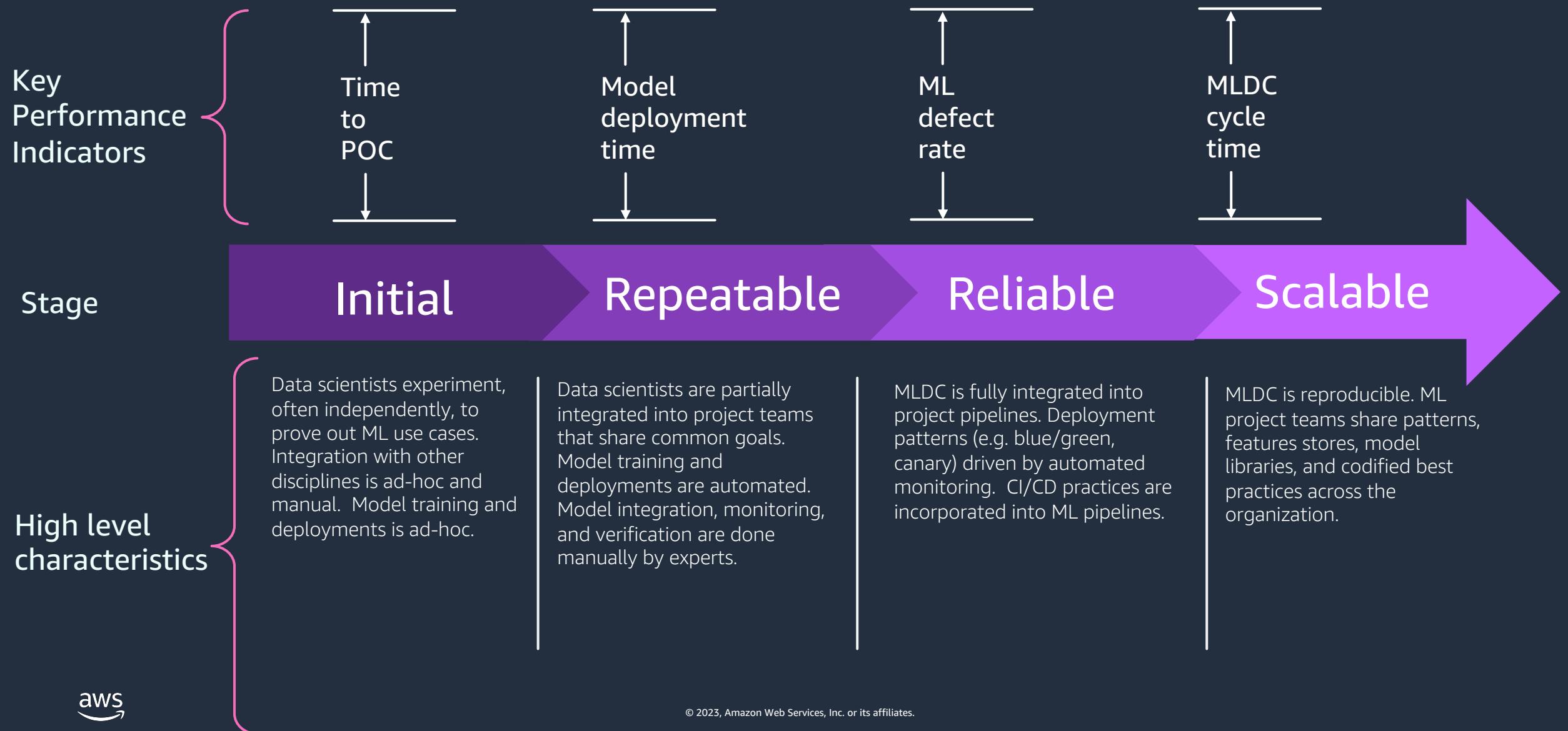
# MLOps Maturity Framework

## Scaling ML through your MLOps journey



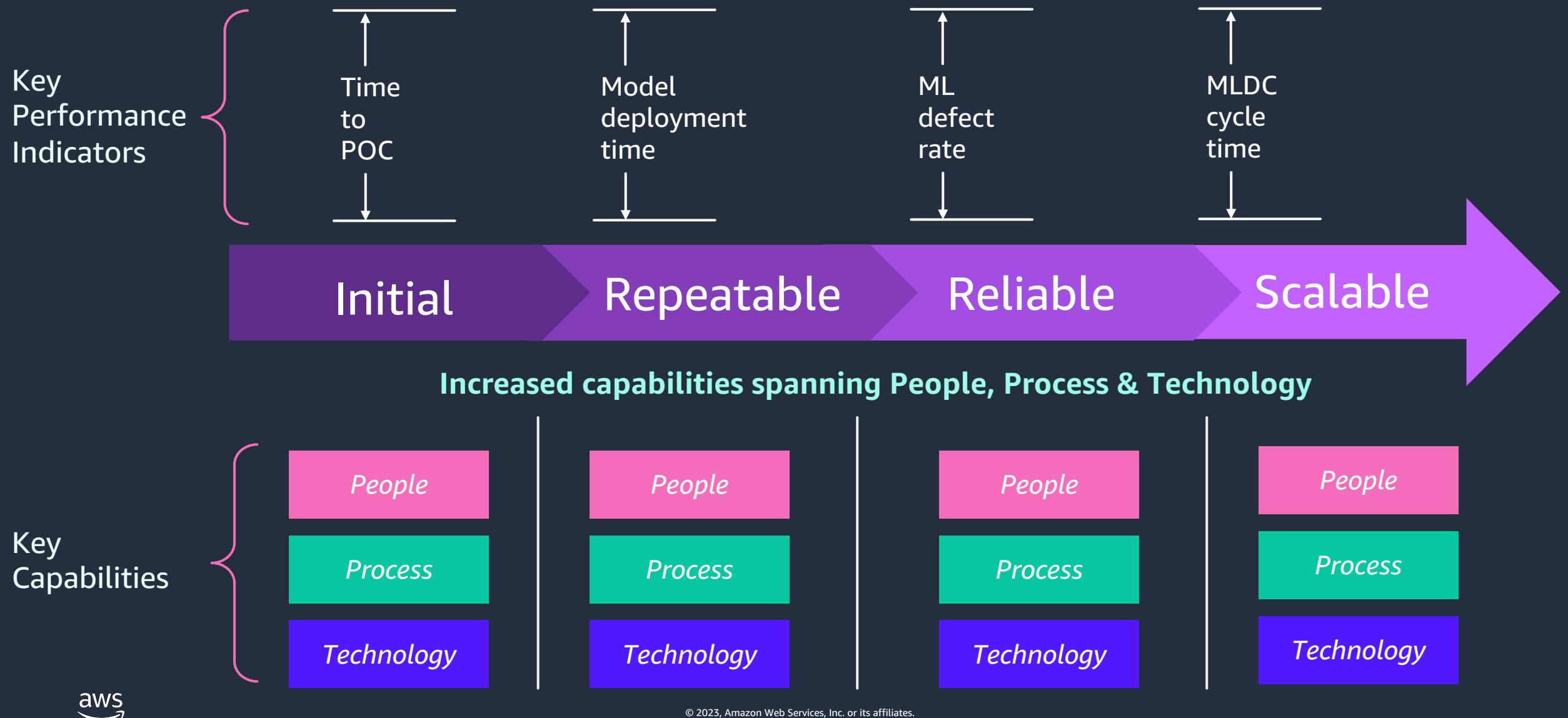
# MLOps Maturity Framework

## Scaling ML through your MLOps journey



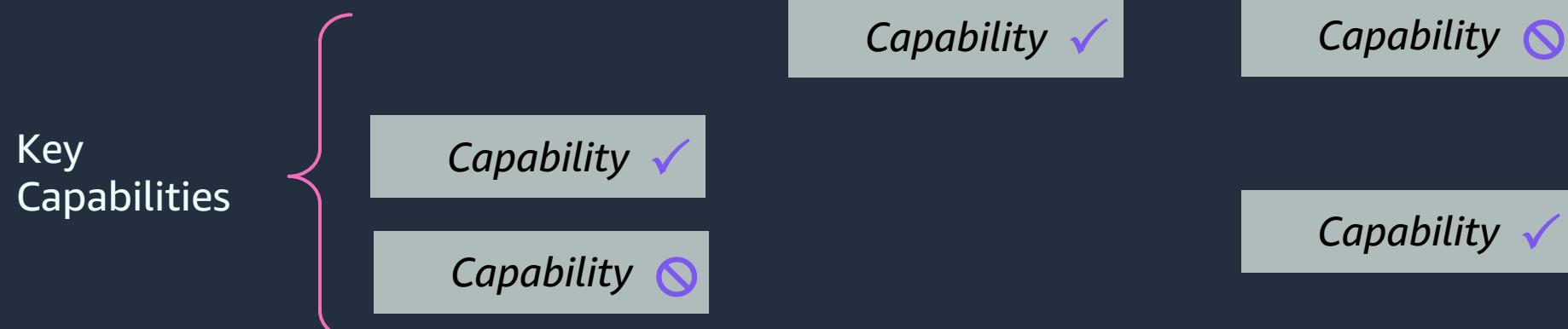
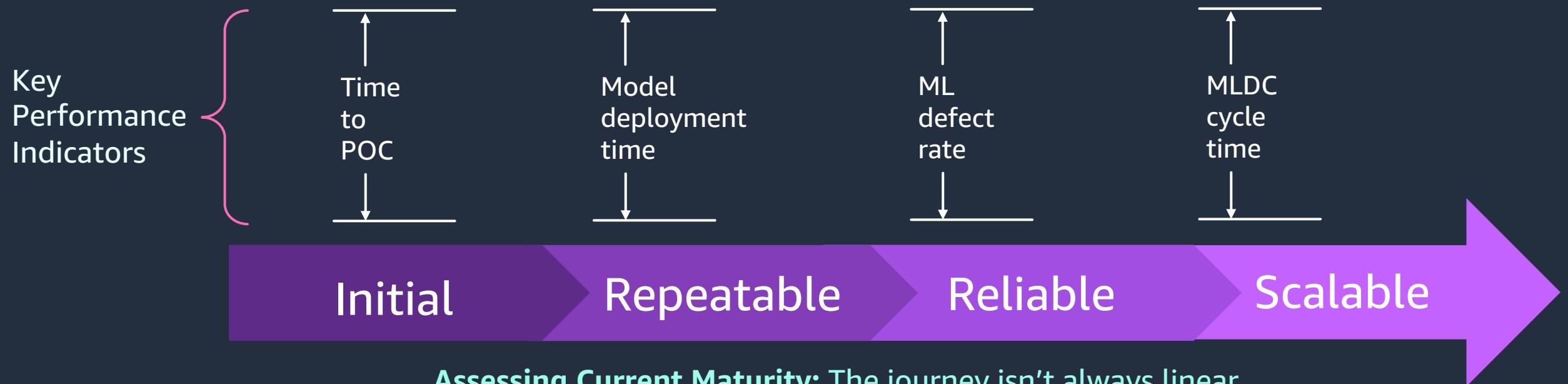
# MLOps Maturity Framework

## Assessing MLOps maturity



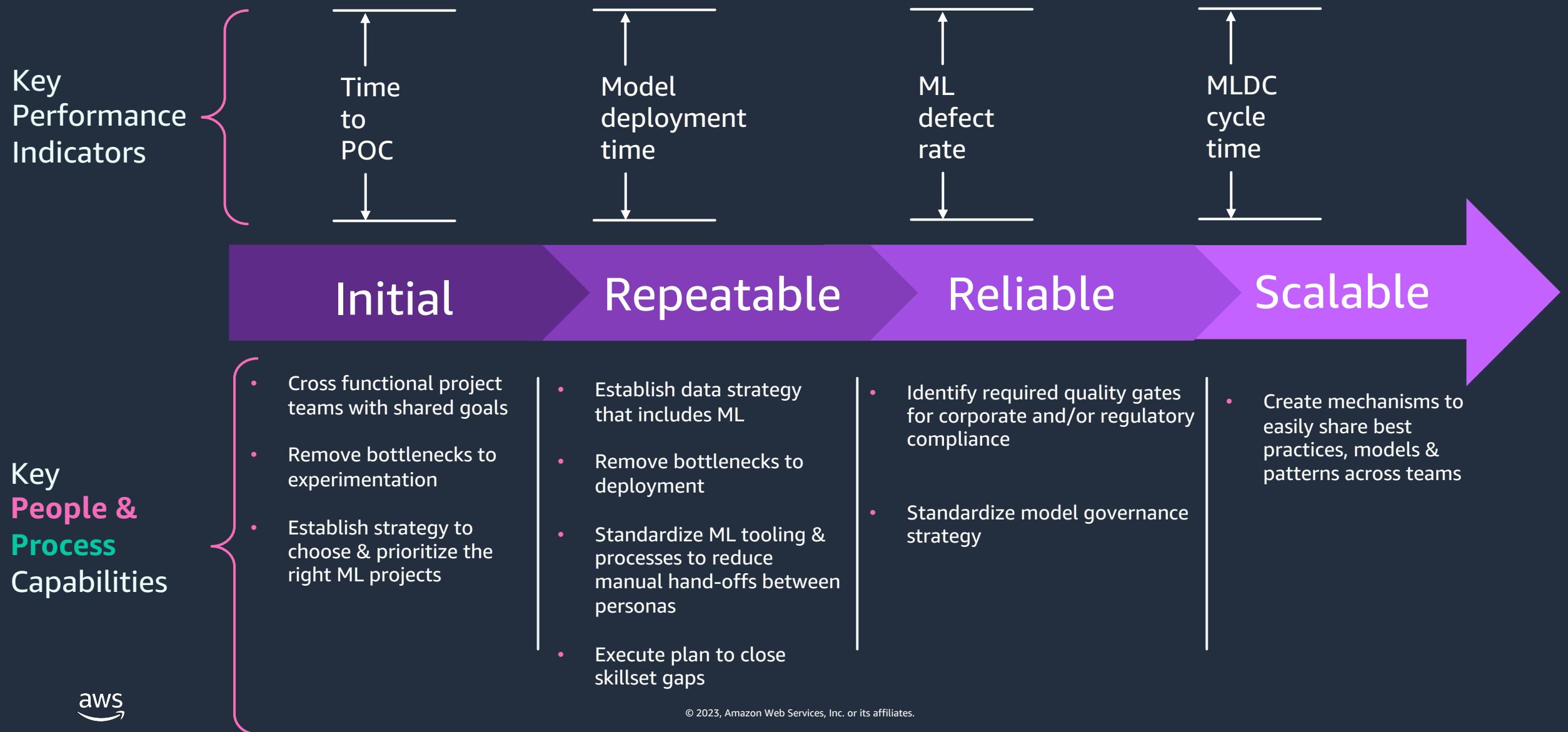
# MLOps Maturity Framework

## Assessing MLOps maturity



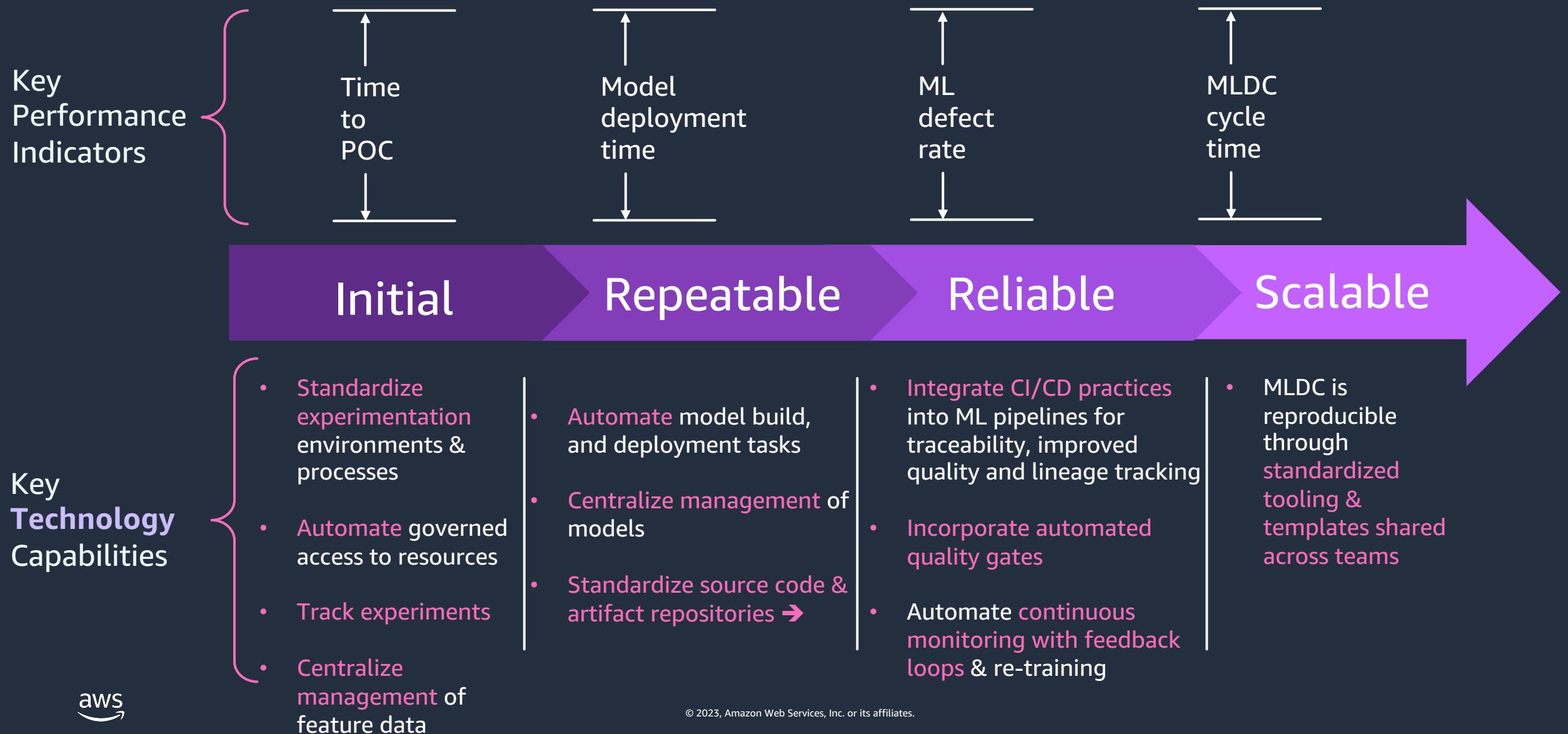
# MLOps Maturity Framework

## Scaling ML through your MLOps journey



# MLOps Maturity Framework

## Scaling ML through your MLOps journey



# Implementing MLOps Practices using Amazon SageMaker

## Today's Focus

### PREPARE DATA AND BUILD, TRAIN, DEPLOY AND MANAGE ML MODELS FOR ANY USE CASE

#### PREPARE →

##### Geospatial

Visualize geospatial data

##### Ground Truth

Create high quality datasets for ML

##### Data Wrangler

Aggregate and prepare data for ML

##### Processing

Built-in Python, BYO R/Spark

##### Feature Store

Store, catalog, search, and reuse features

##### Clarify

Detect bias and understand model predictions

#### BUILD →

##### Studio Notebooks & Notebook Instances

Fully managed Jupyter notebooks with elastic compute

##### Studio Lab

Free ML development environment

##### Built-in Algorithms

Integrated tabular, NLP, and vision algorithms

##### JumpStart

UI based discovery, training, and deployment of models, solutions, and examples

##### Autopilot

Automatically create ML models with full visibility

##### Bring Your Own

Bring your own container and algorithms

##### Local Mode

Test and prototype on your local machine

#### TRAIN & TUNE →

##### Fully Managed Training

Broad hardware options, easy to setup and scale

##### Distributed Training Libraries

High performance training for large datasets and models

##### Training Compiler

Faster deep learning model training

##### Automatic Model Tuning

Hyperparameter optimization

##### Managed Spot Training

Reduce training cost by up to 90%

##### Debugger and Profiler

Debug and profile training runs

##### Experiments

Track, visualize, and share model artifacts across teams

##### Customization Support

Integrate with popular open source frameworks and libraries

#### DEPLOY & MANAGE →

##### Fully Managed Deployment

Ultra low latency, high throughput inference

##### Real-Time Inference

For steady traffic patterns

##### Serverless Inference

For intermittent traffic patterns

##### Asynchronous Inference

For large payloads or long processing times

##### Batch Transform

For offline inference on batches of large datasets

##### Multi-Model Endpoints

Reduce cost by hosting multiple models per instance

##### Multi-Container Endpoints

Reduce cost by hosting multiple containers per instance

##### Shadow Testing

Validate model performance in production

##### Inference Recommender

Automatically select compute instance and configuration

##### Model Monitor

Maintain accuracy of deployed models

##### Kubernetes Operators & Components

Manage and monitor models on edge devices

##### Edge Manager

Manage and monitor models on edge devices

##### MLOps: Pipelines | Projects | Model Registry

Workflow automation, CI/CD for ML, central model catalog

Canvas  
Generate accurate machine learning predictions—no code required

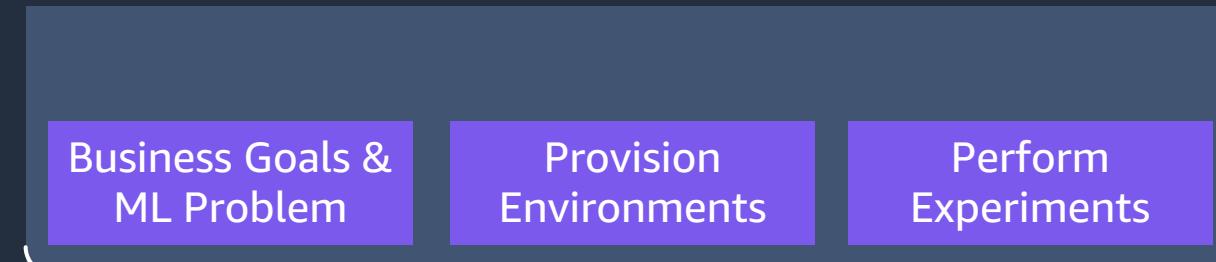
Studio | RStudio  
Integrated development environment (IDE) for ML

Governance  
Model Cards | Dashboard | Permissions

# Implementing MLOps Practices using Amazon SageMaker

## SageMaker Benefits

### Develop Candidate Model(s)



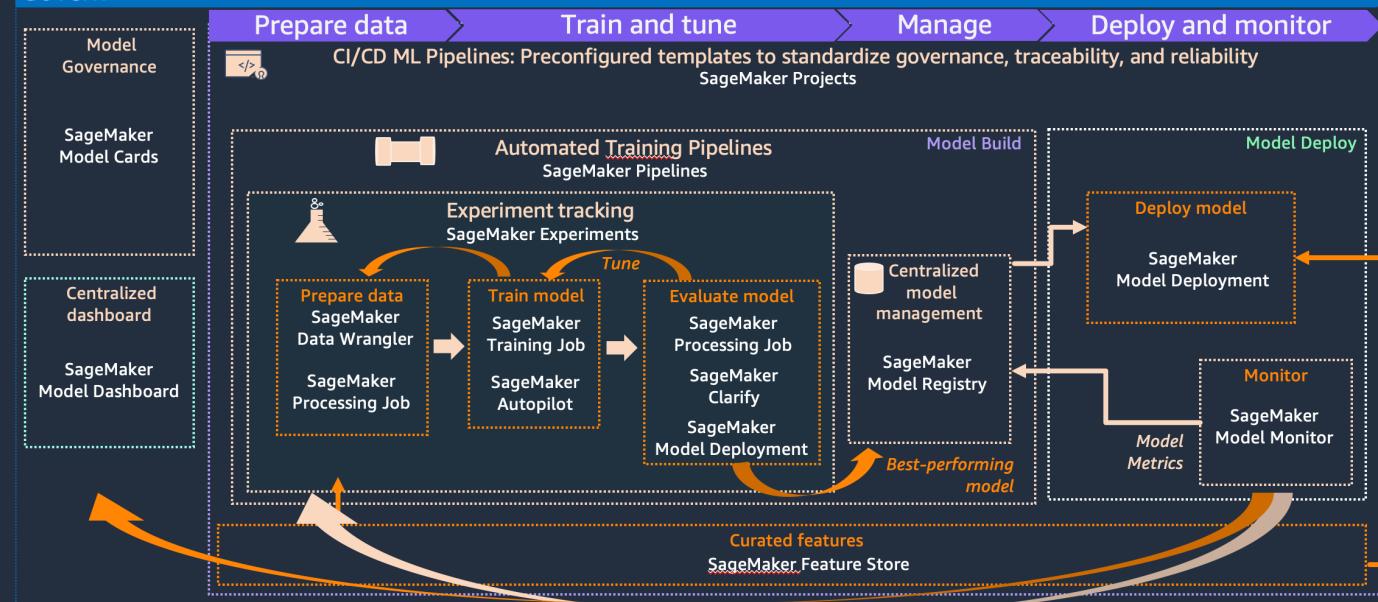
### Deploy Candidate Model(s)



### MLOps Maturity Framework in Practice

#### ML Project Phases – SageMaker View

##### Govern



- End-to-End **integrated features**
- **No servers to manage** for your MLOps tooling
- **Flexible options** across ML use cases
- **Purpose built** for Machine Learning

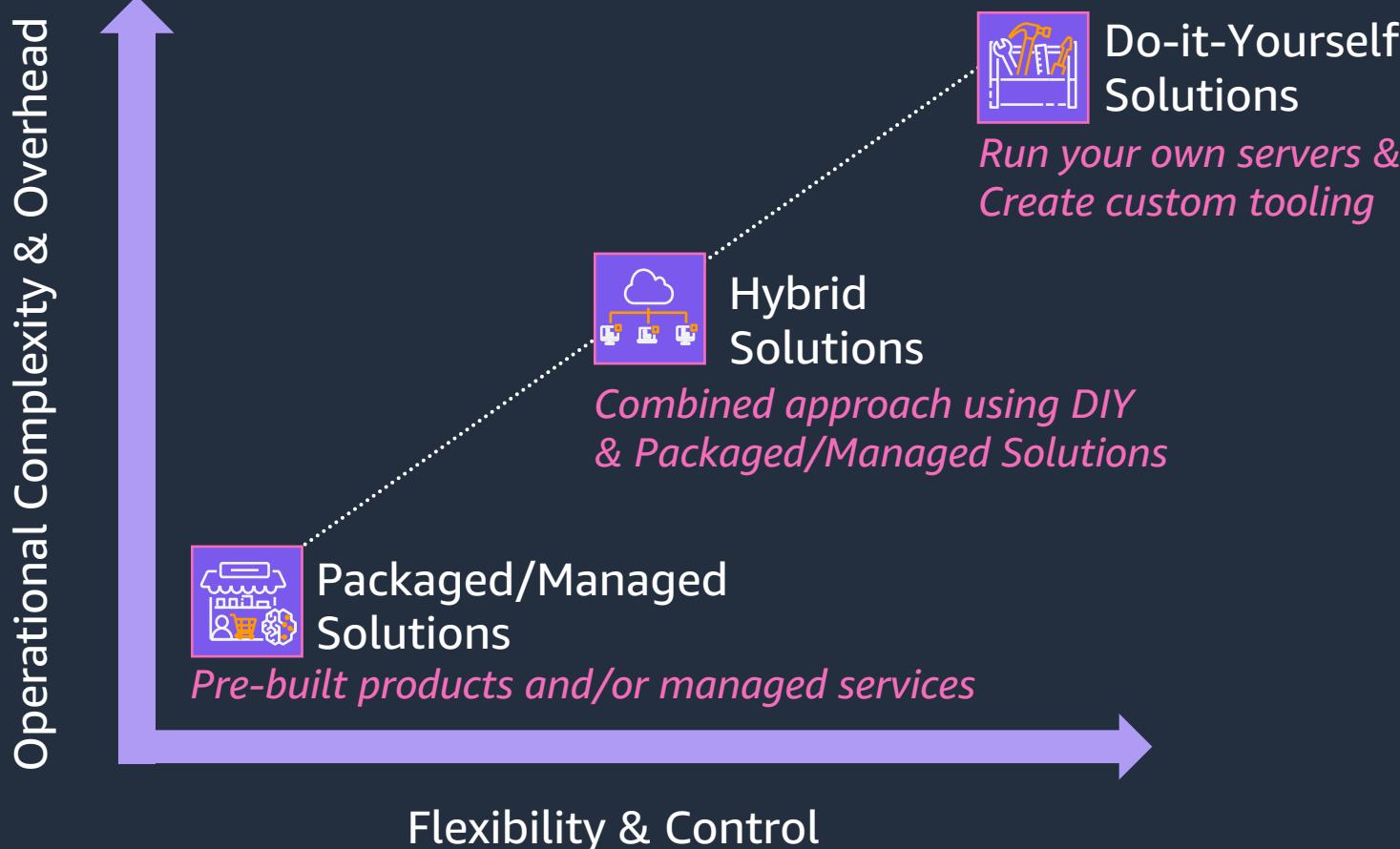


© 2022, Amazon Web Services, Inc. or its affiliates.

© 2022, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

# MLOps Technical Implementation Choices

## Understand Your Tradeoffs



1

Start With Process

Key  
Lessons Learned



1 Start With Process

2

You Can't Scale With People Issues

Key  
Lessons Learned



## Key Lessons Learned

1 Start With Process

2 You Can't Scale With People Issues

3 Technical Implementation(s) – One Size Does Not Fit All



## Key Lessons Learned



- 1 Start With Process
- 2 You Can't Scale With People Issues
- 3 Technical Implementation(s) – One Size Does Not Fit All
- 4 Standardize on Managed Services & Standard Components

## Key Lessons Learned



- 1 Start With Process
- 2 You Can't Scale With People Issues
- 3 Technical Implementation(s) – One Size Does Not Fit All
- 4 Standardize on Managed Services & Standard Components
- 5 Don't Ignore Unique Aspects of the ML Development Lifecycle

## Key Lessons Learned



- 1 Start With Process
- 2 You Can't Scale With People Issues
- 3 Technical Implementation(s) – One Size Does Not Fit All
- 4 Standardize on Managed Services & Standard Components
- 5 Don't Ignore Unique Aspects of the ML Development Lifecycle
- 6 Start Small & Iterate

# Today's Agenda

Time	Type	Title
9:00 – 9:30	Presentation	MLOps Overview & Strategy
9:30 – 10:15	Presentation	Model Development & Experimentation
10:15 – 10:45	Lab	Utilizing SageMaker Experiments during model development
10:45 – 11:00	Break	
11:00 – 11:30	Presentation	Implementing repeatable mechanisms to automate and govern your model build and model deployment tasks
11:30 – 12:15	Lab	Automate model build process & manage model versions at scale
12:15 – 1:00	LUNCH	
1:00 – 1:30	Partner Session	
1:30 – 2:00	Service Team Panel	
2:00 – 2:45	Presentation	Implementing reliable mechanisms to deliver, monitor, and manage models
2:45 – 3:30	Lab	Reliably deliver and manage models at scale
3:30 – 3:45	Break	
3:45 – 4:00	Presentation	SageMaker Third-Party Integrations
4:00 – 4:45	Presentation	ML Platform Deep Dive
4:45 – 5:00	Closeout & Next Steps	

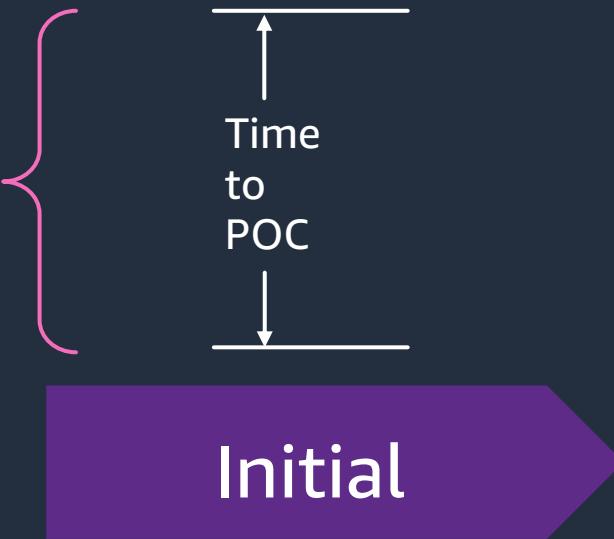
**Initial Stage:**

**Model Development & Experimentation**

# MLOps Maturity Framework in Practice

## Initial: Development & Experimentation with manual workflows

Key Performance Indicators



Initial

Key Technology Capabilities

- Standardize experimentation environments & processes
- Automate governed access to resources
- Track experiments
- Centralize management of feature data



# MLOps Maturity Framework in Practice

## Initial: Development & Experimentation with manual workflows

### Develop Candidate Model(s)

Business Goals  
& ML Problem

Provision  
Environments

Perform  
Experiments

### Operationalize Workflow & Candidate Model(s)

Feedback and Retraining

Automate Model  
Training

Package and  
Test Models

Deploy and  
Monitor



# MLOps Maturity Framework in Practice

## Initial: Development & Experimentation with manual workflows

### Develop Candidate Model(s)

Business Goals  
& ML Problem

Provision  
Environments

Perform  
Experiments

### Operationalize Workflow & Candidate Model(s)

Feedback and Retraining

Automate Model  
Training

Package and  
Test Models

Deploy and  
Monitor

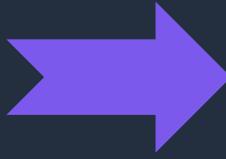


# Reduce time to Proof-of-Concept

## Challenge #1: Enabling innovation while balancing IT needs



Data Scientist wants to prove value of a new use case

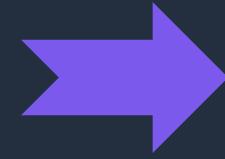


Perform iterative experimentation

- Understand data sources
- Try different data transformations
- Try different ML approaches
- Compare model performance



Development is mainly in a notebook or scripts



Access to model development resources



Data



Compute



Storage

### Key Technical Challenge

- Enable innovation while balancing the needs of IT

### Technical Solutions

- Provide access to governed data science environments with required resources using Infrastructure-as-Code (IaC), Configuration-as-Code (CaC), and Policy-as-Code (PaC)

### On AWS

- Amazon SageMaker Role Manager
- AWS Service Catalog



# Enable Innovation while balancing the needs of IT

## Initial: Standardize SageMaker Experimentation Environments

New!

## Amazon SageMaker Role Manager

### Key Benefits



Simplify permissions  
for ML activities



Guided workflows  
for role creation



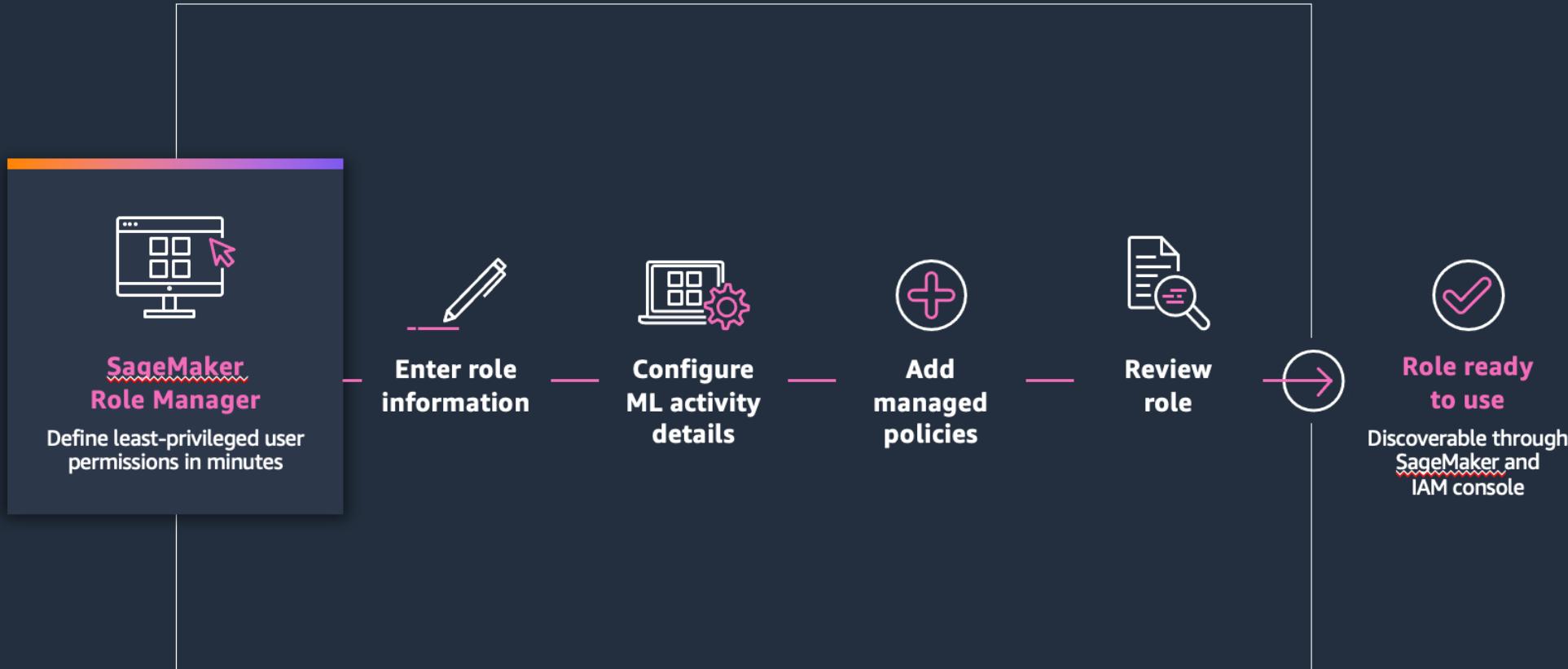
Accelerate user  
onboarding

# Enable Innovation while balancing the needs of IT

## Initial: Standardize SageMaker Experimentation Environments

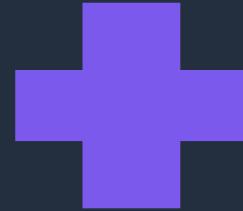
New!

### Amazon SageMaker Role Manager



# Enable Innovation while balancing the needs of IT

## Initial: Standardize SageMaker Experimentation Environments



### Amazon SageMaker Role Manager

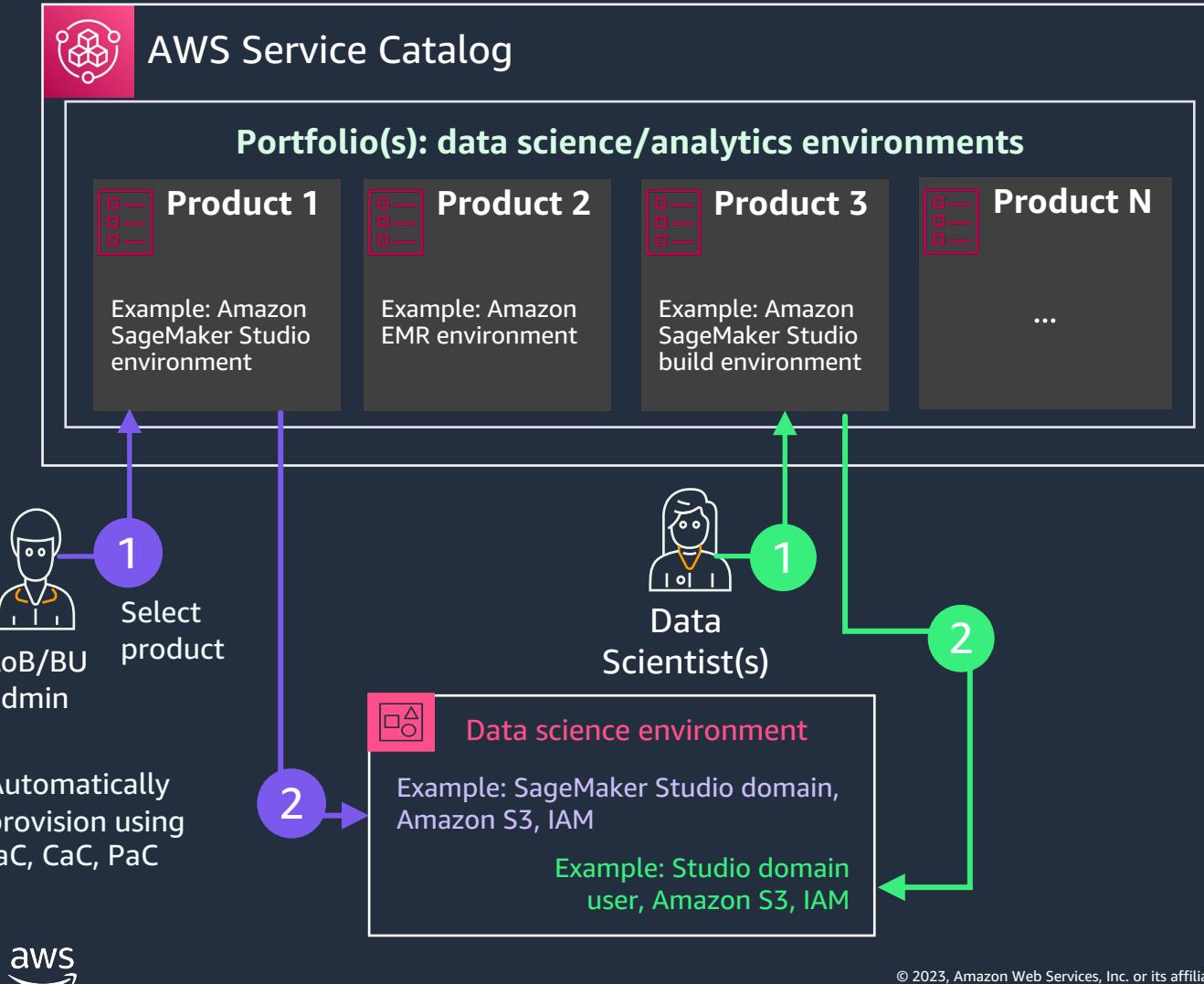
Define least-privileged user  
permissions in minutes

### AWS Service Catalog

Standardize resource  
configuration

# Enable Innovation while balancing the needs of IT

## Initial: Standardize SageMaker Experimentation Environments



## On AWS

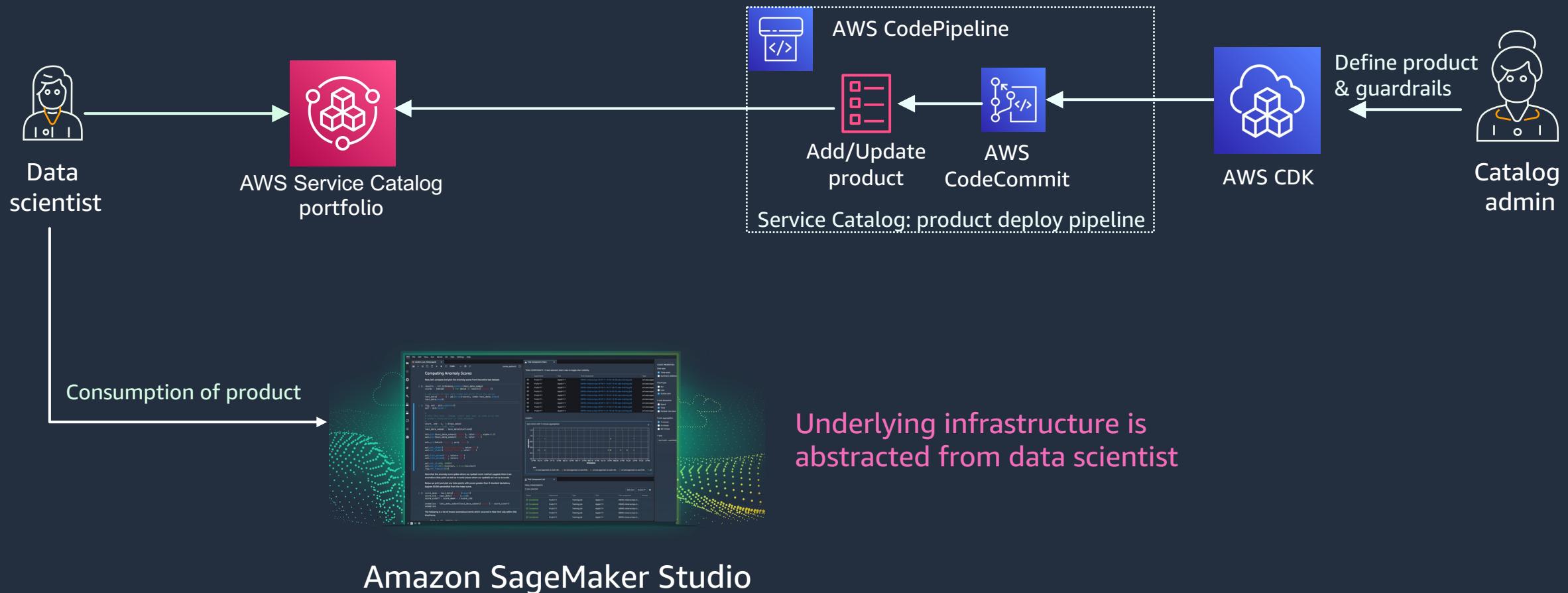
Use **AWS Service Catalog** to standardize the set up and provisioning of model build environments using IaC, CaC, PaC

## Benefits

- Reduce bottlenecks to innovation
- Ensure repeatability and consistency
- Preconfigure security properties
- Implement guardrails via code

# Enable Innovation while balancing the needs of IT

## Initial: Standardize SageMaker Experimentation Environments



# MLOps Maturity Framework in Practice

Initial Phase Goal: Standardize environments & experiment tracking

Govern

Environment Setup  
SageMaker Role Manager  
AWS Service Catalog

Prepare data



Train and tune

Manage

Deploy and monitor

Model Build

Model Development  
SageMaker Studio



Prepare data  
SageMaker Data Wrangler  
SageMaker Processing Job

Train model  
SageMaker Training Job  
SageMaker Autopilot

Evaluate model  
SageMaker Processing Job  
SageMaker Clarify  
SageMaker Model Deployment

Tune

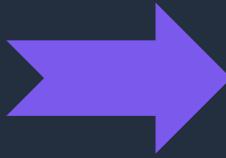


# Reduce time to Proof-of-Concept

## Challenge #2: Reduce feature engineering cycles & Increase consistency



Data Scientist wants to prove value of a new use case



Perform iterative experimentation

- Understand data sources
- Try different data transformations
- Try different ML approaches
- Compare model performance

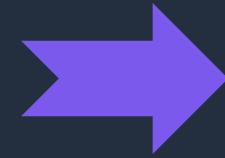


Development is mainly in a notebook or scripts

Access to discover & use curated features for model building



Feature Data



Key Technical Challenge

- Redundant & inconsistent feature engineering

Technical Solutions

- Implement a feature store for discovery, reuse and consistency

On AWS

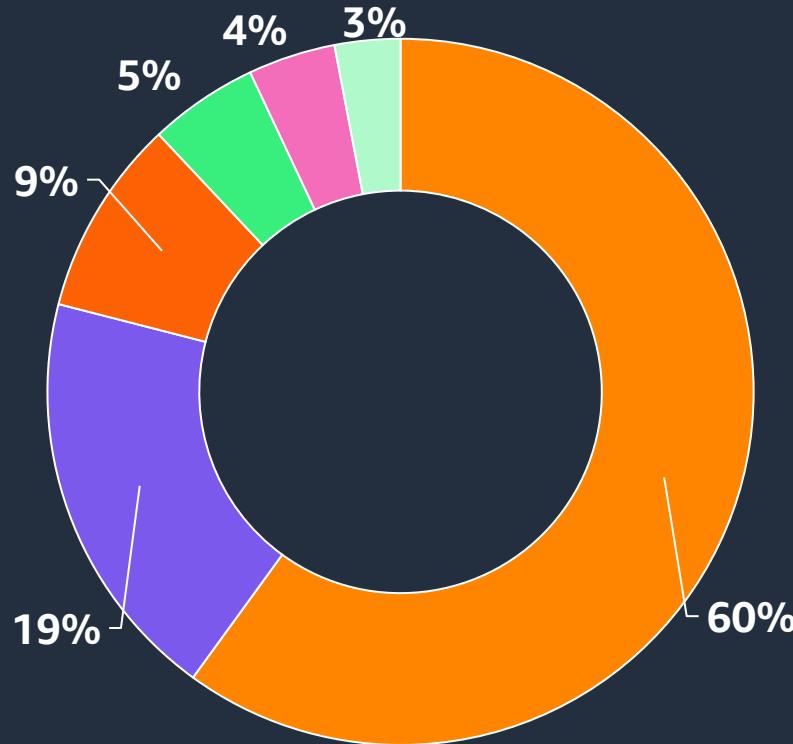
- Amazon SageMaker Feature Store



# Feature discovery, reuse, and consistency

Initial: Standardize Centralized feature stores

**60% of time spent on data preparation**



**What data scientists spend the most time doing**

- Cleaning and organizing data
- Collecting data sets
- Mining data for patterns
- Other
- Refining algorithms
- Building training sets

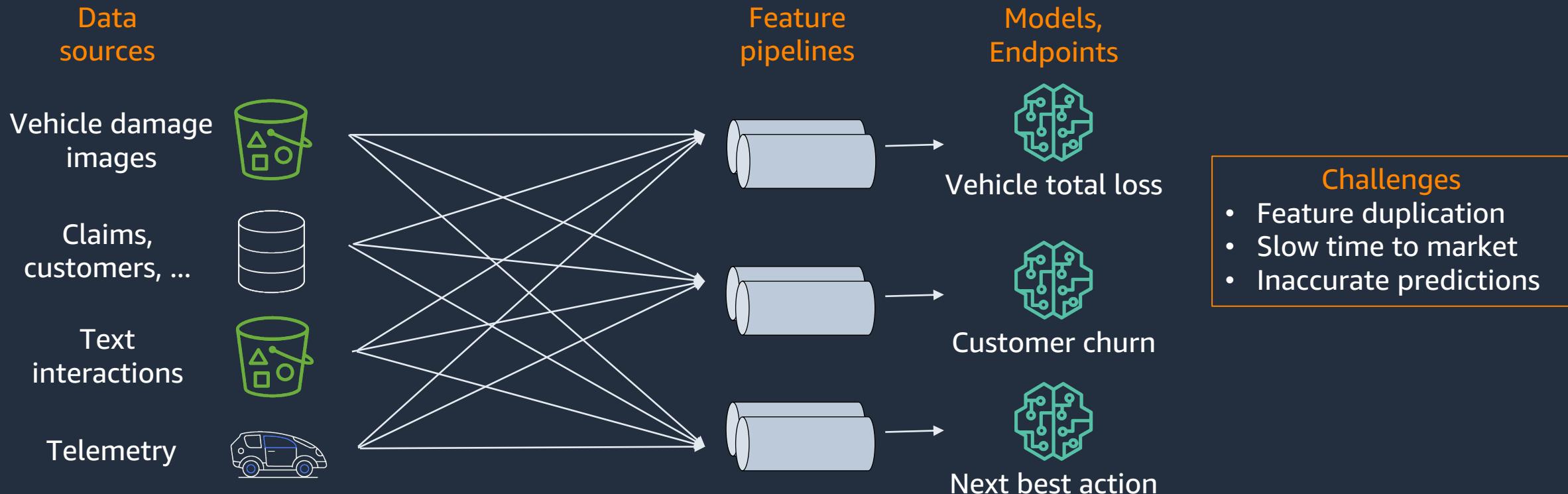
Source: [Forbes survey of 80 data scientists, March 2016](#)



# Feature discovery, reuse, and consistency

Initial: Standardize Centralized feature stores

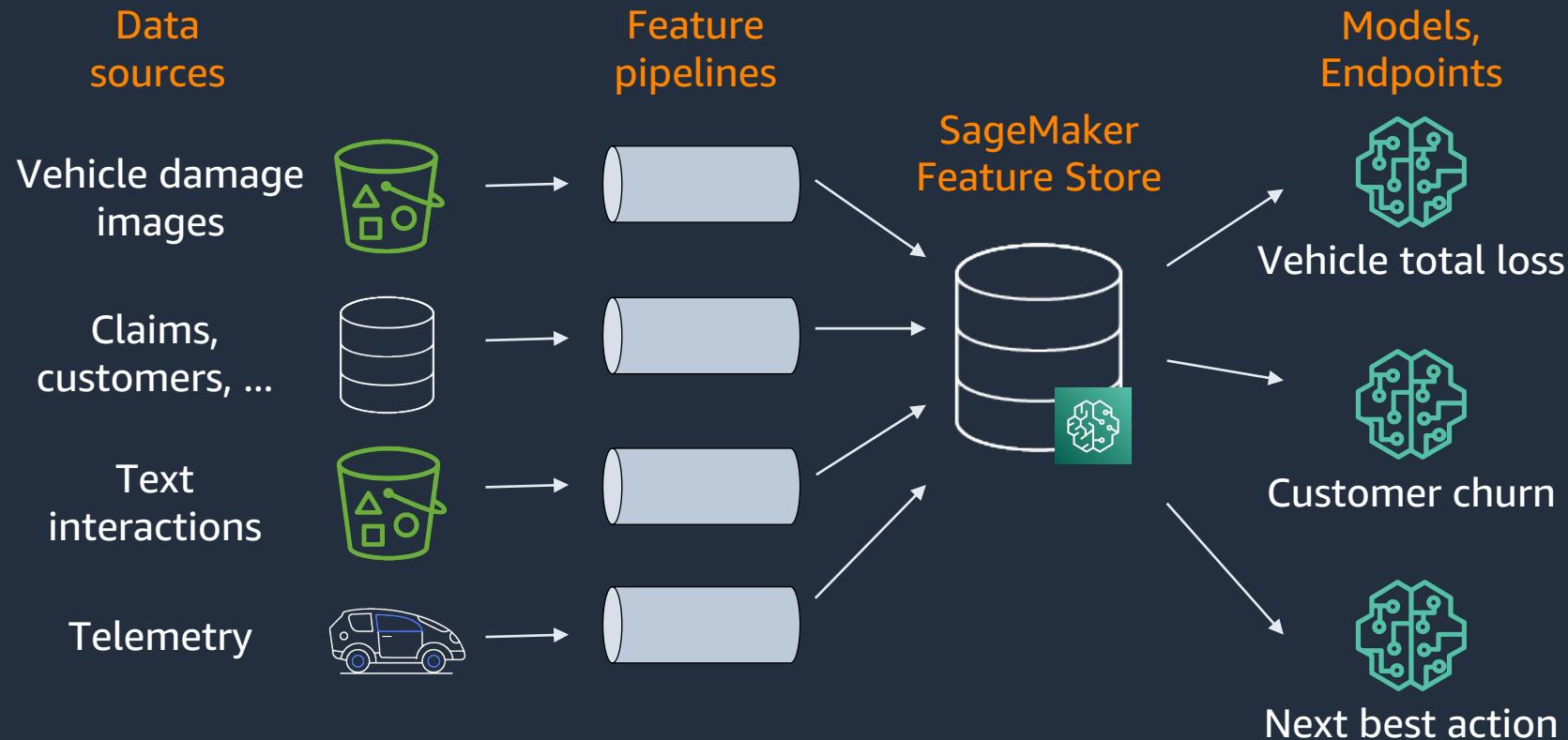
**... and teams too often start from scratch**  
Standalone feature engineering for each new model



# Feature discovery, reuse, and consistency

Initial: Standardize Centralized feature stores

## With SageMaker Feature Store... Build features once, reuse them across teams and models

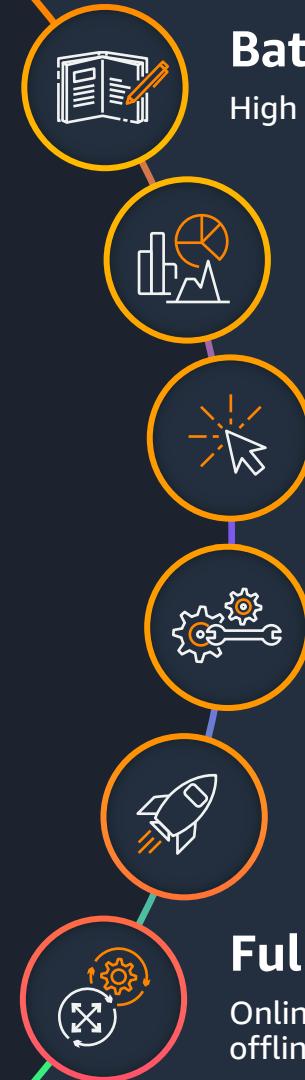


### Benefits

- Feature reuse
- Reproducible features
- Accurate training datasets
- Low latency inference
- Consistent features for training and inference

# SageMaker Feature Store

**Accelerate machine learning with a feature store**



## Batch and streaming ingestion

High throughput writes for ingesting features

## Online and offline features

Online features for real-time prediction, and offline features for historical data for model training and batch prediction

## Feature metadata and data cataloging

Store metadata for features and leverage automatic data cataloging to easily query and extract feature data

## Feature discovery and reuse

Search for features to reuse, before starting new development

## Security and access control

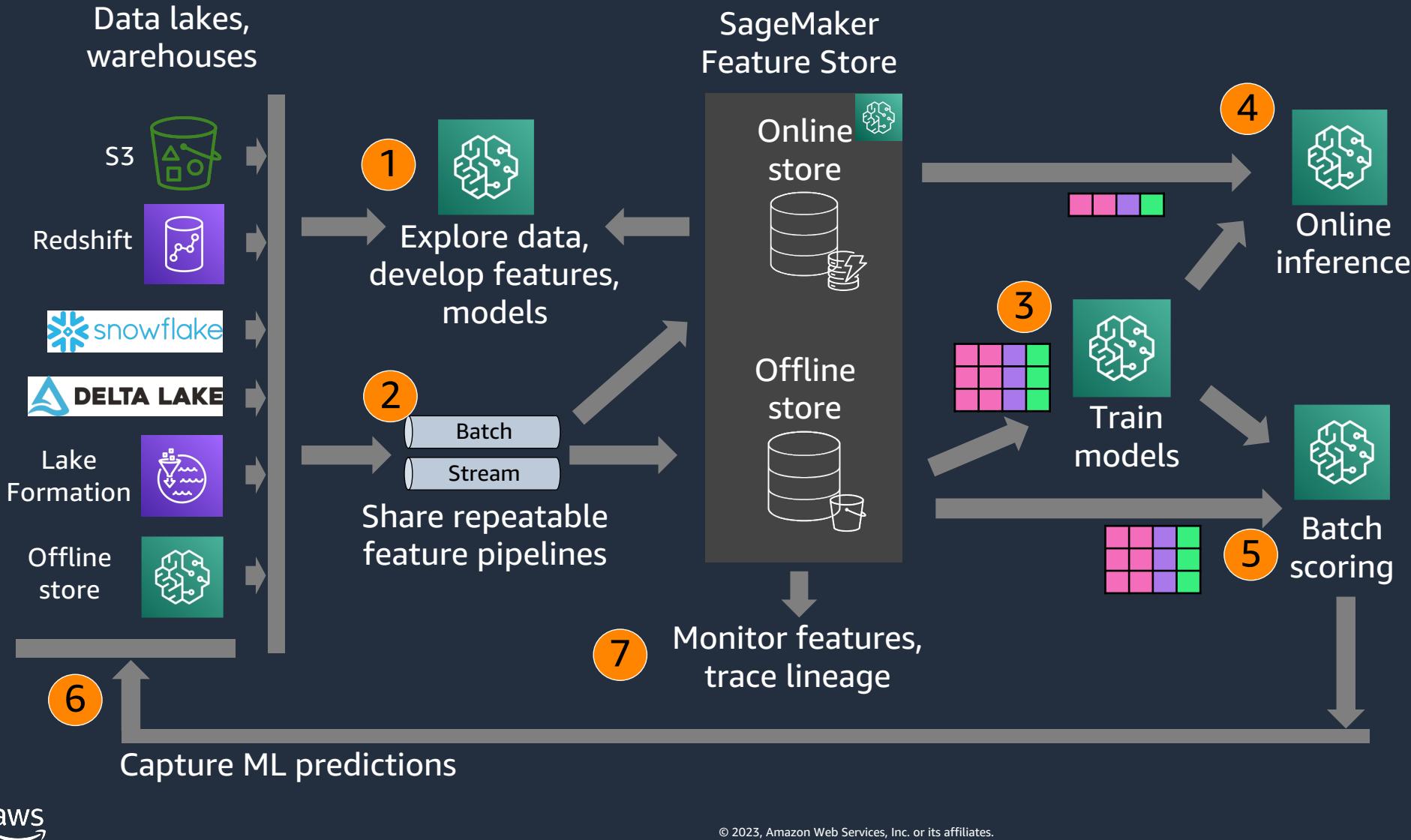
Access control for feature data and feature metadata, and support for encryption at rest, Amazon VPC, and AWS PrivateLink

## Fully managed

Online features cached in low-latency store; maintain consistency between online and offline store to avoid train-infer skew

# Amazon SageMaker Feature Store

## Using Feature Store with data lakes, warehouses



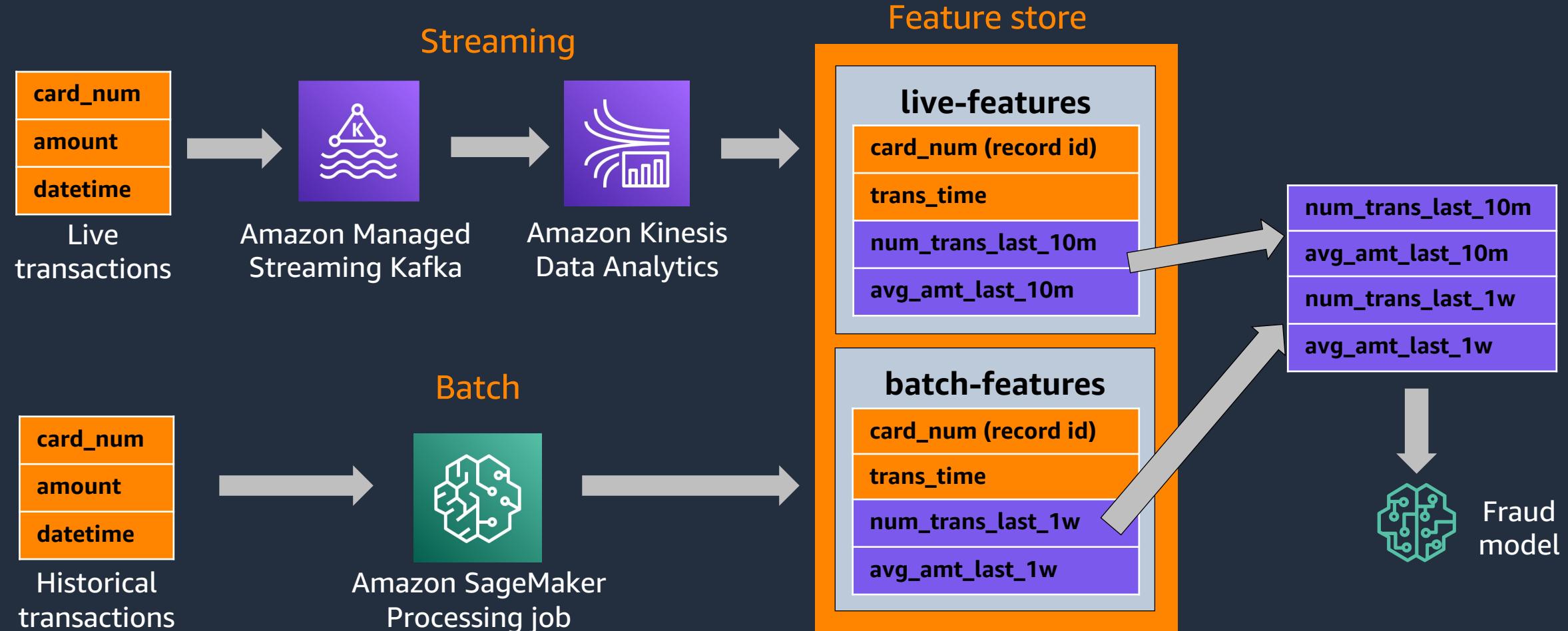
### Steps

- 1 Use Studio notebooks or Data Wrangler to explore existing data, experiment with features and models
- 2 Create repeatable feature pipelines, share features
- 3 Train models from available features, with point-in-time feature accuracy
- 4 Make low latency online predictions without train/serve skew
- 5 Make batch predictions with latest feature values
- 6 Capture prediction results to drive downstream use cases
- 7 Monitor features, and trace lineage

# Amazon SageMaker Feature Store

Features are materialized at different frequencies

Online models can pull latest features with low latency



# Amazon SageMaker Feature Store

## Feature pipeline automation



# MLOps Maturity Framework in Practice

Initial Phase Goal: Standardize environments & experiment tracking

Govern

Environment Setup

SageMaker Role Manager

AWS Service Catalog

Prepare data

Train and tune

Manage

Deploy and monitor



Model Build

Prepare data  
SageMaker Data Wrangler  
SageMaker Processing Job

Train model  
SageMaker Training Job  
SageMaker Autopilot

Evaluate model  
SageMaker Processing Job  
SageMaker Clarify  
SageMaker Model Deployment

Tune

Model Development  
SageMaker Studio



Offline Store

Curated features

SageMaker Feature Store

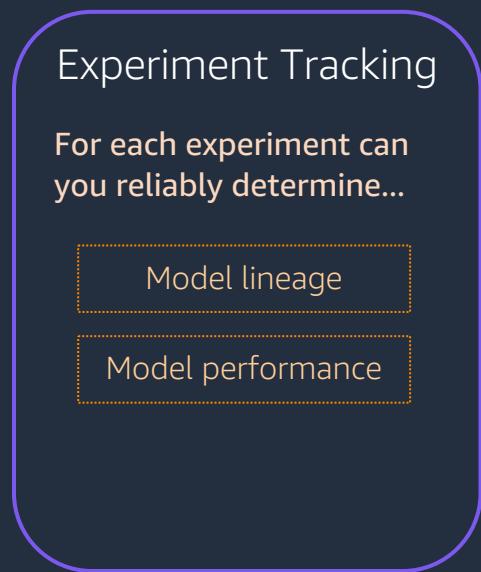
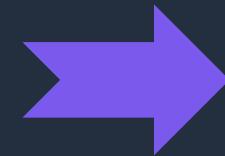
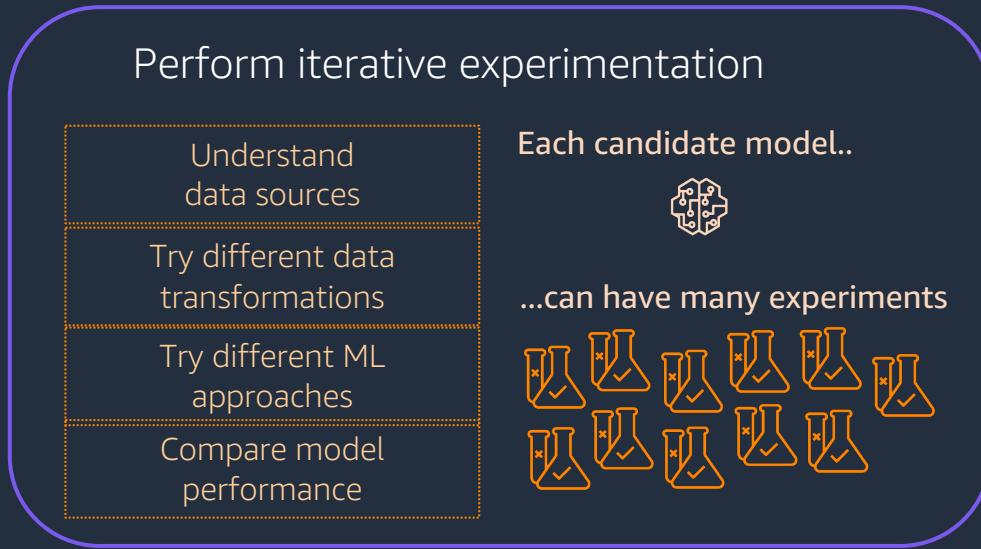
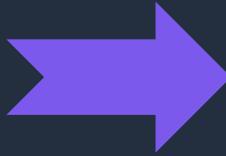


# Reduce time to Proof-of-Concept

## Challenge #3: Efficient experiment evaluation and reliable experiment tracking



Data Scientist wants to prove value of a new use case



### Key Technical Challenge

- Reliable experiment tracking for rapid iteration, management & deployment hand-off

### Technical Solutions

- Manage experiments reliably and at scale with experiment tracking tooling

### On AWS

- Amazon SageMaker Experiments



# Amazon SageMaker Experiments

## Create, manage, analyze, compare and track experiments at scale

Recent update  
(Dec 2022)

The image shows two screenshots of the Amazon SageMaker Studio interface. The top screenshot displays a Jupyter-style notebook titled 'pytorch\_experiment.ipynb'. The notebook content includes instructions for tracking experiments using the PyTorch model, details about experiments and runs, and a runtime section indicating a 20-minute execution time. The bottom screenshot shows the 'Launcher' interface, which provides quick access to various resources like notebooks, code consoles, and image terminals.



### Tracking at scale

Track parameters and metrics across experiments and users



### Custom organization

Organize experiments by teams, goals, and hypotheses



### Visualization

Easily visualize experiments and compare



### Metrics and logging

Log custom metrics using the Python SDK and APIs



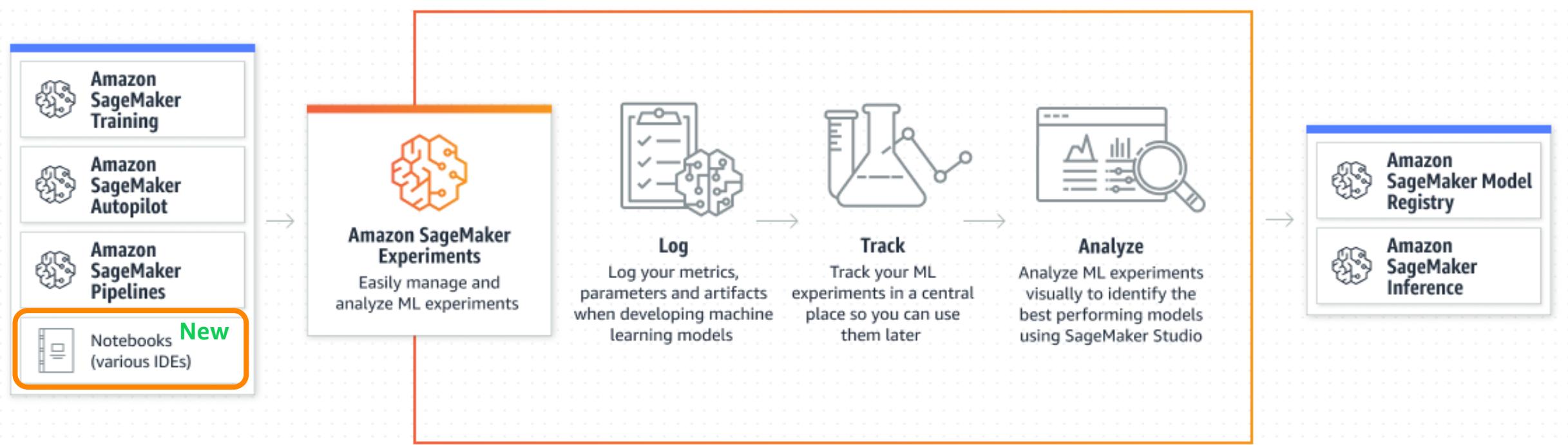
### Fast iteration

Quickly go back and forth, and maintain high-quality

# Amazon SageMaker Experiments

Create, manage, analyze, compare and track experiments at scale

## How it works



# Amazon SageMaker Experiments

Create, manage, analyze, compare and track experiments at scale

## Key Concepts

- 1. Experiment:** An experiment is a collection of runs. When you initialize a run in your training loop, you include the name of the experiment that the run belongs to. Experiment names must be unique within your AWS account.
- 2. Run:** A run consists of all the inputs, parameters, configurations, and results for one iteration of model training. Initialize an experiment run for tracking a training job with `Run()`.
- 3. Run Group (optional):** A run group is a unique identifier for groups of associated runs. For example, if the experiment is a pipeline experiment, then runs associated with a specific pipeline execution all belong to the same run group.



# Amazon SageMaker Experiments

## Create, manage, analyze, compare and track experiments at scale

## Integration with SageMaker Clarify

```
with Run(
    experiment_name='clarify-experiment',
    run_name="joint-run",
    sagemaker_session=sagemaker_session,
) as run:
    xgb.fit({"train": train_input}, log_level=logging.INFO)
    clarify_processor.run_bias(
        data_config=bias_data_config,
        bias_config=bias_config,
        model_config=model_config,
        model_predicted_label_config=preprocess_label_config,
        pre_training_methods="all",
        post_training_methods="all",
    )
    clarify_processor.run_explainability(
        data_config=explainability_data_config,
        model_config=model_config,
        explainability_config=shap_config,
    )
```

Bias Report

The screenshot shows the SageMaker Clarify Bias Report interface. The left sidebar has sections like Overview, Metrics, Charts, Output Artifacts, Bias Reports (which is selected), Explanations, Debugs, Settings, Parameters, Configurations, Input Artifacts, Details, Information, and Metadata. The main area is titled 'Reports' and contains a table of bias metrics. The table has columns for 'Bias metric', 'Bias value', and 'Description'. The metrics listed include Conditional Demographic Disparity in Labels (CDDL), Class imbalance (CI), Difference in Positive Proportions in Labels (DPL), Jensen-Shannon Divergence (JS), Kullback-Liebler Divergence (KL), Kolmogorov-Smirnov Distance (KS), L<sub>p</sub>-Norm (LP), Total Variation Distance (TVD), Accuracy Difference (AD), Conditional Demographic Disparity in Predicted Labels (CDPL), Difference in Acceptance Rates (DAR), Difference in Conditional Acceptance (DCA), Difference in Conditional Outcomes (DCR), Disparate (Adverse) Impact (DI), and Difference in Positive Proportions in Predicted Labels (DPPL). The 'Sex' column is selected in the dropdown, and the '0' threshold is selected in the other dropdown.

Bias metric	Bias value	Description
Conditional Demographic Disparity in Labels (CDDL)	0.21	The metric examines whether, in the training data, the disadvantaged class ...
Class imbalance (CI)	0.35	Detects if the advantaged group is represented in the dataset at a substantia...
Difference in Positive Proportions in Labels (DPL)	0.2	Detects if one class has a significantly higher proportion of desirable (or, alte...
Jensen-Shannon Divergence (JS)	0.031	JS measures how much the label distributions of different classes diverge fro...
Kullback-Liebler Divergence (KL)	0.14	In a binary case, a relative entropy measure of how much the label distribut...
Kolmogorov-Smirnov Distance (KS)	0.2	This metric is equal to the maximum divergence in a label across the classes ...
L <sub>p</sub> -Norm (LP)	0.28	This measure of distance in label distributions is the normed direct distance ...
Total Variation Distance (TVD)	0.2	This measure of distance in label distributions is half the Hamming distance ...
Accuracy Difference (AD)	-0.11	This metric examines whether the classification by the model is more accurat...
Conditional Demographic Disparity in Predicted Labels (CDPL)	0.2	The metric examines whether the model predicted a bigger proportion of rel...
Difference in Acceptance Rates (DAR)	-0.0074	The difference in the rates of positive predicted outcomes across the advant...
Difference in Conditional Acceptance (DCA)	-0.23	This metric compares the actual labels to the predicted labels from the model...
Difference in Conditional Outcomes (DCR)	0.13	This metric compares the actual labels to the predicted labels from the model...
Disparate (Adverse) Impact (DI)	0.33	This metric examines whether the model predicts outcomes differently for e...
Difference in Positive Proportions in Predicted Labels (DPPL)	0.092	This metric examines whether the model predicts outcomes differently for a ...



# Amazon SageMaker Experiments

## Create, manage, analyze, compare and track experiments at scale

## Integration with SageMaker Clarify

```
with Run(
    experiment_name='clarify-experiment',
    run_name="joint-run",
    sagemaker_session=sagemaker_session,
) as run:
    xgb.fit({"train": train_input}, log_level=logging.INFO)
    clarify_processor.run_bias(
        data_config=bias_data_config,
        bias_config=bias_config,
        model_config=model_config,
        model_predicted_label_config=preprocess_label_config,
        pre_training_methods="all",
        post_training_methods="all",
    )
    clarify_processor.run_explainability(
        data_config=explainability_data_config,
        model_config=model_config,
        explainability_config=shap_config,
    )
```

The screenshot shows the SageMaker Clarify interface with the 'Bias Reports' tab selected. The left sidebar includes sections for Overview, Metrics, Charts, Output Artifacts, Bias Reports (selected), Explanations, Debugs, Settings, Parameters, Configurations, Input Artifacts, Details, Information, and Metadata. The main content area displays a list of bias metrics: Predicted column: Target, Predicted value or threshold: 1, Column analyzed for bias: Sex, Column value or threshold analyzed for b: 0. Below this, a list of other metrics is provided: Conditional Demographic Disparity in Labels (CDDL), Class imbalance (CI), Difference in Positive Proportions in Labels (DPL), Jensen-Shannon Divergence (JS), Kullback-Liebler Divergence (KL), Kolmogorov-Smirnov Distance (KS), L<sub>p</sub>-Norm (LP), Total Variation Distance (TVD), Accuracy Difference (AD), Conditional Demographic Disparity in Predicted Labels (CDPL), Difference in Acceptance Rates (DAR), Difference in Conditional Acceptance (DCA), Difference in Conditional Outcomes (DCR), Disparate (Adverse) Impact (DI), and Difference in Positive Proportions in Predicted Labels (DPPL).

### Bias Report

The screenshot shows the SageMaker Clarify interface with the 'Explanations' tab selected. The left sidebar includes sections for Overview, Metrics, Charts, Output Artifacts, Bias Reports, Explanations (selected), Debugs, Settings, Parameters, Configurations, Input Artifacts, Details, Information, and Metadata. The main content area displays a bar chart titled 'FEATURE IMPORTANCE' under the heading 'Explaining your model's predictions'. The chart plots SHAP values for various features: Country (~0.045), Education-Num (~0.035), Relationship (~0.035), Age (~0.035), Capital Gain (~0.03), Marital Status (~0.025), Occupation (~0.02), Hours per week (~0.02), Intltgt (~0.015), and Ethnic group (~0.01). A legend indicates that darker shades of blue represent higher SHAP values.

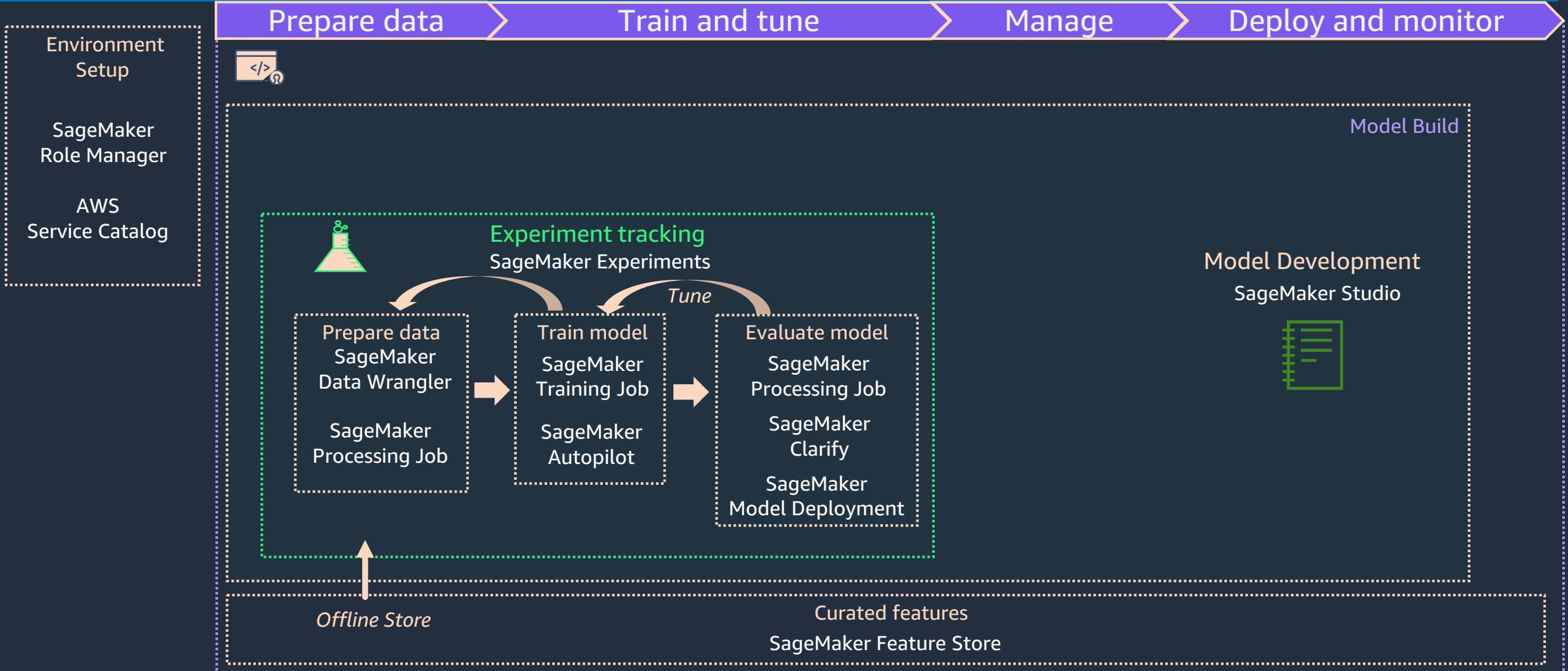
### Model Explainability



# MLOps Maturity Framework in Practice

Initial Phase Goal: Standardize environments & experiment tracking

Govern



## Initial Stage



## Key Technical Takeaways

- 1 Standardize data science environments for rapid experimentation with access to needed resources
- 2 Implement a feature store for feature discovery, reuse, and consistency
- 3 Implement a reliable mechanism/tooling for experiment tracking

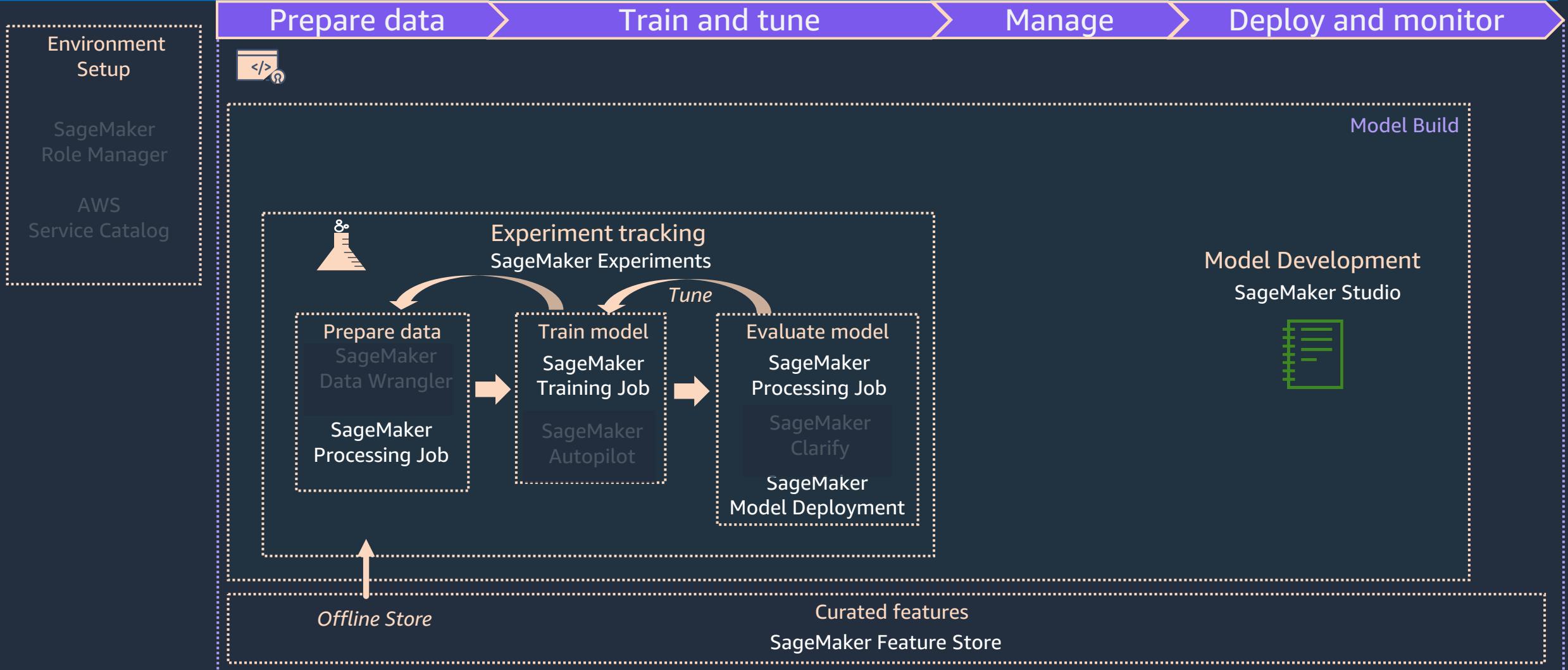
# Lab1: Initial

**Utilizing SageMaker Experiments  
during model development**



# Lab 1: Amazon SageMaker - Features Used

Govern



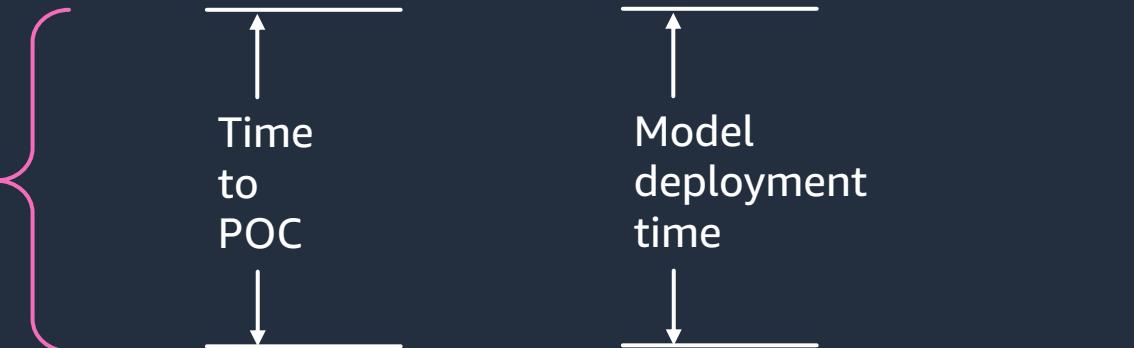
# Repeatable Stage:

**Implementing repeatable mechanisms to automate and govern your model build and model deployment tasks**

# MLOps Maturity Framework

## Scaling ML through your MLOps journey

Key Performance Indicators



Initial

Repeatable

Key Technology Capabilities

- Standardize experimentation environments & processes
  - Automate governed access to resources
  - Track experiments
  - Centralize management of feature data
- 
- Automate model build, and deployment tasks
  - Centralize management of models
  - Standardize source code & artifact repositories ➔



# MLOps Maturity Framework in Practice

Repeatable: Creating a repeatable path to production

**Develop** Candidate Model(s)

Business Goals  
& ML Problem

Provision  
Environments

Perform  
Experiments

**Operationalize** Workflow & Candidate Model(s)

Automate Model  
Training

Package and  
Test Models

Deploy



# Moving from Initial to Repeatable Maturity Level

## Repeatable: Creating a repeatable path to production



What if you have 50 models?  
100 models? Thousands of  
models?



How do you **document the model** and provide  
**model transparency** across stakeholders?



How do you **automate model build tasks** for  
faster iteration, feature re-use, automated  
lineage tracking & model re-training?



Establish **minimum levels of quality**?



Efficiently **manage model  
versions**?



Automate the path to production  
deployment?

# Reduce time to deployment

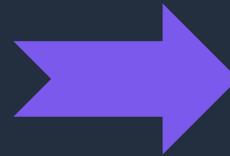
## Challenge #1: Reliable mechanism to document & share model information



Data  
Scientist



ML  
Engineer



### Model Documentation

- Evaluation Details
- Intended Use
- Training Details
- Risk Ratings



Data Owners



Business Stakeholder  
Product Owners



ML Risk Officer

### Key Technical Challenge

- Lack of model transparency across stakeholders impacting the ability to effectively manage ML use cases and models at scale
- Time consuming process to accurately capture model information

### Technical Solutions

- Centralize, standardize, and automate model documentation covering the complete ML use case lifecycle for scaling and governance

### On AWS

- Amazon SageMaker Model Cards

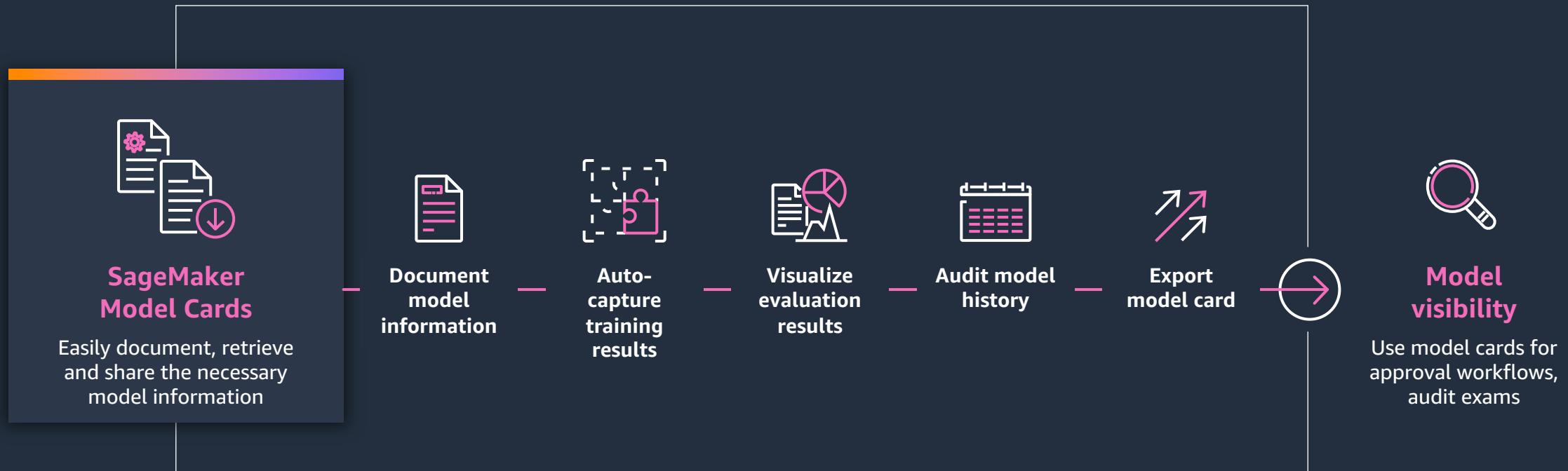


# Amazon SageMaker Model Cards

## Create single source of truth for model information

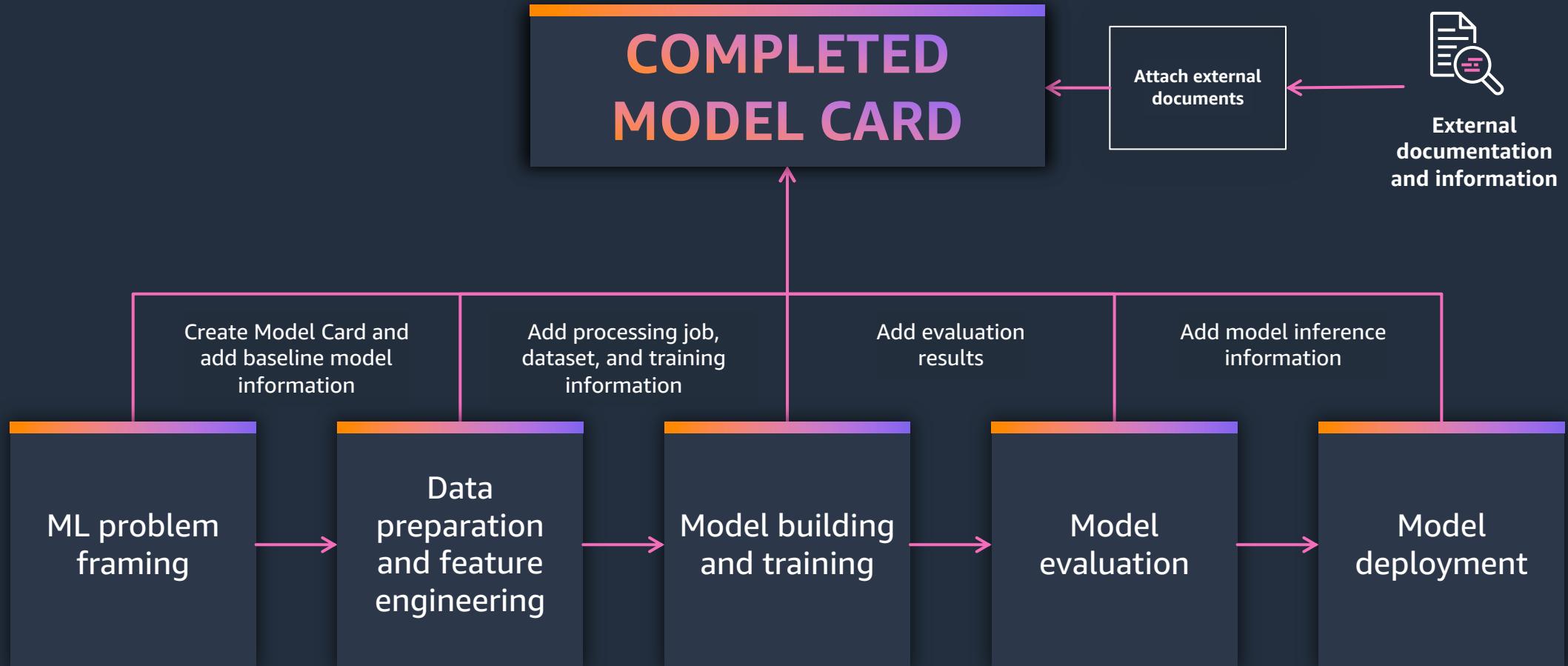
New!!

### How it works...



# Amazon SageMaker Model Cards

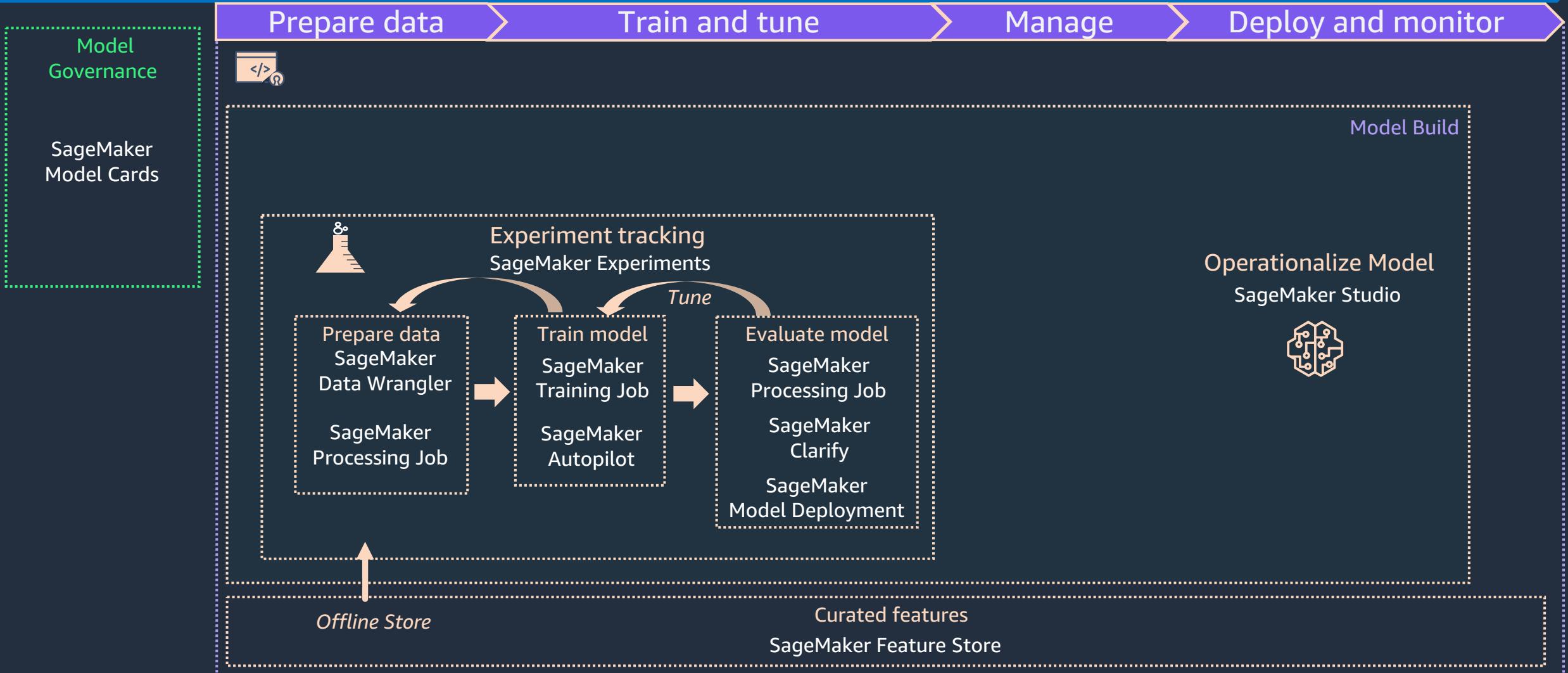
## Create single source of truth for model information



# MLOps Maturity Framework in Practice

Repeatable Phase Goal: Create a repeatable path to production

Govern



# Reduce time to deployment

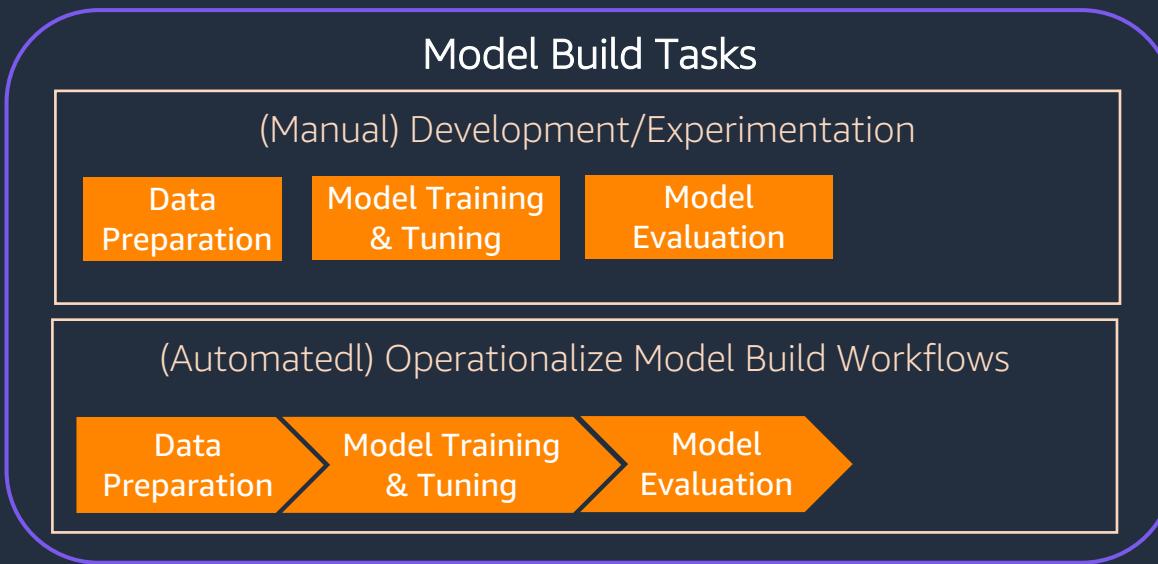
## Challenge #2: Rapidly experiment, deploy, and re-train



Data  
Scientist



ML  
Engineer



### Key Technical Challenges

- Enable rapid experimentation across workflow steps and tasks for model tuning
- Operationalize model build workflows for model re-training

### Technical Solutions

- Utilize workflow orchestration tooling to automate the steps in your model build workflows

### On AWS

- Amazon SageMaker Pipelines



# Amazon SageMaker Pipelines

Create automated model  
build workflows using the  
Pipelines python SDK



## Experiment management

Automatically log pipeline executions to SageMaker Experiments to track model version metadata & performance



## Debug and profile pipeline runs

Easy debugging of pipeline steps



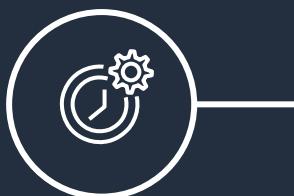
## Step Caching

Eliminate cycles recomputing previous pipeline steps



## Integrated Pipeline Steps

Native integration with common SageMaker features, AWS Services, and the ability to integrate with non-native steps



## Ability to schedule & trigger pipeline runs

Integration with EventBridge to easily trigger and schedule pipelines



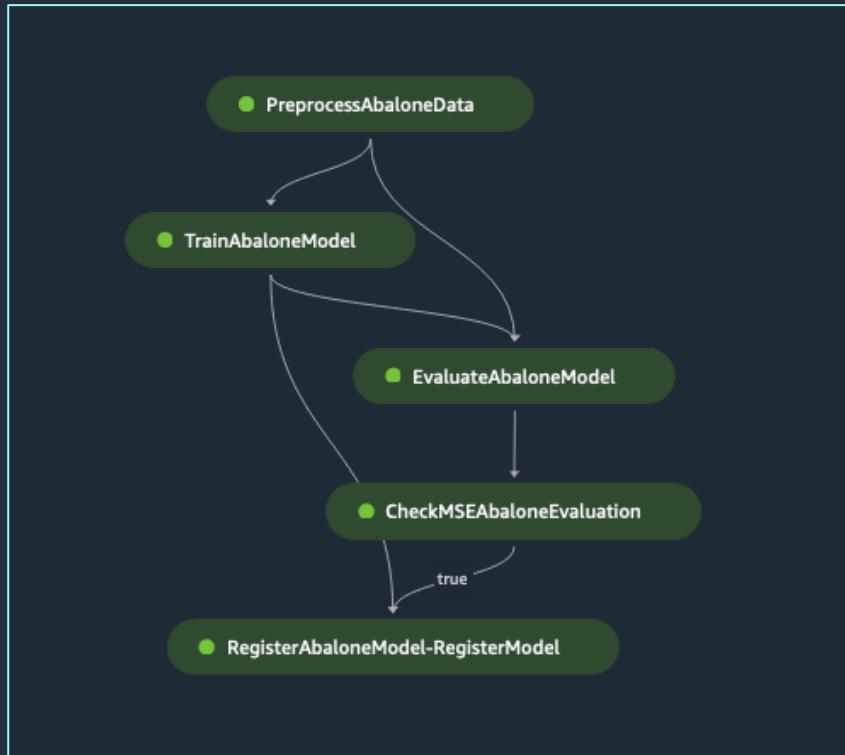
## No servers to manage

Pipeline capabilities offered without any additional servers to manage reducing operational overhead

# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

Programmatically create pipeline DAGs



### Supported Steps:

- **Processing:** Data Processing & Model Evaluation
- **Training:** Model Training using SageMaker Training Jobs
- **Tuning:** Hyperparameter Tuning Jobs
- **Condition:** Conditional step execution
- **Transform:** Batch predictions using SageMaker Batch Transform Jobs
- **Model:** Create and register model resource
- **Callback:** Incorporate additional tasks or AWS services into your workflow
- **Lambda:** Run a Lambda function
- **EMR:** Process tasks on a running EMR cluster
- **ClarifyCheck:** Run baseline drift checks against previous baselines.
- **QualityCheck:** Run baseline suggestions and drift checks against a previous baseline



# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

**Automatically collected input/output metadata for model lineage**

The screenshot illustrates the Amazon SageMaker Pipelines interface, demonstrating how it automatically collects input and output metadata for model lineage.

**Pipeline Execution Graph:** On the left, the "Graph" tab shows a workflow with five steps: PreprocessAbaloneData, TrainAbaloneModel, EvaluateAbaloneModel, CheckMSEAbaloneEvaluation, and RegisterAbaloneModel-RegisterModel. The "TrainAbaloneModel" step is highlighted with a blue border, indicating it is currently selected or active.

**Step Details:** A modal window for the "TrainAbaloneModel" step provides detailed information. It includes a summary bar at the top with the status (green dot), start time (3/17/2023, 12:00 PM), and elapsed time (11m22s). Below this, tabs for Input, Output, Logs, and Information are present, with "Input" being the active tab. The "Input" tab displays the following parameters and their values:

Parameters	Value
AlgorithmSpecification.TrainingInput	File
AlgorithmSpecification.TrainingOutput	683313688378.dkr.ecr.us-east-1.amazonaws.com
OutputDataConfig.S3OutputPath	s3://sagemaker-project-p-0qzy...
StoppingCondition.MaxRuntimeInSeconds	86400
ResourceConfig.InstanceCount	1
ResourceConfig.InstanceType	ml.m5.xlarge
ResourceConfig.VolumeSizeInGB	30
RoleArn	arn:aws:iam::962046451325:role/SageMaker-...
InputDataConfig.0.DataSource	S3Prefix
InputDataConfig.0.DataSource....	FullyReplicated
InputDataConfig.0.ContentType	text/csv
InputDataConfig.0.ChannelName	train
InputDataConfig.1.DataSource....	S3Prefix

**Output Tab:** The "Output" tab shows the following metrics and their values:

Metrics	Value
train:rmse	1.59429
validation:rmse	2.2268

**Files Tab:** The "Files" tab lists the output file and its source:

Files	Source
model.tar.gz	s3://sagemaker-project-p-0qzy...



# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

### Easily Debug Steps

less than a minute ago

### execution-1679832543

Graph    Parameters    Settings

The screenshot shows the execution details for 'execution-1679832543'. At the top, there's a summary card with the status (green circle), started time (3/17/2023, 12:00 PM), and elapsed time (11m22s). Below this is a 'Graph' tab showing a workflow with five steps: 'PreprocessAbaloneData', 'TrainAbaloneModel', 'EvaluateAbaloneModel', 'CheckMSEAbaloneEvaluation', and 'RegisterAbaloneModel-RegisterModel'. The 'TrainAbaloneModel' step is highlighted with a blue border. A search bar is also present. On the right, a detailed view of the 'TrainAbaloneModel' step is shown with tabs for Input, Output, Logs (which is selected and highlighted with an orange border), and Information. The 'Logs' tab displays CloudWatch log entries:

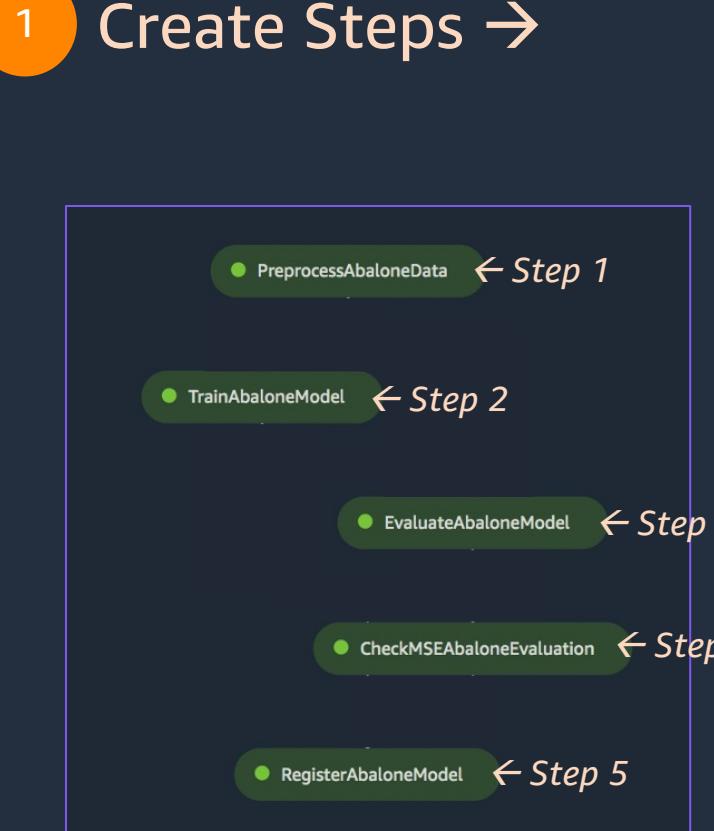
- 1679076461540 | [ 2023-03-17T18:07:41.540Z ] [49]#011train-rmse:1.59262#011validation-rmse:2.22809
- 1679076461540 | [ 2023-03-17T18:07:41.540Z ] [48]#011train-rmse:1.59429#011validation-rmse:2.22680
- 1679076461540 | [ 2023-03-17T18:07:41.540Z ] [47]#011train-rmse:1.61018#011validation-rmse:2.23053
- 1679076461540 | [ 2023-03-17T18:07:41.540Z ] [46]#011train-rmse:1.61224#011validation-rmse:2.23440
- 1679076461540 | [ 2023-03-17T18:07:41.540Z ] [45]#011train-rmse:1.61964#011validation-rmse:2.23879
- 1679076461540 | [ 2023-03-17T18:07:41.540Z ] [44]#011train-rmse:1.62642#011validation-rmse:2.23876



# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

### How it works



# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

### How it works...

#### 1 Create Steps →

```
from sagemaker.inputs import TrainingInput
from sagemaker.workflow.steps import TrainingStep

step_train = TrainingStep(
    name="TrainAbaloneModel",
    estimator=xgb_train,
    inputs={
        "train": TrainingInput(
            s3_data=step_process.properties.ProcessingOutputConfig.Outputs[
                "train"
            ].S3Output.S3Uri,
            content_type="text/csv"
        ),
        "validation": TrainingInput(
            s3_data=step_process.properties.ProcessingOutputConfig.Outputs[
                "validation"
            ].S3Output.S3Uri,
            content_type="text/csv"
        )
    }
)
```



Example:  
Training Step

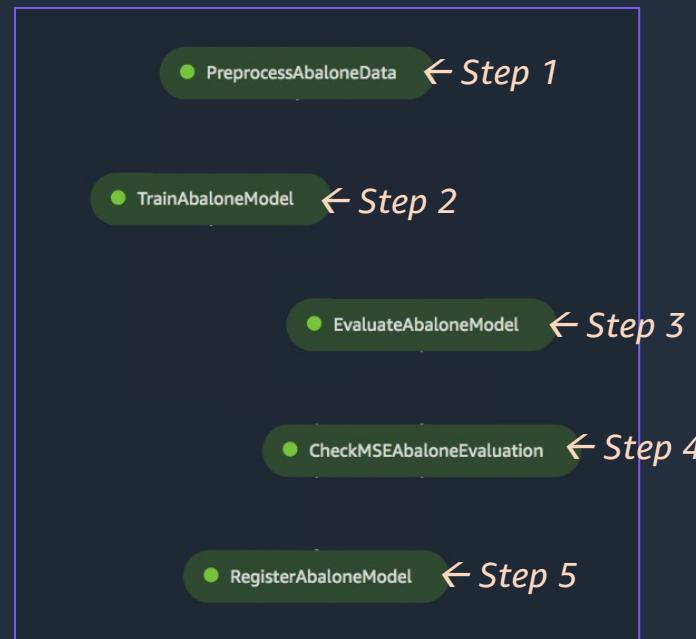


# Amazon SageMaker Pipelines

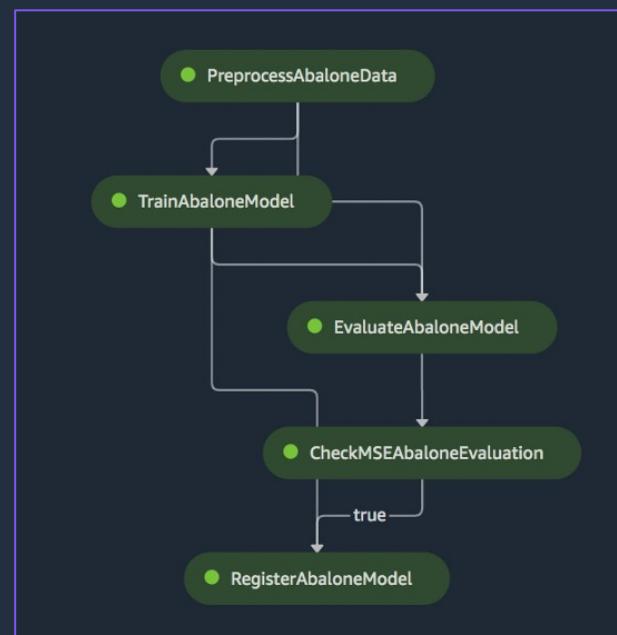
## Orchestrate steps & automate tasks in your ML workflows

### How it works...

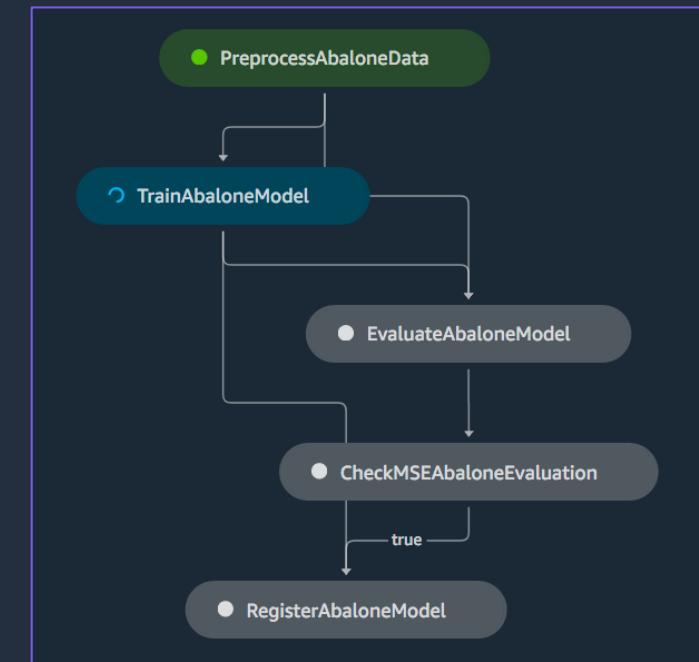
1 Create Steps →



2 Define Pipeline →



3 Start Pipeline →



# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

### How it works...

#### 2 Define Pipeline →

```
from sagemaker.workflow.pipeline import Pipeline

pipeline_name = f"AbalonePipeline"
pipeline = Pipeline(
    name=pipeline_name,
    parameters=[
        processing_instance_count,
        model_approval_status,
    ],
    steps=[step_process, step_train, step_eval, step_cond],
)

#Submit the pipeline definition
pipeline.upsert(role_arn=role)
```

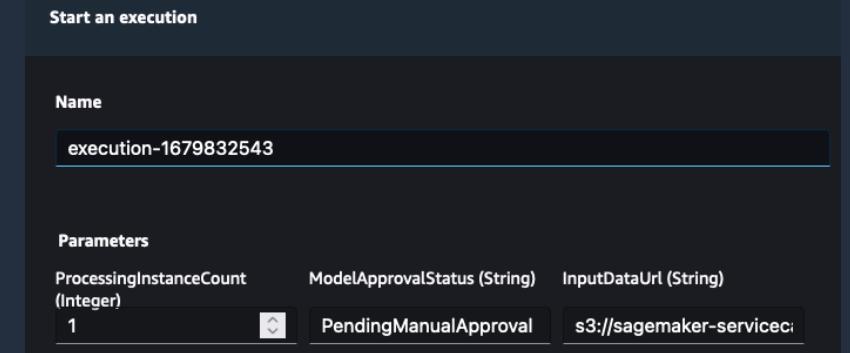
#### 3 Start Pipeline →

*programmatically*

```
execution = pipeline.start()
```

~OR~

*SageMaker Studio*

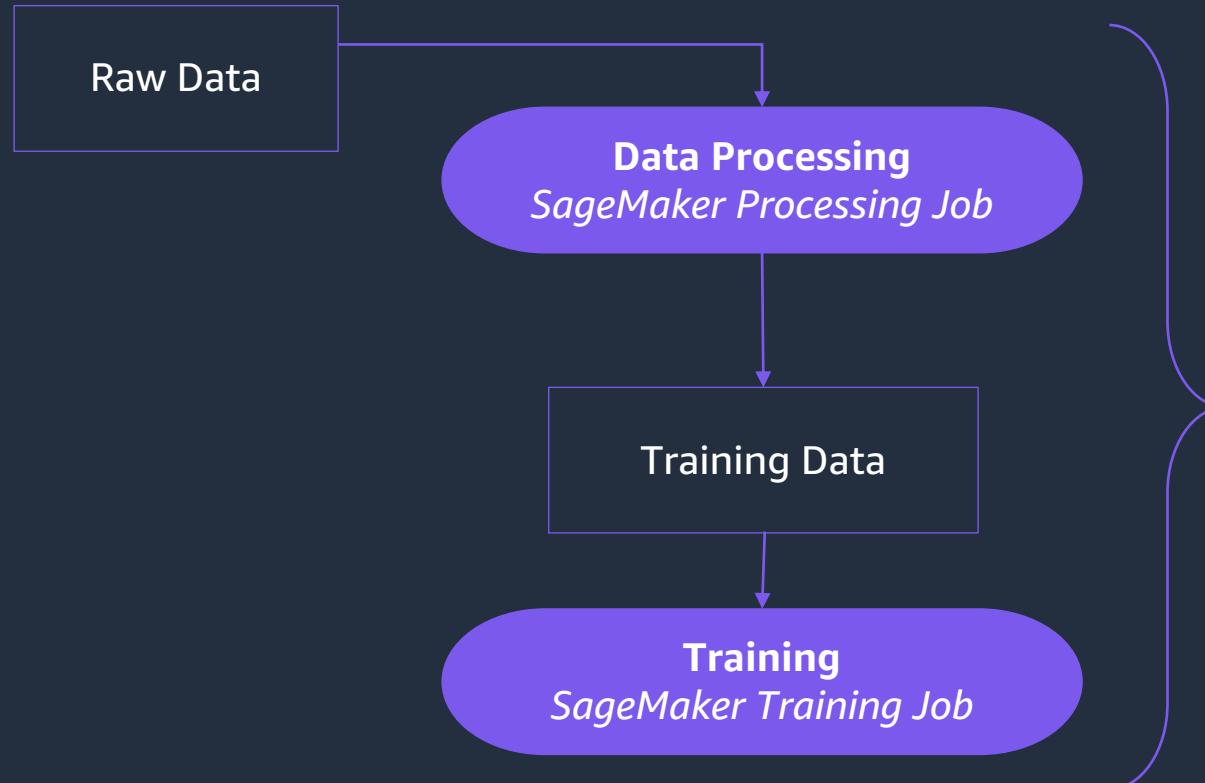


# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

### Feature Highlights: Step Caching

Cache successful results from previous step to avoid redundant processing →



What if you want to:  
- tune hyperparameters?  
- modify training code?

# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

### Feature Highlights: Runtime Parameters

Change parameters at runtime without changing your pipeline code →

1 Configure your parameter →

```
from sagemaker.workflow.parameters import (
    ParameterInteger,
    ParameterString,
    ParameterFloat
)

processing_instance_count = ParameterInteger(
    name="ProcessingInstanceCount",
    default_value=1
)
```

2

Pass in parameter on pipeline create →

```
pipeline = Pipeline(
    name=pipeline_name,
    parameters=[
        processing_instance_count
    ],
    steps=[step_process]
)
```

3

Optionally, pass a non-default value for pipeline execution →

```
execution = pipeline.start(
    parameters=dict(
        ProcessingInstanceCount=4
    )
)
```



# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

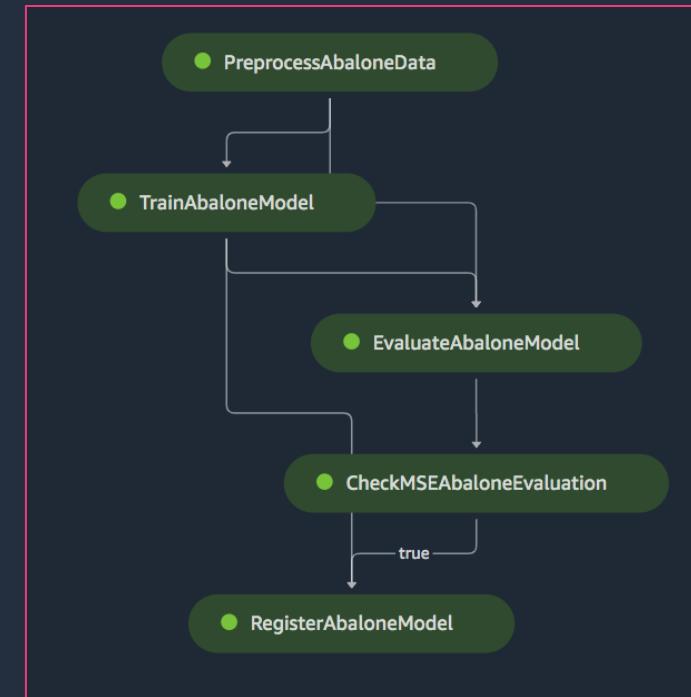
### Feature Highlights: Amazon EventBridge Integration

Setup automatic pipeline triggers based on events or schedules ➔



Event

Trigger Pipeline



# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

### Feature Highlights: Experiments Integration

Automatically capture each pipeline run as an experiment →

```
Pipeline(  
    name="MyPipeline",  
    parameters=[...],  
    pipeline_experiment_config=PipelineExperimentConfig(  
        ExecutionVariables.PIPELINE_NAME,  
        ExecutionVariables.PIPELINE_ID  
    ),  
    steps=[...]  
)
```

The screenshot shows the SageMaker console interface. On the left, there is a code editor window displaying Python code for defining a pipeline. On the right, the SageMaker resources page is shown, specifically the 'Experiments and trials' section. It lists a trial component named 'abalone-btd-1-p-xcfkqunms7sn' with a status of 'Completed'. Below this, a detailed view of the trial component is provided, including its creation time ('30 minutes ago') and status ('Completed'). At the bottom, a metrics table is displayed for the trial component, showing performance metrics like 'train:rmse' and 'validation:rmse' with their respective values.

Name	Minimum	Maximum	Standard Deviation	Average	Final value
train:rmse	1.5859	8.0792	1.2992065703803182	2.255249	1.59099
validation:rmse	2.04849	8.23054	1.2384372986921863	2.5348372000000006	2.08794



# Amazon SageMaker Pipelines

## Orchestrate steps & automate tasks in your ML workflows

### Deployment Support

For Real-Time/Near Real-Time Inference →

SageMaker Pipelines Option

Use Lambda  
or  
Callback Step

~OR~

Other Option

Deploy through CD tooling using IaC/CaC for  
rollback, consistency, and simplified resource  
management

For Batch Inference →

SageMaker Pipelines Option #1

Use Case: Training & Inference share the same lifecycle

Run batch inference as part of your  
training pipeline ~or~ as a separate  
pipeline

SageMaker Pipelines Option #2

Use Case: Training & Inference have a different lifecycle

Run training & batch inference in separate  
pipelines



# Amazon SageMaker Pipelines

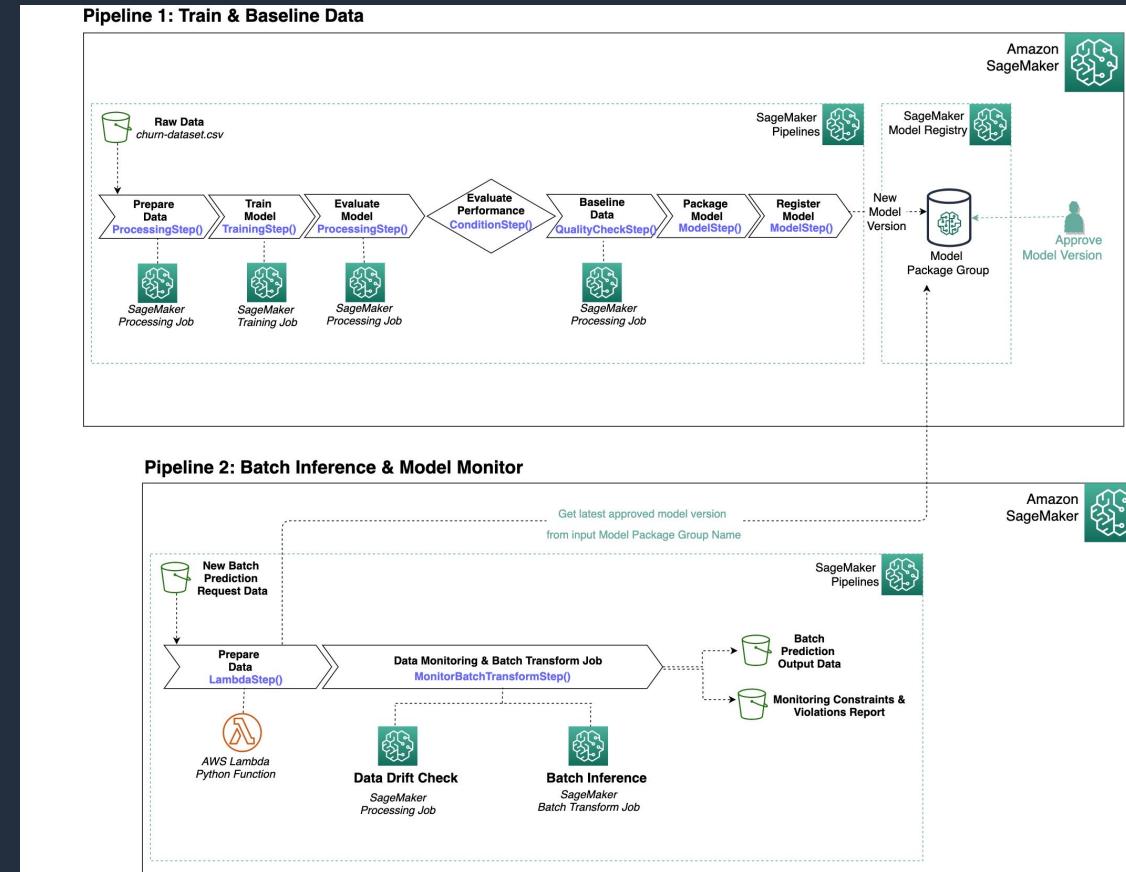
## Orchestrate steps & automate tasks in your ML workflows

### Using SageMaker Pipelines for Batch Inference

#### Option #2: Example

#### Use Case:

- Training & Inference have a different lifecycle
- Model is trained monthly
- Inference is performed daily



# Amazon SageMaker Pipelines

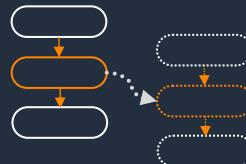
Orchestrate steps & automate tasks in your ML workflows

## What's new?



### Local mode

*Utilize local mode to quickly test your training, processing, and inference scripts*



### Cross account pipeline sharing

*Share pipeline entities across AWS accounts and access shared pipelines through API calls*



### SageMaker Pipelines SDK & SageMaker Studio UI Enhancements

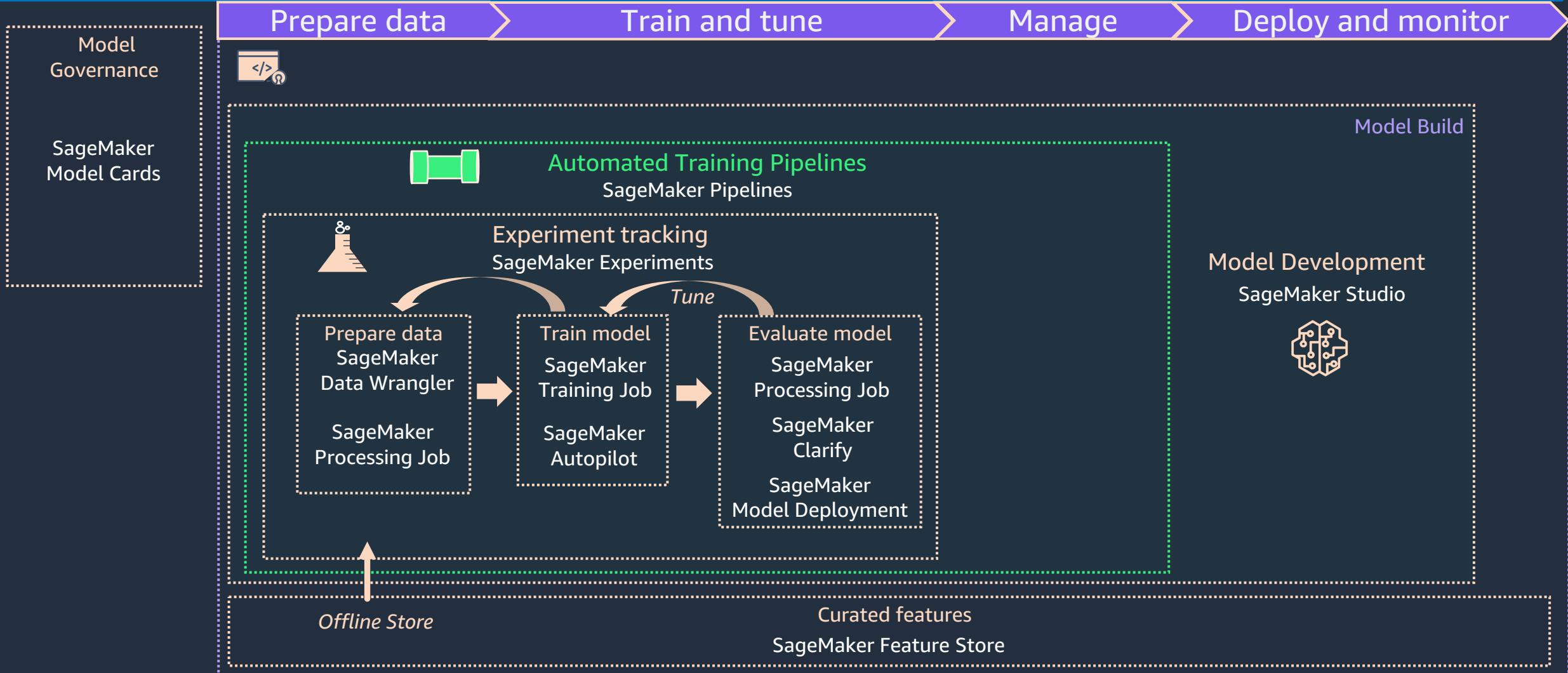
*Updates providing new features, enhanced ease-of-use, and DAG visualization*



# MLOps Maturity Framework in Practice

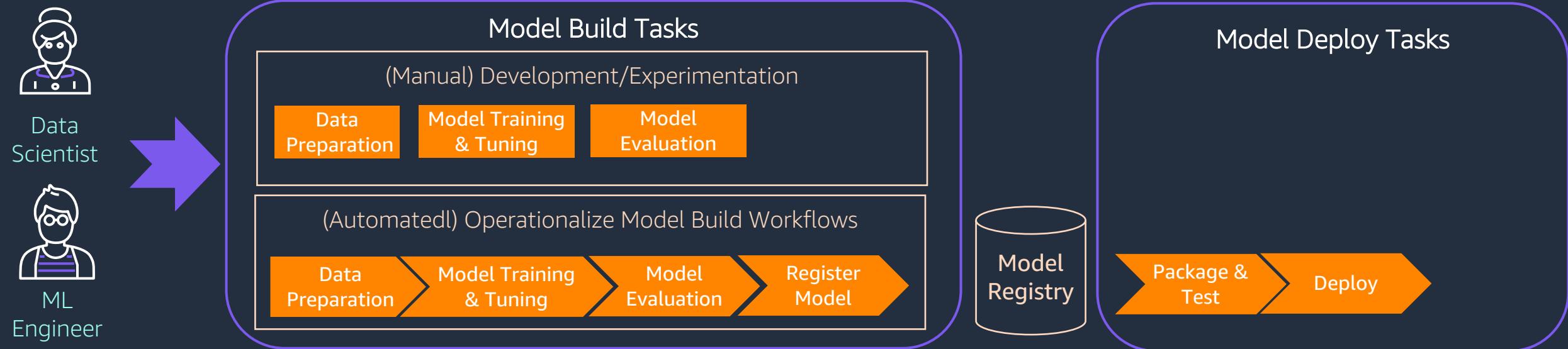
Repeatable Phase Goal: Create a repeatable path to production

Govern



# Reduce time to deployment

## Challenge #3: Ability to manage model versions at scale & streamline deployment



### Key Technical Challenges

- Manage candidate model versions at scale
- Manage the hand-off between model build and model deploy tasks/personas

### Technical Solutions

- Utilize a model registry to manage model versions (candidates for production deployment)
- Utilize a model registry to automate deployment workflows & approvals

### On AWS

- Amazon SageMaker Model Registry

# Amazon SageMaker Model Registry

## Centrally catalog & manage model versions

Model registry

Show introduction

+ Create model version + Create model group

Model group name	Description	Status	Created on	Created by
abalone-btdm-builtin-p-qdv3j9welp9x	--	✓ Completed	3 hours ago	--
abalone-btd-builtin-p-0qzy6hduz67s	--	✓ Completed	3 hours ago	--
aim321-customer-churn	--	✓ Completed	4 months ago	--

End of results

- Associate metadata with a model
- Track model performance metrics

VERSION 1

Status	Pipeline	Execution	Project	Last Stage	Model group	Share	Update status
Pending	abalone-btd-builtin-p-...	execution-167906962...	abalone-btd-builtin	staging	abalone-btd-builtin-p-...		

Activity    Model quality    Explainability    Bias report    Inference recommender    Load test    Settings

Model metric	Metric value	Standard deviation
mse	4.573587555999638	2.135589655997731



# Amazon SageMaker Model Registry

## Centrally catalog & manage model versions

The screenshot shows the Amazon SageMaker Model Registry interface. At the top, there's a navigation bar with tabs for 'Status' (Pending), 'Pipeline' (abalone-btd-builtin-p...), 'Execution' (execution-1679832543), 'Project' (abalone-btd-builtin), 'Model group' (abalone-btd-builtin-p...), 'Share', and 'Update status'. Below the navigation bar, there are several tabs: 'Activity' (selected), 'Model quality', 'Explainability', 'Bias report', 'Inference recommender', 'Load test', and 'Settings'. Under the 'Activity' tab, there's a table with columns: Event type, Event, Comment, Modified by, Last modified, and Actions. A single row is shown: Approval, Status updated to Pend..., (comment), 2 hours ago, ..., and an ellipsis for actions.

- Manage the approval status of a model
- Deploy models to production using Continuous Delivery or Continuous Deployment (with Projects)

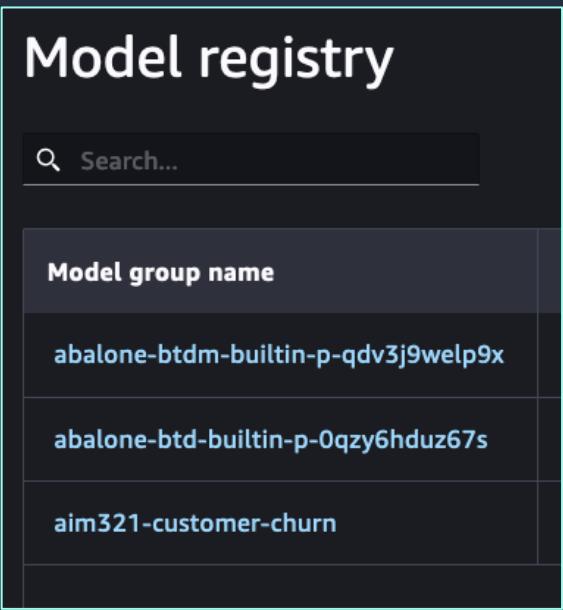
The dialog box is titled 'Update model version status'. It contains a descriptive text: 'Update the model status and add comments. If this model group has a deployment pipeline, the new model version is deployed after it's approved'. Below this is a 'Status' dropdown menu. The menu has four options: 'Approved', 'Rejected', and 'Pending' (which is selected). At the bottom of the dialog are 'Cancel' and 'Update status' buttons.



# Amazon SageMaker Model Registry

Centrally catalog & manage model versions

## Key Concepts



### Hierarchy →



v1



v2



v1

# Amazon SageMaker Model Registry

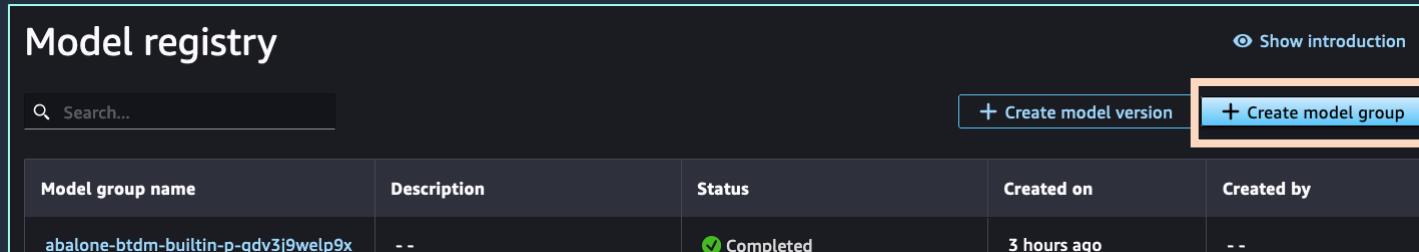
Centrally catalog & manage model versions

## How it works

1

Create Model Group →

*SageMaker Studio*



The screenshot shows the 'Model registry' section of the SageMaker Studio interface. At the top right, there is a 'Show introduction' link and two buttons: '+ Create model version' and '+ Create model group'. The '+ Create model group' button is highlighted with a blue box. Below these buttons is a table with one row. The table columns are: 'Model group name', 'Description', 'Status', 'Created on', and 'Created by'. The row contains the value 'abalone-btdm-builtin-p-qdv3j9welp9x' for 'Model group name', '--' for 'Description', '✓ Completed' for 'Status', '3 hours ago' for 'Created on', and '--' for 'Created by'.

Model group name	Description	Status	Created on	Created by
abalone-btdm-builtin-p-qdv3j9welp9x	--	✓ Completed	3 hours ago	--

~OR~

*programmatically*

```
create_model_package_group_response = sm_client.create_model_package_group(**model_package_group_input_dict)
print('ModelPackageGroup Arn : {}'.format(create_model_package_group_response['ModelPackageGroupArn']))
```



# Amazon SageMaker Model Registry

## Centrally catalog & manage model versions

### How it works...

2

Register a Model Version →

*SageMaker Studio*

The screenshot shows the SageMaker Studio Model registry interface. At the top, there is a search bar and two buttons: '+ Create model version' (highlighted with a red box) and '+ Create model group'. Below is a table with columns: Model group name, Description, Status, Created on, and Created by. One row is visible: 'Model group name' is 'abalone-btdm-builtin-p-qdv3j9welp9x', 'Description' is '--', 'Status' is 'Completed' with a green checkmark, 'Created on' is '3 hours ago', and 'Created by' is '--'. There is also a 'Show introduction' link.

Model group name	Description	Status	Created on	Created by
abalone-btdm-builtin-p-qdv3j9welp9x	--	✓ Completed	3 hours ago	--

~OR~

*programmatically*

```
create_model_package_response = sm_client.create_model_package(**create_model_package_input_dict)
model_package_arn = create_model_package_response["ModelPackageArn"]
print('ModelPackage Version ARN : {}'.format(model_package_arn))
```

~OR~

*within a Pipeline, using RegisterModel Step:*

```
step_register = RegisterModel(
    name="RegisterAbaloneModel",
    estimator=xgb_train,
    model_data=step_train.properties.ModelArtifacts.S3ModelArtifacts,
    content_types=["text/csv"],
    response_types=["text/csv"],
    inference_instances=["ml.t2.medium", "ml.m5.large"],
    transform_instances=["ml.m5.large"],
    model_package_group_name=model_package_group_name,
    approval_status=model_approval_status,
    model_metrics=model_metrics,
)
```



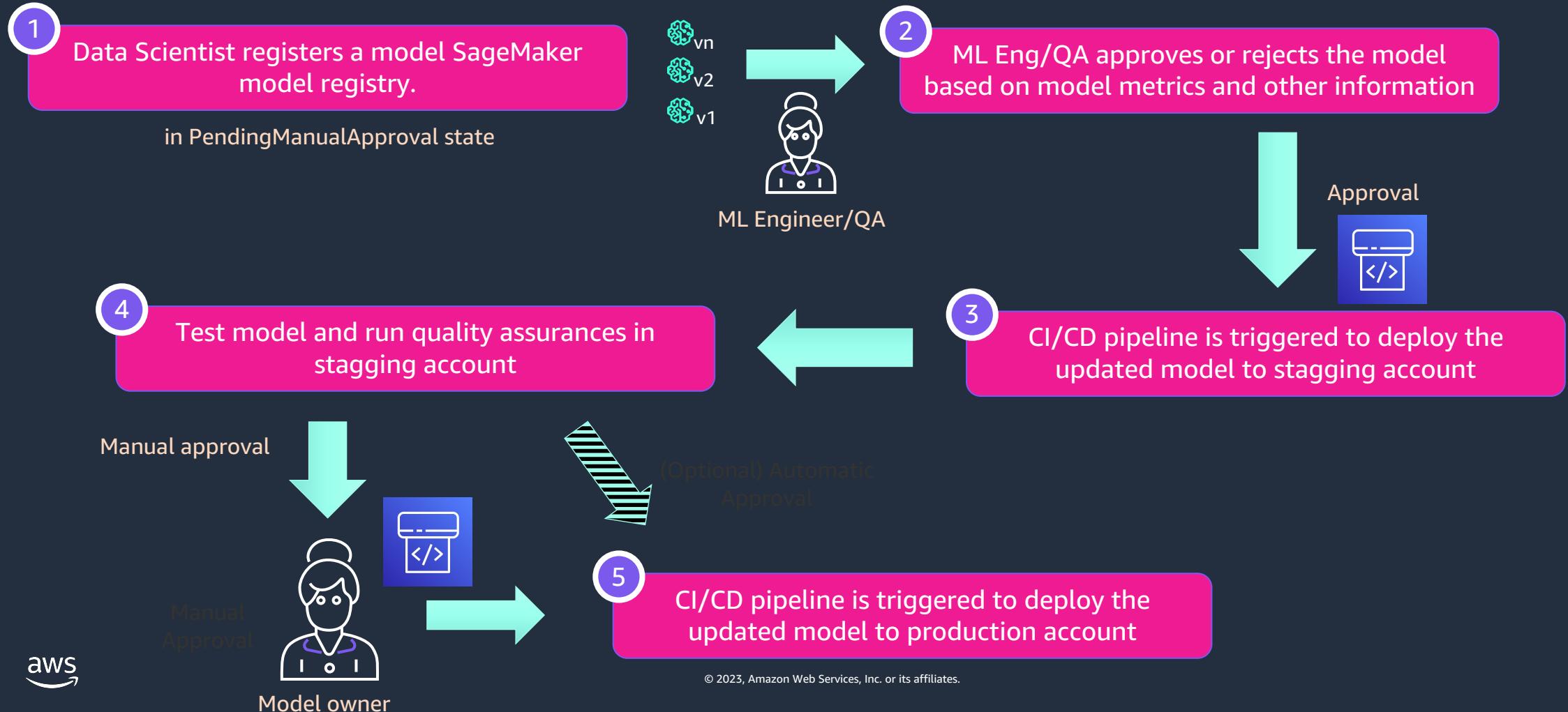
● RegisterAbaloneModel



# Amazon SageMaker Model Registry

## Centrally catalog & manage model versions

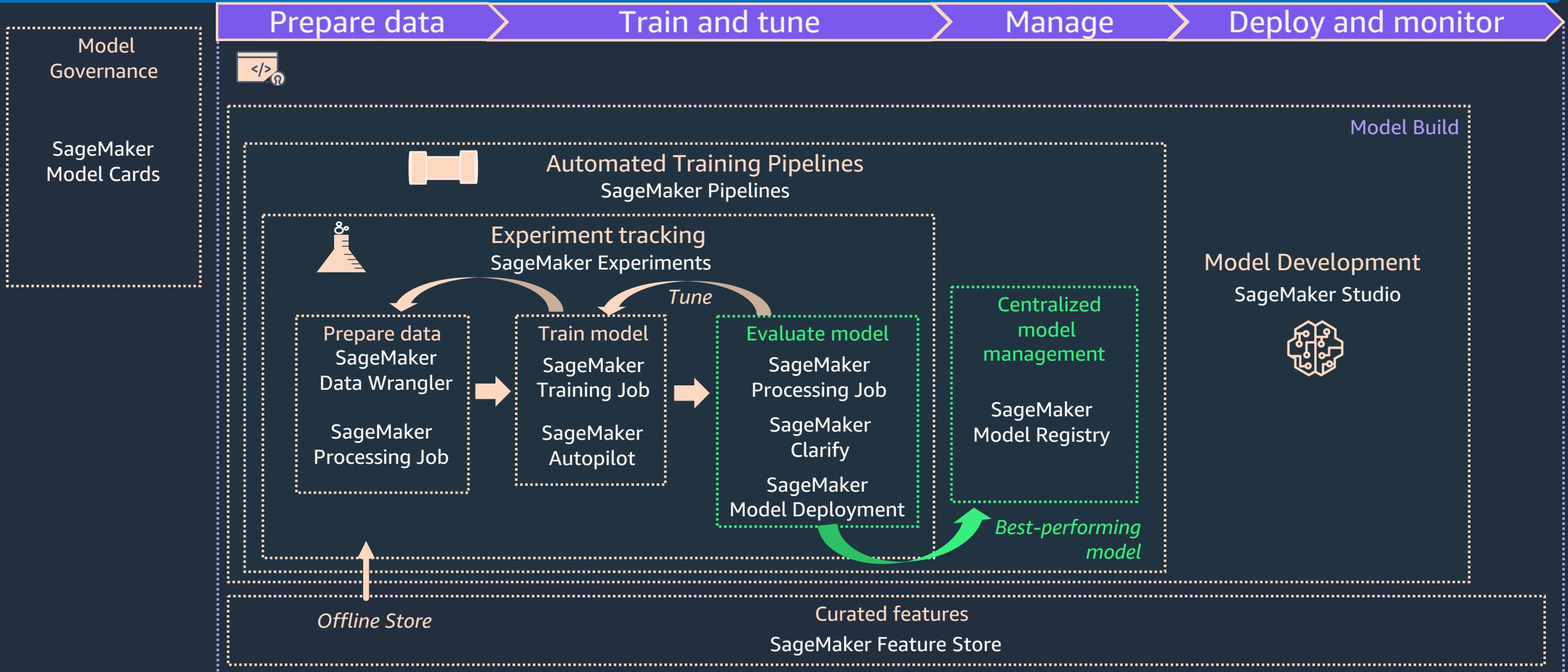
### Example Workflow



# MLOps Maturity Framework in Practice

Repeatable Phase Goal: Create a repeatable path to production

Govern



## Repeatable Stage



## Key Technical Takeaways

- 1
- 2
- 3

Consider a repeatable mechanism to document model information that is transparent across stakeholders

Automate model build workflows for rapid development iterations, and model re-training

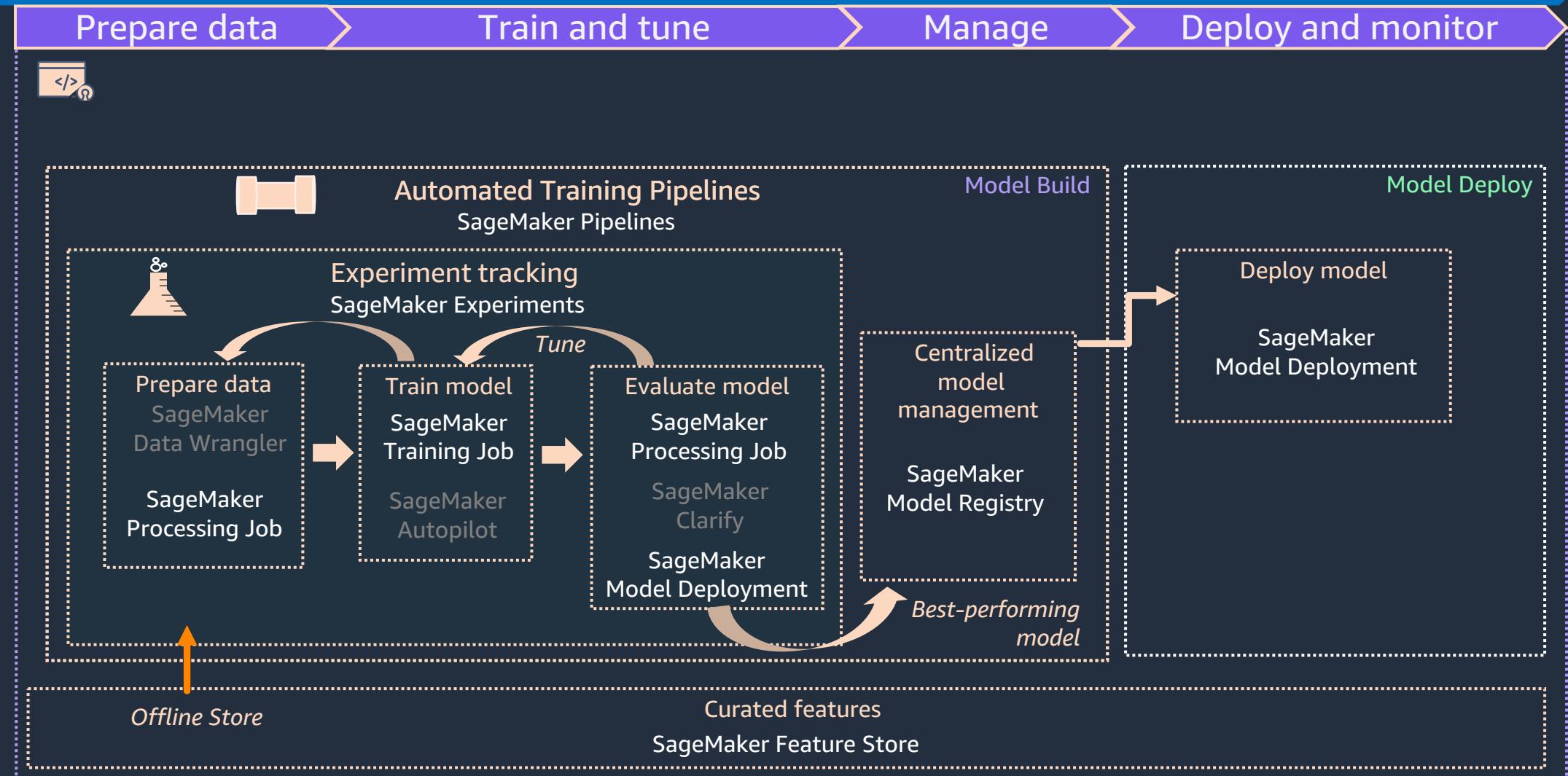
Implement a model registry to manage model versions at scale and bridge the gap between model build and model deploy activities

# Lab2: Repeatable

**Automate model build workflows  
& Manage model versions at scale**

# Lab 2: Amazon SageMaker - Features Used

Govern

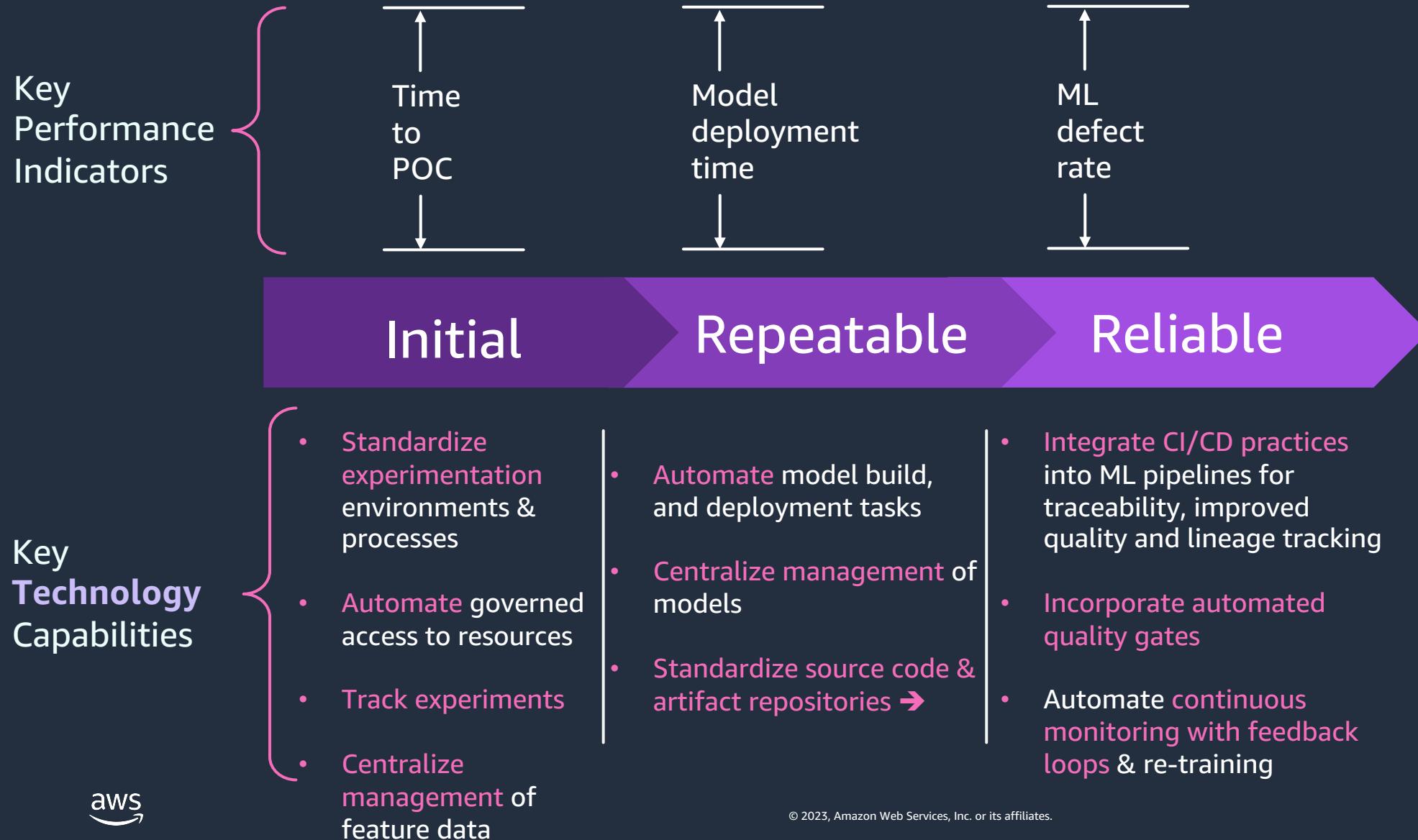


# **Reliable Stage:**

**Implementing reliable mechanisms to build, deliver, and manage models**

# MLOps Maturity Framework

## Scaling ML through your MLOps journey



# MLOps Maturity Framework in Practice

Reliable: Creating managed workflows optimized for quality & traceability

## Develop Candidate Model(s)

Business Goals  
& ML Problem

Provision  
Environments

Perform  
Experiments

## Operationalize Workflow & Candidate Model(s)

Feedback and Retraining

Automate Model  
Training

Package and  
Test Models

Deploy  
& Monitor

End-to-end Quality Gates & Traceability

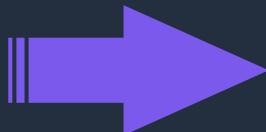


# Moving from Repeatable to Reliable Maturity Level

Reliable: Creating managed workflows optimized for quality & traceability



How can I ensure quality of  
ML workflows and models  
running in production  
environments?



Enforce **version control** your code, data and models



Implement **continuous deployment/delivery pipelines**



Consider **advanced deployment strategies** to validate model performance and reduce risk



Create **automated quality gates** and **automate business process** spanning end-to-end workflows

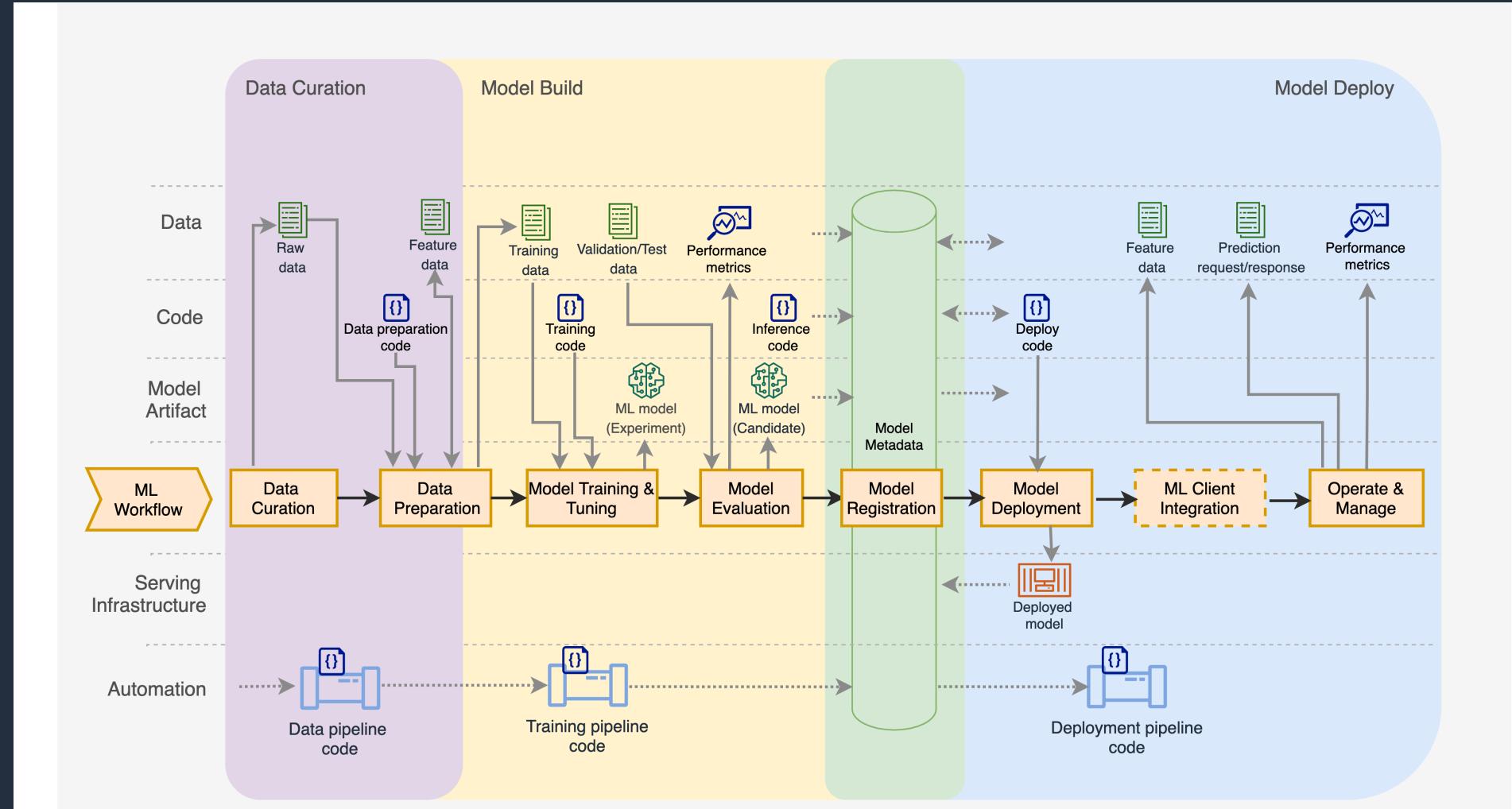


Implement **continuous monitoring** and dashboards for transparency into model performance

# Moving from Repeatable to Reliable Maturity Level

Reliable: Creating managed workflows optimized for quality & traceability

Without standardization, automation, and CI/CD, this is A LOT to track and manage to be able to reliably reproduce and manage

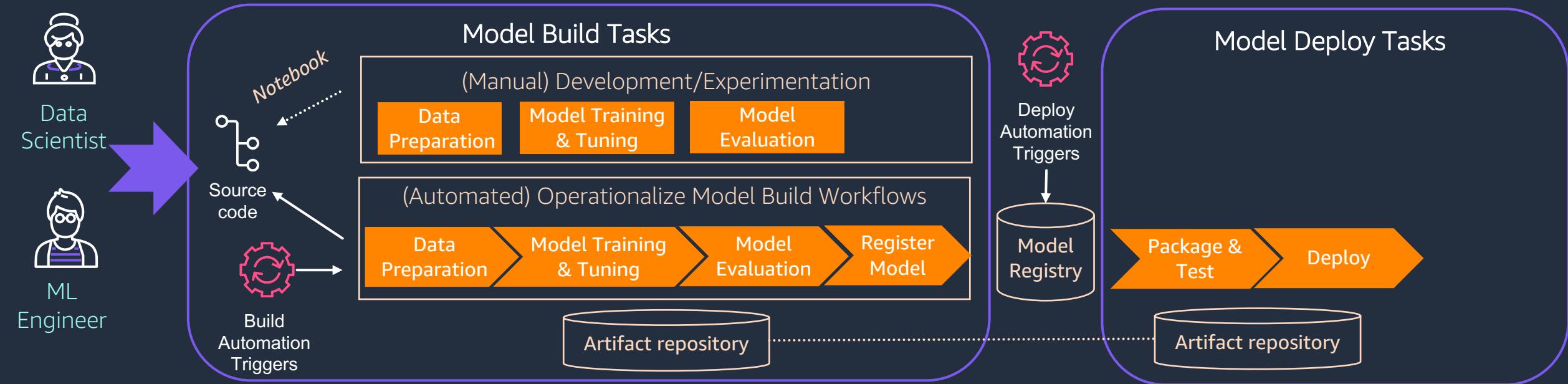


# **Reliable Stage:**

## **Standardizing ML Projects**

# Creating a reliable path to production

## Challenge #1: Standardize ML workflows and reduce manual hand-offs



### Key Technical Challenges

- Identifying, implementing, and enforcing required quality gates & business processes across ML lifecycle
- Automate the hand-off between model build and model deployment tasks

### Technical Solutions

- Implement CI/CD practices that consider implications unique to machine learning

### On AWS

- Amazon SageMaker Projects

# Amazon SageMaker Projects

**Standardize machine learning workflows through MLOps project templates**



## Incorporate CI/CD practices

Automatically create standardized workflows that incorporate CI/CD practices through preconfigured templates



## Standardize workflows across teams

Create custom workflows specific to your own processes and standards

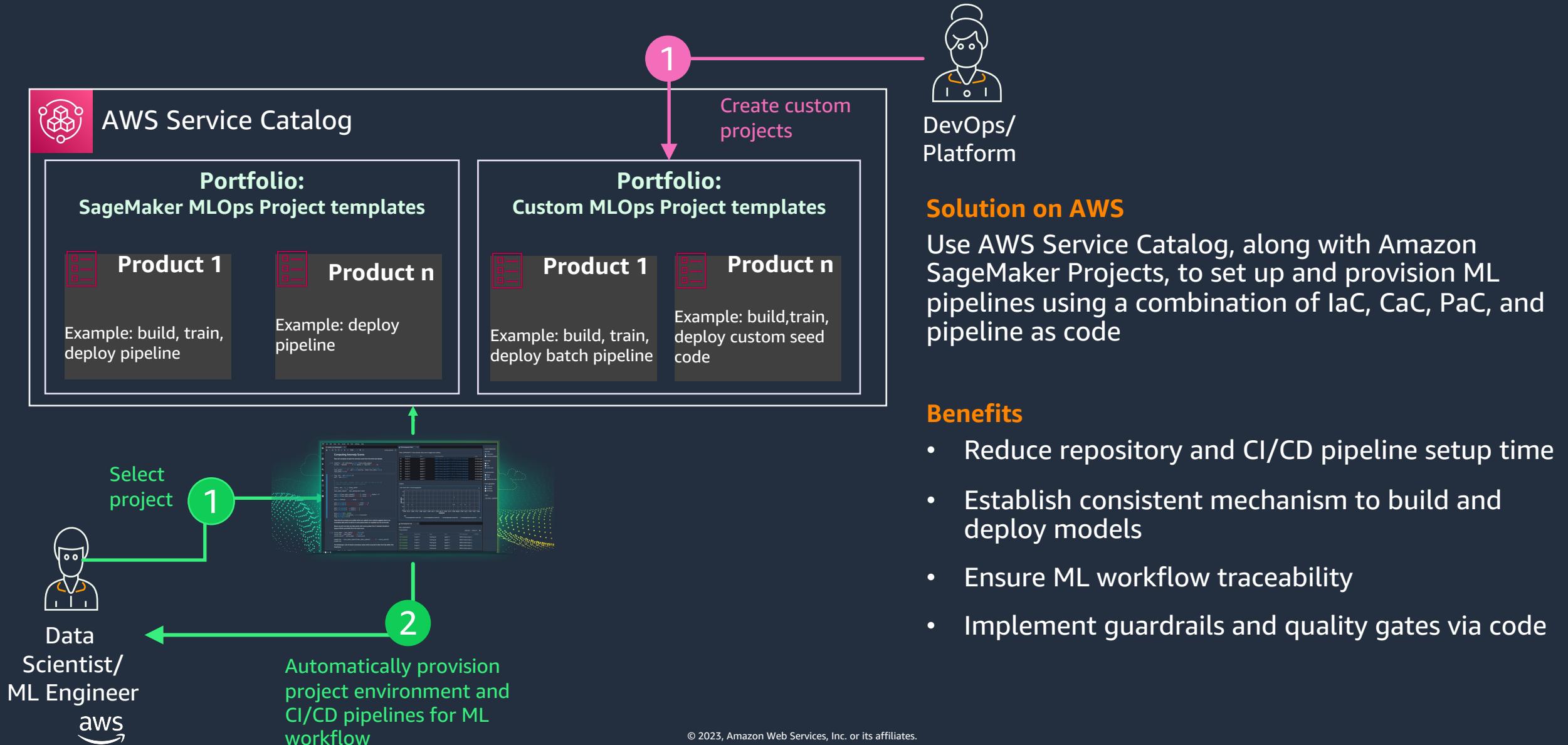


## Integrate with 3<sup>rd</sup> party tooling

Ability to integrate with common 3<sup>rd</sup> party tooling used for CI/CD

# Amazon SageMaker Projects

## Standardize ML workflows using MLOps templates



# Amazon SageMaker Projects

## Standardize ML workflows using MLOps templates

### Template Options

#### Built-In: SageMaker Managed

1. Build, Train, Deploy
  - Using AWS Developer Tools (AWS CodeCommit, AWS CodePipeline)
  - Using AWS CodePipeline & 3<sup>rd</sup> Party Git Repos
  - Using Jenkins & 3<sup>rd</sup> Party Git Repos
2. Build (Image Build), Train, Deploy
3. Build, Train
4. Deploy

#### Organization: Custom

- Create custom project to implement your own organizational standards for common workflow patterns
- Common reasons to create custom projects:
  - Workflows spanning multiple AWS accounts
  - Customize workflows for common patterns
  - Standardize source code repositories
  - Standardize workflows using custom seed code

TIP

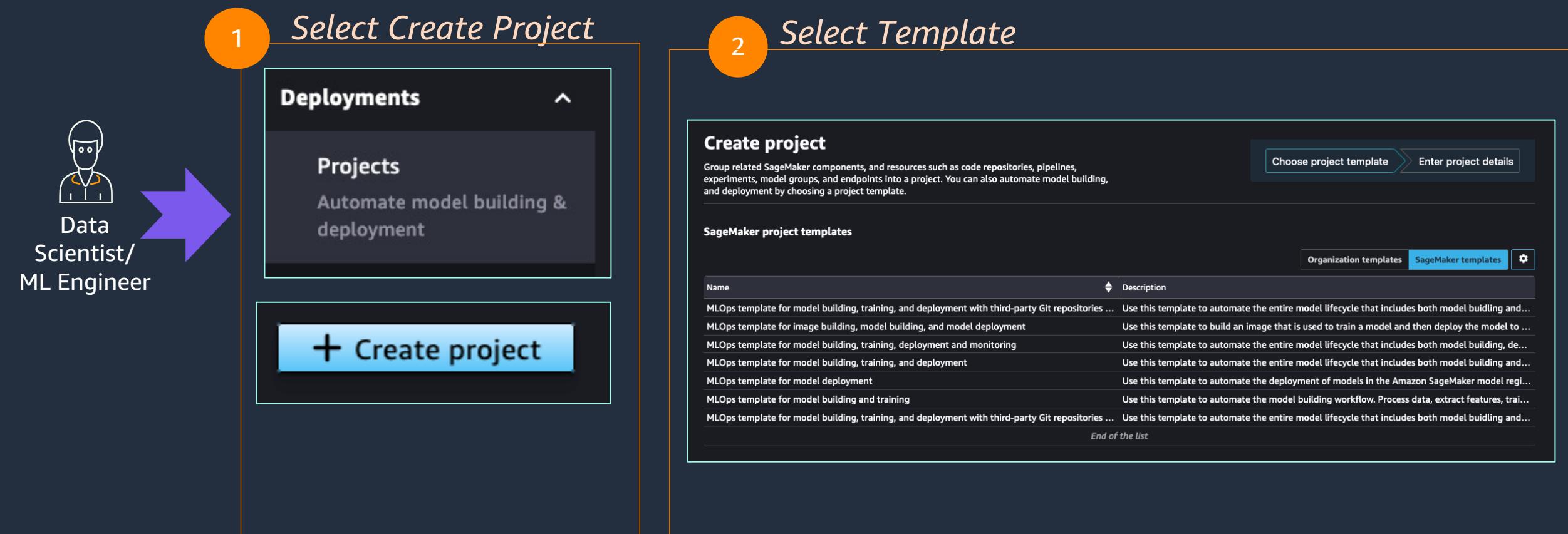
When creating custom Project templates, start with one of the built-in templates, or check out common custom project template examples here:  
<https://github.com/aws-samples/sagemaker-custom-project-templates>



# Amazon SageMaker Projects

## Standardize ML workflows using MLOps templates

### How it works



# Amazon SageMaker Projects

## Standardize ML workflows using MLOps templates

### How it works...

3

#### Enter Project Details

##### Create project

Group related SageMaker components, and resources such as code repositories, pipelines, experiments, model groups, and endpoints into a project. You can also automate model building, and deployment by choosing a project template.

##### Project details

Please provide the following details for your project:

###### Name

builtin-buildtraindeploy-abalone



- Totally unique, and never used before in your AWS account, or AWS region.
- Does not contain spaces. Uses hyphens (-) instead. Can contain lowercase letters, uppercase letters, and numbers (e.g. a - z, A - Z, 0 - 9).
- Between 1 - 32 characters in length.

###### Description - optional

Instantiate project using built-in template

###### Tags - optional

Key	Value	Remove
team	marine-life	

- Ability to add any custom tags
- Tag propagation: Resources provisioned through built-in templates will have project-name & project-id Tags assigned. Example(s):
  - AWS CodePipeline → Pipelines
  - AWS CodeBuild → Build Projects
  - SageMaker Pipelines → Pipeline
  - → Training & Processing Jobs (via AWS CodeBuild)

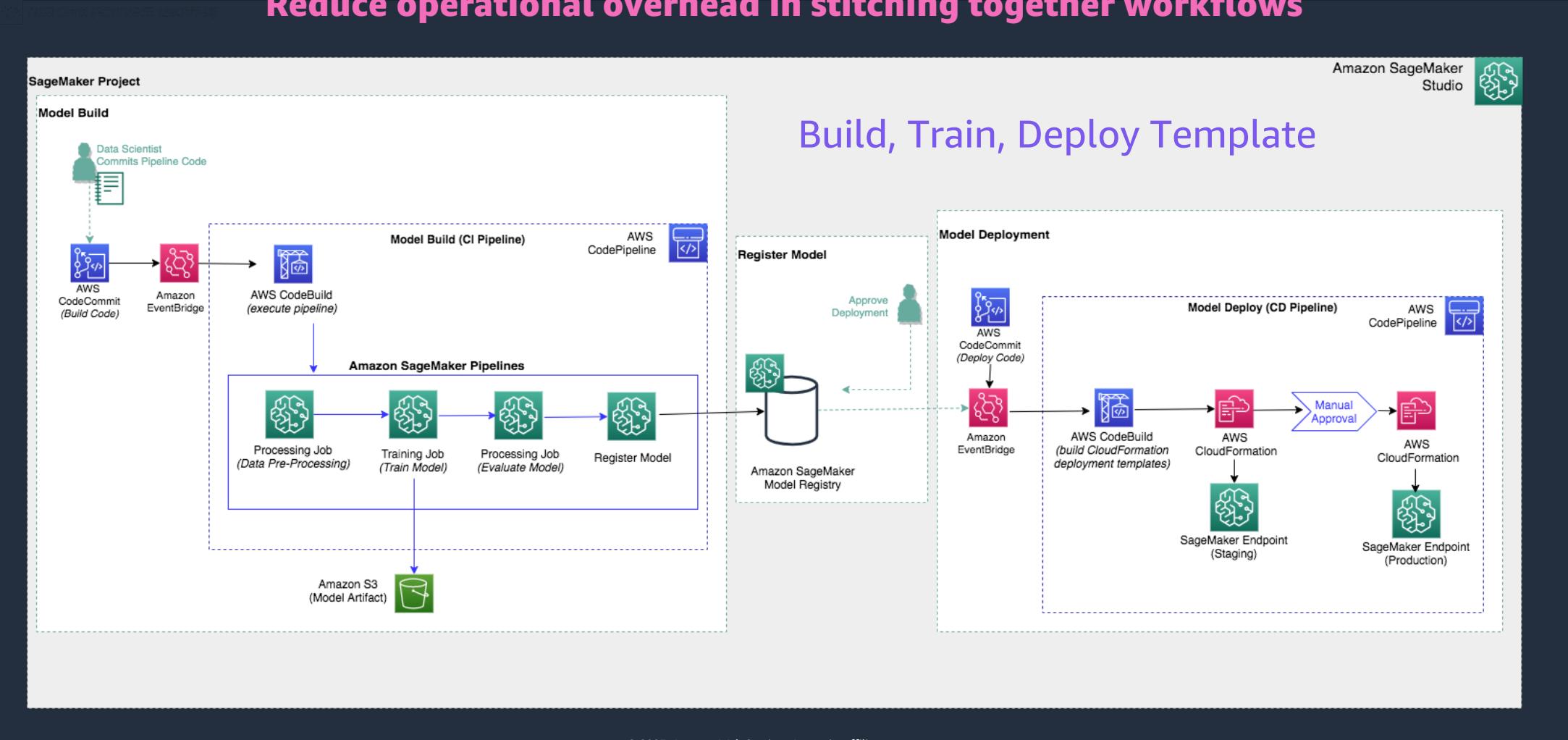
Everything needed to create a CI/CD Pipeline for Machine Learning gets automatically provisioned & configured for you.....



# Amazon SageMaker Projects

## Standardize ML workflows using MLOps templates

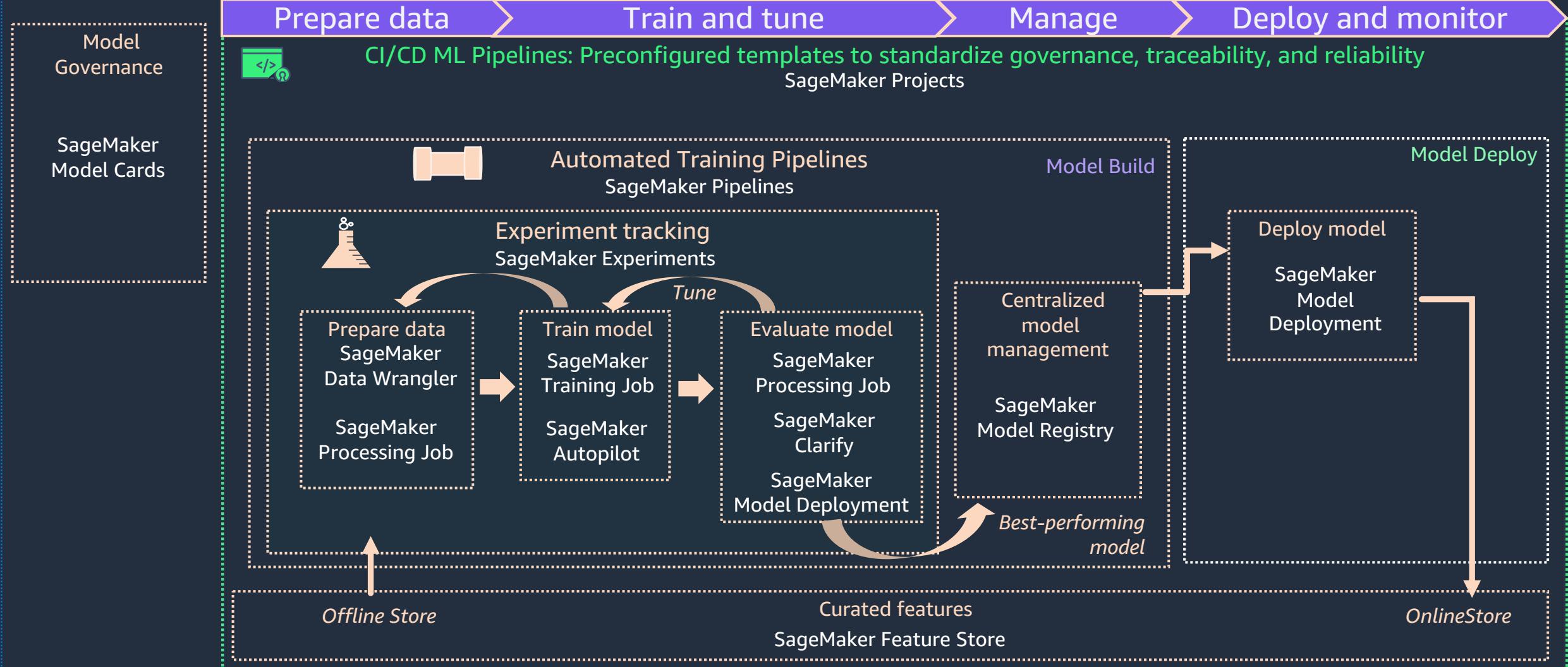
**Automatic creation & configuration of CI/CD ML Pipeline**  
**Reduce operational overhead in stitching together workflows**



# MLOps Maturity Framework in Practice

Reliable Phase Goal: Creating managed workflows optimized for quality & traceability

Govern

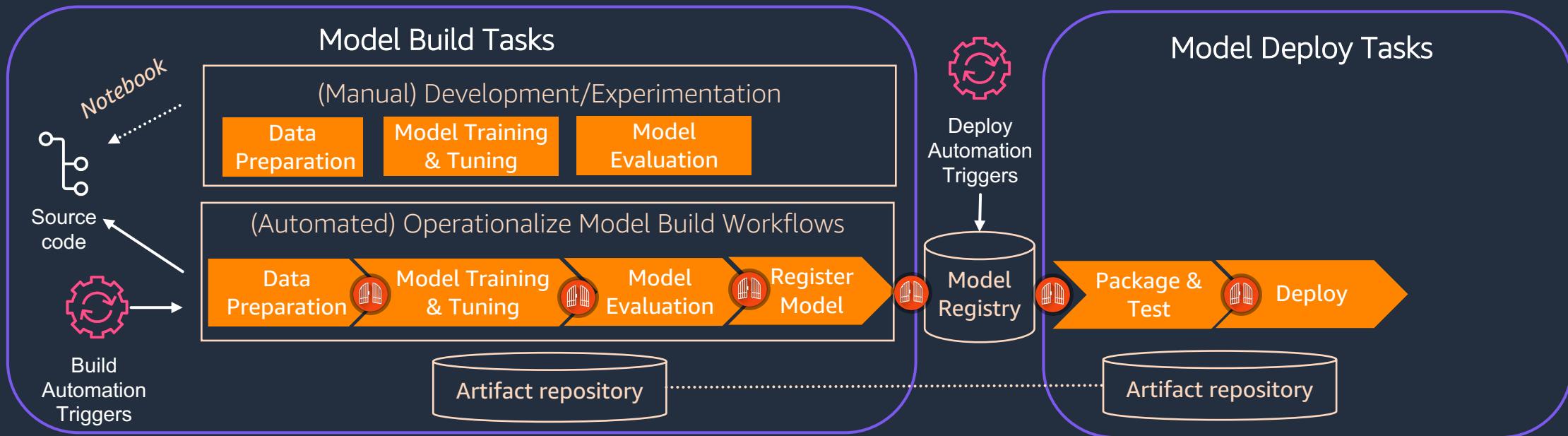


# **Reliable Stage:**

# **MLOps Quality Gates**

# Implement automated quality gates across MLDC

## Challenge #2: Identify and implement quality gates across the MLDC



### Key Technical Challenges

- Identifying, implementing, and enforcing required quality gates & business processes across ML lifecycle

### Technical Solutions

- Incorporate additional quality gates in your CI/CD pipeline

### On AWS

- Amazon SageMaker Projects, Pipelines with built-in and custom steps
- Amazon SageMaker Clarify



# Considerations for quality gates

Reliable: Creating managed workflows optimized for quality

## Data Quality Assurance

Before the model is trained, consider automated data validation for ➔



---

Data format - file format,  
structure, naming



Data Quality – missing values,  
feature distribution, feature  
completeness



---

Label annotation and data  
feasibility

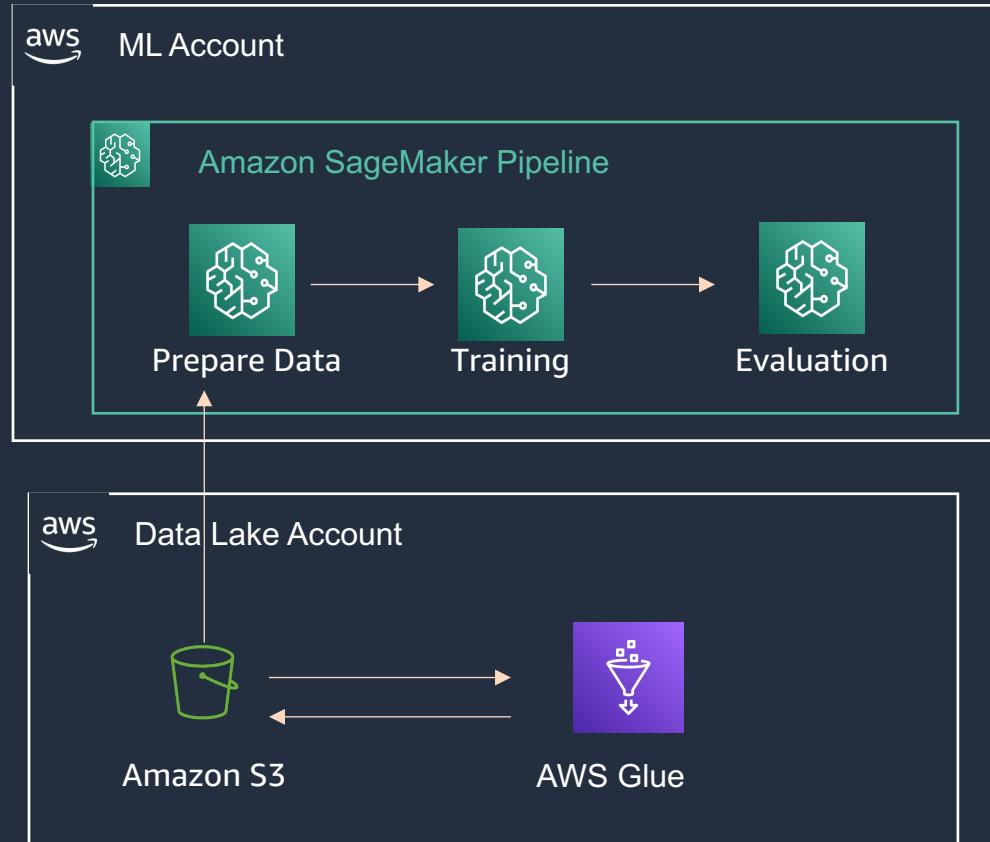
# Considerations for quality gates

Reliable: Creating managed workflows optimized for quality

New!!

## AWS Glue Data Quality continuous monitoring

DELIVER HIGH QUALITY DATA ACROSS YOUR DATA LAKES AND DATA PIPELINES



- Automatic data quality rule recommendations based on your data
- Keep data quality high with ongoing data analysis and quality checking
- Data quality for datasets in your data lake and data pipelines
- Cost-effective to scale with pay-as-you-go billing, with no lock-in

In - Preview



# Considerations for quality gates

Reliable: Creating managed workflows optimized for quality and traceability

## Package & Code Quality Assurance

Scan packages for common vulnerabilities & promote high coding standards to avoid bugs ➔



**Vulnerability Scans**  
Scan containers and  
libraries for common  
vulnerabilities



**Code style**  
indentation,  
formatting,  
documentation



**Testing**  
unit test, integration  
test, end-to-end test



**Performance**  
compilation of code,  
stress test, resource  
utilization optimization



# Considerations for quality gates

Reliable: Creating managed workflows optimized for quality and traceability

## Model & System Testing

Validate data, model and infrastructure for your ML workflow →



### Data & Feature Tests

- Data validation
- Feature Importance
- Policy-compliant Data and feature Pipeline
- Unit tests for feature engineering process



### Model Tests

- ML Metrics
- A/B experiments
- Performance Validation
- Fairness/Bias checks



### ML System Tests

- Latency and throughput testing
- Compute resource utilization & right-sizing



### Integration Tests

- Data, Modeling, Orchestration, 3<sup>rd</sup> Party tools integrations
- ML consuming client integration tests

# Amazon SageMaker Clarify

DETECT MODEL STATISTICAL  
BIAS AND UNDERSTANDING  
MODEL PREDICTIONS



## Identify imbalances in data

Detect bias during data preparation



## Check your trained model for statistical bias

Evaluate the degree to which various types of bias are present in your model



## Explain overall model behavior

Understand the relative importance of each feature to your model's behavior



## Explain individual predictions

Understand the relative importance of each feature for individual inferences



## Detect drift in statistical bias and model behavior over time

Provide alerts and detect drift over time due to changing real-world conditions



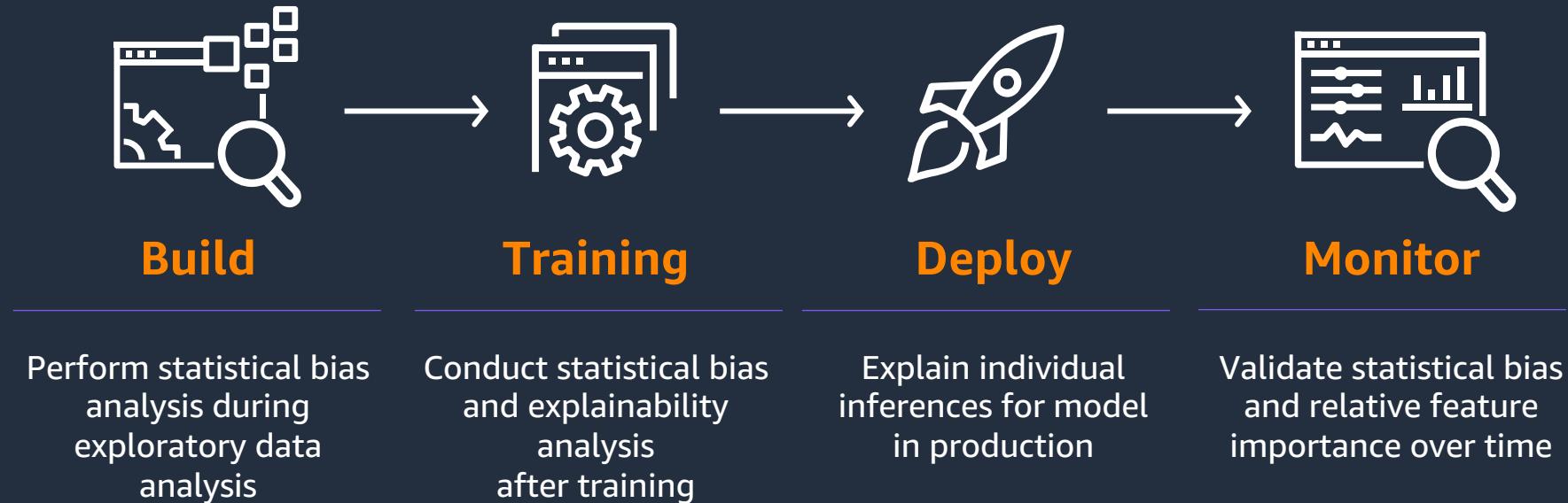
## Generated automated reports

Produce reports on bias and explanations to support internal presentations

# Amazon SageMaker Clarify

## Incorporate model explainability and statistical bias detection

### Model explainability & statistical bias detection

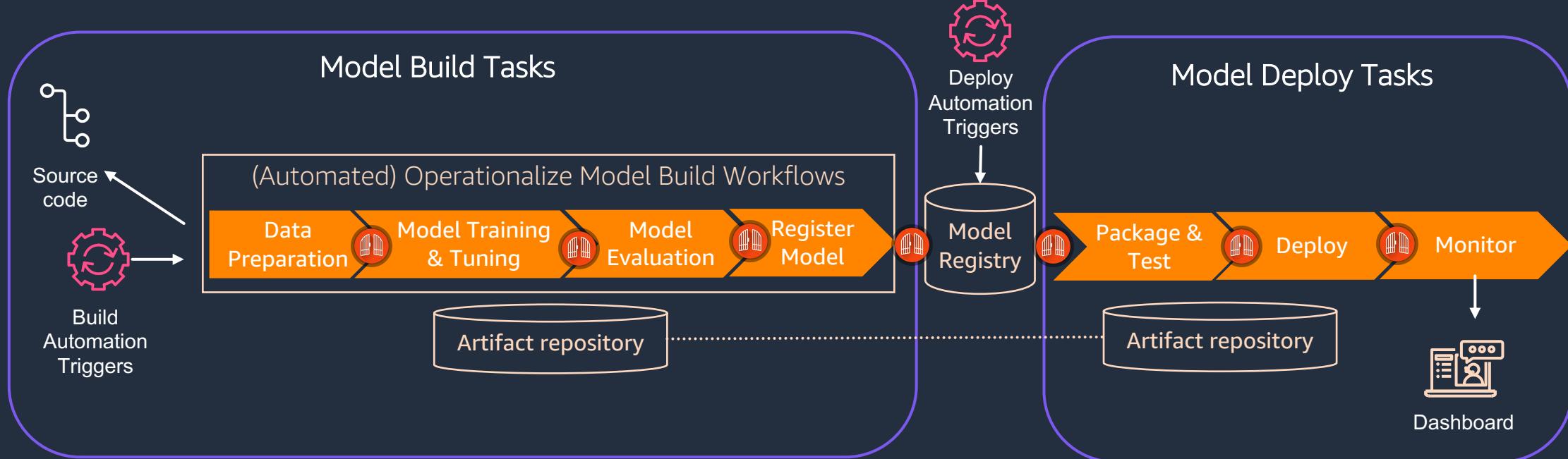


# **Reliable Stage:**

## **Deployment Quality Gates**

# Minimize deployment risk & Continuous monitoring for ML

## Challenge #3: Minimize risk with new model versions & monitoring ML workloads



### Key Technical Challenges

- Monitoring and measuring model performance
- Deploying new model versions with reduced risk

### Technical Solutions

- Incorporate advanced deployment strategies into CD pipelines
- Provide transparency into model performance across key stakeholders

### On AWS

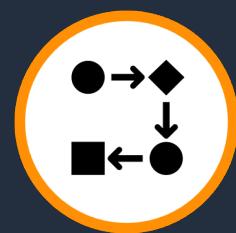
- [Amazon SageMaker Model Monitor](#)
- [Amazon SageMaker Model Dashboard](#)
- [Amazon SageMaker Model Deployment](#)

# Considerations for monitoring

Reliable: Creating managed workflows optimized for quality

## Continuous Monitoring & Logging

Consider all monitoring dimensions →



Pipeline



Model



System



Business

Alerts & Dashboards

Stakeholder transparency



Platform Team



Data Engineer



Data Scientist



Business Sponsor



Security and  
Compliance  
Officer



MLOps  
Engineer

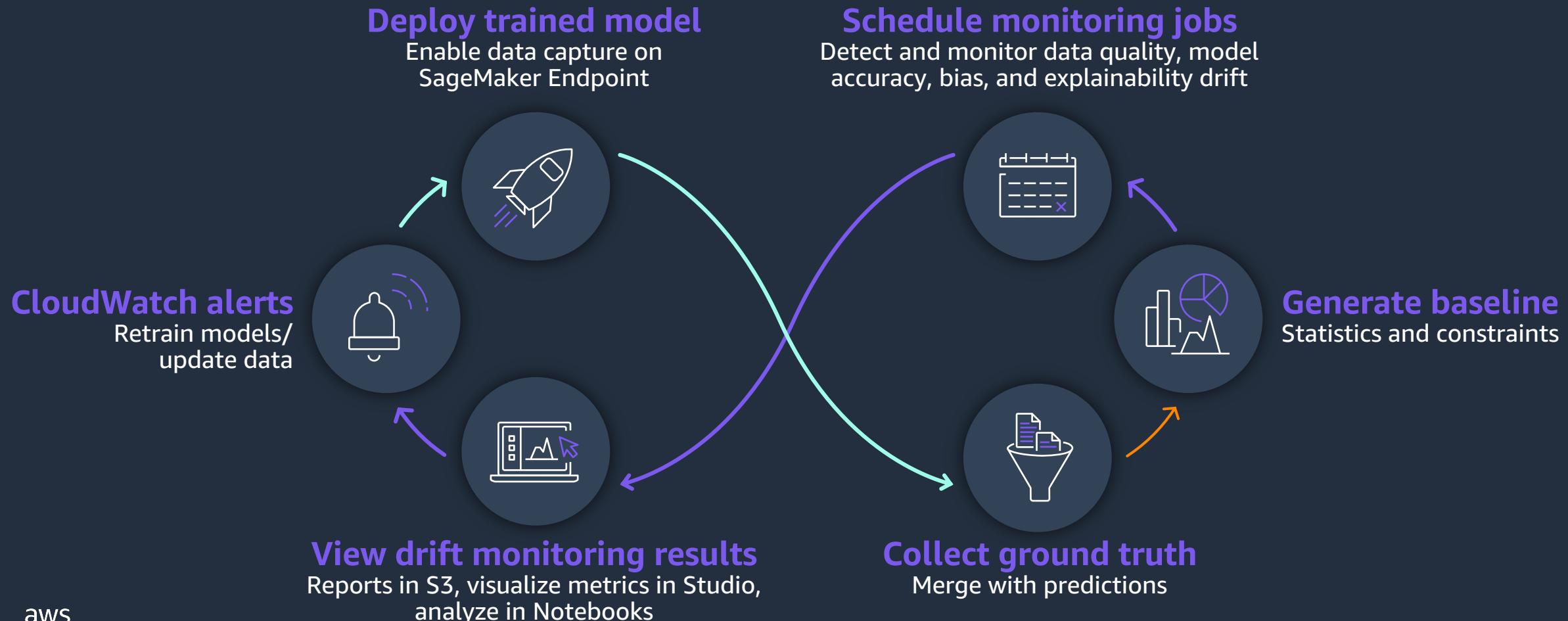
...



# Amazon SageMaker Model Monitor

## Automatically detect & mitigate drift

Easily incorporate model monitoring into your ML workflows  
Concept & Data Drift



# Amazon SageMaker Model Monitor

Automatically detect & mitigate drift

## Supported Monitor Types



Data  
Quality

***Detect divergence in data***

Real time data capture from endpoints

Define Baseline

Pre-built container for analysis

Support custom analysis

Type, Num Present, Num Missing

Mean, Sum, Std Dev

KLL Sketch



Model  
Quality

***Detect quality degradation over time***

Merge predictions with ground truth

Compare predictions to ground truth

Generate reports and violations

MAE, RMSE, F1



Model  
Bias

***Track model balance***

Overfitting

Underfitting

Class Imbalance

DPL

KL Divergence

LP-Norm



Model  
Explainability

***Feature Attribution***

How much each feature contributed to predictions

Shapley Values

# Amazon SageMaker Model Dashboard

## Provide transparency into model performance

Recent release  
(re:invent 2022)

## Unified monitoring across all your models in production



Single pane of glass view for models, endpoints, batch transform jobs, and monitoring jobs



Insights into model performance at endpoints and batch transform jobs



View alerts for violations, missing, and inactive monitors



Supports model quality, data quality, bias drift, and feature attribution drift

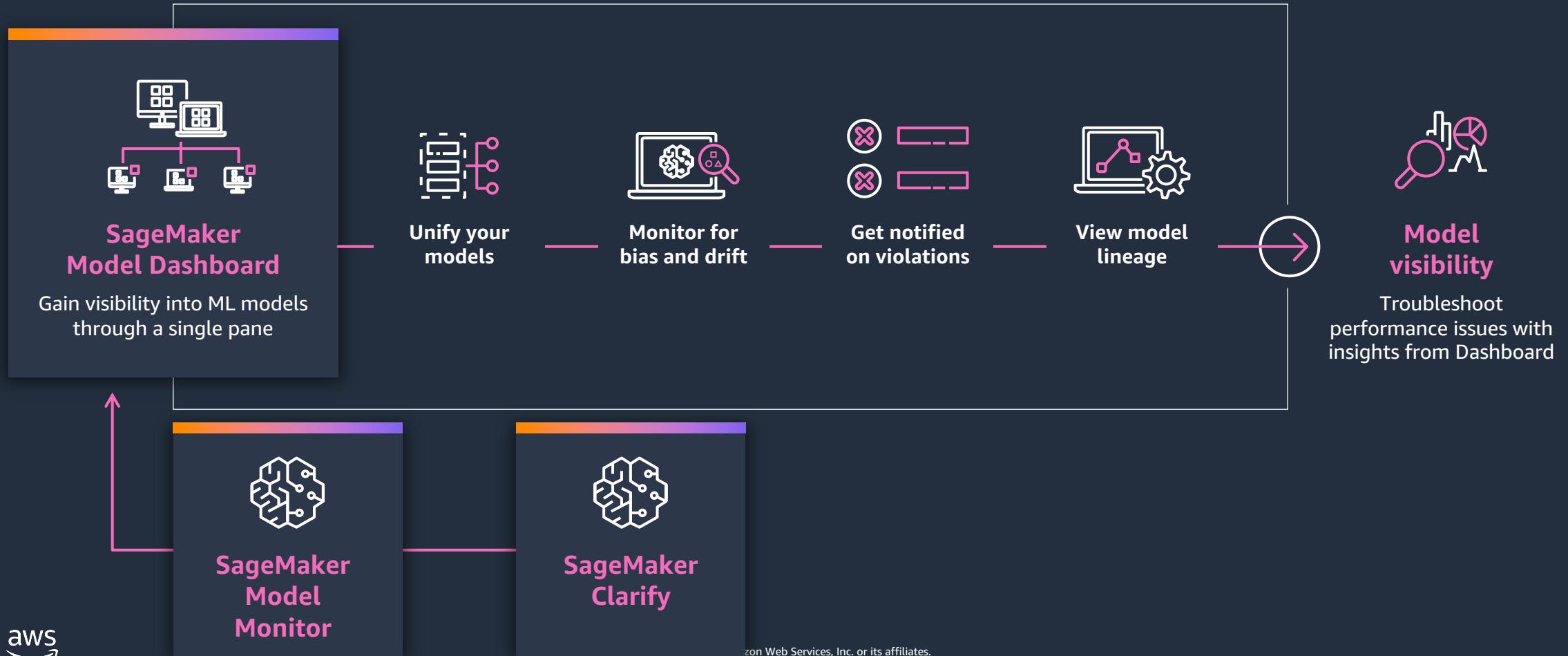
Model dashboard <a href="#">Info</a>								
Models <a href="#">Info</a>								
Filter models or endpoints by property or value								
[Filter 01] <a href="#">X</a>	and <a href="#">▼</a>	[Filter 02] <a href="#">X</a>						
Model name	Risk rating	Model quality	Data quality	Bias drift	Feature attribution drift	Endpoints	Last updated	Last checked
[ Model-01] <a href="#">...</a>	High	⚠ 1/1/2022 5:30PM	⚠ 1/1/2022 5:30PM	⊖ Inactive	⚠ 1/1/2022 5:30PM	Endpoint_Name	[B]	[B]
[ Model-02] <a href="#">...</a>	Unknown	⚠ 1/1/2022 5:30PM	⚠ 1/1/2022 5:30PM	⚠ 1/1/2022 3:00PM	⚠ 1/1/2022 10:00AM	Endpoint_Name	[B]	[B]
[ Model-03] <a href="#">...</a>	High	⚠ 1/1/2022 5:30PM	⌚ Scheduled	⌚ Scheduled	⚠ 1/1/2022 10:00AM	Endpoint_Name & 3...	[B]	[B]
[ Model-04] <a href="#">...</a>	High	ⓘ None	ⓘ None	ⓘ None	ⓘ None	Endpoint_Name	[B]	[B]
[ Model-05] <a href="#">...</a>	High	ⓘ None	ⓘ None	ⓘ None	ⓘ None	Endpoint_Name & 3...	[B]	[B]
[ Model-06] <a href="#">...</a>	Medium	-	-	-	-	-	-	-
[ Model-07] <a href="#">...</a>	Medium	⊖ Inactive	⊖ Inactive	⊖ Inactive	⊖ Inactive	Endpoint_Name	[B]	[B]
[ Model-08] <a href="#">...</a>	Medium	⊖ Inactive	⊖ Inactive	⊖ Inactive	⊖ Inactive	Endpoint_Name & 3...	[B]	[B]
[ Model-09] <a href="#">...</a>	Medium	⊖ Inactive	⊖ Inactive	⊖ Inactive	⊖ Inactive	Endpoint_Name	[B]	[B]
[ Model-10] <a href="#">...</a>	Low	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	Endpoint_Name	[B]	[B]
[ Model-11] <a href="#">...</a>	Low	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	Endpoint_Name & 3...	[B]	[B]
[ Model-12] <a href="#">...</a>	Low	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	Endpoint_Name	[B]	[B]
[ Model-13] <a href="#">...</a>	Low	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	Endpoint_Name	[B]	[B]
[ Model-14] <a href="#">...</a>	Low	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	Endpoint_Name & 3...	[B]	[B]
[ Model-15] <a href="#">...</a>	Low	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	⌚ Scheduled	Endpoint_Name	[B]	[B]

# Amazon SageMaker Model Dashboard

## Provide transparency into model performance

Recent release  
(re:invent 2022)

## How it works



# **Reliable Stage: MLOps Deployment Strategies**

# Considerations for model deployment strategies

Reliable: Creating managed workflows optimized for quality and traceability

## Model Deployment Strategies

Validate new model versions with advanced deployment strategies to minimize deployment risk ➔



### Shadow Test

Validate model versions in production without returning its outputs to users



### A/B Tests

Validate model versions in production serving both outputs to users



### Blue/Green Deployments

Securely update models in production with quick rollback functionality

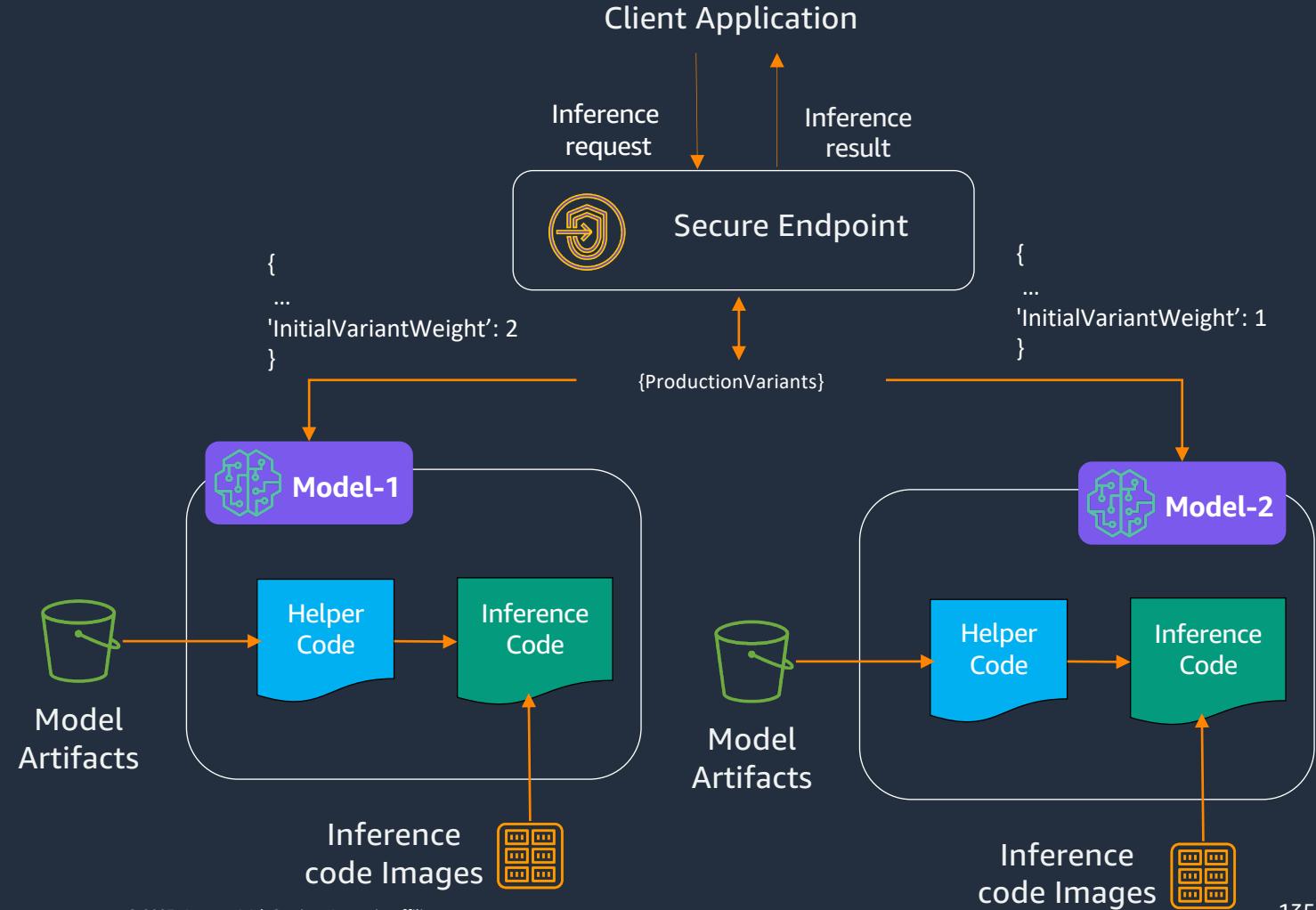
# Amazon SageMaker Model Deployment

Automate advanced model deployment strategies

- 1-10 Production Variants (Model Versions)
- All models must have the same I/O schema
- Endpoint Modification w/o service disruption
- Perform traffic distribution per variant based on weight or invoke a specific variant per request



## A/B Testing Compare model versions



# Amazon SageMaker Model Deployment

## Automate advanced model deployment strategies

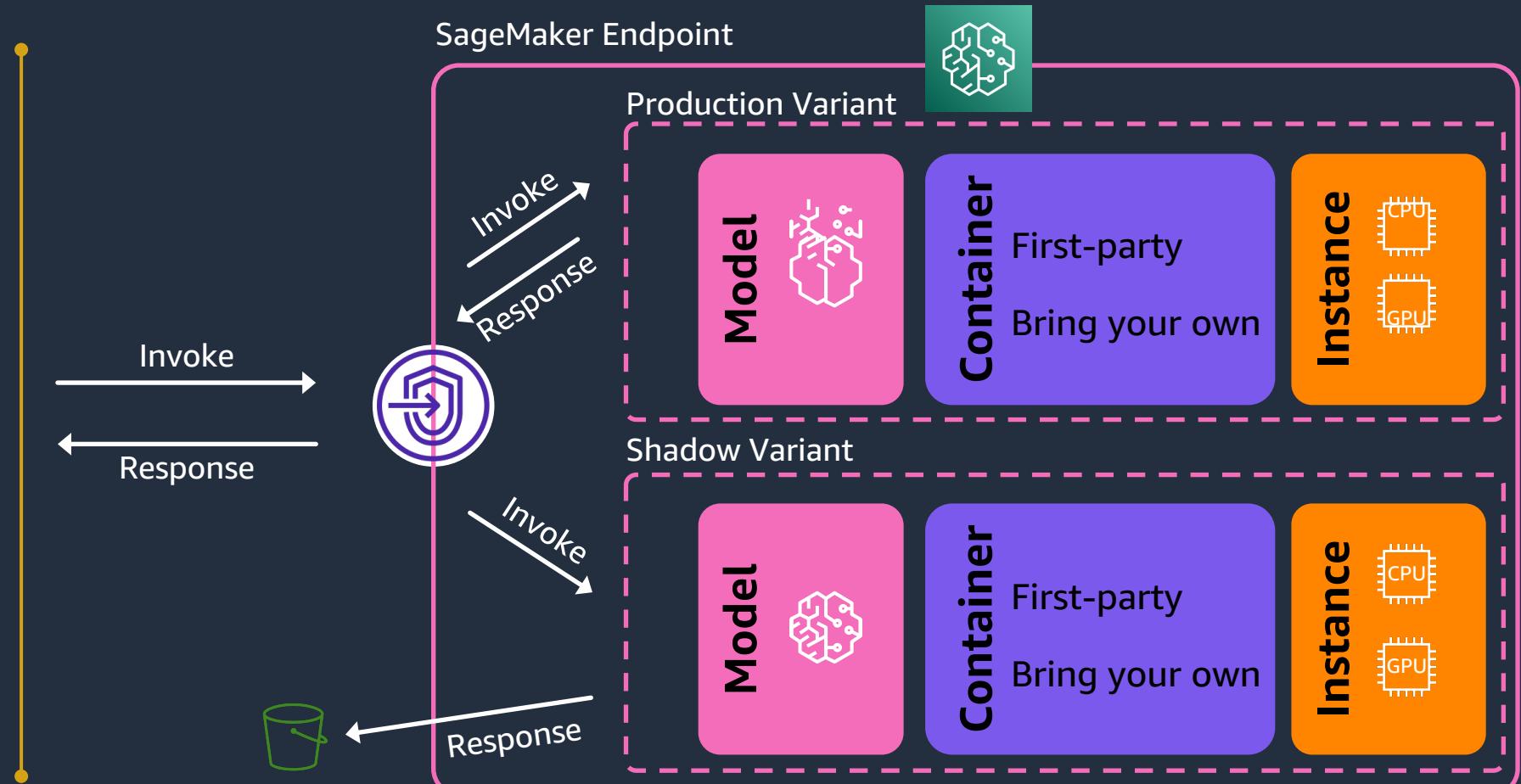
### Shadow Testing

Validate new models in parallel to the model serving production traffic

SageMaker takes care of mirroring requests

Start small and dial up to control costs

Accessible through AWS console, CLI, APIs



# Amazon SageMaker Model Deployment

## Automate advanced model deployment strategies

### Shadow Testing

#### Tools to manage shadow tests



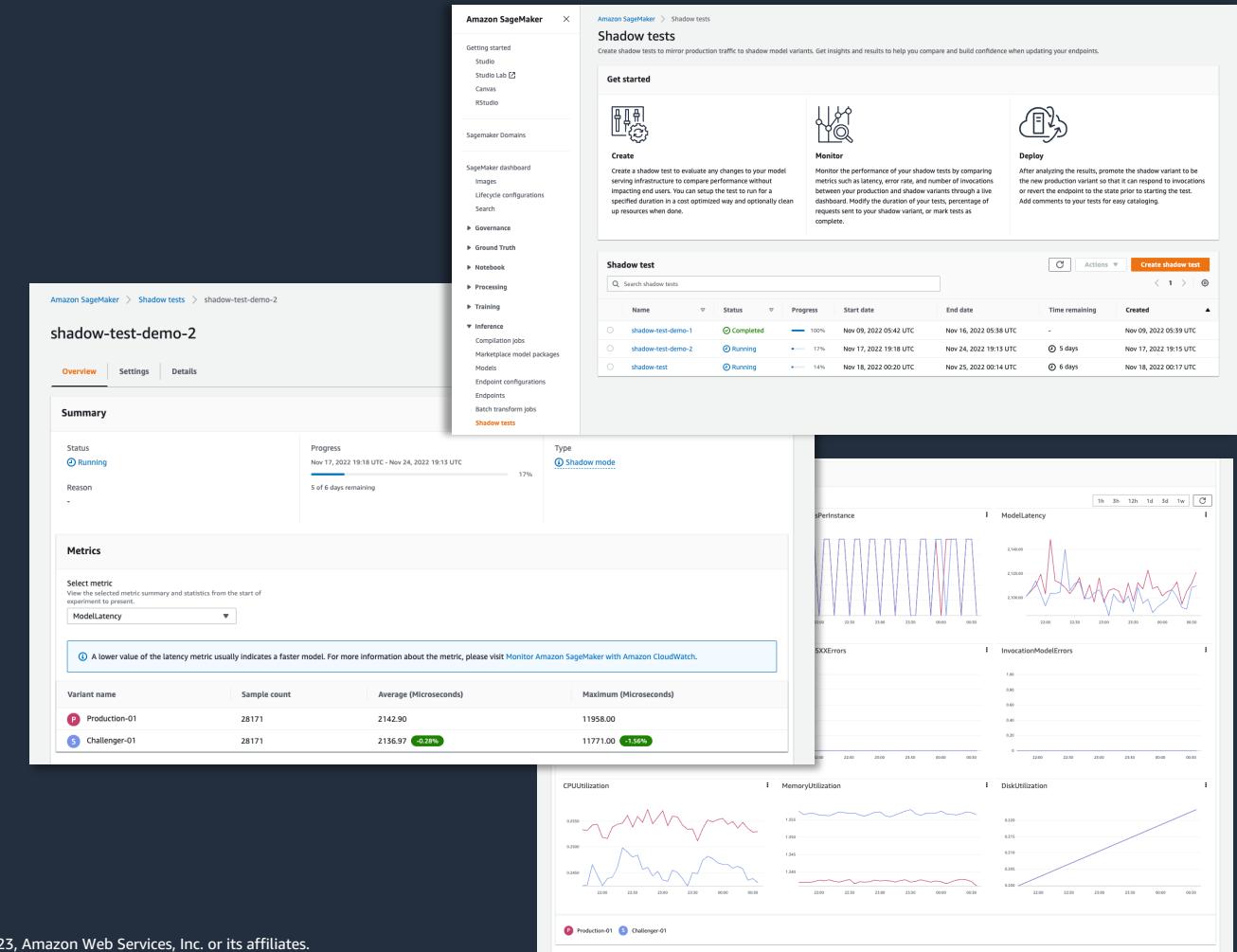
Set up a test for a predefined duration



Monitor operational performance through a live dashboard



Deploy models with confidence or rollback



# Amazon SageMaker Model Deployment

Automate advanced model deployment strategies

## Deployment Guardrails

Safe deployments into production environments

VALIDATE YOUR MODEL BEFORE PRODUCTION UPDATE

Blue/Green Deployment  
Traffic shifting modes

All At Once  
1-step traffic shift

Canary  
2-step traffic shift

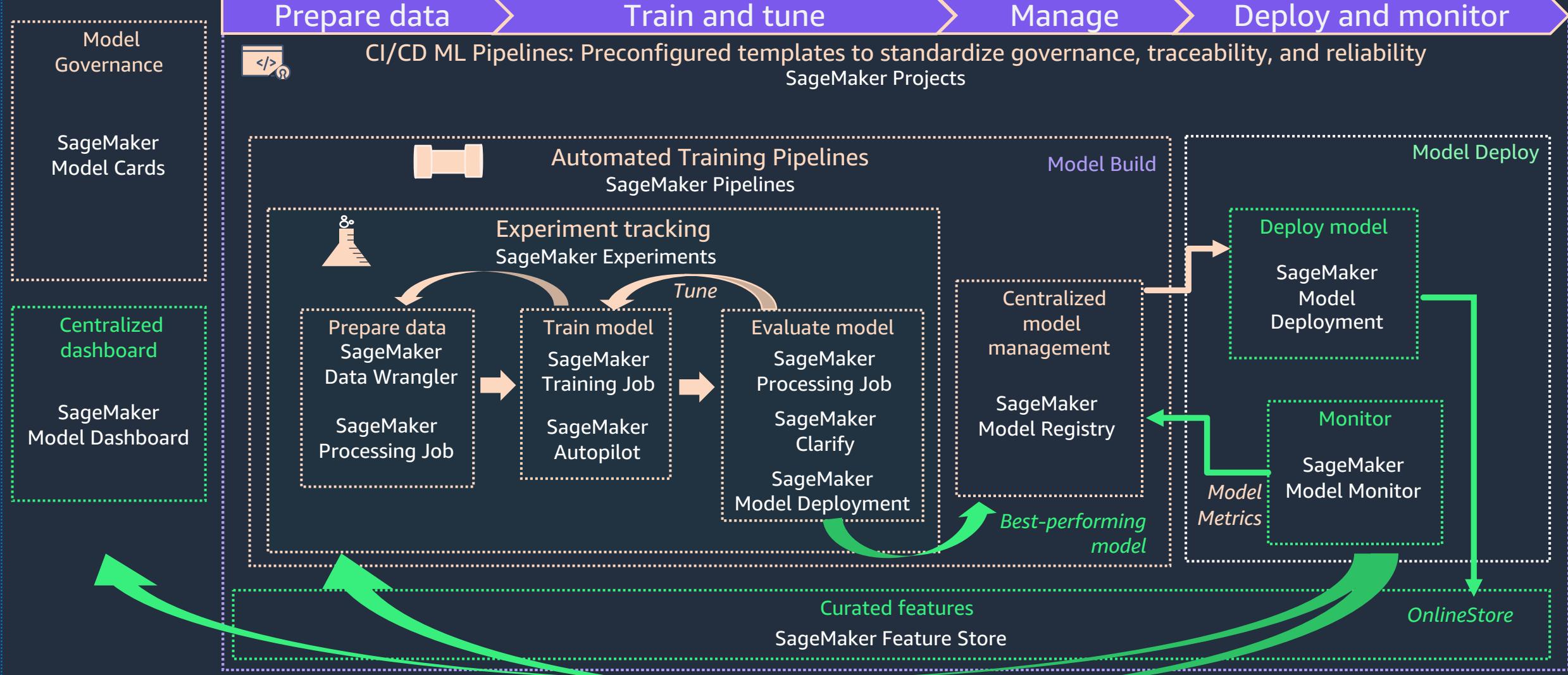
Linear  
n-step traffic shift

- Traffic shifting modes, such as canary and linear
- Deployment safety while updating production environments
- Built-in safeguards such as auto-rollsbacks
- Fully managed deployment
- Visibility: Track the progress of your deployment

# MLOps Maturity Framework in Practice

Reliable Phase Goal: Creating managed workflows optimized for quality & traceability

Govern



## Reliable Stage



## Key Technical Takeaways

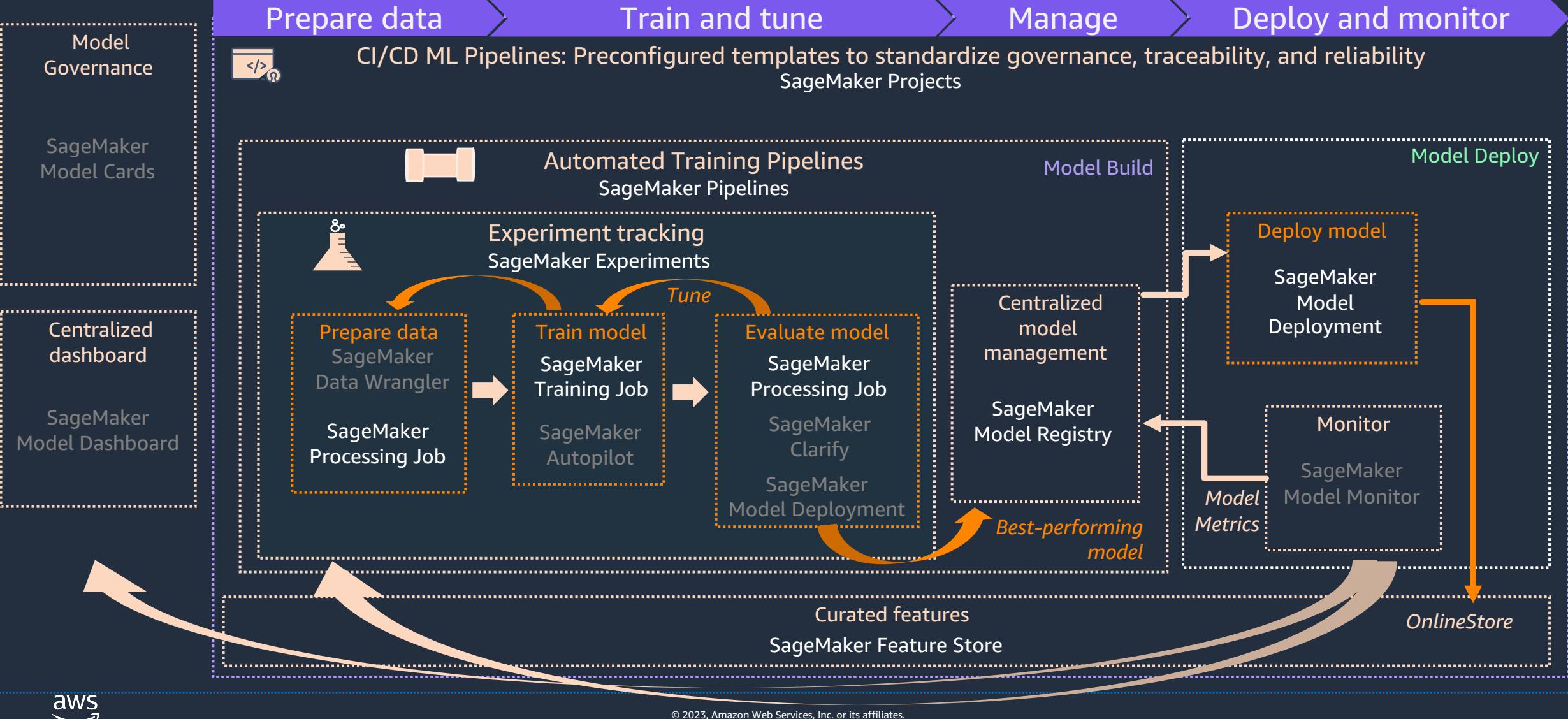
- 1 Implement CI/CD practices to incorporate versioning of all key inputs and artifacts for model lineage and reproducibility and to streamline deployment
- 2 Identify & implement automated (or manual when necessary) quality gates that span the end-to-end workflow
- 3 Implement model monitoring to detect signals of data or concept drift
- 4 Consider advanced deployment strategies to minimize risk when deploying a new model version

# Lab3: Reliable

**Reliably build, deliver, and  
manage models at scale**

# Lab 3: Amazon SageMaker - Features Used

Govern

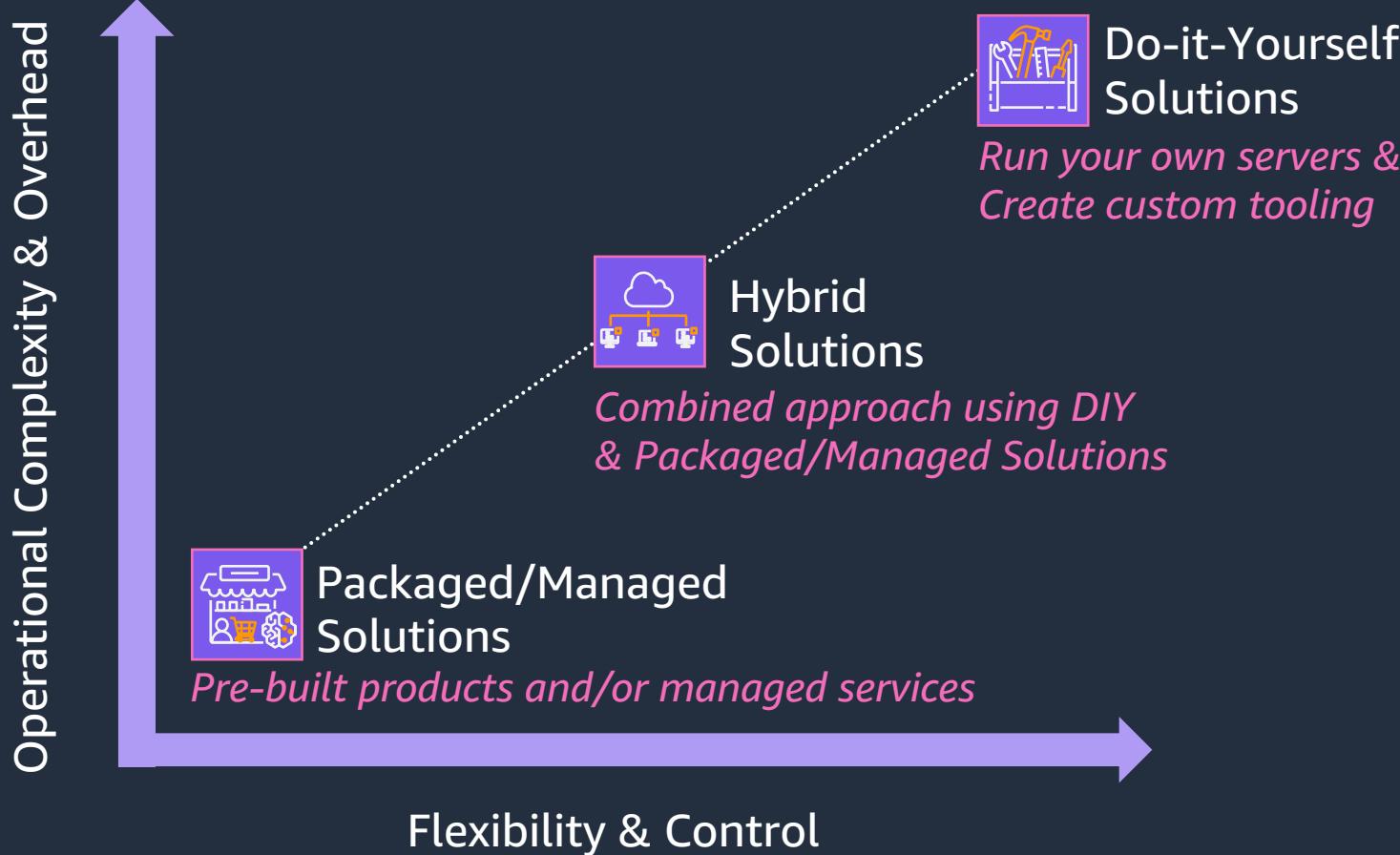


# AWS & Third-party integrations



# MLOps Technical Implementation Choices

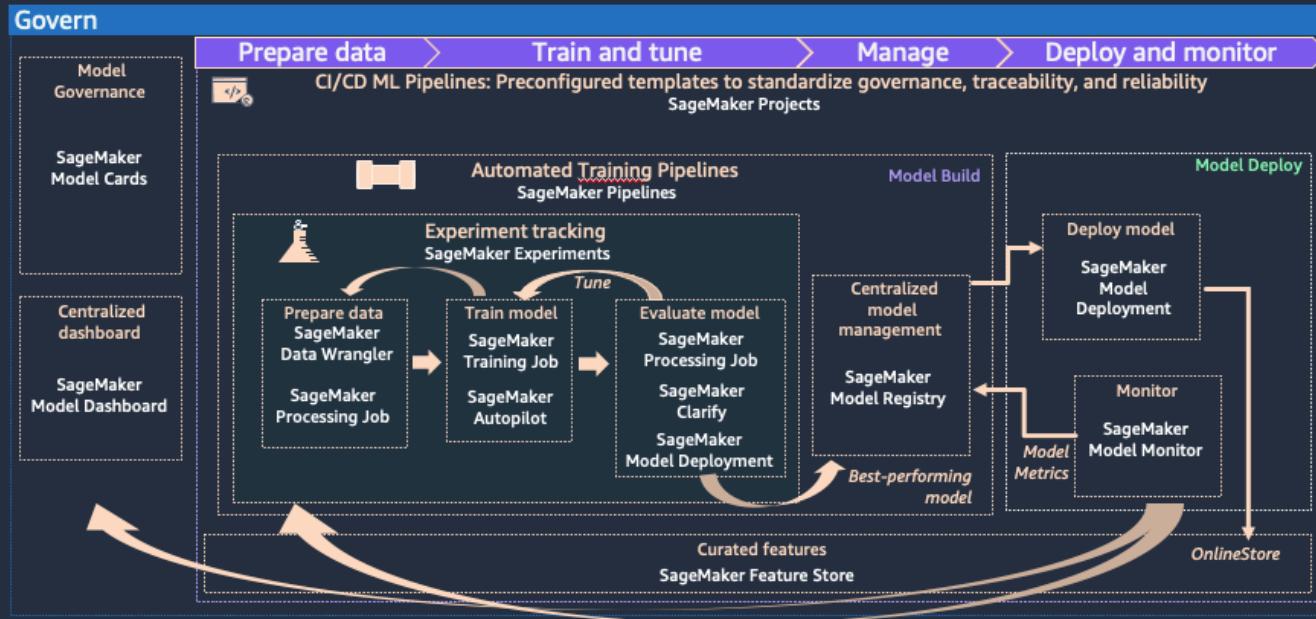
## Understand Your Tradeoffs



# Amazon SageMaker

## Flexible across use cases & requirements

SageMaker offers the flexibility to utilize native features end-to-end....



...Or integrate with existing tooling as needed based on existing investments, skillsets, or organizational requirements

# Tooling Considerations

## Integrating SageMaker Jobs into MLOps Tooling



*\*\*Note: This list is not inclusive of every do-it-yourself combination*

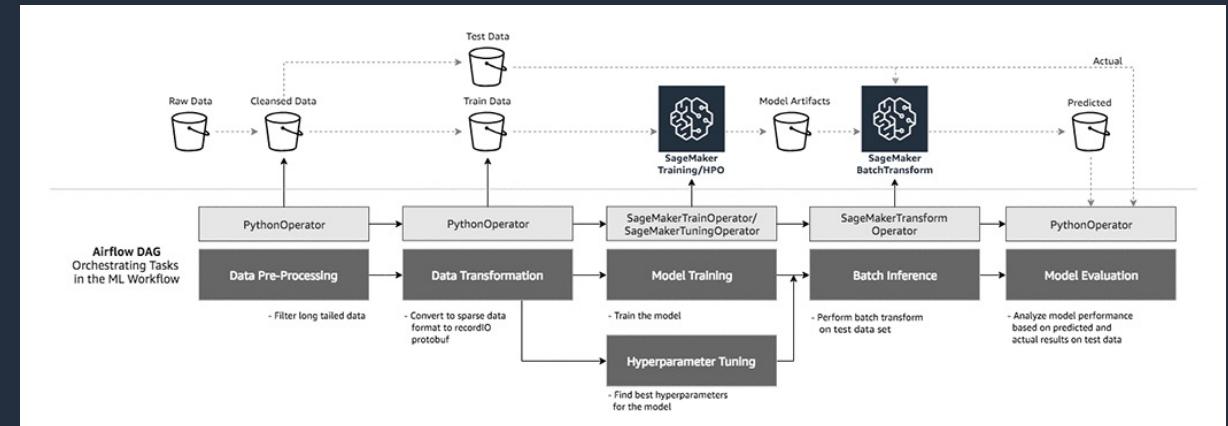


# Tooling Considerations: Airflow

## Integrating SageMaker Jobs into MLOps Tooling



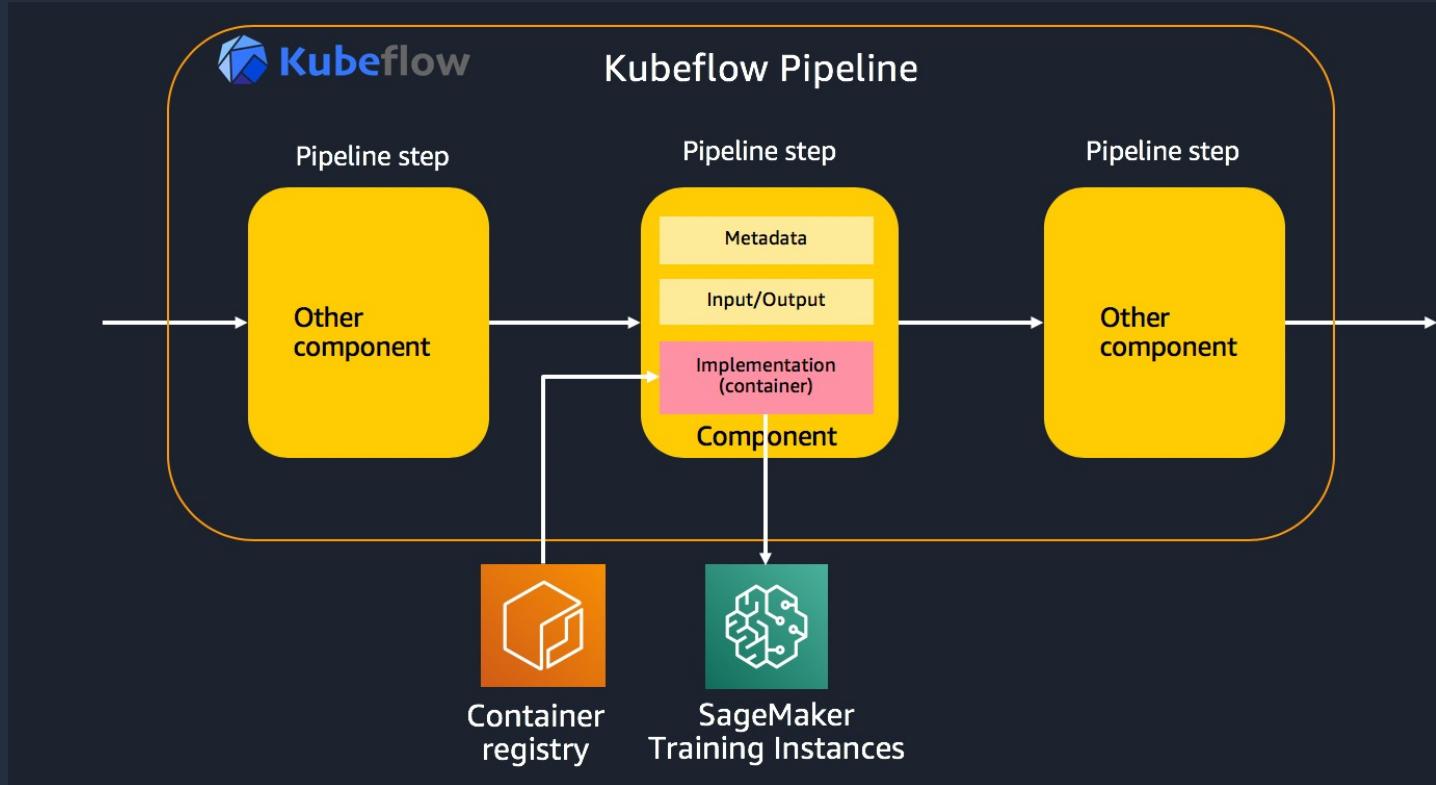
- **SageMaker Operators**
  - SageMakerTrainingOperator
  - SageMakerTuningOperator
  - SageMakerModelOperator
  - SageMakerTransformOperator
  - SageMakerEndpointConfigOperator
  - SageMakerEndpointOperator
  - SageMakerAutoMLOperator
  - SageMakerCreateExperimentOperator
  - SageMakerStartPipelineOperator
  - SageMakerRegisterModelVersionOperator



# Tooling Considerations: Kubeflow

## Integrating SageMaker Jobs into MLOps Tooling

### SageMaker Components for Kubeflow Pipelines



#### Key Features:

- Run production model pipelines from open source in Kubeflow Pipelines
- Take advantage of faster and cheaper training with SageMaker using managed spot training, distributed training, and security compliance
- Run inference using SageMaker Batch or Hosting



## Key Technical Takeaways

1

Consider the number of integrations you will need to manage

2

Standardize on core components such as Model Registry or Feature Store

3

Aim to utilize purpose built tooling

4

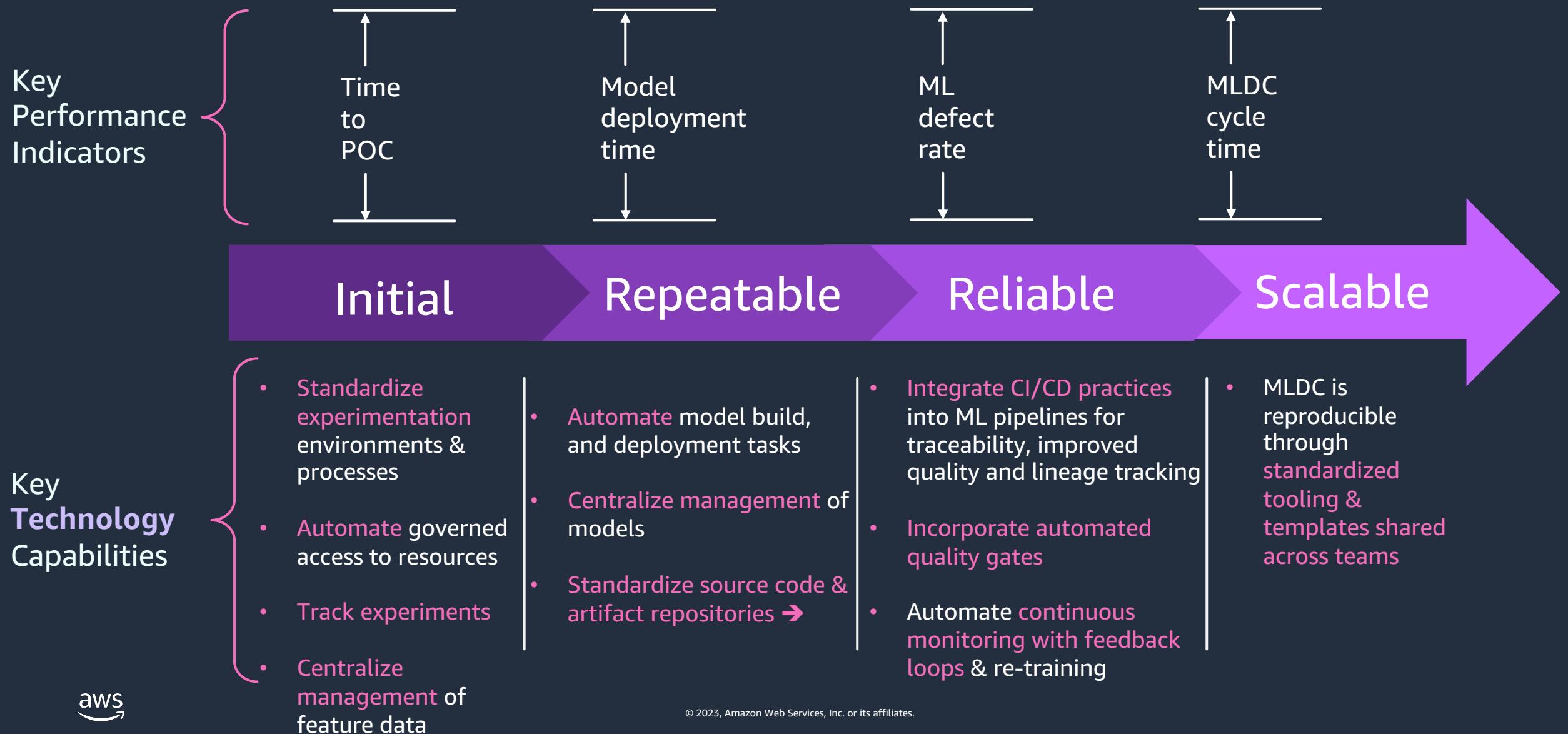
Utilize managed tooling to reduce operational overhead and server management

# Scalable Stage:

**Implementing standardized  
ML platforms & workflows  
across teams**

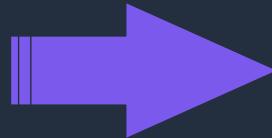
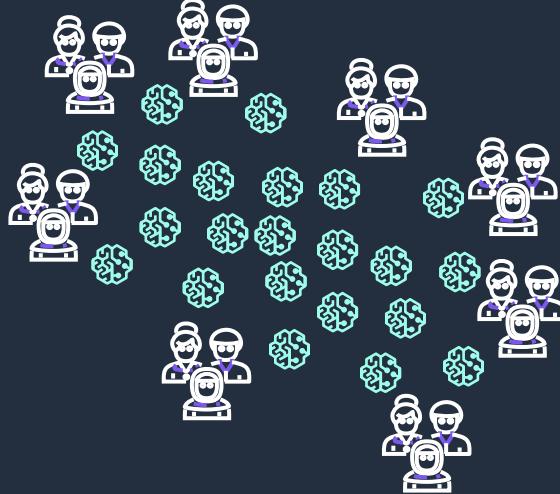
# MLOps Maturity Framework

## Scaling ML through your MLOps journey



# Moving from Repeatable to Scalable Maturity Level

Scalable: Shared best practices & provide a reproducible MLDC



What if you're supporting  
100's/1000s of use cases?  
teams? tooling management  
& integrations?



Select platform and tooling choices that  
**minimize operational overhead and total cost**  
while **maintaining IT controls**



Consider **standardized tooling**



Consider **readily available templates** to  
streamline workflow

# Governed ML Platform on AWS

**Build & Govern ML workloads  
at scale**

# Agenda

- Why build Governed ML Platform?
- Governed ML Platform Capabilities
- Governed ML Platform Reference Architecture
- Resources
- Q & A



# ML and Path to Production

MACHINE LEARNING CODE AND DATA SCIENCE NOTEBOOKS...

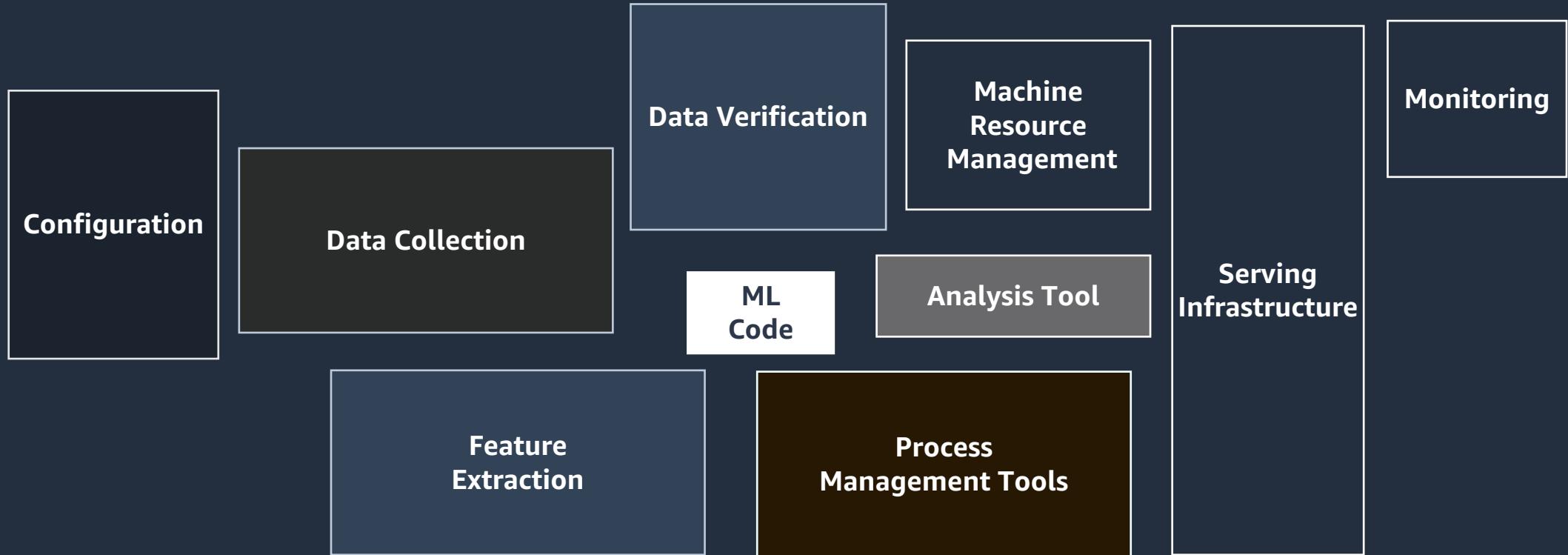
ML  
Code



© 2023, Amazon Web Services, Inc. or its affiliates.

# ML and Path to Production

...IS ONE SMALL PART OF THE OVERALL DEPLOYMENT PICTURE



***"Only a small fraction of real-world ML systems is composed of the ML code"***

source: Hidden Technical Debt in Machine Learning Systems [D. Sculley, & al.] – 2015

<https://papers.nips.cc/paper/5656-hidden-technical-debt-in-machine-learning-systems.pdf>

# Challenges



**Abstraction** – Abstract infra setup and access controls



**Scalability** – Being able to scale resources for teams with isolation



**Reliability** – Perform intended function correctly and consistently at scale



**Auditability** – Logs, versions and dependencies of artifacts



**Reproducibility** – Need to recreate past experiments



**Governability** – Ability to centrally govern resources and manage risk

# Governed ML Platform Capabilities

## Build ML Foundations

Operating Model

Multi-Account Foundations

Data Foundations

Central Feature Store

## Scale ML Operations

Self-Service Onboarding

Guided ML Workflow

CI/CD Build & Deployment

Central Model Registry

## Observable ML

Experiment Tracking & Lineage

Centralizing Logging & Monitoring

Model Governance

Cost Control & Reporting

## Secure ML

Security Foundations

Authentication & Authorization

Data & Model Security

Differential Privacy



# 1. Build Governed ML Platform Foundations

# Operating Model

Team Structure	Definition
 Centralized	Data science activities are centralized within a single team or organization
 Decentralized	Data science activities are distributed across different business functions or divisions with no central governance
 Federated	Shared services and governance functions managed by centralized team and each business unit team manages their data science functions

# Accounts are the building blocks

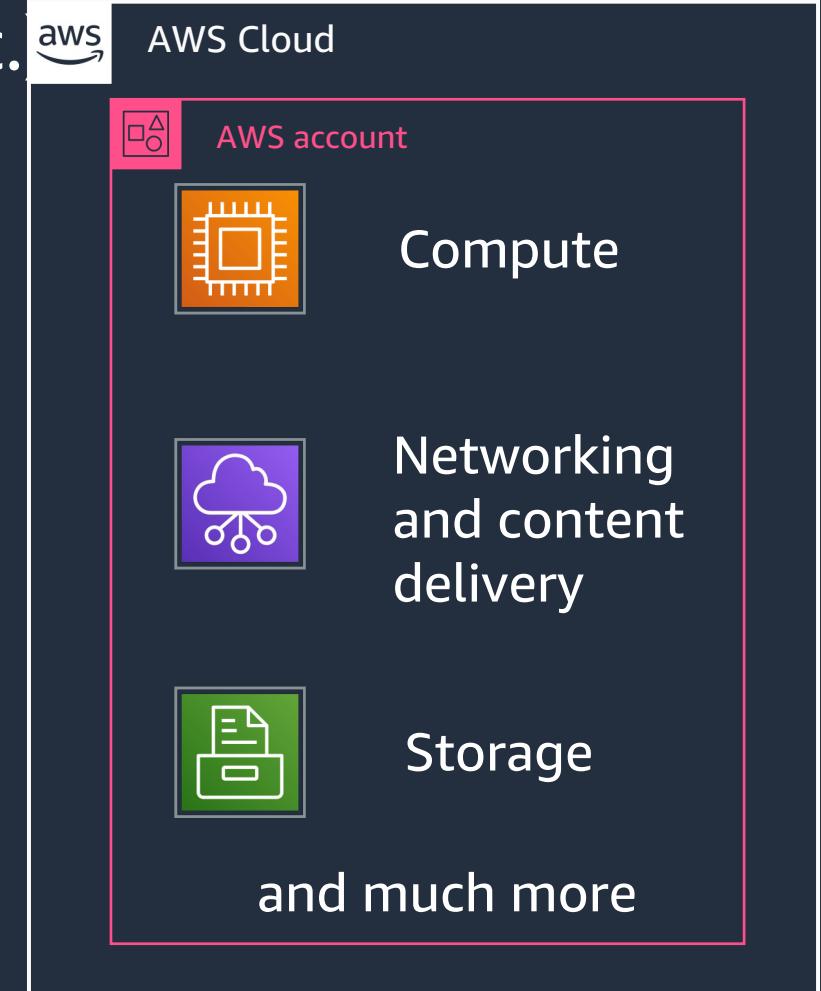
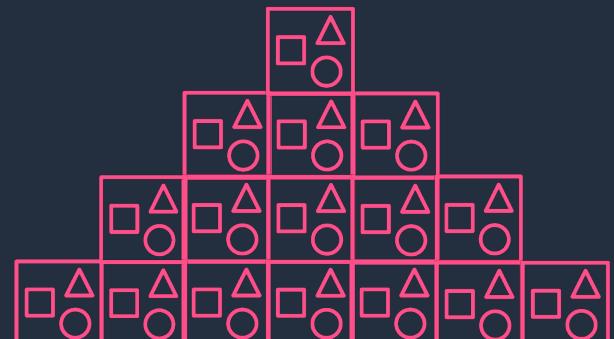
Think of AWS resources (instances, storage, etc.) as material and accounts as bricks

Scale across accounts

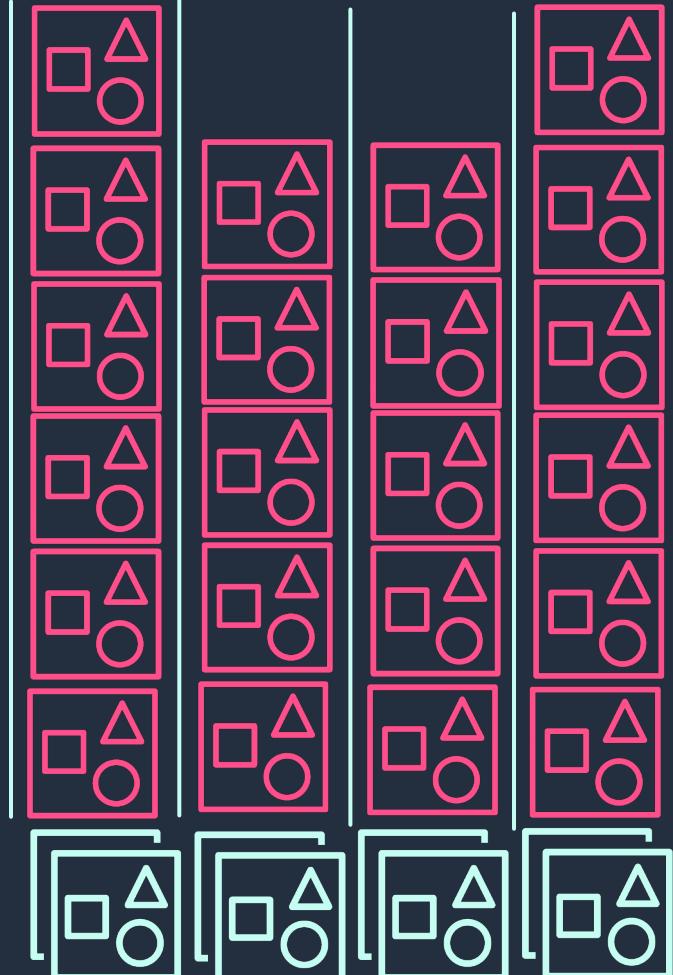
**Account limits**  
Quotas

**Security**  
Natural boundaries, isolation

**Compliance/  
business processes**  
Billing, custom requirements



# Framework for accounts: Organization Unit(OU)



To build out structure, you need supports, trusses, or frames in place to define the pattern for bricks

Organizational units (OUs) are the frames

Design your OU structure around security and operational needs

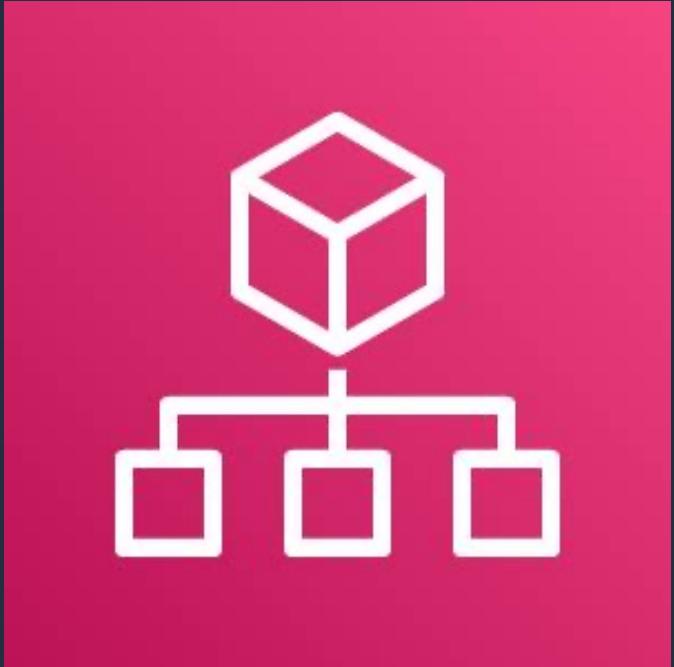


Start small and broad, expand as needed



Guardrails/ policies are assigned to OUs

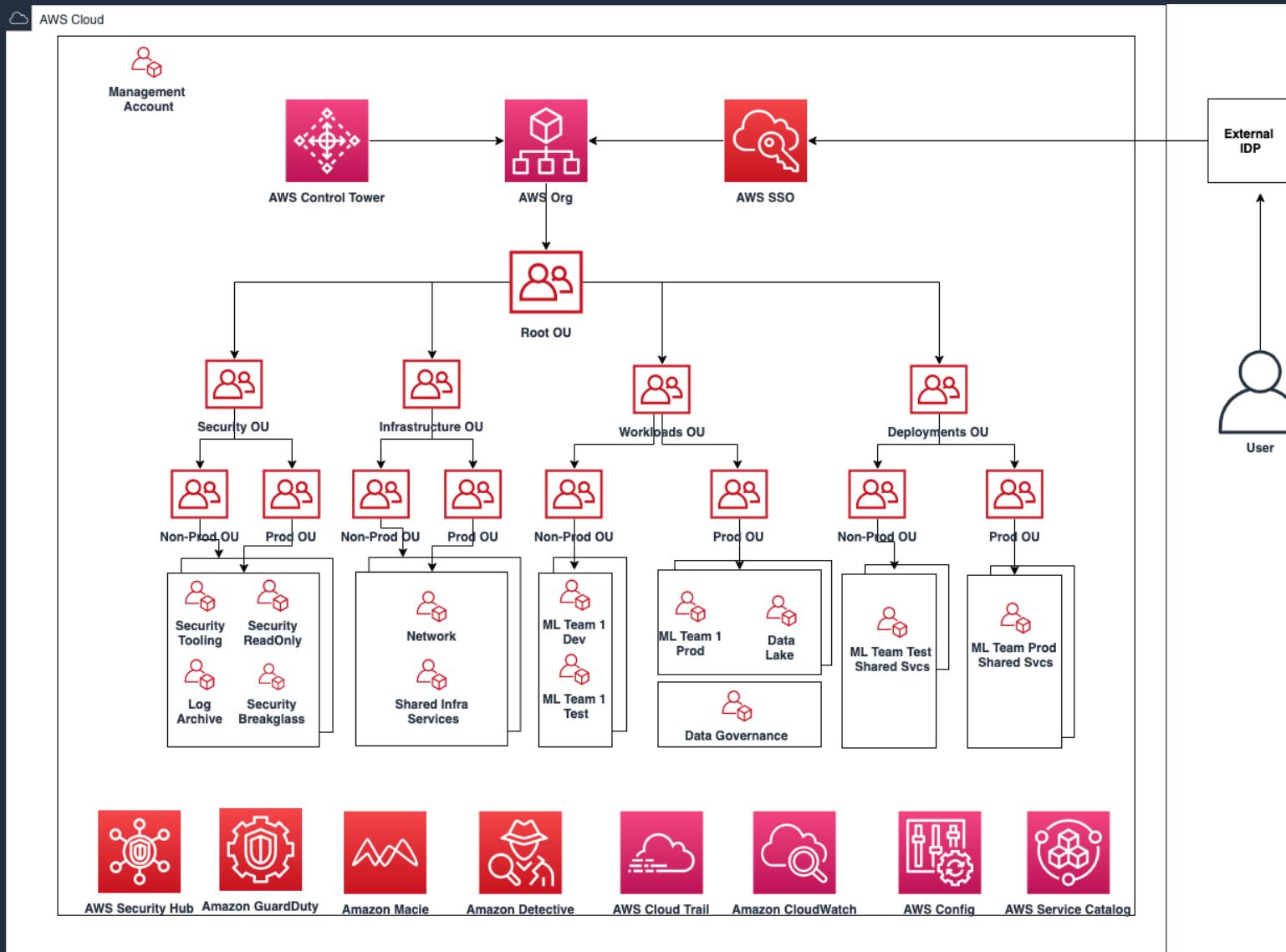
# AWS Organizations



Provides you tools to centrally govern and manage your cloud environment

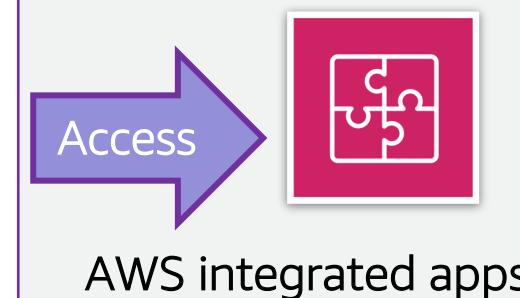
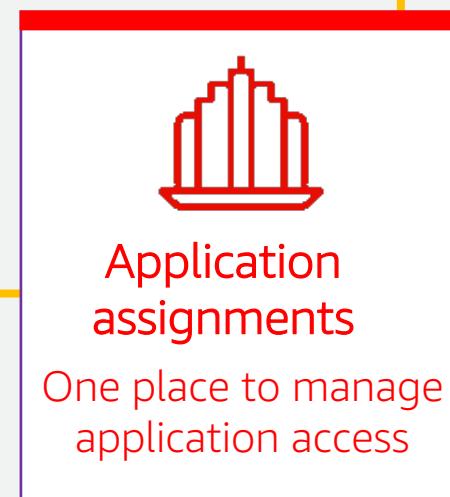
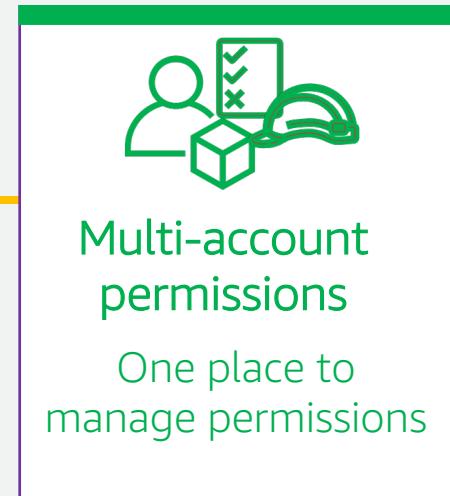
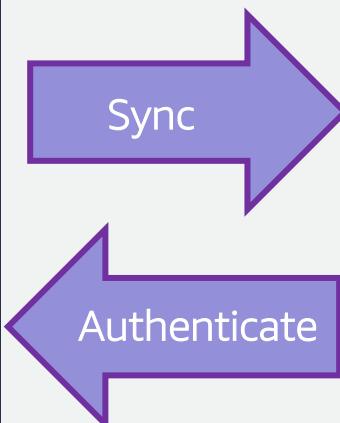
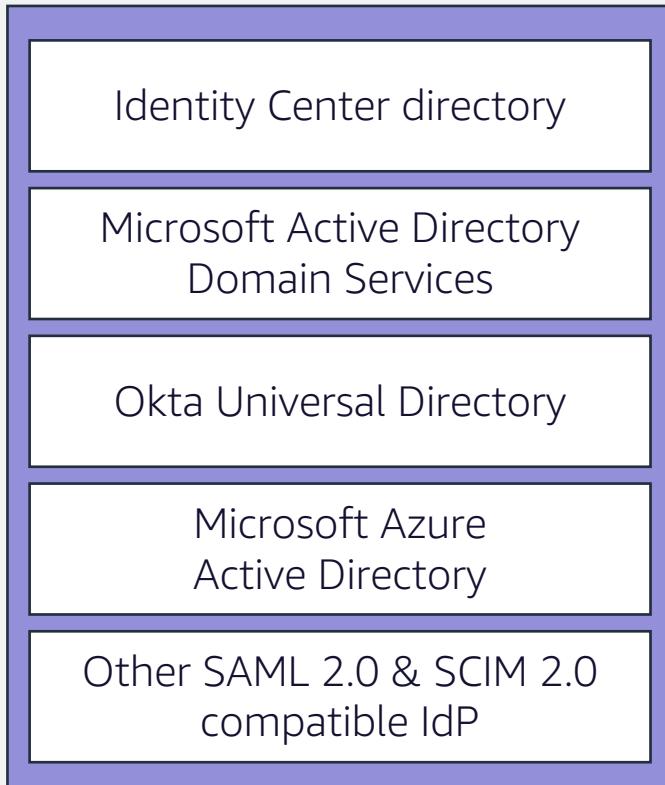
- Quickly scale by creating accounts and allocate resources
- Customize your environment by applying governance policies
- Secure and audit your environment
- Share resources and control access
- Manage costs and identify cost-saving measures

# AWS Organization Structure



# IAM Identity Center (formerly AWS SSO)

## Choose your identity source



SAML 2.0 apps

# Overview of AWS Control Tower



## AWS Control Tower

Easily set up and manage a secure multi-account environment

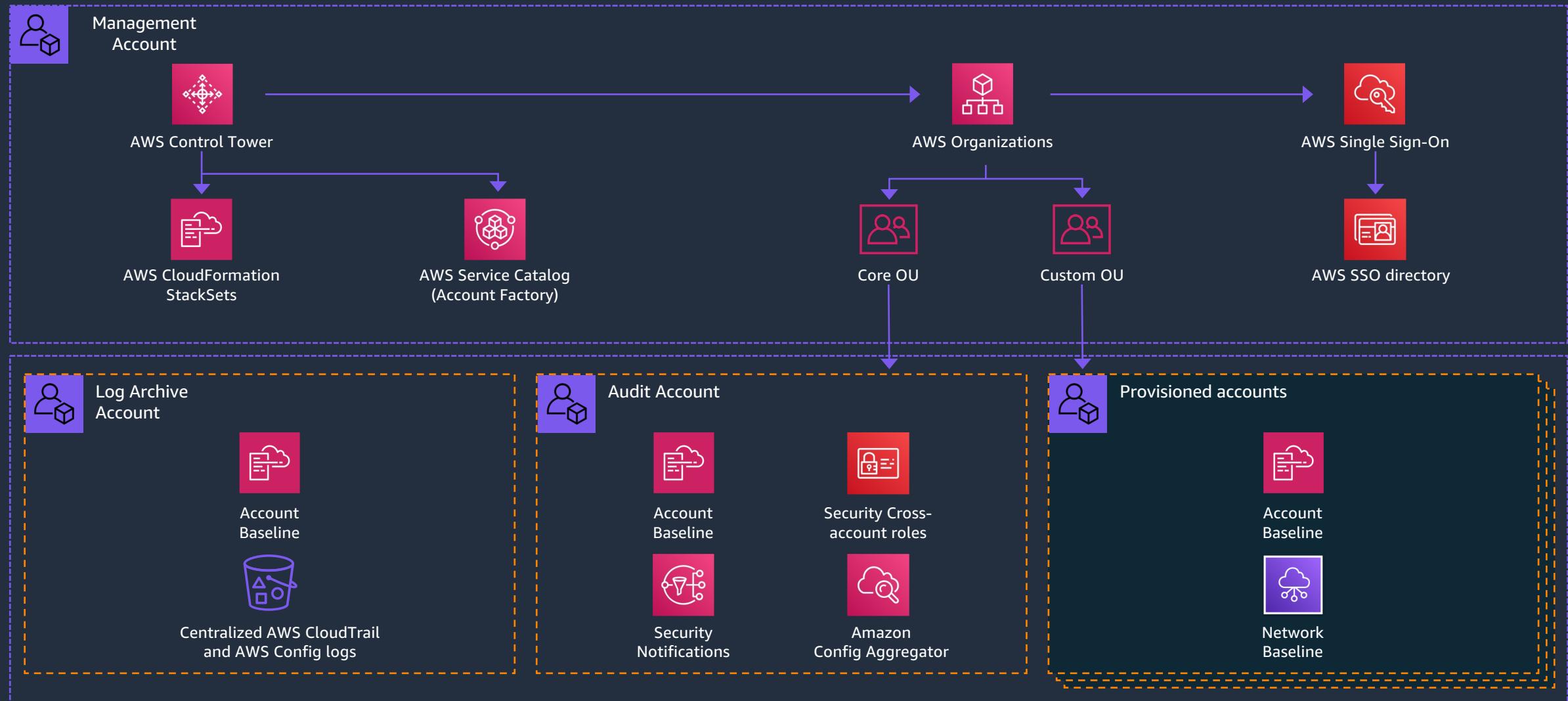


### Service overview

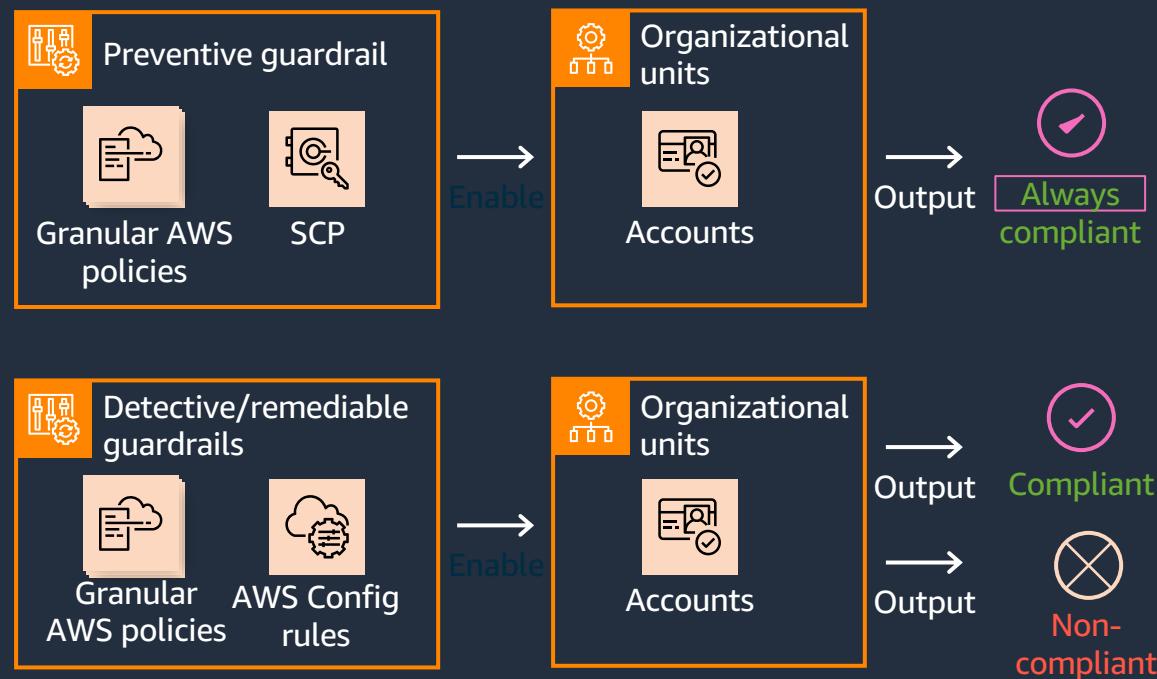
- Build an AWS management foundation based on best practices
  - Deploy Landing Zone using AWS Organizations, AWS CloudTrail, AWS IAM, etc.
- Install guardrails
  - Pre-packaged “guardrails” of security, operations, and compliance requests across the enterprise or only to specific accounts
- Free of charge  
(but incurs the cost of each AWS service required to configure the Landing Zone)



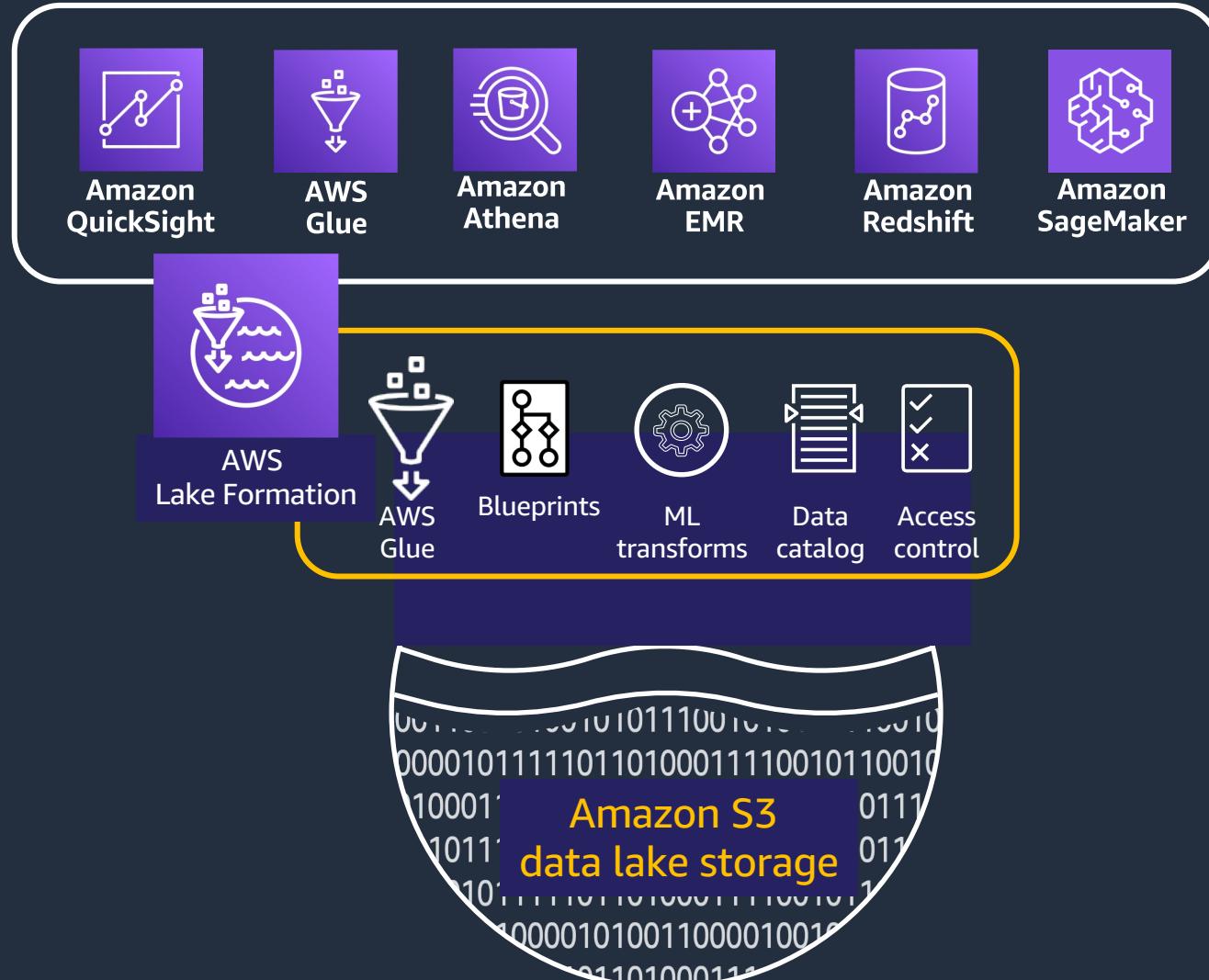
# Landing Zone provisioned by AWS Control Tower



# Establish guardrails



# AWS Lake Formation: The foundation for data lakes



Single place to manage access

Open file formats

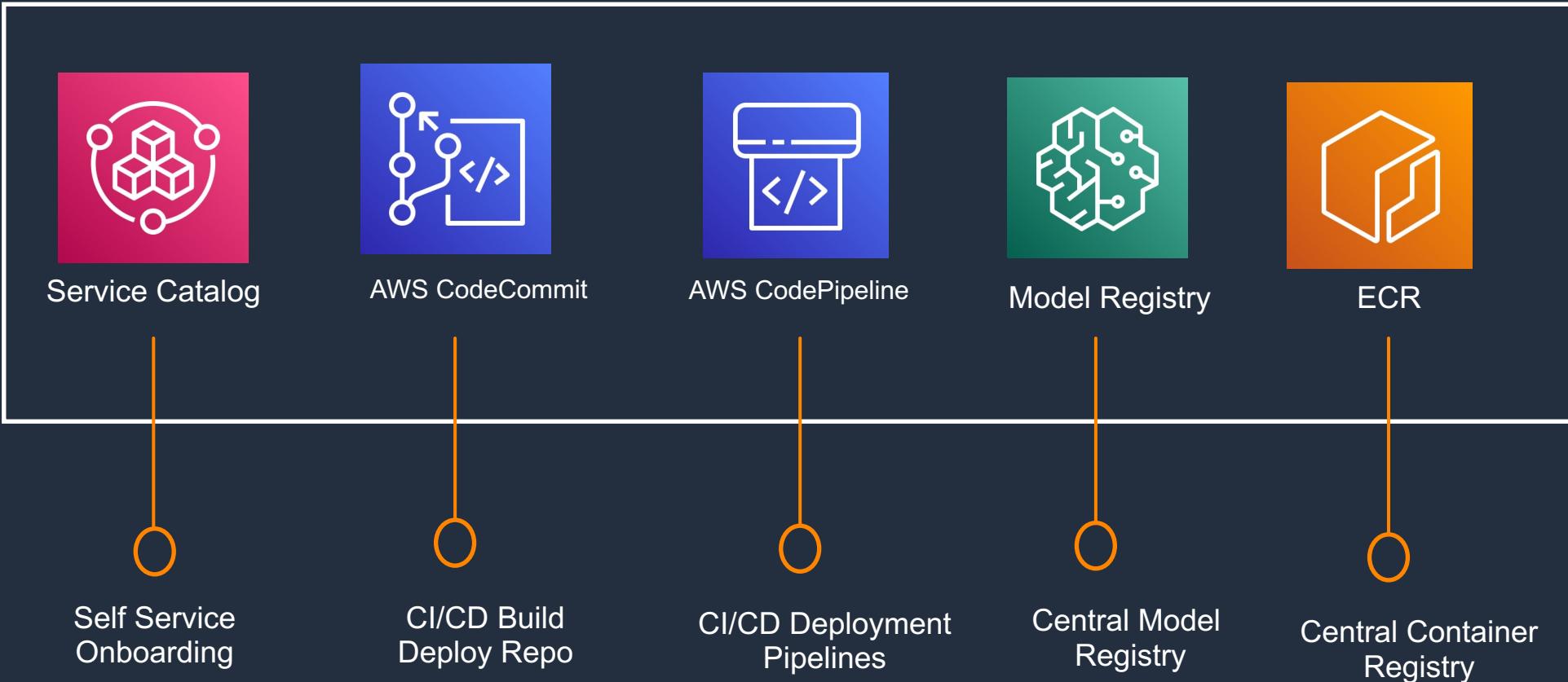
Efficient sharing

Ecosystem of integrated tools

Cost effective

## 2. Scale Governed ML Platform Operations

# ML Platform Shared Services



# Onboarding: ML Platform Key Personas and Roles

Advance Analytics Team  
Data Lake



## Data Engineer

Prepare & Ingest data building ETL pipelines for the ML use cases

Data Science Team  
Experimentation & MLOps



## Data Scientist

Create and Maintain the best ML models to solve business problems



## Data Owners

Manage data sharing and provide access to other teams



## ML Engineer

Collaborate with DS to productionize their code and models following development best practices

Platform Team  
Secure Cloud/Data/ML Platform



## Platform/MLOps Engineer

Standardize CI/CD, Container creation, model consumption, testing and deployment methodology based on business requirements



## Security

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



## Architects/ SysOps Engineer

Standardize account infrastructure e.g. VPC and endpoints, subnets, security groups, connectivity, user roles implementation

Business  
Viz Dashboards, ML Adoption, & ROI



## Business Stakeholder Product Owners

Define business problem, ML use case business KPIs, interact with other BUs, and make business decisions

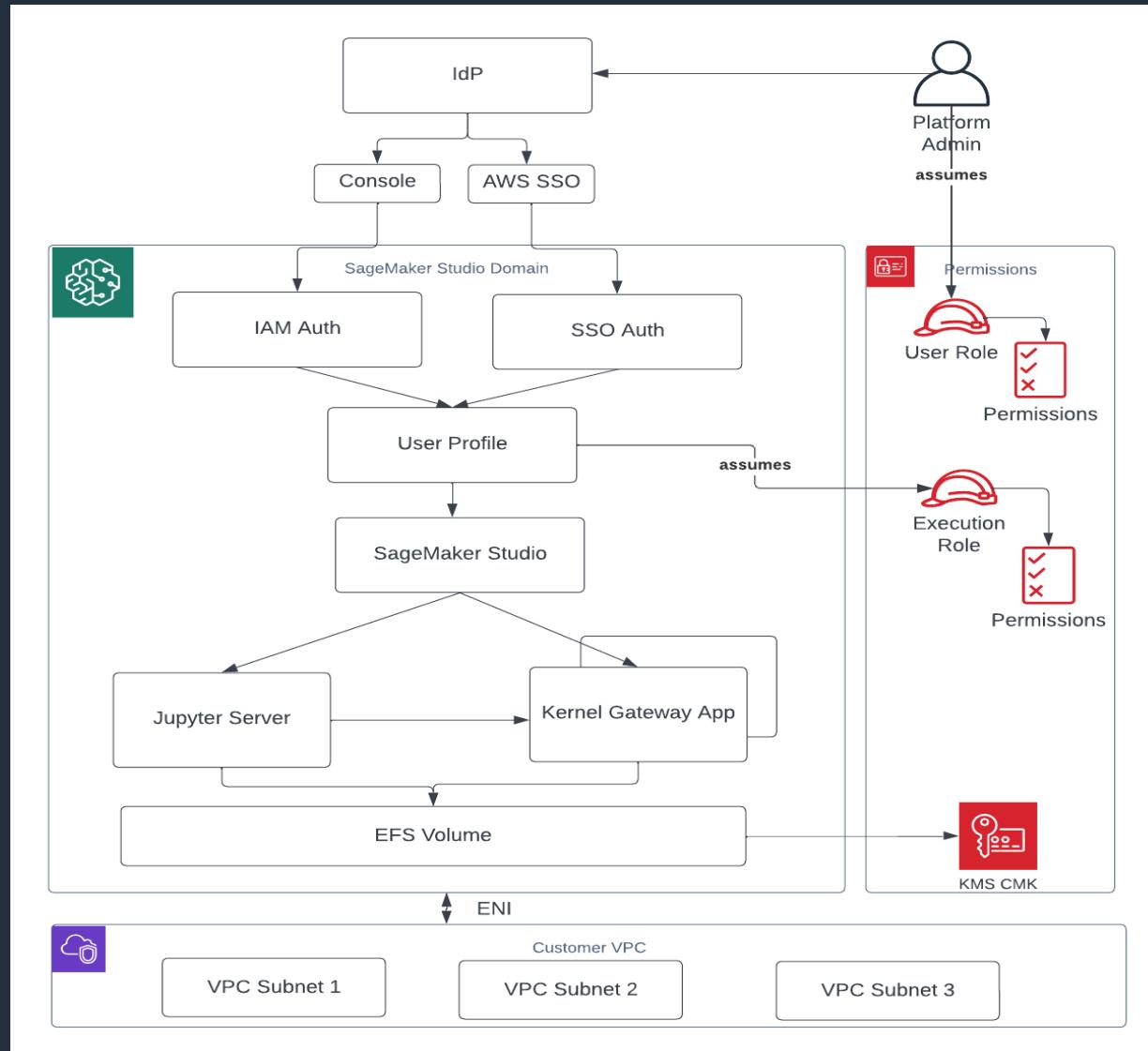


## Business Stakeholder ML Consumers

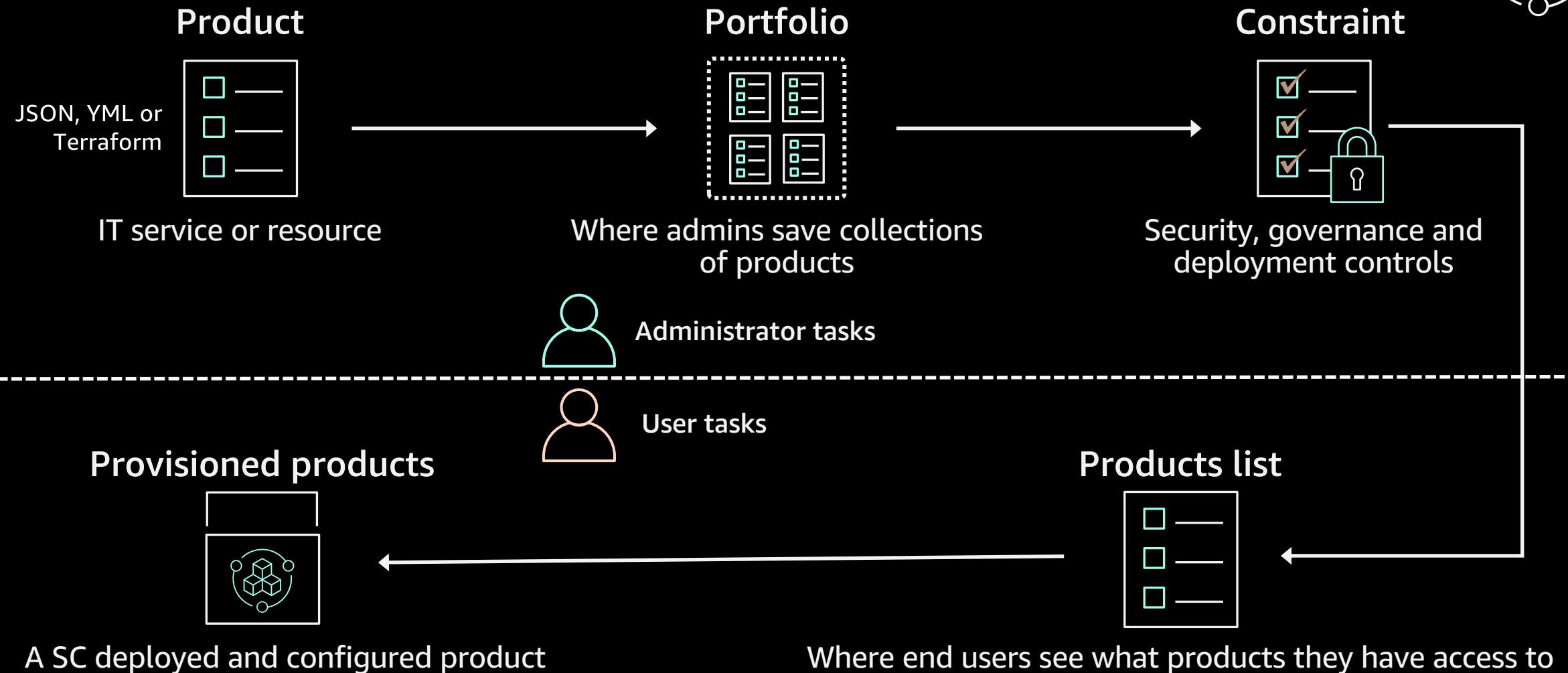
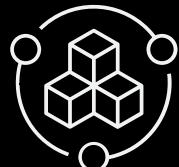
Consumers of ML results from other BUs, driving business decision making



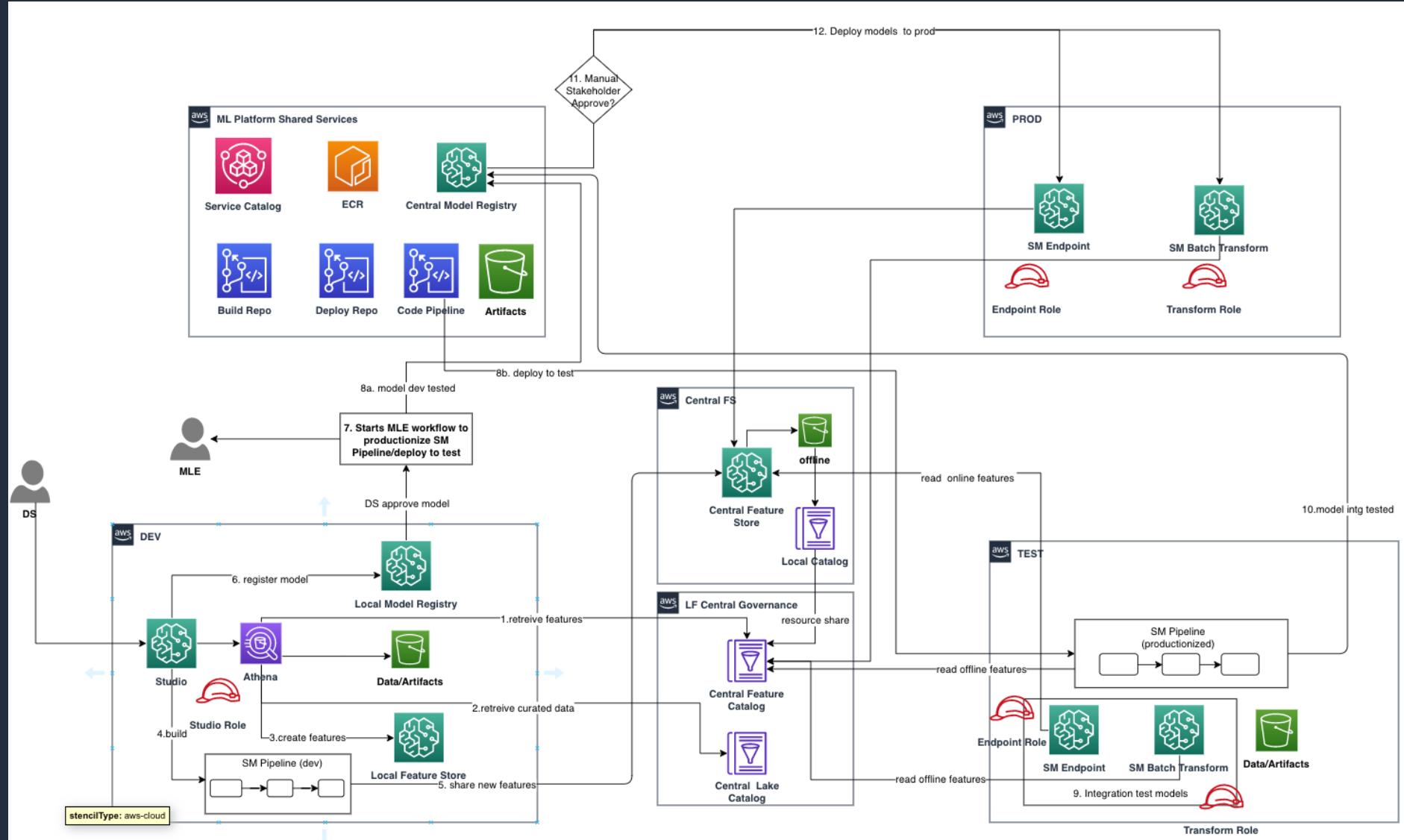
# Onboarding: SageMaker Domain Management



# AWS Service Catalog



# Central Model Registry and Feature Store Workflow



# 3. Observable ML

# Experiment Tracking with SageMaker Experiments

Runs 7

Search

Name	Id	Run Group	Type	Test:loss	Train:accuracy	Test:accuracy	Train:loss
exp-12-12-2022-18-39-07-run-6	Default-Run-Group-exp-12-1...	Default-Run-Group-exp-12-1...	--	0.04963957405090332	98.46	98.47	0.051136534420897564
exp-12-12-2022-18-39-07-run-5	Default-Run-Group-exp-12-1...	Default-Run-Group-exp-12-1...	--	0.10452271881103516	97.17	97.05	0.09755590740218759
exp-12-12-2022-18-39-07-run-4	Default-Run-Group-exp-12-1...	Default-Run-Group-exp-12-1...	--	0.058010439300537106	98.126666666666667	98.21	0.059526144051055115
exp-12-12-2022-18-39-07-run-3	Default-Run-Group-exp-12-1...	Default-Run-Group-exp-12-1...	--	0.09766765747070312	96.956666666666666	97.26	0.10599174875694638
exp-12-12-2022-18-39-07-run-2	Default-Run-Group-exp-12-1...	Default-Run-Group-exp-12-1...	--	0.053717627716064455	98.126666666666667	98.32	0.06060071001363297
exp-12-12-2022-18-39-07-run-1	Default-Run-Group-exp-12-1...	Default-Run-Group-exp-12-1...	--	0.11910663757324219	96.83333333333333	96.78	0.1146532233208418
exp-12-12-2022-18-39-07-run-0	Default-Run-Group-exp-12-1...	Default-Run-Group-exp-12-1...	--	0.14192571716308594	95.465	95.64	0.1484463992357254

Refresh Runs per page 10 Go to page 1 Page 1 of 1 < >

CHARTS 2

Compact View + Add Chart

Test:loss\_last

Step	run-0	run-1	run-2	run-3	run-4	run-5	run-6
1.0	0.14	0.18	0.20	0.19	0.17	0.16	0.15
2.0	0.13	0.15	0.17	0.16	0.14	0.13	0.12
3.0	0.12	0.14	0.16	0.15	0.13	0.12	0.11
4.0	0.11	0.13	0.15	0.14	0.12	0.11	0.10
5.0	0.10	0.12	0.14	0.13	0.11	0.10	0.09
6.0	0.11	0.13	0.15	0.14	0.12	0.11	0.10
7.0	0.12	0.14	0.16	0.15	0.13	0.12	0.11
8.0	0.11	0.13	0.15	0.14	0.12	0.11	0.10
9.0	0.10	0.12	0.14	0.13	0.11	0.10	0.09
10.0	0.14	0.12	0.10	0.11	0.09	0.08	0.07

Test:loss\_last

run	Test:loss
run-0	0.14
run-1	0.12
run-2	0.05
run-3	0.10
run-4	0.06
run-5	0.10
run-6	0.05



# Model inventory using SageMaker Model Registry

Model registry

Show introduction

+ Create model version + Create model group

Model group name	Description	Status	Created on	Created by
abalone-btdm-built-in-p-qdv3j9welp9x	--	✓ Completed	3 hours ago	--
abalone-btd-built-in-p-0qzy6hduz67s	--	✓ Completed	3 hours ago	--
aim321-customer-churn	--	✓ Completed	4 months ago	--

End of results

VERSION 1

Status	Pipeline	Execution	Project	Last Stage	Model group	Share	Update status
Pending	abalone-btd-built-in-p...	execution-167906962...	abalone-btd-builtin	staging	abalone-btd-built-in-p...		

Activity      Model quality      Explainability      Bias report      Inference recommender      Load test      Settings

Model metric	Metric value	Standard deviation
mse	4.573587555999638	2.135589655997731



# SageMaker Model Cards

Amazon SageMaker > Model cards > sentiment-analysis-model-card

## Model card - sentiment-analysis-model-card

Edit Clone Actions ▾

### Model card overview

Model card version	4	KMS encryption key	-
Model card status	Draft	Model card ARN	<a href="#">arn:aws:sagemaker:us-east-2:██████████:model-card/sentiment-analysis-model-card</a>
Created date	11/14/2022, 10:17:18 PM		

### Model overview

Model name	Sentiment-Analysis-Model	Inference environment	<a href="#">257758044811.dkr.ecr.us-east-2.amazonaws.com/sagemaker-xgboost:1.3-1</a>
Model description	the model is updated.	Problem type	Binary Classification
Model versions	-	Algorithm type	Logistic Regression
Model arn	<a href="#">arn:aws:sagemaker:us-east-2:██████████:model/sentiment-analysis-model</a>	Model creator	DEMO-user
Model artifacts	<a href="#">s3://sagemaker-studio-us-east-2:██████████/xgboost-churn/output/workshop-xgboost-customer-churn-2022-11-10-23-13-38-509/output/model.tar.gz</a>	Model owner	DEMO-user

Intended uses Training details Evaluation results Additional details Version history

### Intended uses

Intended uses of the model  
Not used except this test.

Factors affecting model efficacy  
No.

Risk rating  
Low

Explanations for risk rating  
Just an example.



# SageMaker Model Dashboard

Amazon SageMaker > Model dashboard > Sentiment-Analysis-Model

## Sentiment-Analysis-Model [Info](#)

[Edit Model Card](#)

### Model overview [Info](#)

Model card

[sentiment-analysis-model-card](#)

Model lineage

[View lineage](#)

Additional model details

[Sentiment-Analysis-Model](#)

Model card risk rating

Low

### Endpoints [Info](#)

Endpoint name

[Sentiment-Analysis-Model-Endpoint](#)

Endpoint status

 In Service

Creation Date

Nov 15, 2022 03:58 UTC

Last modification time

Nov 15, 2022 04:01 UTC

### Monitor schedule [Info](#)

[Activate/ Deactivate monitor schedule](#)[Edit alert](#)

Schedule name

[sentiment-data-monitor-schedule-2022-11-15-04-10-39](#)

Endpoint name

[Sentiment-Analysis-Model-Endpoint](#)

Monitor type

DataQuality

Monitor frequency

Every hour

Schedule status

 Scheduled

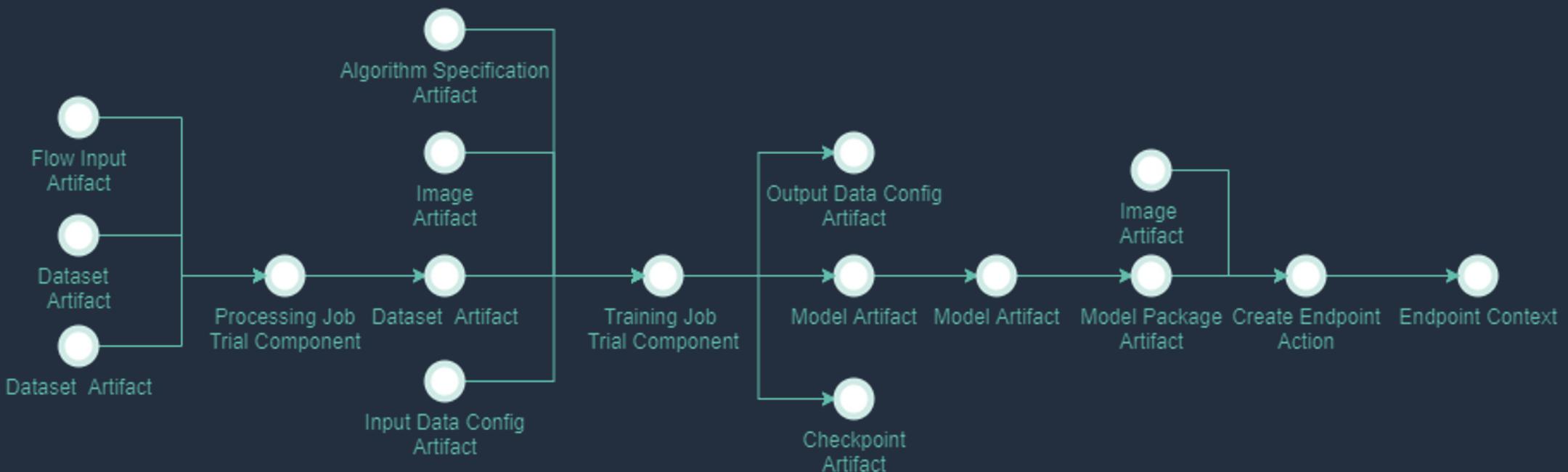
Alert



# SageMaker ML Lineage Tracking

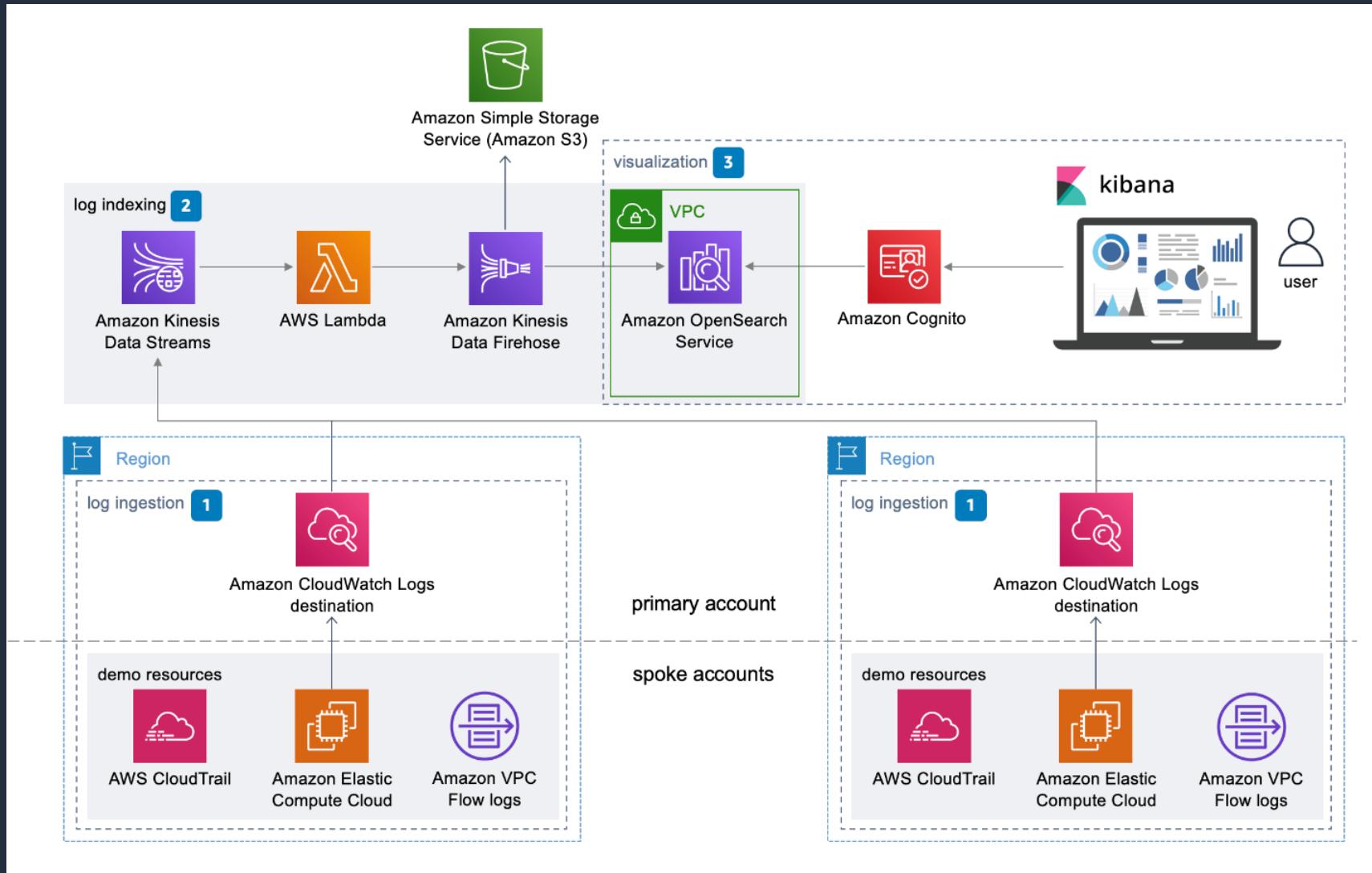
## Lineage Metadata

SageMaker automatically creates a connected graph of lineage entity metadata tracking your workflow.



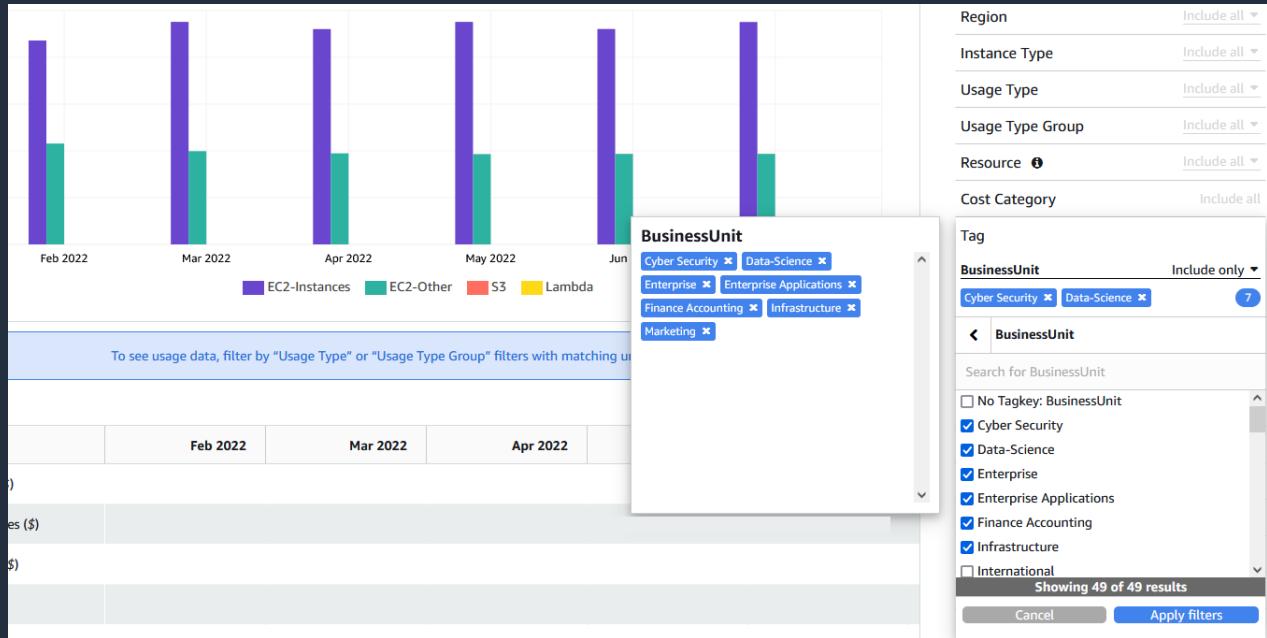
**Use the SageMaker Python SDK to query lineage and manually track entities**

# Centralized Logging and Monitoring



# Cost Control and Reporting in SageMaker

- Cost allocation for ML environments and workloads using resource tagging in SageMaker
- Enforce tagging using IAM policies
  - Use *Tag policies* to standardize tagging at the organization level
- View ML spend through **AWS Cost Explorer** and **AWS Cost and Usage Report (CURs)**
- Use **AWS Budgets** to provide early warning when ML costs exceed (or forecasted to exceed) a threshold

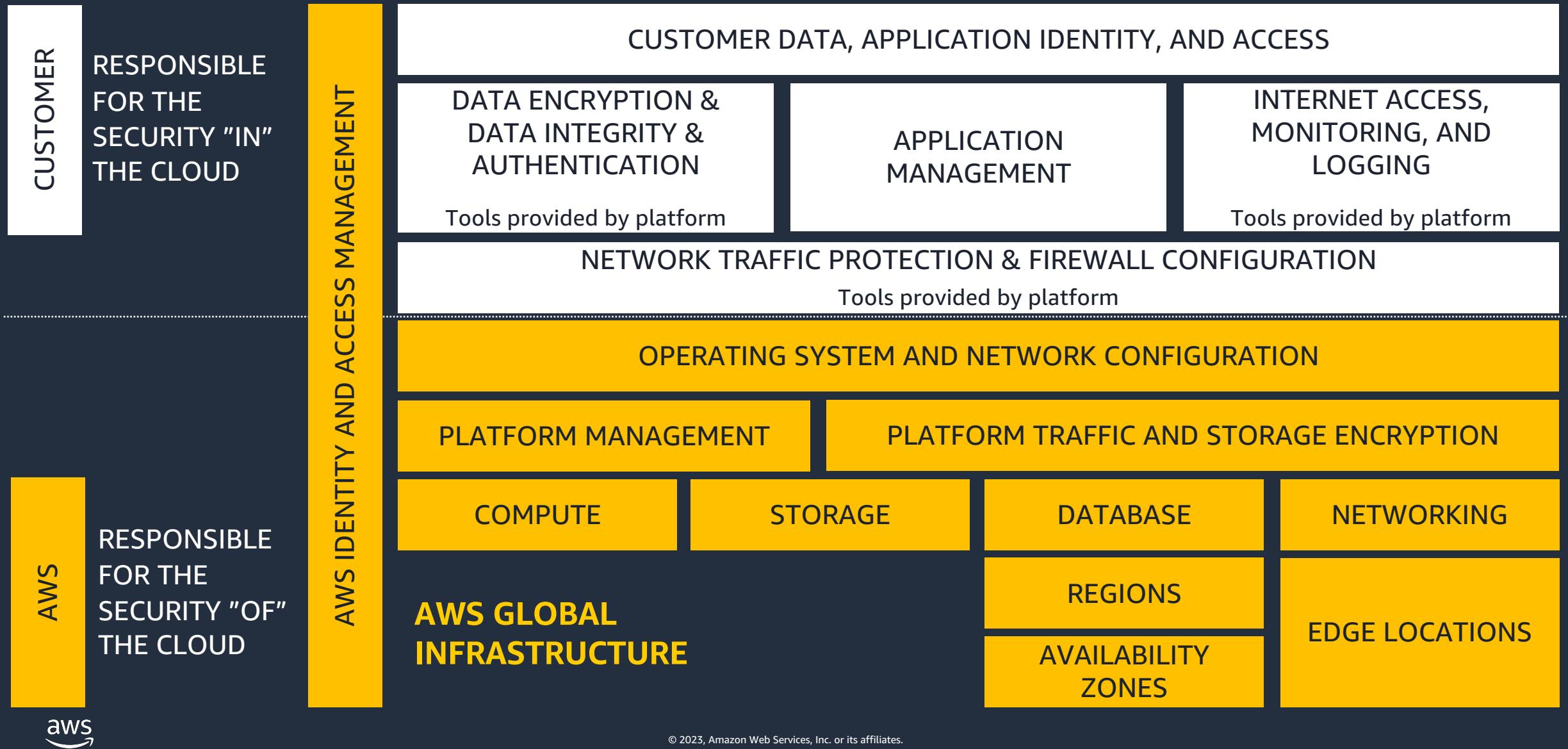


AG	AH	AI	AJ	AK	AL	AM	AN	AO
product/region	product/serviceName	product/usageType	product/vcpu	pricing/public	pricing/term	pricing/unit	resourceTags/user:cost-center	
us-west-2	Amazon SageMaker	USW2-Host:ml.m5.xlarge	4	0.23	OnDemand	Hrs	TF2WorkflowEndpoints	
us-west-2	Amazon SageMaker	USW2-Host:ml.m5.xlarge	4	0.23	OnDemand	Hrs	TF2WorkflowEndpoints	
us-west-2	Amazon SageMaker	USW2-Host:ml.m5.xlarge	4	0.23	OnDemand	Hrs	TF2WorkflowEndpoints	
us-west-2	Amazon SageMaker	USW2-Host:ml.m5.xlarge	4	0.23	OnDemand	Hrs	TF2WorkflowEndpoints	
us-west-2	Amazon SageMaker	USW2-Host:ml.m5.xlarge	4	0.23	OnDemand	Hrs	TF2WorkflowEndpoints	
us-west-2	Amazon SageMaker	USW2-Host:ml.m5.xlarge	4	0.23	OnDemand	Hrs	TF2WorkflowEndpoints	

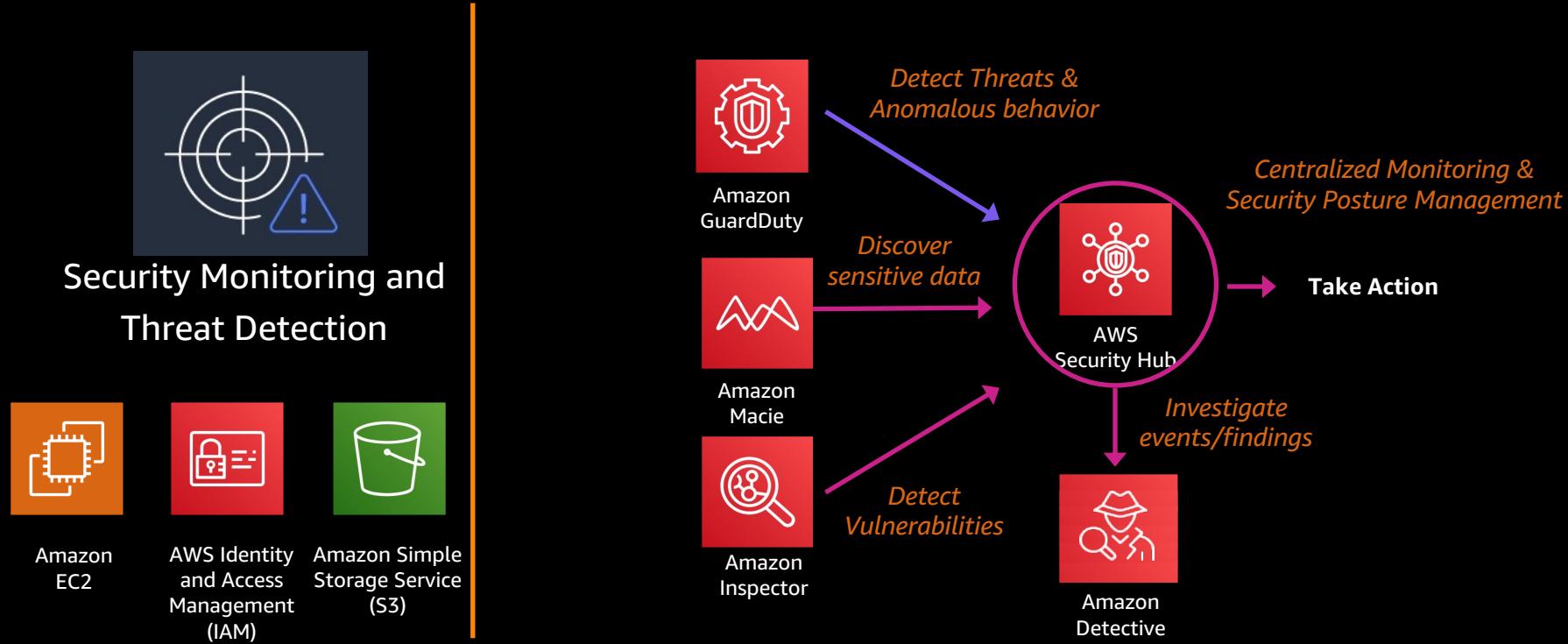


# 4. Secure ML

# SHARED RESPONSIBILITY MODEL: SAGEMAKER STUDIO



# AWS Security Hub: Centralize Security Monitoring



# Amazon SageMaker security features help you go from idea to production faster



## Infrastructure and network isolation

Control data traffic across SageMaker components over a private network, and ensure appropriate ingress/egress with single-tenancy

## Authentication and authorization

Define, enforce, and audit who can be authenticated and authorized to use Amazon SageMaker resources

## Data protection

Ensure automatic data encryption at rest and in transit with flexibility to bring your own keys

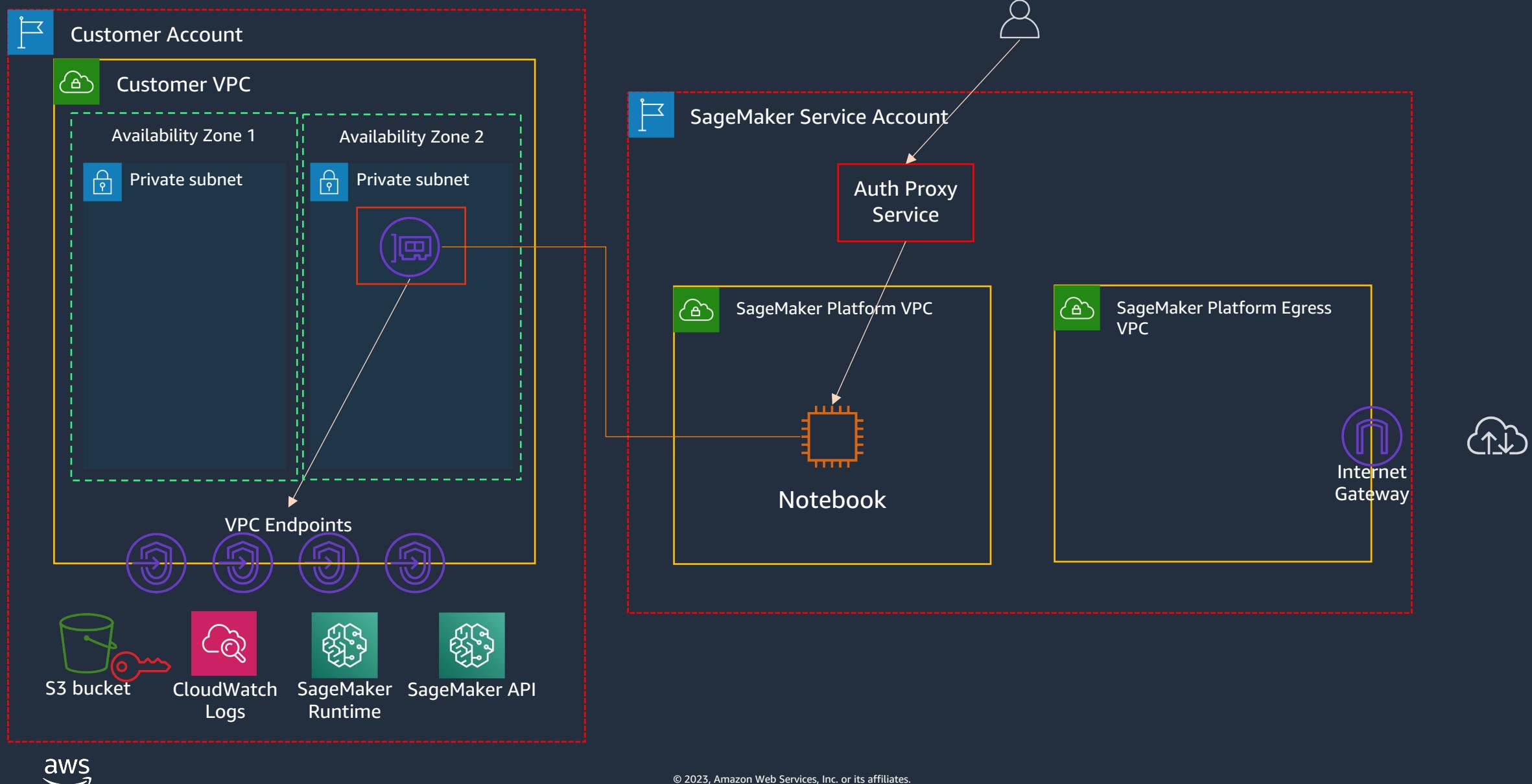
## Monitoring and auditability

Track, trace, and audit all API calls, events, data access, or interactions down to the user and IP level to ensure quick remediation

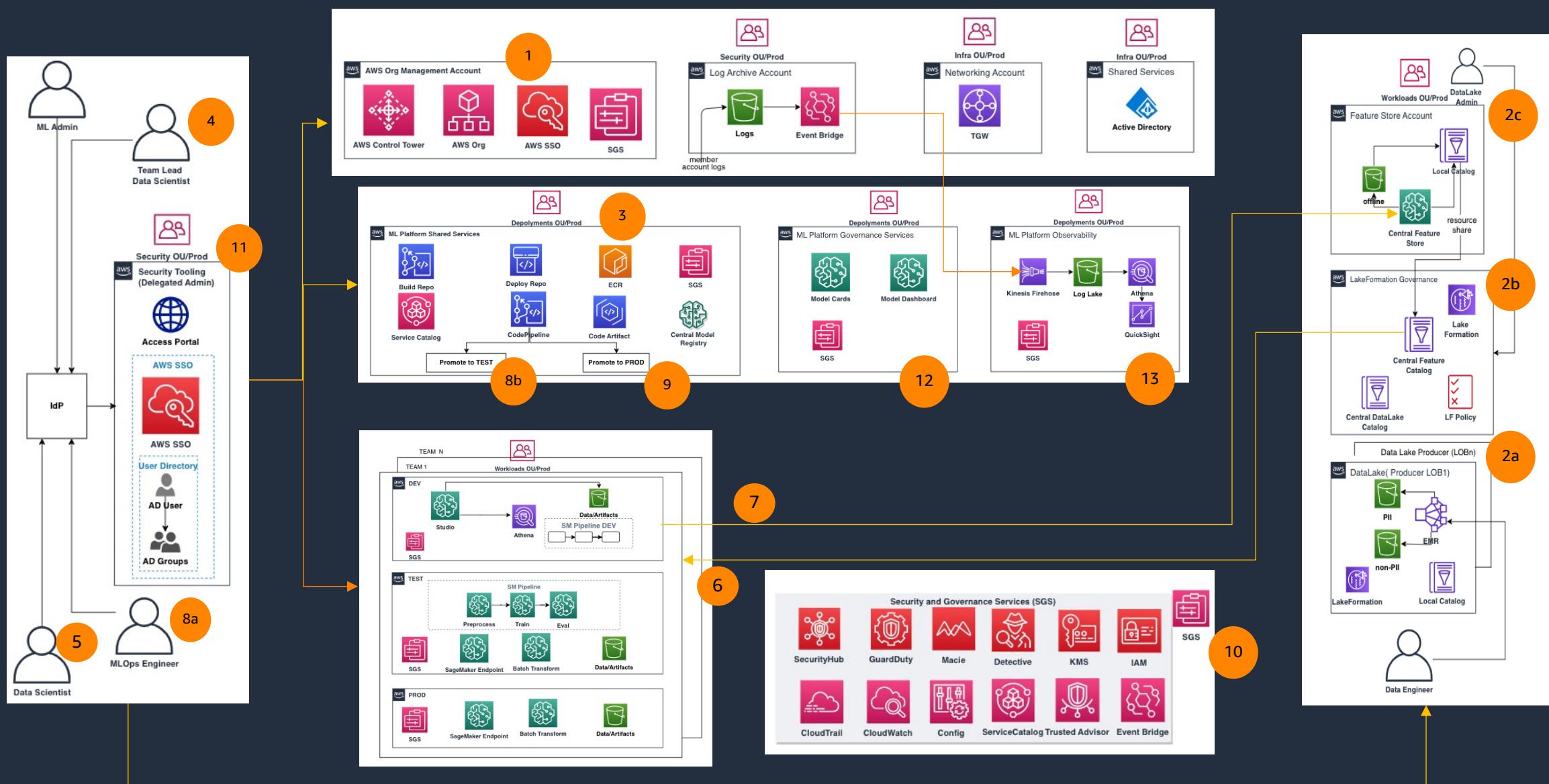
## Compliance certifications

Inherit the most comprehensive compliance controls, and easily abide by your industry's legislation

# SageMaker Studio - VPC connectivity (recommended)



# Governed ML Platform Reference Architecture



# Key Takeaways

---

01

## Build Strong ML Platform Foundations

Establish right operating model with multi-account and data strategies

---

02

## Scale ML Platform Operations

Centralize Shared Services and Scale operations with platform level services e.g. on-boarding, CI/CD

---

03

## Create an Observable ML Platform

Centralize logging, monitoring, experiment tracking and leverage ML governance services

---

04

## Secure your ML Platform

Centralize Identity, Security, Network services, and leverage preventative & detective controls



# Resources

- [SageMaker Studio Best Practices WhitePaper](#)
- [Building Secure ML Platform Whitepaper](#)
- [Onboarding SageMaker Studio with AWS SSO and Okta](#)
- [ML Platform Setup Blogs](#)
- [Macie Sensitive Data Detection](#)
- [SageMaker workshops for Studio](#)
  - [SageMaker End-to-End Workshop](#)
  - [SageMaker Immersion Day](#)
  - [Secure Data science with SageMaker Studio](#)