



Introduction to MLOps

Bringing DevOps and Automation to Machine Learning

Hei Chow
Solutions Architect

Current state of AI/ML

State of machine learning

• Today

- 53% of POCs make it into production
- Average 9 months
- Gartner

Last decade

- Focusing mostly on building ML models
- Operationalization was an afterthought

By end of 2024

- 75% of organizations will shift from piloting to operationalizing AI
- Gartner

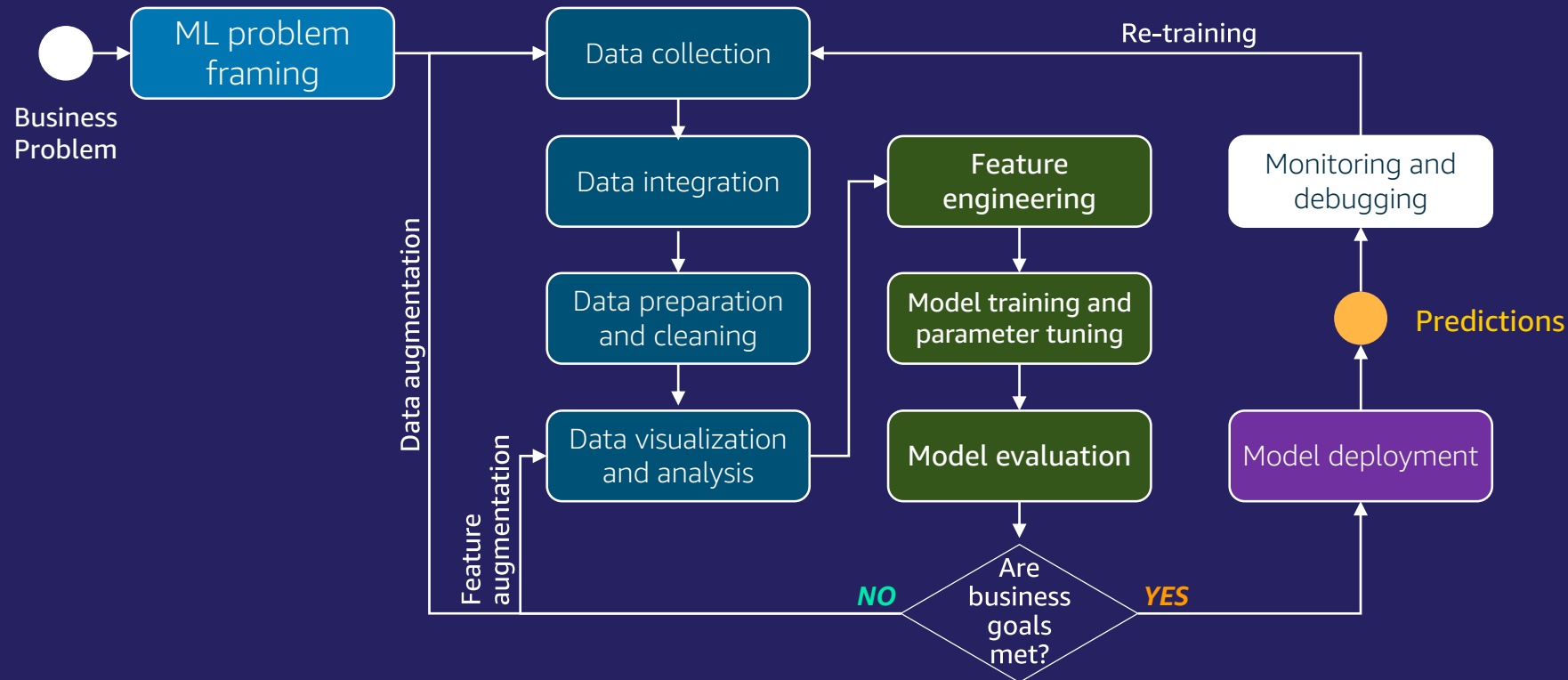
<https://www.idgconnect.com/article/3583467/gartner-accelerating-ai-deployments-paths-of-least-resistance.html>

Main Challenges

- Publishing a ML model is not enough.
- Managing the published ML models is as important as developing them.
- *"IT leaders responsible for AI are discovering '**AI pilot paradox**', where launching pilots is deceptively easy but deploying them into production is notoriously challenging."*
- **Chirag Dekate**, Vice President Analyst, Gartner

From DevOps to MLOps

The ML process



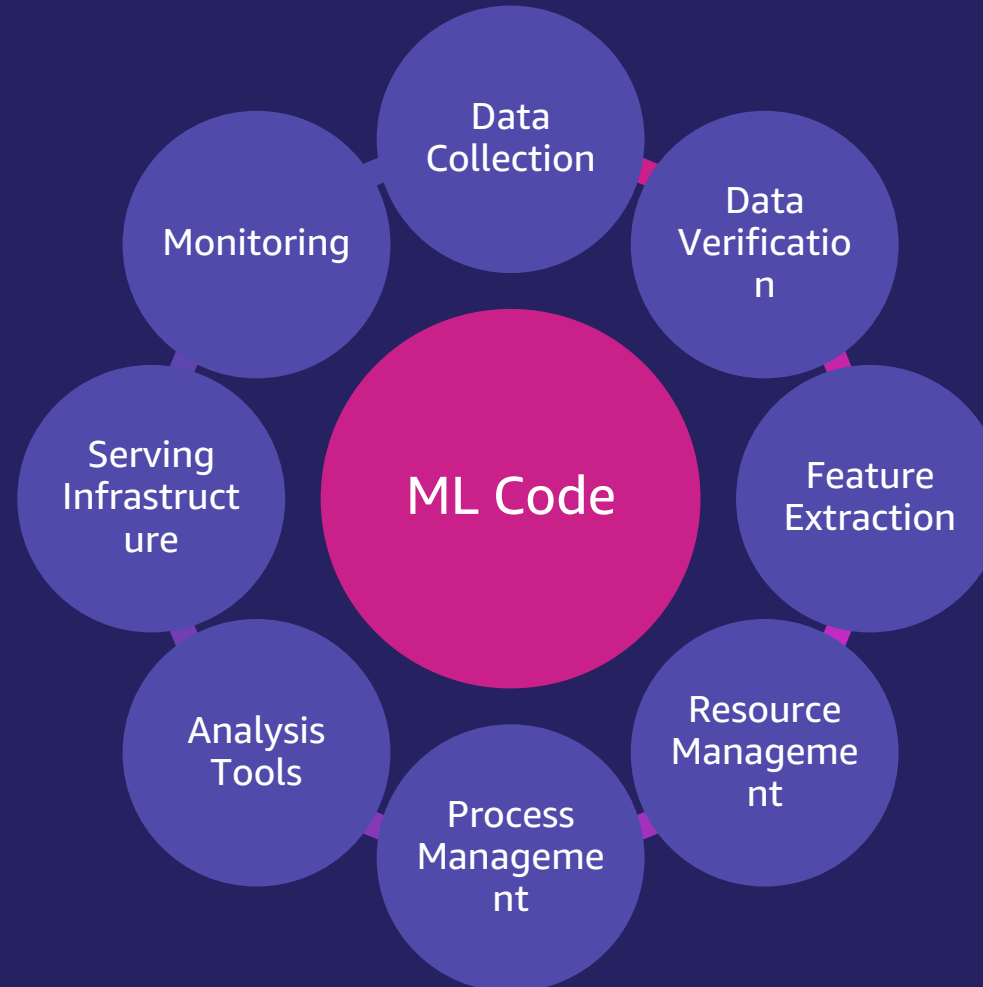
Phase 1: Research/Experiment

Question: "Can we use ML to solve this?"

- *"Is it possible to ... ?"*
- *"Can we use this data to solve the following problem?"*
- *"Surely we must be able to ..."*

Typical scenarios

- Scientific projects
- Proof-of-concepts (PoCs)



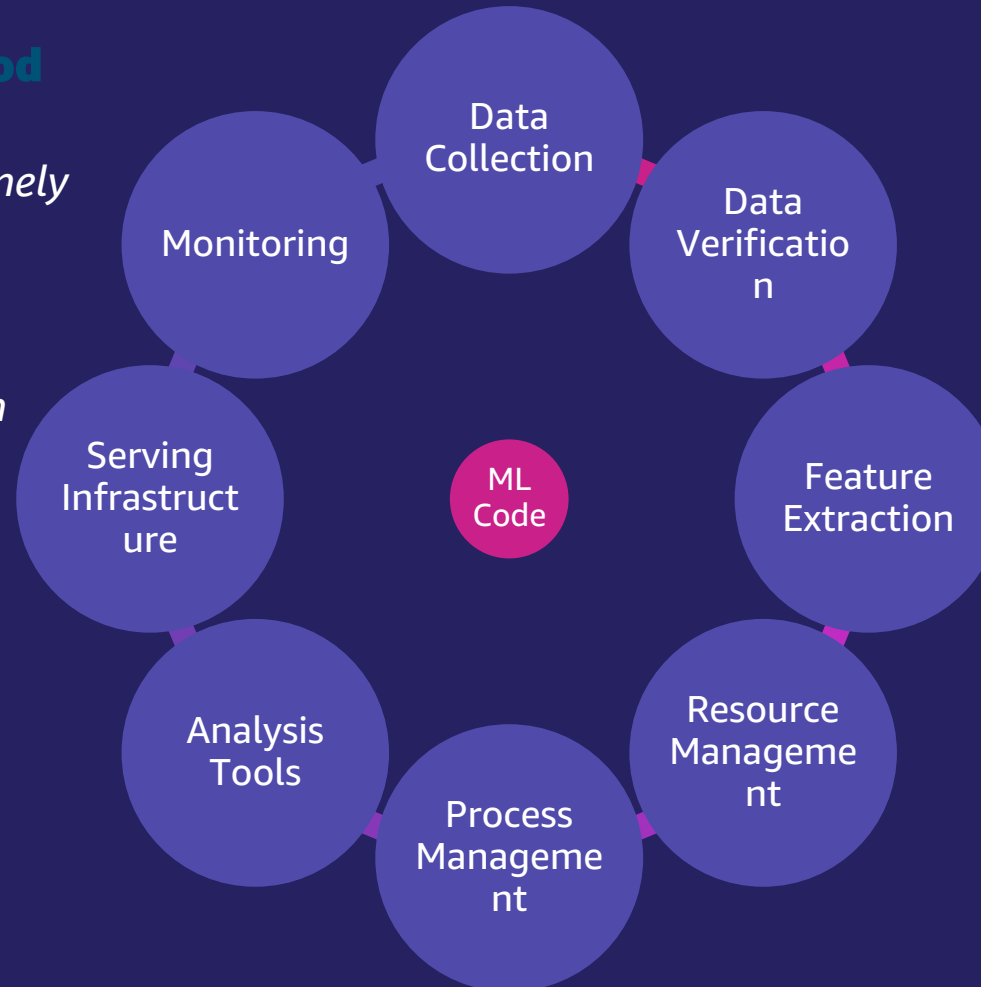
Phase 2: Operational

Question: “How do we implement this method at scale?”

- *How do we pipe the data into the model in a timely fashion?*
- *How do we collect, store and transform data so models can be retrained consistently?*
- *How do we build an A/B testing environment, in order to test future model iterations?*

Typical scenarios

- After PoC, bringing your ML models to production
- Migration of existing models into ML platform



MLOps – Why?



Agility

- Continuous and faster deliveries
- Faster modifications
- Faster bug-fixing



Experiments

- Faster and Controlled Experiments
- Faster integration of successful experiments to other environments



Scalability

- Ease integration of new ML model
- Standarization of code
- Lower operational costs



Time to Market

- Reduced time-to-market
- Faster planning and delivery expectations



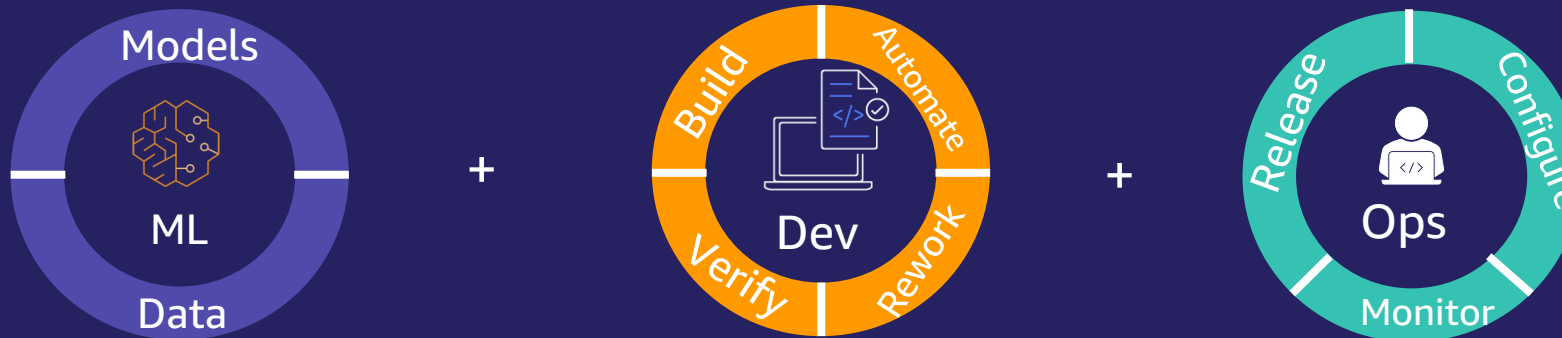
Business Owners

- Strong collaboration
- Improve iterations

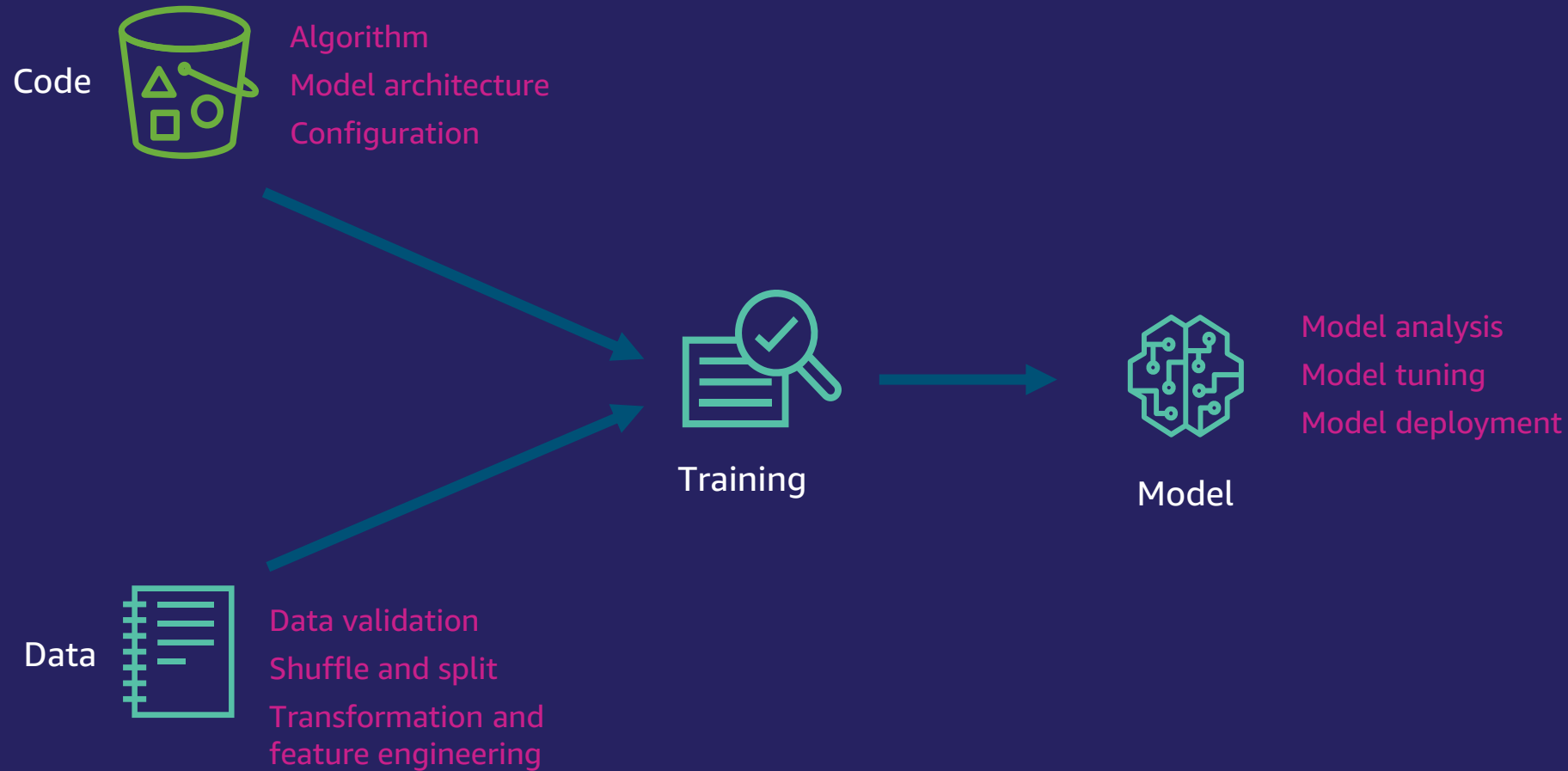
MLOps – What?

ML + Dev + Ops = MLOps

Collaborative and experimental in nature | Automate as much as possible |
Continuous improvement of ML Models | Standardize and Scale



ML Code and Data are Independent



How is MLOps different from DevOps?

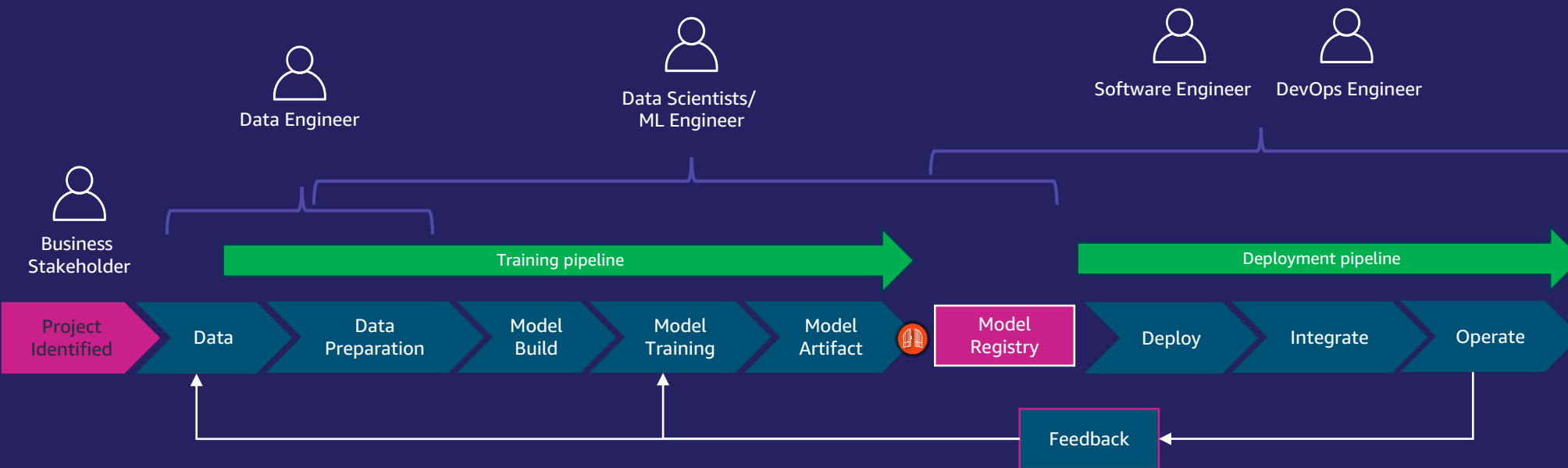
	DevOps	MLOPS
Code versioning	✓	✓
Compute environment	✓	✓
Continuous integration/delivery (CI/CD)	✓	✓
Monitoring in production	✓	✓
Data provenance		✓
Datasets		✓
Models		✓
Hyperparameters		✓
Metrics		✓
Workflows		✓

MLOPS

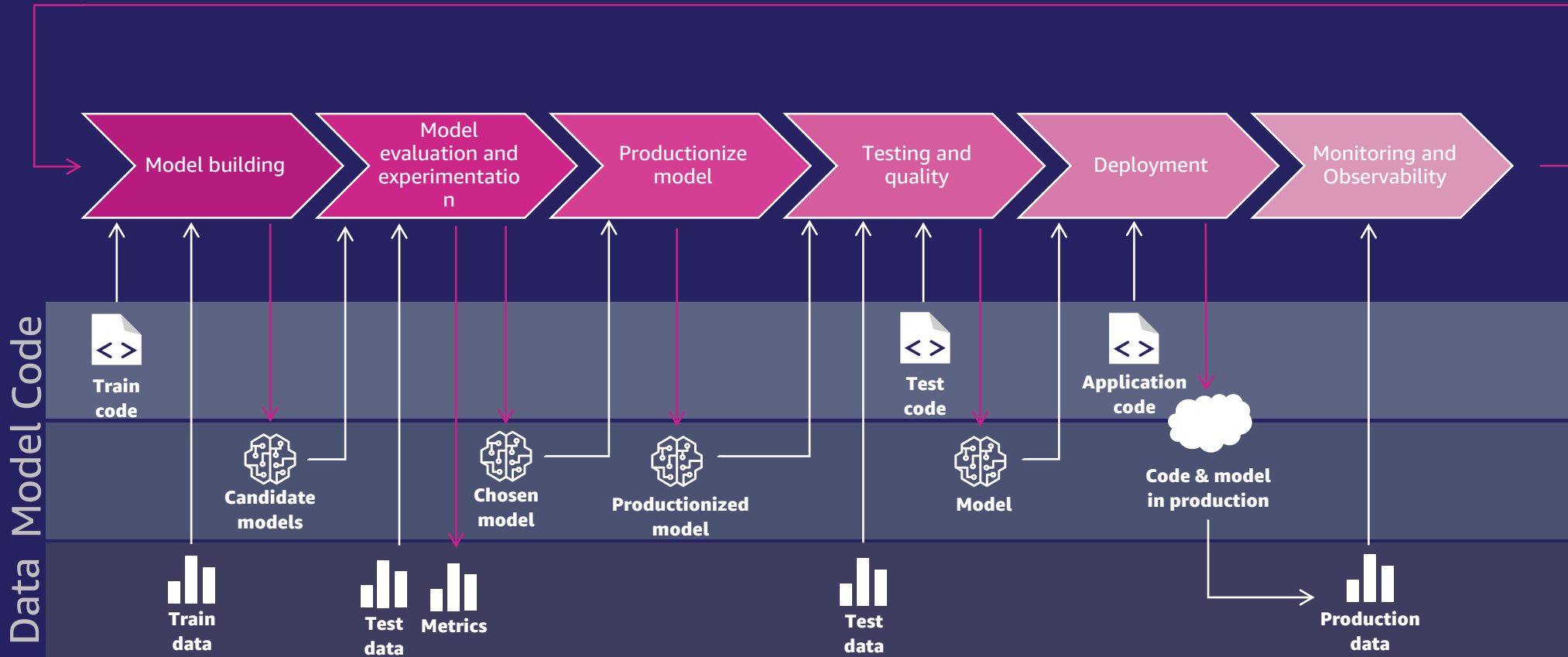
End-to-end ML
lifecycle
management

<https://medium.com/analytics-vidhya/mlops-the-epoch-of-productionizing-ml-models-4eec06d93623>

MLOps practices



ML lifecycle management



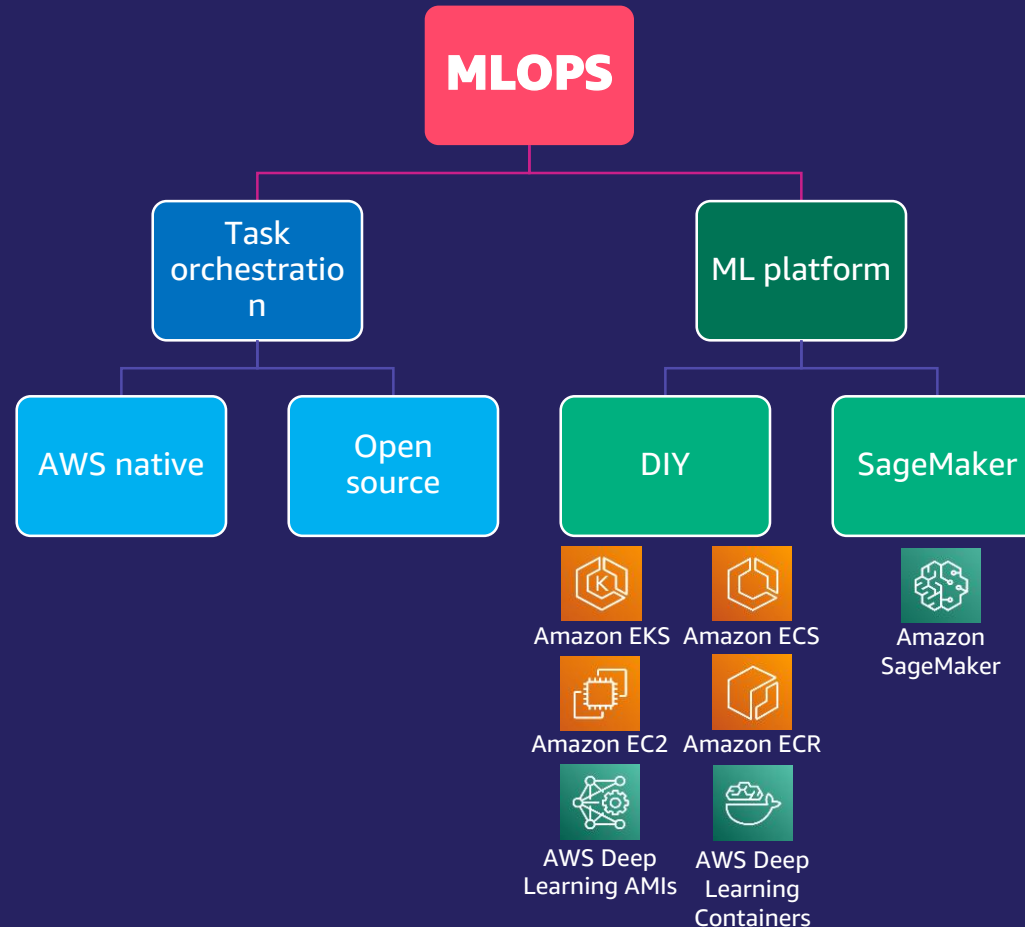
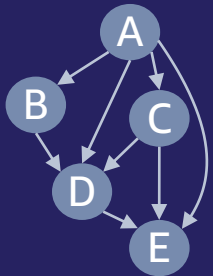
Automating ML Workflows

using SageMaker



Technology components in MLOps

- Create and manage workflows
- Automate ML steps & pipelines
- Implement CI/CD
- Form a Directed Acyclic Graph (DAG)



- ML development, experimentation, collaboration
- Compute/training environment
- Model registry
- Feature store
- Model deployment
- Monitoring in production
- Hyperparameter optimization
- Dataset management

Amazon SageMaker

Most complete, end-to-end ML service

Integrated Workbench

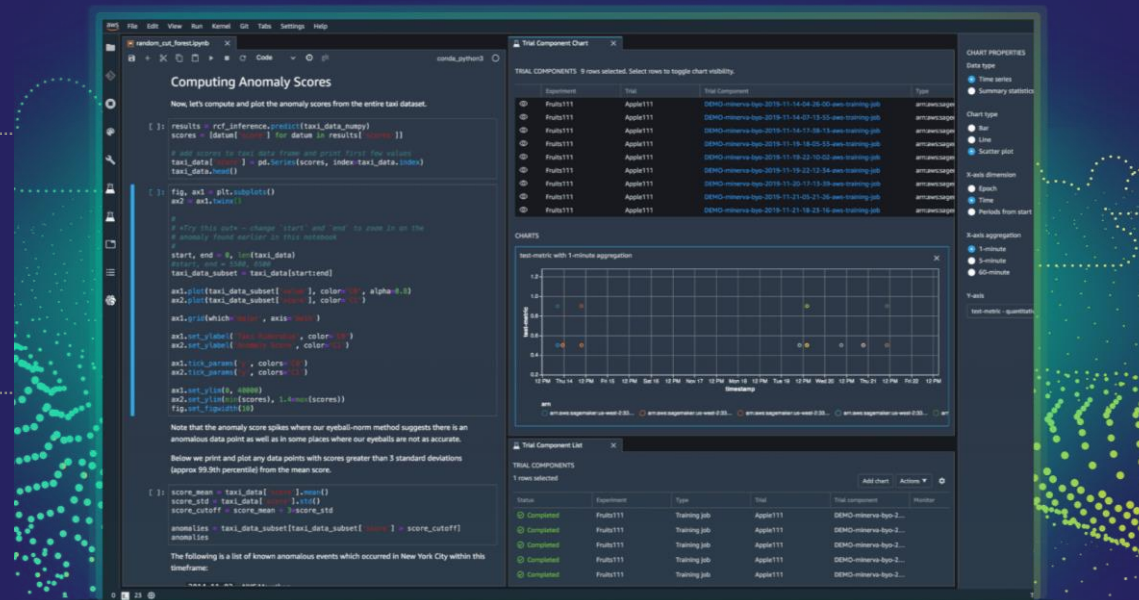
Capabilities designed specifically for ML, data preparation, experiment management, and workflows

Managed Infrastructure

Designed for ultra low latency and high throughput, automatic scaling, and distributed training

Managed Tooling

Purpose-built from the ground up to work together including auto ML, collaboration, debugger, profiler, bias analyzer, and explainability

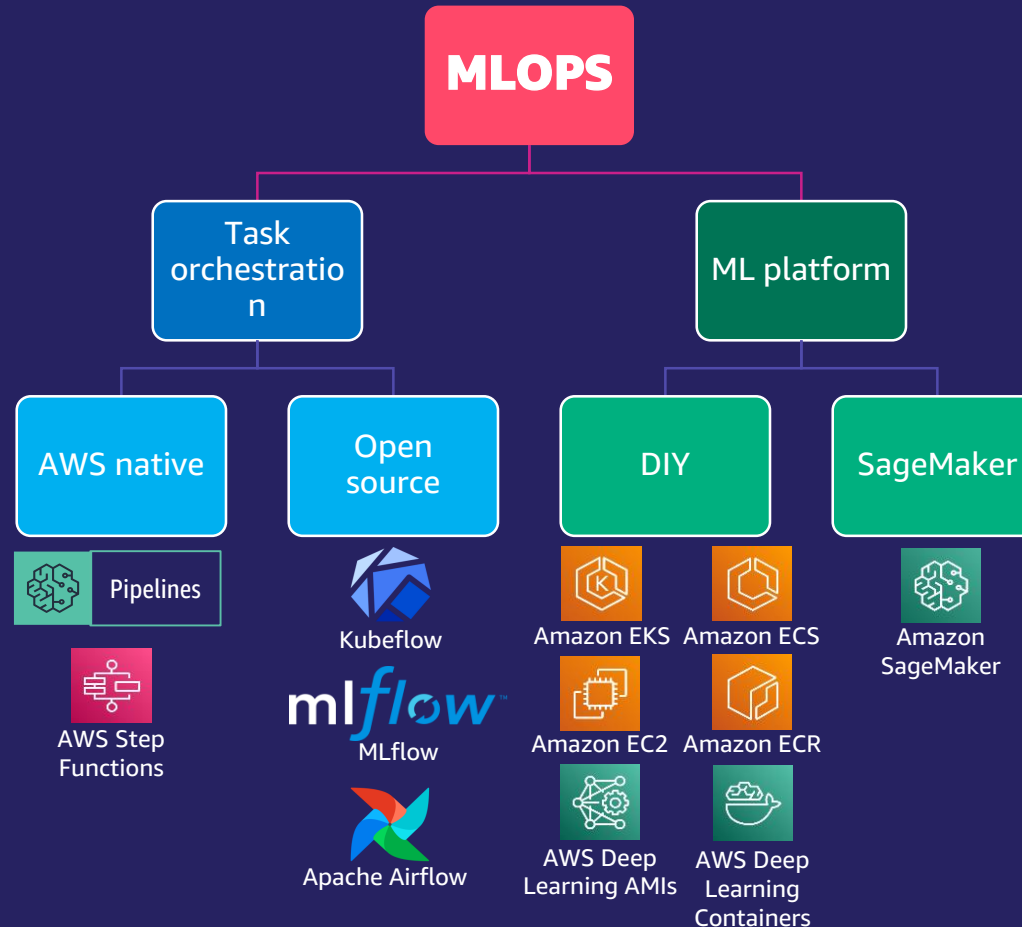
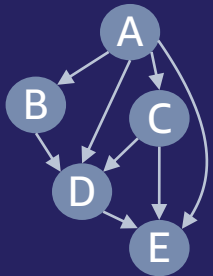


Amazon SageMaker Features



Technology components in MLOps

- Create and manage workflows
- Automate ML steps & pipelines
- Implement CI/CD
- Form a Directed Acyclic Graph (DAG)



- ML development, experimentation, collaboration
- Compute/training environment
- Model registry
- Feature store
- Model deployment
- Monitoring in production
- Hyperparameter optimization
- Dataset management

Task orchestration

Open source 3rd party options



MLflow

Open source platform for the ML lifecycle



Apache Airflow

Platform to author, schedule and monitor workflows



Kubeflow

ML toolkit for Kubernetes

Native AWS options



AWS Step Functions

Serverless pipeline orchestration




Amazon SageMaker Pipelines

Managed ML pipelines in SageMaker Studio

Native integration with SageMaker

Apache Airflow

- SageMaker Operators in Apache Airflow
-  Amazon Managed Workflows for Apache Airflow

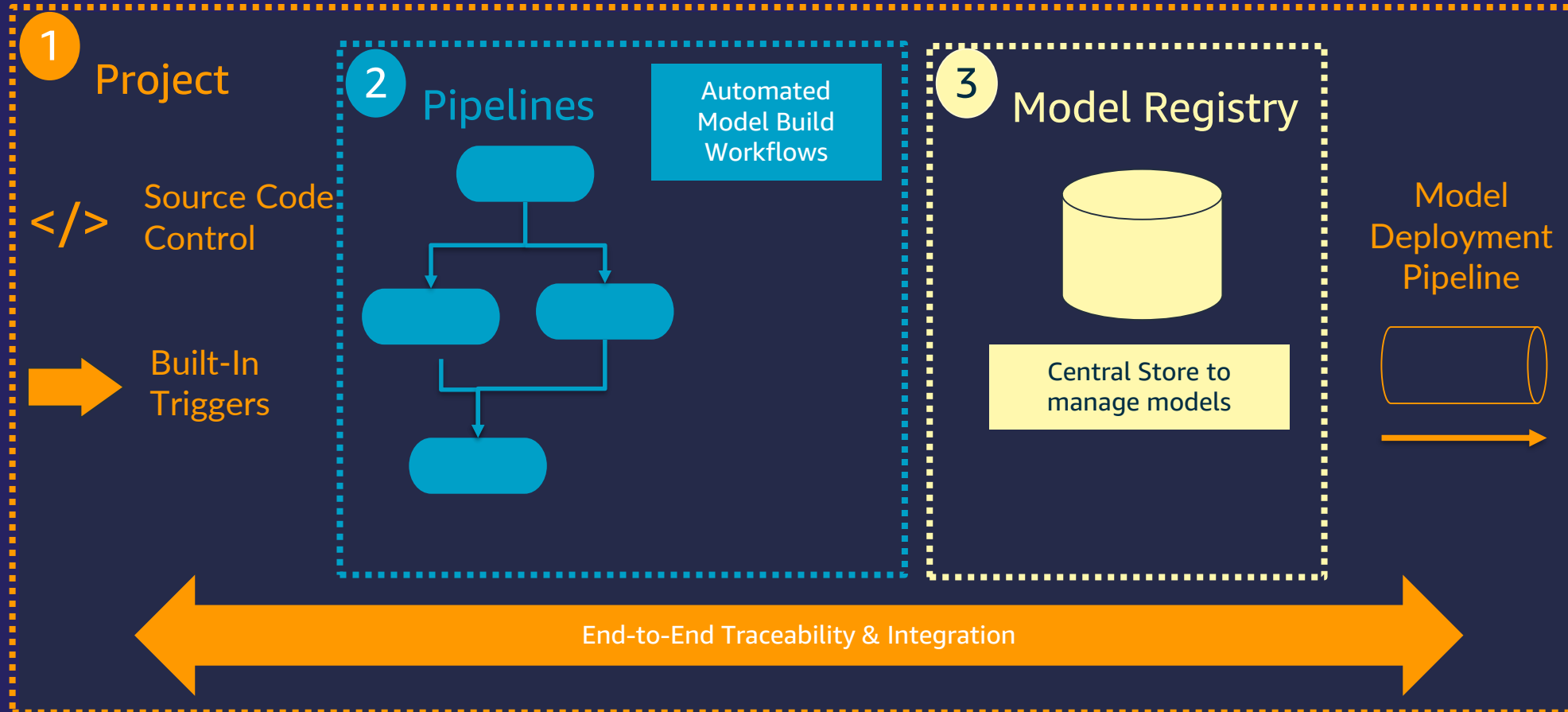
(managed Airflow service)

Kubeflow & Kubernetes

- SageMaker Components for Kubeflow Pipelines
- SageMaker Operators for Kubernetes

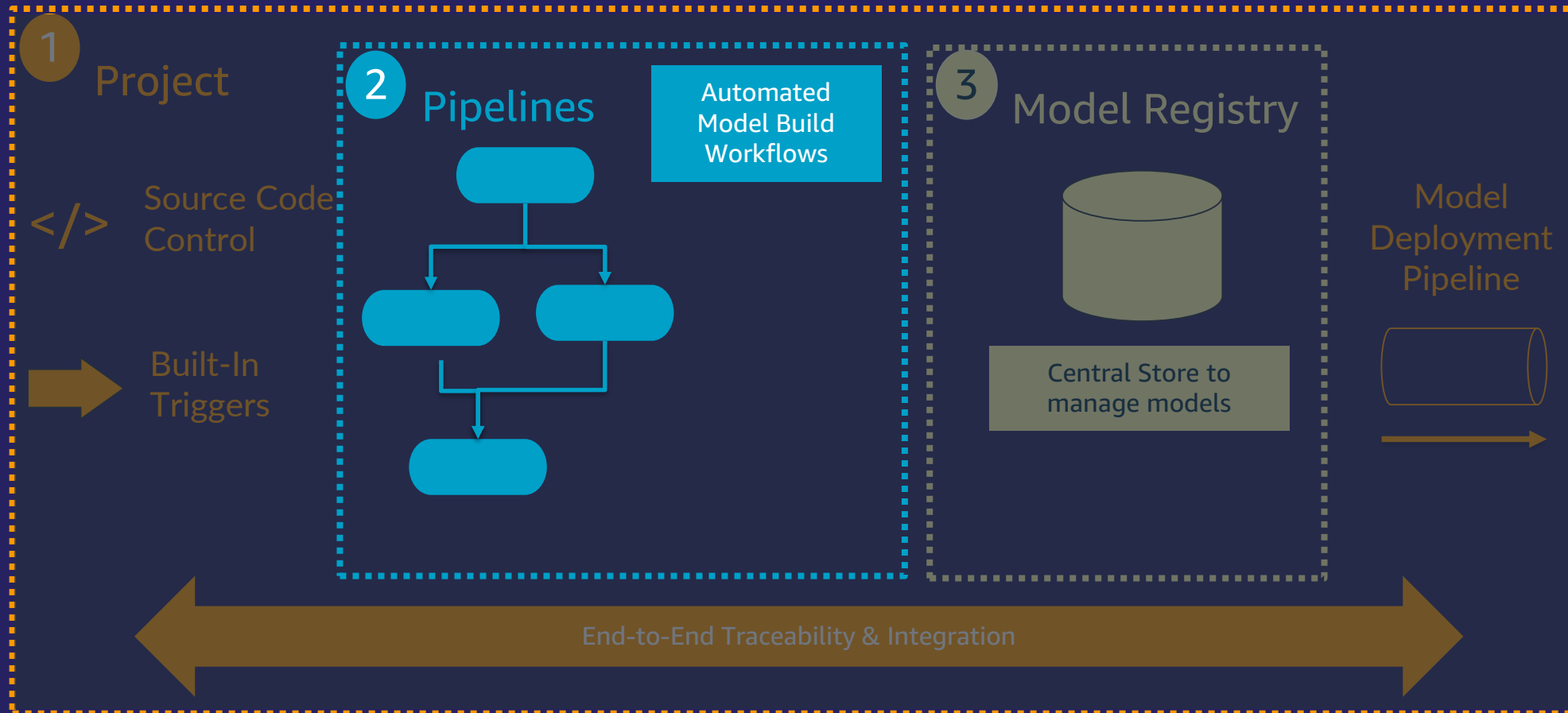


Amazon SageMaker Pipelines Components



Amazon SageMaker Pipelines

Components – Pipelines



Amazon SageMaker Pipelines

Components – Pipelines

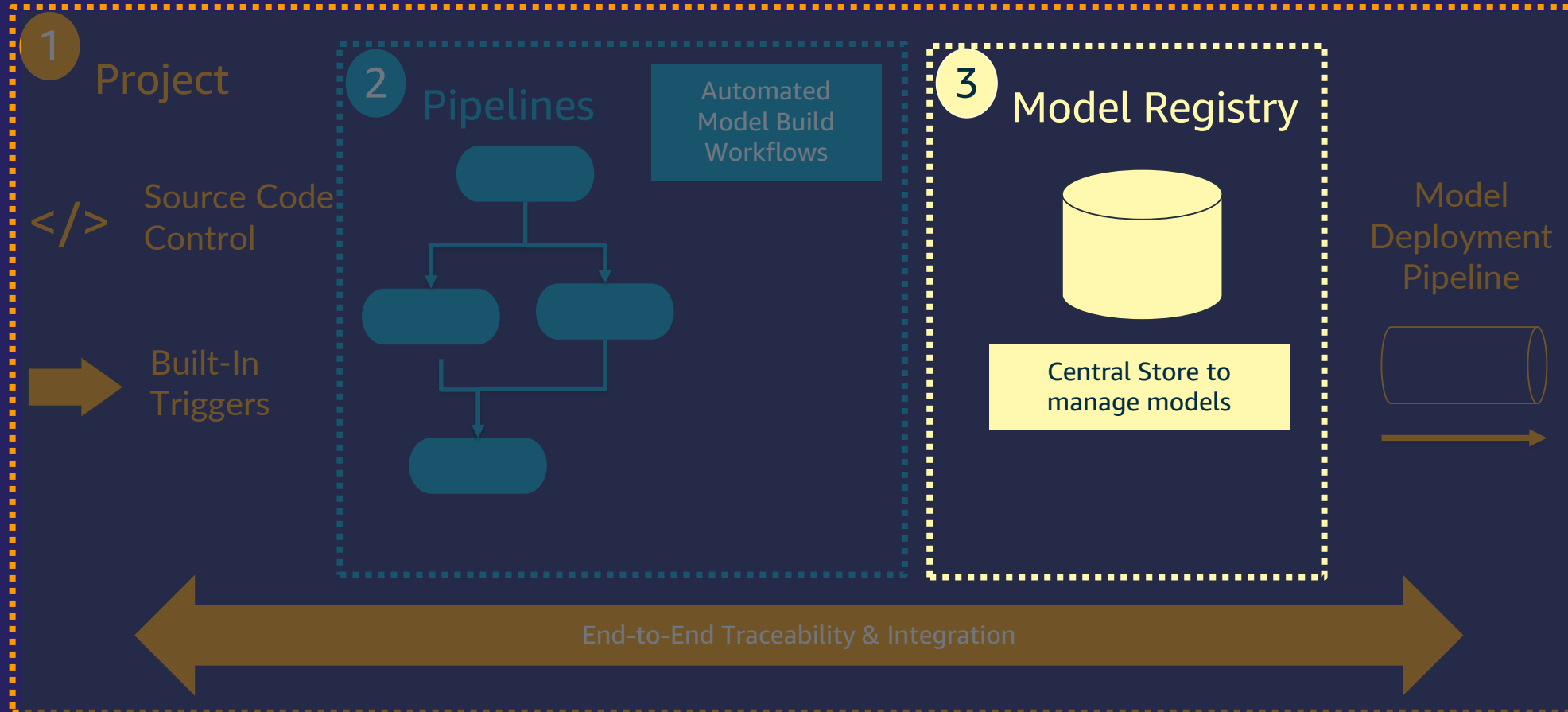
The screenshot displays the Amazon SageMaker Pipelines console interface. On the left, a sidebar titled 'Components and registries' shows a list of pipelines, including 'shelbee-demoagain-p...', 'shelbee-btd-abalone-p...', 'shelbee-demo-p-rbxf5...', 'shelbee-bt-p-lg35zjypy...', and 'shelbee-btd-p-ocmm2...'. The main panel shows the details of a specific pipeline execution, 'execution-1607360596599', which is in a 'Status' of 'Completed' (indicated by a green dot). The execution started on '12/7/2020, 10:03 AM' and took '12m10s' to complete. Below this, a 'Graph' tab displays a flowchart of the pipeline steps: 'PreprocessAbaloneData' leads to 'TrainAbaloneModel', which leads to 'EvaluateAbaloneModel', then 'CheckMSEAbaloneEvaluation', and finally 'RegisterAbaloneModel'. The 'CheckMSEAbaloneEvaluation' step has a 'true' condition leading to 'RegisterAbaloneModel'. The console also shows a list of tabs at the top, including 'untitled.f', 'shelbee-c', 'executior', 'Create pr', and several 'shelbee-c' and 'executior' tabs.

Supported Steps:

- Processing
- Training
- Tuning
- Conditional
- Register Model
- Create Model

Amazon SageMaker Pipelines

Components – Model Registry



Amazon SageMaker Pipelines

Components – Model Registry

The screenshot displays the Amazon SageMaker Model Registry interface. The top section shows a list of model versions with columns for Version, Stage, Status, Short description, Modified by, and Last modified. Below this, a detailed view for 'Version 2' is shown, including a status bar with 'Approved' status, pipeline name, execution name, last stage, and model group, along with an 'Update status' button. The 'Metrics' tab is active, showing a table of model metrics.

Version	Stage	Status	Short description	Modified by	Last modified
3	None	Pending			
2	prod	Approved		workshop-user	
1	staging	Approved		shelbee-iggy	

Version 2		
Status	Pipeline	Execution
Approved	shelbee-demoagain-p-...	workflow-2
Last Stage	Model group	
prod	shelbee-demoagain-p-...	

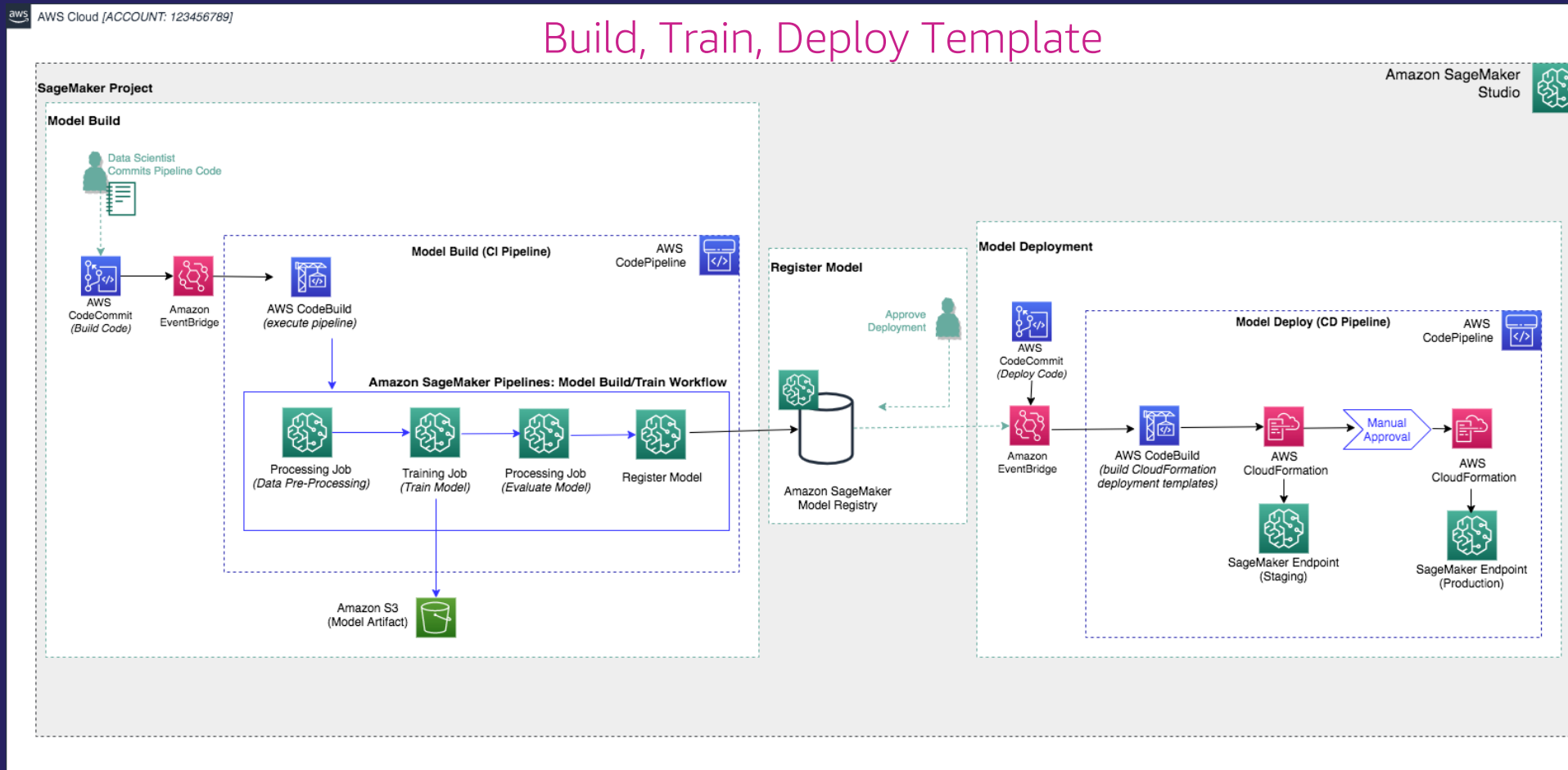
Model metrics		
Model metric	Metric value	Standard deviation
mse	4.823176167079859	2.1960237865043672

- Catalog models for production
- Manage model versions
- Associate metadata with a model
- Manage the approval status of a model
- Deploy models to production (with Projects)

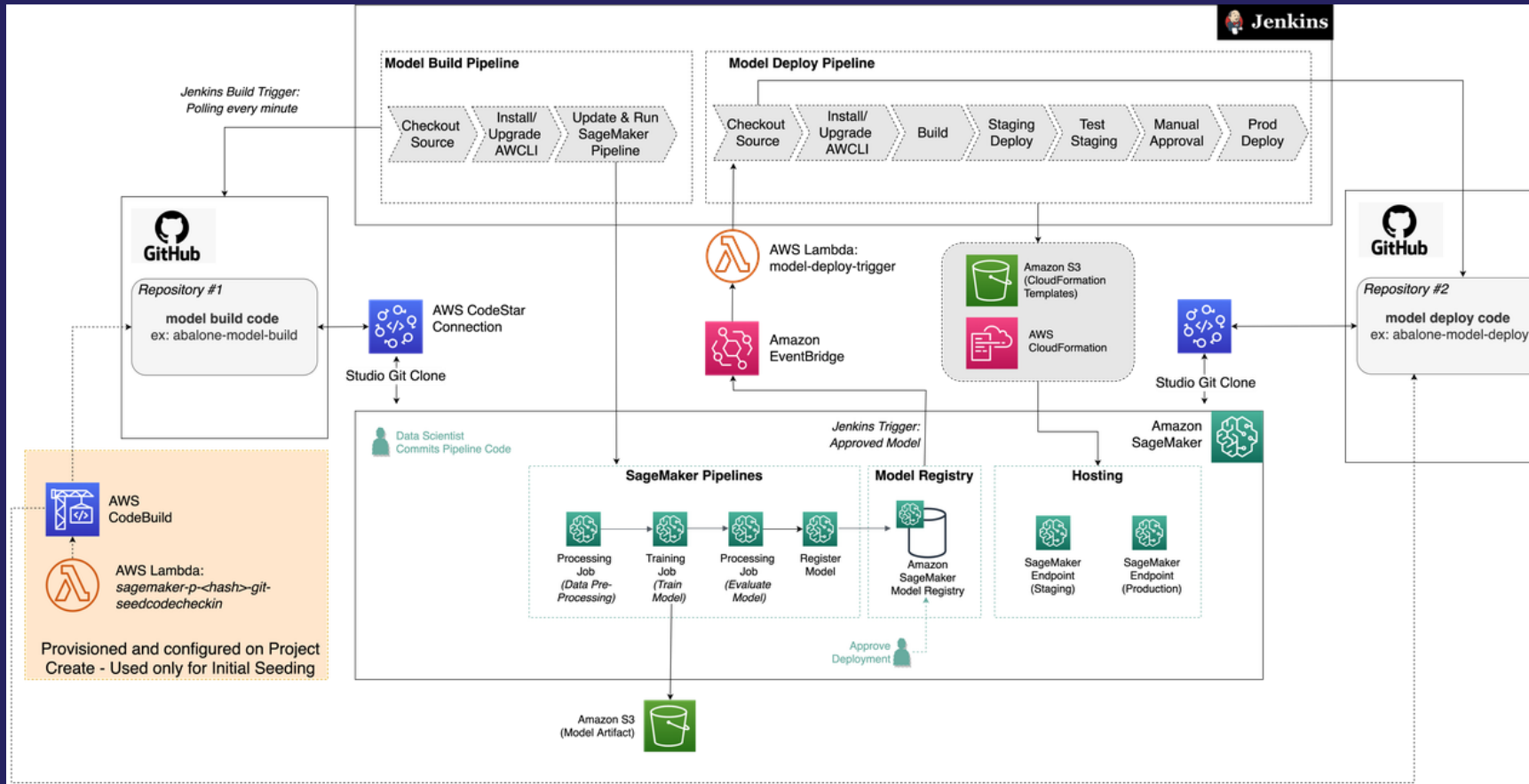
- Track model performance metrics

Amazon SageMaker Projects

High Level Services View



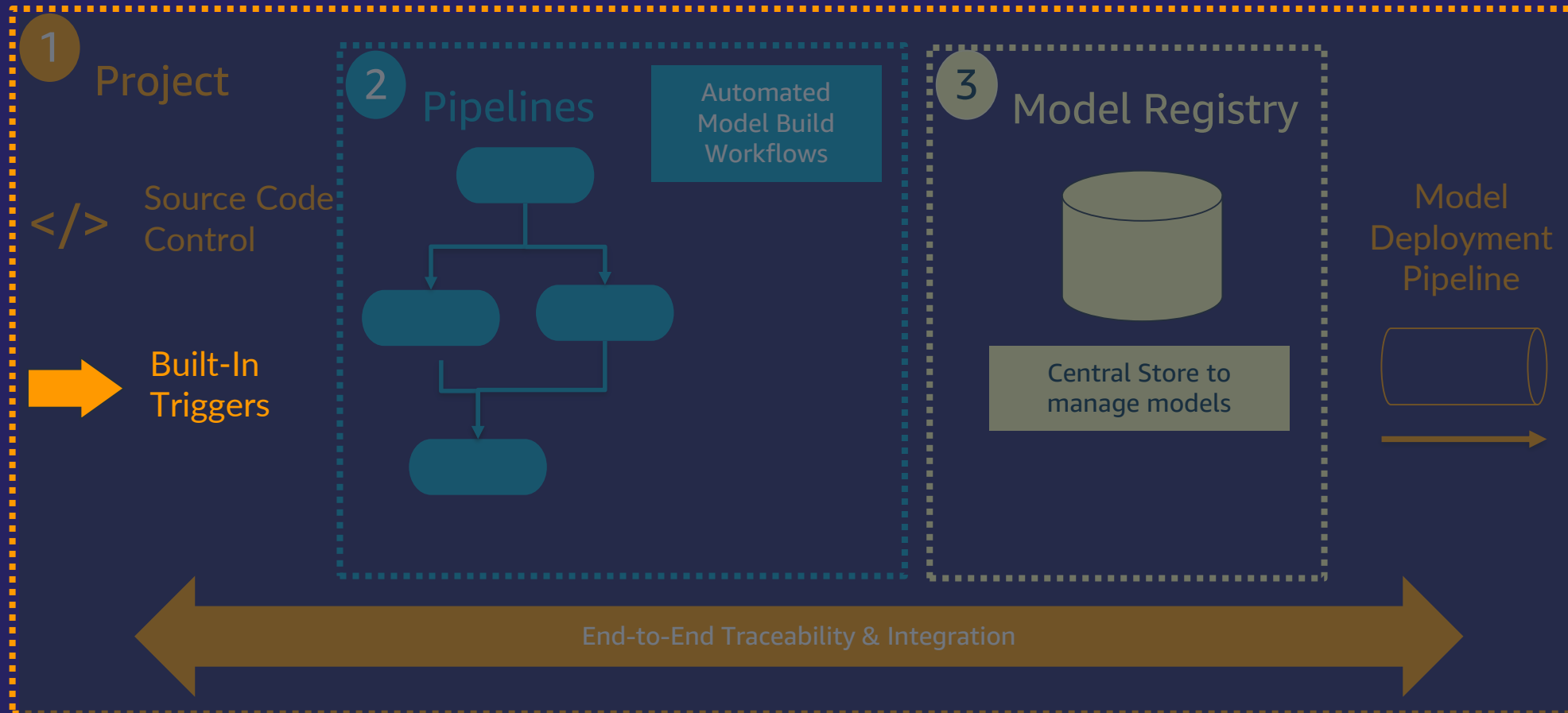
Amazon SageMaker Projects using third-party source control and Jenkins



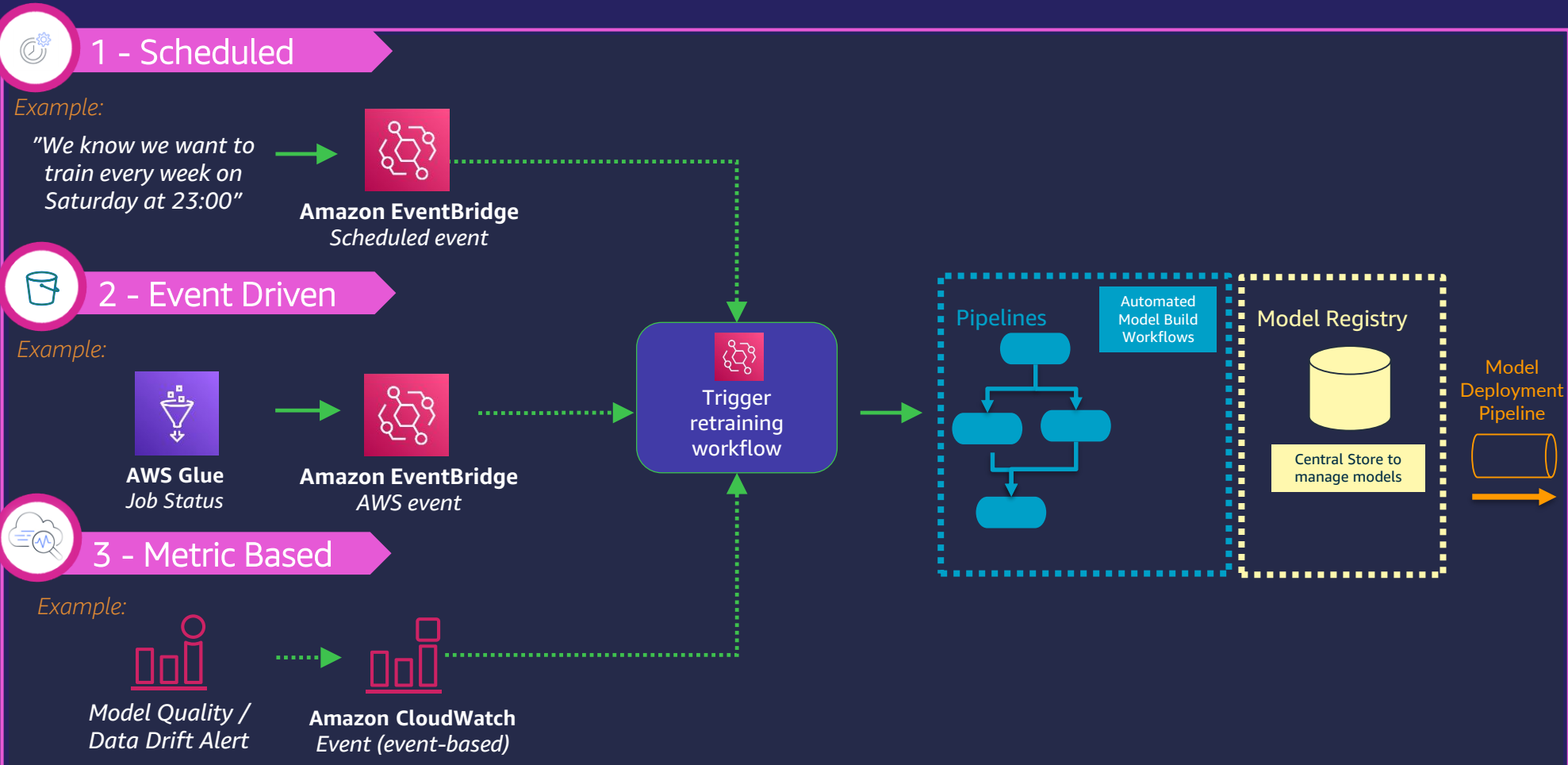
<https://aws.amazon.com/blogs/machine-learning/create-amazon-sagemaker-projects-using-third-party-source-control-and-jenkins/>

Amazon SageMaker Pipelines

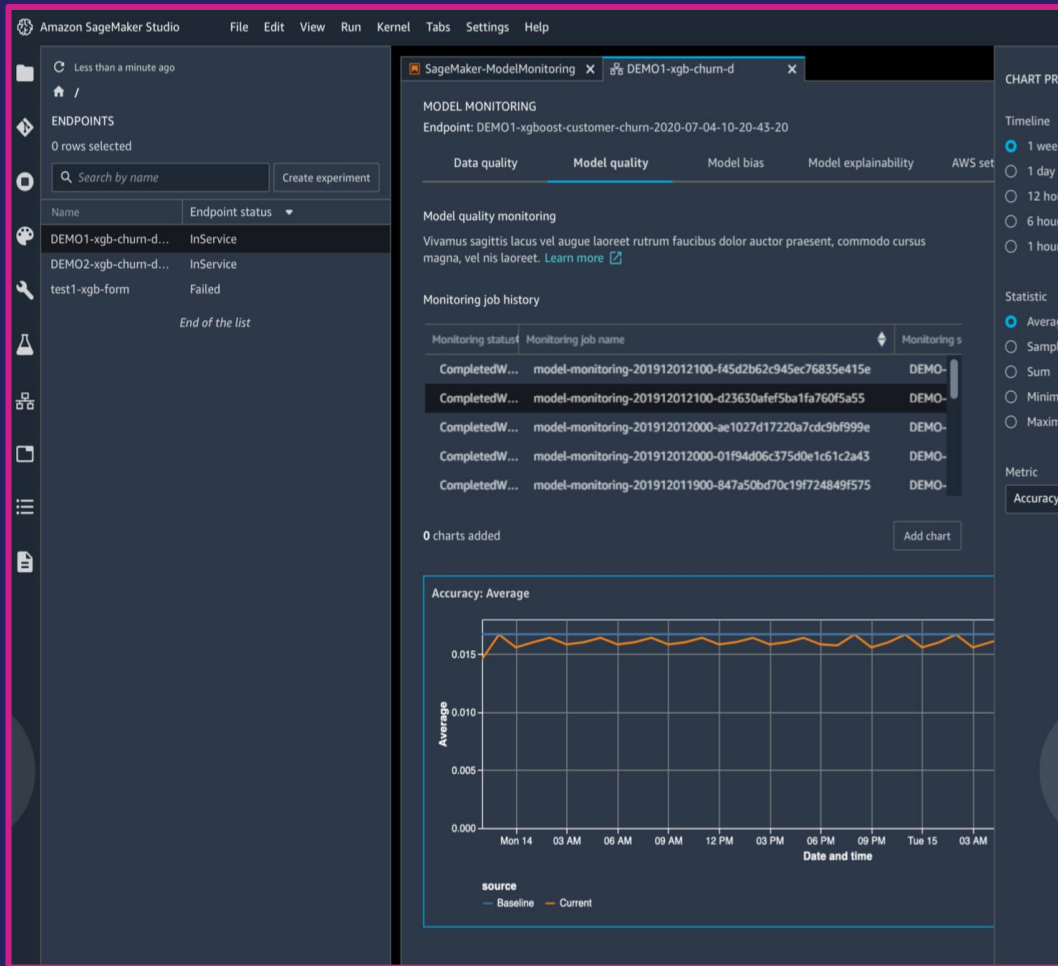
Built-In Triggers



Creating retraining strategies



Amazon SageMaker Model Monitor



Supported Features:

- Automatic data collection
- Continuous monitoring
- Flexible Monitoring Rules
- Visual data analysis
- CloudWatch integration

Amazon SageMaker MLOps

Streamline the ML lifecycle



Automate ML workflows to scale model development



Build CI/CD pipelines for ML to accelerate model deployment



Catalog model versions, metadata, metrics, and approvals for traceability and reusability



Track lineage for troubleshooting and compliance



Maintain accuracy of predictions after models are deployed



Enhance governance and security

Getting Started



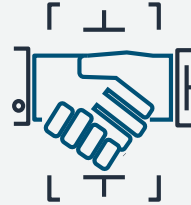
Getting started: Next steps



**Discovery and
Get Hands Dirty**



**Proof of
Concepts (PoC)**



**AWS Partner
Network (APN)**



**Training and
Certification**



**Thank you! Fill in the event survey
and get USD 25 AWS Credits**