



Architecting for ML, On AWS

Day 3

Phi Nguyen
Machine Learning Solutions Architect

Agenda

Day 1

AI/ML on AWS
Intro lab

Team up
Define problem

Write-up

Day 2

Feature engineering
Model evaluation

Build

Working model

Day 3

Moving to
production
Build

Present

Solution
architecture

Advanced Sagemaker Topics

Architecture Patterns, Security and DevOps

Survey link

<http://bit.ly/2M3QF8j>

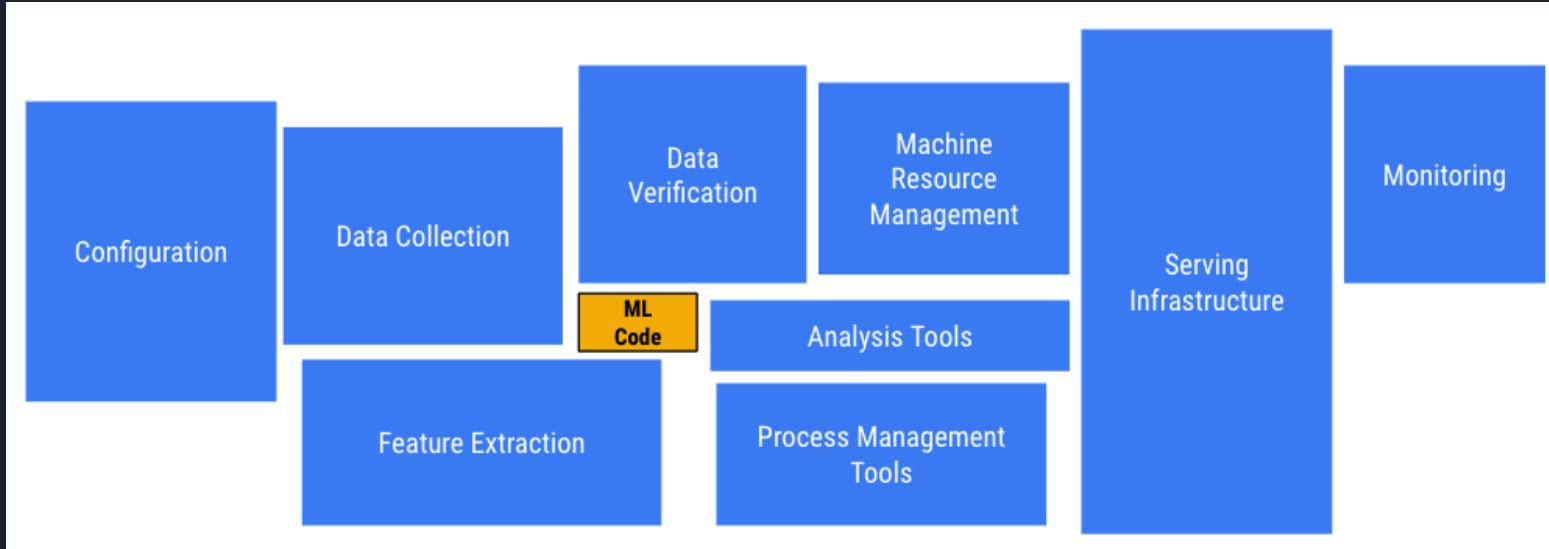
Presentations PDF

<https://bit.ly/2kRXfVZ>

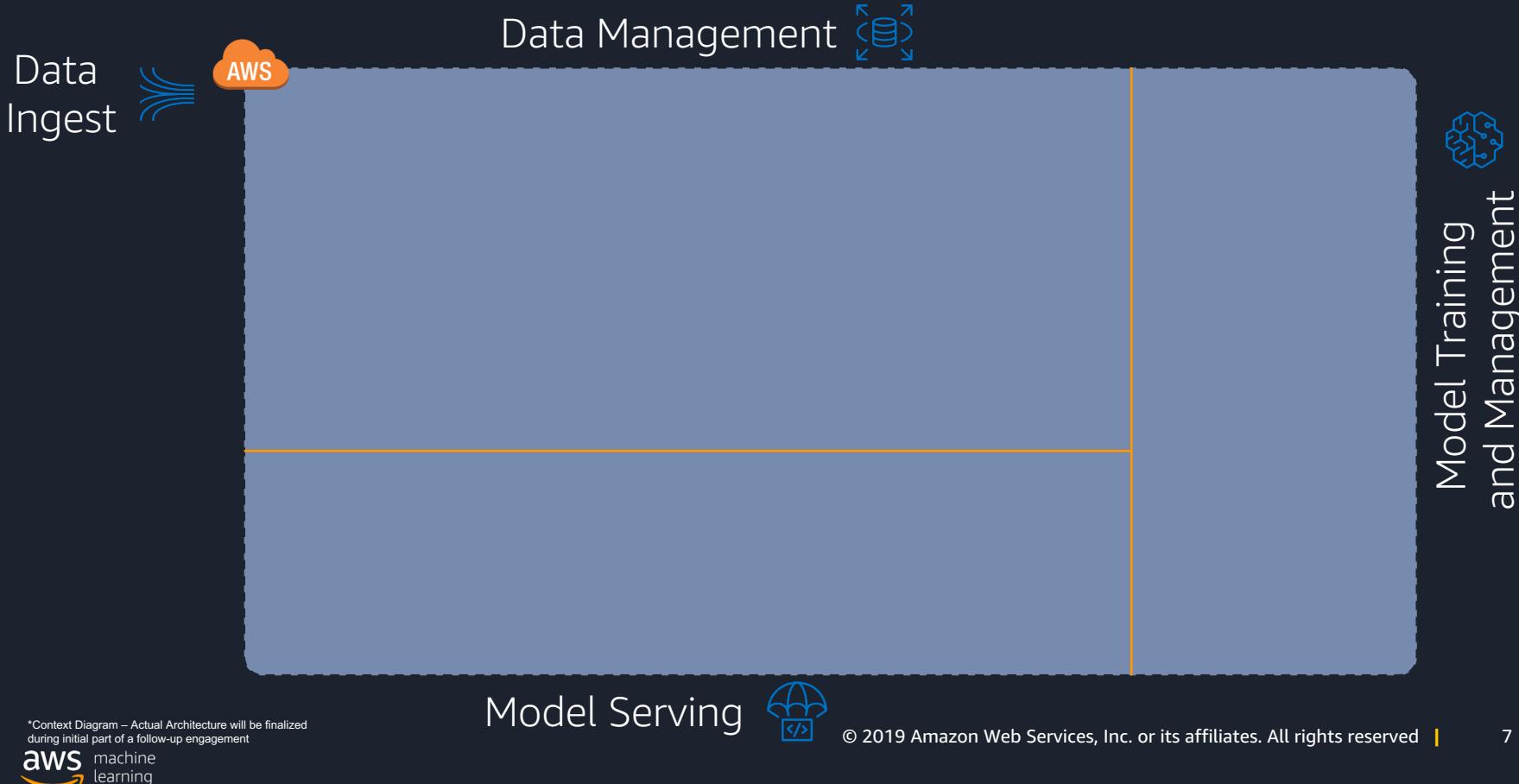
Architecture Patterns

From MVP to production

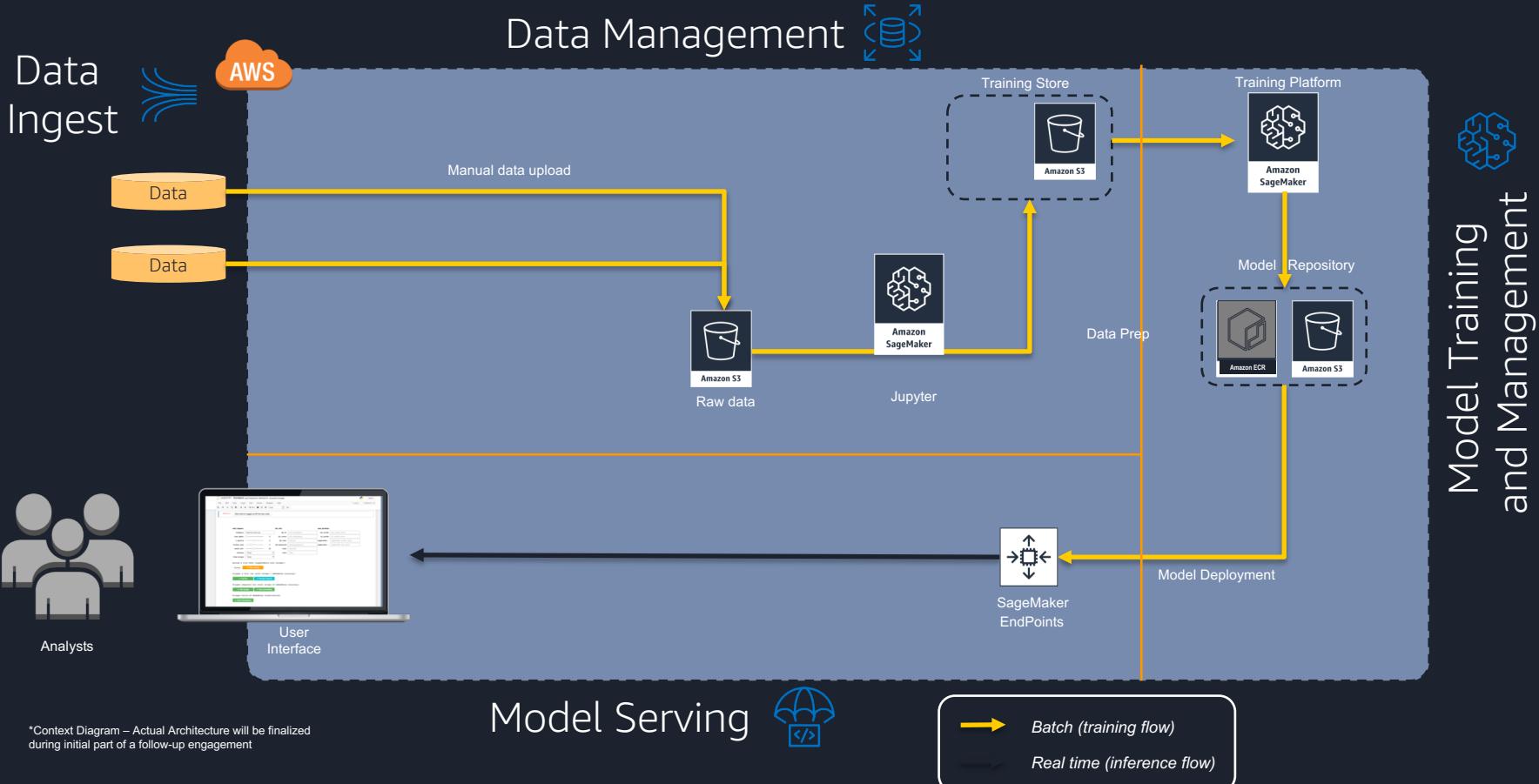
Hidden technical debt of Machine Learning



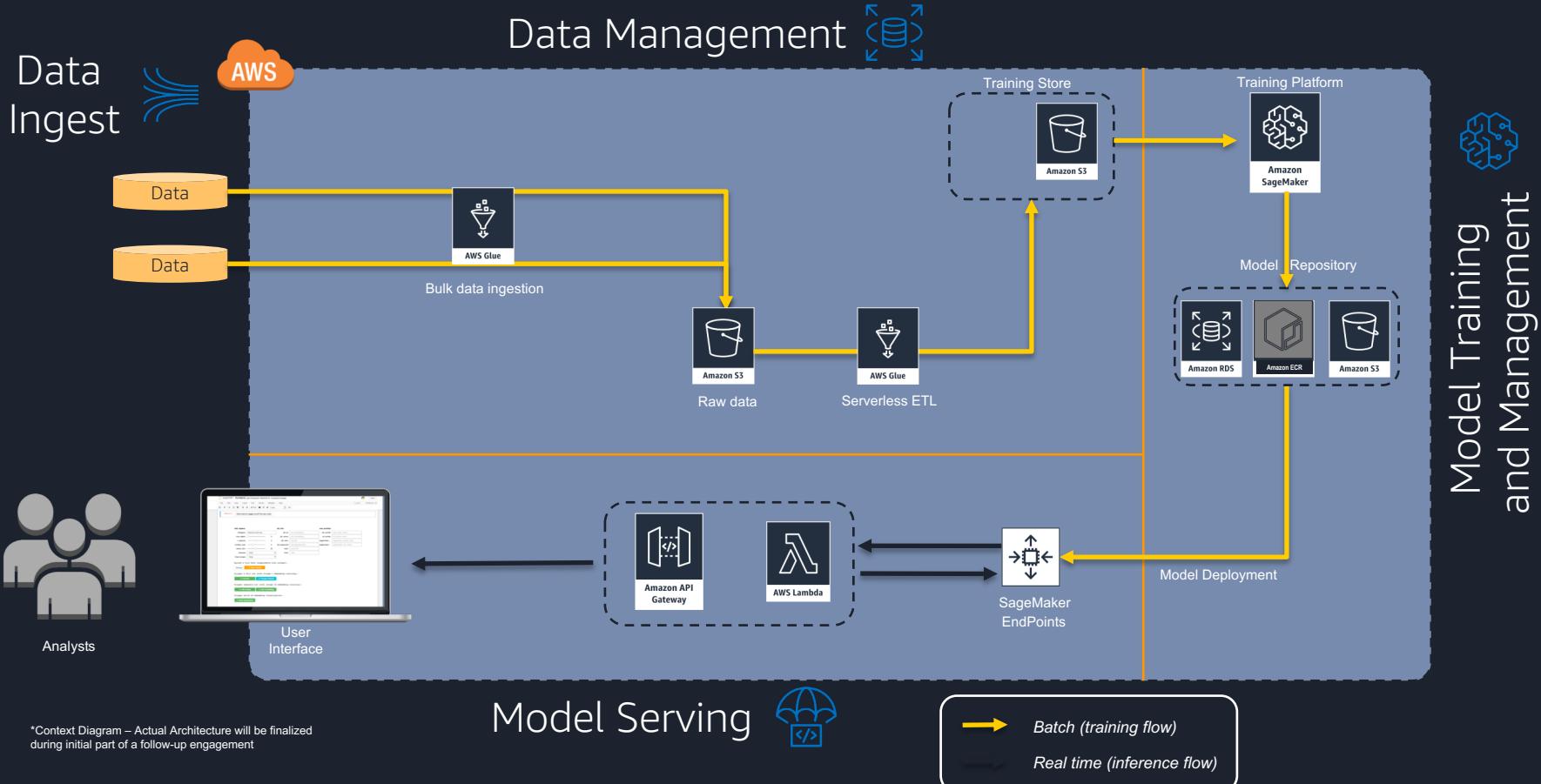
Architecture building blocks



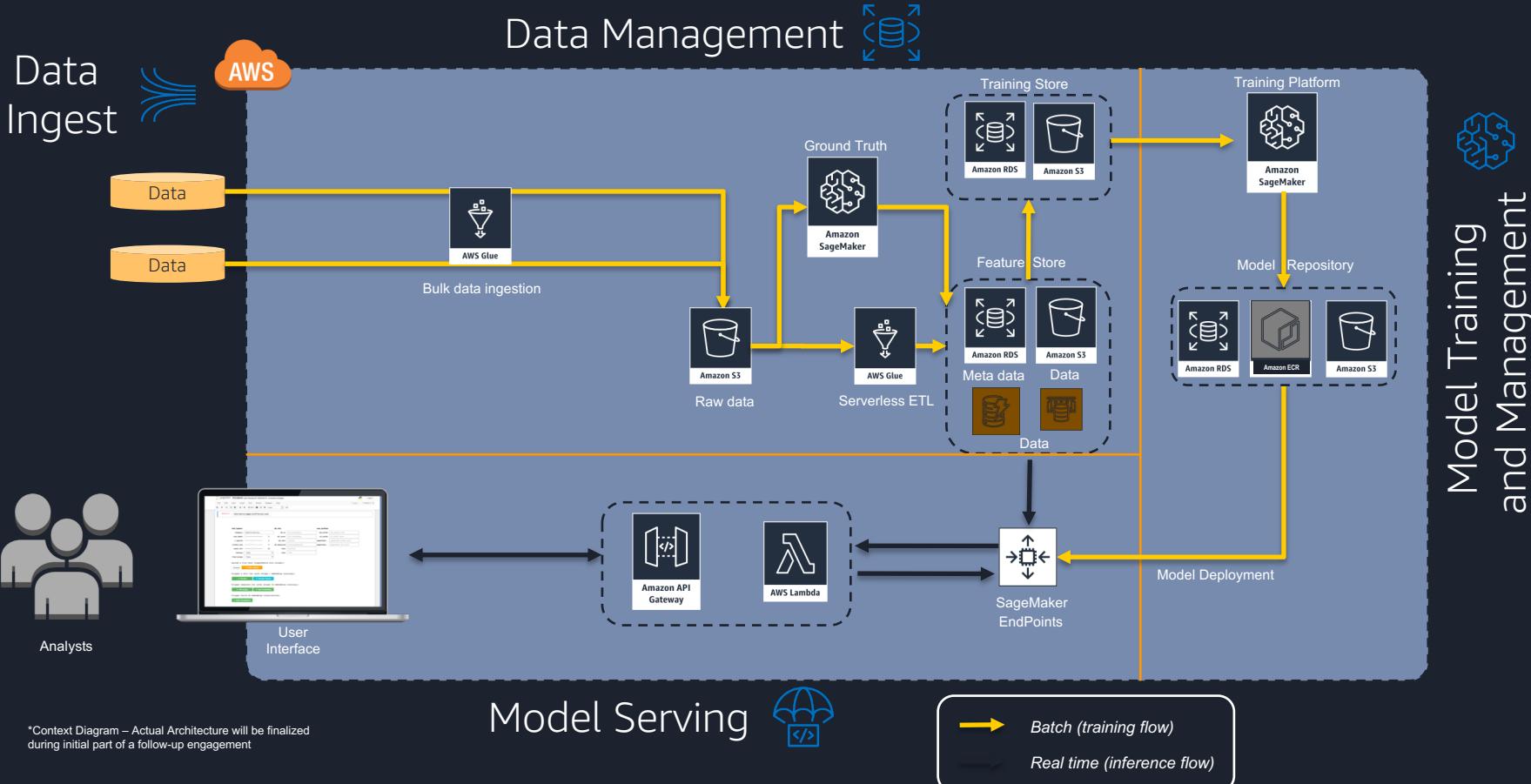
POC/MVP



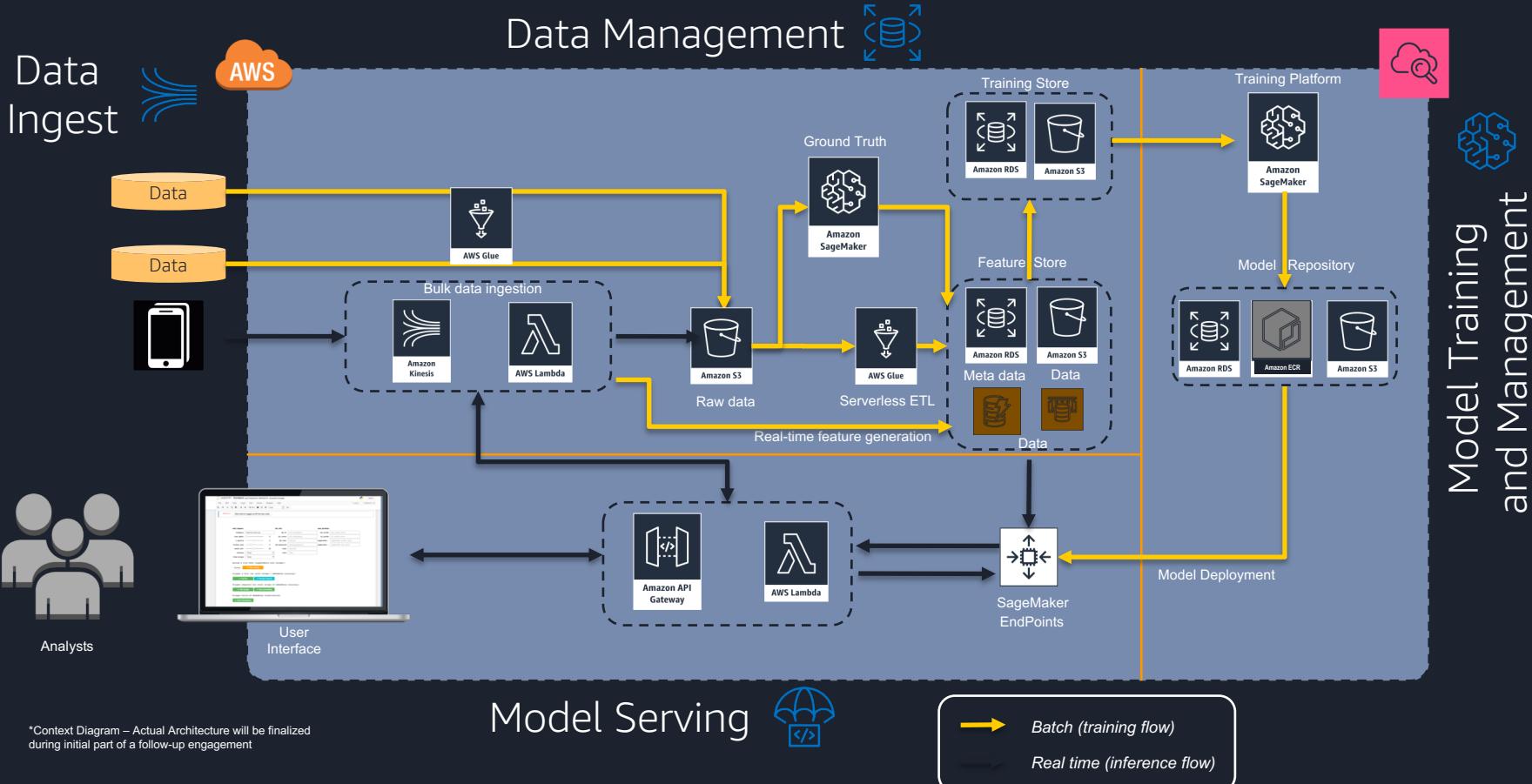
Model management enhancement



Feature and training store management

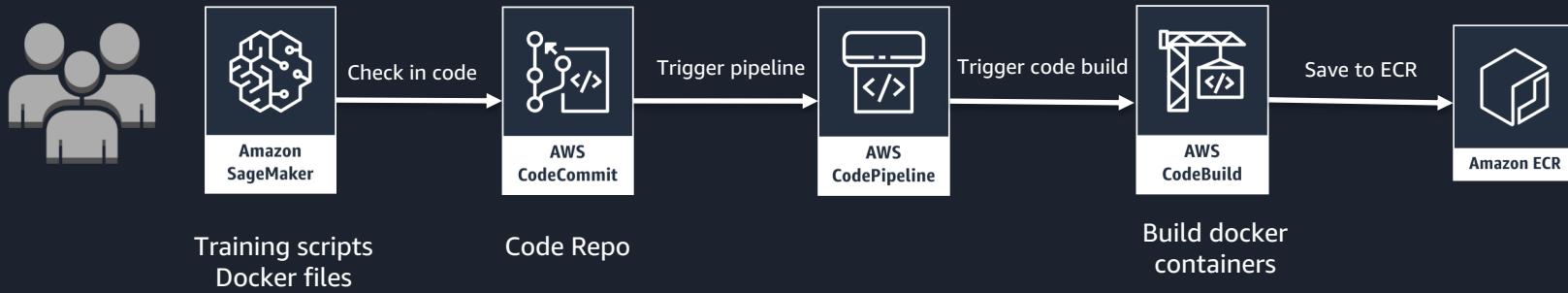


Realtime capability and Monitoring



MLOps

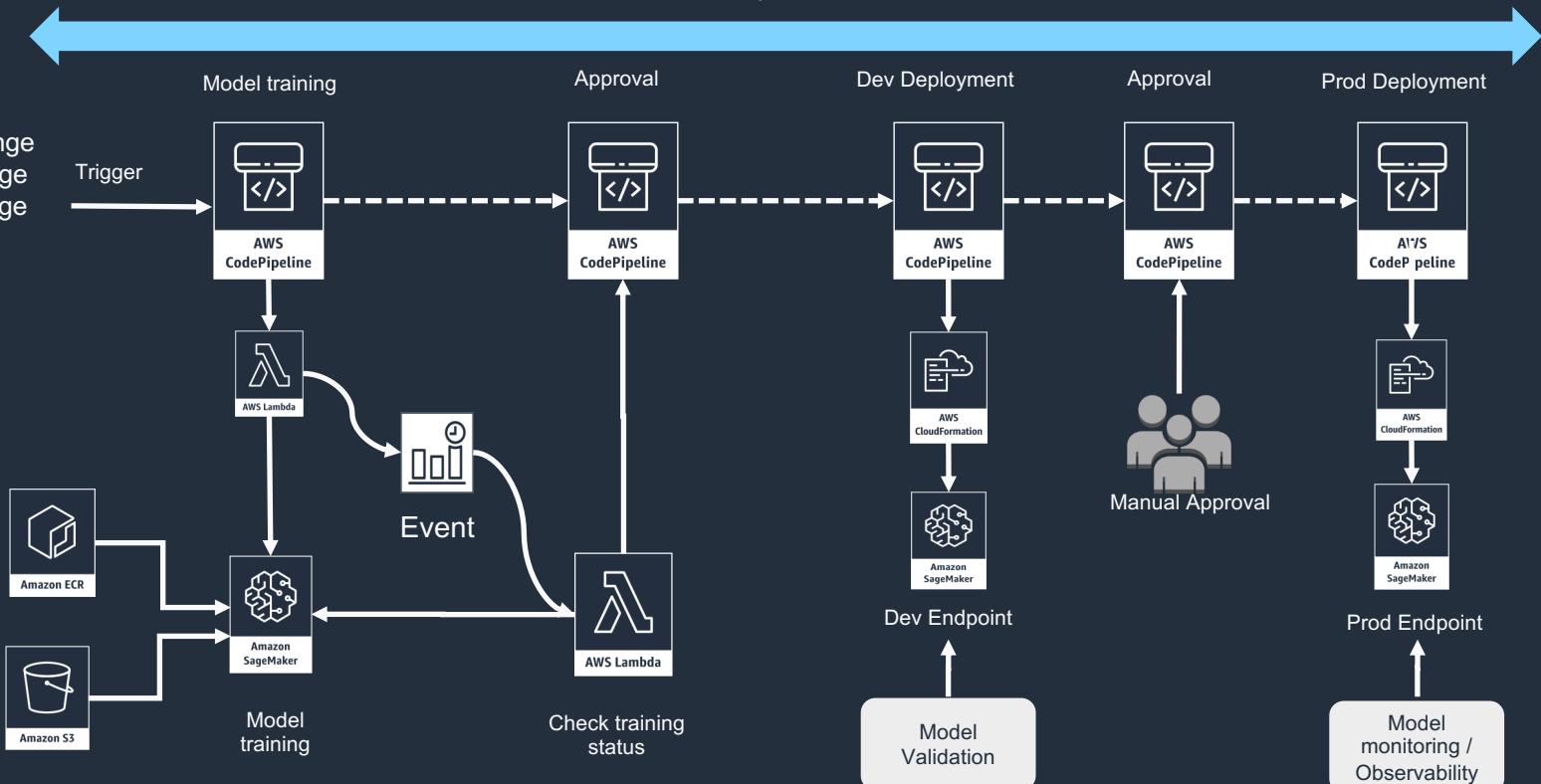
ML DevOps – Custom Docker Image Build Pipeline



Sample Training/Deploy Pipeline

AWS CodePipeline Automation

- Code change
- ECR change
- Data change
- Schedule
- Manual



How do I orchestrate my ML pipeline?

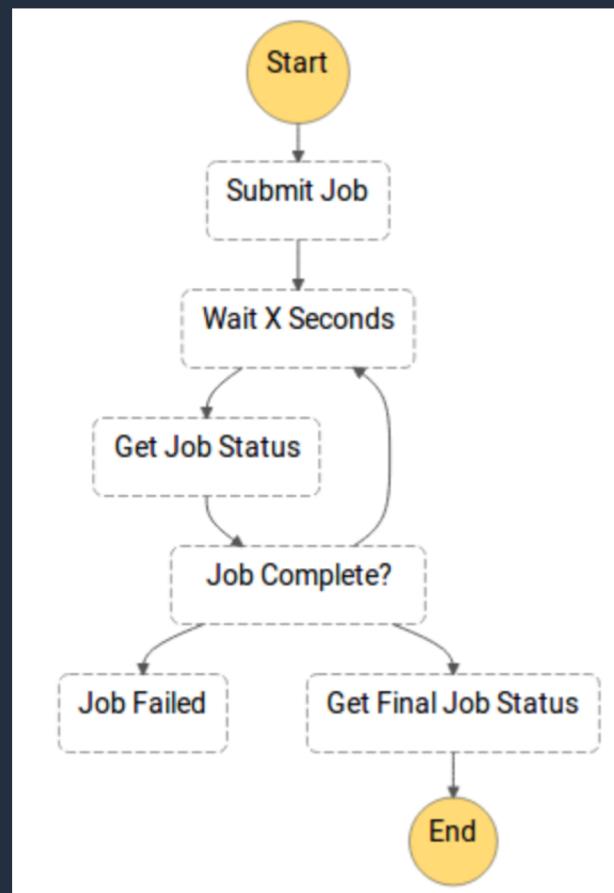
- Step functions
- Cloud Development Kit (CDK)
- Apache Airflow
- Kubeflow

Step Functions

Supported Amazon SageMaker APIs :

[CreateTrainingJob](#)

[CreateTransformJob](#)



CDK example

API gateway > Lambda > SM



```
lambdaFn = lambda_.Function(  
    self,  
    "callsmlambda",  
    code=lambda_.InlineCode(handler_code),  
    handler="index.lambda_handler",  
    timeout=core.Duration.seconds(300),  
    runtime=lambda_.Runtime.PYTHON_3_7,  
    environment={"endpoint_name":endpoint_name, # CHANGE TO YOUR ENDPOINT NAME!!  
                 "content_type":"text/csv"}  
)  
  
lambdaFn.add_to_role_policy(aws_iam.PolicyStatement(actions=['sagemaker:InvokeEndpoint'],  
                                                 resources = ['arn:aws:sagemaker:{}:{}:endpoint/{}'.format(my_region,my_acc_id,endpoint_name),]))
```

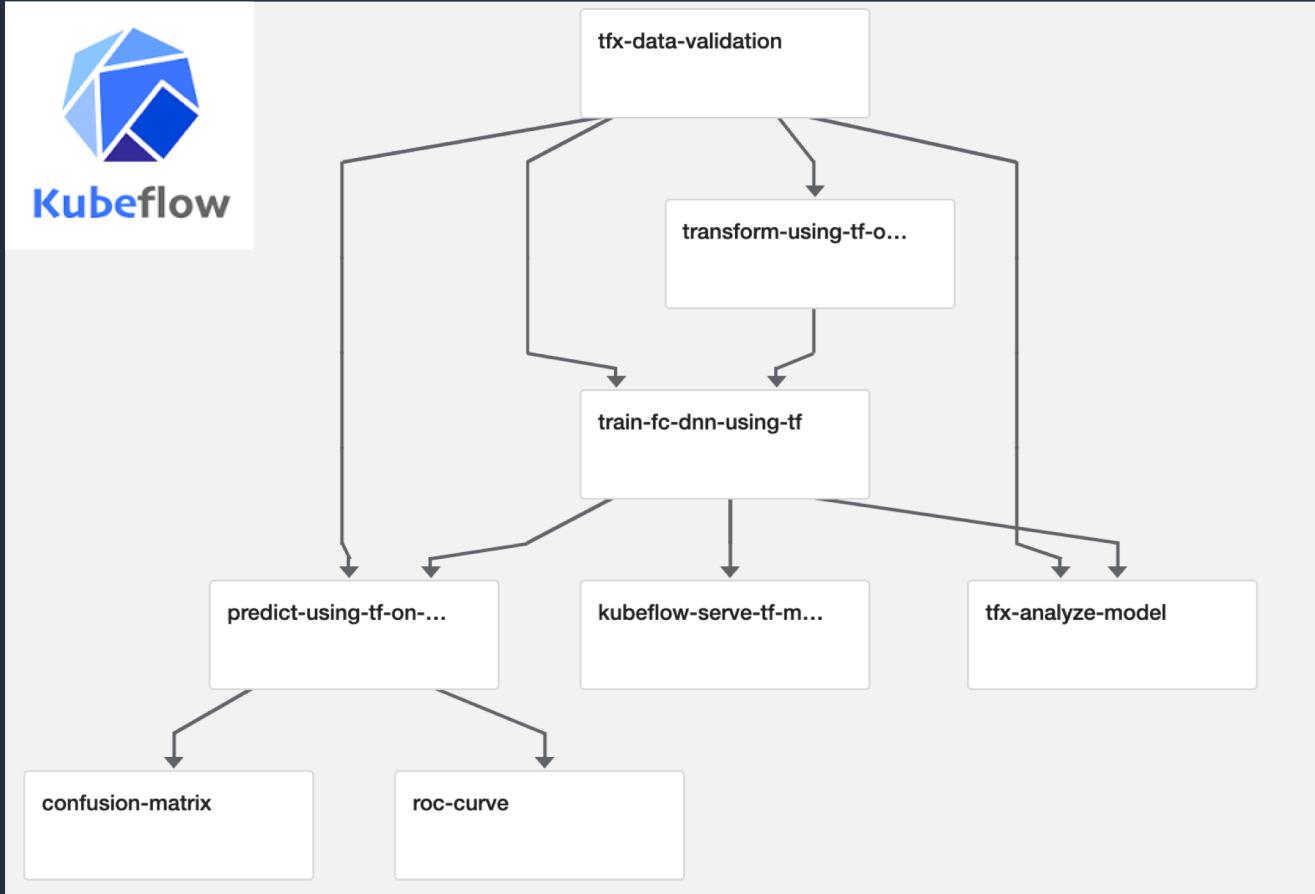
Airflow with SageMaker on AWS

Currently, the following SageMaker operators are supported:

- SageMakerTrainingOperator
- SageMakerTuningOperator
- SageMakerModelOperator
- SageMakerTransformOperator
- SageMakerEndpointConfigOperator
- SageMakerEndpointOperator

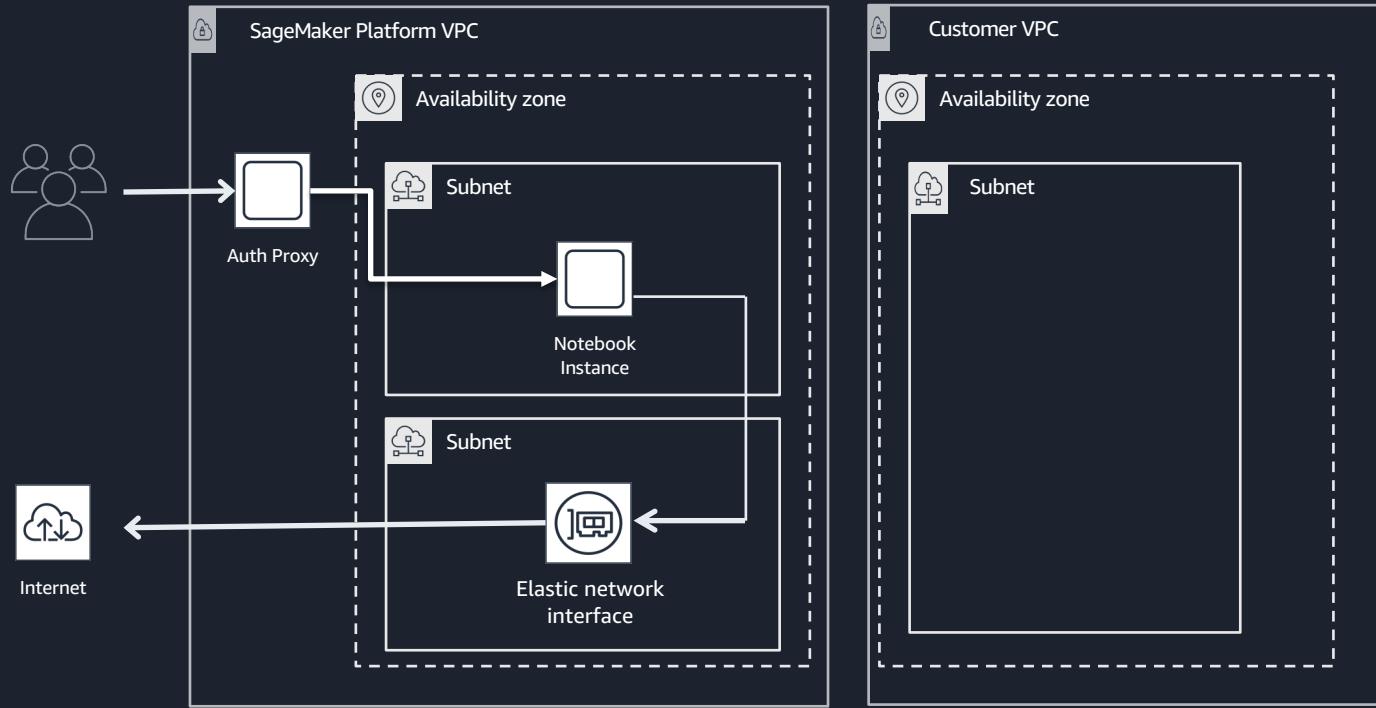
https://sagemaker.readthedocs.io/en/stable/using_workflow.html

Taxi Tip Prediction Model Trainer using Tensorflow DNN

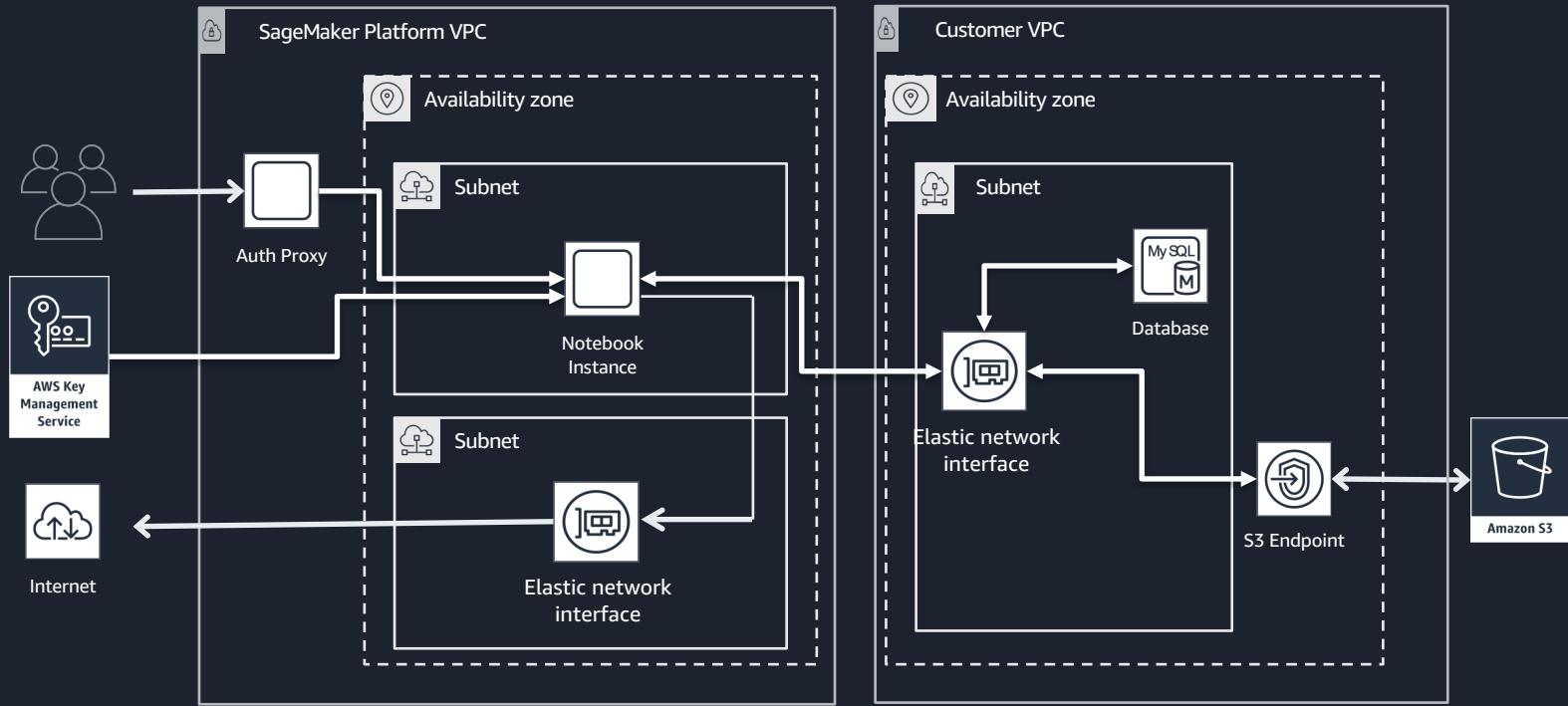


Notebook Security

Jupyter notebooks - Default



Secure Jupyter notebooks + VPC



Deployment Best Practices

Deployment best practices

- Deploy to at least 2 availability zones, multiple regions for DR
- Use auto-scaling to accommodate traffic spikes
- Consider Lambda and API Gateway
- Consider batch inference or edge deployment
- Use Elastic Inference Accelerators to reduce inference costs
- Automate training and deployment pipeline
- Use blue/green deployment techniques
- Clearly define re-training triggers
- Enable model versioning, implement MLDevOps pipeline
- Monitor performance

Project Presentations

- Project goal, business context
- ETL strategy
- Modeling strategy
- Solution architecture
- Goal for going into production
- 20 minutes to present, can use slides