



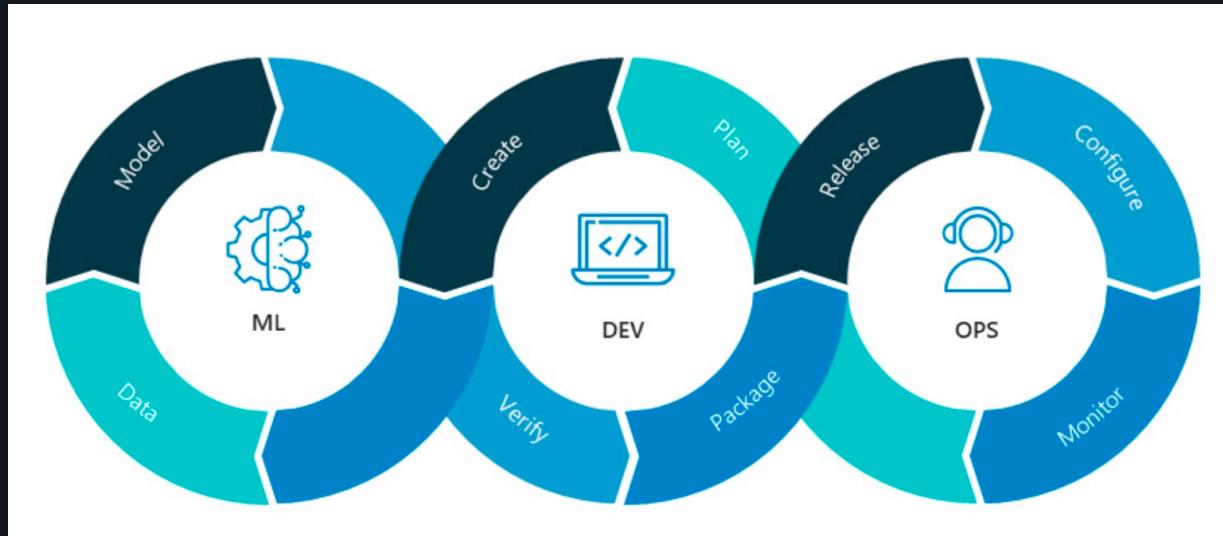
MLOps Model Deployment & Infrastructure

Sergey Pugachev
Solutions Architect, AWS



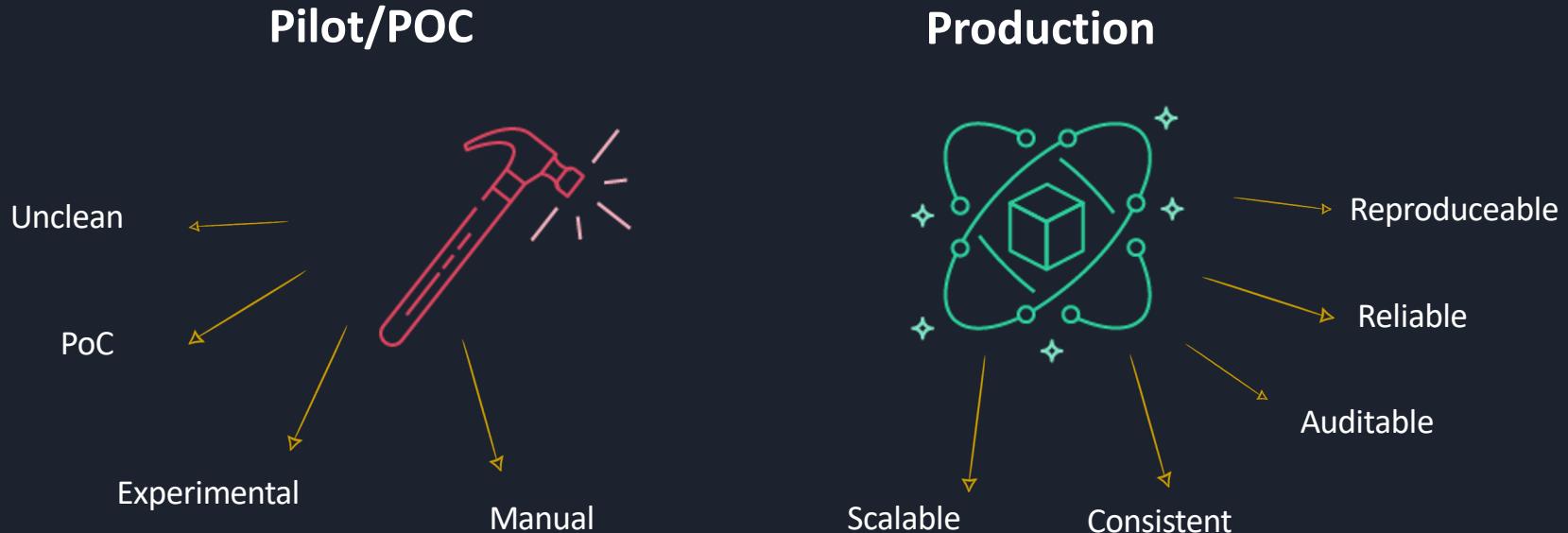
What is MLOps?

MLOps, a close relative to DevOps, is a combination of philosophies and practices designed to enable data science and IT teams to rapidly develop, deploy, maintain, and scale out Machine Learning models.



<https://blogs.nvidia.com/blog/2020/09/03/what-is-mlops/>

AI/ML only brings value in production



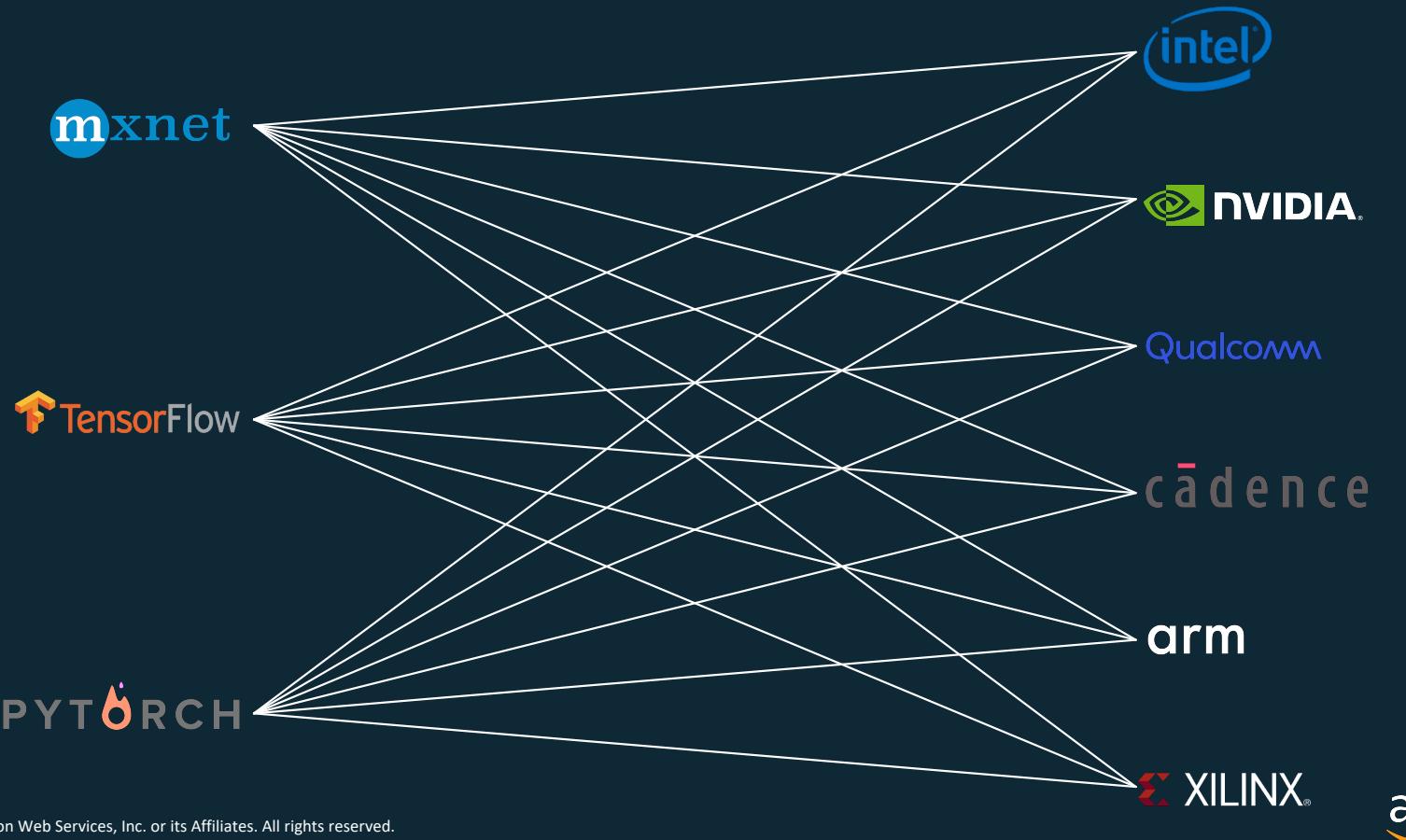
Agenda

1. What to run?
2. Where to run?
3. How to run?

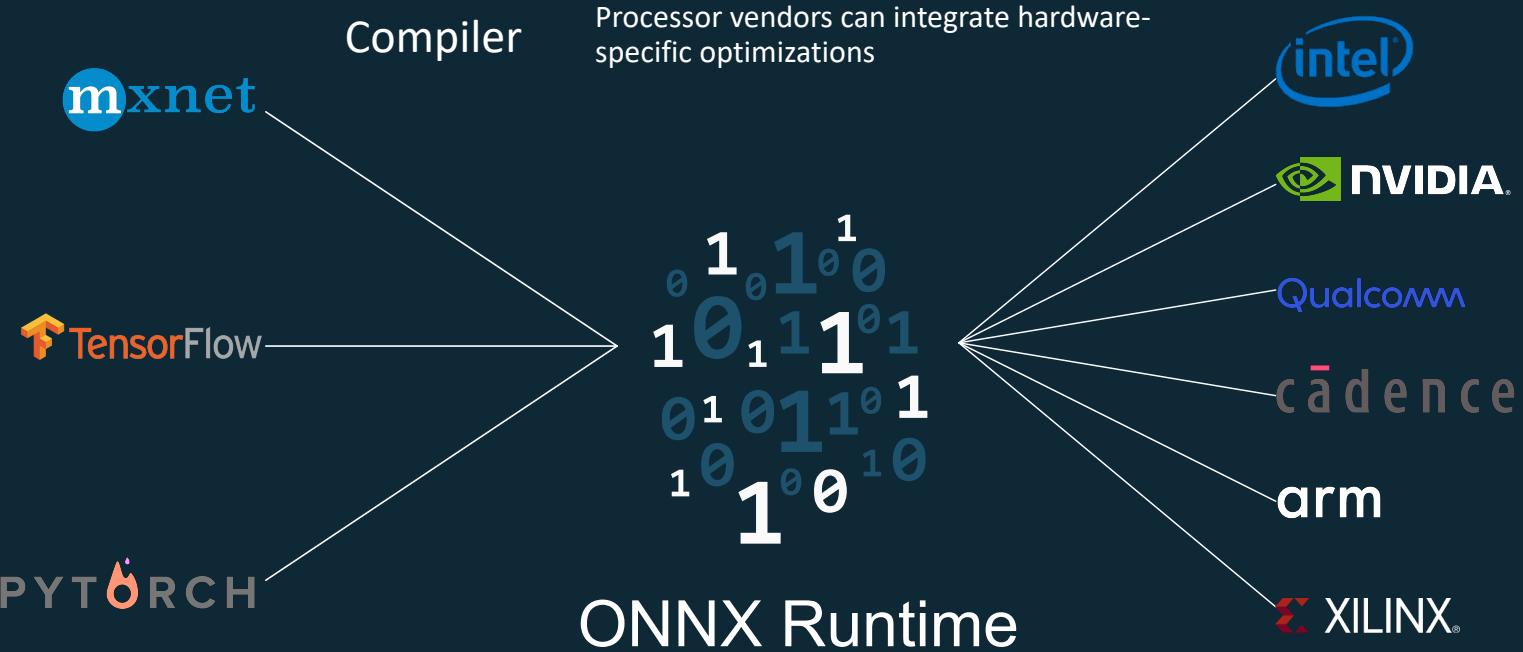
1. What?

Save to file: pickle and unpickle

Problem: Framework optimization is complicated



ONNX runtime: Train once, run anywhere



<https://onnx.ai/supported-tools.html>

2. Where?

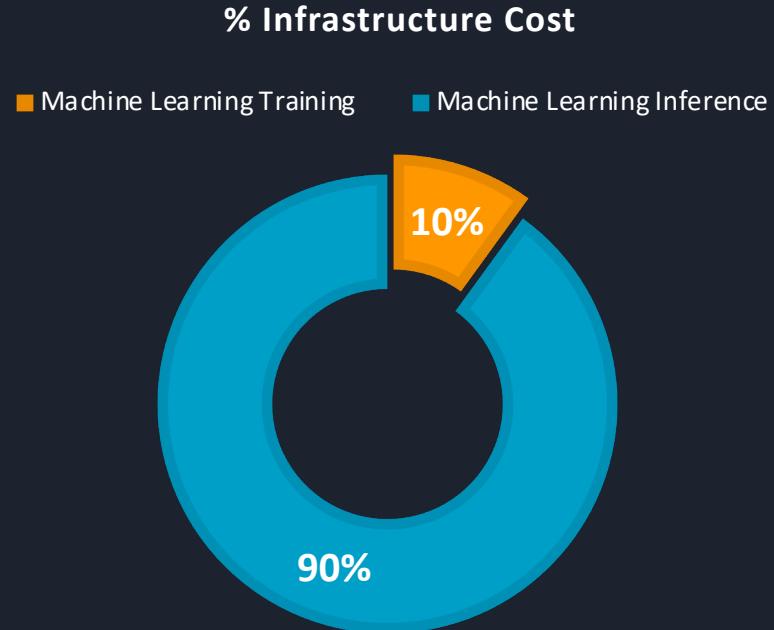
Machine Learning Inference Dominates ML Infra. Costs

1

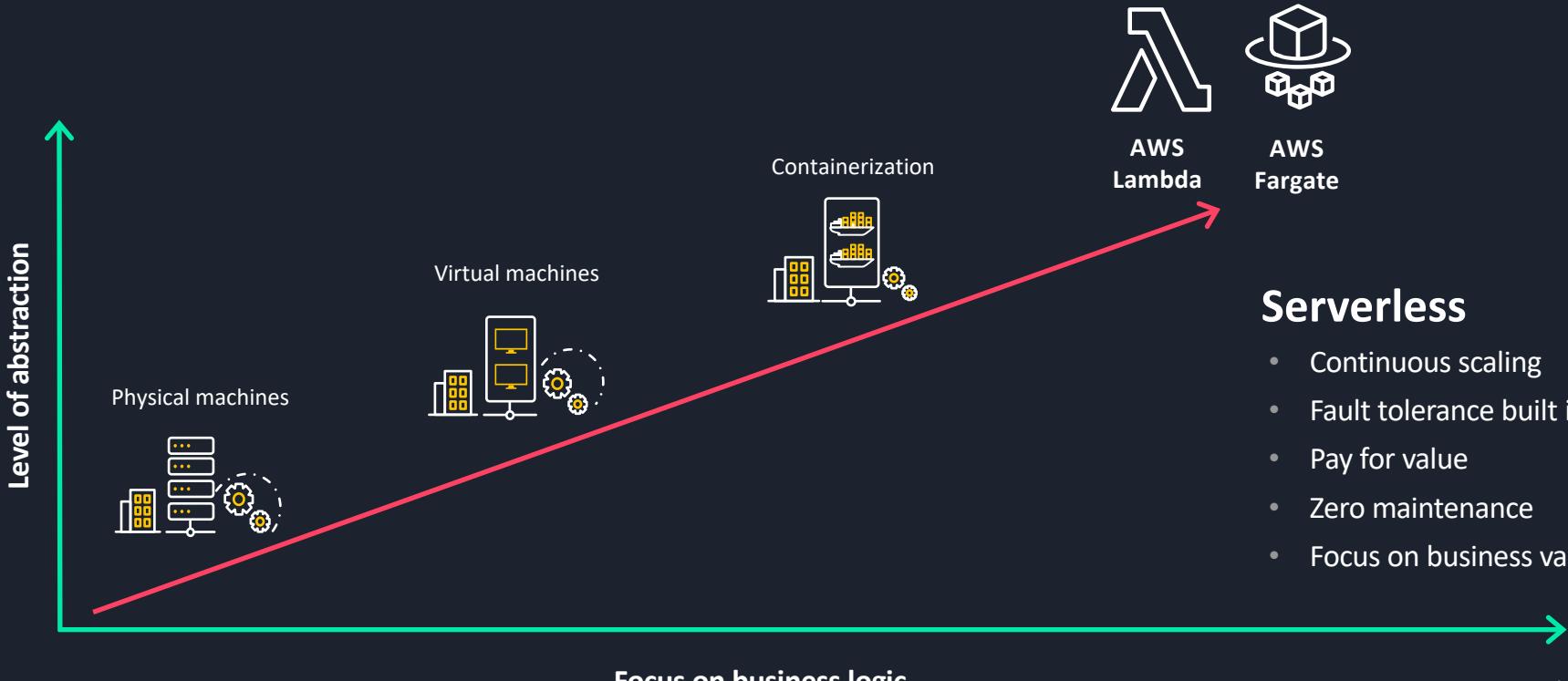
Cost to run ML models in production dominates the overall infrastructure costs to develop and run ML applications

2

ML practitioners are creating increasing complex models that are prohibitively more expensive to run in production



We are witnessing a paradigm shift



Serverless means...



No servers to provision
or manage



Scales with usage

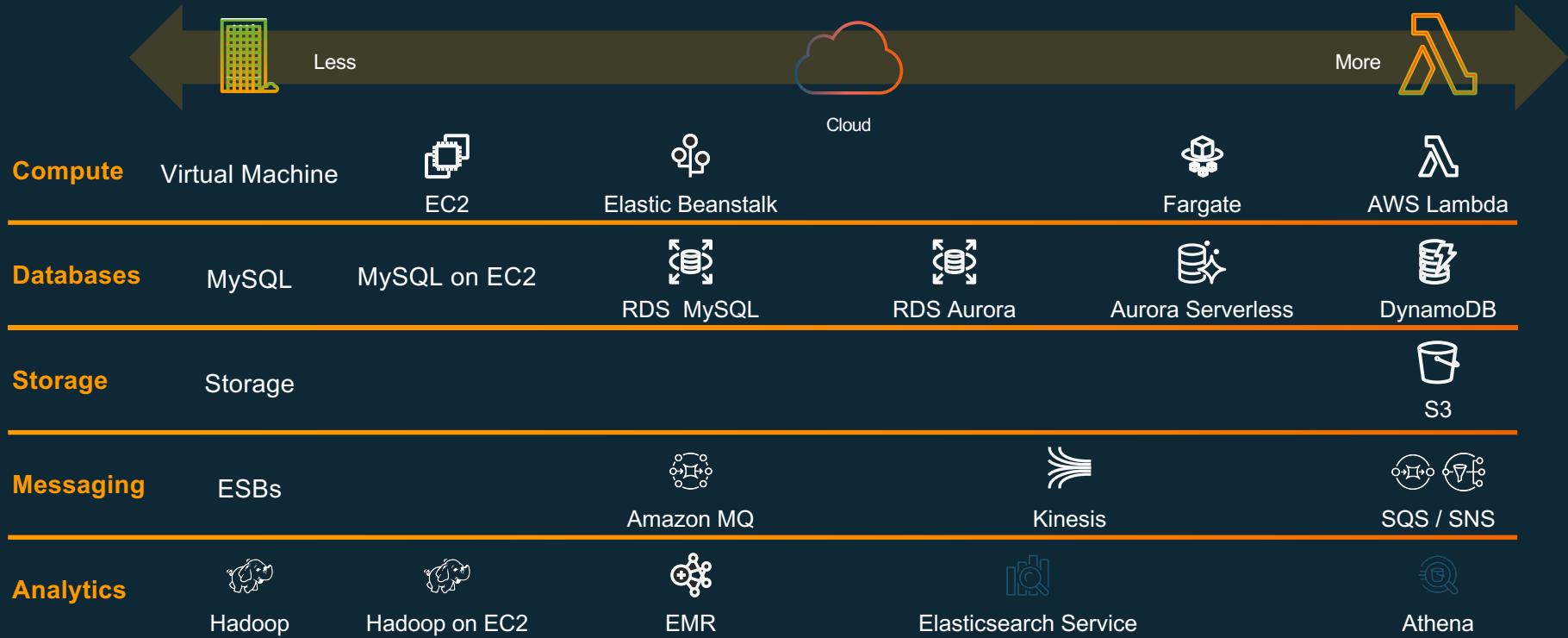


Never pay for idle



Availability and fault
tolerance built in

Cloud Services



Optimizing ML performance

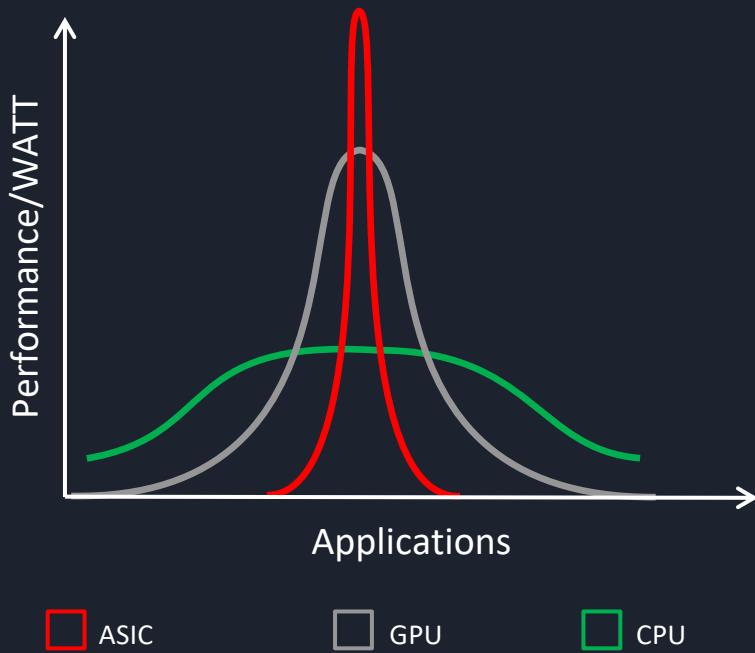
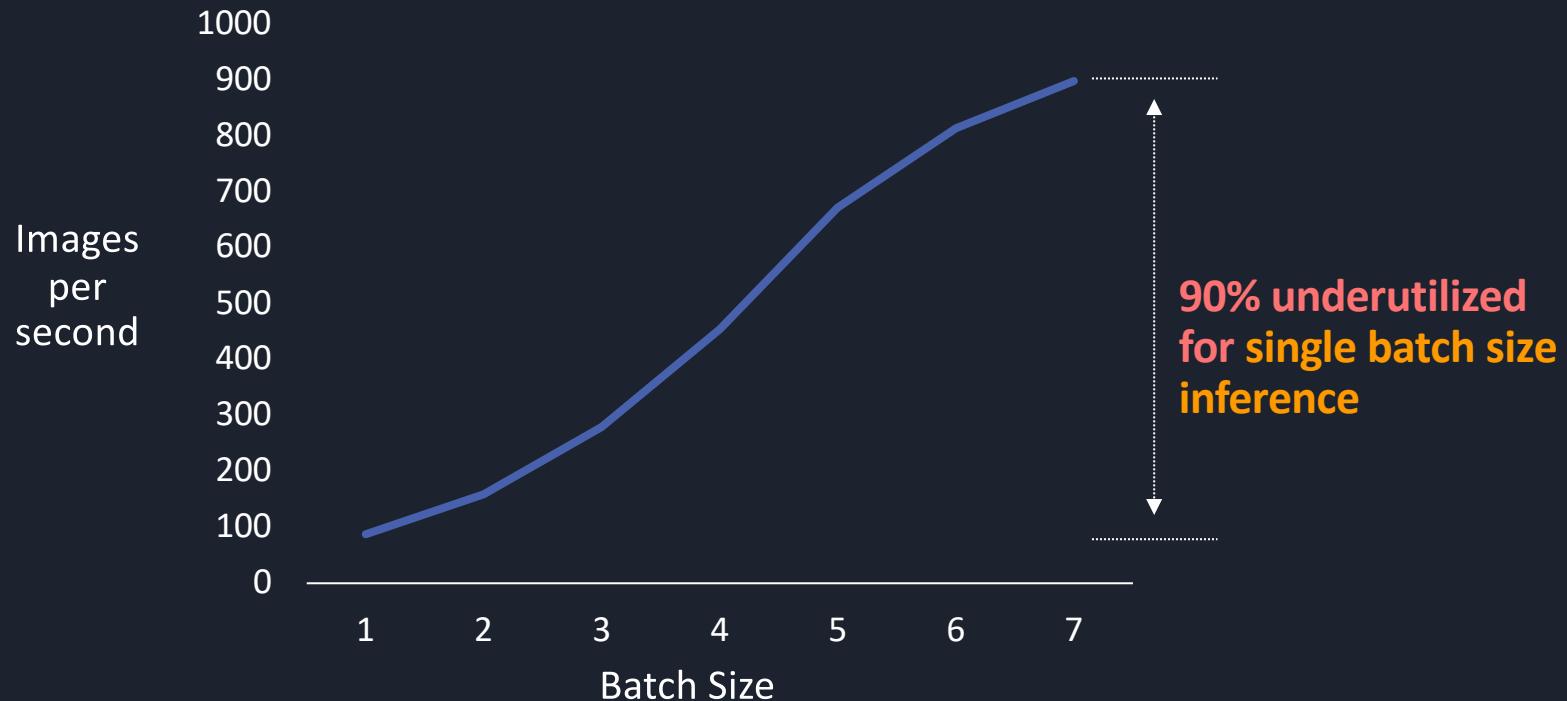


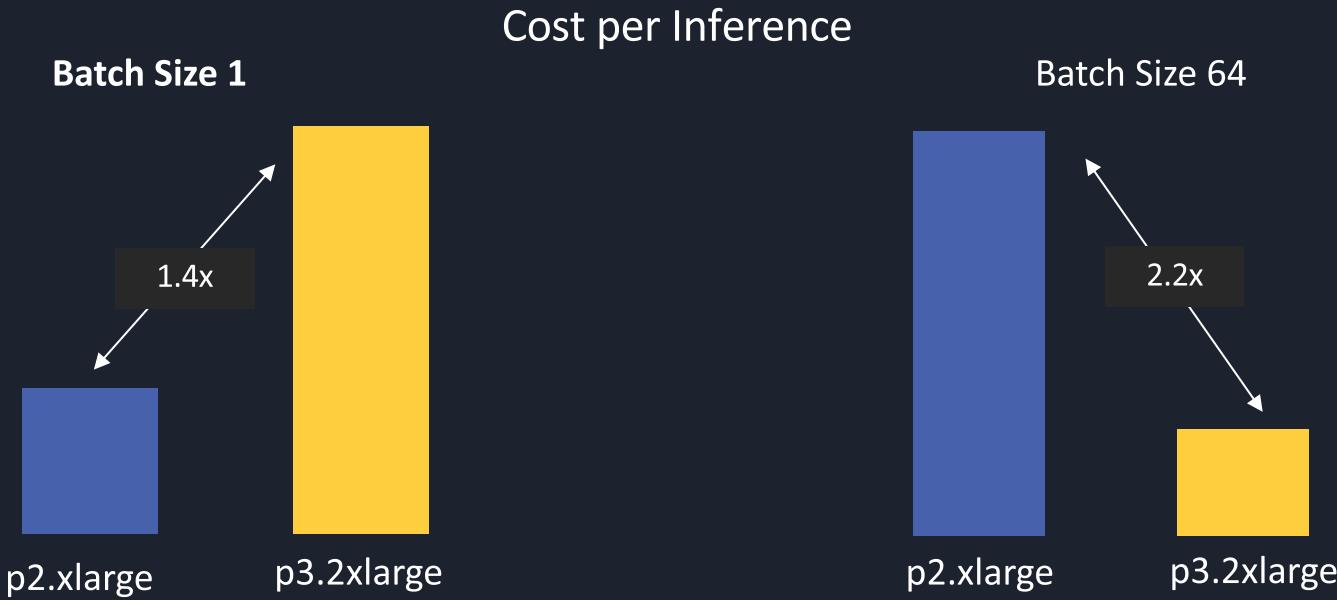
Chart for illustrative purposes

GPU utilization for inference: latency and throughput



How cost effective are GPU instances for inference?

Smaller P2 instances are more effective for real time inference with small batch sizes



Lower numbers are better

Inception-v3 FP32 MXNet

Processors



Intel® Xeon Scalable
processors



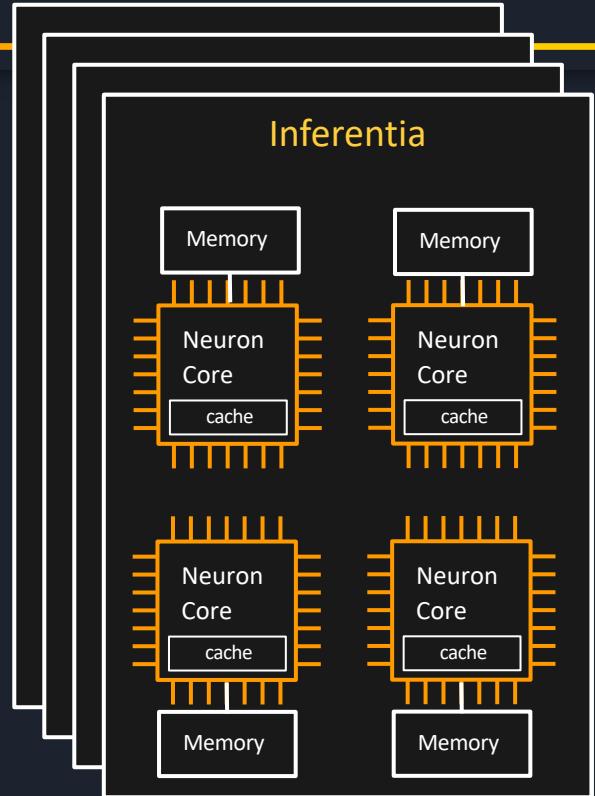
AMD EPYC
processors



AWS Graviton
processors

AWS Inferentia ASIC

- 4 Neuron Cores
- Up to 128 TOPS
- 2-stage memory hierarchy
 - Large on-chip cache and commodity DRAM
- Supports FP16, BF16, INT8 data types
- Fast chip-to-chip interconnect

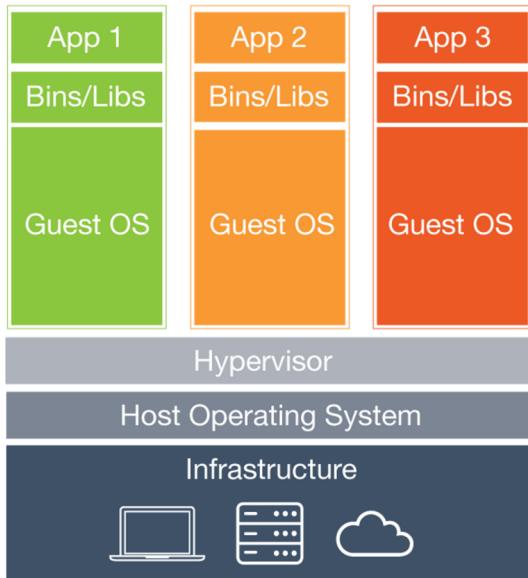




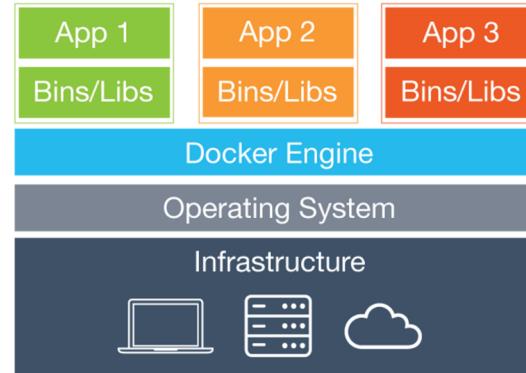
3. How?

Containers & Docker

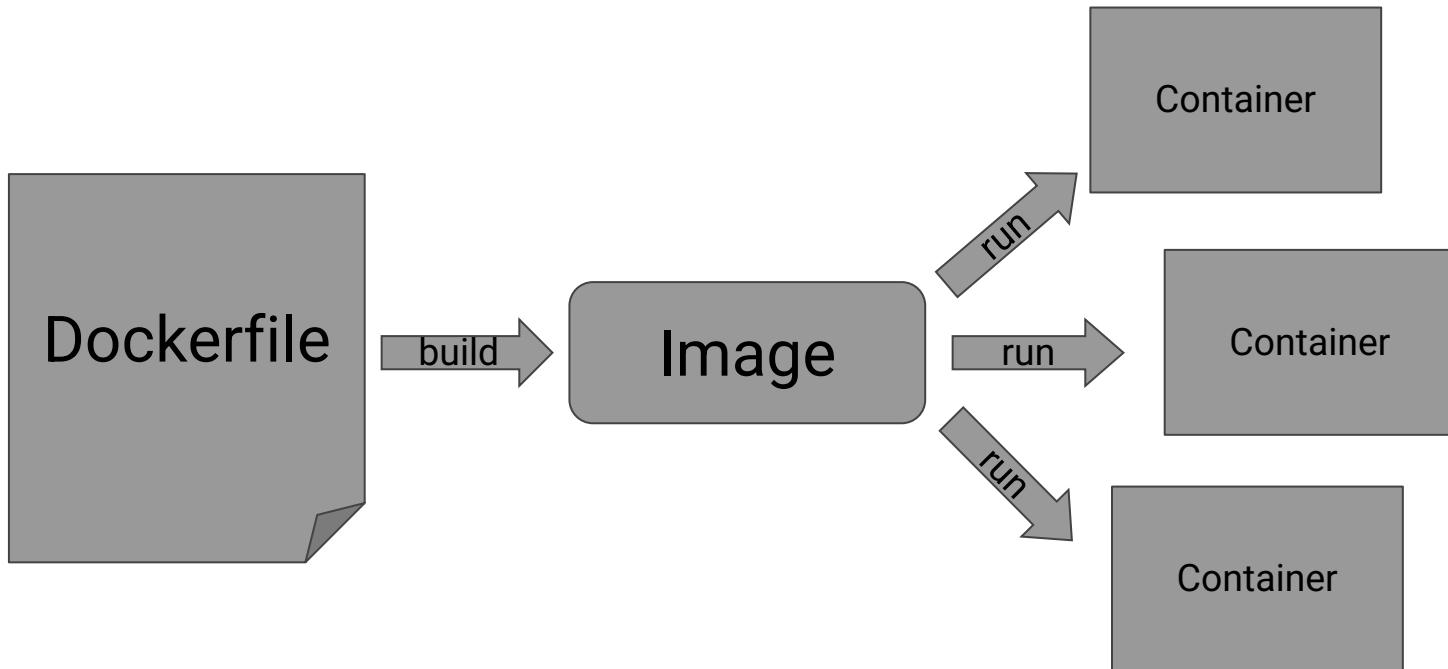
A container is a standard **unit of software** that packages up **code and all its dependencies** so the application **runs quickly and reliably** from one computing environment to another.



Virtual Machines



Containers



Infrastructure as Code / Terraform

```
resource "aws_apprunner_service" "example" {
    service_name = "example"

    source_configuration {
        image_repository {
            image_configuration {
                port = "8000"
            }
            image_identifier      = "public.ecr.aws/jg/hello:latest"
            image_repository_type = "ECR_PUBLIC"
        }
    }

    tags = {
        Name = "example-apprunner-service"
    }
}
```

Infrastructure as Code / CloudFormation

27 lines (27 sloc) | 915 Bytes

Raw Blame ⌂ ⌄ ⌒ ⌓

```
1 # Copyright Amazon.com Inc. or its affiliates. All Rights Reserved.
2 # SPDX-License-Identifier: Apache-2.0
3 AWSTemplateFormatVersion: 2010-09-09
4 Description: CloudFormation template that deploys hello-app-runner app
5 Parameters:
6   ServiceName:
7     Type: String
8     Default: hello-app-runner
9     Description: Name for your App Runner service.
10 Resources:
11   Service:
12     Metadata:
13       'aws:apprunner:description': 'hello-app-runner example service'
14     Type: AWS::AppRunner::Service
15     Properties:
16       ServiceName: !Ref ServiceName
17       SourceConfiguration:
18         AutoDeploymentsEnabled: false
19         ImageRepository:
20           ImageIdentifier: public.ecr.aws/aws-containers/hello-app-runner:latest
21           ImageRepositoryType: ECR_PUBLIC
22           ImageConfiguration:
23             Port: 8000
24     Outputs:
25       Endpoint:
26         Description: "The endpoint of the App Runner service."
27         Value: !GetAtt Service.ServiceUrl
```

Infrastructure as Code / CDK

30 lines (24 sloc) | 841 Bytes

Raw Blame   

```
1 import { App, CfnOutput, Stack, StackProps } from 'aws-cdk-lib';
2 import * as apprunner from 'aws-cdk-lib/aws-apprunner';
3 import { Construct } from 'constructs';
4
5 export class AppRunnerStack extends Stack {
6     constructor(scope: Construct, id: string, props: StackProps = {}) {
7         super(scope, id, props);
8
9         const service = new apprunner.CfnService(this, 'service', {
10             serviceName: 'hello-apprunner',
11             sourceConfiguration: {
12                 autoDeploymentsEnabled: false,
13                 imageRepository: {
14                     imageIdentifier: 'public.ecr.aws/aws-containers/hello-app-runner:latest',
15                     imageRepositoryType: 'ECR_PUBLIC',
16                 },
17             },
18         });
19
20         new CfnOutput(this, 'url', {
21             value: 'https://' + service.attrServiceUrl,
22         });
23     }
24 }
25
26 const app = new App();
27
28 new AppRunnerStack(app, 'hello-apprunner');
29
30 app.synth();
```

github.com

TensorFlow Serving

Ubuntu passing Ubuntu TF Head passing Docker CPU Nightly passing Docker GPU Nightly failing

TensorFlow Serving is a flexible, high-performance serving system for machine learning models, designed for production environments. It deals with the *inference* aspect of machine learning, taking models after *training* and managing their lifetimes, providing clients with versioned access via a high-performance, reference-counted lookup table. TensorFlow Serving provides out-of-the-box integration with TensorFlow models, but can be easily extended to serve other types of models and data.

To note a few features:

- Can serve multiple models, or multiple versions of the same model simultaneously
- Exposes both gRPC as well as HTTP inference endpoints
- Allows deployment of new model versions without changing any client code
- Supports canarying new versions and A/B testing experimental models
- Adds minimal latency to inference time due to efficient, low-overhead implementation
- Features a scheduler that groups individual inference requests into batches for joint execution on GPU, with configurable latency controls
- Supports many *servables*: Tensorflow models, embeddings, vocabularies, feature transformations and even non-Tensorflow-based machine learning models

Serve a Tensorflow model in 60 seconds

```
# Download the TensorFlow Serving Docker image and repo
docker pull tensorflow/serving

git clone https://github.com/tensorflow/serving
# Location of demo models
TESTDATA="$(pwd)/serving/tensorflow_serving/servables/tensorflow/testdata"
```

No packages published

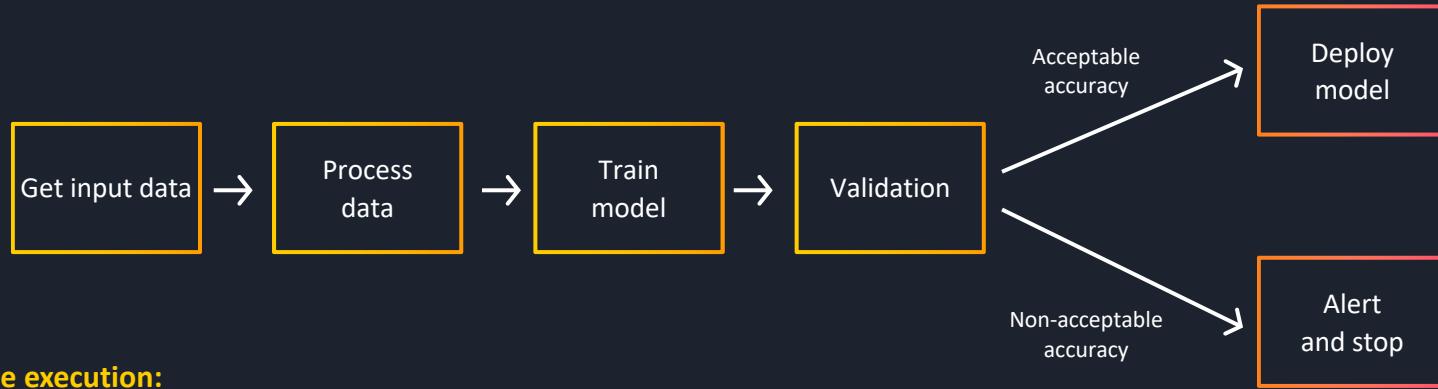
Contributors 181

+ 170 contributors

Languages

Language	Percentage
C++	87.8%
Starlark	6.4%
Python	4.5%
Other	1.3%

ML Pipelines



Start pipeline execution:

- Manually
- New data uploads
- On schedule
- Upon code check-in

Creating and managing ML workflows



Compose and manage workflows

Create your ML workflows using the Python. Define the parameters and steps

Visualize the workflows

Include steps such as data transformations, training, debugging, and optimizing models

```
# Example processing step for feature engineering

step_process = ProcessingStep(
    name="PreprocessAbaloneData",
    processor=sklearn_processor,
    outputs=[
        ProcessingOutput(output_name="train", source="/opt/ml/processing/train"),
        ProcessingOutput(output_name="validation", source="/opt/ml/processing/validation"),
        ProcessingOutput(output_name="test", source="/opt/ml/processing/test"),
    ],
    code=os.path.join(BASE_DIR, "preprocess.py"),
    job_arguments=["--input-data", input_data],
)

# Example training step

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",
        ),
    },
)
```

Task orchestration options

Open source 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

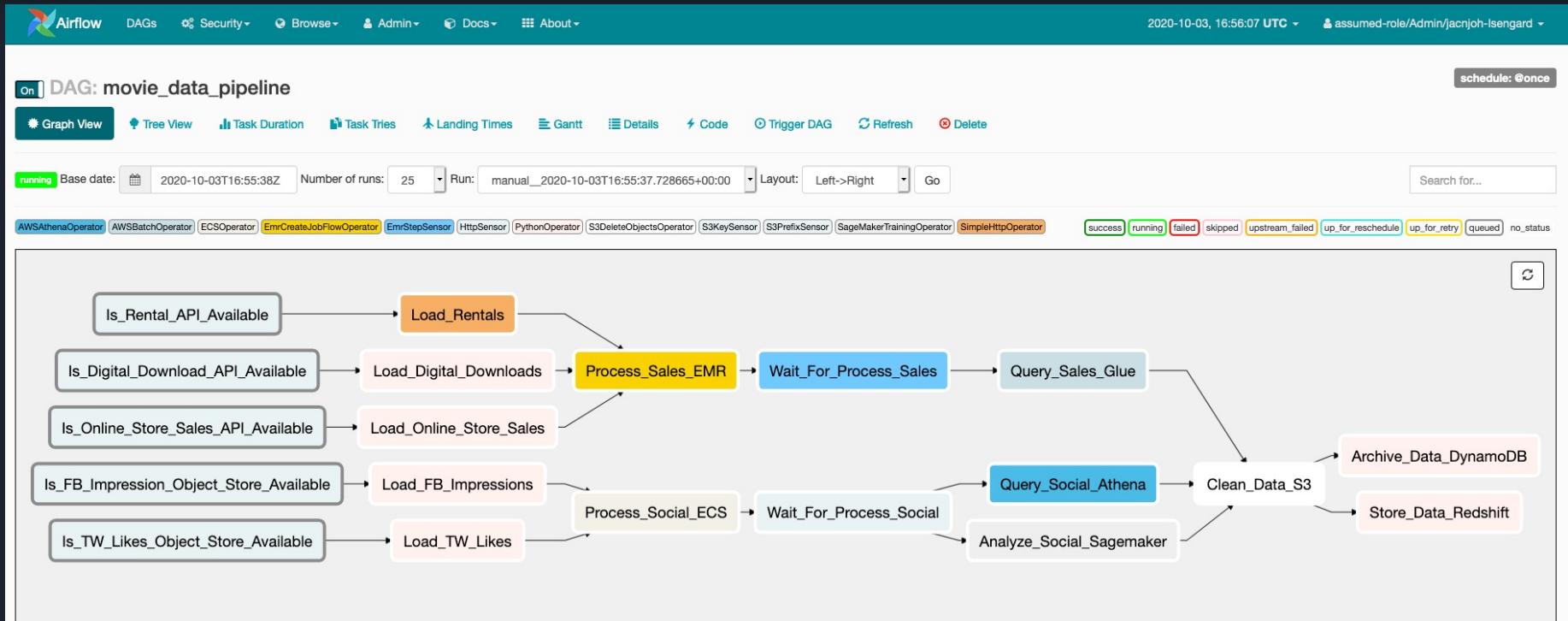
Serverless pipeline orchestration



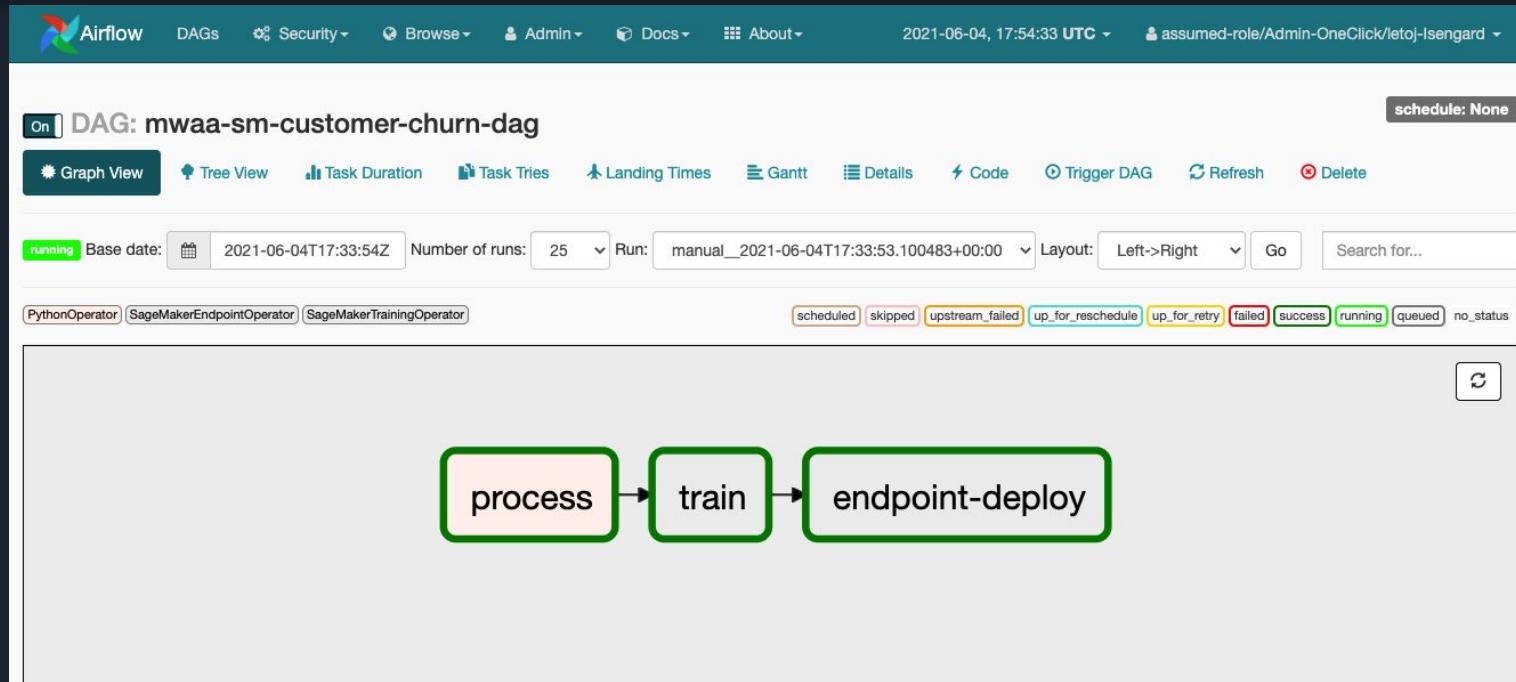
Amazon SageMaker Pipelines

Managed ML pipelines in SageMaker Studio

Apache Airflow



Apache Airflow



<https://aws.amazon.com/blogs/machine-learning/orchestrate-xgboost-ml-pipelines-with-amazon-managed-workflows-for-apache-airflow/>

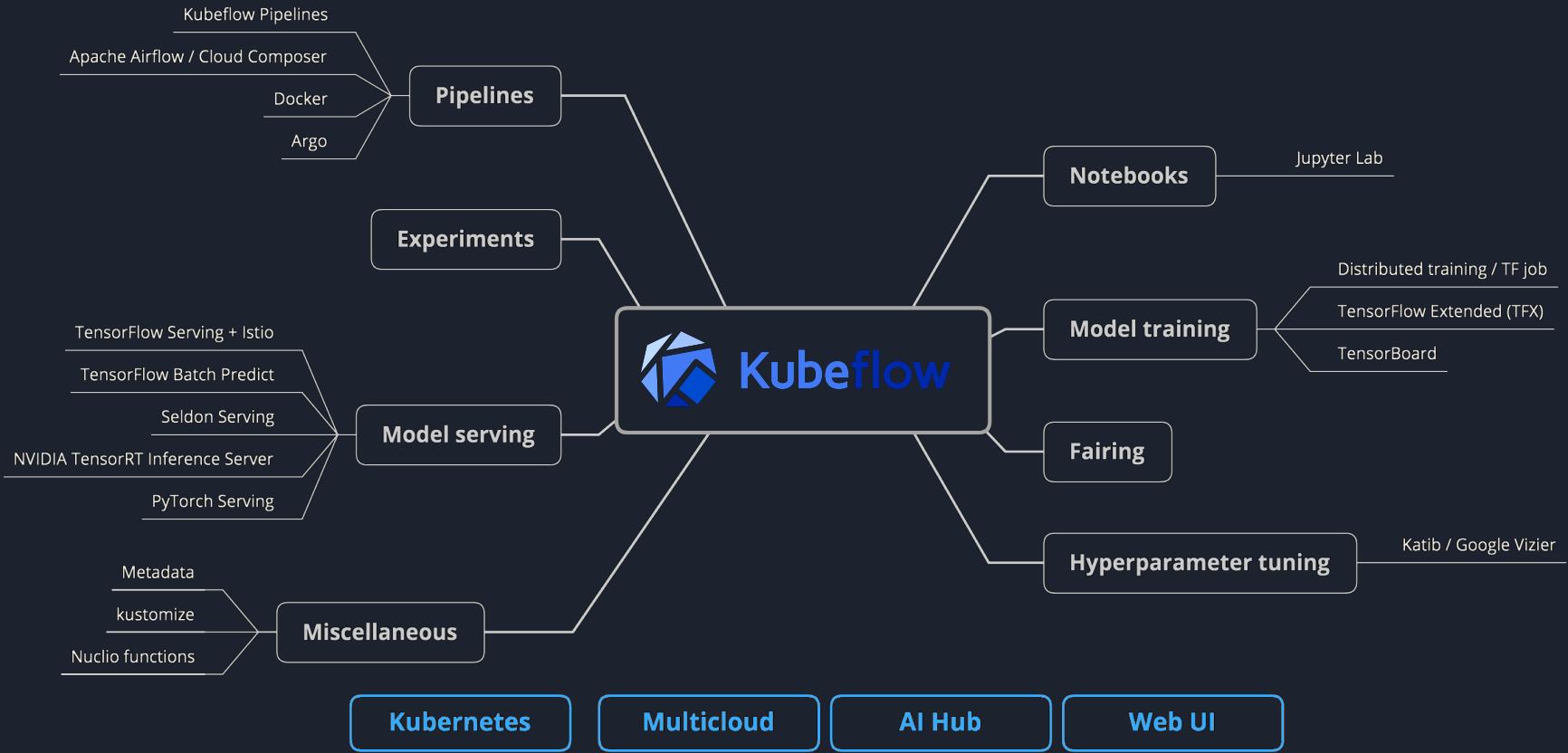
Monitoring Pipeline

*Machine Learning lifecycle does not end after you have deployed your Machine Learning model.
Part of the lifecycle is about tracking the model performance in production.*

There are multiple sources of errors that could degrade your ML model performance, such as:

Model Drift

Data Drift



The AWS ML Stack

Broadest and most complete set of Machine Learning capabilities

AI SERVICES

VISION



Amazon
Rekognition

SPEECH



Amazon
Polly
+Medical

TEXT



Amazon
Comprehend
+Medical

SEARCH



Amazon
Translate



Amazon
Textract

SEARCH



Amazon
Kendra

CHATBOTS



Amazon
Lex

PERSONALIZATION



Amazon
Personalize

FORECASTING



Amazon
Forecast

FRAUD



Amazon
Fraud Detector

DEVELOPMENT



Amazon
CodeGuru

CONTACT CENTERS



Contact Lens
For Amazon Connect

ML SERVICES



Amazon
SageMaker

Ground
Truth

ML
Marketplace

Built-in
algorithms

Notebooks

Experiments

Model
training &
tuning

SageMaker Studio IDE

Debugger

Autopilot

Model
hosting

Model Monitor

Neo

Augmented
AI

ML FRAMEWORKS & INFRASTRUCTURE

TensorFlow

mxnet

PYTORCH

GLUON

MXNet
gluon

Keras

DeepGraphLibrary

Deep Learning
AMIs & Containers

GPUs &
CPUs

Elastic
Inference

Inferentia

FPGA



Questions?
DevOps & MLOps

> Thank you! ■