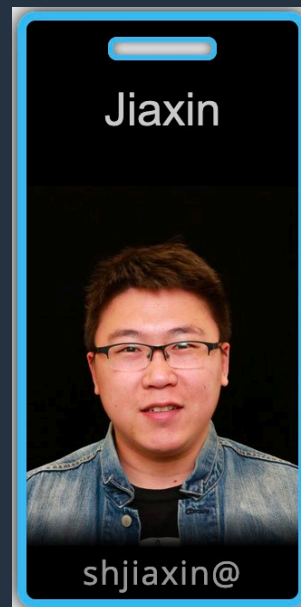# Machine Learning with Kubernetes

Yaniv Donenfeld
AI/ML Solutions
Container Services, AWS

Jiaxin Shan
Software Development Engineer
Container Services, AWS

"Cloud has removed so many of the barriers to experimenting and innovating with AI that even risk-adverse businesses are making it part of their strategies."

*- Yaniv Donenfeld, just now.*

**40%** of digital transformation initiatives supported by AI in 2019 *—IDC 2018*

aws

# Our mission at AWS

---

Put machine learning in the hands
of every developer

aws

# The AWS ML Stack

## Broadest and deepest set of capabilities

### AI Services

| VISION | | | SPEECH | | LANGUAGE | | CHATBOTS | FORECASTING | RECOMMENDATIONS |
|---|---|---|---|---|---|---|---|---|---|
| REKOGNITION IMAGE | REKOGNITION VIDEO | TEXTRACT | POLLY | TRANSCRIBE | TRANSLATE | COMPREHEND & COMPREHEND MEDICAL | LEX | FORECAST | PERSONALIZE |

### ML Services

| Amazon SageMaker | Ground Truth | Notebooks | Algorithms + Marketplace | Reinforcement Learning | Training | Optimization | Deployment | Hosting |
|---|---|---|---|---|---|---|---|---|

### ML Frameworks + Infrastructure

| FRAMEWORKS | | INTERFACES | INFRASTRUCTURE | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TensorFlow  mxnet  PYTORCH | | GLUON  Keras | EC2 P3 & P3DN | EC2 G4 EC2 C5 | FPGAS | DL CONTAINERS & AMIs | ELASTIC CONTAINER SERVICE | ELASTIC KUBERNETES SERVICE | GREENGRASS | ELASTIC INFERENCE | INFERENTIA |

aws

# Why Machine Learning on Kubernetes?

ON-PREMISES

CLOUD

Composability

Portability

Scalability

# Use Case #1: Large Scale ML

aws

Autonomous Vehicles Workloads

# Typical Autonomous Vehicle Development Workflow



2 Data Ingestion

3 Data Pre-processing

4 Labeling

5 Model Training

6 Model Simulation (SIL/HIL)

7 Evaluation & Validation

1 Data Acquisition

8 Model Deployment and CI/CD

aws

# Typical Autonomous Vehicle Development Workflow



1 Data Acquisition

2 Data Ingestion

3 Data Pre-processing

4 Labeling

5 Model Training

6 Model Simulation (SIL/HIL)

7 Evaluation & Validation

8 Model Deployment and CI/CD

aws

- Distributed Training Challenges

- Single GPU code ➔ multiple | Horovod + MPIJob (or TFJob)

- Dataset Copying time

Use FSx Lustre / EFS

- Dataset Sharing and Reuse

👉 Built-in CSI driver with S3 integration

aws

# Want to Run Distributed Training on EKS?

Ajay

ajayvohr@

Distributed TensorFlow
training using Kubeflow on
Amazon EKS

Ajay Vohra
Principal SA -
Vision/AI/ML

aws

Typical Autonomous Vehicle Development Workflow

2 Data Ingestion

3 Data Pre-processing

4 Labeling

1 Data Acquisition

5 Model Training

6 Model Simulation (SIL/HIL)

7 Evaluation & Validation

8 Model Deployment and CI/CD

aws

Can you run my workload?

Concurrent CPUs

# Total Core Hours / Year

# Simulations Architecture

# Simulations Architecture

# TOP500 – Top 10 Supercomputers in June 2019

| Rank / Name | Rmax / Rpeak (Petaflops) |
|---|---|
| 1. Summit | 148.600 / 200.795 |
| 2. Sierra | 94.640 / 125.712 |
| 3. Sunway Tahihu Light | 93.015 / 125.436 |
| 4. Tianhe-2A | 61.445 / 100.679 |
| 5. Frontera | 23.516 / 38.746 |
| 6. Piz Daint | 21.230 / 27.154 |
| 7. Trinity | 20.159 / 41.461 |
| 8. AI Bridging Cloud Infrastructure | 19.880 / 32.577 |
| 9. SuperMUCNG | 19.477 / 26.874 |
| 10. Lassen | 18.200 / 23.047 |

aws

We're helping our customers run at Supercomputer Scale, targeting the equivalent of **one of the Top 10 largest supercomputers** in the world.

aws

# Use Case #2: ML Development Platform

aws

Kubeflow

aws

# Jupyter Notebook / JupyterHub

- Build, deploy, and train ML models

- Live code, equations, visualizations, and narrative text

- 40+ programming languages

- Sharing and collaboration

👉 EFS for reusing training data and results

👉 Built-in AWS CLI and ECR support

# Kubeflow KFServing

- Simple and pluggable platform for ML inference
- Intuitive and consistent experience
- Serving models on arbitrary frameworks
-   e.g. TensorFlow, XGBoost, SciKitLearn
- Encapsulates GPU auto-scaling, canary rollouts

Credits @ellis-bigelow (Kubeflow slack)

aws

# Kubeflow KFServing

# Pluggable Interface

```yaml
apiVersion: "serving.kubeflow.org/v1alpha1"
kind: "InferenceService"
metadata:
  name: "sklearn-iris"
spec:
  default:
    sklearn:
      storageUri: "gs://kfserving-samples/models/sklearn/iris"
```
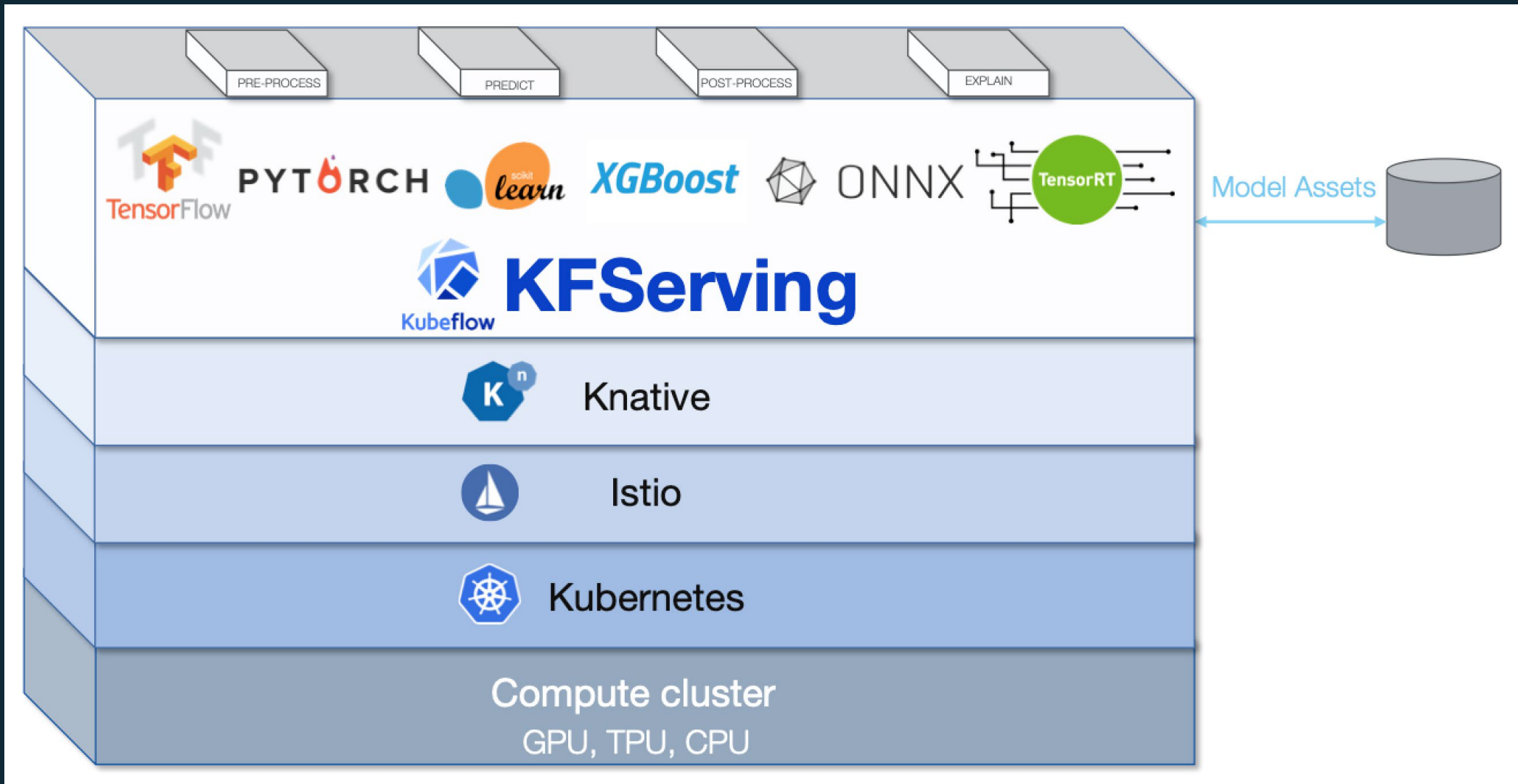
```yaml
apiVersion: "serving.kubeflow.org/v1alpha1"
kind: "InferenceService"
metadata:
  name: "flowers-sample"
spec:
  default:
    tensorflow:
      storageUri: "gs://kfserving-samples/models/tensorflow/flowers"
```

```yaml
apiVersion: "serving.kubeflow.org/v1alpha1"
kind: "KFService"
metadata:
  name: "pytorch-cifar10"
spec:
  default:
    pytorch:
      storageUri: "gs://kfserving-samples/models/pytorch/cifar10"
      modelClassName: "Net"
```

scikit learn

TensorFlow

PYTORCH

aws

# Kubeflow Pipelines

- A user interface (UI) for managing and tracking experiments, jobs, and runs.

- An engine for scheduling multi-step ML workflows.

- An SDK for defining and manipulating pipelines and components.

# Kubeflow Pipelines Component



Kubeflow Pipeline

Pipeline step

Metadata

Input/Output

Implementation
(container)

Component

Pipeline step

Metadata

Input/Output

Implementation
(container)

Component

Pipeline step

Metadata

Input/Output

Implementation
(container)

Component

Container registry

aws

# Creating a pipeline



Pipeline decorator

Pipeline function

Pipeline component

Compile pipeline

```
@dsl.pipeline(
  name='Sample Trainer',
  description=''
)

def sample_train_pipeline(... ):

    create_cluster_op = CreateClusterOp('create-cluster', ...)

    analyze_op = AnalyzeOp('analyze', ...)

    transform_op = TransformOp('transform', ...)

    train_op = TrainerOp('train', ...)

    predict_op = PredictOp('predict', ...)

    confusion_matrix_op = ConfusionMatrixOp('confusion-matrix', ...)

    roc_op = RocOp('roc', ...)

kfp.compiler.Compiler().compile(sample_train_pipeline , 'my-pipeline.zip')
```
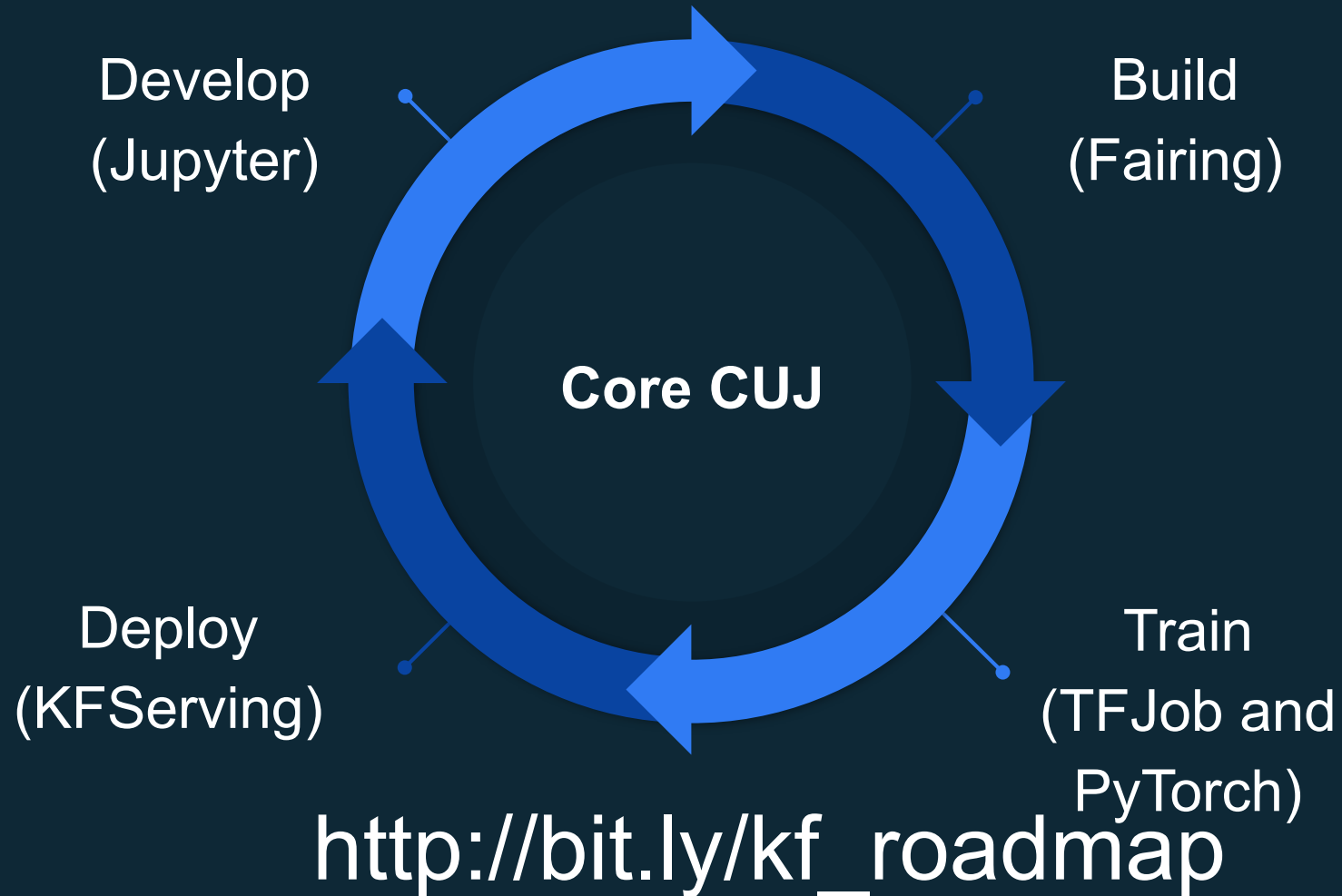
# Kubeflow 1.0 Arriving January 2020



Develop (Jupyter)

Build (Fairing)

Core CUJ

Deploy (KFServing)

Train (TFJob and PyTorch)

http://bit.ly/kf_roadmap

aws

# Kubeflow 1.0 – Main components

- Graduating 1.0

  - kfctl for deployment and upgrades

  - TFJob and PyTorch for distributed training (already 1.0)

  - Jupyter notebook controller and web app

  - Profile controller and UI for multiuser management

- Beta

  - Katib for hyper-parameter tuning

  - Fairing SDK to facilite use of notebooks for build-train-deploy

  - Metadata SDK, UI, and backend

  - KFServing for model deployment and inference

# Kubeflow 1.0 – AWS Support

- Multi user support

  - Kubeflow pipelines

  - Managed contributors

- IAM Roles for Service Accounts integration with notebooks

# Want to Dive Deeper on Kubeflow?

## Now

2:30PM    Kubeflow Workshop
(Workshop Room Harborside)

## Later

https://eksworkshop.com/kubeflow/

aws

# Join the kubeflow#aws Slack channel !

aws

# DEMO

aws