# Serving Tensorflow Models with Kubernetes

Ron Bodkin, Technical Director, Applied AI

Brian Foo, Senior Software Engineer, Applied AI

**Step One: Get An Account**
**Go To goo.gl/sgZ2Qp to start**

Google Cloud
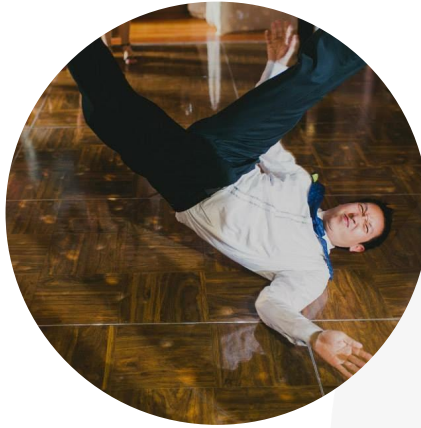
# Your Instructors

Ron
Bodkin

Brian
Foo

Cassie
Kozyrkov

Set up an account here: **goo.gl/sgZ2Qp**

# Your Instructors



Ron
Bodkin

Brian
Foo

Cassie
Kozyrkov

Set up an account here: **goo.gl/sgZ2Qp**

# Overview

- Tensorflow
  - Building and Running Tensorflow graphs & APIs
  - Tensorflow Serving for Online Predictions

- Docker and Kubernetes
  - A brief history of containers
  - Kubernetes: pods and services

- The Marriage of TF Serving and Kubernetes

# Building and Running

# Tensorflow Graphs

# Build a Tensorflow Graph

```python
import tensorflow as tf

x = tf.placeholder(tf.float32, shape=(100))

# Some preloaded model weights and biases

w = tf.get_variable('weights', shape=(100))

b = tf.get_variable('bias', shape=[])

y = tf.tensordot(w, x, 1) + b
```

# Run the Tensorflow Graph

```python
import tensorflow as tf

x = tf.placeholder(tf.float32, shape=(100))

# Some preloaded model weights and biases

w = tf.get_variable('weights', shape=(100))

b = tf.get_variable('bias', shape=[])

y = tf.tensordot(w, x, 1) + b
```
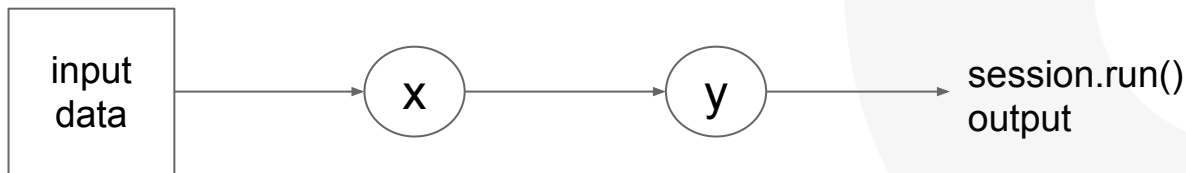
```python
with tf.Session() as sess:

    print(sess.run(y, feed_dict={x: input_data}))
```
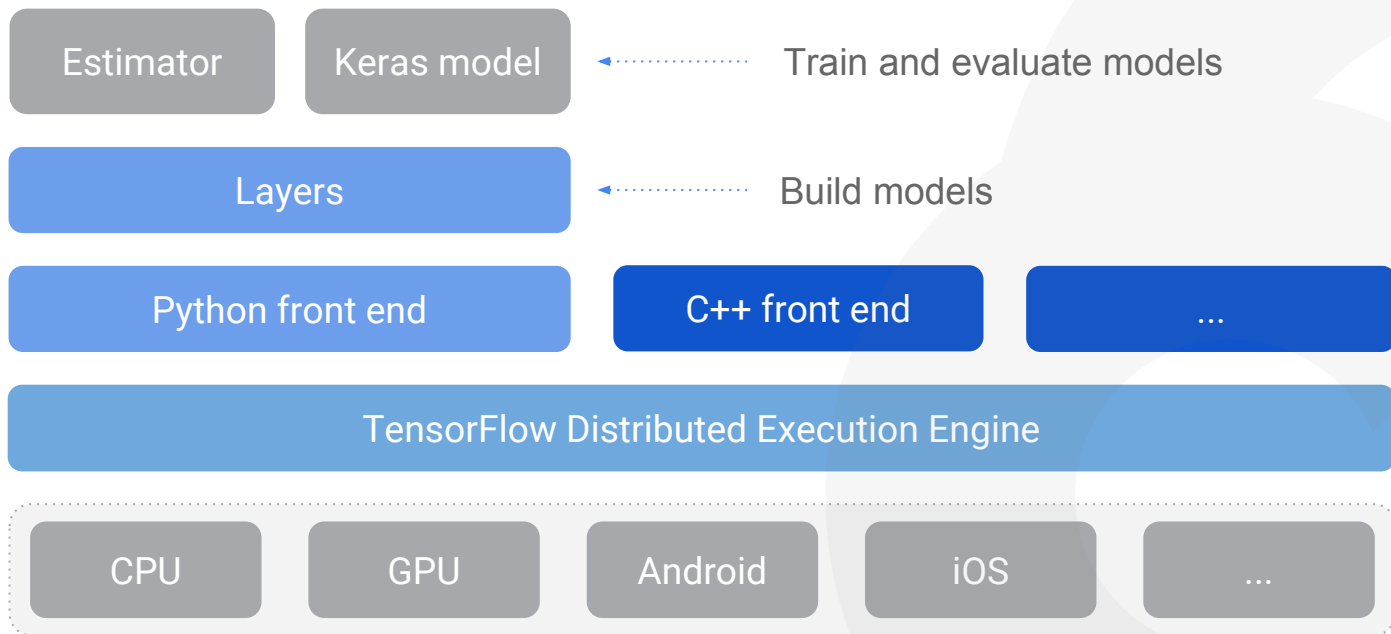
Data entrypoint for predictions

input data → x → y → session.run() output

# Different APIs for TF

| Estimator | Keras model | ◁┈┈┈┈┈ Train and evaluate models |

| Layers | ◁┈┈┈┈┈ Build models |

| Python front end | C++ front end | ... |

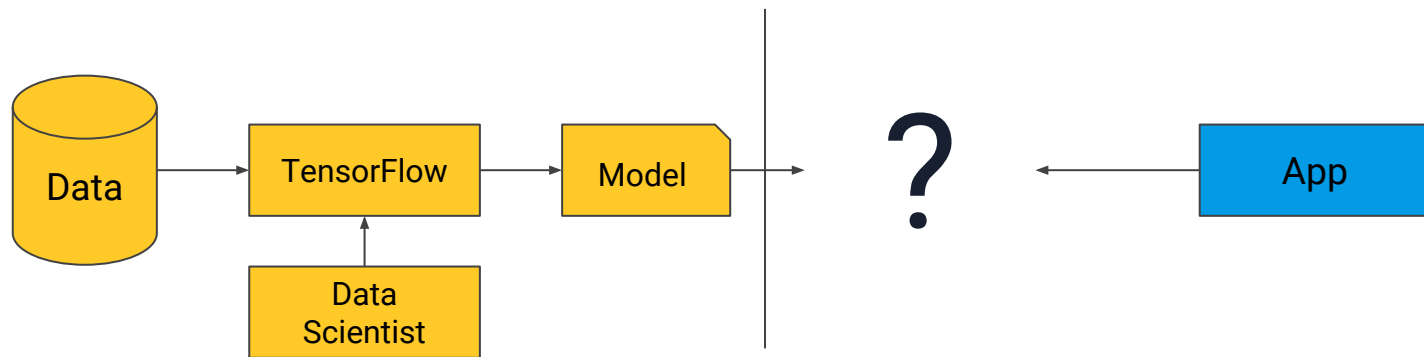| TensorFlow Distributed Execution Engine |

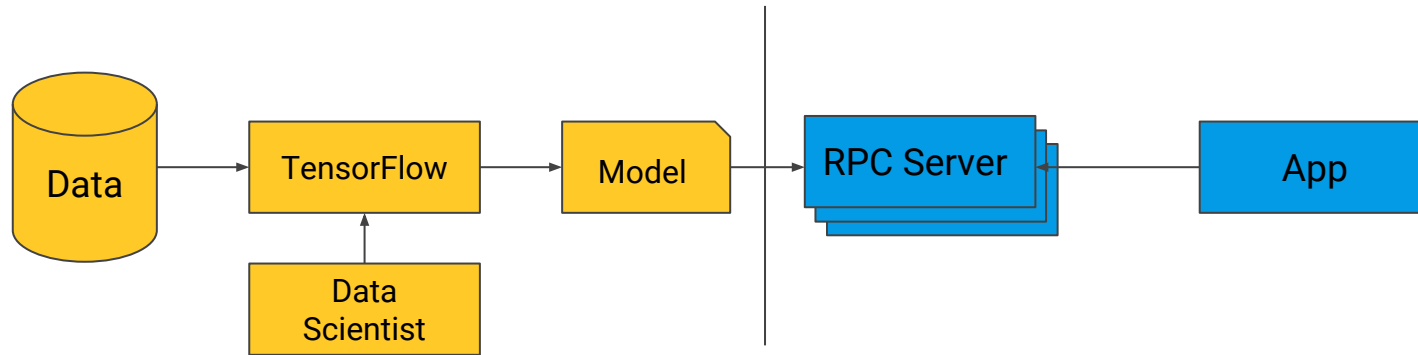| CPU | GPU | Android | iOS | ... |

# What is Serving?

Serving is how you *apply* a ML model, *after* you've trained it

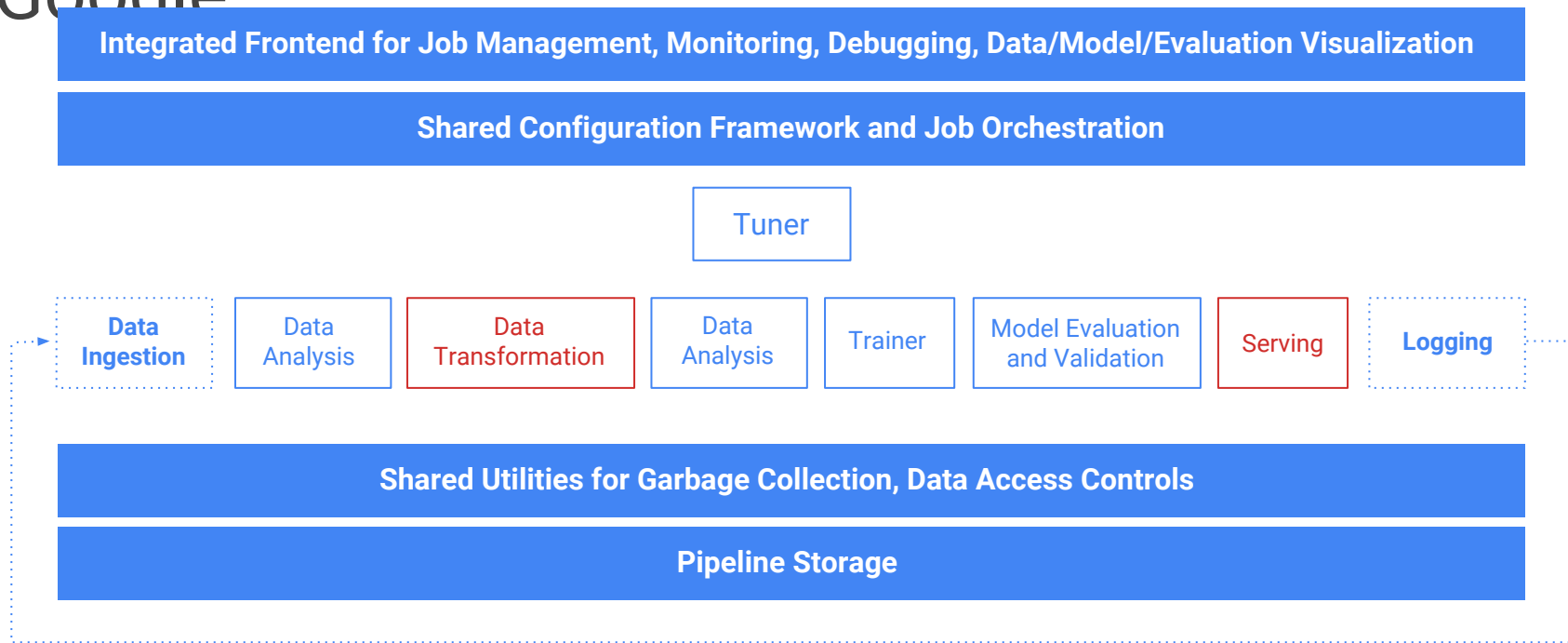# What is Serving?

# What is Serving?

# TensorFlow Serving Goals

- Online, low latency
- Multiple models in a single process
- Multiple *versions* of a model loaded over time
- Compute cost varies in real-time to meet product demand
  - auto-scale with CloudML, Docker & K8s
- Aim for the efficiency of mini-batching at training time …
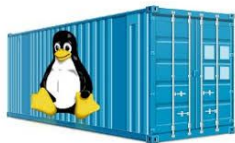  - except with requests arriving asynchronously

# TFX: TensorFlow-Based ML deployment at Google

**Integrated Frontend for Job Management, Monitoring, Debugging, Data/Model/Evaluation Visualization**

**Shared Configuration Framework and Job Orchestration**

Tuner

| **Data Ingestion** | Data Analysis | Data Transformation | Data Analysis | Trainer | Model Evaluation and Validation | Serving | **Logging** |

**Shared Utilities for Garbage Collection, Data Access Controls**

**Pipeline Storage**
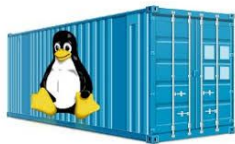
# Docker and Kubernetes

# A Brief History of Containers



2008

In 2008, Linux introduced containers.
- Isolated environment for running applications, except…
- All applications must have a common OS Kernel (e.g. Ubuntu, Debian, etc.)

# Enter Docker...



2008

2013

In 2013, Docker found a way around the shared kernel problem.
- Application containers have their own OS kernel.
- Flexible resource requirements.
- Perfect for cloud computing!
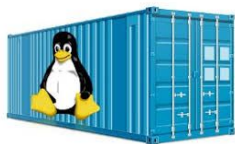
Google Cloud

# Enter Kubernetes

2013

2014

2008

In 2014, Google open sourced Kubernetes.
- Deploy Docker containers to any number of machines
- Create load balancing and front-end services to handle external requests.
- Automatically restart backend containers when they fail.

# And Many More...


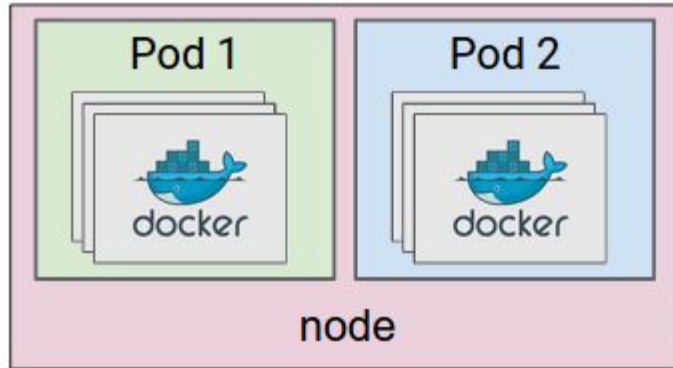
```
2008                    2013      2014    2015...
```

Example projects built on top of Kubernetes:
- Monitoring (Heapster)
- Deployment languages (KSonnet)
- Deployment automation (Kubeflow for ML)
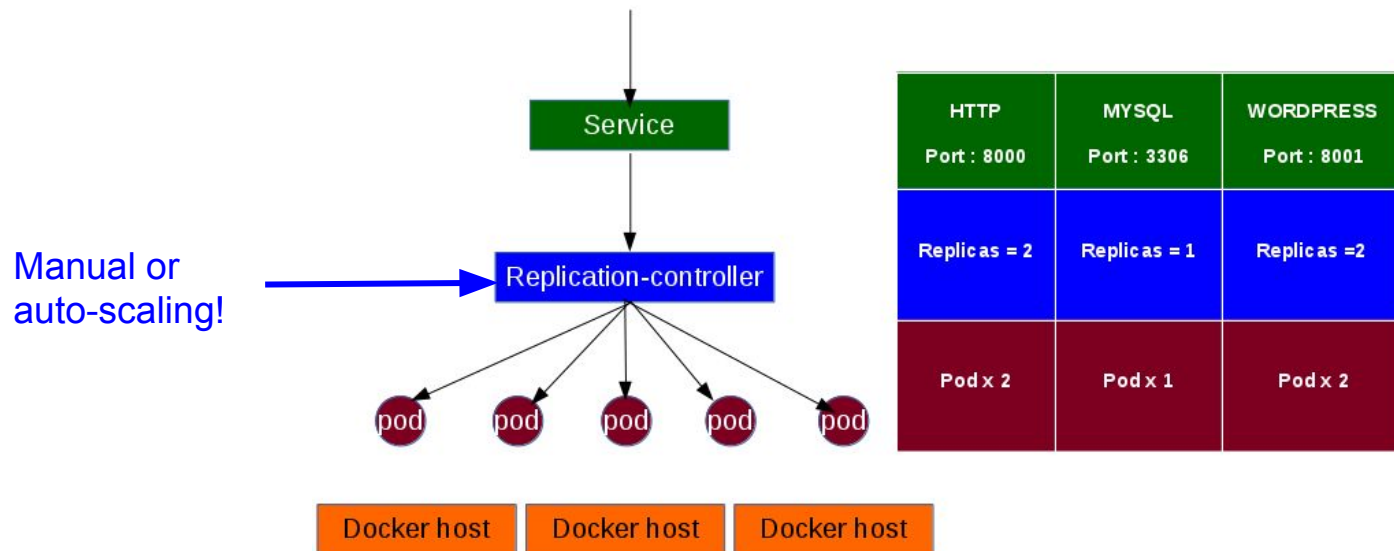- and many more!

# Kubernetes in a Nutshell



**Kubernetes Pods**
*collections of containers that are co-scheduled*

Pod 1 — docker

Pod 2 — docker

node

# Kubernetes in a Nutshell

**How Kubernetes Works?**

Service

Manual or
auto-scaling!  →  Replication-controller

pod  pod  pod  pod  pod

Docker host    Docker host    Docker host

| HTTP | MYSQL | WORDPRESS |
|---|---|---|
| Port : 8000 | Port : 3306 | Port : 8001 |
| Replicas = 2 | Replicas = 1 | Replicas =2 |
| Pod x 2 | Pod x 1 | Pod x 2 |

Google Cloud

# Tensorflow Serving

# using Kubernetes

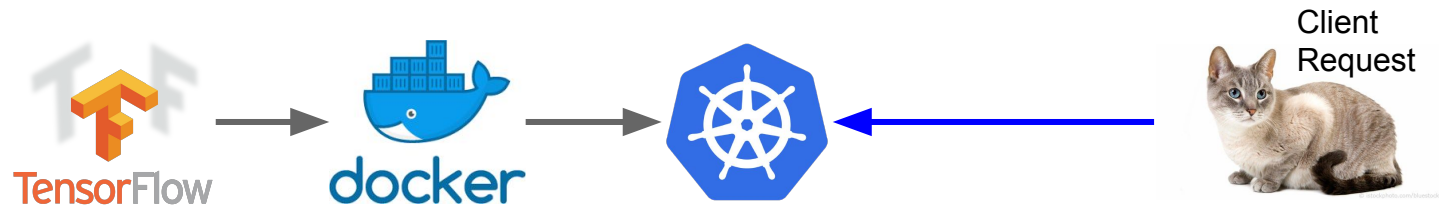# TF Serving on Kubernetes Workflow



Label: cat
Prob: 97.6%

## What do we want?

- A prediction service that can handle multiple client requests
- Load-balancing across TF model servers
- Ability to scale up

Google Cloud

# TF Serving on Kubernetes Workflow
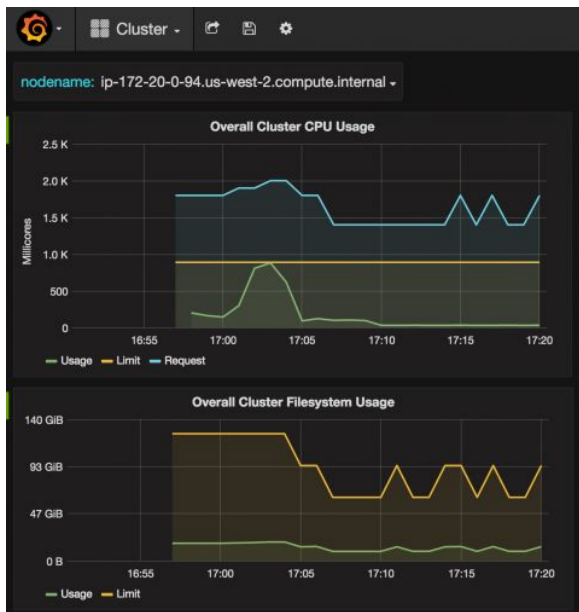## Exercises



## How do we get there?

- Convert TF training code to model for serving.
- Package model in Docker container and upload to a registry.
- Use **Kubernetes** to:
  - Deploy container on multiple back-end pods.
  - Deploy a front-end service to send client requests to a backend pod.
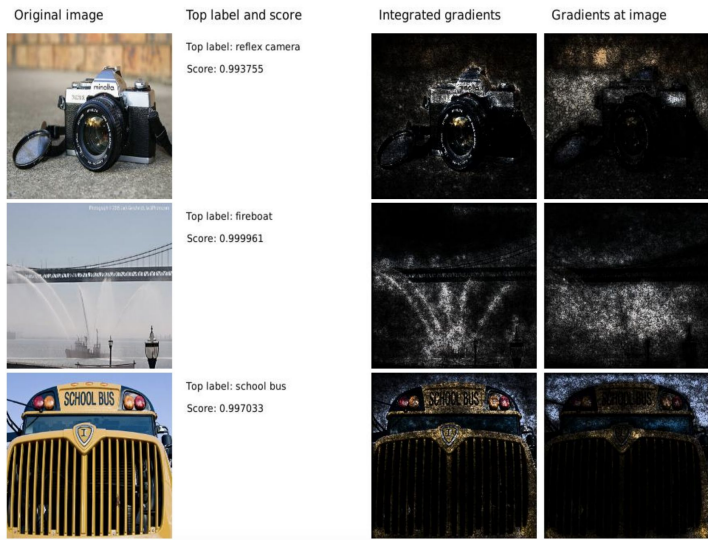- Send protobufs (encoded JSONs) containing images to kubernetes cluster.
- Load testing.

Google Cloud

# Monitoring, Interpretation, and Keras
## Bonus Exercises

### Heapster Grafana Dashboard
### (Pod and Cluster Resource Usage)



### Model Understanding and Visualization
### using Integrated Gradients

# Codelab Time!

Open a Chrome incognito window.

Log in at events.qwiklabs.com

If you don't have an account register at **goo.gl/sgZ2Qp**

# Recap and Demos

# Acknowledgements and Additional Resources

Special thanks to:

- [Kubeflow](#): providing Docker images and templates for TF Serving on Kubernetes
  David Aronchick, Jeremy Lewi, Vishnu Kannan (Google); Peng Yu (Shopify)

- [Google Cloud ML](#): GPU batch profiling work using Beam and Tensorflow

Reference - Model Visualization:

- Sundararajan, Taly, Yan. *Axiomatic Attribution for Deep Networks.* ICML 2017. Link:
  [arxiv.org/pdf/1703.01365.pdf](#)

Google Cloud

# Thank you

1. Please leave feedback
2. Resources at [goo.gl/Sg6ecA](goo.gl/Sg6ecA)
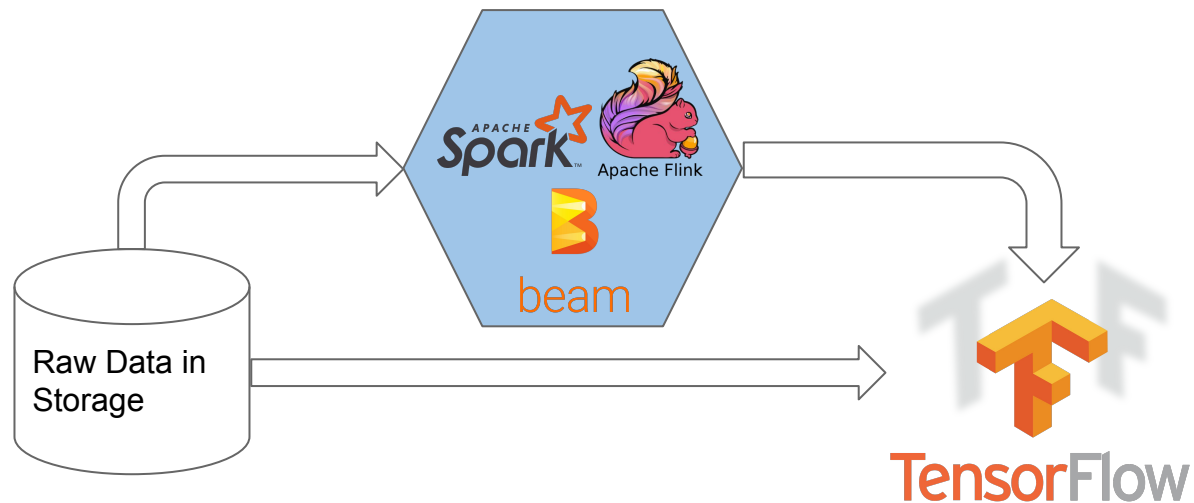3. Save any work you want to keep

# Appendix: pipelines versus

# Client-Server Architectures

# Pipeline Architecture: Batch/Online Processing

- Read offline data from local/HDFS/Google Storage/AWS
- Preprocess (clean, filter, aggregate) using Spark/Beam/Flink
- Create batches to run through a TF graph
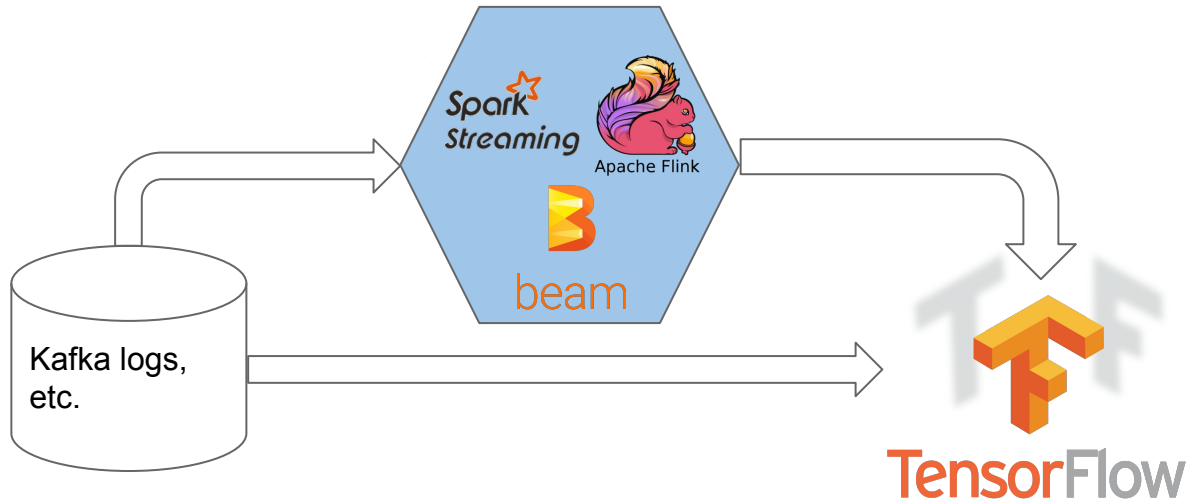- Update model params (training) / Collect inference results (serving)



Google Cloud

# Pipeline Architecture: Batch/Online Processing

- **Benefits**:
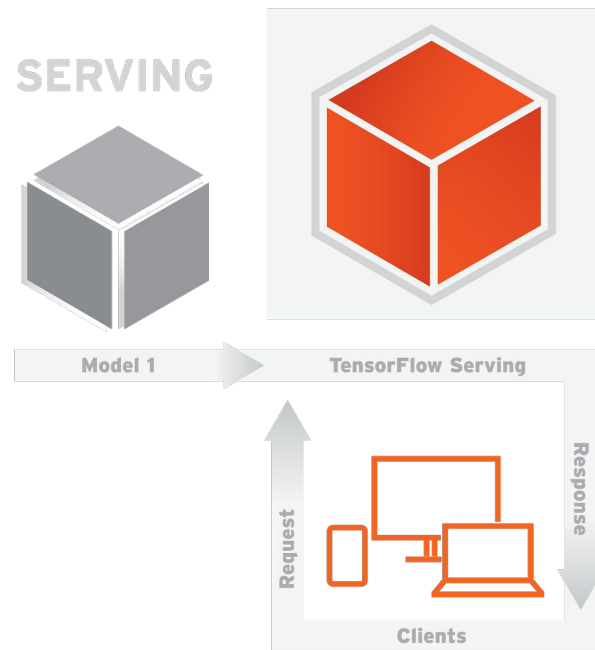  - Full control over pipeline application and model!
- **Limitations**:
  - **Language Dependency:** Requires Python, or Java JNI to C++
  - **No Proprietary Models:** Requires graph and model params to be exposed in code.
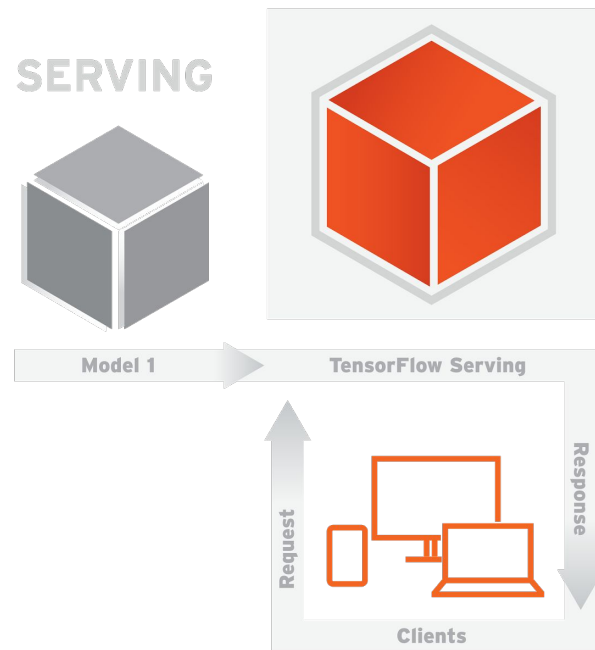  - **Experience:** Months to years of expertise to build, debug, and manage pipelines effectively.



Google Cloud

# Client-Server Architecture: Tensorflow Serving

- Asynchronous and Streaming Model Serving
- Efficient implementation in c++
- Server build can be optimized for native environment
    - CPUs or GPUs
    - Just-in-time (JIT) compilation
    - etc.



Google Cloud

# Client-Server Architecture: Tensorflow Serving

- Asynchronous and Streaming Model Serving
- Efficient implementation in c++
- Build can be optimized for the environment (CPUs or GPUs)
- **Language independent Protobufs!**
  - RESTful API calls using serialized dictionaries
  - Send dictionary of data
  - Receive dictionary of prediction results

# Client-Server Architecture: Tensorflow Serving

- Asynchronous and Streaming Model Serving
- Efficient implementation in c++
- Build can be optimized for the environment (CPUs or GPUs)
- **Language independent Protobufs!**
    - RESTful API calls using serialized dictionaries
    - Send dictionary of data
    - Receive dictionary of prediction results
- How do we guarantee identical serving environments?
- How do we scale?
- How do we handle failures gracefully?



SERVING

Model 1 → TensorFlow Serving

Request → Clients → Response