

TensorFlow on Google Cloud

May / 2020



1

Context

2

Why TensorFlow
on GCP?

3

Getting Started
with Tensorflow

4

Tensorflow
Extended

5

Hands On
Qwiklabs

Context

ML on Google Cloud



Our model +
Our data

API's

ML on Google Cloud



Our model +
Our data

API's

Our model +
Your data

AutoML

ML on Google Cloud



Our model +
Our data

API's

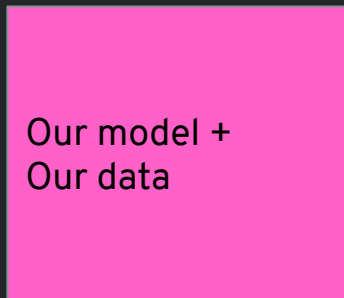
Our model +
Your data

AutoML

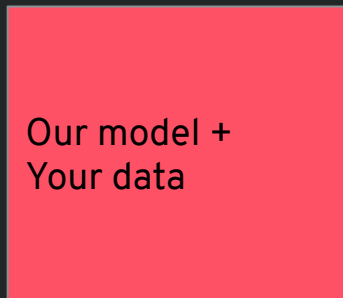
Your model +
Your data

AI Platform

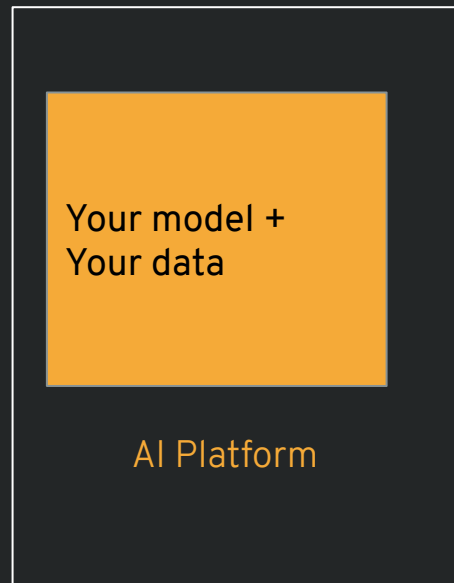
ML on Google Cloud



API's



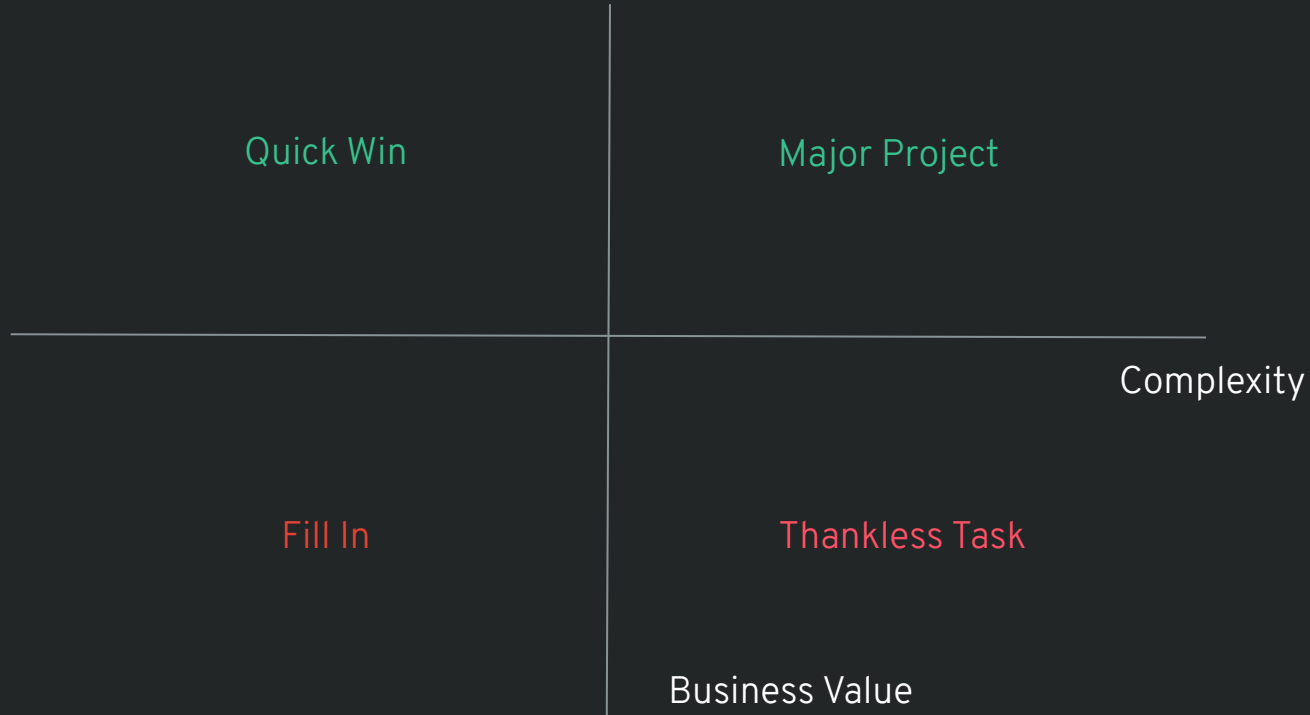
AutoML



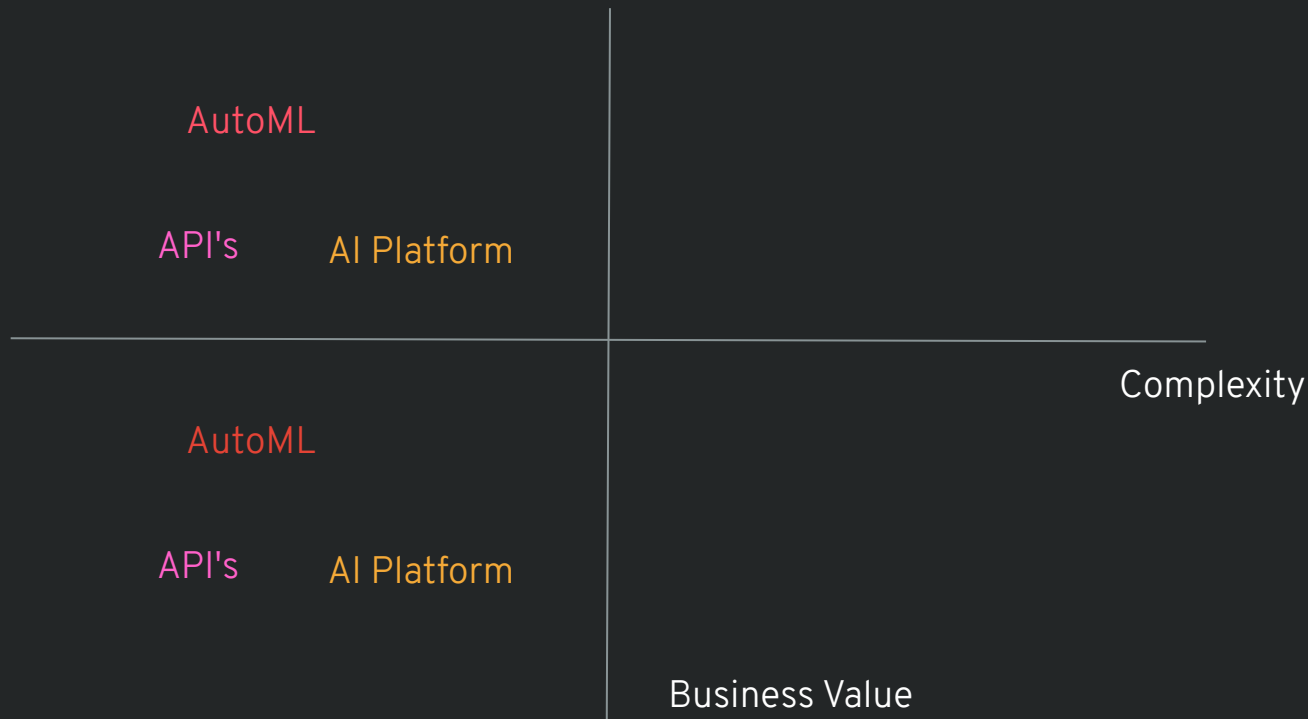
AI Platform

Today
Presentation

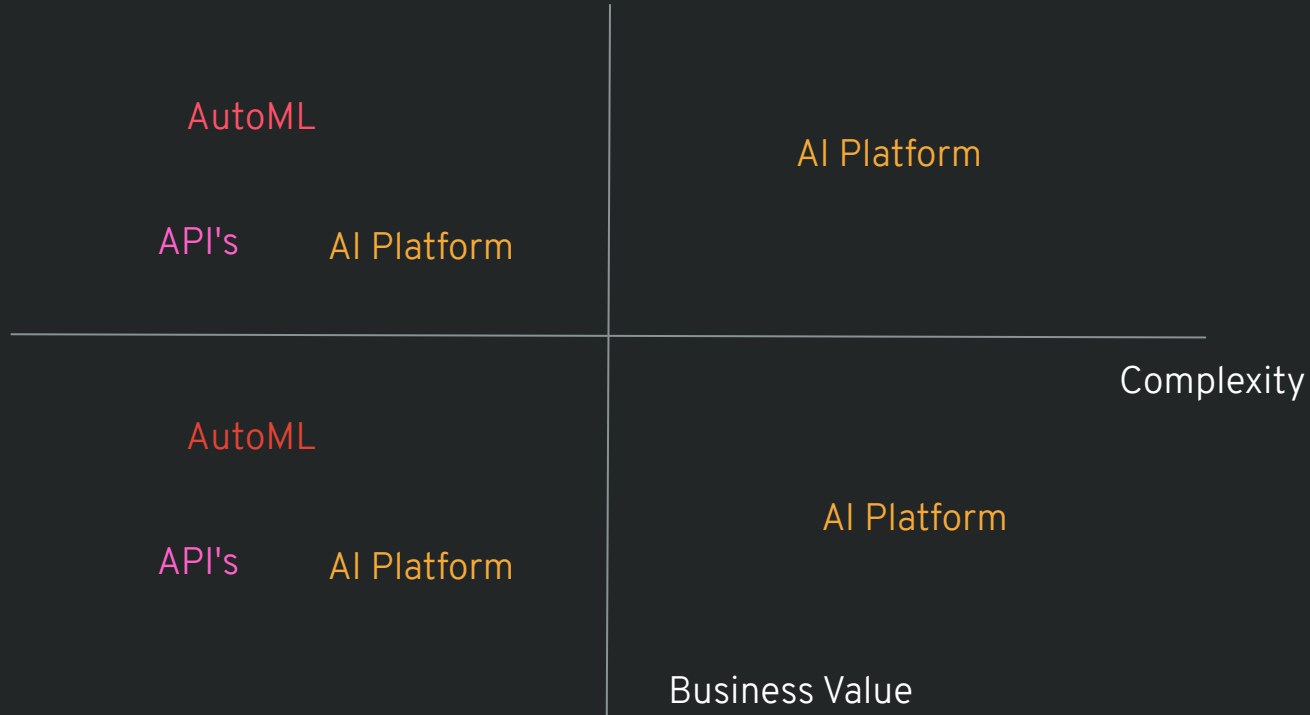
AvenueCode Process



AvenueCode Process



AvenueCode Process



Problem Example



Problem Example

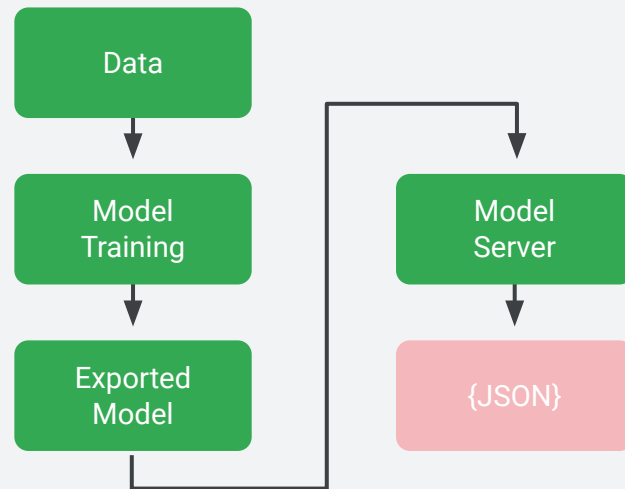
Is this a plastic bag?



Our model +
Our data

API's

Machine Learning APIs



Problem Example

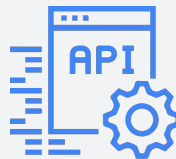
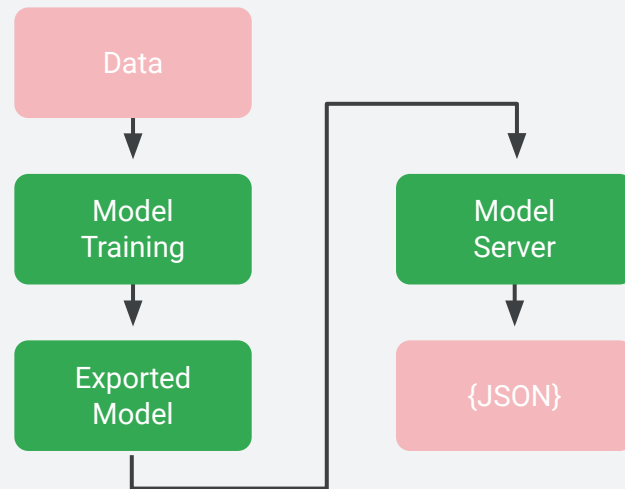
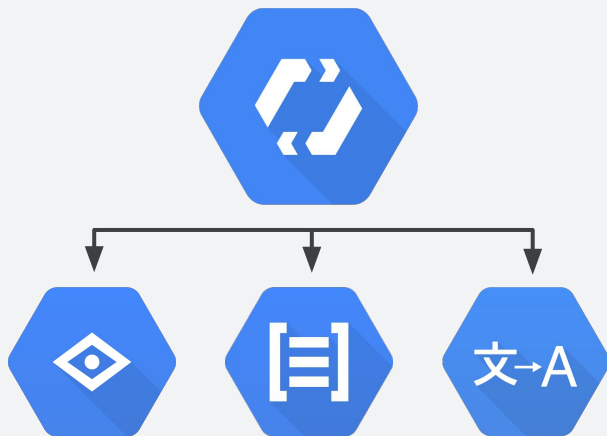
What's the plastic type?



Our model +
Your data

AutoML

AutoML



Problem Example

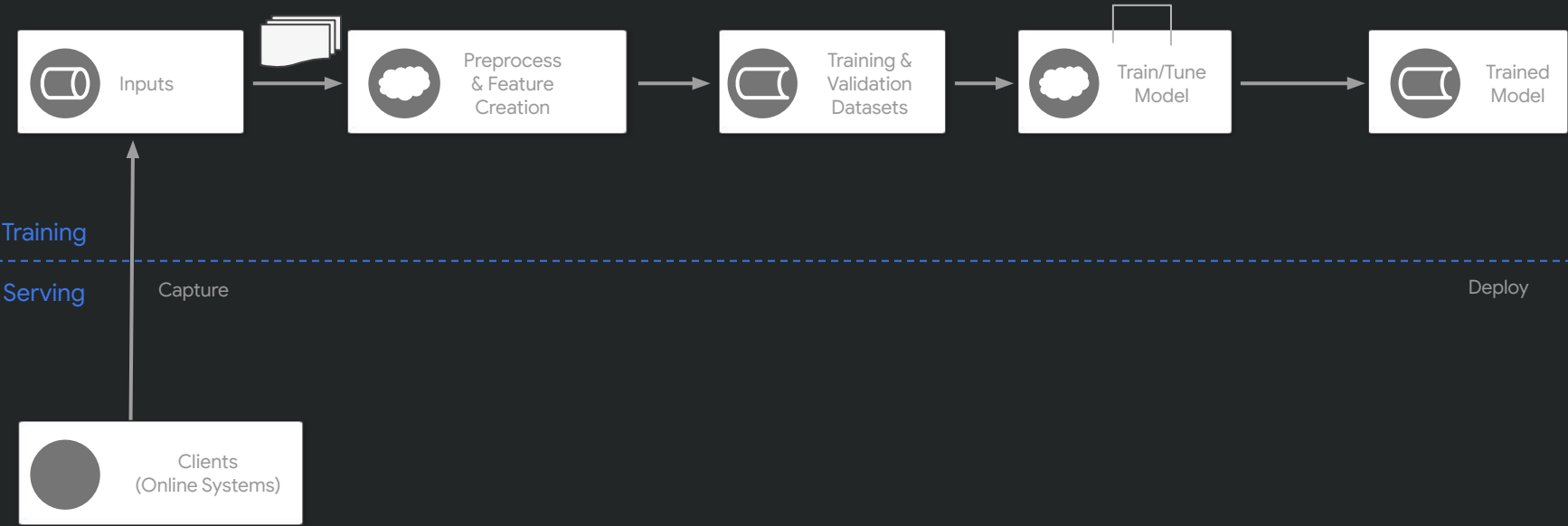
Will this plastic bag be recycled? From today how much time will take ?



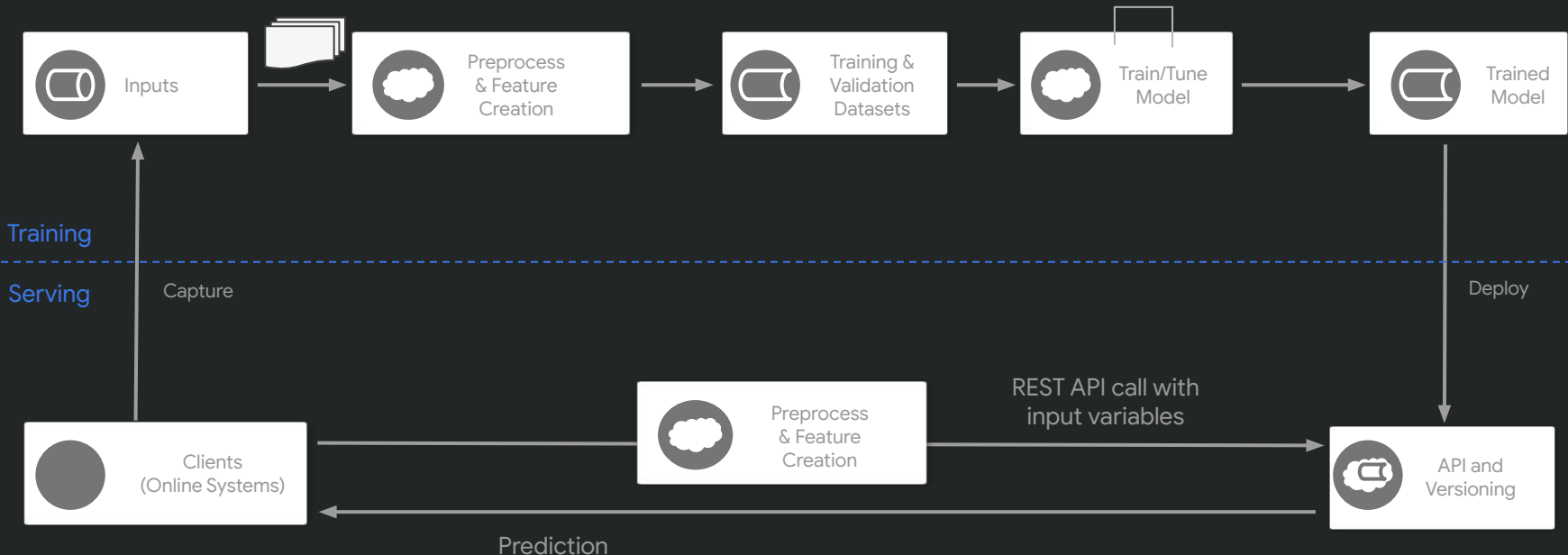
Your model +
Your data

AI Platform

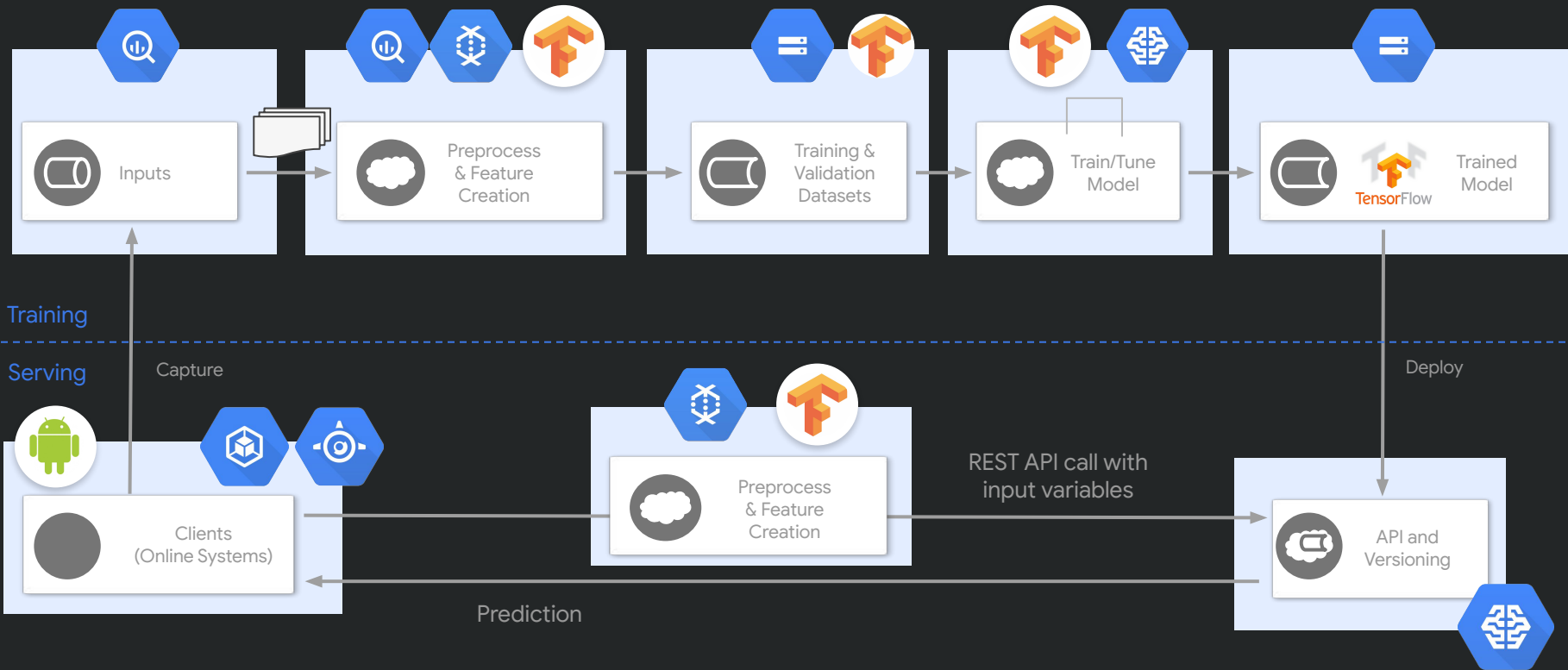
Machine learning pipeline @ GCP



Machine learning pipeline @ GCP



Machine learning pipeline @ GCP





Why TensorFlow on GCP?

“More than 87% of data science projects never make it into production”

- Multiple studies and surveys -



Hidden Technical Debt in Machine Learning Systems

D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips
`{dsculley, gholt, dgg, edavydov, toddphillips}@google.com`
Google, Inc.

Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-François Crespo, Dan Dennison
`{ebner, vchaudhary, mwyoung, jfcrespo, dennison}@google.com`
Google, Inc.

Why ML deployment is hard?

Deploying ML models at scale is very challenging

In IT System

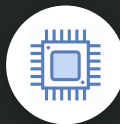
Behavior of system is defined by **Code**

In AI/ML System

Behavior of system is defined by **Data**



Organization - Managing different ecosystems like programming languages



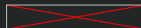
Compute - Continuous and reliable compute resources (dedicated servers or cloud only option)

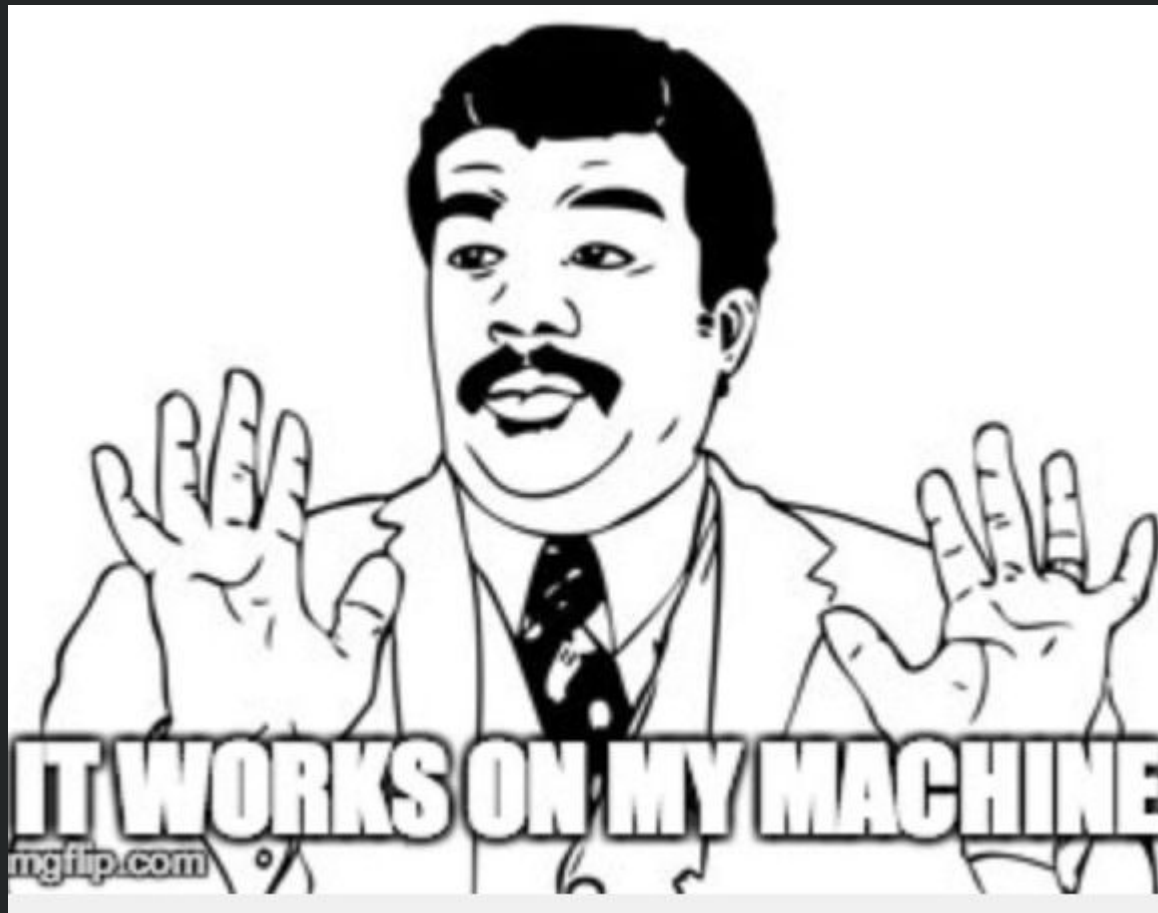


Portability - Huge dependencies on legacy systems



Seasonality - ML workloads works in patches; this need auto scaling capabilities





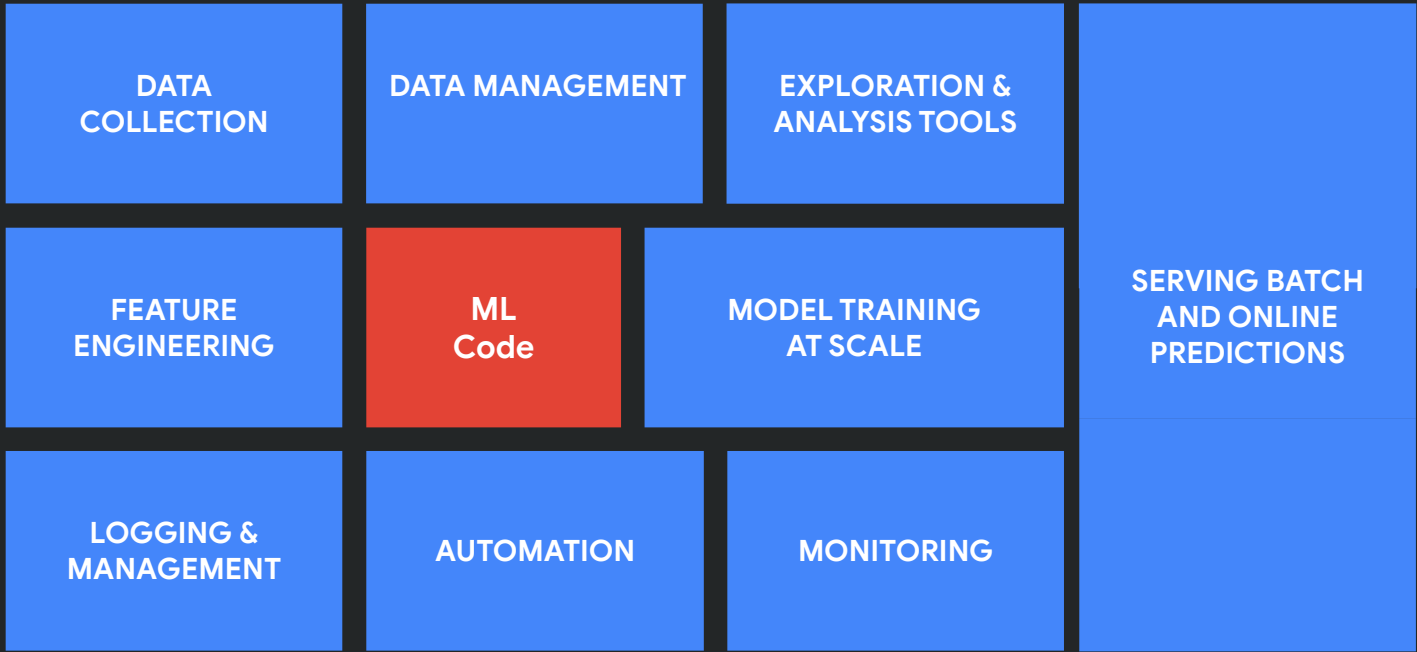
Operational ML : What you first think



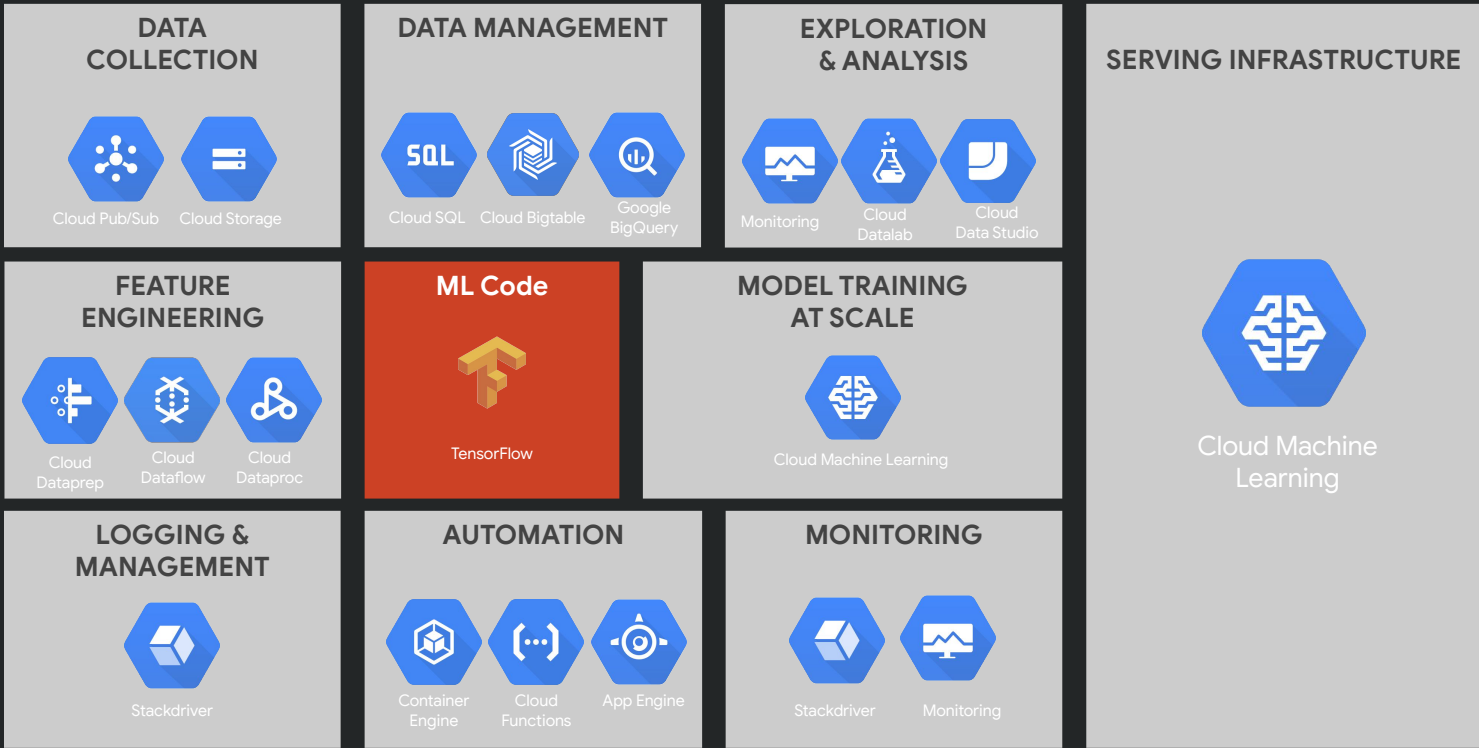
ML
Code



Operational ML : What really is



Operational ML - end-to-end ML solution on GCP





**Excuse me
sir, what
about
TensorFlow?**

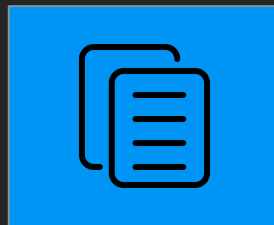
TensorFlow deals with all kinds of data



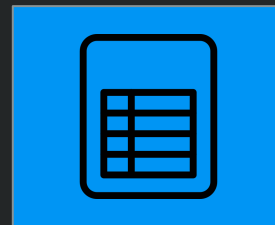
Audio



Image



Text



Tabular



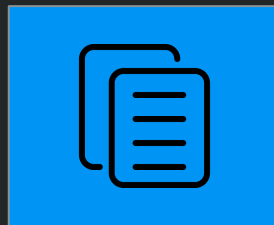
TensorFlow deals with main machine learning tasks



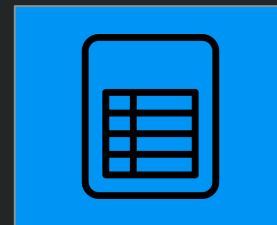
Audio



Image



Text

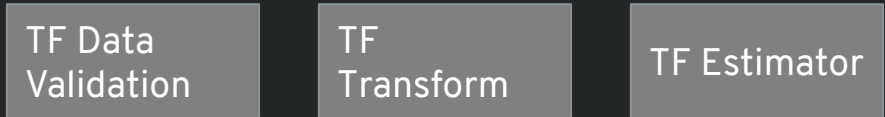


Tabular

Supervised Learning
Unsupervised Learning
Reinforcement Learning
Transfer Learning



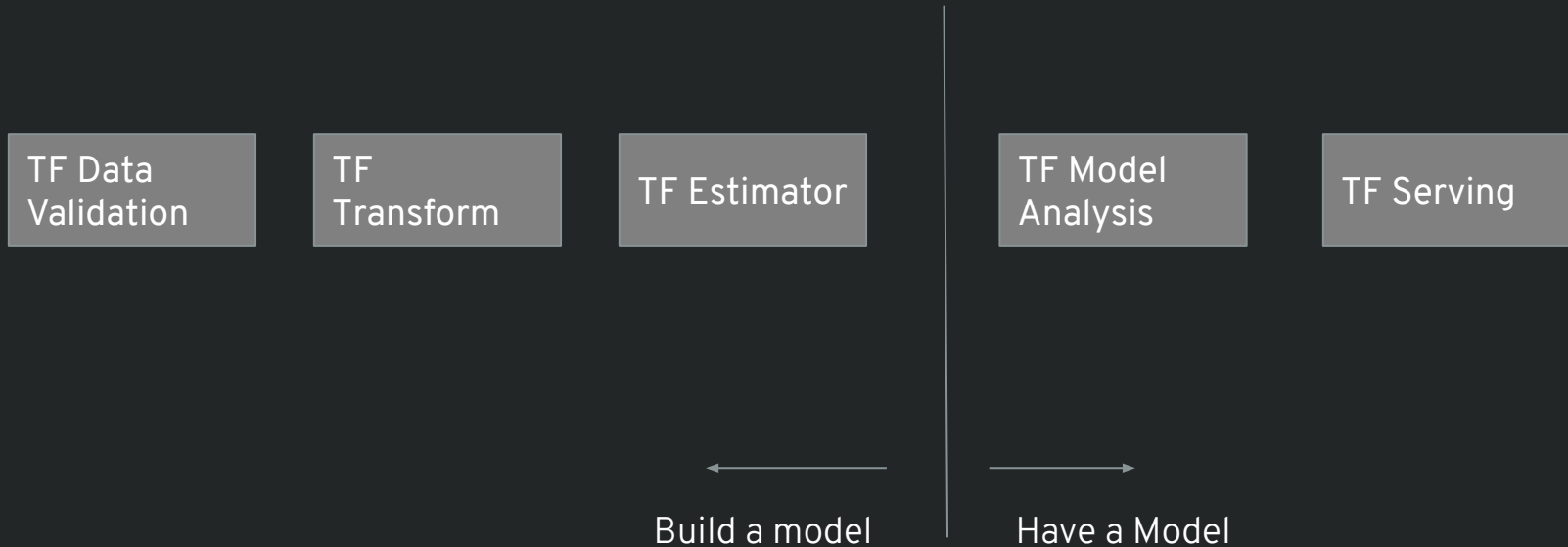
TensorFlow provides a end-to-end solution

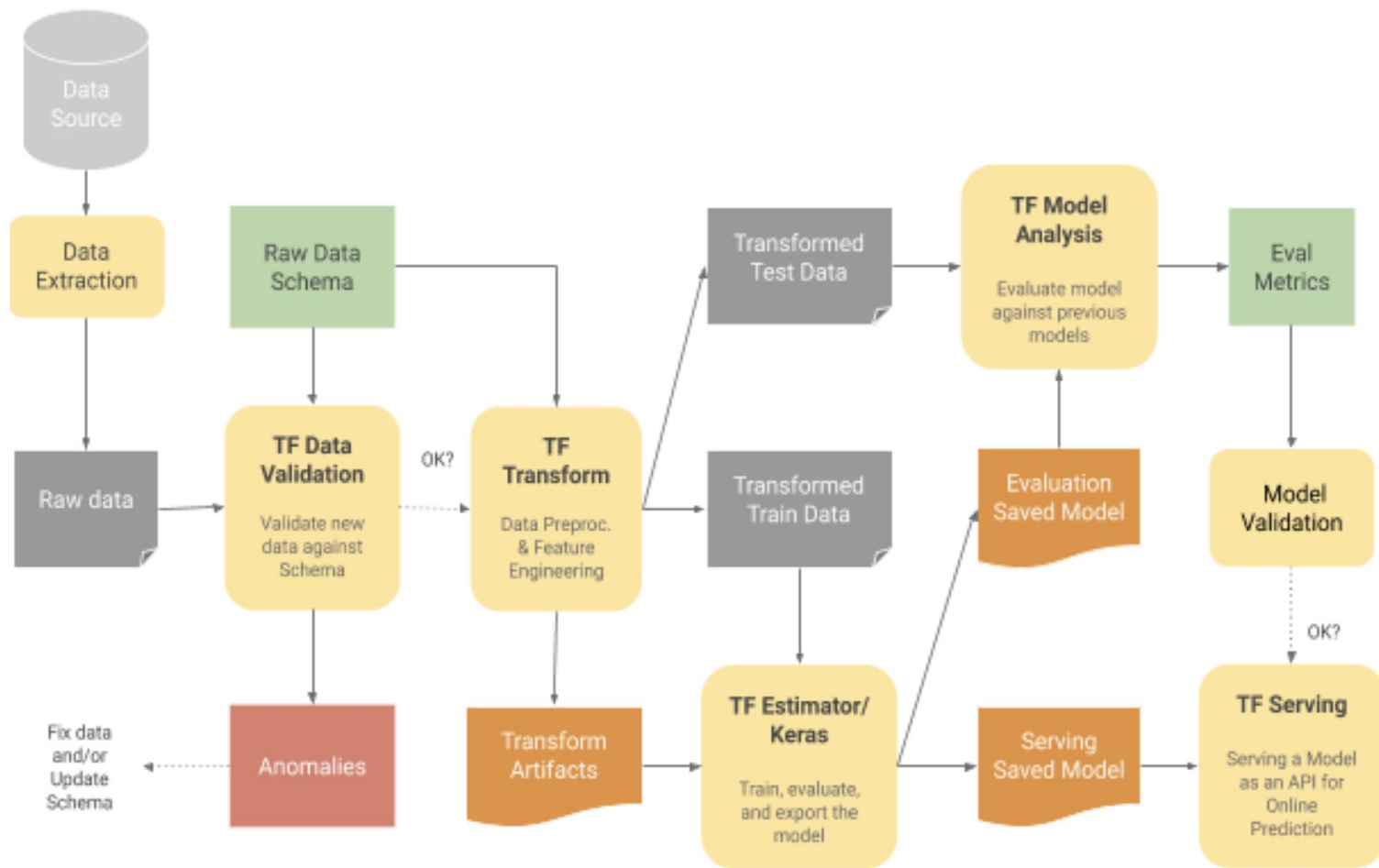


Build a model



TensorFlow provides a end-to-end solution





Typical TFX-based ML system



Getting Started with TensorFlow

What are tensors?

- Pattern of data in TensorFlow
- N-dimensional array.
- Beyond the numeric representation, tensors carry informations as data types and array's shape.

A tensor is an N-dimensional array of data



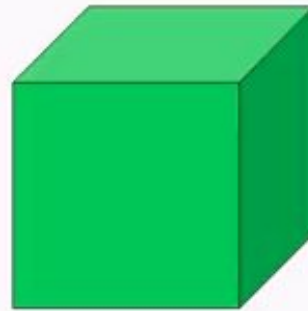
Rank 0
Tensor
scalar



Rank 1
Tensor
vector



Rank 2
Tensor
matrix



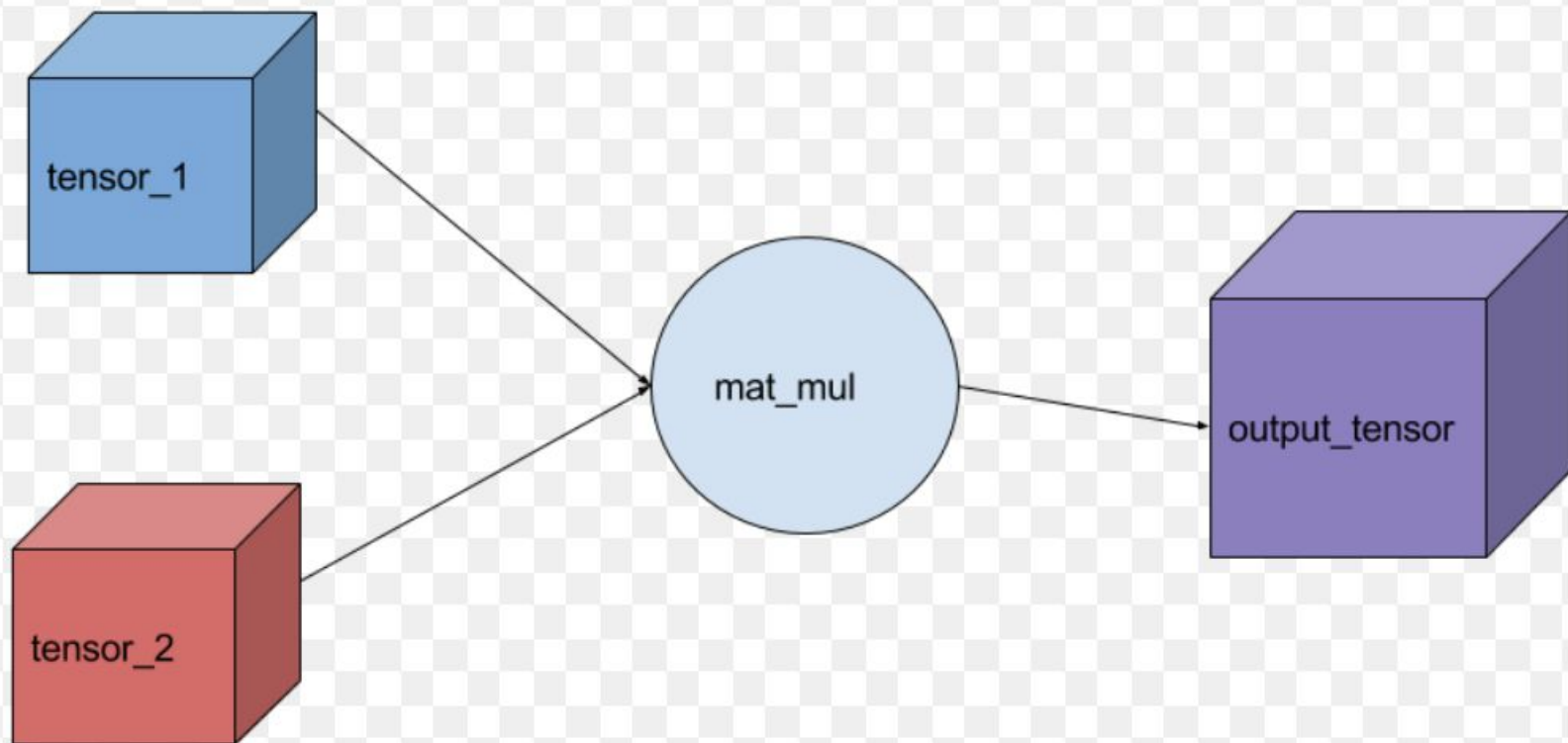
Rank 3
Tensor



Rank 4
Tensor

Graphs : Brief Introduction

- Math structure to compute relationship between objects. .
- Tensorflow generate graphs every time a operation is computed.



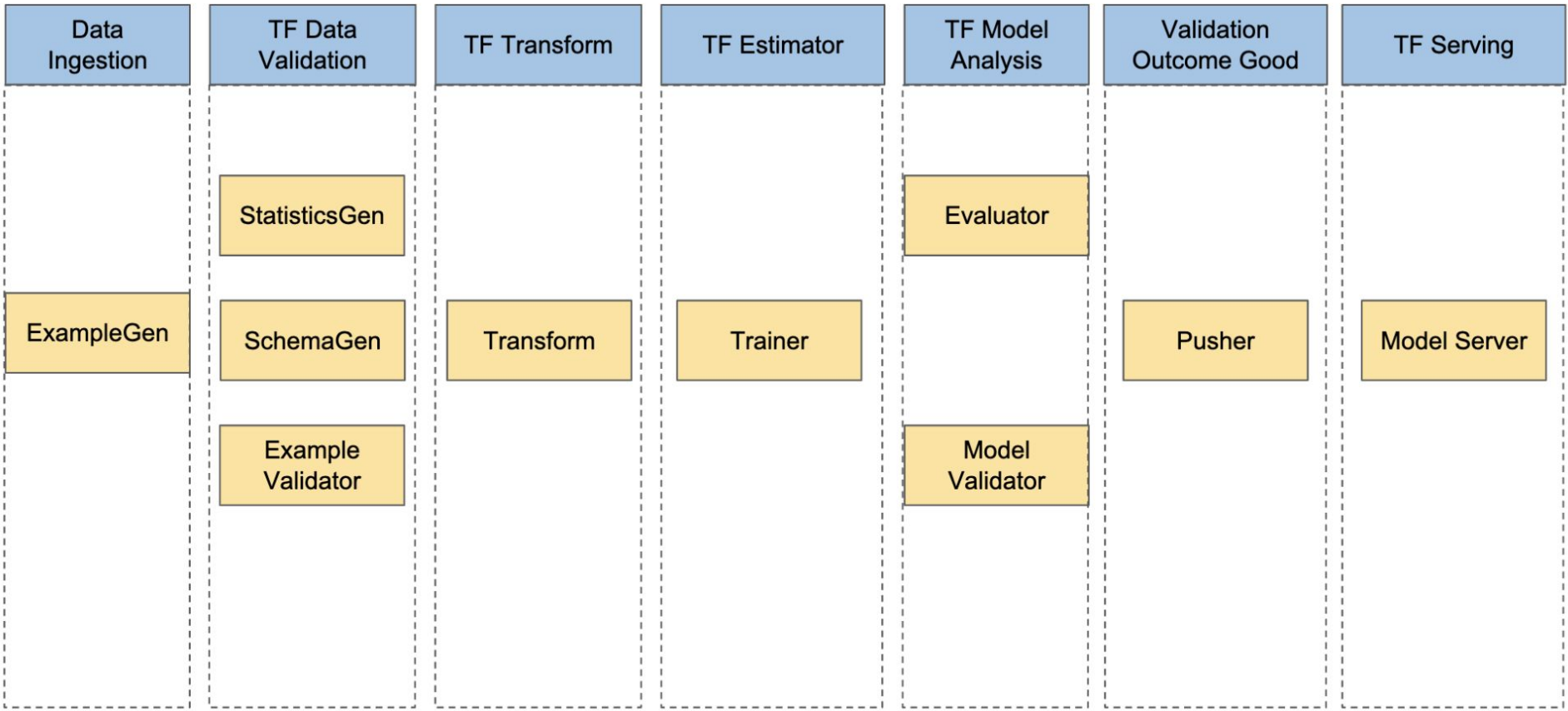
Code
(Link)



TensorFlow Extended

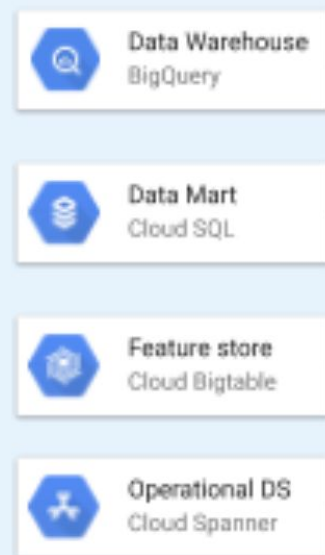
What is TensorFlow Extended?

- General purpose open source machine learning platform implemented by google.
- A set of glueable machine learning components in one platform simplifying machine learning deployment.
- Aims to reduce time to put machine learning in production and cover technical debts in those practices. .



Google Cloud Platform

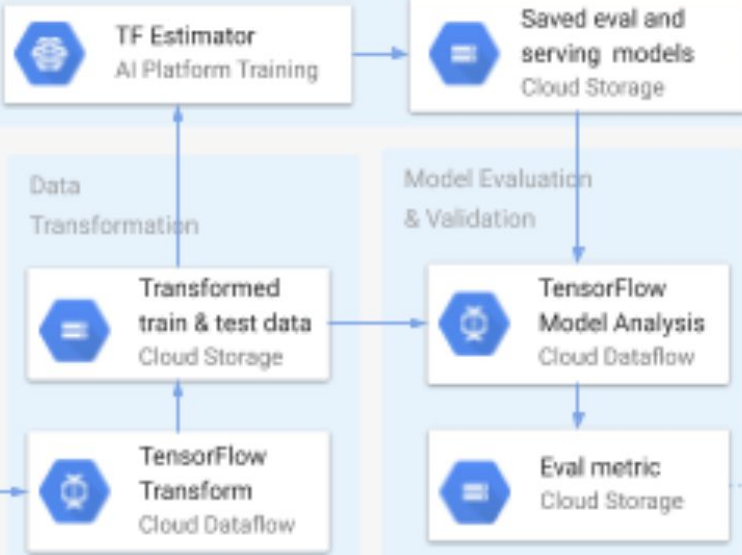
Data sources



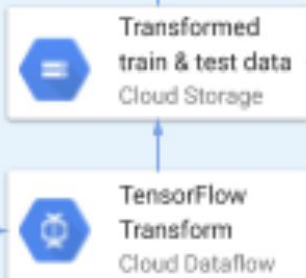
Data Extraction & Validation



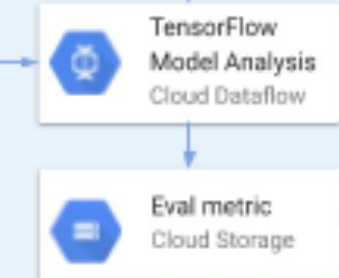
Model Training & Tuning



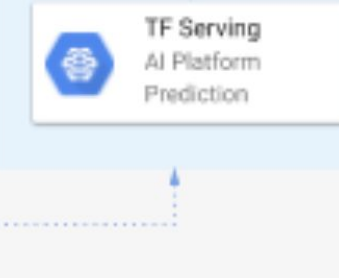
Data Transformation



Model Evaluation & Validation



Model Serving for Predictions



tf.data_validation

Real World Challenges

Challenges

- Schema Validation
- Data contains anomalies
- Label invalid data
- Change feature skew

How TFDV can help?

- `tfdv.schema_validation()`
- `tfdv.display_anomalies()`
- `tfdv.display_anomalies()`
- `tfdv.visualize_statistics()`

Better Data, Better Models	Automated schema generation	Integration with Facets (live demo)	Anomaly detection
Compute summary statistics for train/test data	Input feature value ranges	Identify train/test/validation set skew	Detect missing values
Drift detection	Type imputation	Unexpected feature values	Out of range values
Training-Serving skew detection	Environment specific schema	Feature by feature analysis	Wrong feature types
	Schema validation	Compare statistics across two or more data sets	Correct non-conforming data
		Supports visualization of large datasets responsively	

Code
(Link)

tf.transform

About Transform

- Framework for running batch and streaming data processing jobs.
- Uses Apache Beam.
- Easily integrated with Kubeflow. [\(example\)](#)

Feature Engineering @ Scale	Transformations
Transform data before it goes into a model. Output is exported as a TensorFlow graph, used for both training & serving.	<code>tft.scale_by_min_max()</code> , <code>tft.scale_to_0_1()</code> , <code>tft.scale_to_z_score()</code>
Create feature embeddings & enriching text features	<code>tft.tfidf()</code> , <code>tft.ngrams()</code> , <code>tft.hash_strings()</code>
Builds transformations into TF graph for your model, so same transformations are applied for train & serving.	Convert strings to integers by generating a vocabulary over all input values
Vocabulary generation	<code>tft.compute_and_apply_vocabulary()</code> , <code>tft.string_to_int()</code>
Normalize values & Bucketization	<code>tft.bucketize()</code>

Code
(Link)

tf.estimators

About Estimators

- Easily create Estimators using Keras API.
- Sklearn compatible syntax.
- Train models using CPU/ GPU / TPU.
- Perform distributed training.
- Save Summaries on TensorBoard.
- Model checkpoint and failure recover.

Code
(Link)

tf.model_analysis

About Model Analysis

- Evaluate model performance on large datasets.
- Check performance choosing different metrics and slices of data.
- Track Performance over the time.
- Friendly Visualization tool.

tf.serving

About Serving

- Low Latency,
- Scalable horizontally.
- Reliable & Robust.
- Load/ host multiple versions of a model.
- Built in A/B testing.

5

Qwiklabs
link

Thank you!

Questions?

Tulio Vieira de Souza
tulio.souza@avenuecode.com
Data Science Engineer

