



Michelangelo

Jeremy Hermann, Machine Learning Platform @ Uber

ML at Uber

Uber Data



75 Million Riders
and 3 million drivers



4 Billion Trips
completed worldwide in
2017



65 Countries
and 600+ cities worldwide



15 Million Trips
completed each day

Uber Data

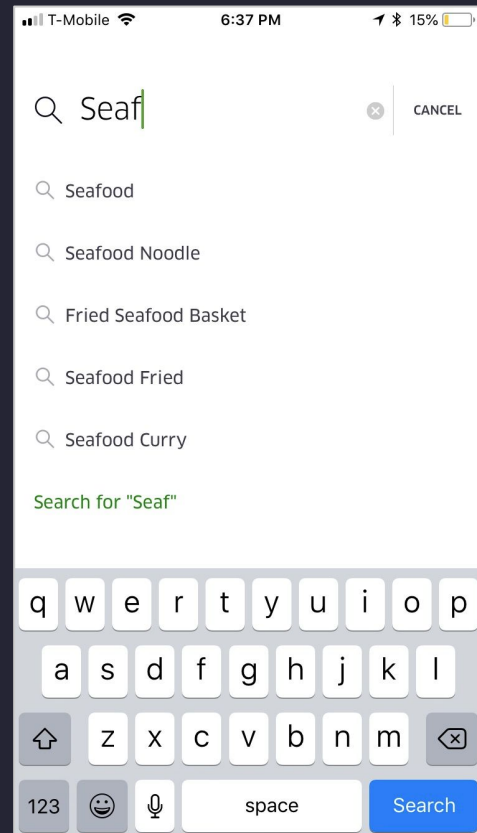
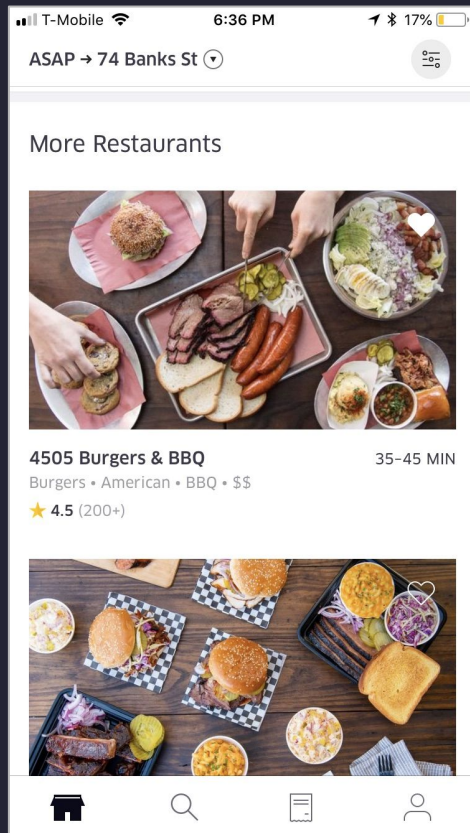


ML at Uber

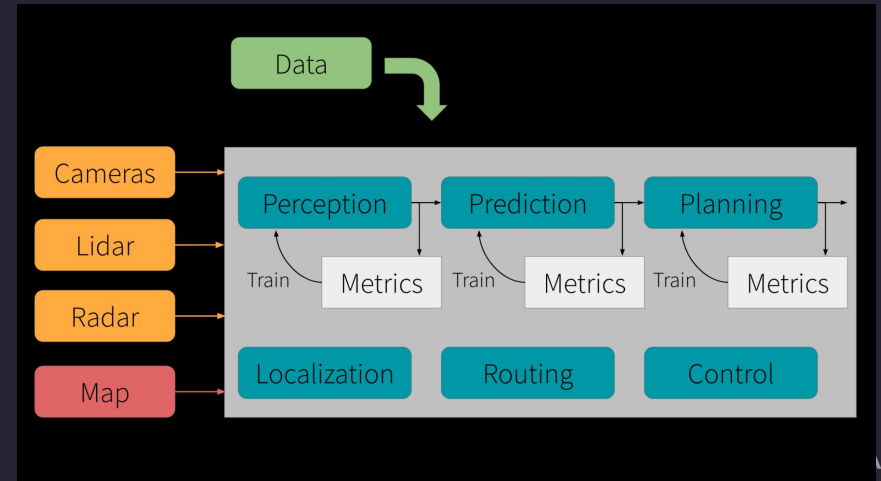
- Uber Eats
- ETAs
- Self-Driving Cars
- Customer Support
- Dispatch
- Personalization
- Demand Modeling
- Dynamic Pricing
- Forecasting
- Maps
- Fraud
- Safety
- Destination Predictions
- Anomaly Detection
- Capacity Planning
- And many more...

ML at Uber - Eats

- Models used for
 - Ranking of restaurants and dishes
 - Delivery times
 - Search ranking
- 100s of ML models called to render Eats homepage

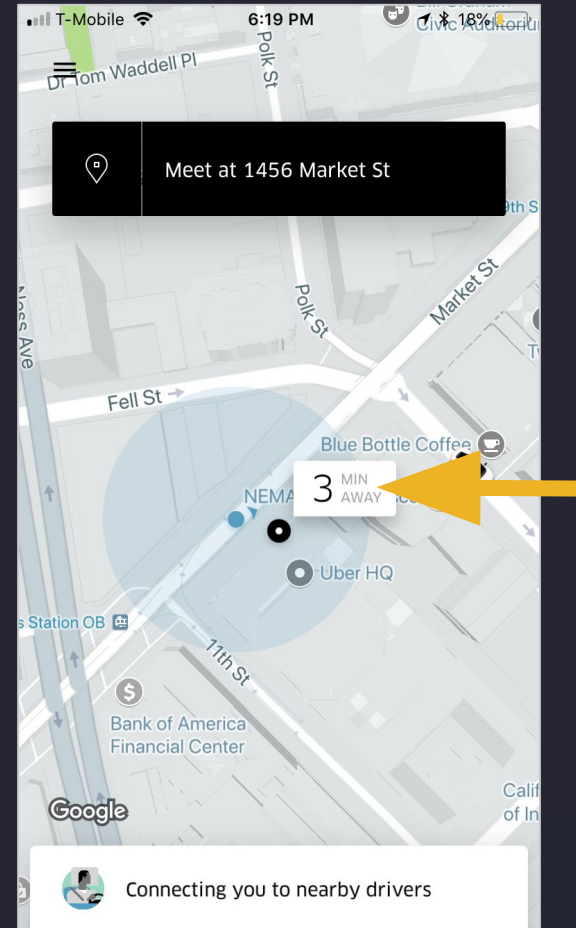


ML at Uber - Self-Driving Cars



ML at Uber - ETAs

- ETAs are core to customer experience
- ETAs used by myriad internal systems
- ETAs are generated by route-based algorithm called Garafu
- ML model predicts the Garafu error
- Use the predicted error to correct the ETA
- ETAs now dramatically more accurate



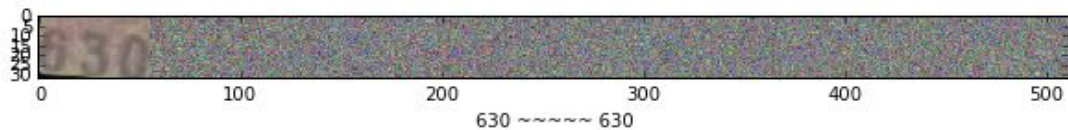
ML at Uber - Map Making



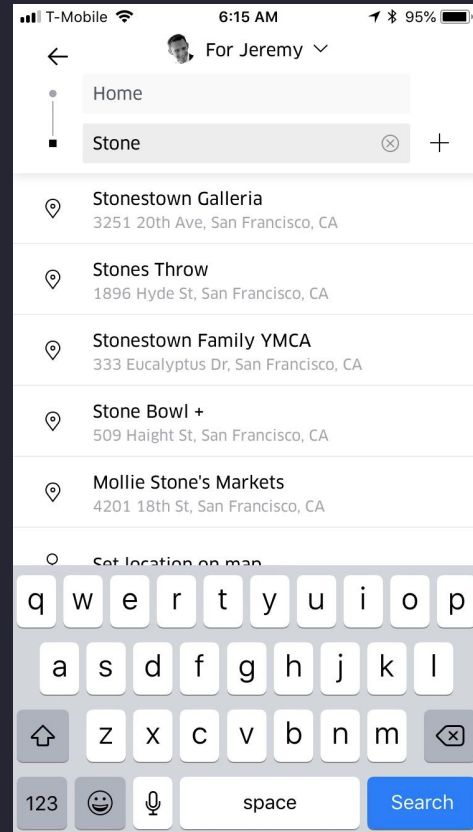
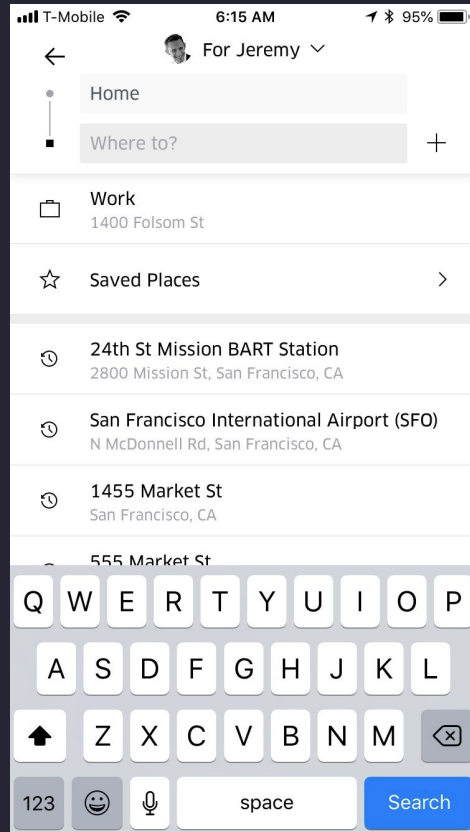
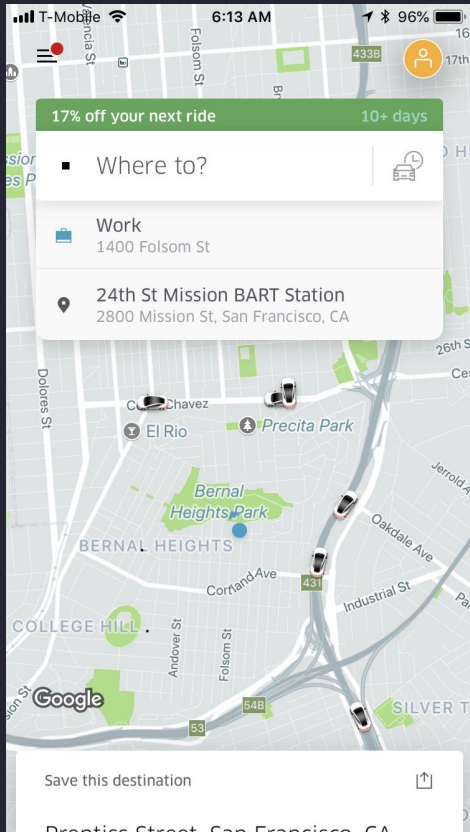
ML at Uber - Map Making



ML at Uber - Map Making



ML at Uber - Destination Prediction



ML at Uber - Spatiotemporal Forecasting

Supply

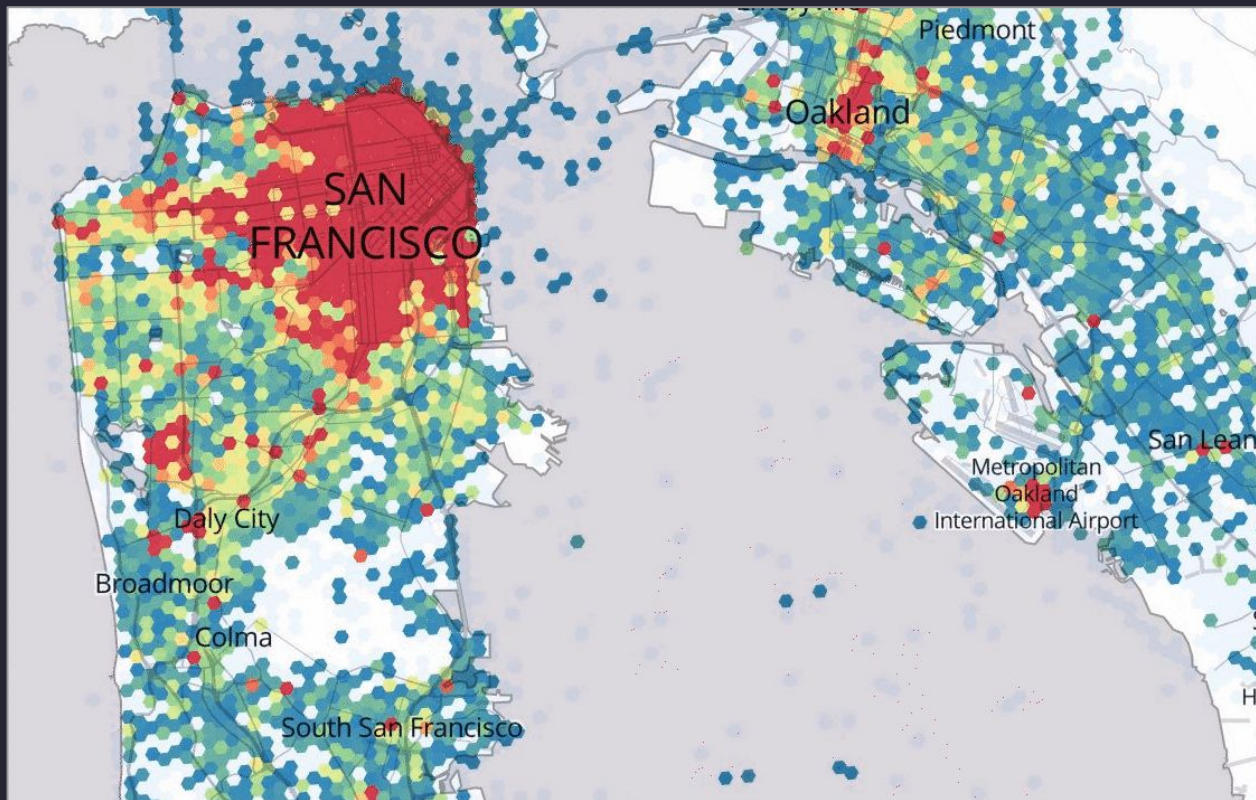
- Available Drivers

Demand

- Open Apps

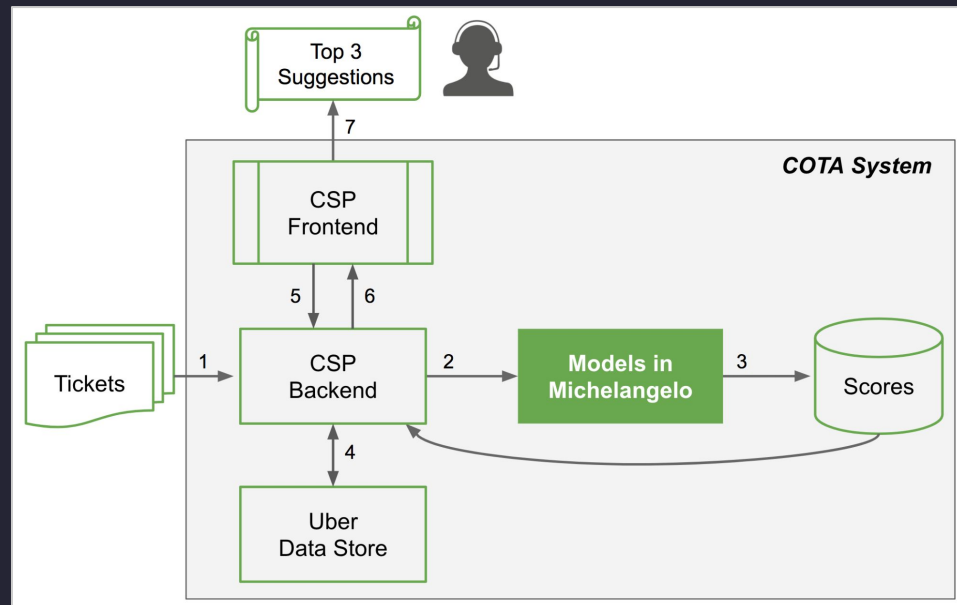
Other

- Request Times
- Arrival Times
- Airport Demand



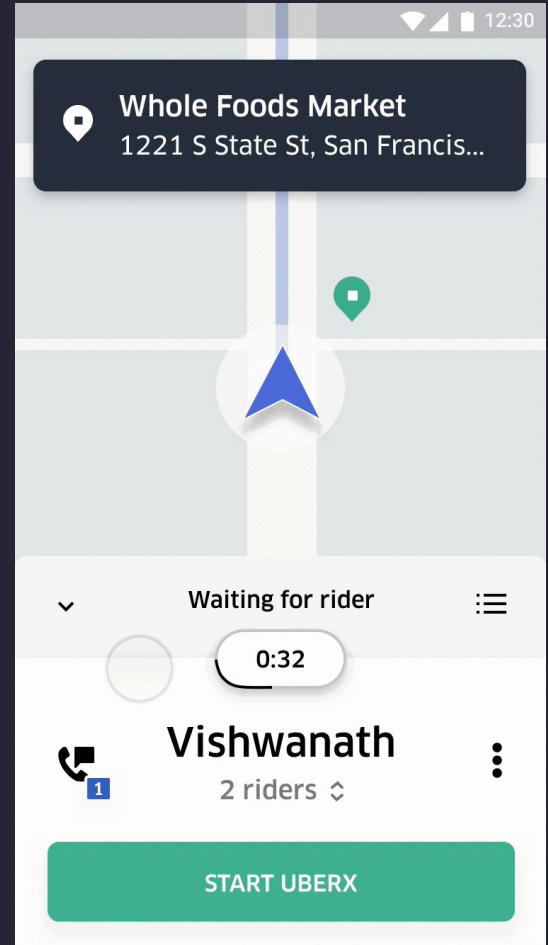
ML at Uber - Customer Support

- 5 customer-agent communication channels
- Hundreds of thousands of tickets surfacing daily on the platform across 400+ cities
- NLP models classify tickets and suggest response templates
- Reduce ticket resolution time by 10%+ with same or higher CSAT



ML at Uber - One Click Chat

- It's important for riders and driver partners to be able to communicate efficiently during pickup
- The one click chat feature streamlines communication between riders and driver-partners
- Uses natural language processing (NLP) models that predict and display the most likely replies to in-app chat messages.

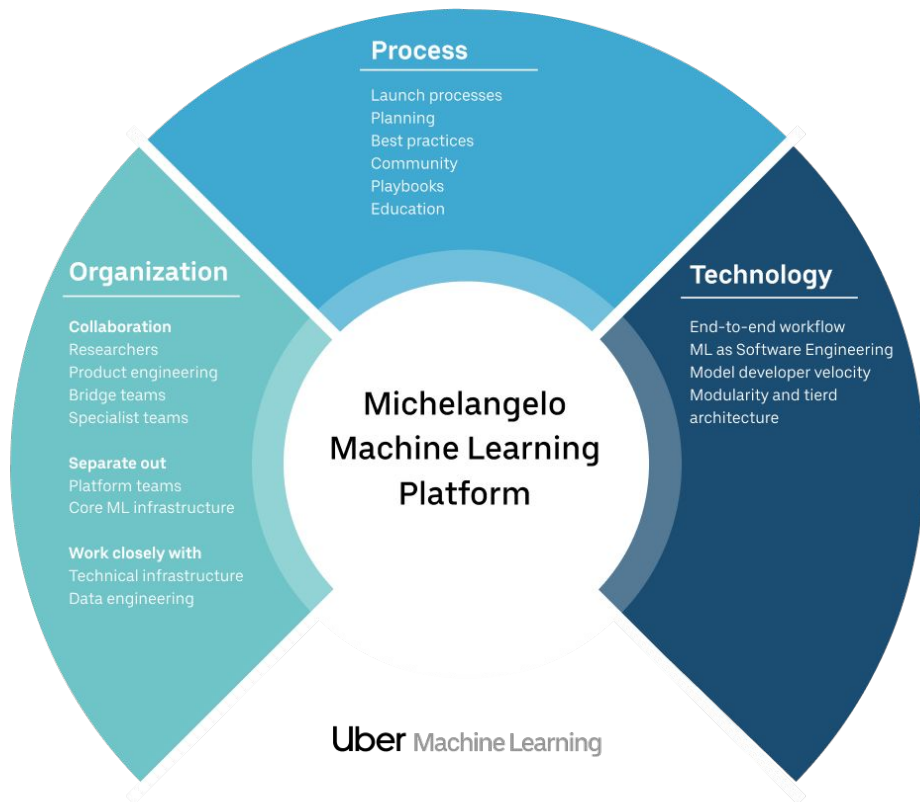


ML Platform

MISSION

Enable engineers and data scientists across the company to easily build and deploy machine learning solutions at scale.

ML Foundations - Organization, Process, and Technology



ML Platform Evolution

V1: Enable ML at Scale

- End-to-end workflow
- High scale training
- High scale model serving
- Feature Store



V2: Accelerate ML

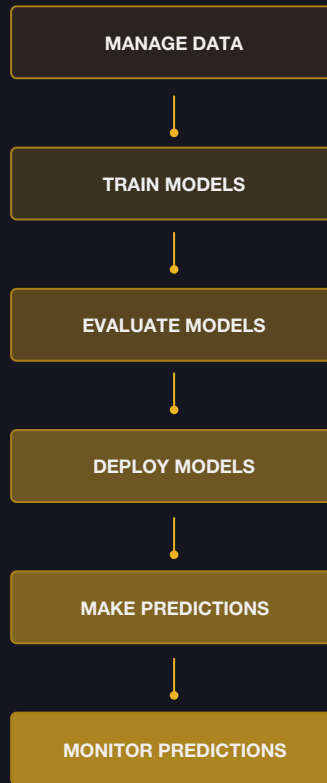
- PyML
- Horovod
- AutoTune
- Manifold Model Viz
- Realtime Model
Monitoring

Enable ML at Scale

Machine Learning Workflow

Same basic ML workflow & system requirements for

- Traditional ML & Deep Learning
- Supervised, unsupervised, & semi-supervised learning
- Online learning
- Batch, online, & mobile deployments
- Time-series forecasting



Enable ML at Scale:

Manage Data

Feature Store (aka Palette)

Problem

- Hardest part of ML is finding good features
- Same features are often used by different models built by different teams

Solution

- Centralized feature store for collecting and sharing features
- Platform team curates core set of widely applicable features
- Modellers contribute more features as part of ongoing model building
- Meta-data for each feature to track ownership, how computed, where used, etc
- Modellers select features by name & join key. Offline & online pipelines are automatically deployed

Enable ML at Scale:

Train Models

Distributed Training of Non-DL Models

Large-scale distributed training (billions of samples)

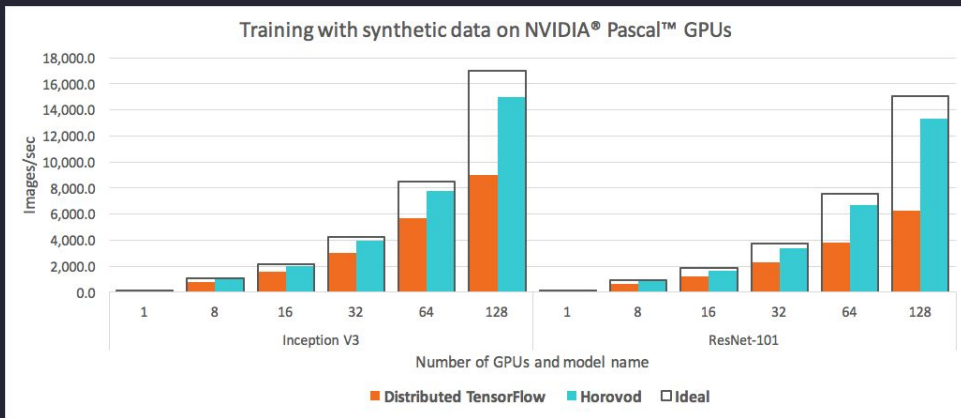
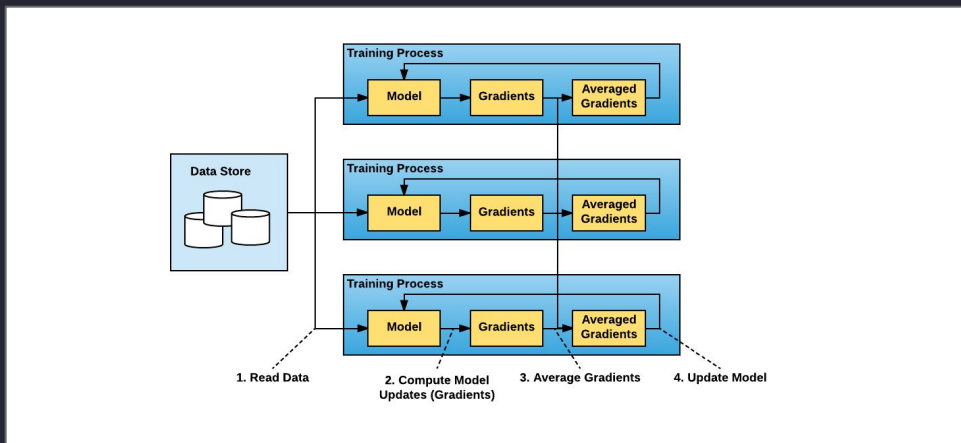
- Decision trees
- Linear and logistic models
- Unsupervised learning
- Time series forecasting
- Hyperparameter search for all model types

Speed and reliability

- Fuse operators into single job for speed
- Break operators into separate jobs to reliability

Distributed Training of Deep Learning Models with Horovod

- Data-parallelism works best when model is small enough to fit on each GPU
- Ring-allreduce is more efficient than parameter servers for averaging weights
- Faster training and better GPU utilization
- Much simpler training scripts
- More details at <http://eng.uber.com/horovod>



Enable ML at Scale:

Manage & Eval Models

Evaluate Models

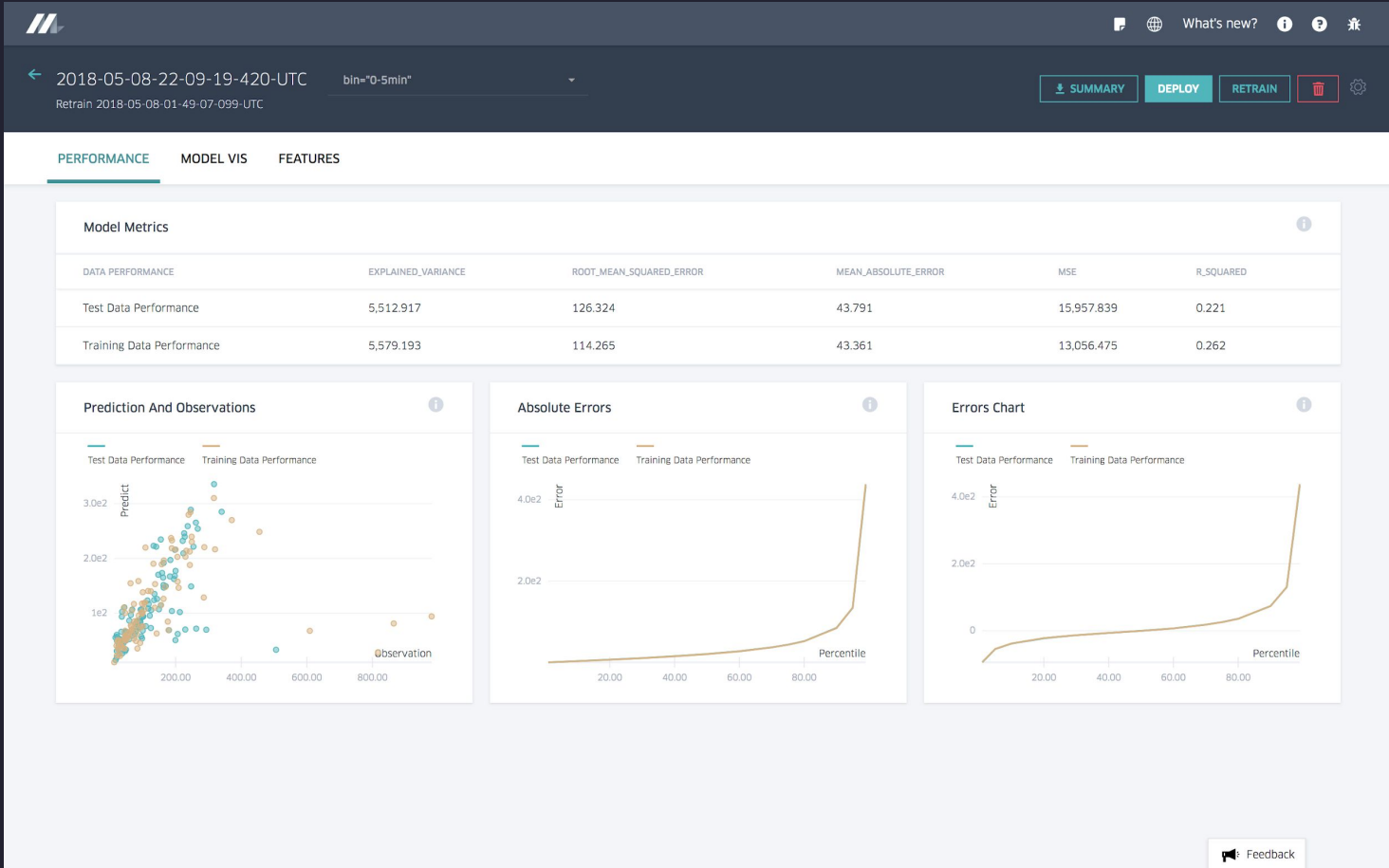
Problem

- It takes many iterations to produce a good model
- Keeping track of how a model was built is important
- Evaluating and comparing models is hard

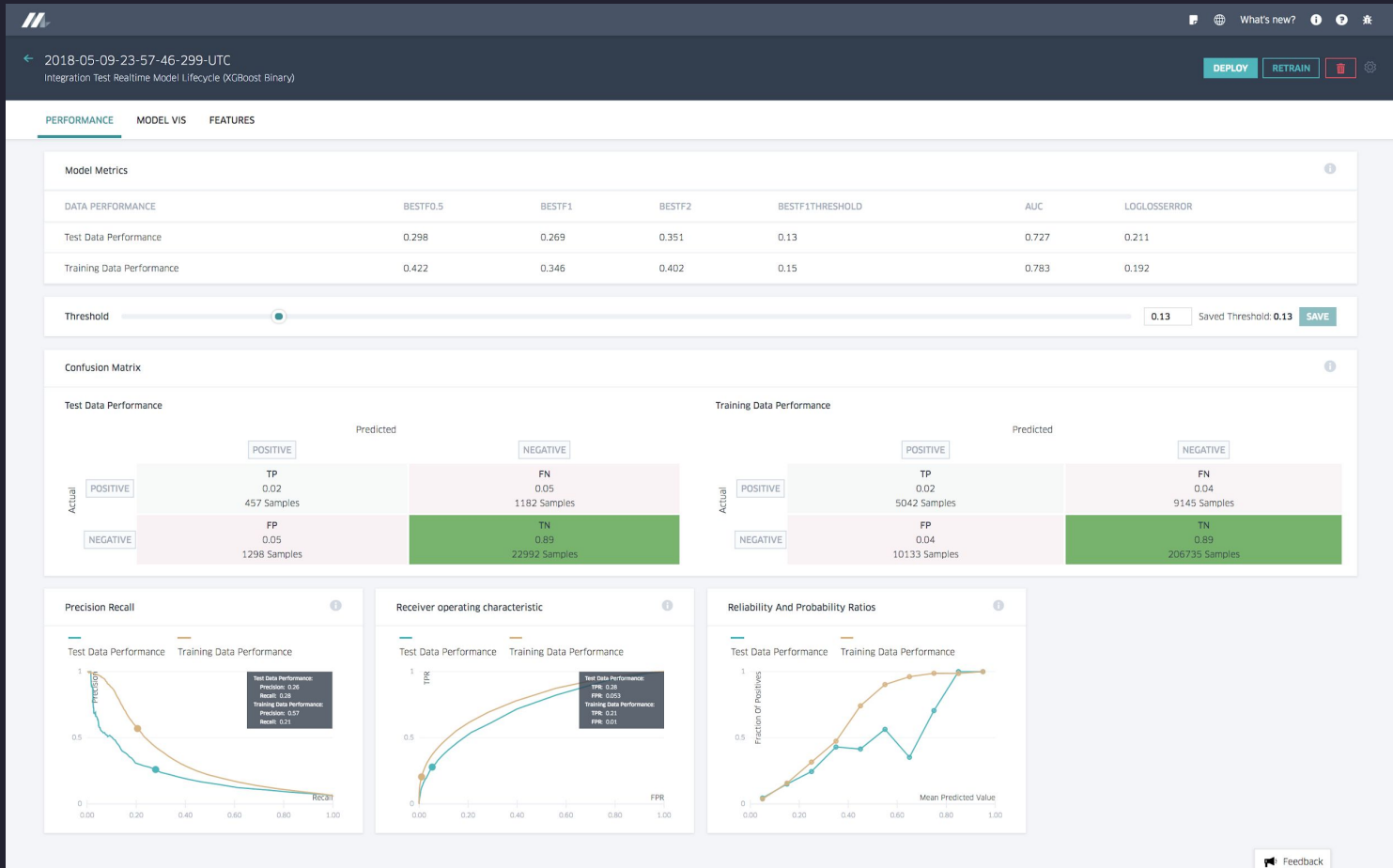
With every trained model, we capture standard metadata and reports

- Full model configuration, including train and test datasets
- Training job metrics
- Model accuracy metrics
- Performance of model after deployment

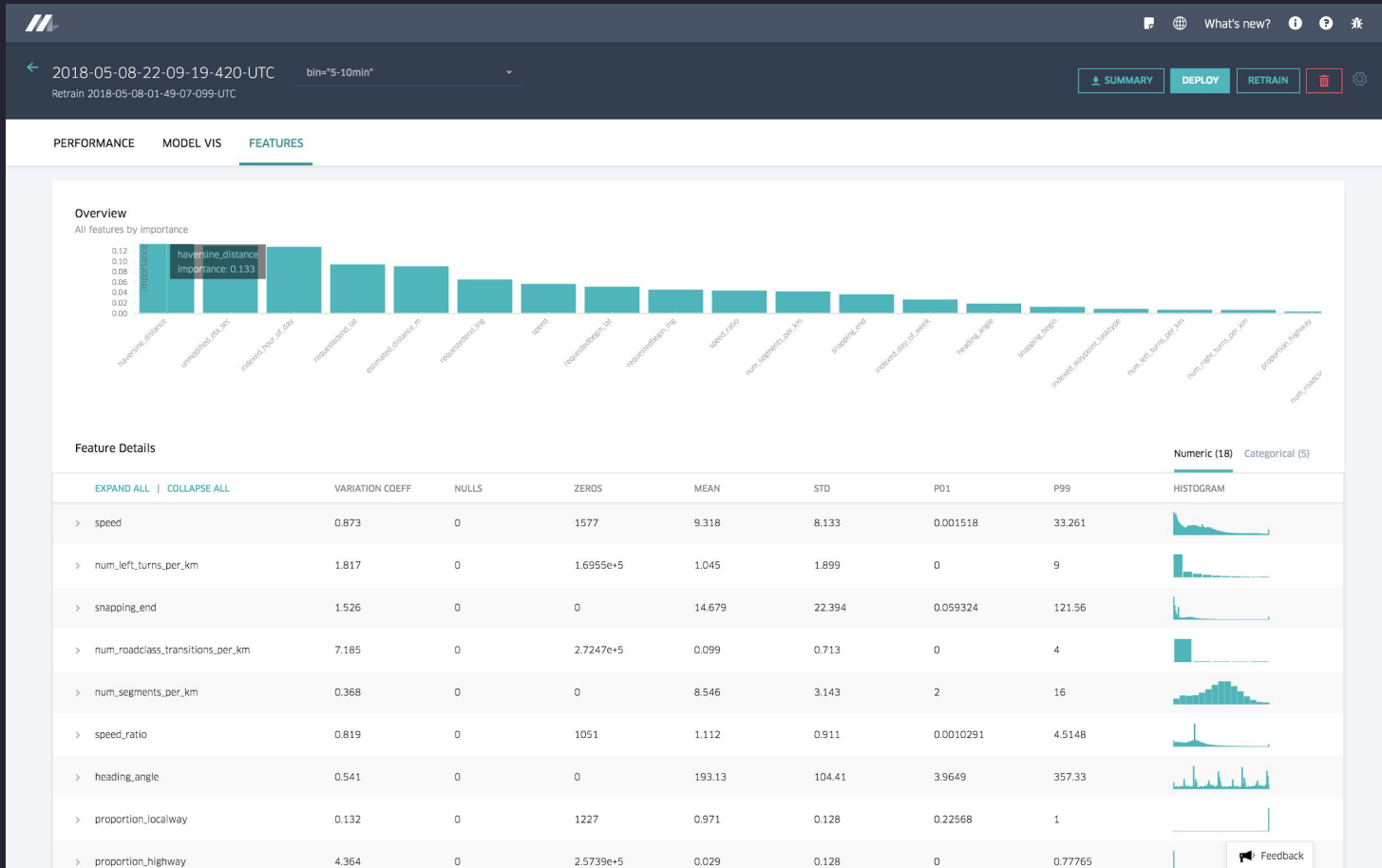
Model Visualization - Regression Model



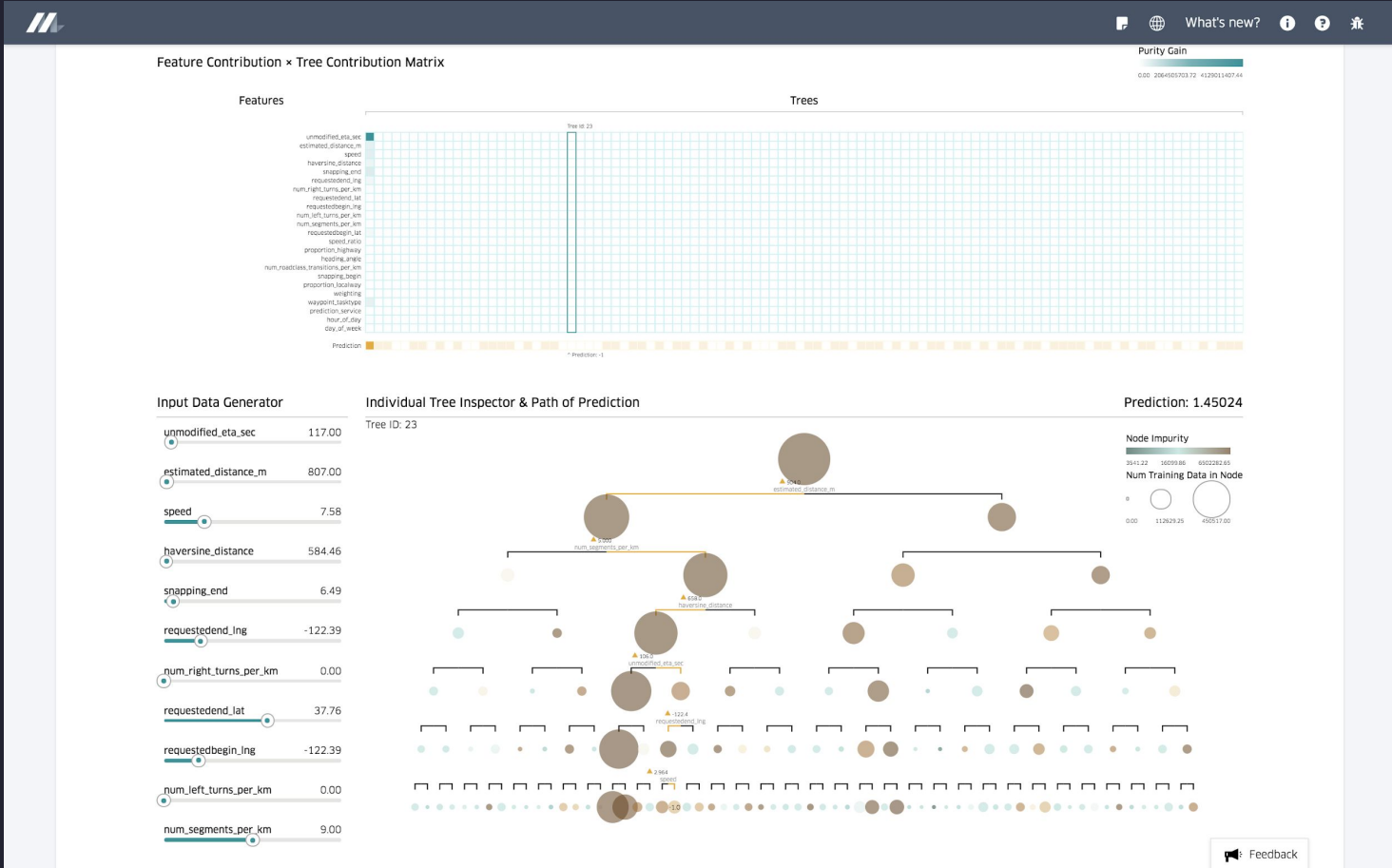
Model Visualization - Classification Model



Model Visualization - Feature Report



Model Visualization - Decision Tree



Enable ML at Scale:

Deployment & Serving

Online Prediction Service

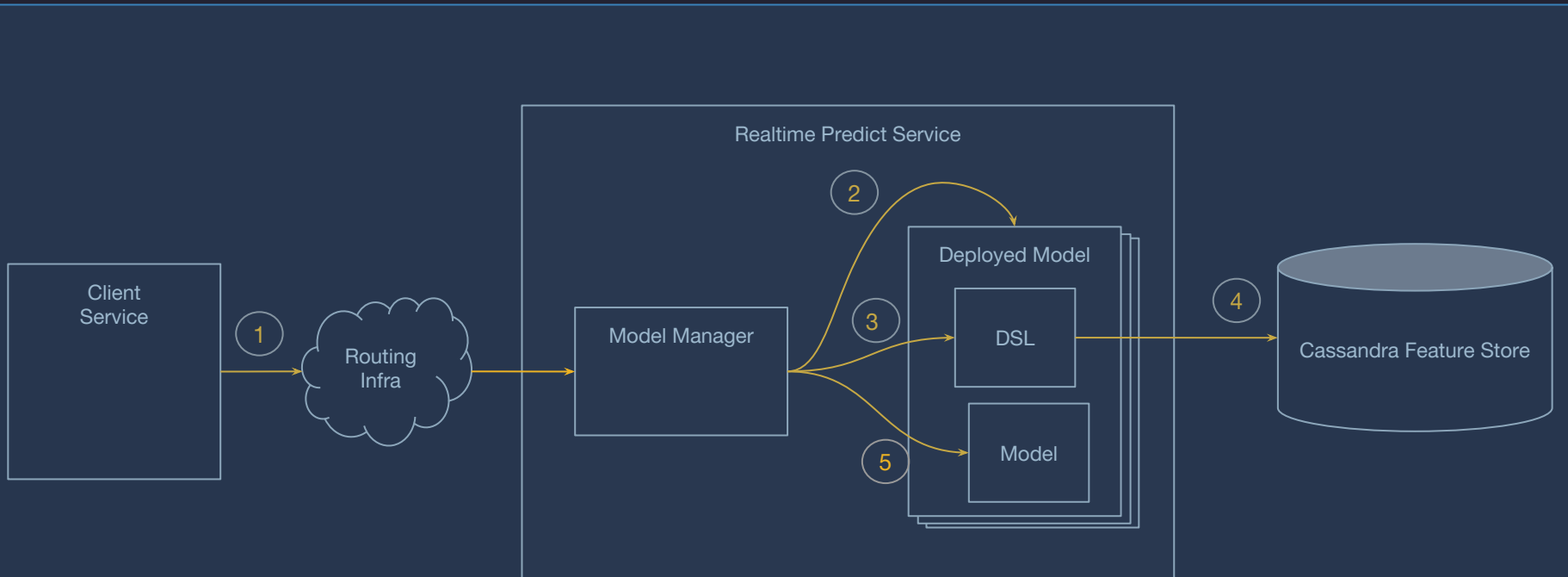
Prediction Service

- Thrift service container for one or more models
- Scale out in Docker on Mesos
- Single or multi-tenant deployments
- Connection management and batched / parallelized queries to Cassandra
- Monitoring & alerting

Deployment

- Model & DSL packaged as JAR file
- One click deploy across DCs via standard Uber deployment infrastructure
- Health checks and rollback

Online Prediction Service



Online Prediction Service

Typical p95 latency from client service

- ~5ms when all features from client service
- ~10ms when joining pre-computed features from Cassandra

Peak prediction volume across current online deployments

- 1M+ QPS

Enable ML at Scale:

Monitor Models in Production

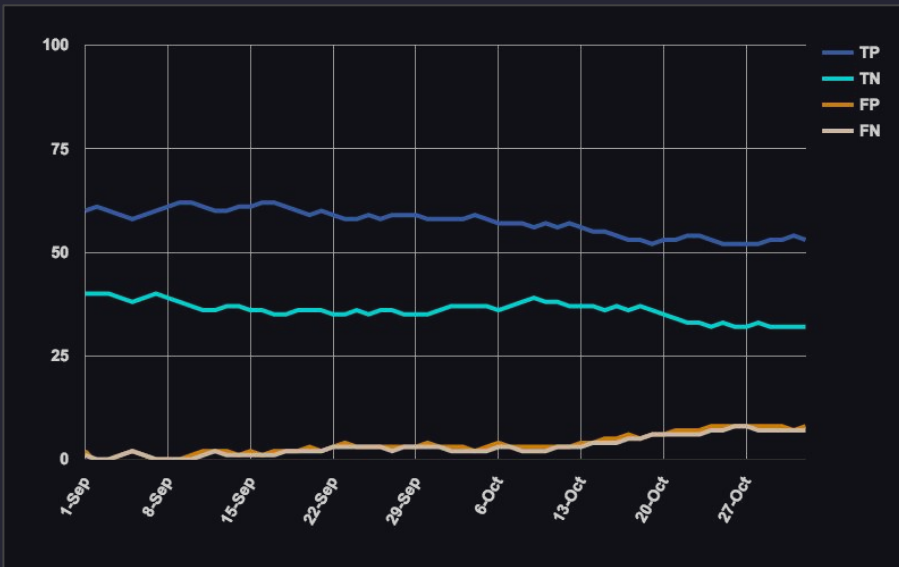
Monitor Predictions

Problem

- Models trained and evaluated against historical data
- Need to ensure deployed model is making good predictions going forward

Solution

- Log predictions & join to actual outcomes
- Publish metrics feature and prediction distributions over time
- Dashboards and alerts



Enable ML at Scale:

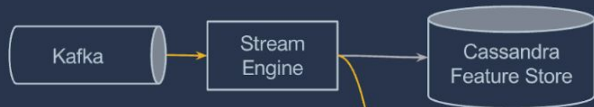
System Architecture

GET DATA

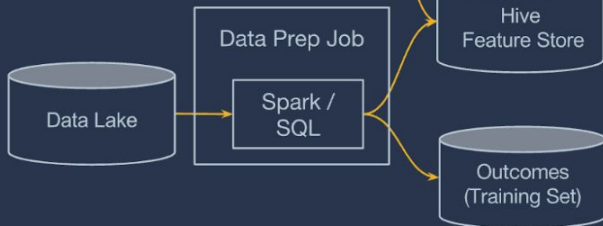
TRAIN MODELS

EVAL MODELS

DEPLOY, PREDICT & MONITOR



ONLINE



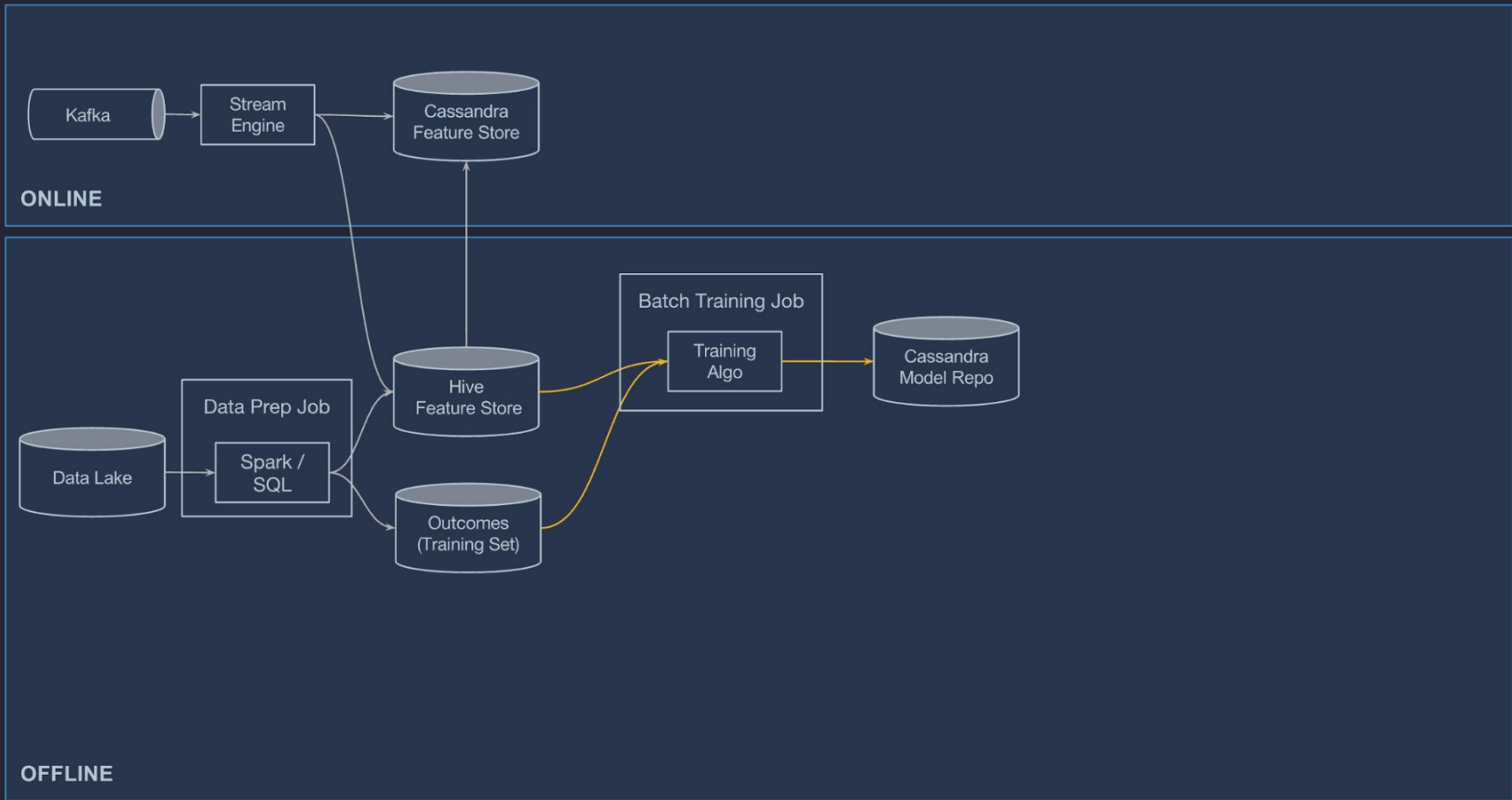
OFFLINE

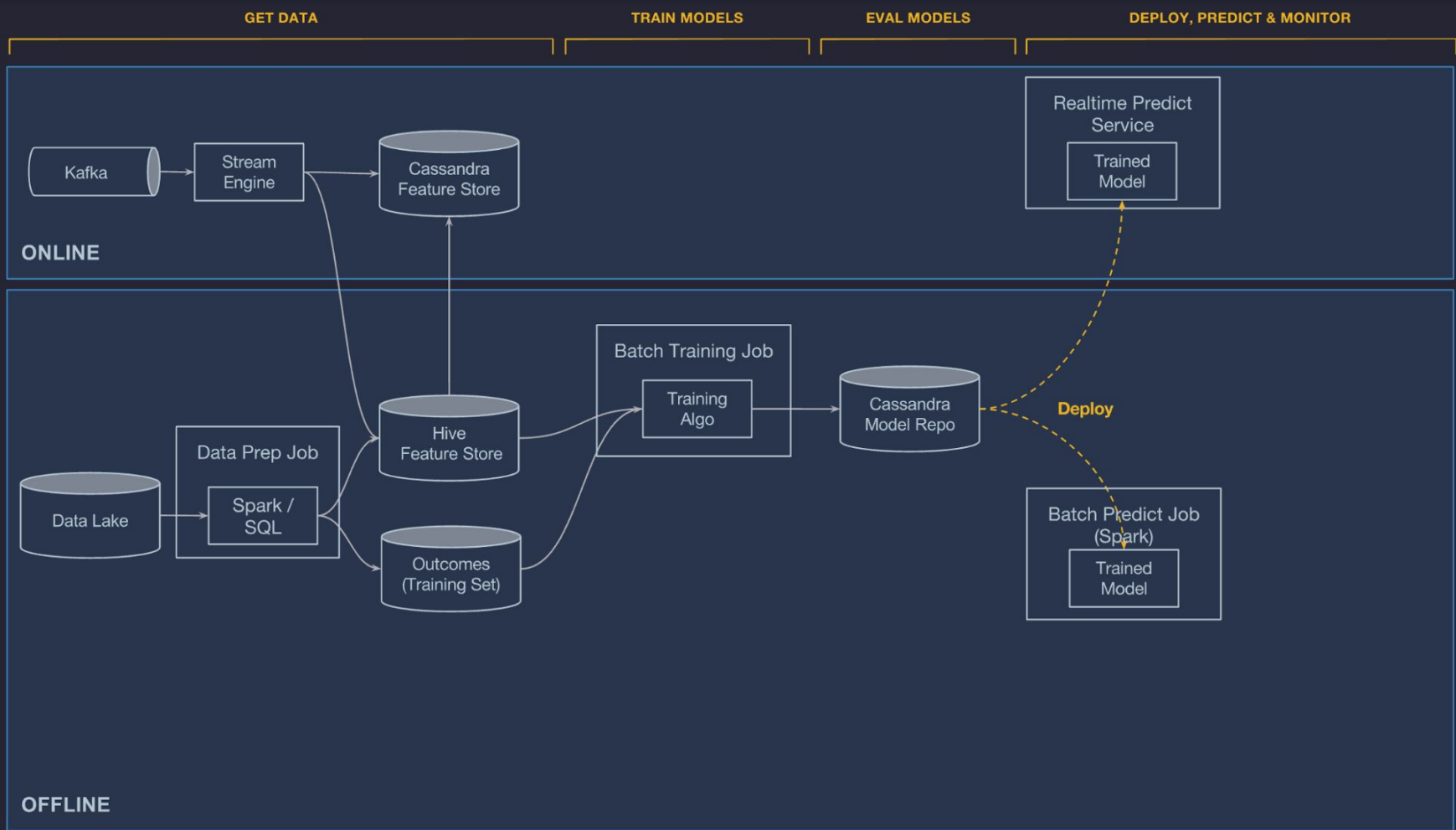
GET DATA

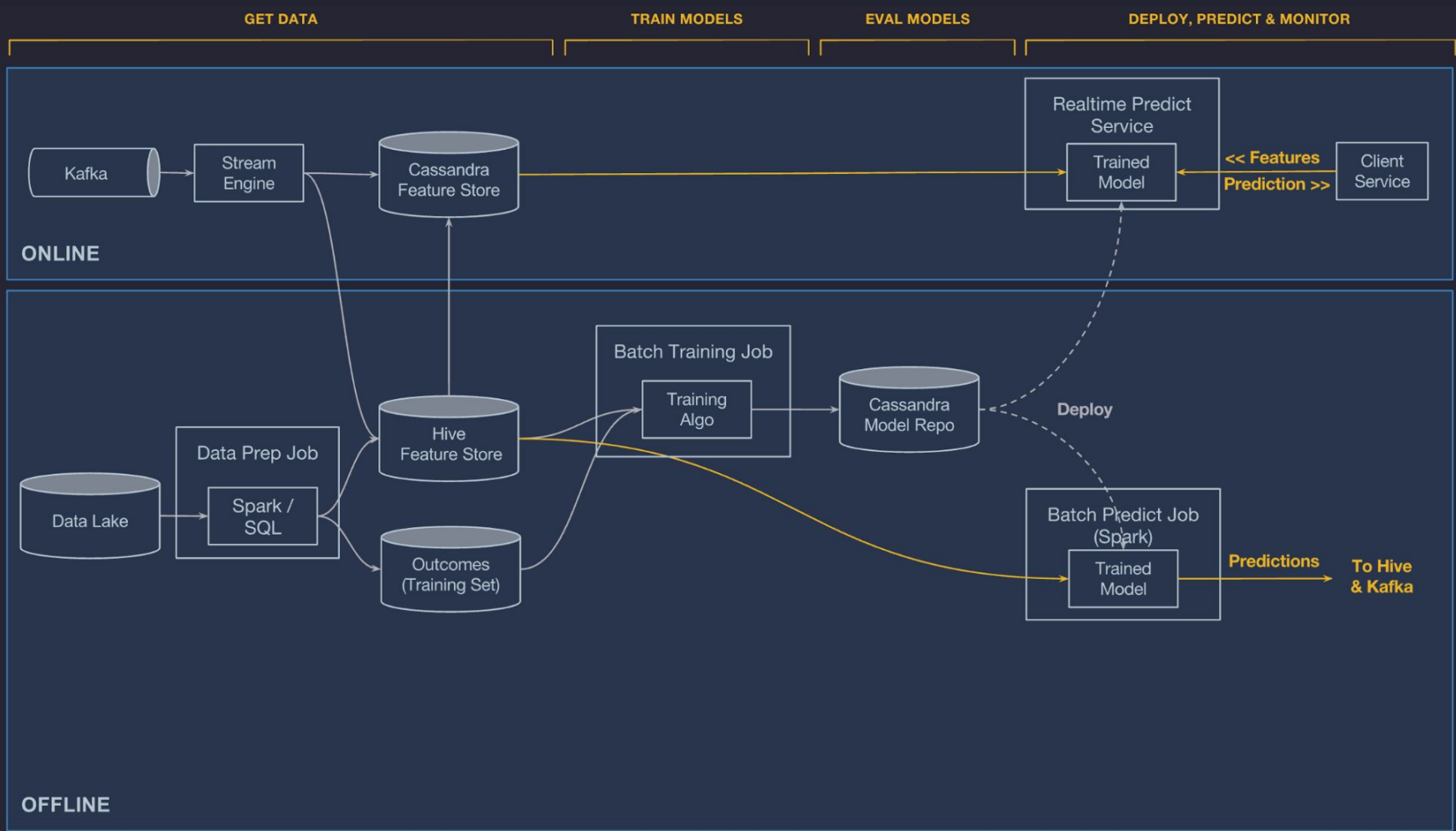
TRAIN MODELS

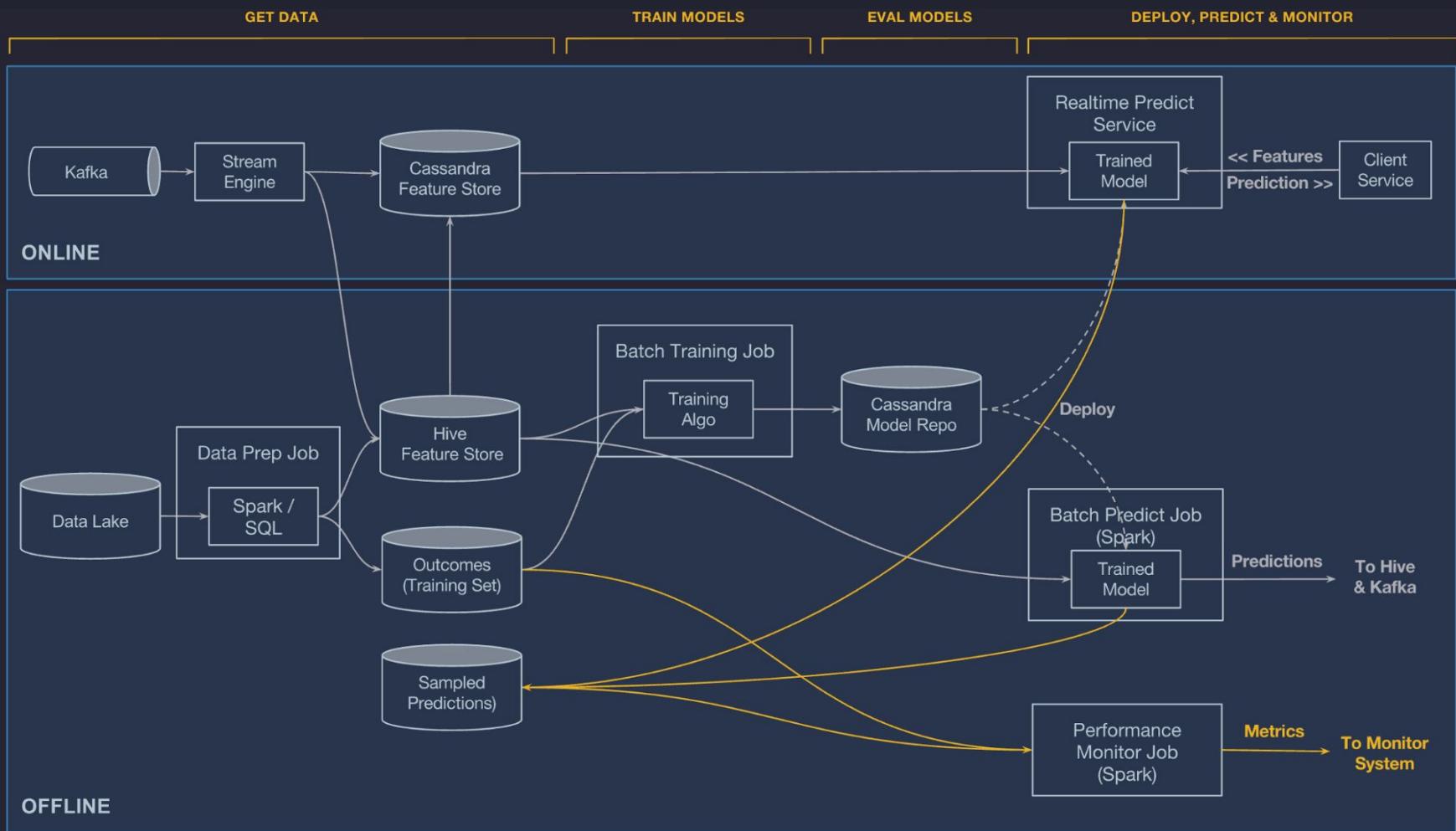
EVAL MODELS

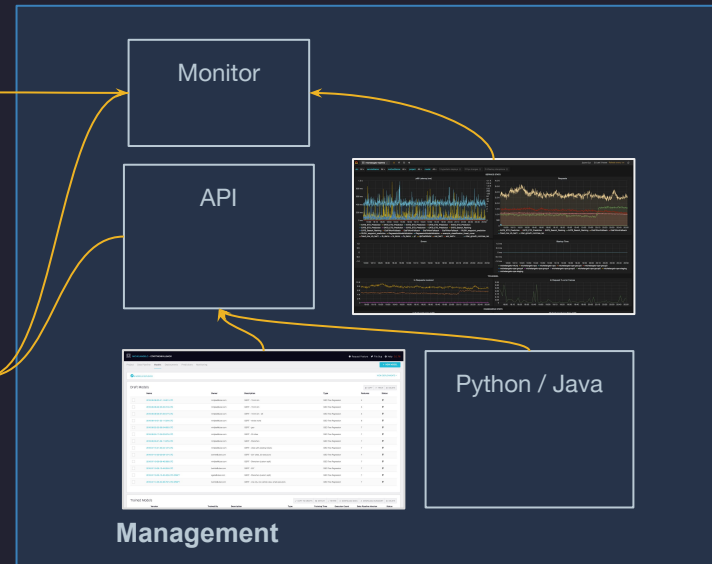
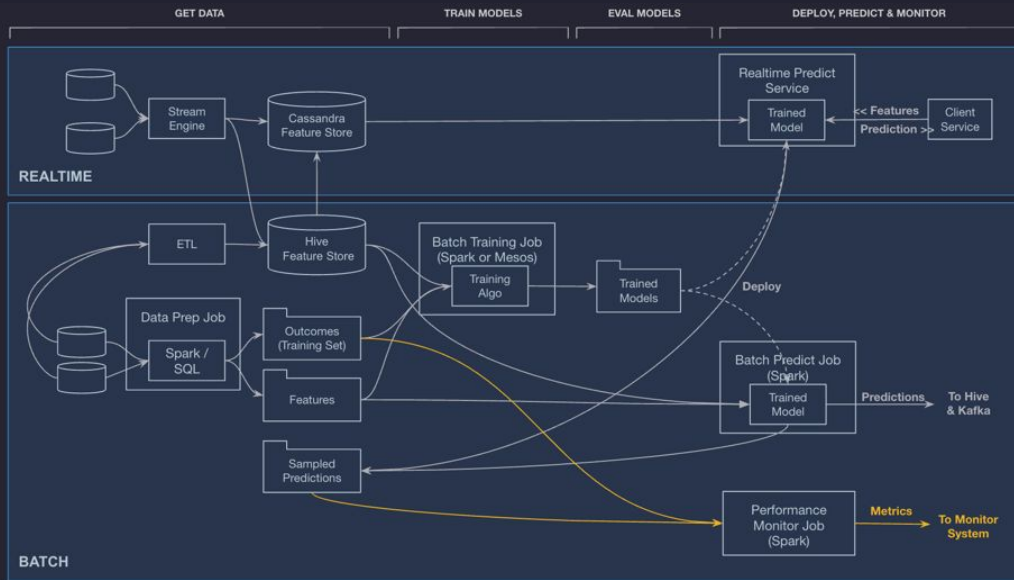
DEPLOY, PREDICT & MONITOR













lm2d180725-174721-ELRYCYFA-KWCORR

An example classification model

Deploy

Retrain

Stop

PERFORMANCE

MODEL VIS

FEATURES

Model Vis

This visualization is generated from the trained model. Model visualizations allow you to see a prediction path for the dataset. To edit, drag the slider to change input feature values. The prediction will change according to your input. Click on the matrix to switch rows.

Feature Contribution + Tree Contribution Matrix

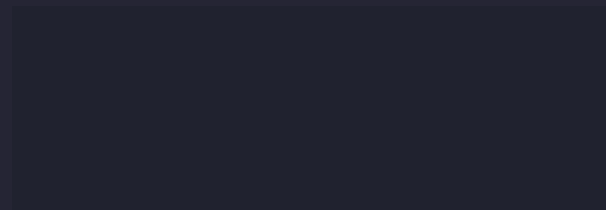
Policy 0.01

Features

Tree



Accelerate ML



ML Platform Evolution

V1: Enable ML at Scale

- End-to-end workflow
- High scale training
- High scale model serving
- Feature Store



V2: Accelerate ML

- PyML
- Horovod
- AutoTune
- Manifold Model Viz
- Realtime Model Monitoring

Keys to High Velocity ML

- Reduce friction at every step of complex, iterative workflow
- End-to-end ownership by modeler - no handoffs
- Bring the tools to the modeler
- Simple, elegant APIs
- Rich visual tools for understanding data and models
- Measure time from idea to model in production

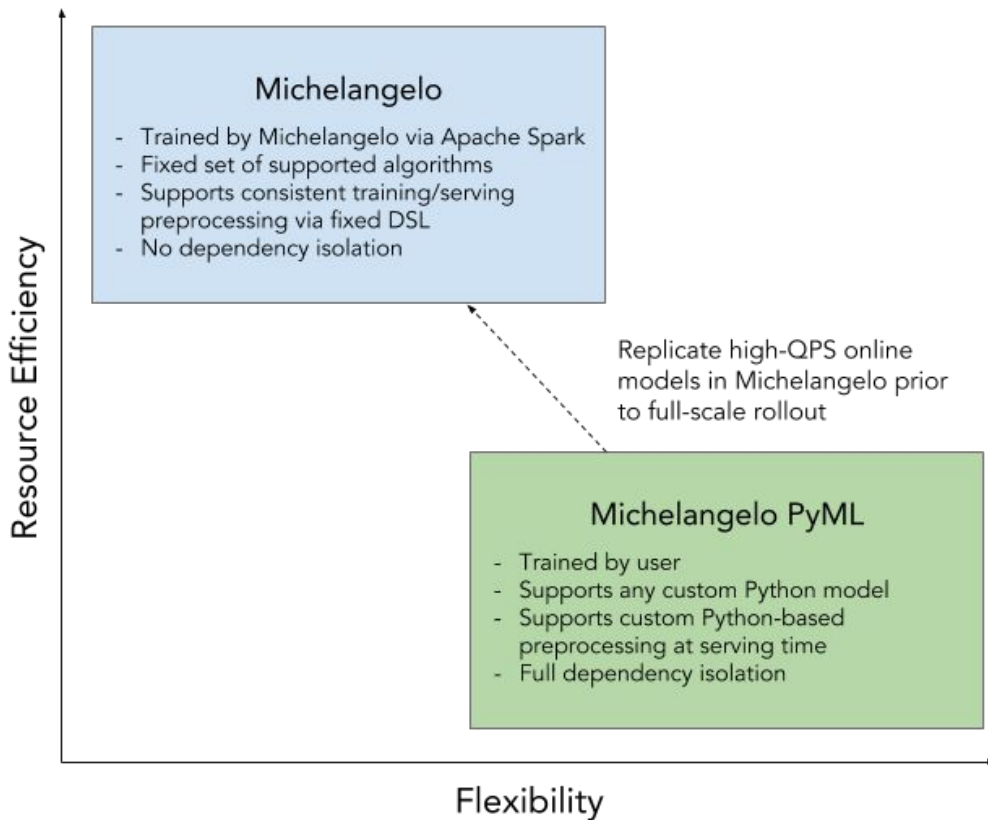
Accelerate ML:

PyML

PyML

- Problem
 - Michelangelo initially targeted high scale use cases
 - Good for first wave of production use cases
 - But, not very easy for early prototyping and limited flexibility
- Solution
 - Support regular Python for lower scale, easy to use and very flexible modeling
 - Import any libraries
 - Train any model
 - Implement serving interface with `predict()` method
 - Call API to package model artifacts and upload to Michelangelo
 - Deploy to production via API or Michelangelo UI

PyML



PyML - 1. Train and Test Model

```
import pandas as pd
import numpy as np
from sklearn.datasets import load_breast_cancer

# Prepare the dataset
dataset = load_breast_cancer()
feature_columns = [name.replace(' ', '_') for name in dataset.feature_names.tolist()]
pandas_df = pd.DataFrame(data= np.c_[dataset.data, dataset.target],
                        columns=feature_columns + ['target'])

# Train logistic regression
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(dataset.data,
                                                    dataset.target,
                                                    random_state=42)

log_reg = LogisticRegression()
log_reg.fit(X_train, y_train)

# Test that model works on first few records
log_reg.predict_proba(pandas_df[feature_columns][:10,0])
```

```
array([[ 1.          ,  0.99999995,  0.99999994,  0.41742828,  0.99998692,
         0.76040329,  0.99999709,  0.97970554,  0.90590637,  0.99841058]])
```

PyML - 2. Save Model

```
# Create prediction_model folder
!mkdir -p prediction_model

# Export the model weights
from sklearn.externals import joblib
joblib.dump(log_reg, 'prediction_model/weights.pkl')

# Export feature columns
import pickle
pickle.dump(feature_columns, open('prediction_model/feature_columns.pkl', 'wb'))
```

PyML - 3. Implement Serving Interface

```
# Create model.py inline via Jupyter command "writefile"
%%writefile prediction_model/model.py
import pandas as pd
import numpy as np
import pickle

from pyml.model.dataframe_model import DataFrameModel
from sklearn.externals import joblib

class LogisticRegressionModel(DataFrameModel):
    """A DataFrameModel takes input as a pandas DataFrame and produces predictions as a Pandas
    DataFrame.
    """
    def __init__(self):
        super(LogisticRegressionModel, self).__init__()

        # Load the model weights and feature columns
        self.clf = joblib.load('weights.pkl')
        self.feature_columns = pickle.load(open('feature_columns.pkl', 'rb'))

    def predict(self, df):
        df['probability'] = self.clf.predict_proba(
            df[self.feature_columns])[:,0]
        return df
```

PyML - 4. Package and Upload Model

```
from pyml import PyMLModel

pyml_model = PyMLModel(model_path="prediction_model/", model_name=example_prediction_model)
pyml_model.predict(pandas_df)[:2]['probability']
```

```
0    0.999995
1    0.999996
```

```
from pyml import Client
client = Client(user_email="kstumpf@uber.com", team_name="michelangelo")

# Upload the model and build the model's Docker image
model_id = client.upload_model(pyml_model)
```


PyML - 5. See Model in Michelangelo UI

The screenshot displays the Michelangelo UI interface for a project named "My Test Models". The page header includes navigation links for "BACK TO PROJECTS", "DOCUMENTATION", "TIER" (set to 5), and "ROLLOUT". A sidebar on the left contains icons for "Data Sources", "Data Snapshots", "Models", and "Deployments". The main content area shows a list of models under the "MODELS" tab. A single model is listed with the ID "example_prediction_model", created on 2018-10-15 at 13:37, of type "PyML", and created by "kstumpf@uber.com". The model's status is "Trained" and "Deployed", both indicated by green checkmarks. A "DEPLOY" button and a menu icon are visible next to the model entry.

My Test Models
Project used for demonstration purposes

DOCUMENTATION TIER 5 ROLLOUT

+ Add

MODELS TEMPLATES + CREATE MODEL

All Trained Deployed Search...

STATUS	MODEL ID	TYPE	CREATED BY	
✓	example_prediction_model Created 2018-10-15 13:37	PyML	kstumpf@uber.com	DEPLOY ...

PyML - 6. Deploy and Test Model

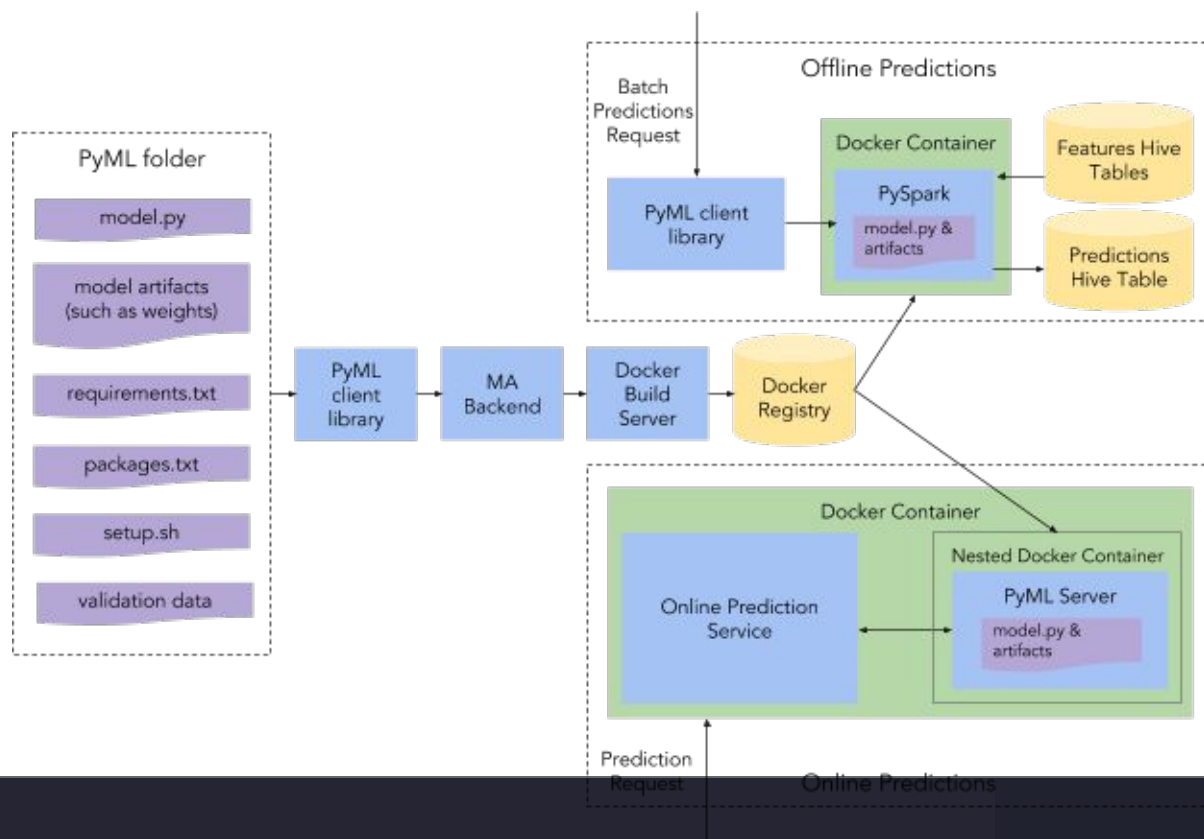
```
client.deploy_model(model_id)
```

```
from pyml import OnlineClient

online_client = OnlineClient(model_id=model_id)
output_df = online_client.predict(pandas_df[:2])
print output_df[['target', 'probability']]
```

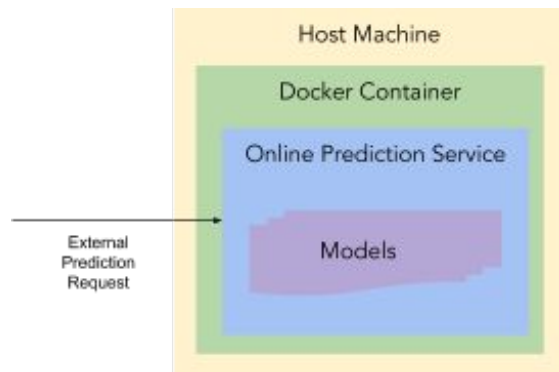
```
target probability
0      0.0      0.999996
1      0.0      0.999995
```

PyML - Architecture

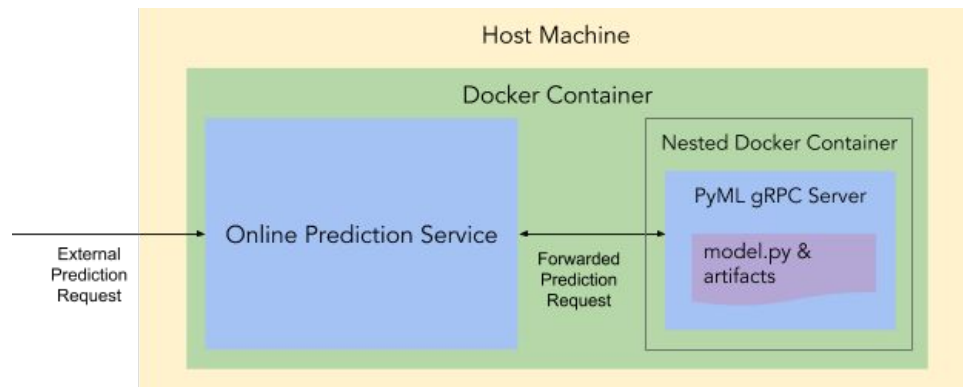


PyML - Serving Architecture

Standard Model



PyML Model



Accelerate ML:

Horovod

Horovod - Intro

- There are many ways to do data-parallel training.
- Some are more confusing than others. UX varies greatly.
- Our goals:
 - Infrastructure people deal with choosing servers, network gear, container environment, default containers, and tuning distributed training performance.
 - ML engineers focus on making great models that improve business using deep learning frameworks that they love.

Horovod - Complex Parameter Server Setup

```
import argparse
import sys

import tensorflow as tf

FLAGS = None

def main():
    ps_hosts = FLAGS.ps_hosts.split(",")
    worker_hosts = FLAGS.worker_hosts.split(",")

    # Create a cluster from the parameter server and worker hosts.
    cluster = tf.train.ClusterSpec({"ps": ps_hosts, "worker": worker_hosts})

    # Create and start a server for the local task.
    server = tf.train.Server(cluster,
                             job_name=FLAGS.job_name,
                             task_index=FLAGS.task_index)

    if FLAGS.job_name == "ps":
        server.join()
    elif FLAGS.job_name == "worker":

        # Assigns ops to the local worker by default.
        with tf.device(tf.train.replica_device_setter(
            worker_device="/job:worker/task:%d" % FLAGS.task_index,
            cluster=cluster)):

            # Build model...
            loss = ...
            global_step = tf.contrib.framework.get_or_create_global_step()

            train_op = tf.train.AdagradOptimizer(0.01).minimize(
                loss, global_step=global_step)

            # The StopAtStepHook handles stopping after running given steps.
            hooks=[tf.train.StopAtStepHook(last_step=100000)]

            # The MonitoredTrainingSession takes care of session initialization,
            # restoring from a checkpoint, saving to a checkpoint, and closing when done
            # or an error occurs.
            with tf.train.MonitoredTrainingSession(master=server.target,
                                                  is_chief=(FLAGS.task_index == 0),
                                                  checkpoint_dir="/tmp/train_logs",
                                                  hooks=hooks) as mon_sess:

                while not mon_sess.should_stop():
                    # Run a training step asynchronously.
                    # See `tf.train.SyncReplicasOptimizer` for additional details on how to
                    # perform *synchronous* training.
                    # mon_sess.run handles AbortedError in case of preempted PS.
                    mon_sess.run(train_op)

if __name__ == "__main__":
    parser = argparse.ArgumentParser()
    parser.register("type", "bool", lambda v: v.lower() == "true")
    # Flags for defining the tf.train.ClusterSpec
    parser.add_argument(
        "--ps_hosts",
        type=str,
        default="",
        help="Comma-separated list of hostname:port pairs"
    )
    parser.add_argument(
        "--worker_hosts",
        type=str,
        default="",
        help="Comma-separated list of hostname:port pairs"
    )
    parser.add_argument(
        "--job_name",
        type=str,
        default="",
        help="One of 'ps', 'worker'"
    )
    # Flags for defining the tf.train.Server
    parser.add_argument(
        "--task_index",
        type=int,
        default=0,
        help="Index of task within the job"
    )
    FLAGS, unparsed = parser.parse_known_args()
```

Image Source: TensorFlow

-- <https://www.tensorflow.org/deploy/distributed>

Horovod - Simple Horovod Setup

```
import tensorflow as tf
import horovod.tensorflow as hvd

# Initialize Horovod
hvd.init()

# Pin GPU to be used
config = tf.ConfigProto()
config.gpu_options.visible_device_list = str(hvd.local_rank())

# Build model...
loss = ...
opt = tf.train.AdagradOptimizer(0.01)

# Add Horovod Distributed Optimizer
opt = hvd.DistributedOptimizer(opt)

# Add hook to broadcast variables from rank 0 to all other processes during initialization.
hooks = [hvd.BroadcastGlobalVariablesHook(0)]

# Make training operation
train_op = opt.minimize(loss)

# The MonitoredTrainingSession takes care of session initialization,
# restoring from a checkpoint, saving to a checkpoint, and closing when done
# or an error occurs.
with tf.train.MonitoredTrainingSession(checkpoint_dir="/tmp/train_logs",
                                       config=config, hooks=hooks) as mon_sess:
    while not mon_sess.should_stop():
        # Perform synchronous training.
        mon_sess.run(train_op)
```

Initialize Horovod

Assign a GPU to each TensorFlow process

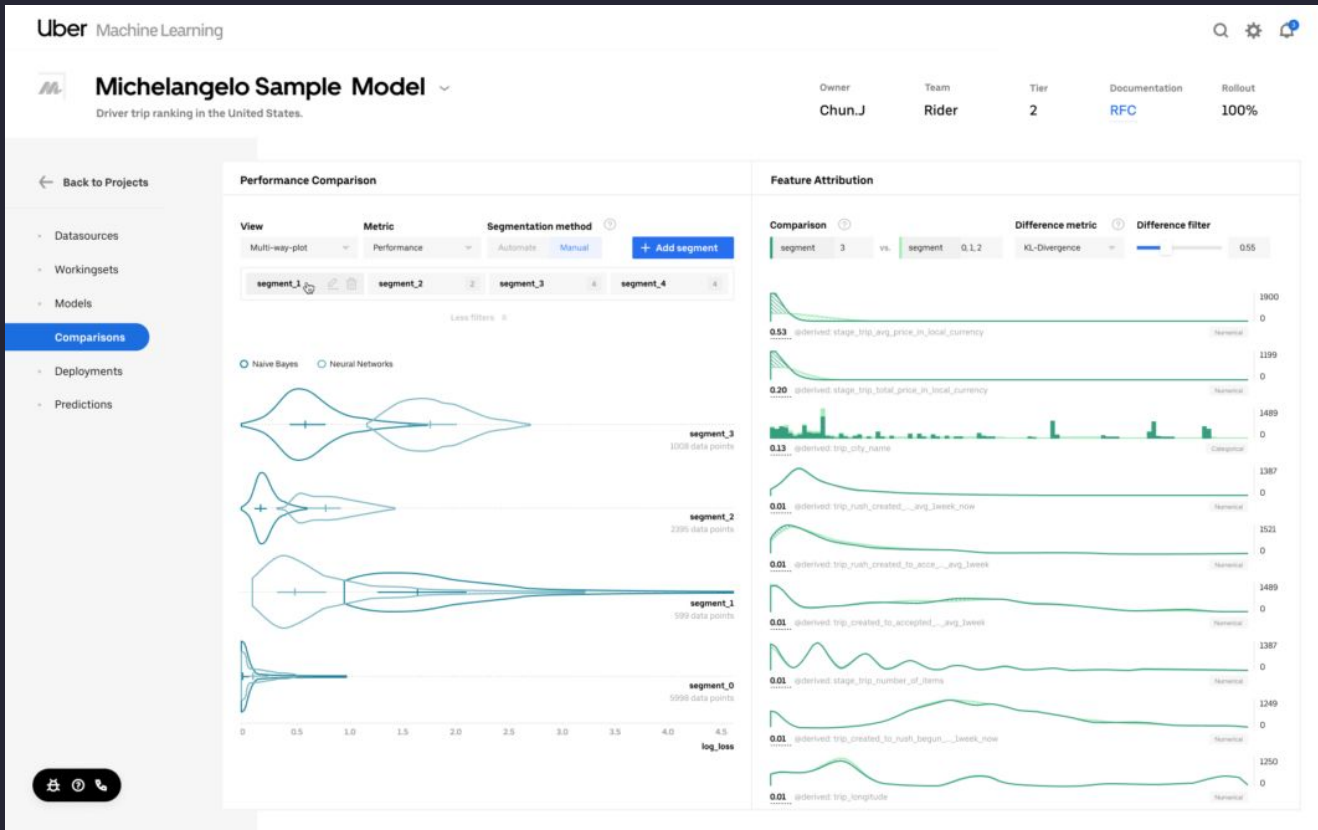
Wrap regular TensorFlow optimizer with Horovod optimizer which takes care of averaging gradients using ring-allreduce

Broadcast variables from the first process to all other processes to ensure consistent initialization

Accelerate ML:

Manifold Viz

Manifold Viz



Accelerate ML:

AutoTune

AutoTune

Model Settings

Description
GBD Tree with AutoTune

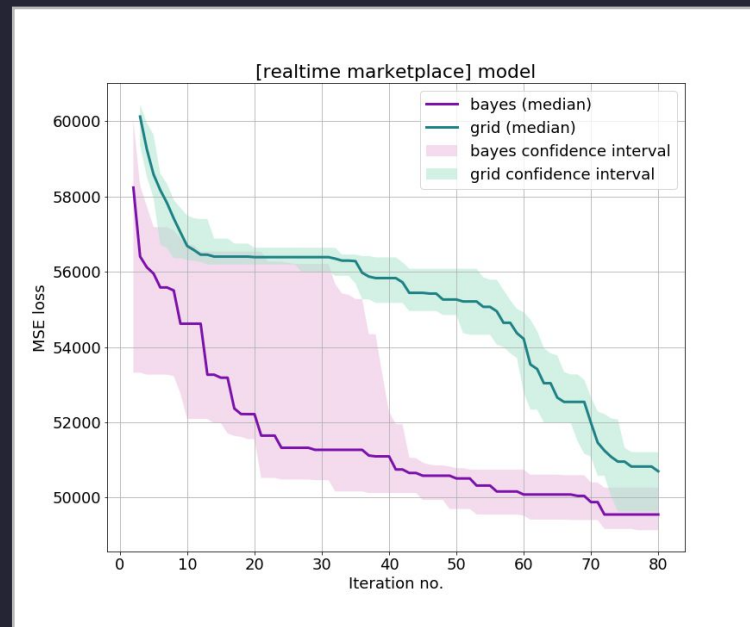
Model Type
GBD Tree Classification [Binary]

HYPERPARAMETERS NEW Use Autotune

Algorithm: Bayesian Optimization with Gaussian Process Kernel
Acquisition Function: Lower Confidence Bound

Max Depth: 6-10 (Example: 6-10 or 6,8,10)
Max Count: 100-200 (Example: 50-120 or 50,100,120)
Max After Bal: 2
Number Of Bins: 10 (Example: 10-50 or 10,30,50)

Min Rows: 100 (Example: 50-200 or 50,100,200)
Learn Rate: 0.1-0.2 (Example: 0.01-0.2 or 0.01,0.1,0.2)
Probability Calibration: None



Key Lessons Learned

Key Lessons Learned

- Let developers use the tools that they want
- Data is the hardest part of ML and the most important piece to get right
- It can take significant effort to make open source and commercial components work well at scale.
- Develop iteratively based on user feedback, with the long-term vision in mind
- Real-time ML is challenging to get right

Thank you!

eng.uber.com/michelangelo

eng.uber.com/horovod

uber.com/careers

Uber

Proprietary and confidential © 2018 Uber Technologies, Inc. All rights reserved. No part of this document may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying, recording, or by any information storage or retrieval systems, without permission in writing from Uber. This document is intended only for the use of the individual or entity to whom it is addressed and contains information that is privileged, confidential or otherwise exempt from disclosure under applicable law. All recipients of this document are notified that the information contained herein includes proprietary and confidential information of Uber, and recipient may not make use of, disseminate, or in any way disclose this document or any of the enclosed information to any person other than employees of addressee to the extent necessary for consultations with authorized personnel of Uber.