



全球软件开发大会【上海站】

Designing Machine Learning Platform

刘彦东 网易传媒技术副总裁



每天10分钟，邀请顶级技术专家，为你传道授业解惑。



扫一扫，试读专栏

主办方 Geekbang® InfoQ_{CONF}
极客邦科技

ArchSummit

全 球 架 构 师 峰 会 2017

12月8-9日 北京 · 国际会议中心



APSEC 2017

24th Asia-Pacific Software Engineering Conference
4-8 December 2017, Nanjing, Jiangsu, China

12月4-8日
中 国 南 京



了解详情

AiCon

全球人工智能技术大会 2018

助力人工智能落地

2018.1.13 - 1.14 北京国际会议中心



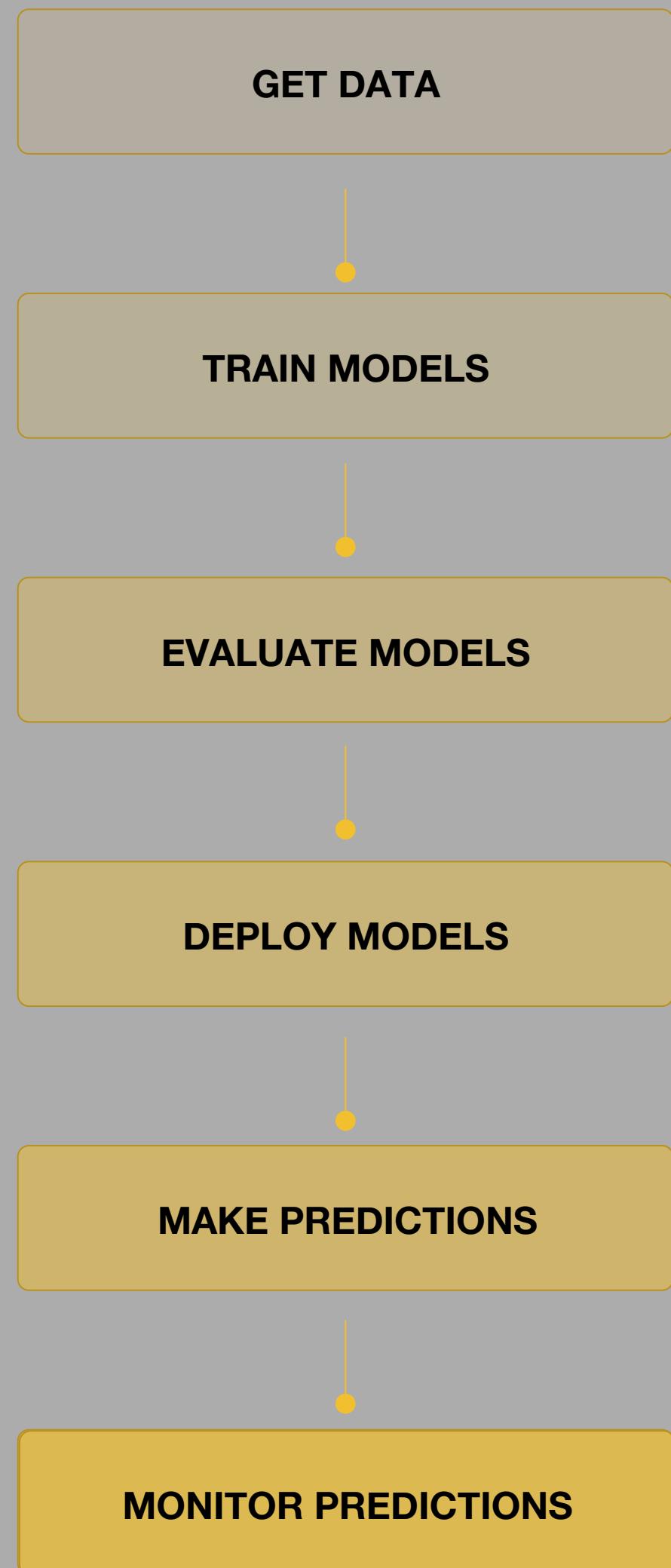
扫描关注大会官网

Agenda

- How ML works
- ML platform
- ML Case Study
- Key Challenges: Deep dive
- System Architectures

The Life Cycle of a Model

- Data
 - Data sources, batch or streaming, data aggregation...
- Training in various environments
 - Fast iteration, traditional ML vs DL, training environments ...
- Model evaluations
 - Standard evaluation vs. customized evaluation
- Model deployment
 - Versioning, production
- Inference
 - Batch predictions vs online predictions, scaling out, SLA ...
- Monitoring
 - Signal selection





ML PLATFORM MISSION

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

Design of ML Platform

- ML as a Service
- Scalable infrastructure for training and serving
- Workflow tools for prototyping, iteration, and productionization
- Model and data serving with full monitoring for batch and real-time
- Scope
 - Traditional ML & Deep Learning
 - Supervised, Unsupervised and Semi-supervised
 - Online learning



Having trouble? →

[GO TO OLD UI](#)[← BACK TO PROJECTS](#)

EATS_ETD_Prediction

OWNER	TEAM	DOCUMENTATION	TIER	ROLLOUT
ntle, alexnik, nico, myz	eats_core	+ Add	3	<div style="width: 100%; background-color: #2e7131; height: 10px;"></div>



DATA SOURCES

[+ ADD DATA SOURCE](#)

Data Sources



Dynamic Data Source

Type:

[DYNAMIC](#)

Source:

SQL Pipeline

Data Type:

spark_sql

[SETTINGS](#)

...



Models



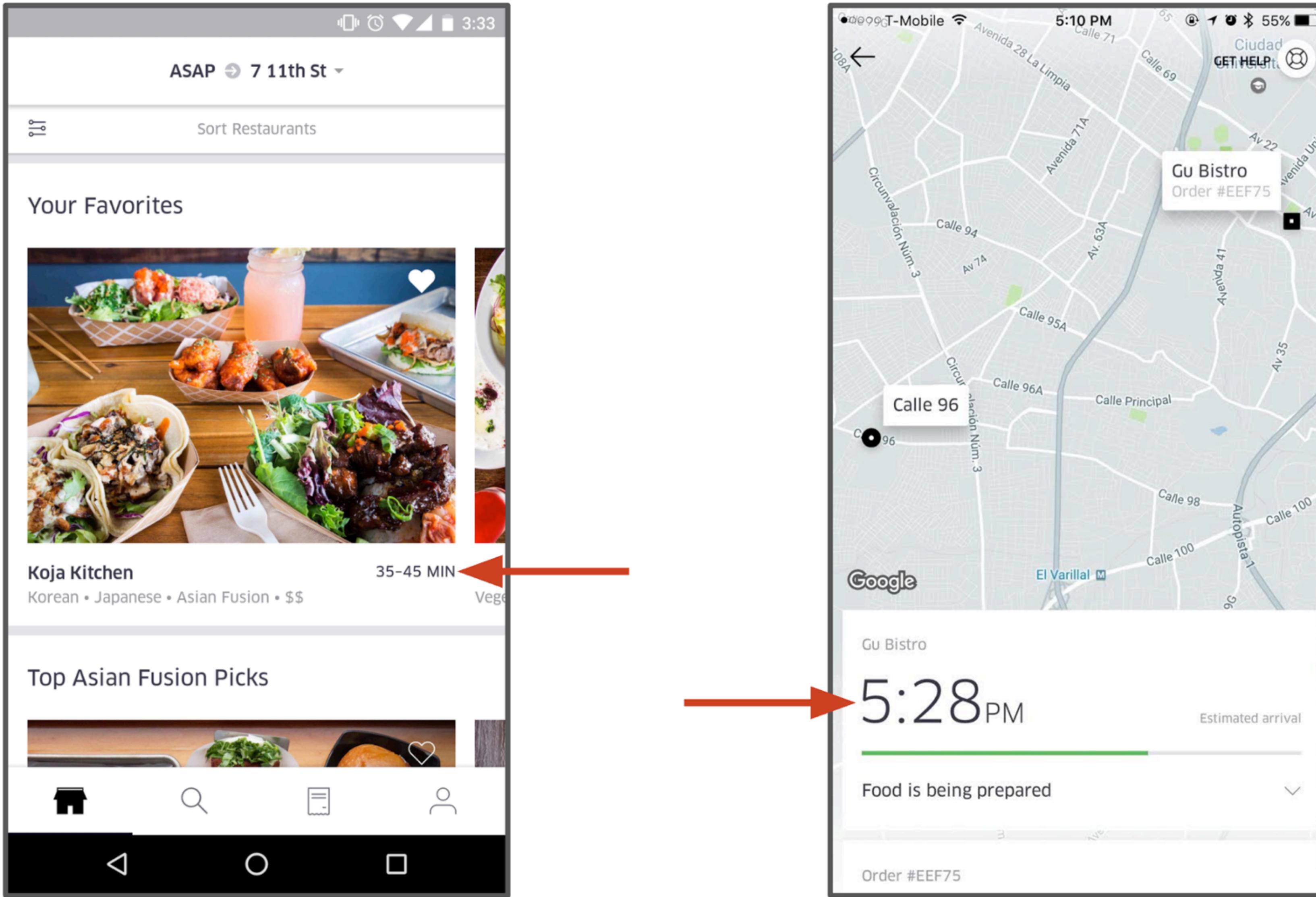
Deployments



Predictions

Example Use Case: UberEATS

MEAL DELIVERY TIME



Uber EATs Delivery Time Models

- Features
 - Curated features
 - Request Level Features – *user's current location*
- Models
 - Several models for different stages of order
 - GBDT Regression
 - Different versions of each for experimentation

Key Challenges

- Guarantee same data for batch training and online scoring
- Train and deploy separate model per city
- One-click deploy & easy scale out
- Live monitoring of model performance

Challenge 1:

Same data for train & predict

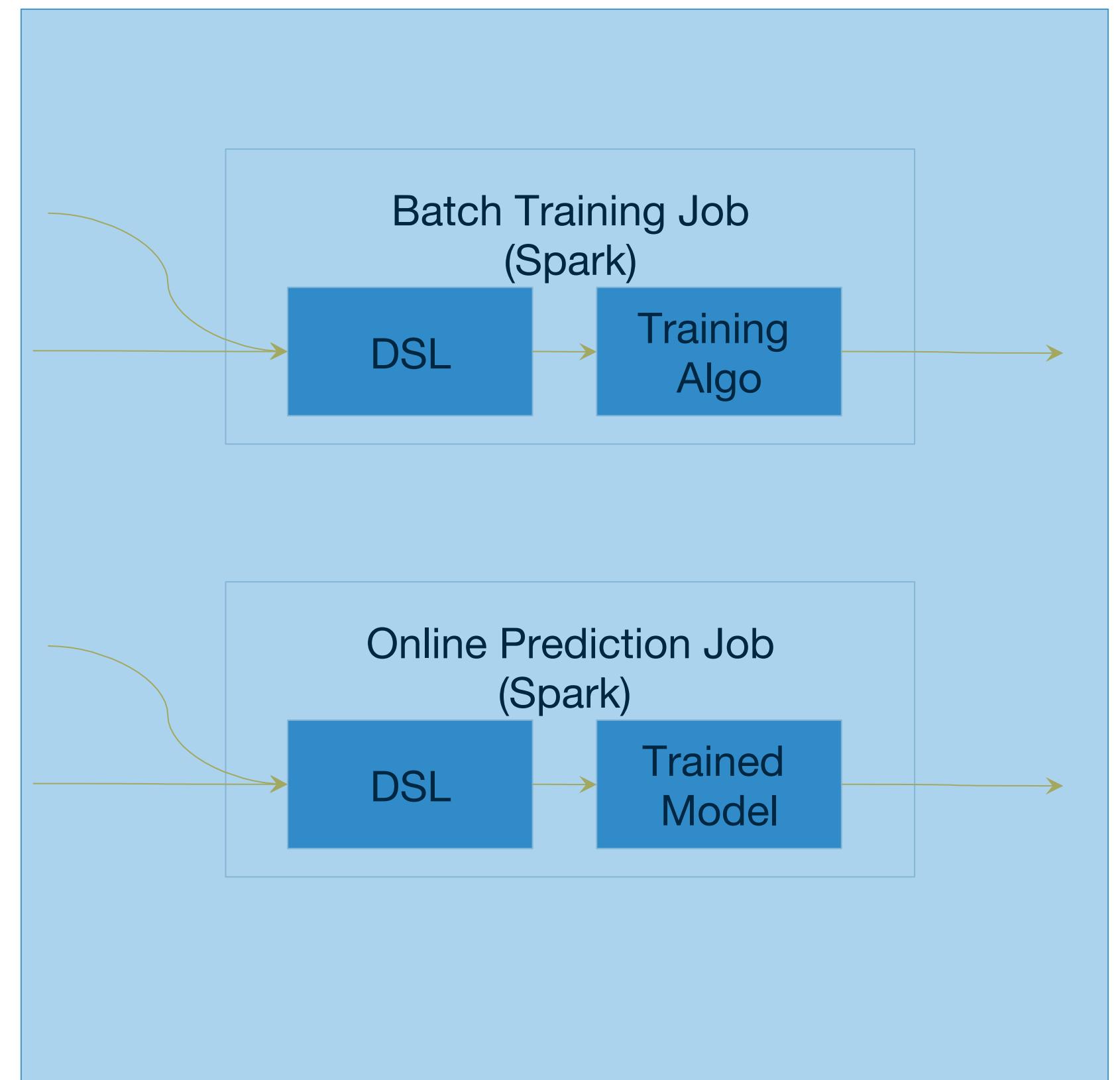
Data Sources: Problems

	Request Level Features	Aggregated Feature	Near Realtime Aggregated Feature
Training	Batch	Batch	Batch
Online scoring	Given by user	Curated in batch Consumed by query	Curated in streaming, consumed by query

Generation Pattern: Batch and streaming
Consuming Pattern: Batch and query

Data Sources (Solutions)

- Data Storage
 - Spark for batch jobs
 - Cassandra for online jobs
 - Streaming jobs
- Data Accessors
 - Own DSL
 - Access basis features, curated features, and column stats
- Data Transformation
 - Standard transformation functions + UDFs



Challenge 2: Partition Model

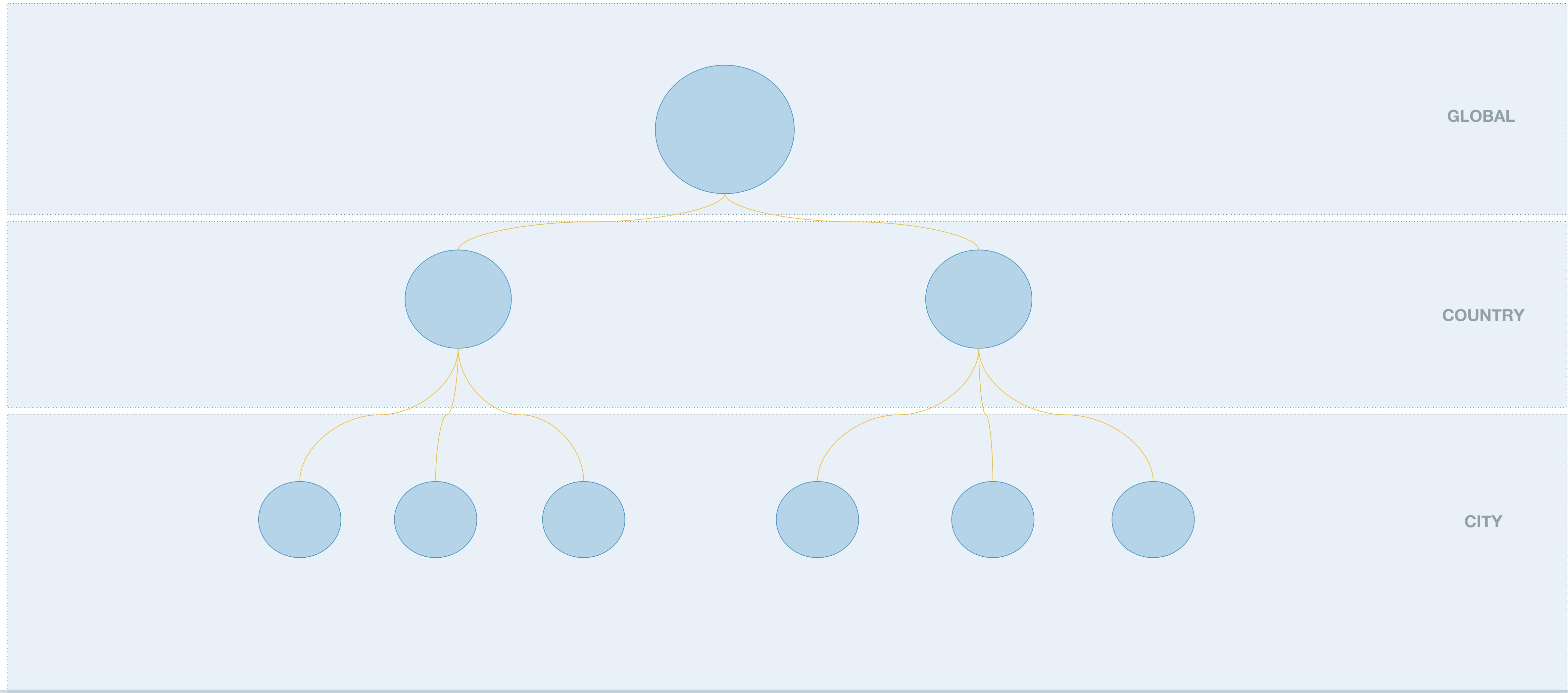
PROBLEM

- Often you want to train a model per city
- Hard to train and deploy a few hundred individual models

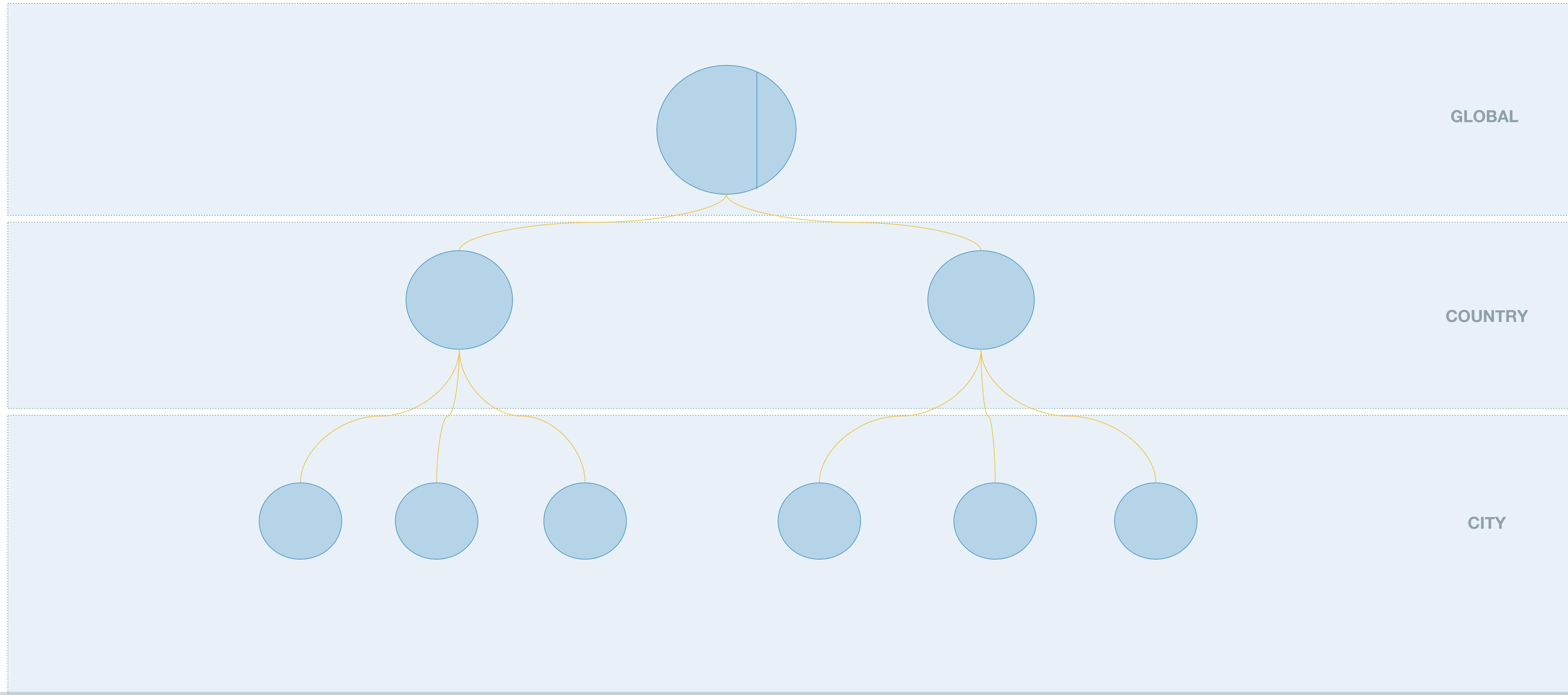
SOLUTION

- Let users define hierarchical partitioning scheme
 - Automatically train model per partition
- Manage and deploy as single logical model

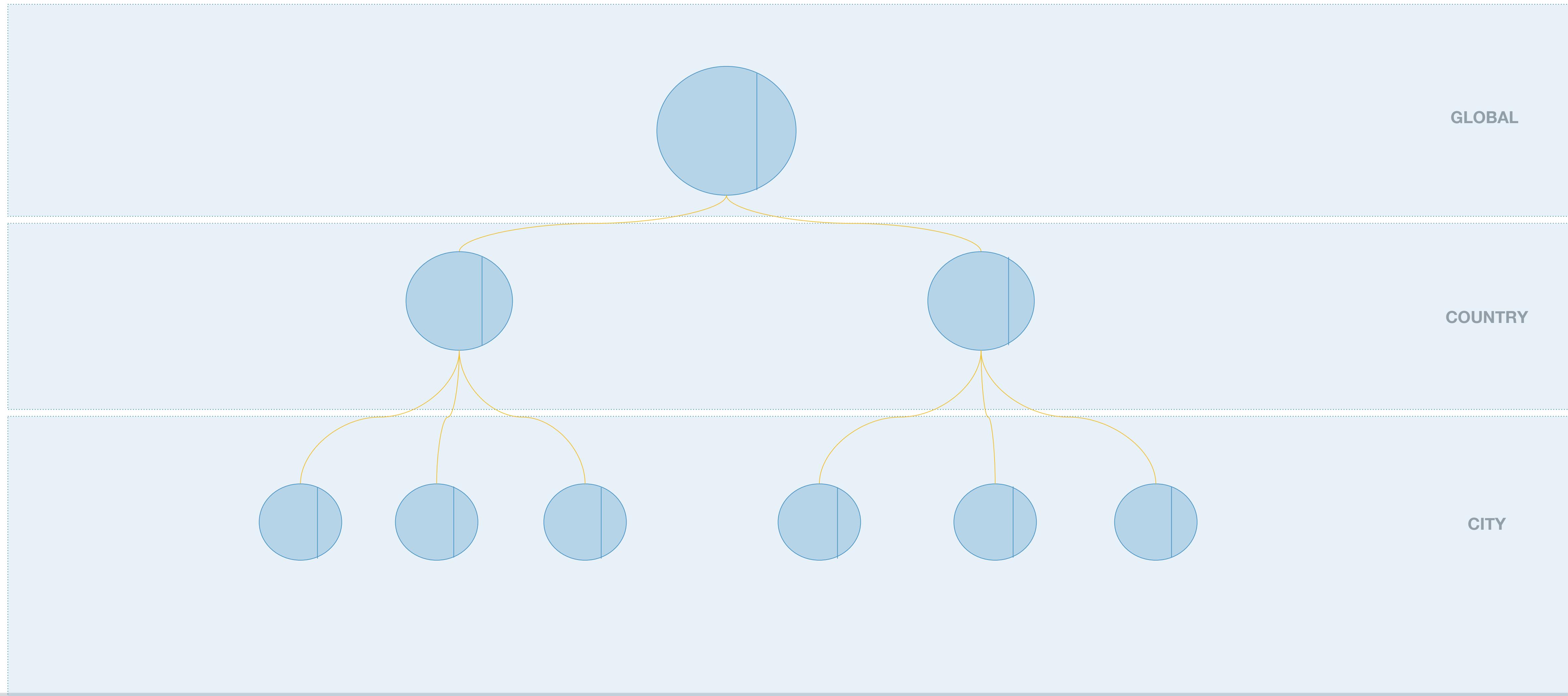
Define partition scheme



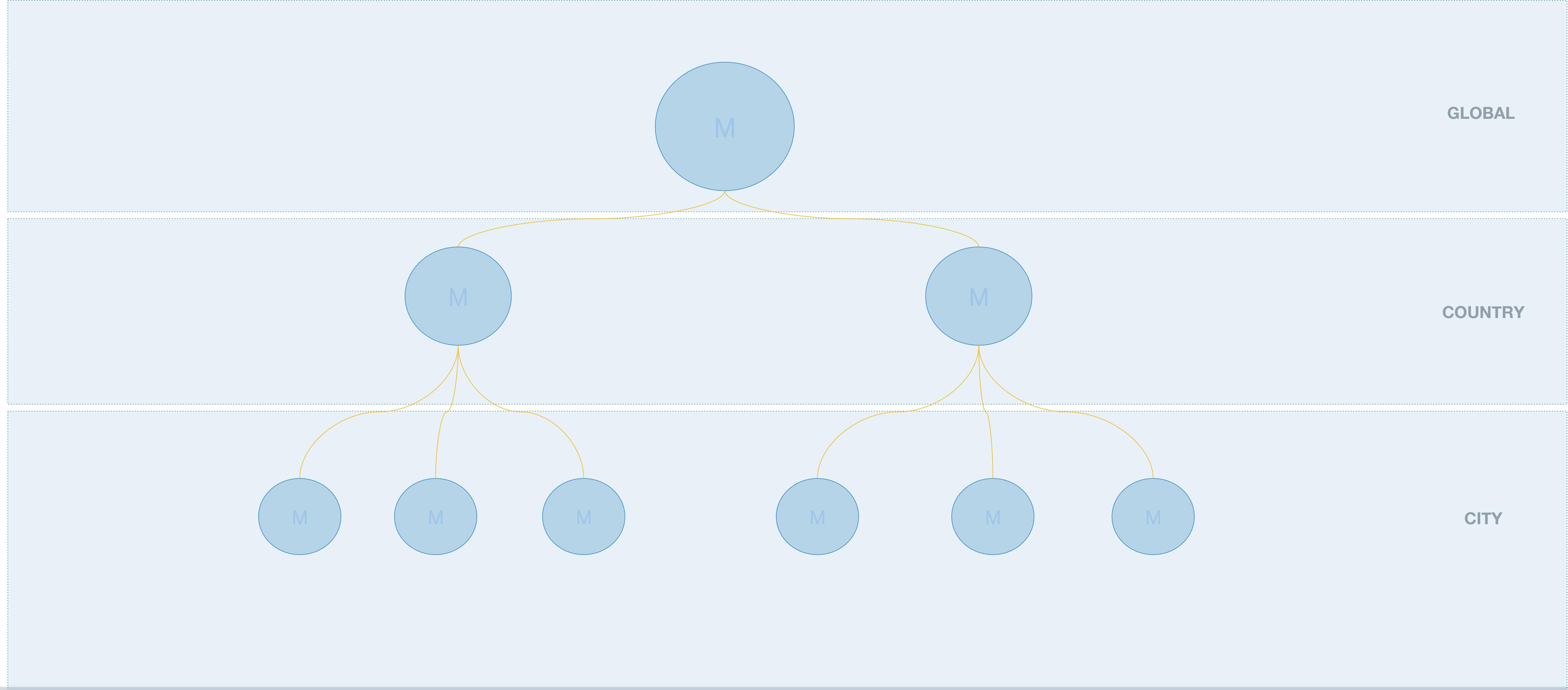
Make train / test split



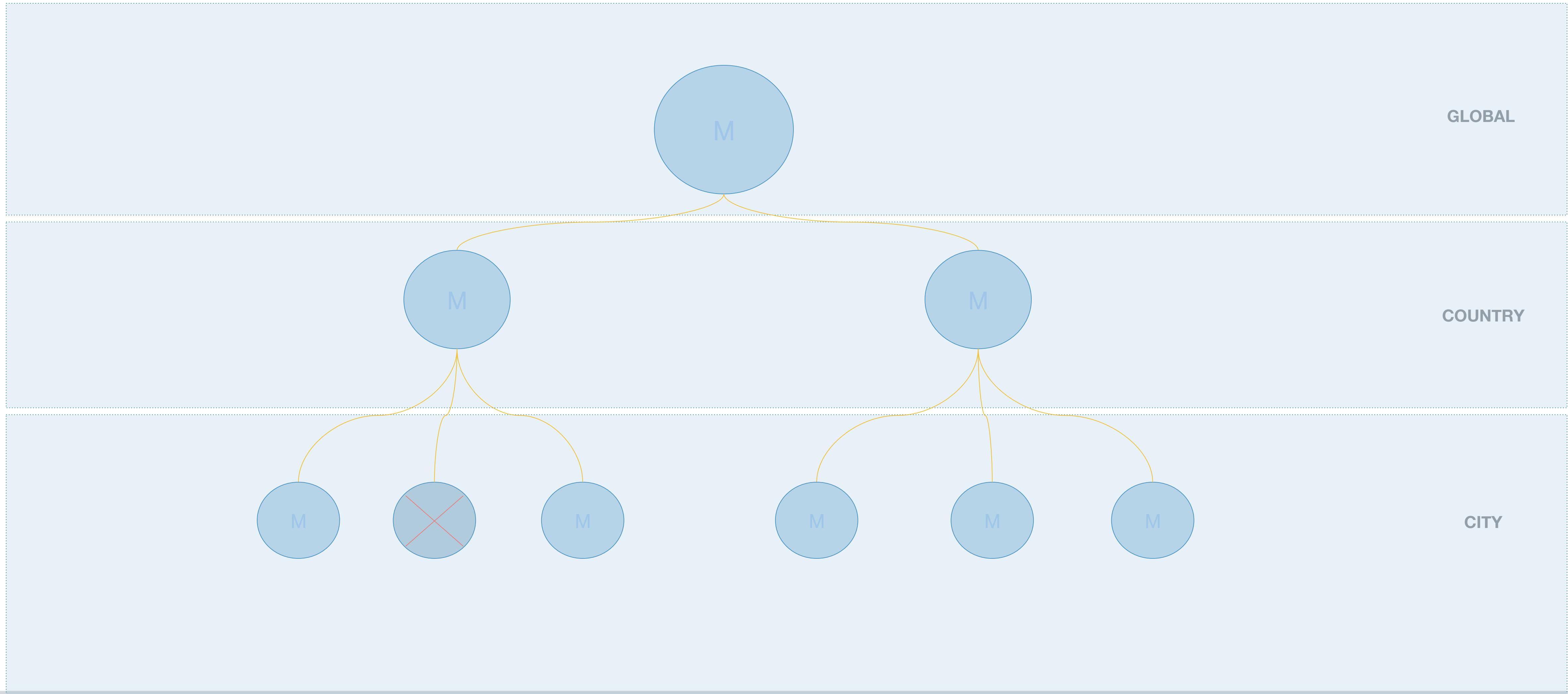
Keep same split and partition for each level



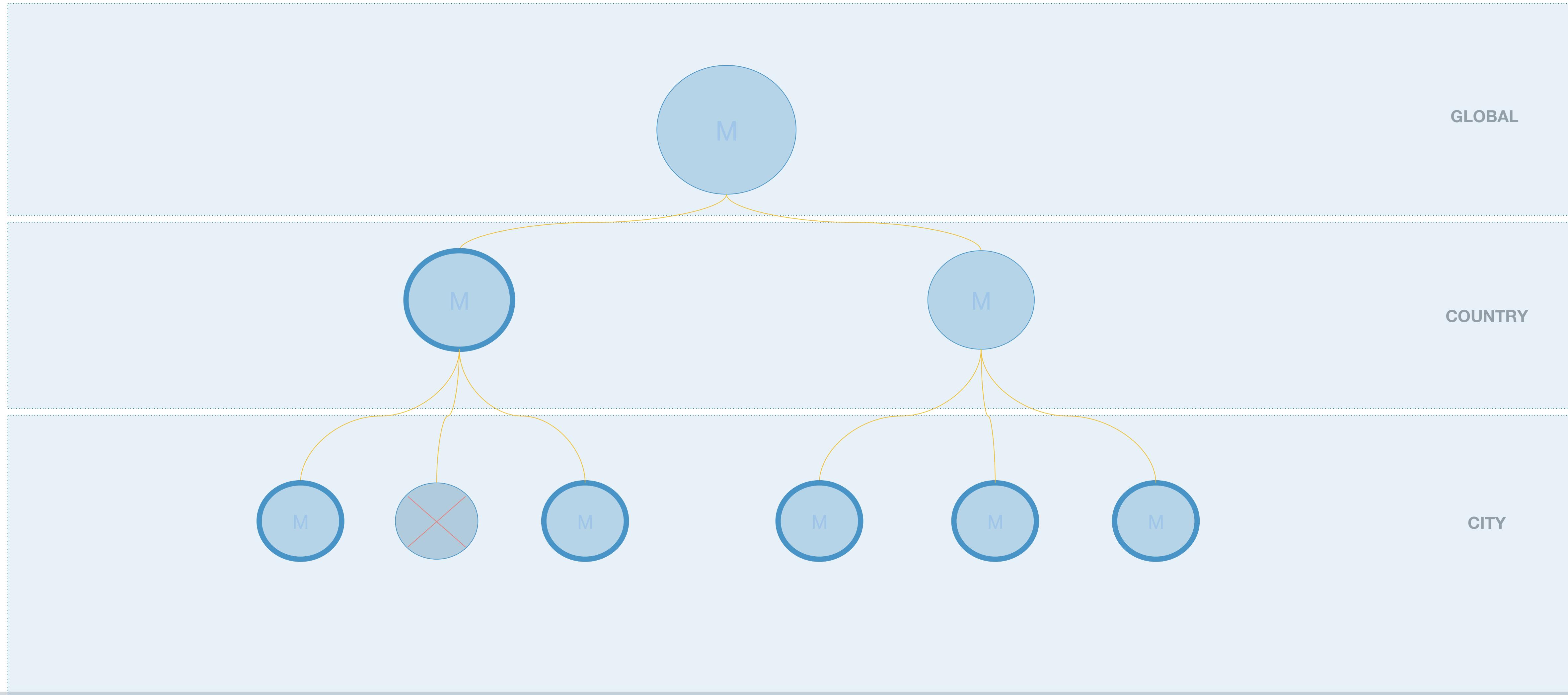
Train model for every node



Prune bad models



At serving time, route to best model for each node



Challenge 3: **One-click deploy and scale out**

REALTIME PREDICT SERVICE

- Predict service
 - RPC 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
 - Each model is packed individually
 - One click deploy across DCs via standard deployment infrastructure
 - Health checks and rollback

Challenge 4: Live model performance monitoring

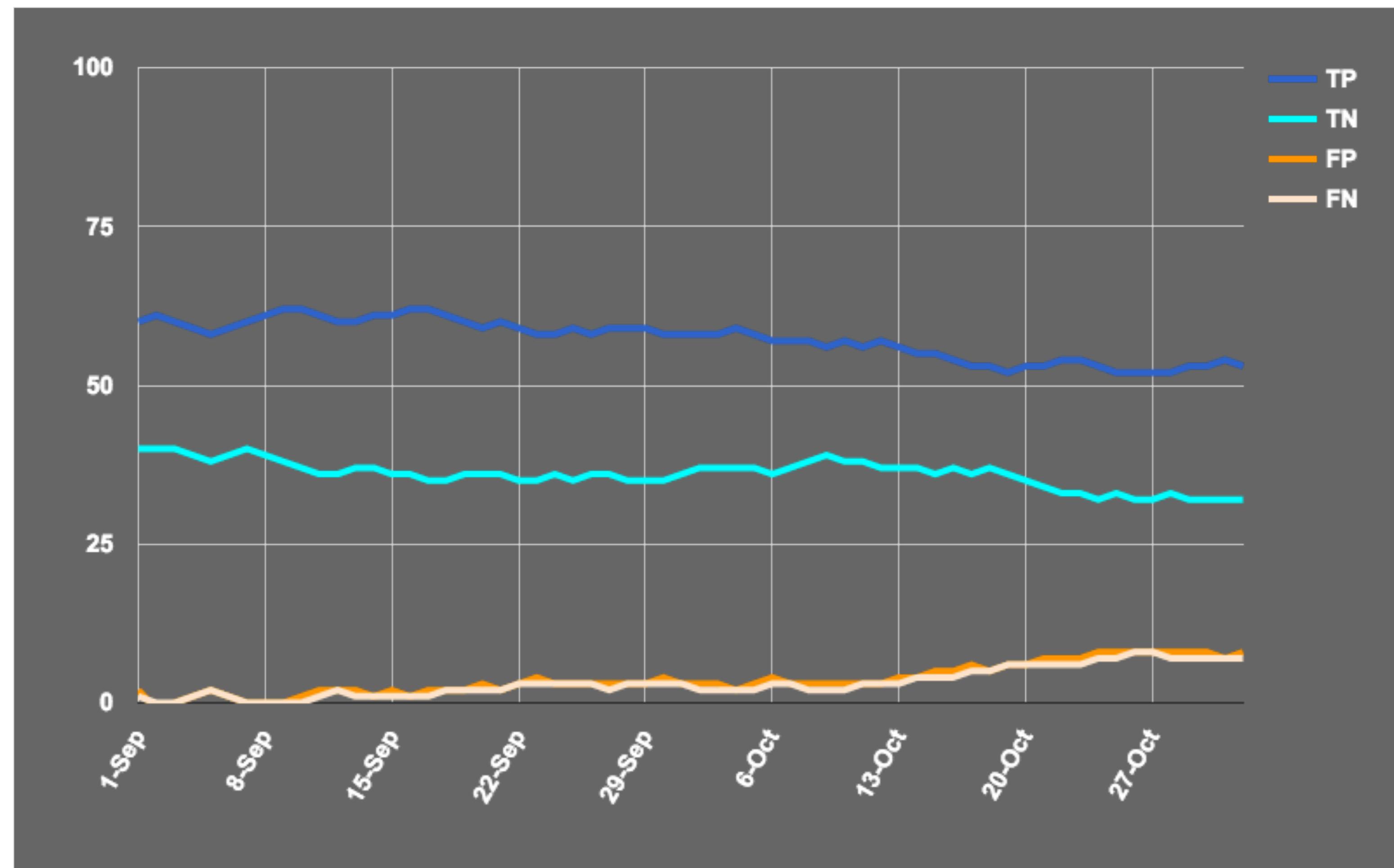
LIVE PREDICTION MONITORING

Problem

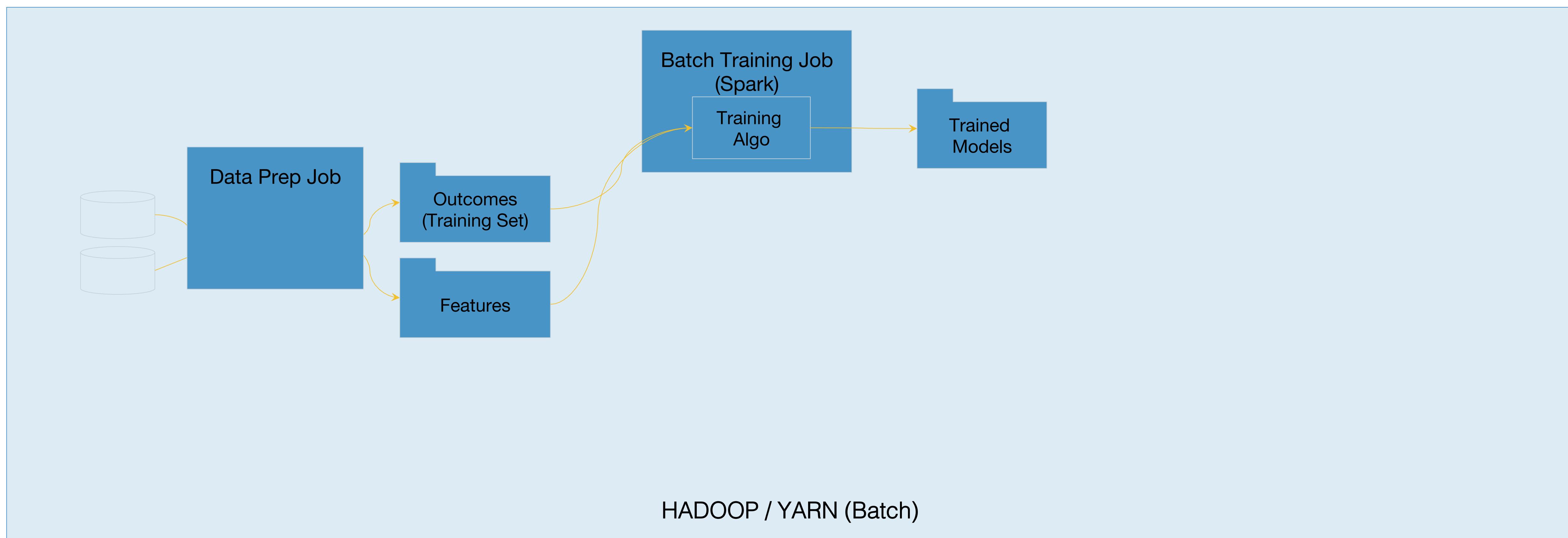
- Ensure deployed model is making good predictions

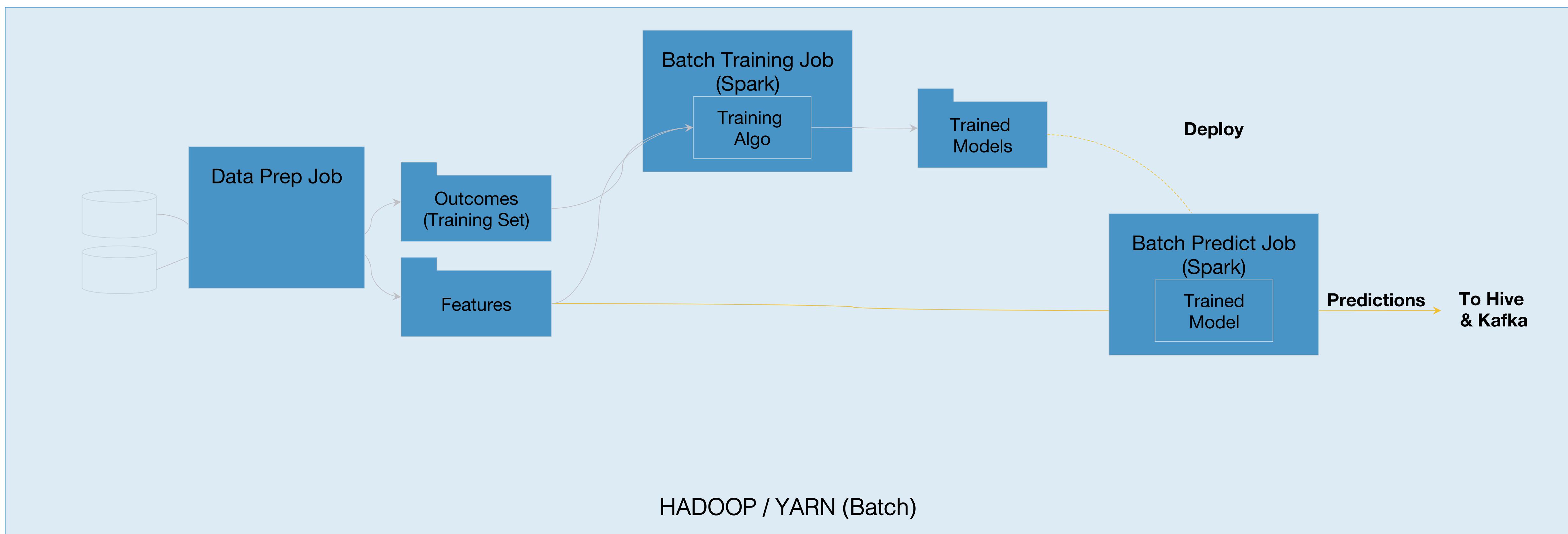
Solution

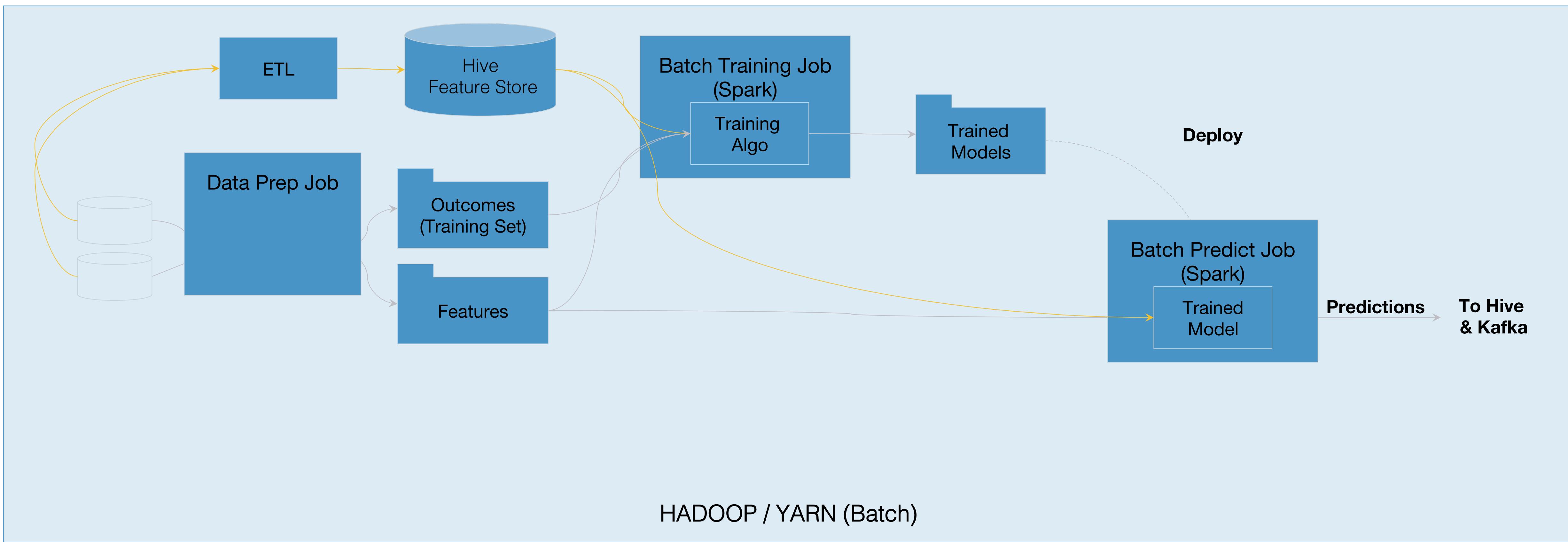
- Log predictions
- Join logged predictions to actual outcomes
- Publish metrics for monitoring and alerting
- Optionally hold back logged predictions



System Architecture





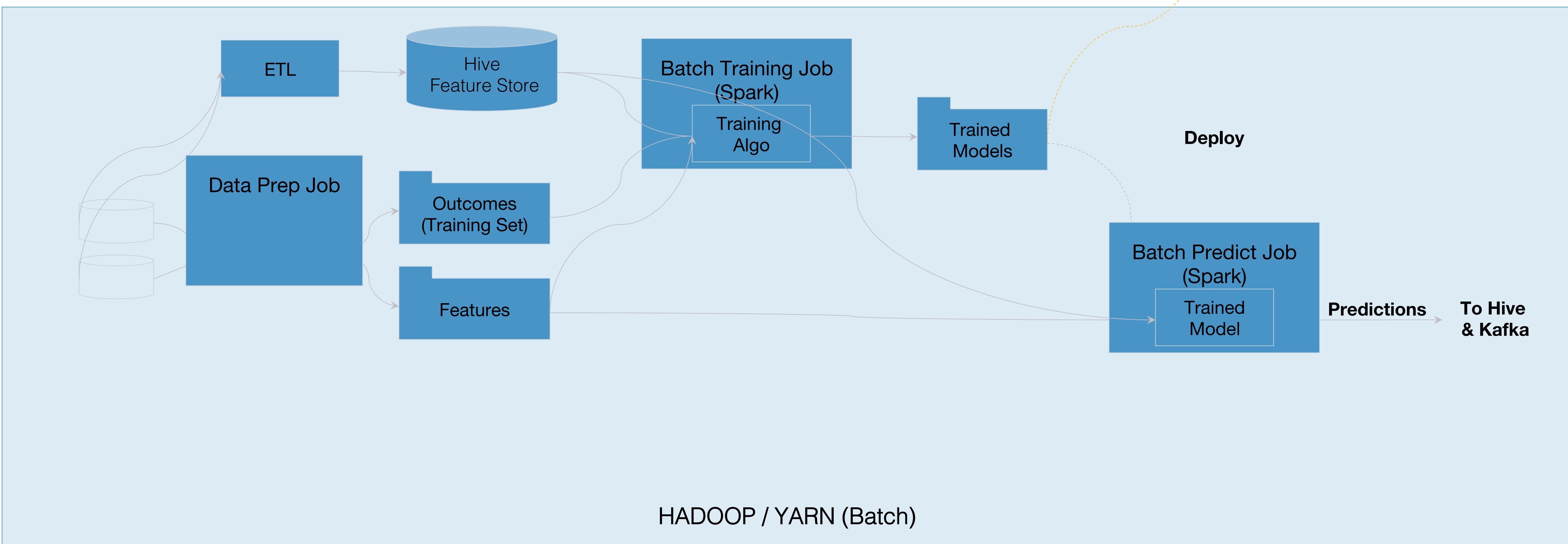
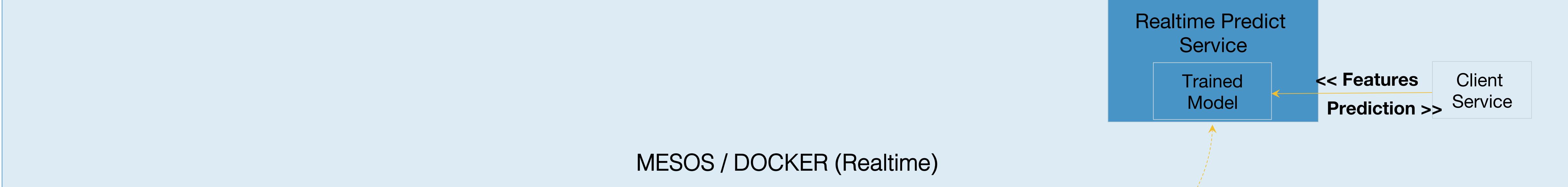


GET DATA

TRAIN MODELS

EVAL MODELS

DEPLOY, PREDICT & MONITOR

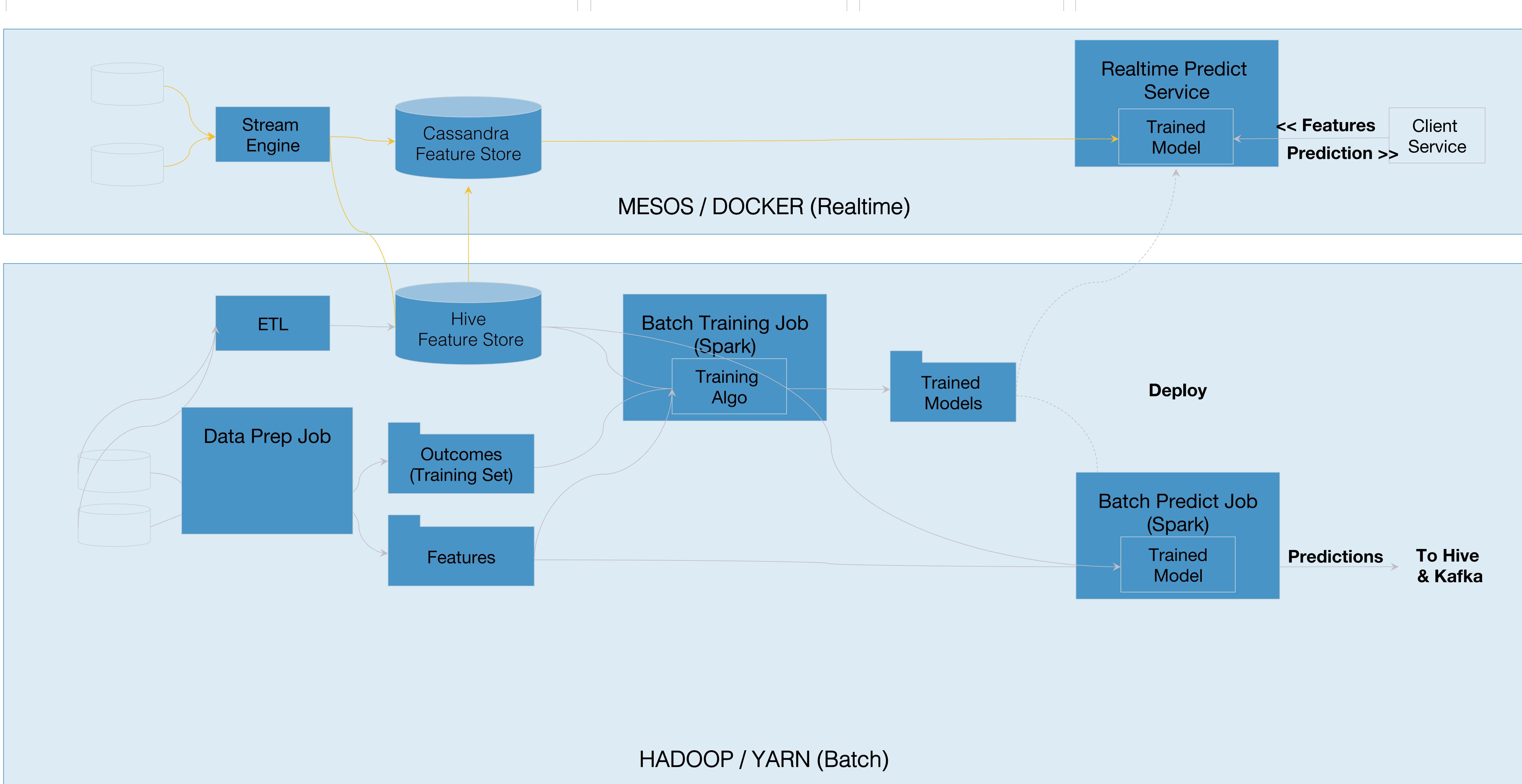


GET DATA

TRAIN MODELS

EVAL MODELS

DEPLOY, PREDICT & MONITOR

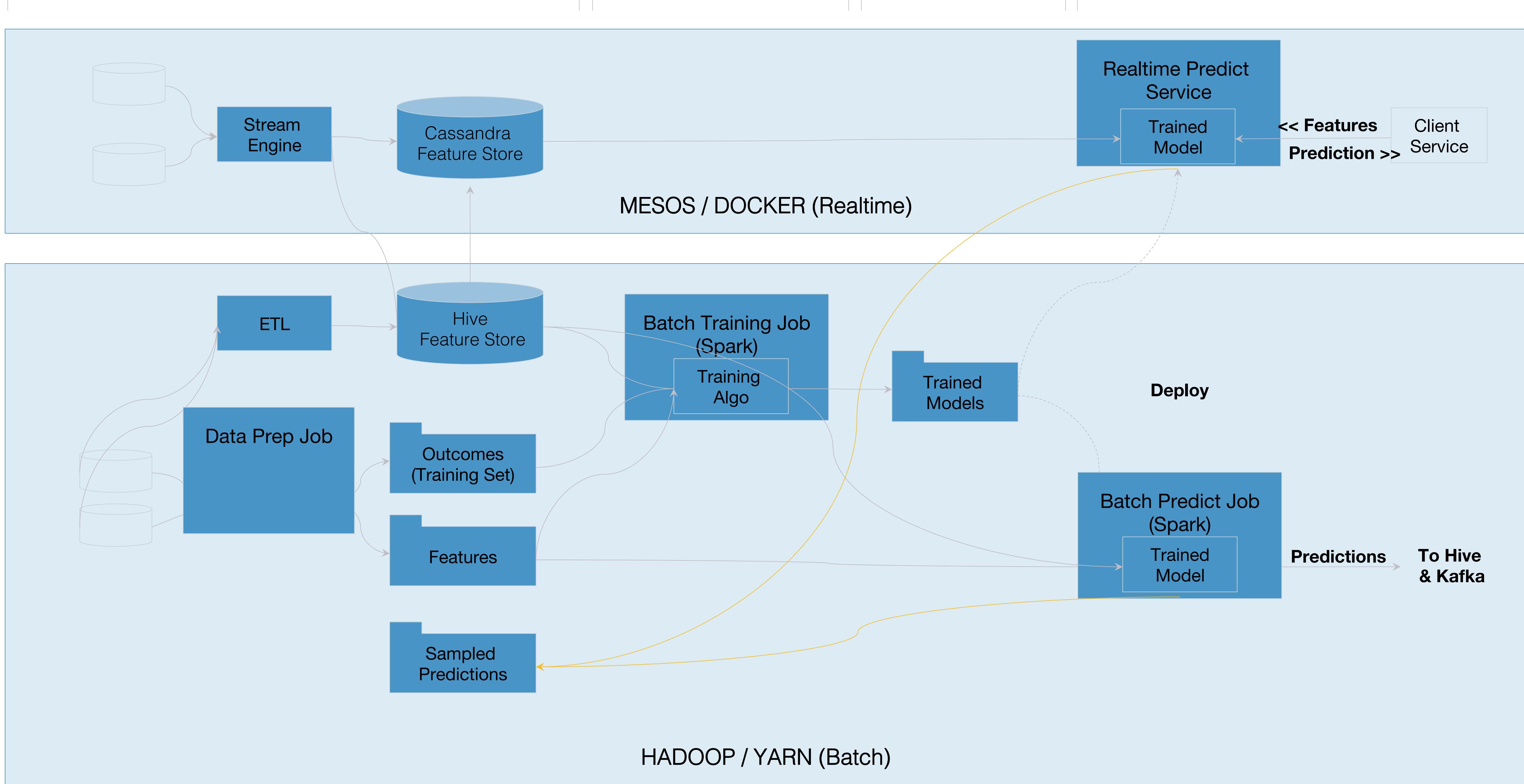


GET DATA

TRAIN MODELS

EVAL MODELS

DEPLOY, PREDICT & MONITOR

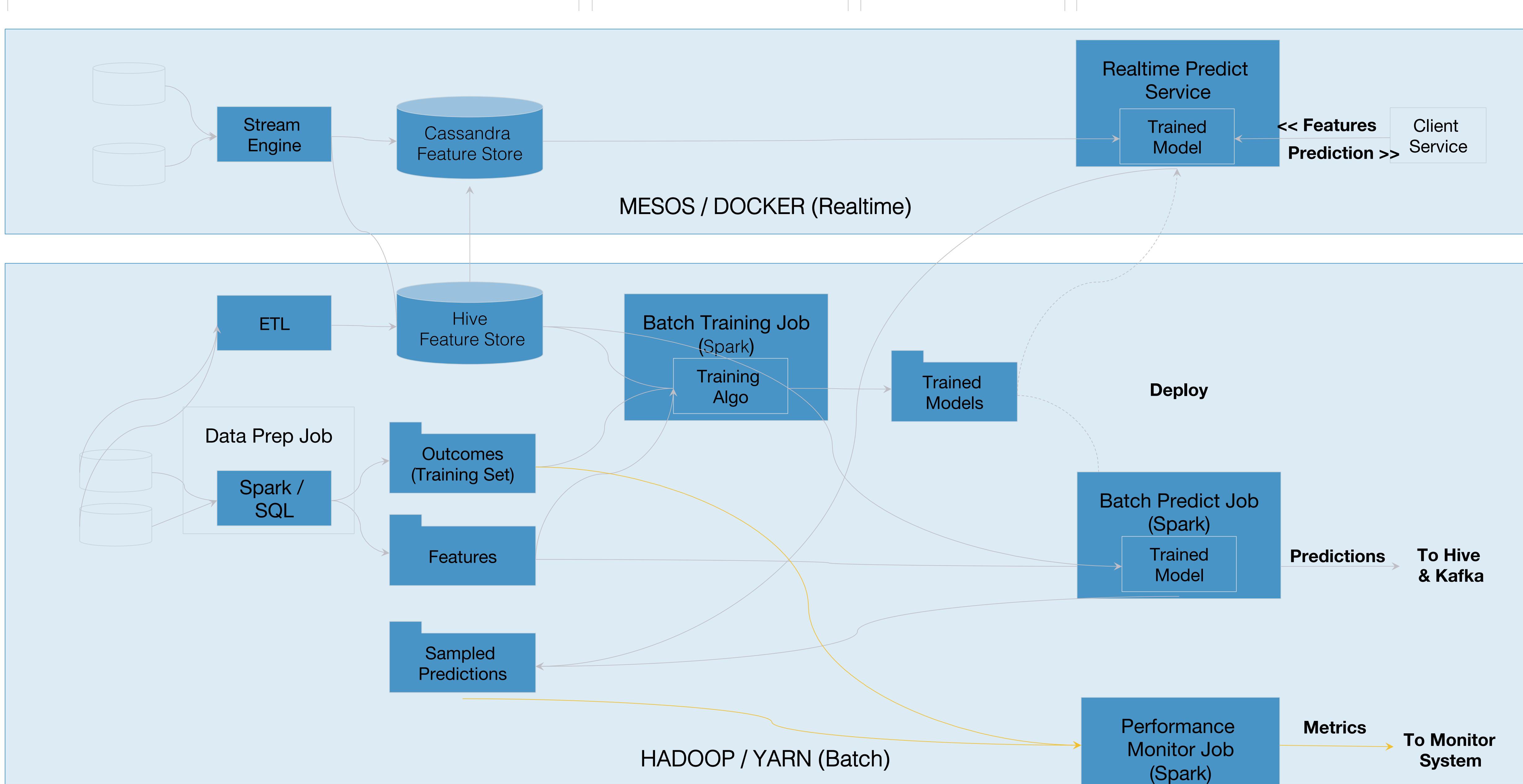


GET DATA

TRAIN MODELS

EVAL MODELS

DEPLOY, PREDICT & MONITOR

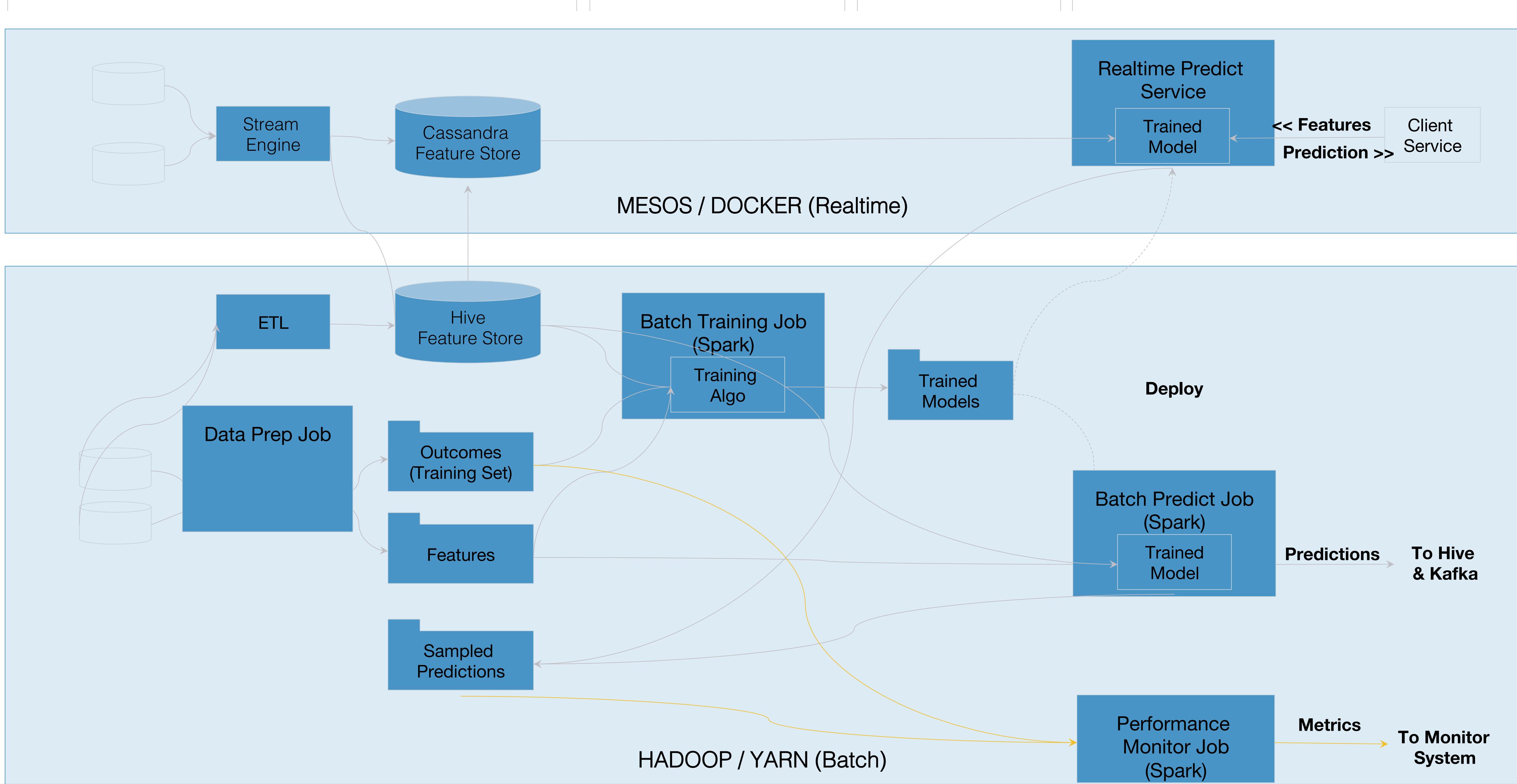


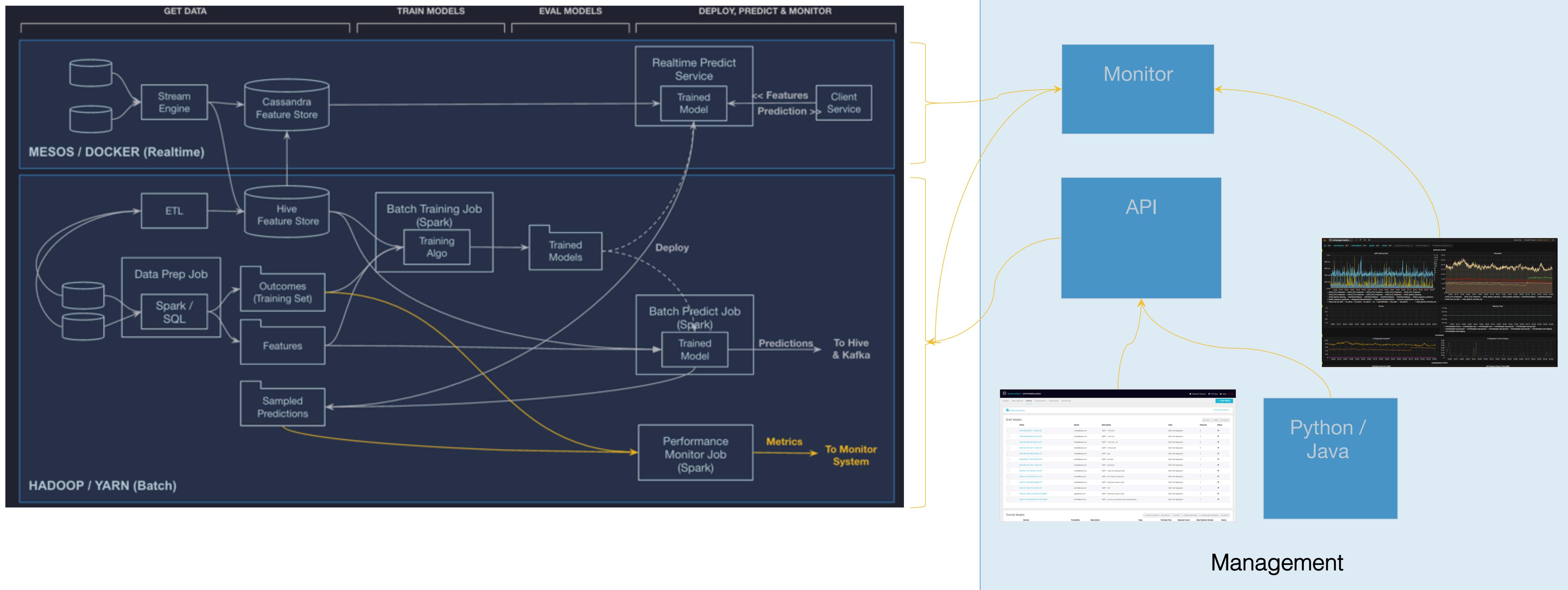
GET DATA

TRAIN MODELS

EVAL MODELS

DEPLOY, PREDICT & MONITOR







关注QCon微信公众号
获得更多干货!

Thanks!

INTERNATIONAL SOFTWARE DEVELOPMENT CONFERENCE

主办方: Geekbang > InfoQ
极客邦科技

