



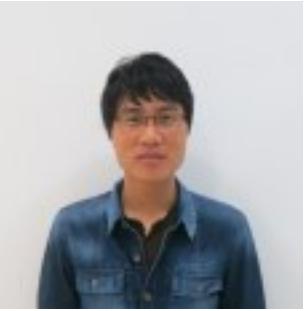
Christian Kadner



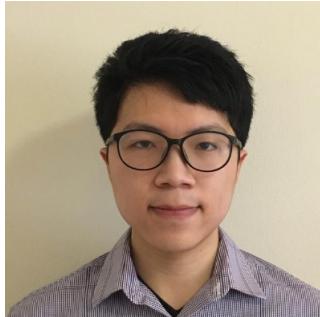
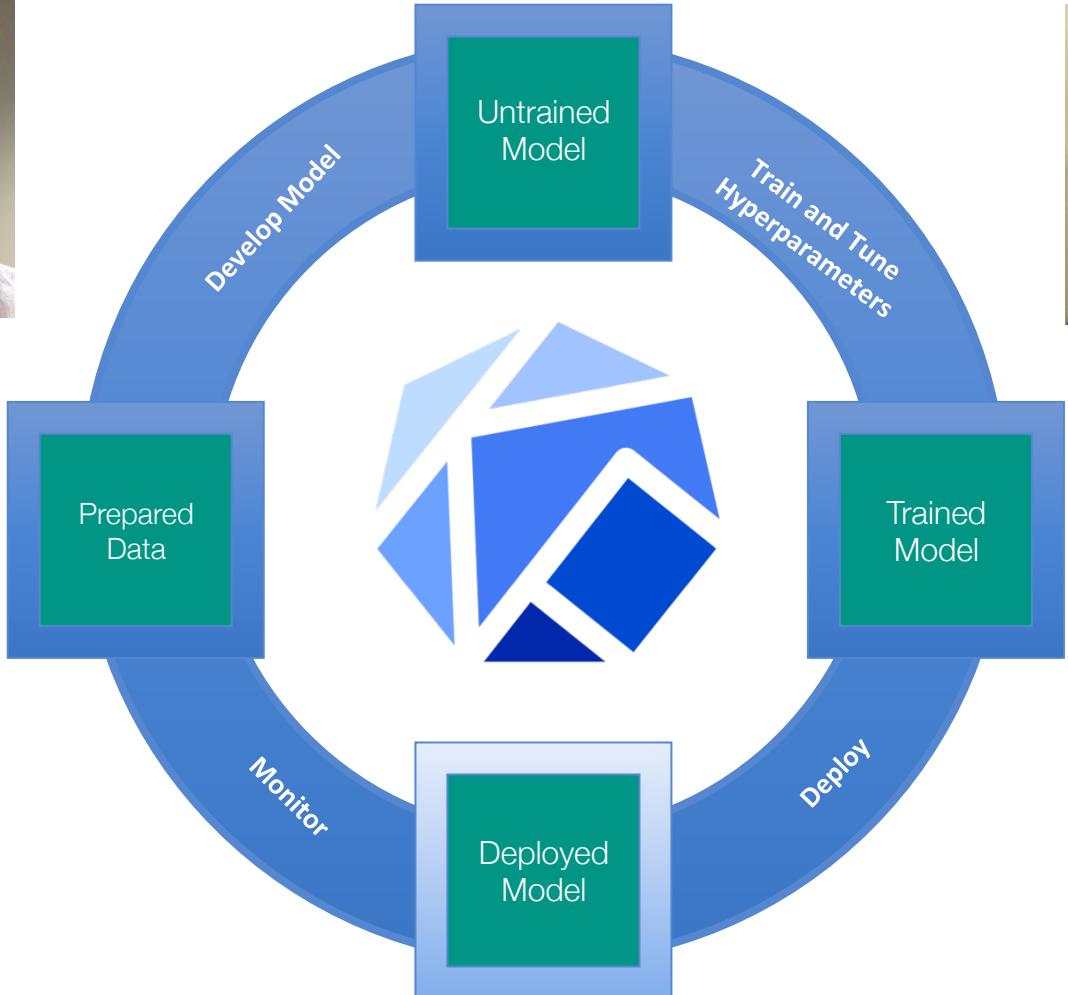
Weiqiang Huang



Jin Chi He



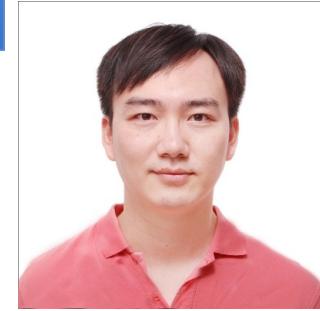
Feng Li



Tommy Li



Andrew Butler



Ke Zhu



Kevin Yu



'Upstream' is about extracting oil and natural gas from the ground; 'midstream' is about safely moving them thousands of miles; and 'downstream' is converting these resources into the fuels and finished products we all depend on.

## Upstream



Upstream has many phases, beginning with the exploratory process. Geologists search on dry land or in oceans for signs of gas or oil.

## Midstream



When a well is producing, oil or gas enters the midstream juncture. The middle part of the process requires multiple cooperation.

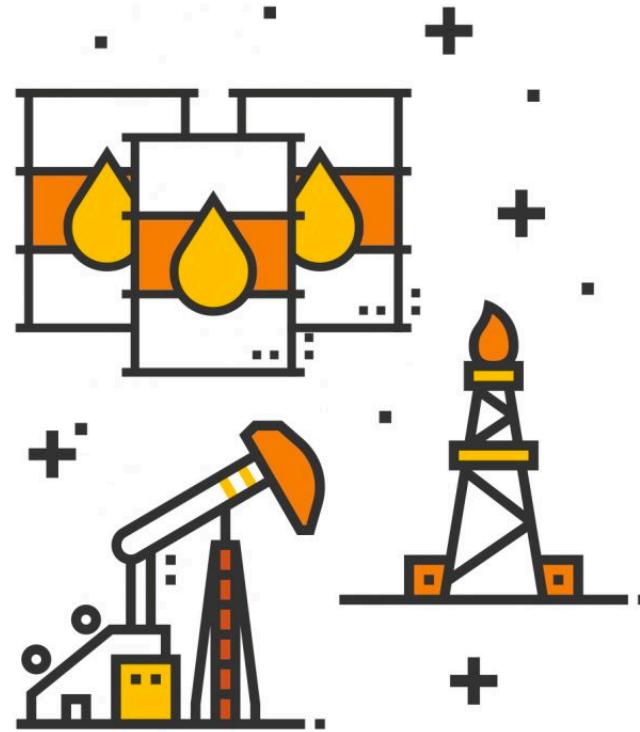
## Downstream



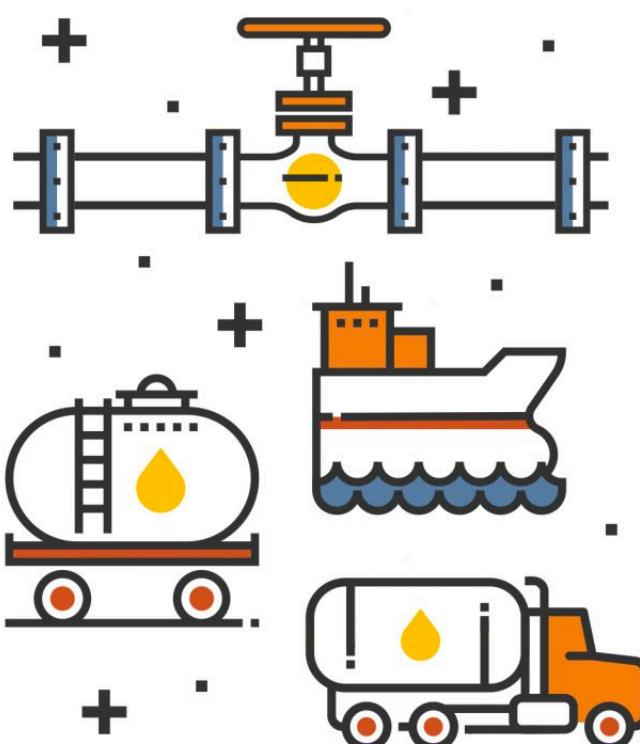
The downstream stage handles processing, selling, marketing and distributing gas or oil. Final products depend upon the initial resource.

'Upstream' is about extracting oil and natural gas from the ground; 'midstream' is about safely moving them thousands of miles; and 'downstream' is converting these resources into the fuels and finished products we all depend on.

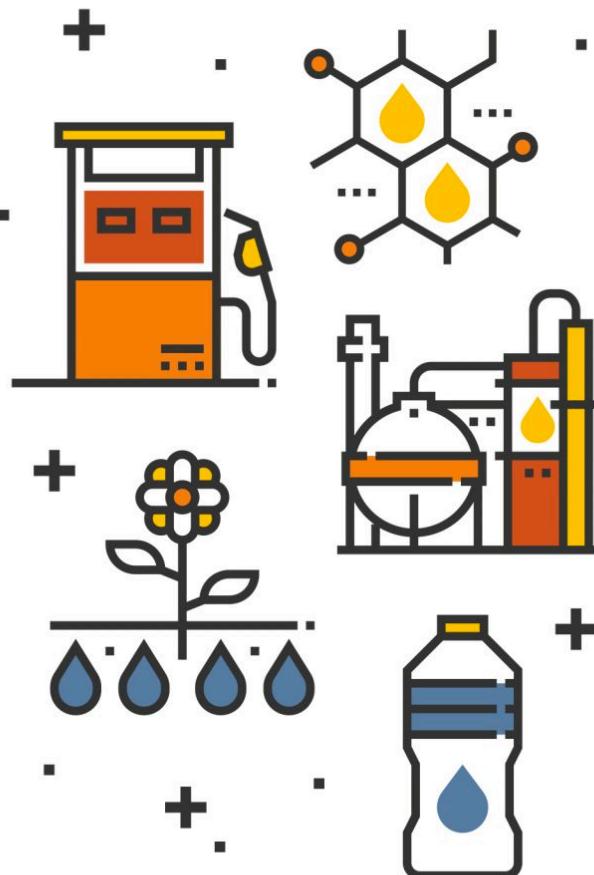
## UPSTREAM



## MIDSTREAM



## DOWNSTREAM



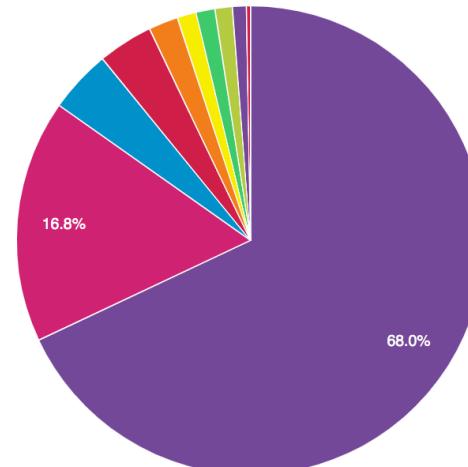
- End to end ML Platform on Kubernetes. Focused on multiple aspects of Model Lifecycle
- Originated at Google, and has grown to have a large community of developers
- Google, IBM, Cisco, RedHat, Intel, Microsoft, AWS and others contributing
- IBM is the 2<sup>nd</sup> contributor in terms of overall commits. IBM maintainers (committers/reviewers) in Katib (HPO+Training), Kubeflow Serving, Manifests, Pipelines etc.

Commits by Company

Show 10 entries Search

| # | Company       | Commits |
|---|---------------|---------|
|   | *Independent  | 3669    |
| 1 | Google        | 904     |
| 2 | IBM           | 234     |
| 3 | Cacloud       | 205     |
| 4 | Alibaba       | 110     |
| 5 | Red Hat       | 71      |
| 6 | Intel         | 65      |
| 7 | Huawei        | 52      |
| 8 | Cisco Systems | 15      |
| 9 | VA Linux      | 12      |

Showing 1 to 10 of 29 entries Previous Next



## Libraries and CLIs - Focus on end users

Arena

kfctl

kubectl

fairing

## Systems - Combine multiple services

katib

pipelines

Model DB

kube  
bench

notebooks

TFX

## Low Level APIs / Services (single function)

TFJob

PyTorchJ  
ob

Pipelines CR

Argo

Jupyter  
CR

MPI CR

Seldon CR

Study Job

Spark  
JobDeveloped By  
KubeflowDeveloped Outside  
Kubeflow

\* Not all components shown

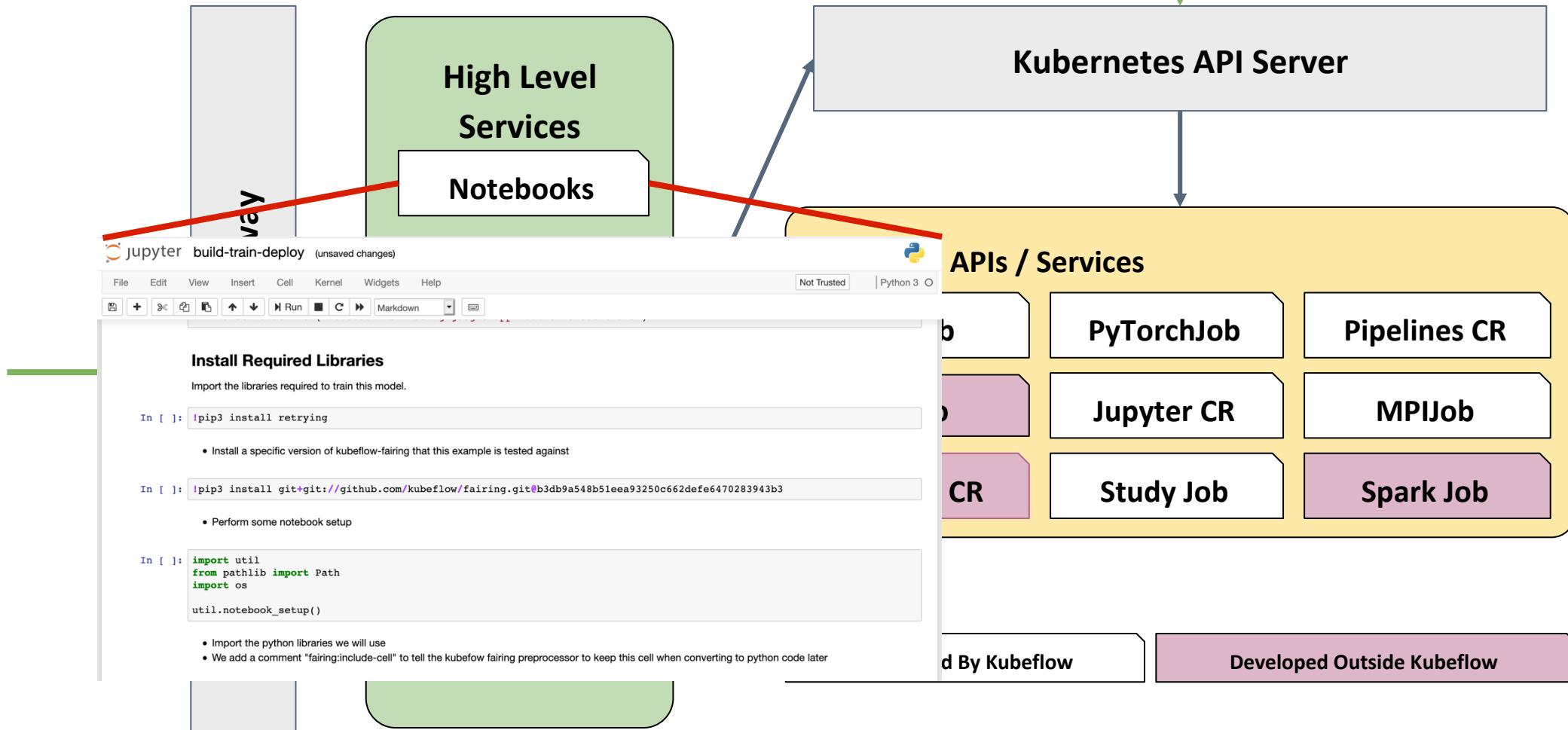
IAM

Orchestration

Scheduling



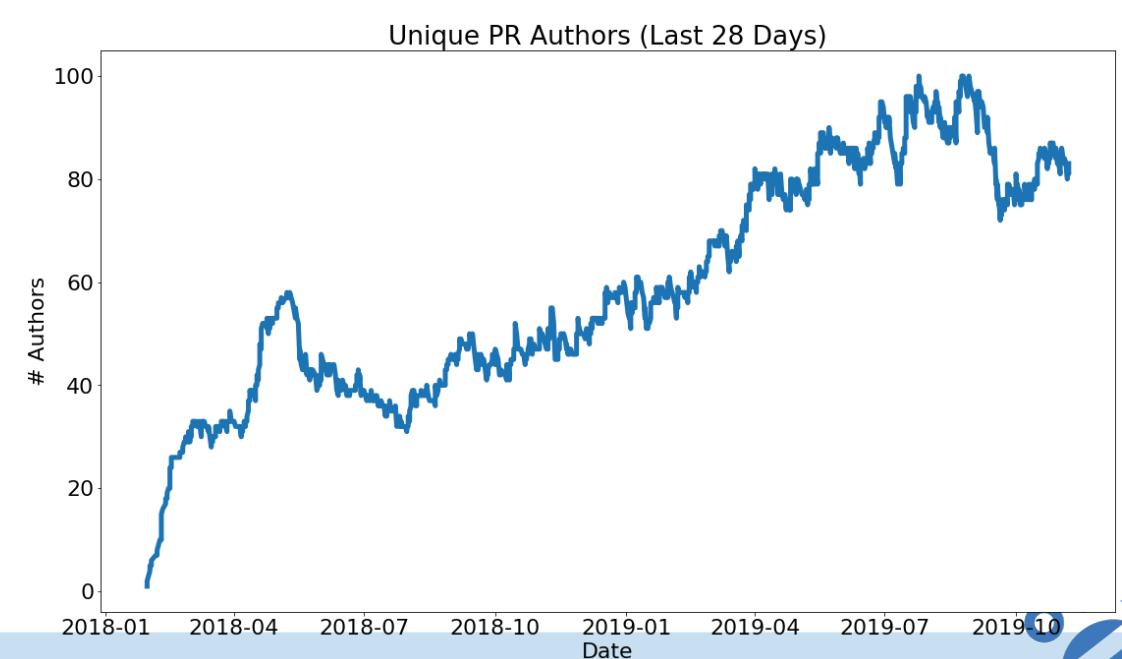
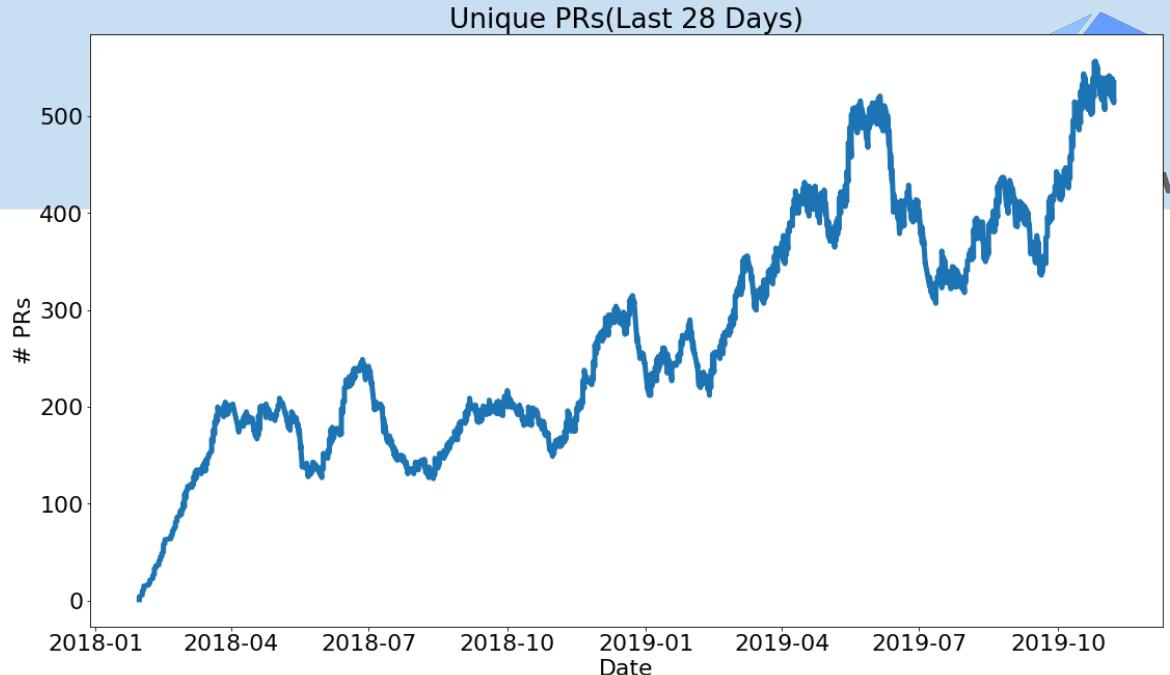
kubectl apply -f tfjob



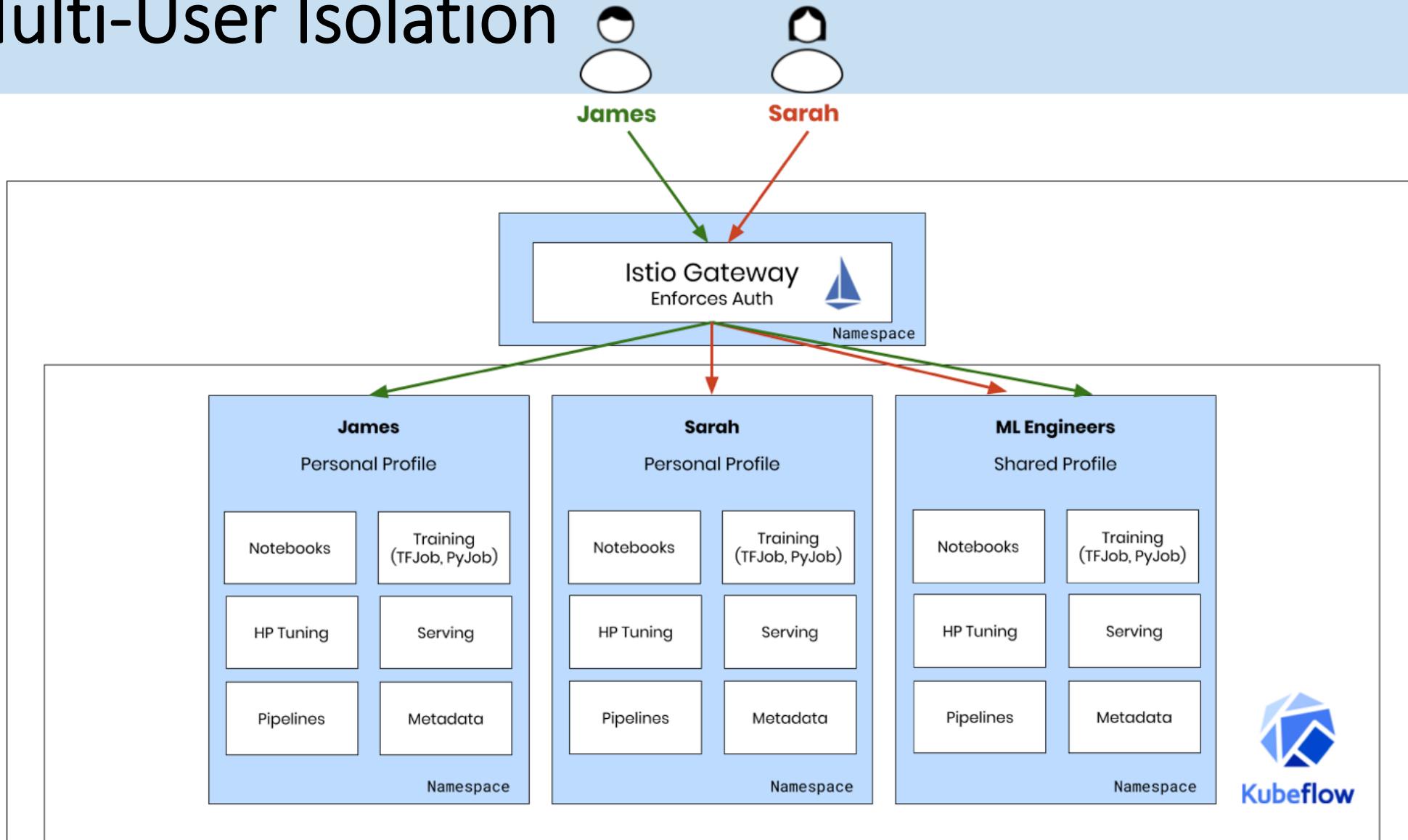
Adapted from Kubeflow Contributor Summit 2019 talk: Kubeflow and ML Landscape (Not all components are shown)



# Community is growing!



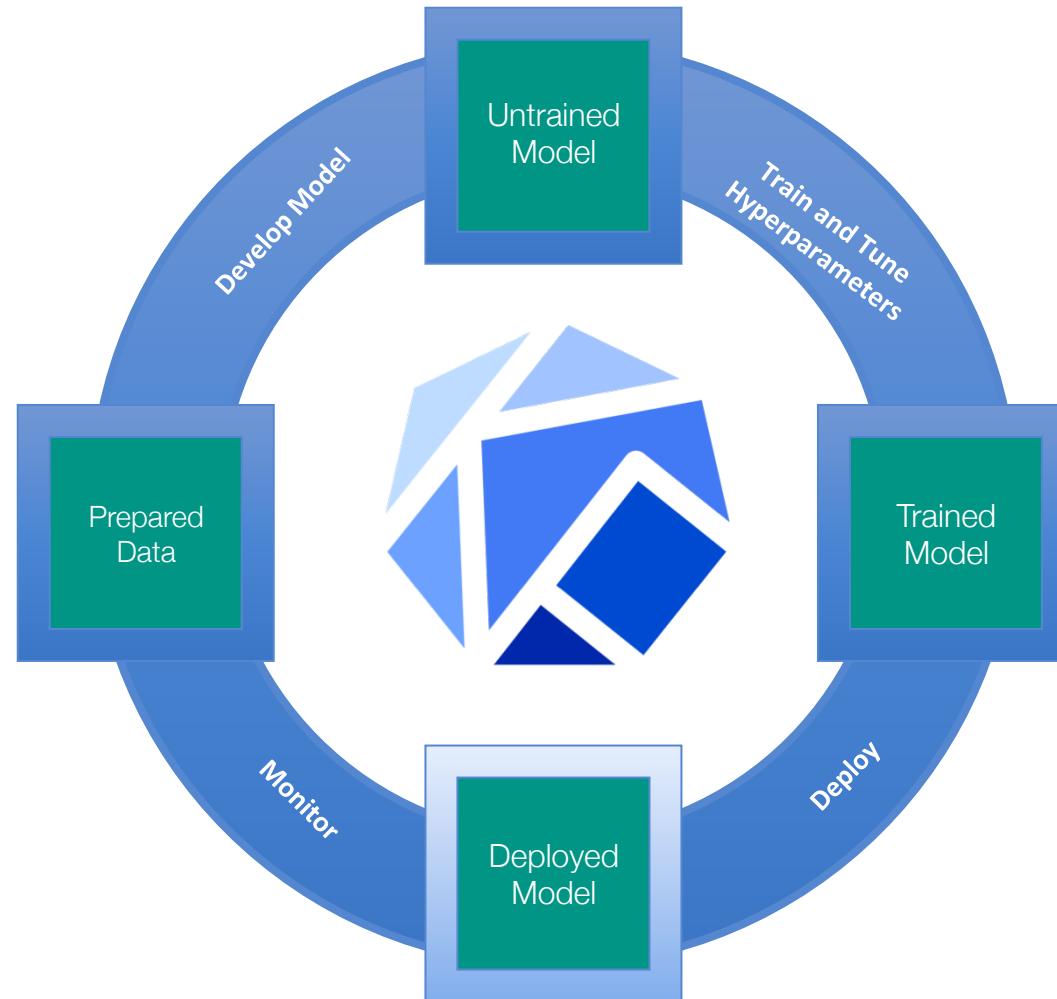
# Multi-User Isolation



kubernetes



# ML Lifecycle: Build: Development, Training and HPO



- Data Scientist
  - Self-service Jupyter Notebooks provide faster model experimentation
  - Simplified configuration of CPU/GPU, RAM, Persistent Volumes
  - Faster model creation with training operators, TFX, magics, workflow automation (Kale, Fairing)
  - Simplify access to external data sources (using stored secrets)
  - Easier protection, faster restoration & sharing of “complete” notebooks
  
- IT Operator
  - Profile Controller, Istio, Dex enable secure RBAC to notebooks, data & resources
  - Smaller base container images for notebooks, fewer crashes, faster to recover



|                                 | TF Operator                                                                                  | PyTorch Operator                                                                            | MPI Operator                                                                                                                            |
|---------------------------------|----------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------|
| Framework Support               |  TensorFlow |  PyTorch | <br>TensorFlow/Keras<br>Apache MXNet/PyTorch/OpenMPI |
| Distribution Strategy & Backend | <code>tf.distribute</code><br>MPI/NCCL/PS/TPU                                                | <code>torch.distributed</code><br>Gloo/MPI/NCCL                                             | <code>horovod</code><br>DistributedOptimizer<br>Gloo/MPI/NCCL                                                                           |



# Distributed Training Operators



## tf-operator

Tools for ML/Tensorflow on Kubernetes.

● Jsonnet ⚙ Apache-2.0 323 ⭐

## pytorch-operator

PyTorch on Kubernetes

● Jsonnet ⚙ Apache-2.0 87 ⭐ 19

## mpi-operator

Kubernetes Operator for Allreduce-style

kubernetes tensorflow mpi dist  
horovod kubeflow

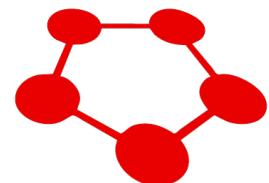
● Go ⚙ Apache-2.0 83 ⭐ 125

## xgboost-operator

Incubating project for xgboost operator

● Go ⚙ Apache-2.0 23 ⭐ 41

# XGBoost



# Chainer

## mxnet-operator

A Kubernetes operator for mxnet jobs

● Go ⚙ Apache-2.0 20 ⭐ 50

## chainer-operator

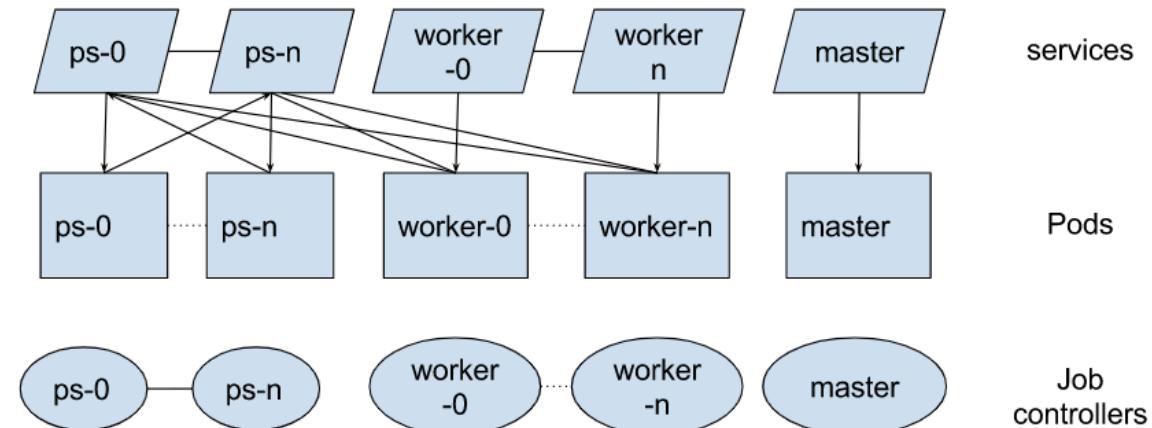
Repository for chainer operator

● Go ⚙ Apache-2.0 9 ⭐ 12



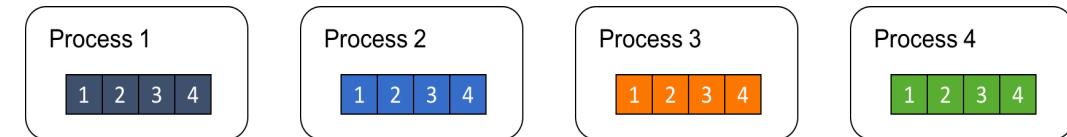
# Distributed Tensorflow Operator

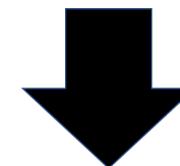
- A distributed Tensorflow Job is collection of the following processes
  - Chief – The chief is responsible for orchestrating training and performing tasks like checkpointing the model
  - Ps – The ps are parameters servers; the servers provide a distributed data store for the model parameters to access
  - Worker – The workers do the actual work of training the model. In some cases, worker 0 might also act as the chief
  - Evaluator - The evaluators can be used to compute evaluation metrics as the model is trained

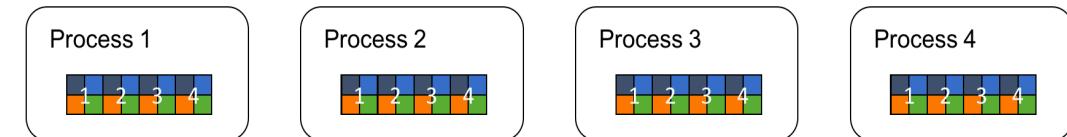


# Distributed MPI Operator - AllReduce

- AllReduce is an operation that reduces many arrays spread across multiple processes into a single array which can be returned to all the processes
- This ensures consistency between distributed processes while allowing all of them to take on different workloads
- The operation used to reduce the multiple arrays back into a single array can vary and that is what makes the different options for AllReduce



 AllReduce



# IBM Hyper Parameter Optimization and Neural Architecture Search - Katib

- Katib: Kubernetes Native System for Automated tuning of machine learning model's Hyperparameter Tuning and Neural Architecture Search.
- Github Repository:  
<https://github.com/kubeflow/katib>
- Hyperparameter Tuning
  - [Random Search](#)
  - [Tree of Parzen Estimators \(TPE\)](#)
  - [Grid Search](#)
  - [Hyperband](#)
  - [Bayesian Optimization](#)
  - [CMA Evolution Strategy](#)
- Neural Architecture Search
  - [Efficient Neural Architecture Search \(ENAS\)](#)
  - [Differentiable Architecture Search \(DARTS\)](#)





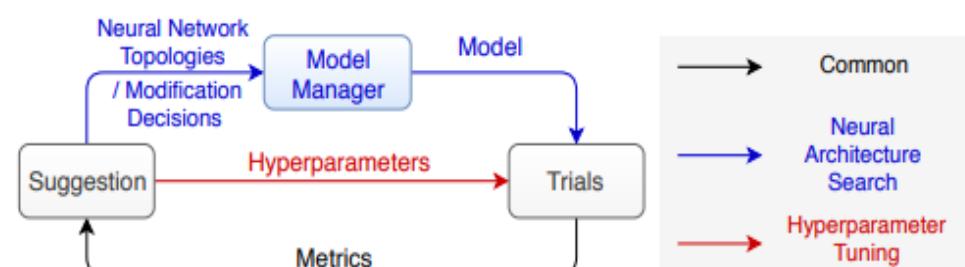
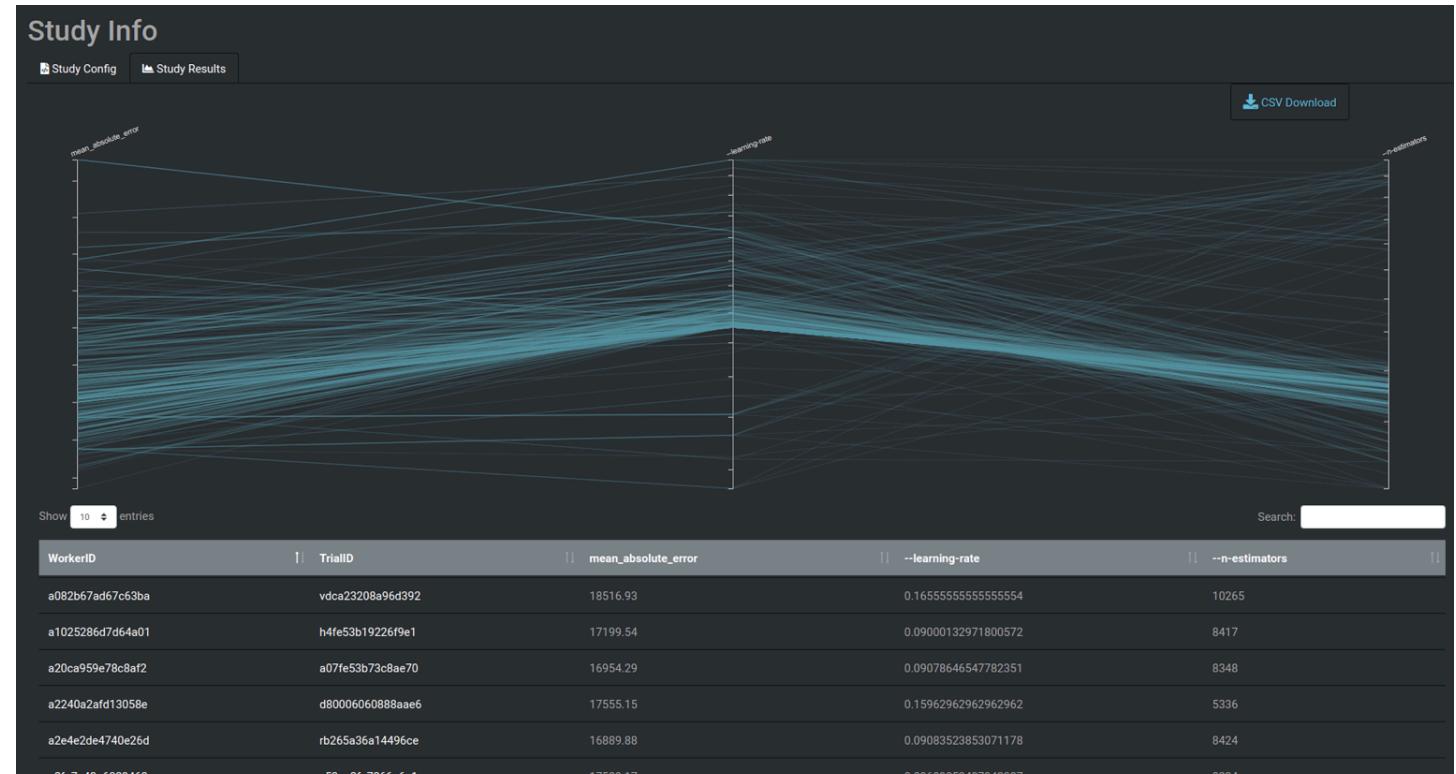
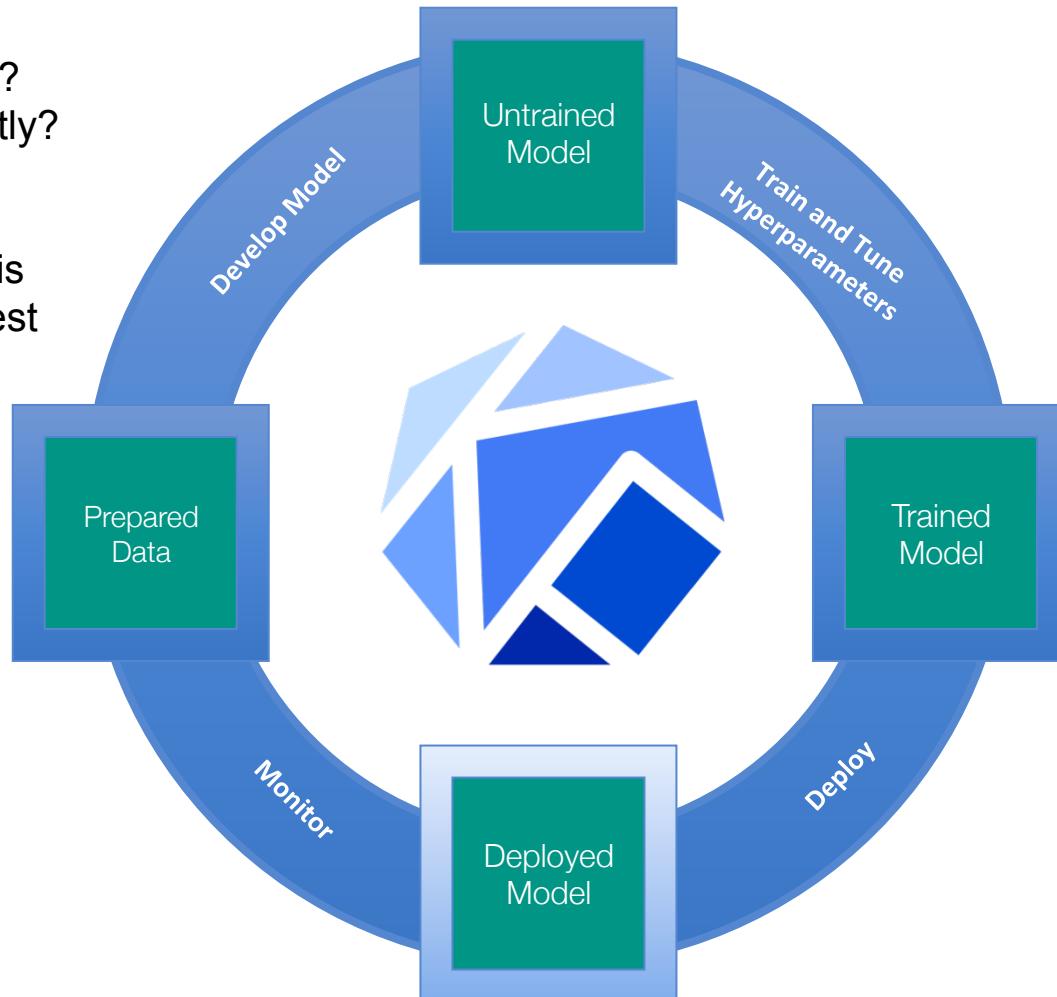


Figure 1: Summary of AutoML workflows

- Cost:  
Is the model over or under scaled?  
Are resources being used efficiently?
- Monitoring:  
Are the endpoints healthy? What is the performance profile and request trace?
- Rollouts:  
Is this rollout safe? How do I roll back? Can I test a change without swapping traffic?
- Protocol Standards:  
How do I make a prediction?  
GRPC? HTTP? Kafka?



- How do I handle batch predictions?
- How do I leverage standardized Data Plane protocol so that I can move my model across MLServing platforms?
- Frameworks:  
How do I serve on Tensorflow?  
XGBoost? Scikit Learn? Pytorch?  
Custom Code?
- Features:  
How do I explain the predictions?  
What about detecting outliers and skew? Bias detection? Adversarial Detection?
- How do I wire up custom pre and post processing



- Seldon Core was pioneering Graph Inferencing.
- IBM and Bloomberg were exploring serverless ML lambdas. IBM gave a talk on the ML Serving with Knative at last KubeCon in Seattle
- Google had built a common Tensorflow HTTP API for models.
- Microsoft Kuberntizing their Azure ML Stack



Bloomberg





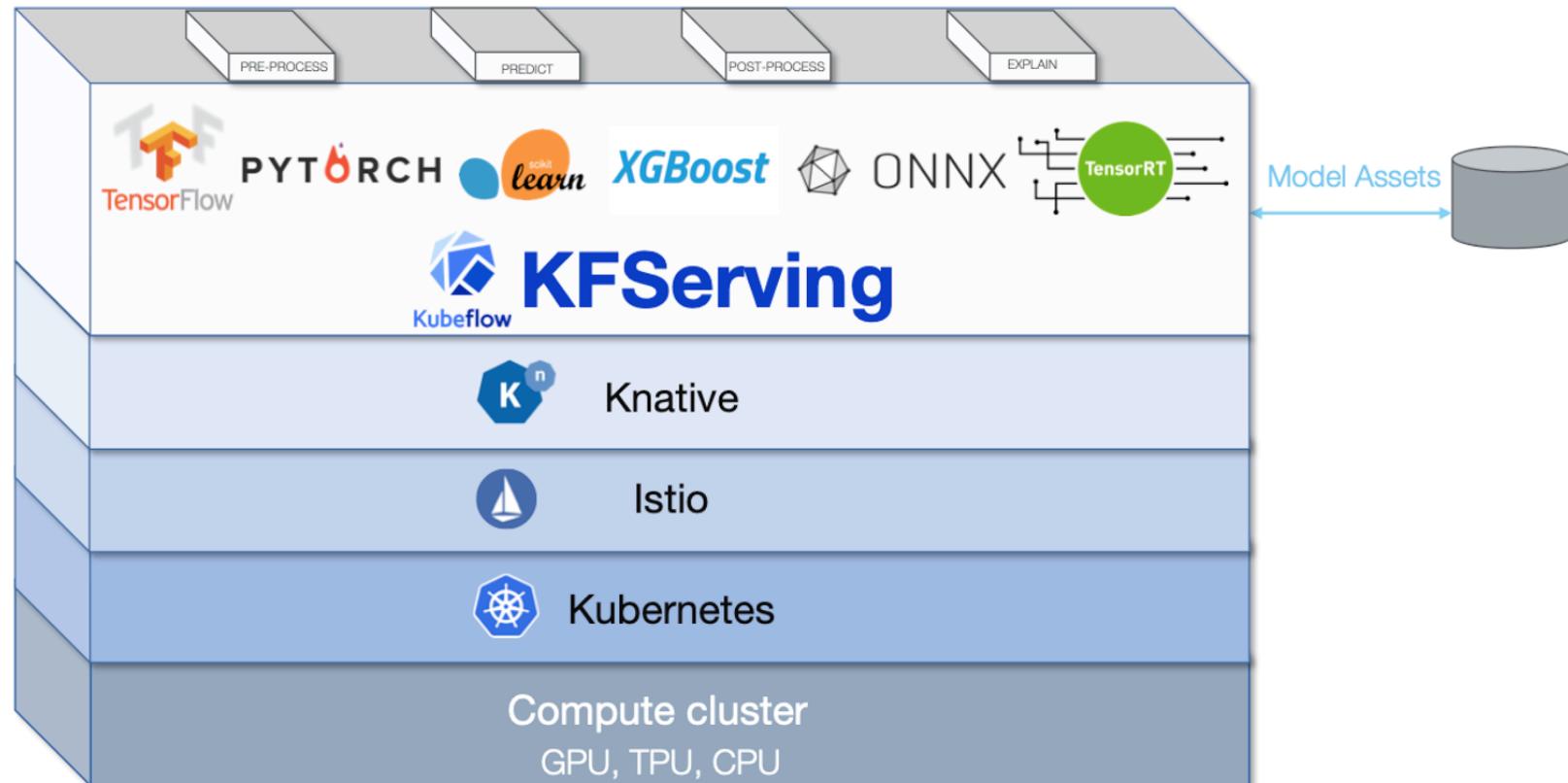
# Putting the pieces together



- Kubeflow created the conditions for collaboration.
- A promise of open code and open community.
- Shared responsibilities and expertise across multiple companies.
- Diverse requirements from different customer segments



- Founded by Google, Seldon, IBM, Bloomberg and Microsoft
- Part of the Kubeflow project
- Focus on 80% use cases - single model rollout and update
- Kfserving 1.0 goals:
  - Serverless ML Inference
  - Canary rollouts
  - Model Explanations
  - Optional Pre/Post processing

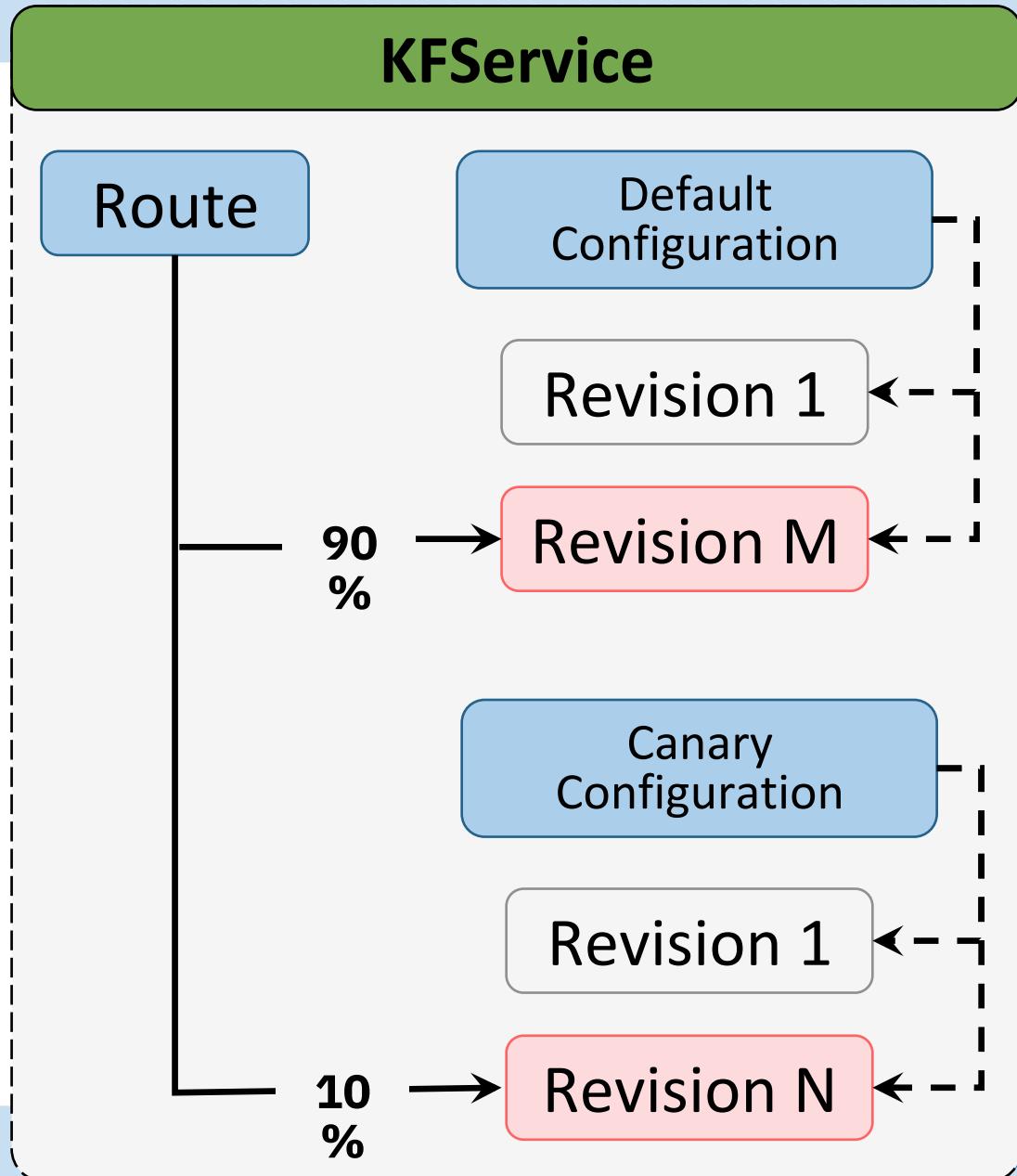




# IBM KFServing: Default and Canary Configurations

Manages the hosting aspects of your models

- **InferenceService** - manages the lifecycle of models
- **Configuration** - manages history of model deployments. Two configurations for default and canary.
- **Revision** - A snapshot of your model version
- **Route** - Endpoint and network traffic management



## Model Servers

- TensorFlow
- Nvidia TRTIS
- PyTorch
- XGBoost
- SKLearn
- ONNX

## Components:

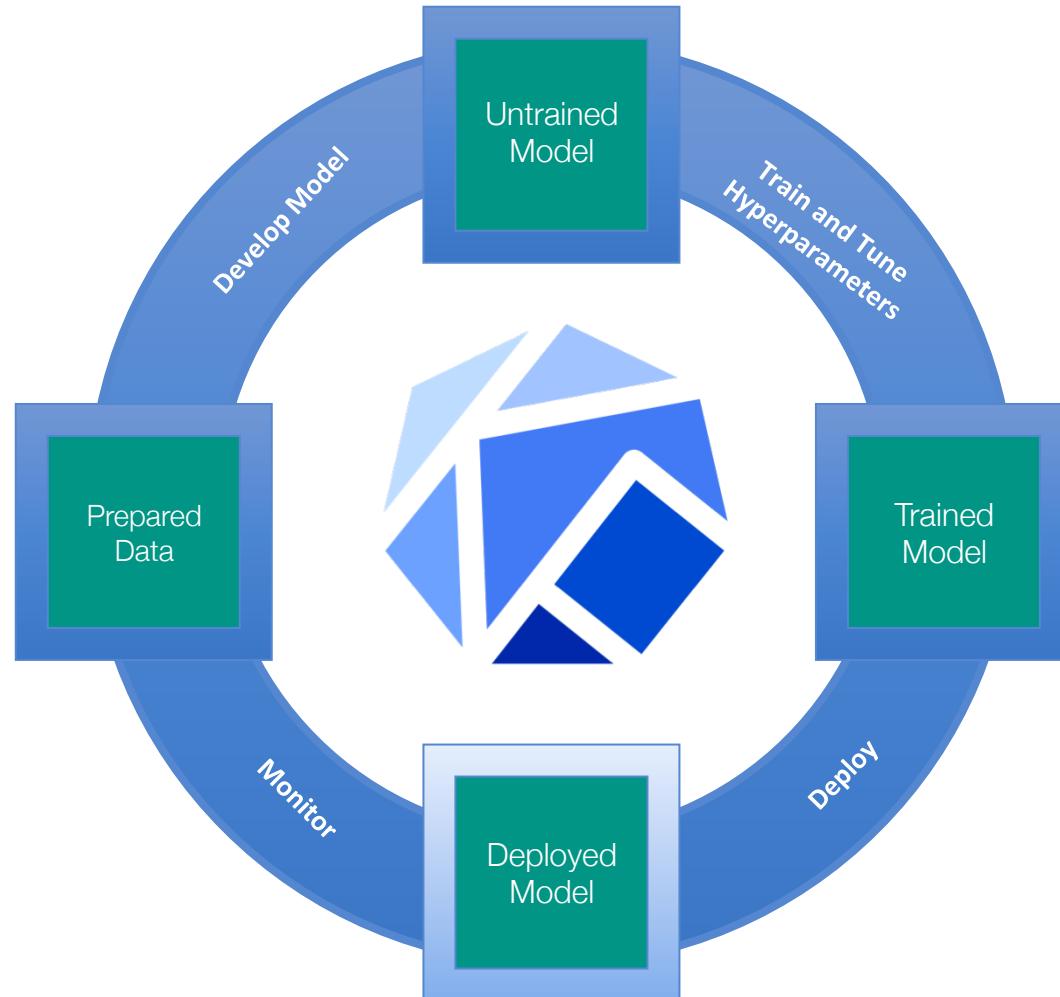
- - Predictor, Explainer, Transformer (pre-processor, post-processor)

## Storage

- AWS/S3
- GCS
- Azure Blob
- PVC

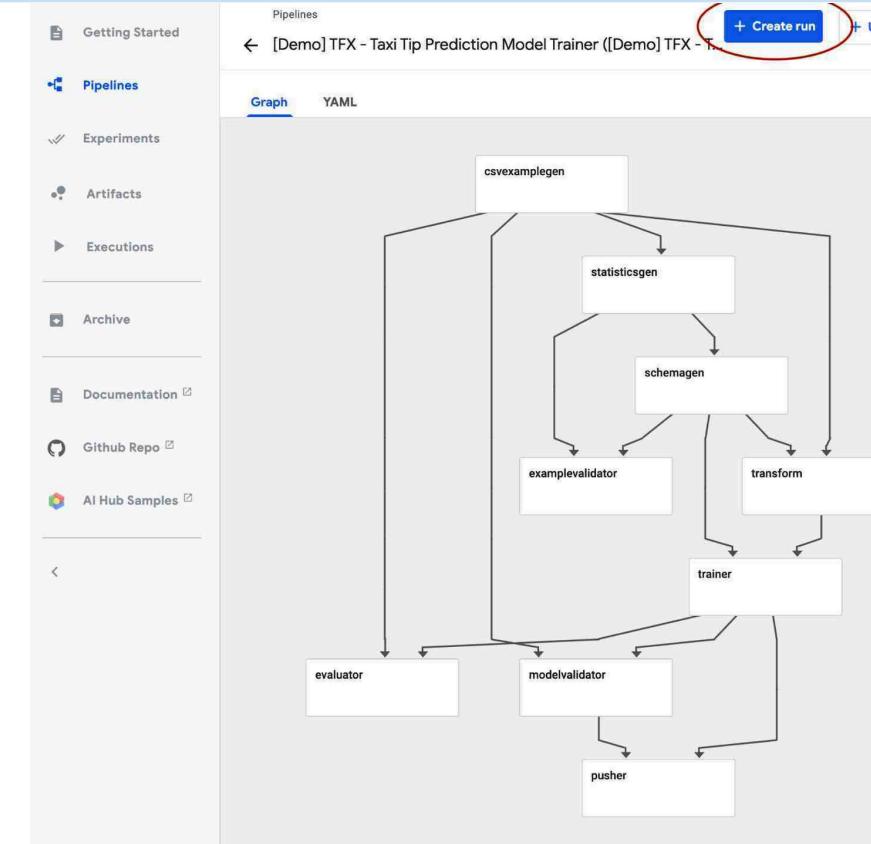


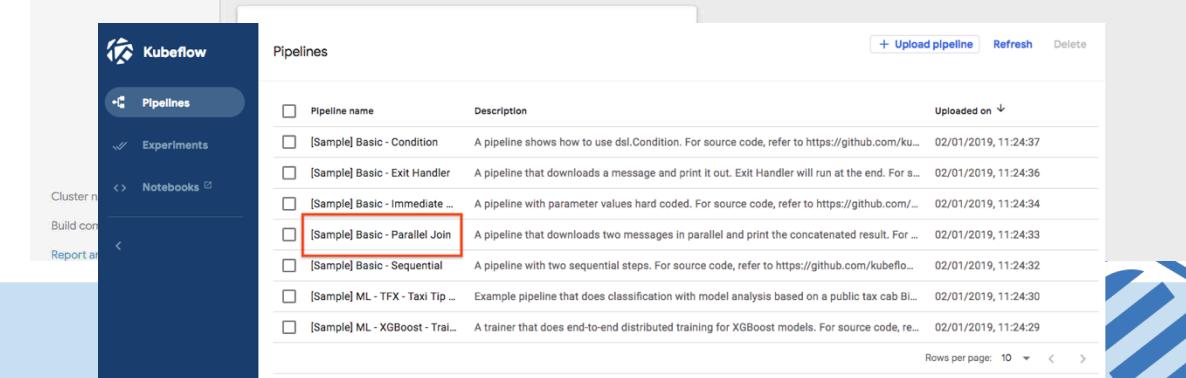
# ML Lifecycle: Orchestrate Build, Train, Validate and Deploy



# Kubeflow Pipelines

- Containerized implementations of ML Tasks
  - Pre-built components: Just provide params or code snippets (e.g. training code)
  - Create your own components from code or libraries
  - Use any runtime, framework, data types
  - Attach k8s objects - volumes, secrets
- Specification of the sequence of steps
  - Specified via Python DSL
  - Inferred from data dependencies on input/output
- Input Parameters
  - A “Run” = Pipeline invoked w/ specific parameters
  - Can be cloned with different parameters
- Schedules
  - Invoke a single run or create a recurring scheduled pipeline





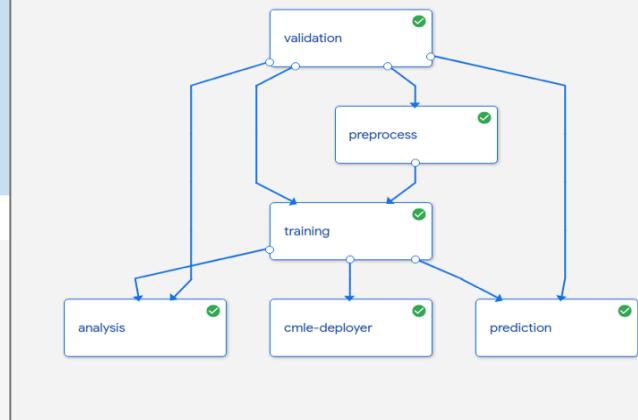
| Pipeline name                    | Description                                                                                     | Uploaded on          |
|----------------------------------|-------------------------------------------------------------------------------------------------|----------------------|
| [Sample] Basic - Condition       | A pipeline shows how to use dsl.Condition. For source code, refer to https://github.com/ku...   | 02/01/2019, 11:24:37 |
| [Sample] Basic - Exit Handler    | A pipeline that downloads a message and print it out. Exit Handler will run at the end. For ... | 02/01/2019, 11:24:36 |
| [Sample] Basic - Immediate ...   | A pipeline with parameter values hard coded. For source code, refer to https://github.com/...   | 02/01/2019, 11:24:34 |
| [Sample] Basic - Parallel Join   | A pipeline that downloads two messages in parallel and print the concatenated result. For ...   | 02/01/2019, 11:24:33 |
| [Sample] Basic - Sequential      | A pipeline with two sequential steps. For source code, refer to https://github.com/kubeflo...   | 02/01/2019, 11:24:32 |
| [Sample] ML - TFX - Taxi Tip ... | Example pipeline that does classification with model analysis based on a public tax cab Bl...   | 02/01/2019, 11:24:30 |
| [Sample] ML - XGBoost - Trai...  | A trainer that does end-to-end distributed training for XGBoost models. For source code, re...  | 02/01/2019, 11:24:29 |



# Define Pipeline with Python SDK

```
@dsl.pipeline(name='Taxi Cab Classification Pipeline Example')
def taxi_cab_classification(
    output_dir,
    project,
    Train_data      = 'gs://bucket/train.csv',
    Evaluation_data = 'gs://bucket/eval.csv',
    Target          = 'tips',
    Learning_rate   = 0.1, hidden_layer_size = '100,50', steps=3000):

    tfdv           = TfdvOp(train_data, evaluation_data, project, output_dir)
    preprocess     = PreprocessOp(train_data, evaluation_data, tfdv.output["schema"], project, output_dir)
    training       = DnnTrainerOp(preprocess.output, tfdv.schema, learning_rate, hidden_layer_size, steps,
                                target, output_dir)
    tfma           = TfmaOp(training.output, evaluation_data, tfdv.schema, project, output_dir)
    deploy         = TfServingDeployerOp(training.output)
```



## Compile and Submit Pipeline Run

```
dsl.compile(taxi_cab_classification, 'tfx.tar.gz')
run = client.run_pipeline(
    'tfx_run', 'tfx.tar.gz', params={'output': 'gs://dpa22', 'project': 'my-project-33'})
```



# Visualize the state of various components

Pipelines  
Experiments **Artifacts**  
Executions  
Archive  
Documentation  
Github Repo  
AI Hub Samples

Cluster name: cluster-4  
Build commit: 743746b  
Report an Issue

Graph Run output Config

csvexamplegen → statisticsgen → schemagen → examplevalidator → evaluator → pusher

resolvernode-lates... → evaluator

train → evaluator

Static HTML

Sort by Feature ▾ Reverse order Feature search (...)

Features:  int(8)  float(7)  string(2)  
 unknown(1)

| Numeric Features (15)           |         |        |         |
|---------------------------------|---------|--------|---------|
| count                           | missing | mean   | std dev |
| dropoff_census_tract<br>3,618   | 28.93%  | 17.0B  | 333k    |
| dropoff_community_area<br>4,905 | 3.65%   | 21.2   | 17.85   |
| dropoff_latitude<br>4,915       | 3.46%   | 41.9   | 0.04    |
| dropoff_longitude<br>4,915      | 3.46%   | -87.65 | 0.06    |

Runtime execution graph. Only steps that are currently running or have a

## Pipelines

[+ Upload pipeline](#)[Refresh](#)[Delete](#)

Filter pipelines



| <input type="checkbox"/> | Pipeline name                                                               | Description                                                                                                                                          | Uploaded on           |
|--------------------------|-----------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------|
| <input type="checkbox"/> | [Tutorial] DSL - Control structures                                         | <a href="#">source code</a> Shows how to use conditional execution and exit handlers. This pipeline will randomly fail to demonstra...               | 2/20/2020, 3:28:12 PM |
| <input type="checkbox"/> | [Tutorial] Data passing in python com...                                    | <a href="#">source code</a> Shows how to pass data between python components.                                                                        | 2/20/2020, 3:28:11 PM |
| <input type="checkbox"/> | [Demo] TFX - Taxi Tip Prediction Mod...                                     | <a href="#">source code</a> <a href="#">GCP Permission requirements</a> . Example pipeline that does classification with model analysis based on ... | 2/20/2020, 3:28:10 PM |
| <input type="checkbox"/> | Version name                                                                |                                                                                                                                                      | Uploaded on           |
| <input type="checkbox"/> | TFX - Taxi Tip Prediction Model Trainer_version_at_2020-03-03T15:44:30.197Z |                                                                                                                                                      | 3/3/2020, 7:55:03 AM  |
| <input type="checkbox"/> | [Demo] TFX - Taxi Tip Prediction Model Trainer                              |                                                                                                                                                      | 2/20/2020, 3:28:10 PM |
|                          |                                                                             |                                                                                                                                                      | Rows per page: 10 < > |
| <input type="checkbox"/> | [Demo] XGBoost - Training with Confu...                                     | <a href="#">source code</a> <a href="#">GCP Permission requirements</a> . A trainer that does end-to-end distributed training for XGBoost models.    | 2/20/2020, 3:28:09 PM |
|                          |                                                                             |                                                                                                                                                      | Rows per page: 10 < > |

Pipelines lets you group and manage multiple versions of a pipeline.



**Artifacts**

| Pipeline/Workspace ↑          | Name             | ID | Type              | URI                                                                                     | Created at           |
|-------------------------------|------------------|----|-------------------|-----------------------------------------------------------------------------------------|----------------------|
|                               |                  | 1  | ExternalArtifact  | <a href="gs://ml-pipeline-playground/tfx_t...">gs://ml-pipeline-playground/tfx_t...</a> | 2020-02-20T15:10:00Z |
| taxi_pipeline_with_parameters | examples         | 2  | Examples          | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | statistics       | 3  | ExampleStatistics | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | schema           | 4  | Schema            | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | anomalies        | 5  | ExampleAnomalies  | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | transform_graph  | 6  | TransformGraph    | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | transformed_e... | 7  | Examples          | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | model            | 8  | Model             | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | evaluation       | 9  | ModelEvaluation   | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | blessing         | 10 | ModelBlessing     | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | pushed_model     | 11 | PushedModel       | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |
|                               | evaluation       | 12 | ModelEvaluation   | <a href="gs://aju-pipelines/tfx_taxi_simpl...">gs://aju-pipelines/tfx_taxi_simpl...</a> | 2020-02-20T15:10:00Z |

**model**

Overview Lineage Explorer

Type: Model

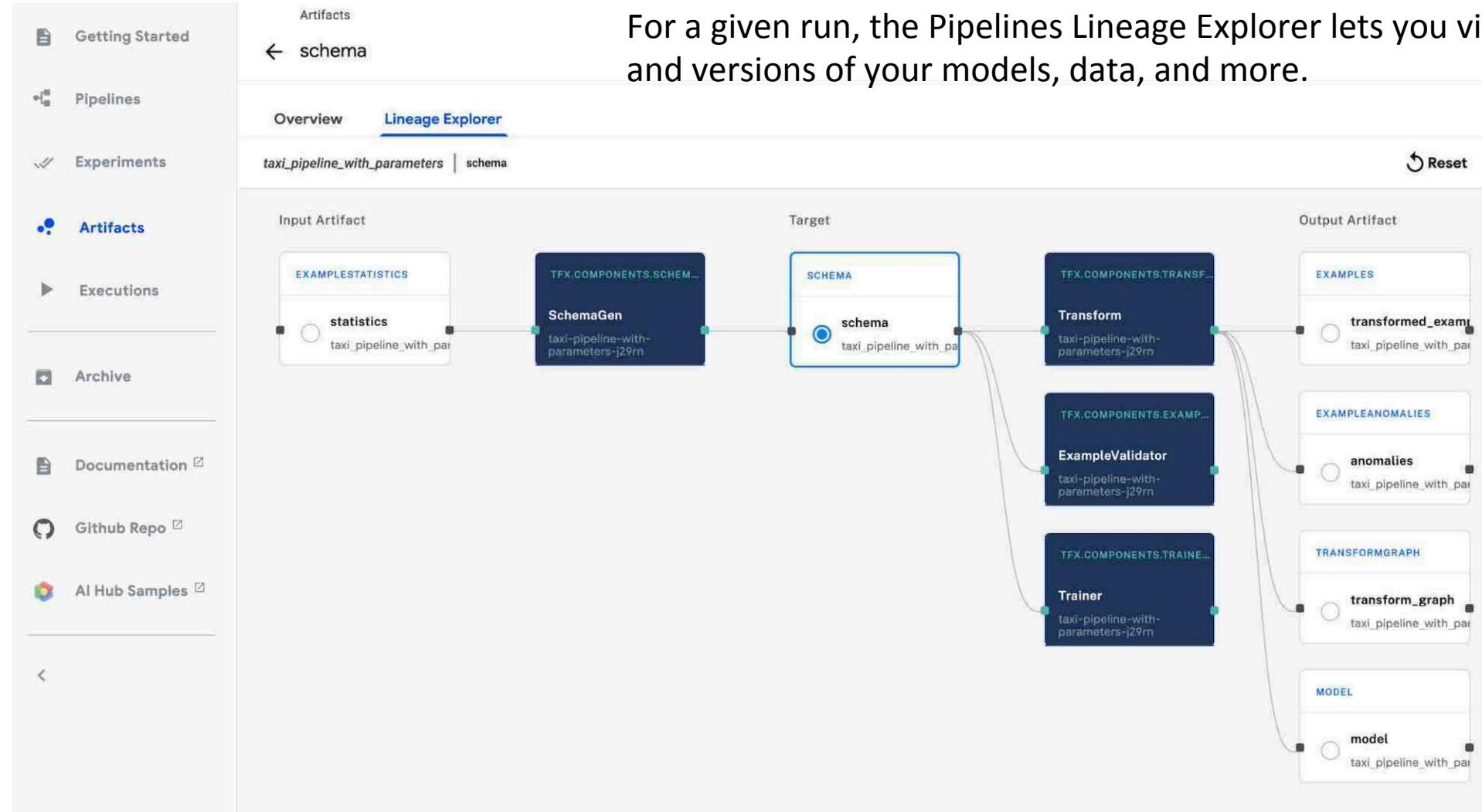
URI: [gs://aju-pipelines/tfx\\_taxi\\_simple/85265540-6a06-4969-a49e-1f65741878be/Trainer/model/7](gs://aju-pipelines/tfx_taxi_simple/85265540-6a06-4969-a49e-1f65741878be/Trainer/model/7)

Properties

Custom Properties

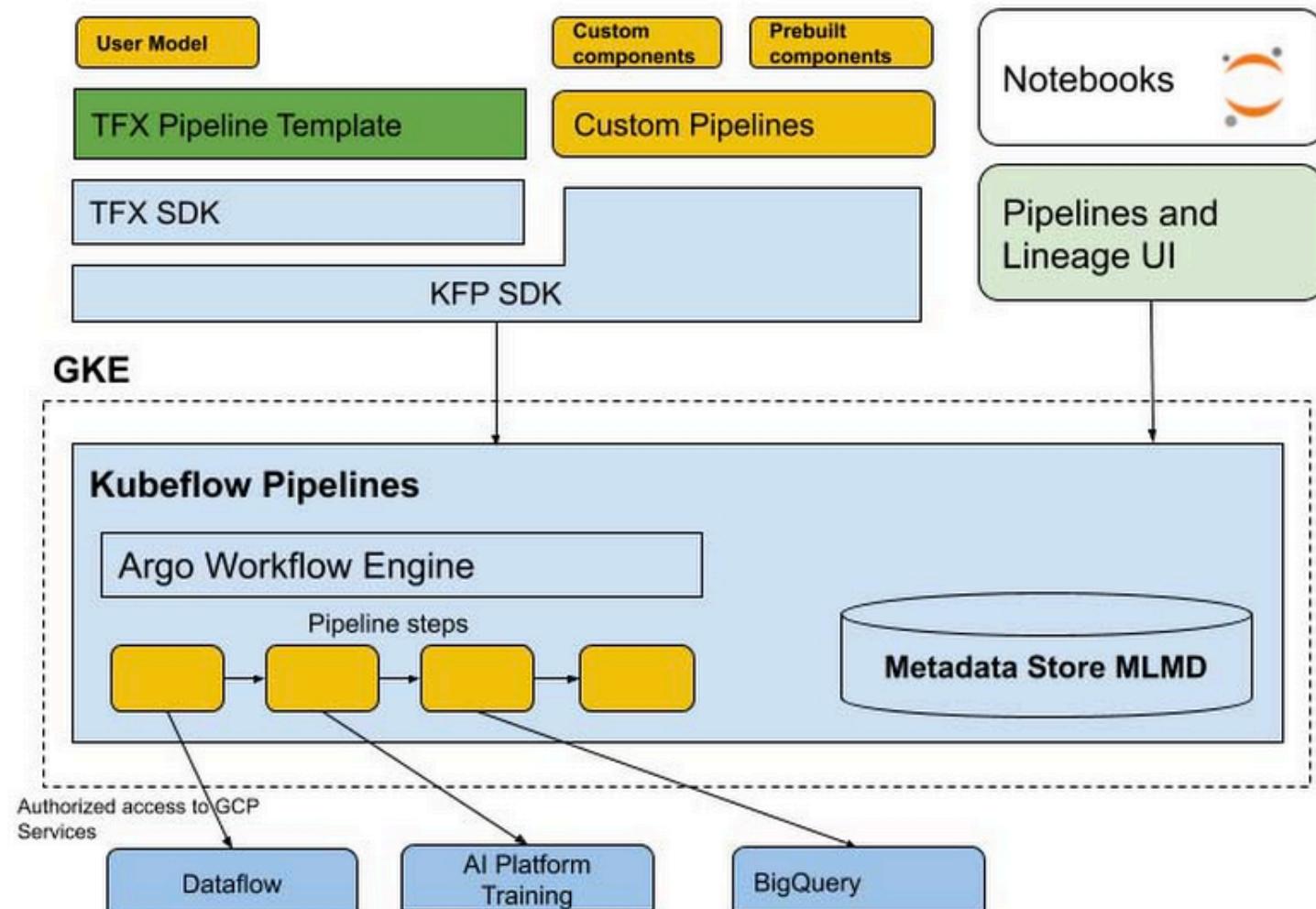
|       |                               |                    |           |
|-------|-------------------------------|--------------------|-----------|
| name  | pipeline_name                 | producer_component | state     |
| model | taxi_pipeline_with_parameters | Trainer            | published |

Artifacts for a run of the “TFX Taxi Trip” example pipeline. For each artifact, you can view details and get the artifact URL—in this case, for the model.



For a given run, the Pipelines Lineage Explorer lets you view the history and versions of your models, data, and more.

# Kubeflow Pipeline Architecture





**seldon**

**Spark**

jupyter

Kubernetes  
Ready

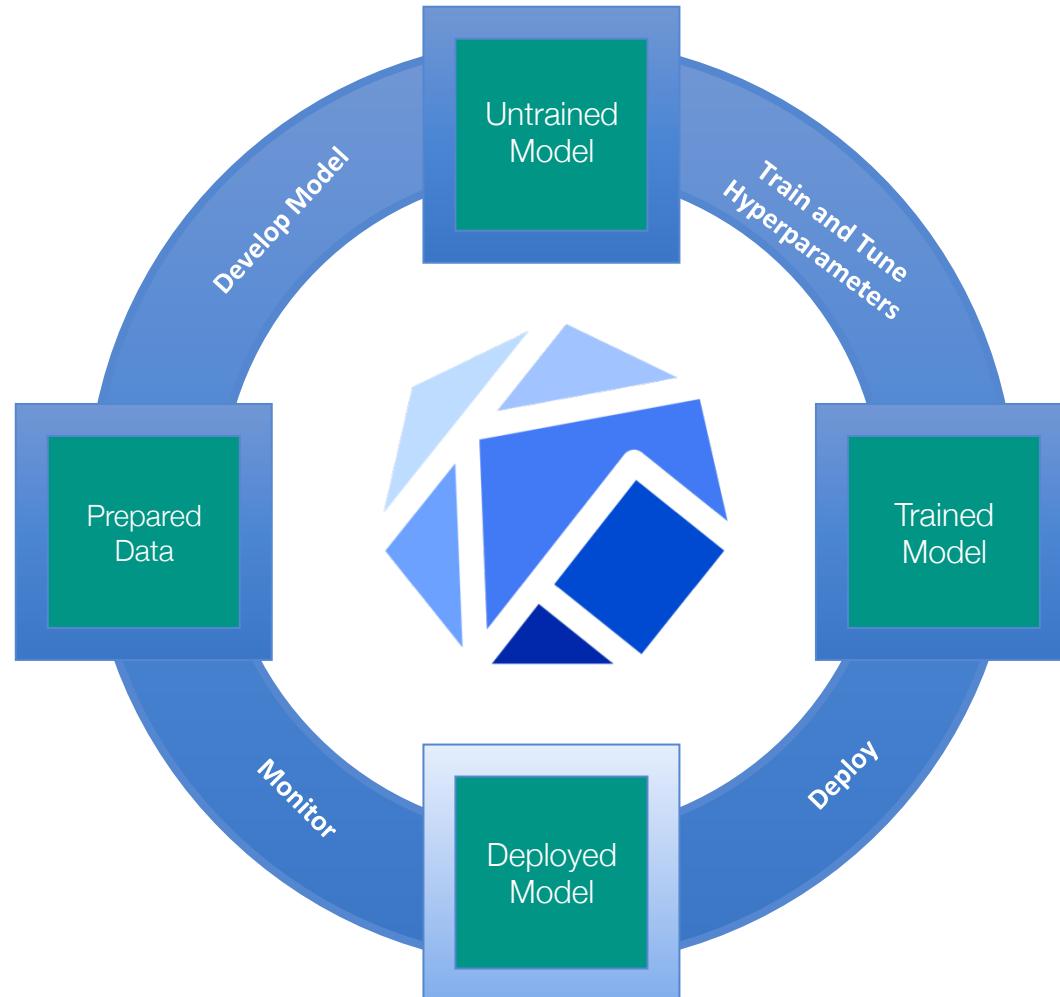


**Kubeflow**

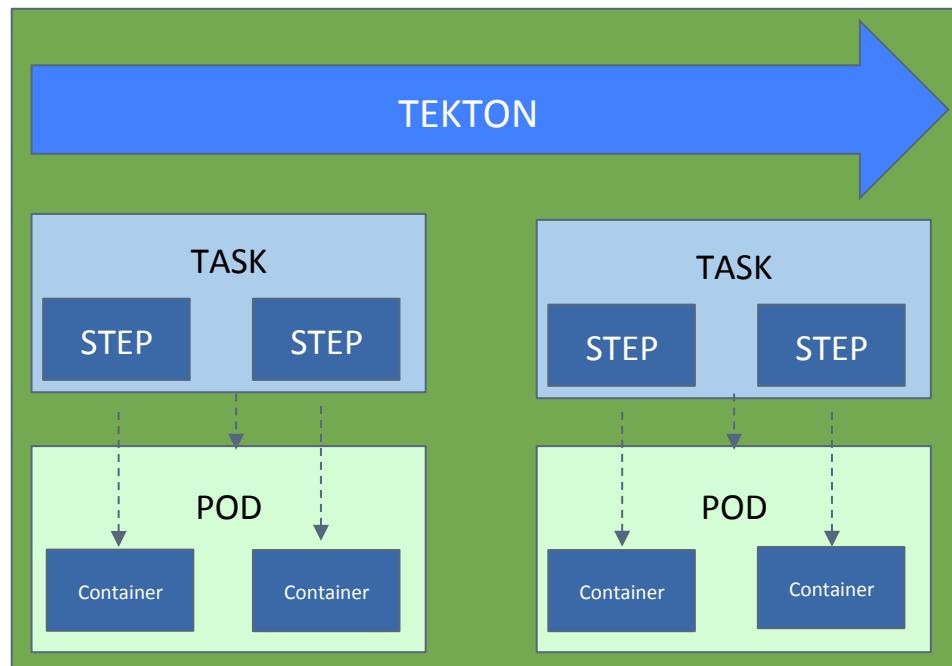
**ML and AI Platform**



# Kubeflow Pipelines with Tekton

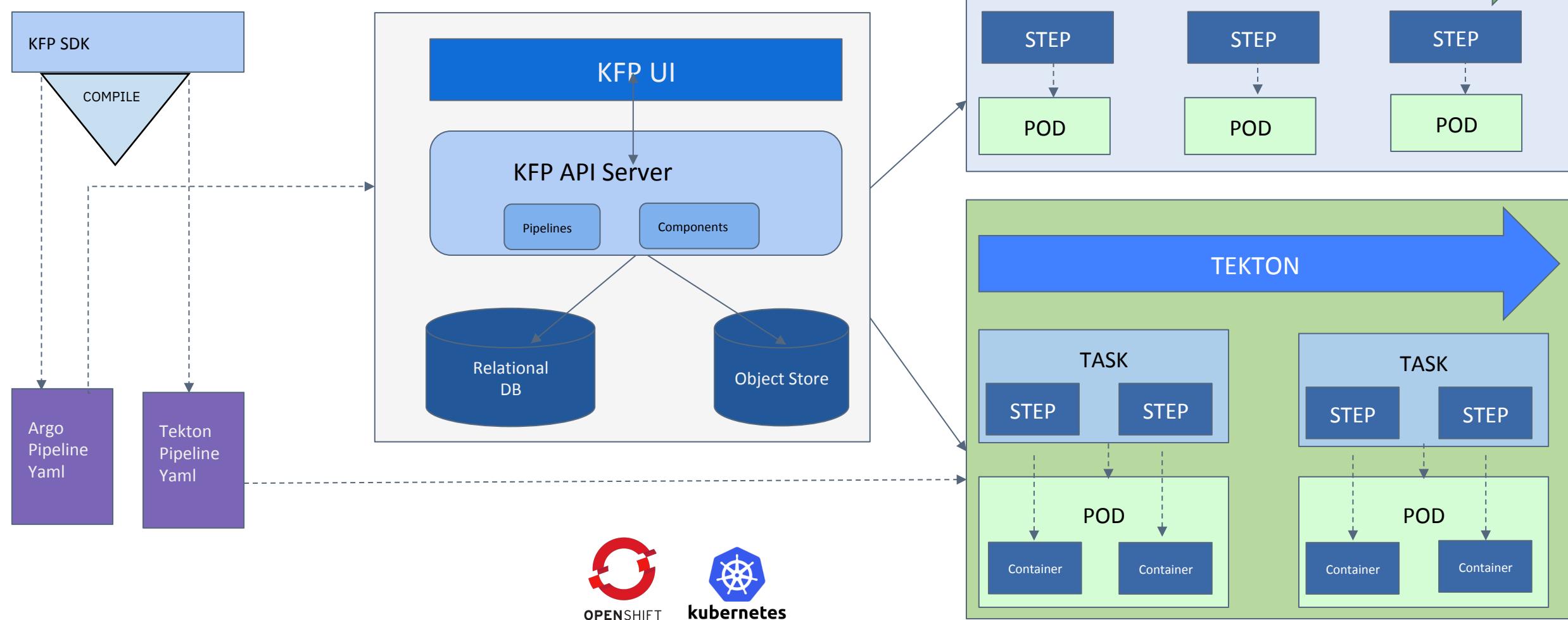


- The Tekton Pipelines project provides Kubernetes-style resources for declaring CI/CD-style pipelines.
- Tekton introduces several new CRDs including Task, Pipeline, TaskRun, and PipelineRun.
- A PipelineRun represents a single running instance of a Pipeline and is responsible for creating a Pod for each of its Tasks and as many containers within each Pod as it has Steps.



- A **PipelineResource** defines an object that is an input (such as a git repository) or an output (such as a docker image) of the pipeline.
- A **PipelineRun** defines an execution of a pipeline. It references the Pipeline to run and the PipelineResources to use as inputs and outputs.
- A **Pipeline** defines the set of Tasks that compose a pipeline.
- A **Task** defines a set of build Steps such as compiling code, running tests, and building and deploying images.

# KFP – Tekton Phase One



Pluggable Components

Spark

Watson Studio

WML

Open Scale

Kubeflow Training

Seldon

AIF360

ART

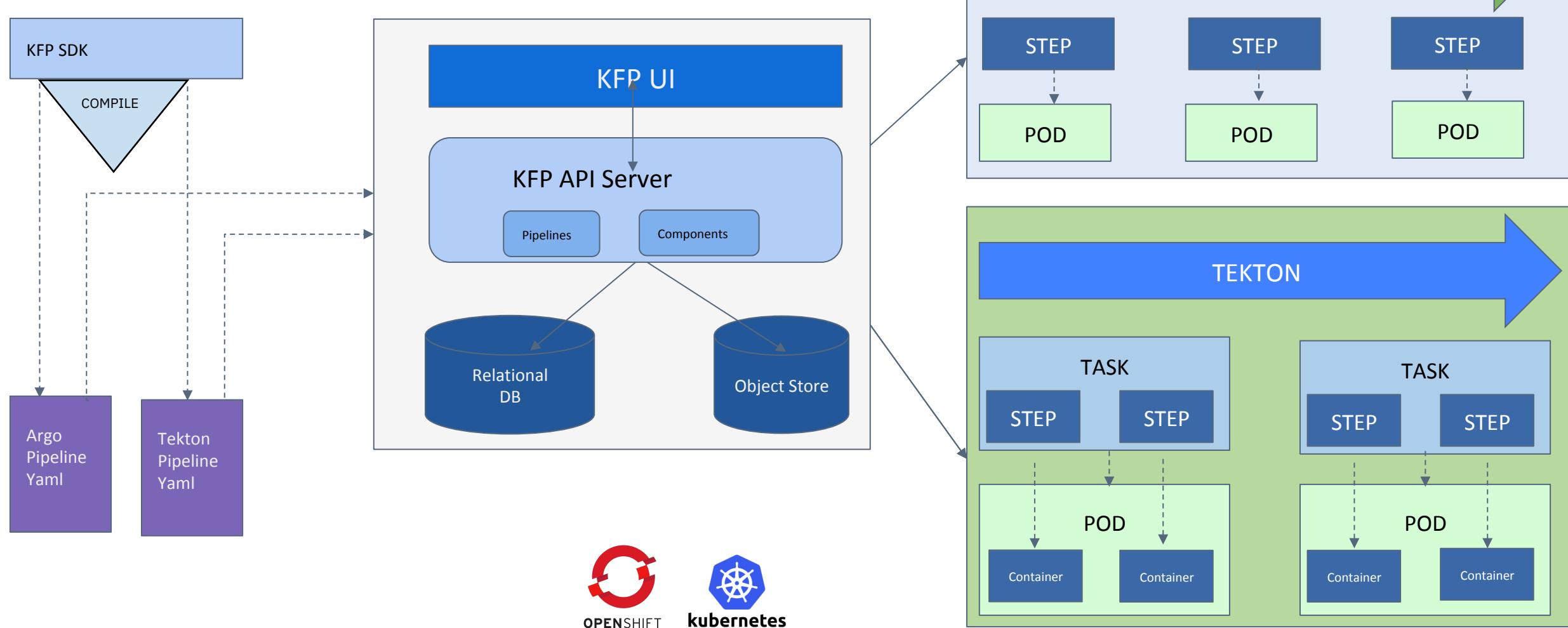
KATIB

KFSERVING

...



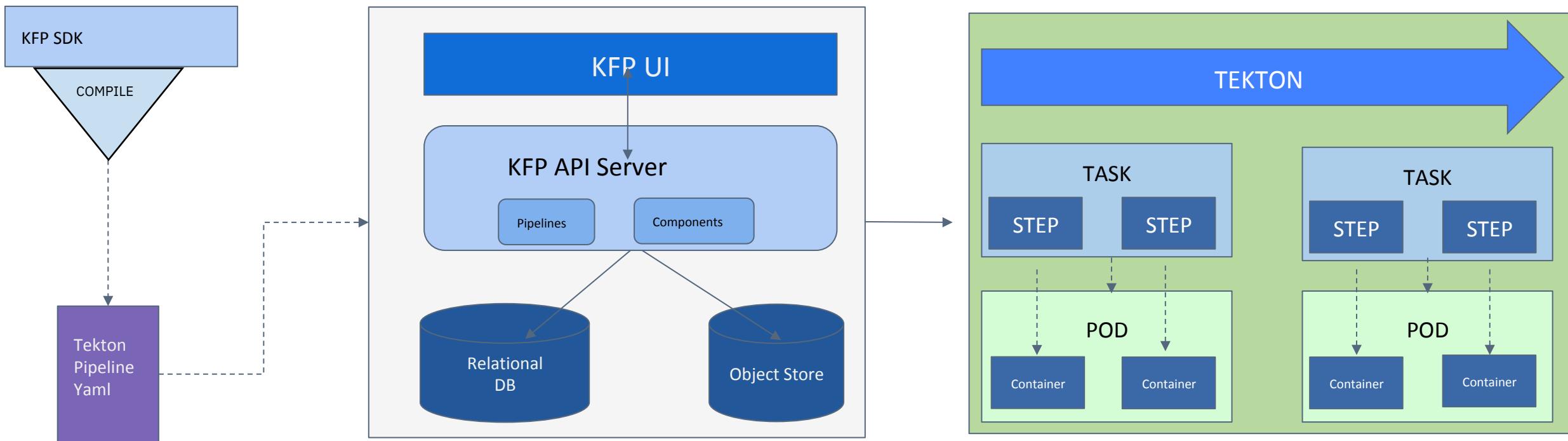
# KFP – Tekton Phase Two



Pluggable Components



# KFP – Tekton: Current Implementation



OPENSIFT



kubernetes



Pluggable Components

Spark

Watson Studio

WML

Open Scale

Kubeflow Training

Seldon

AIF360

ART

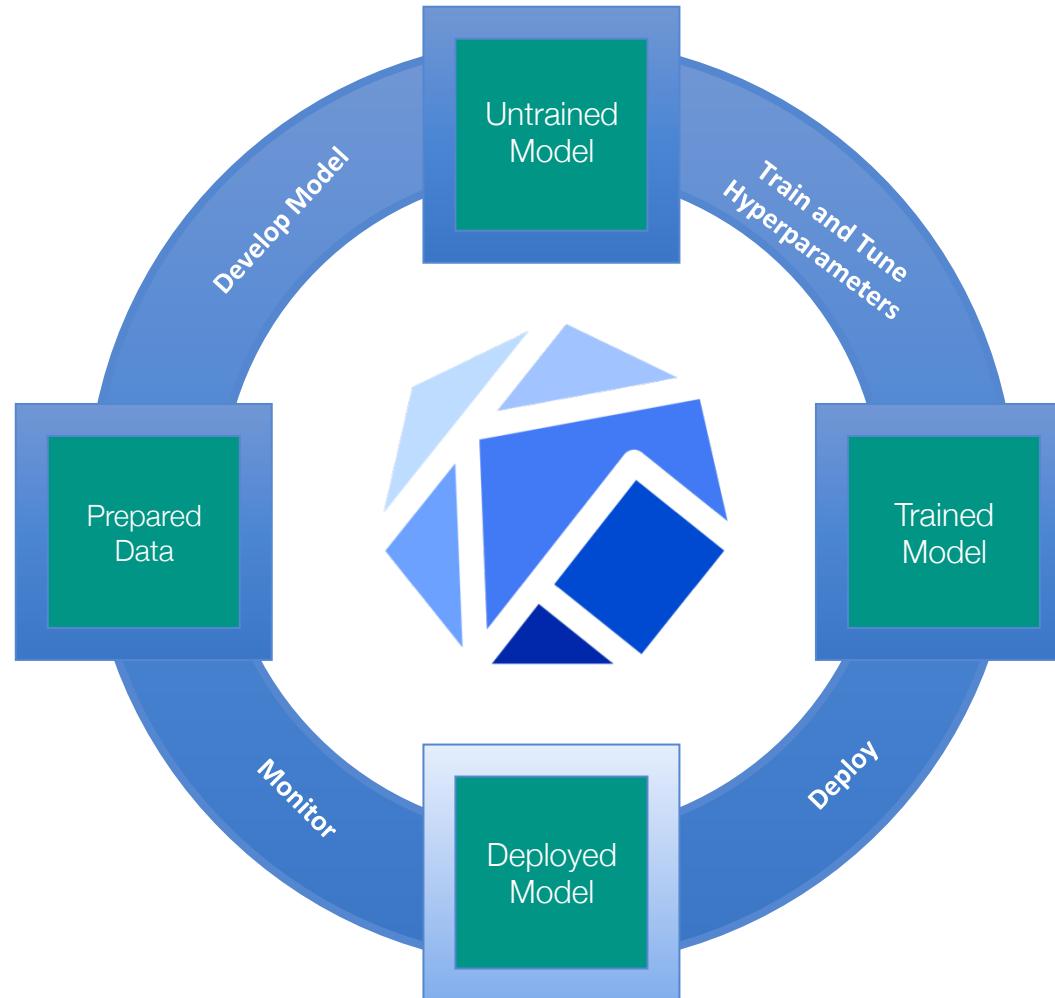
KATIB

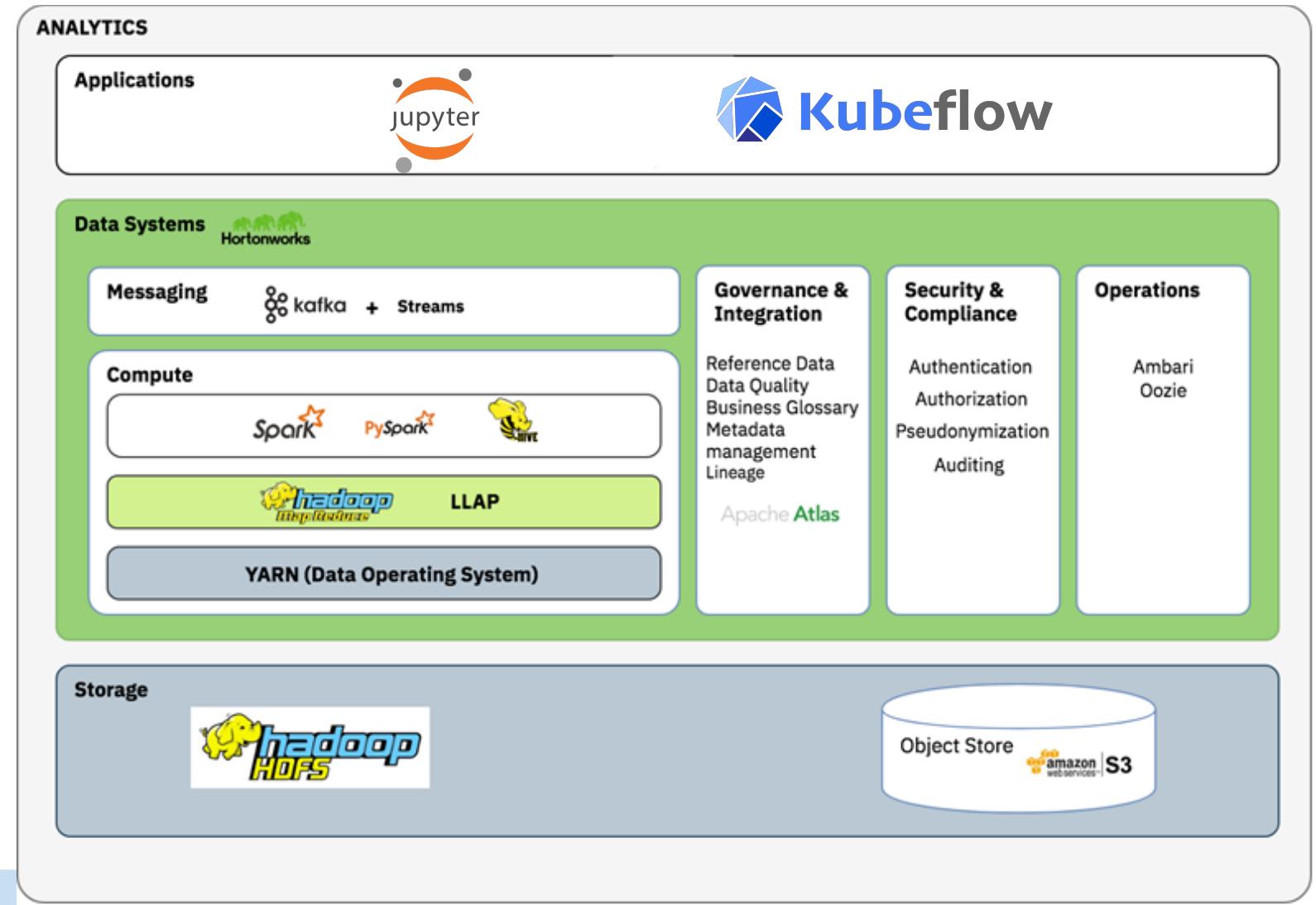
KFSERVING

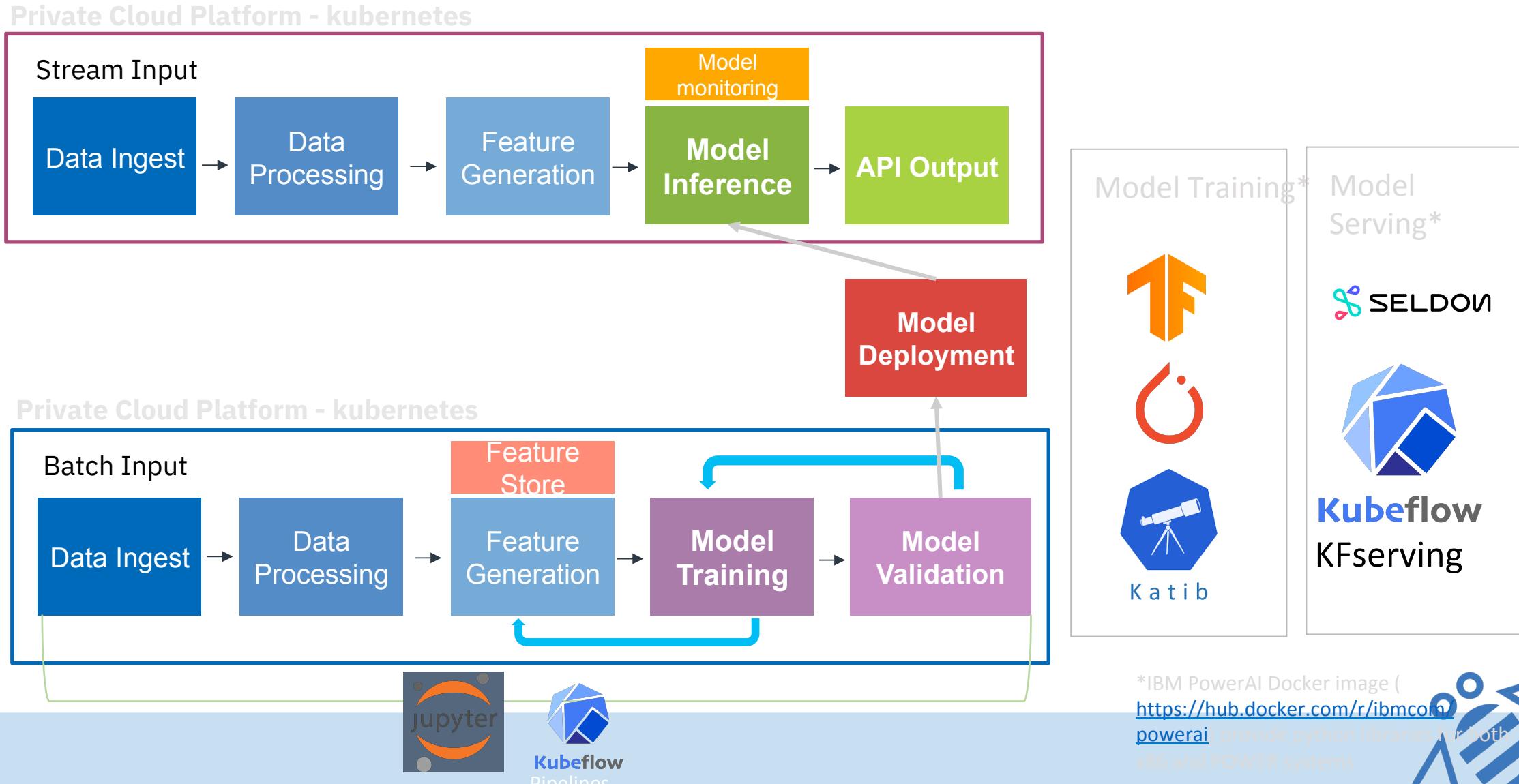
...

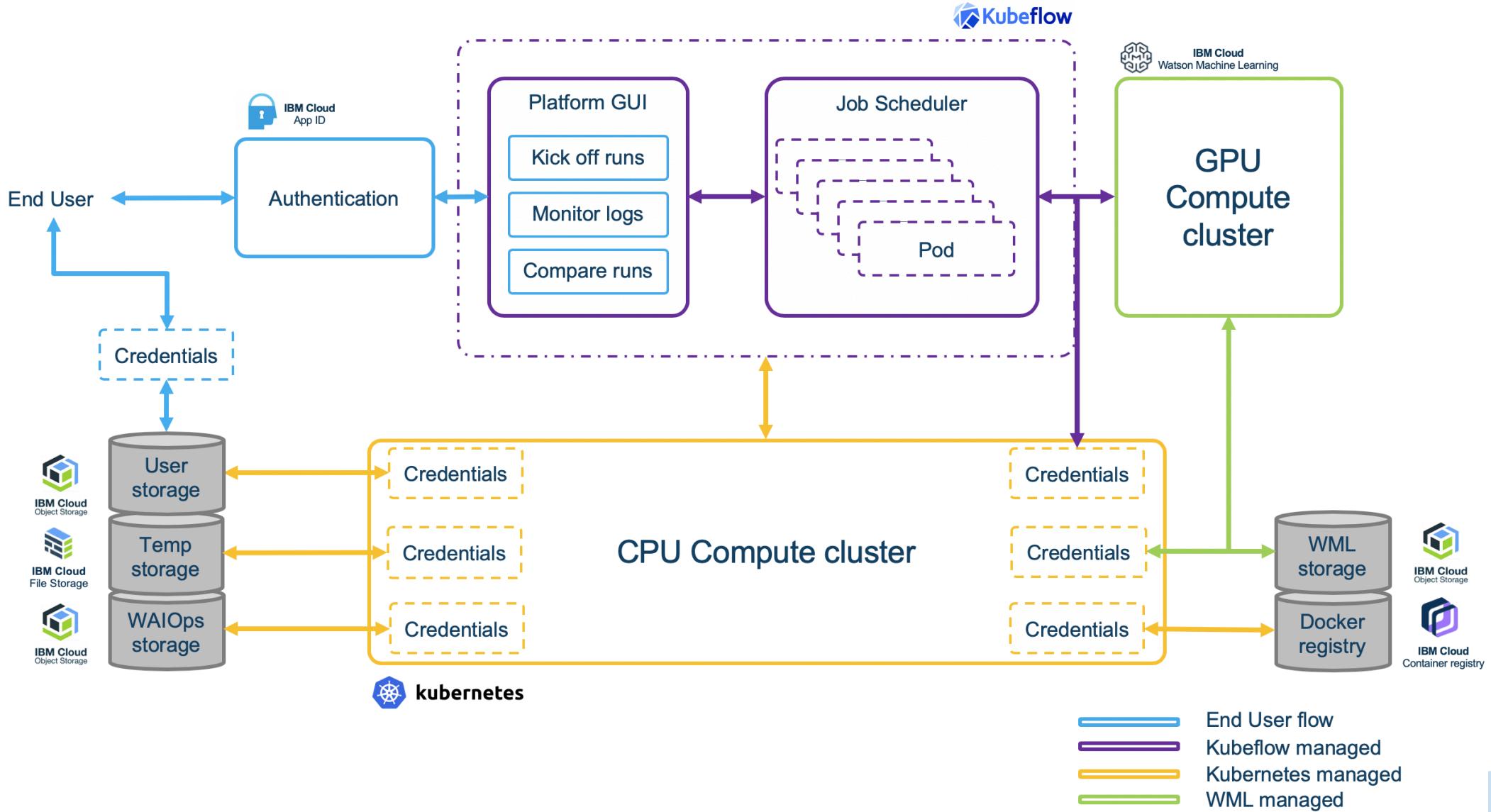


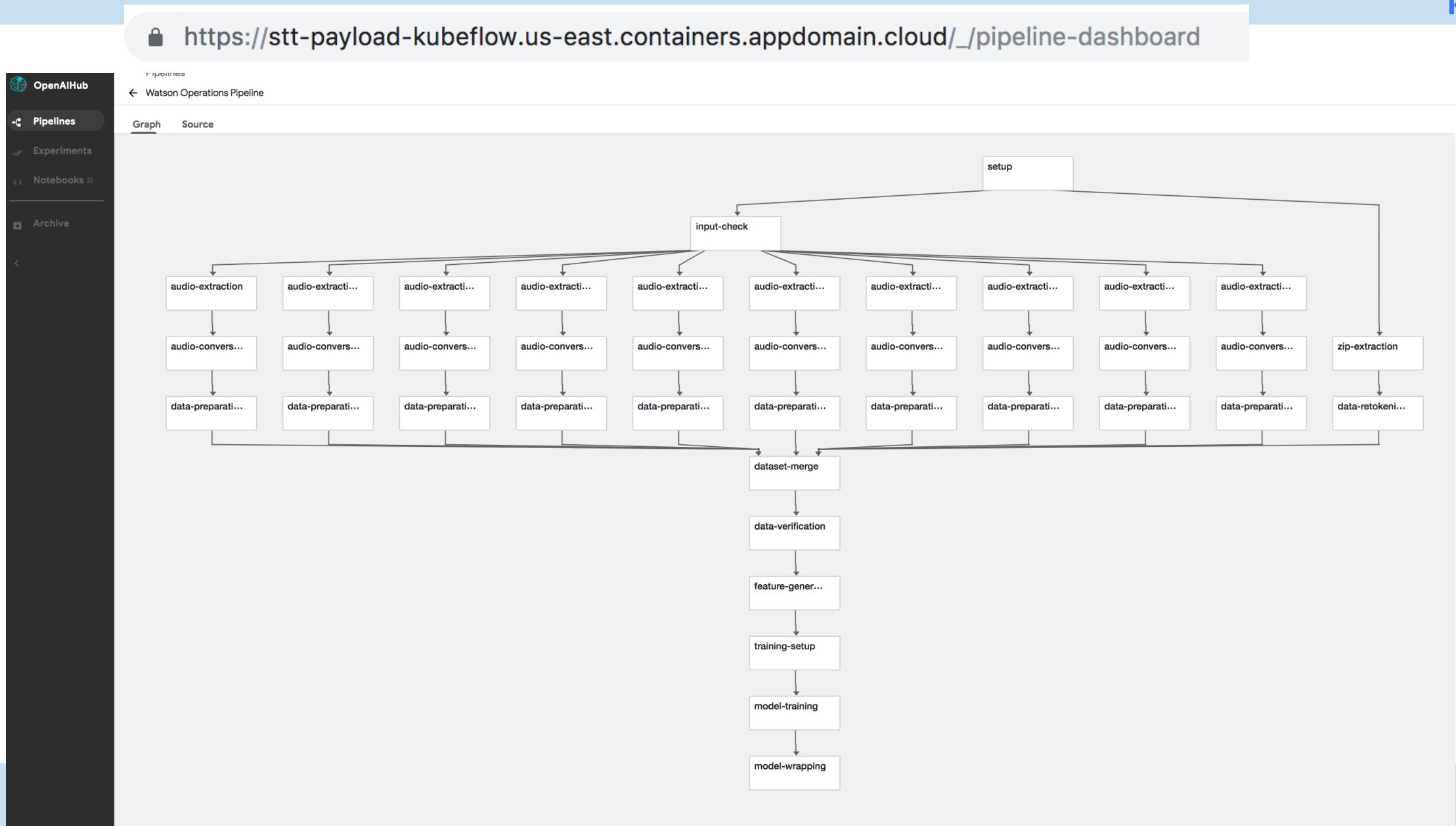
# Kubeflow Adoption: External and Internal

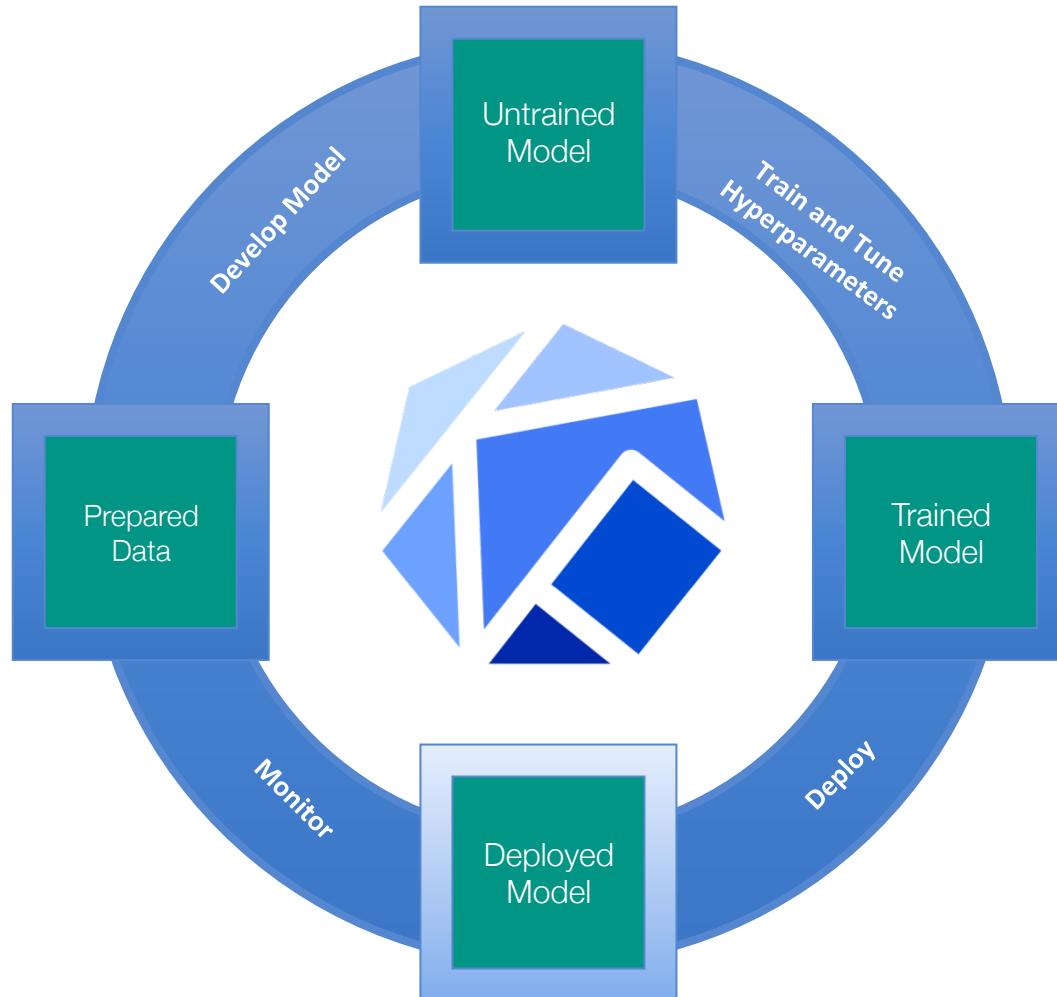


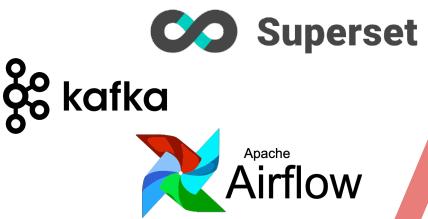












jupyterhub

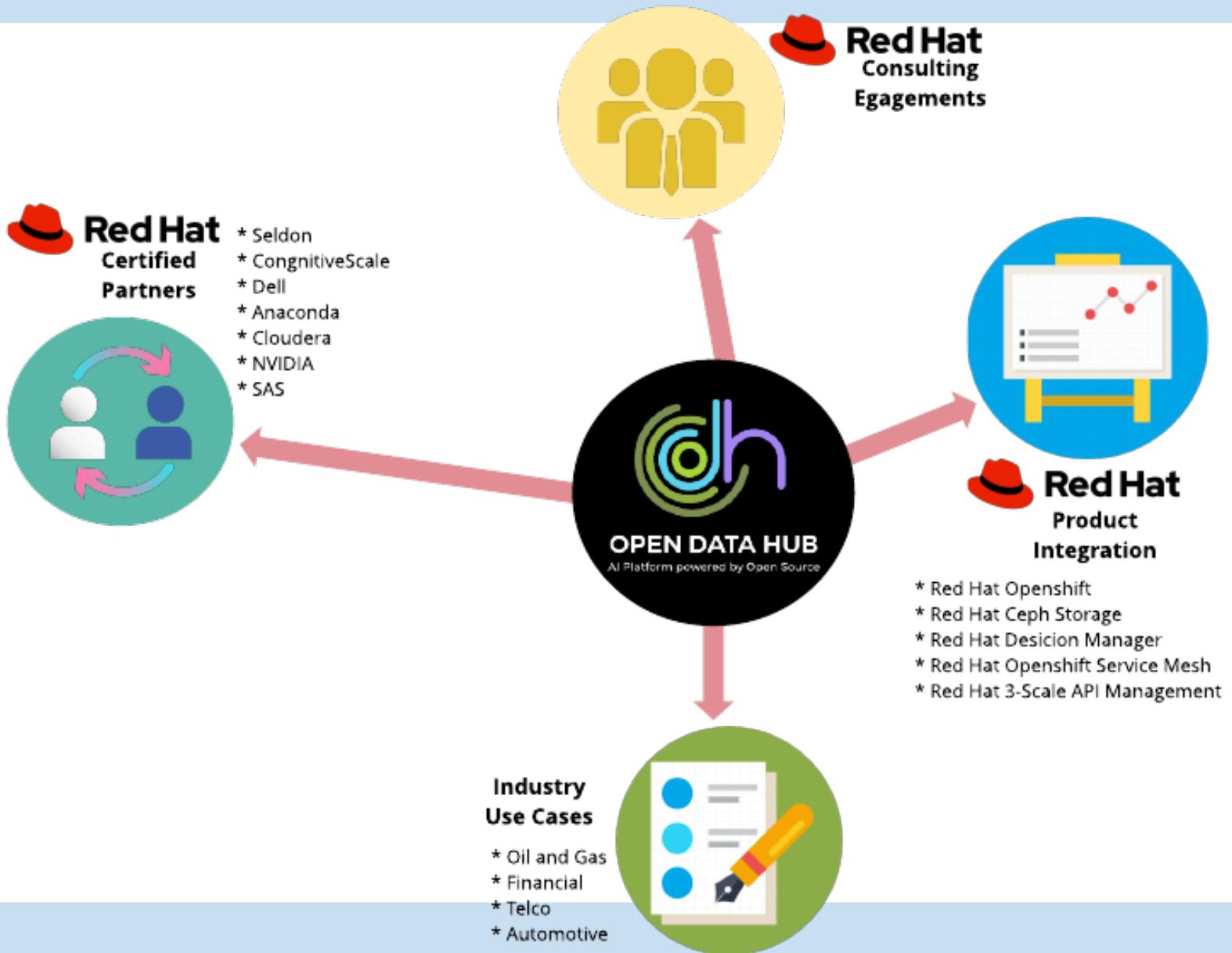


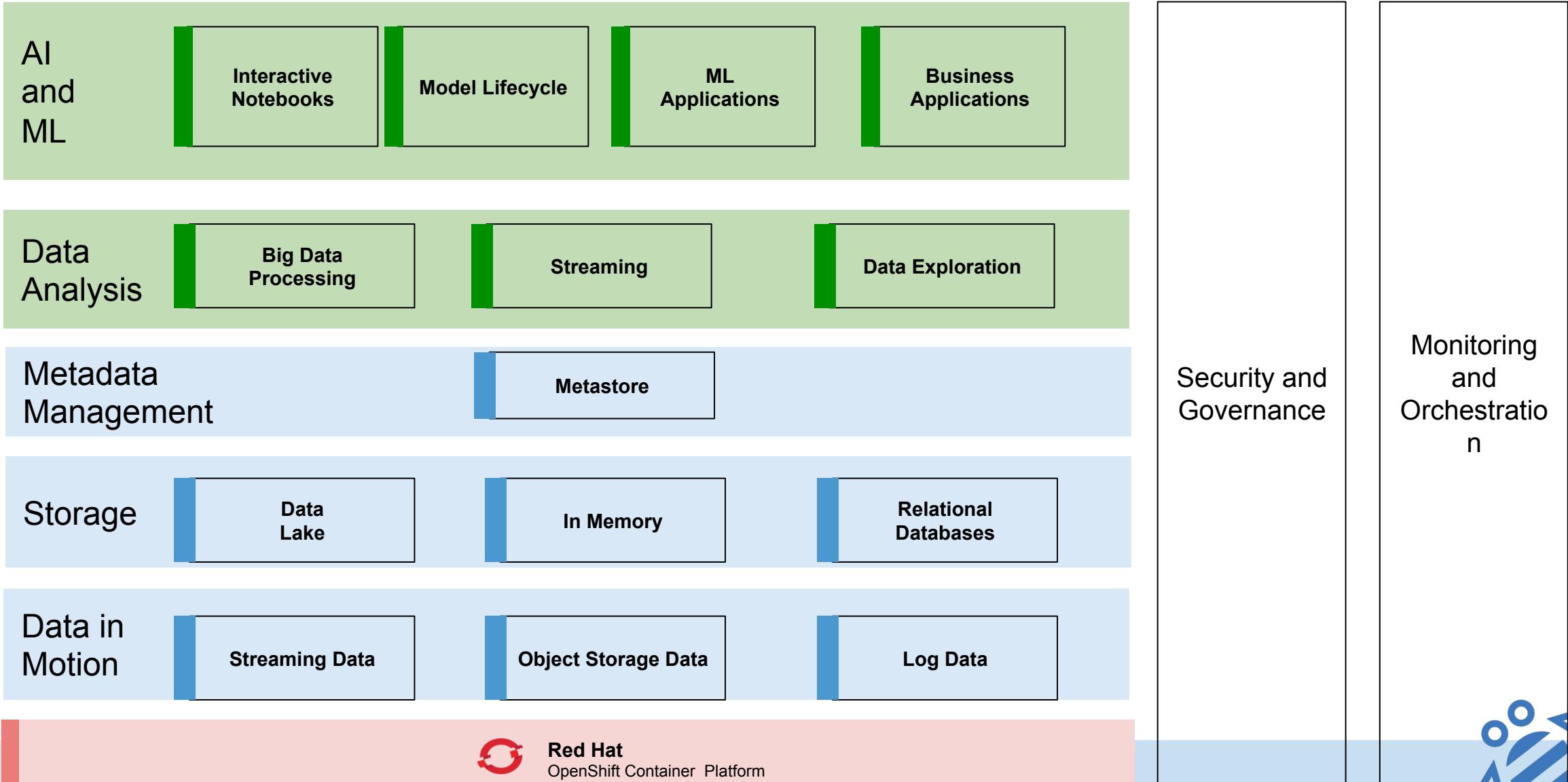
OpenShift  
Ready

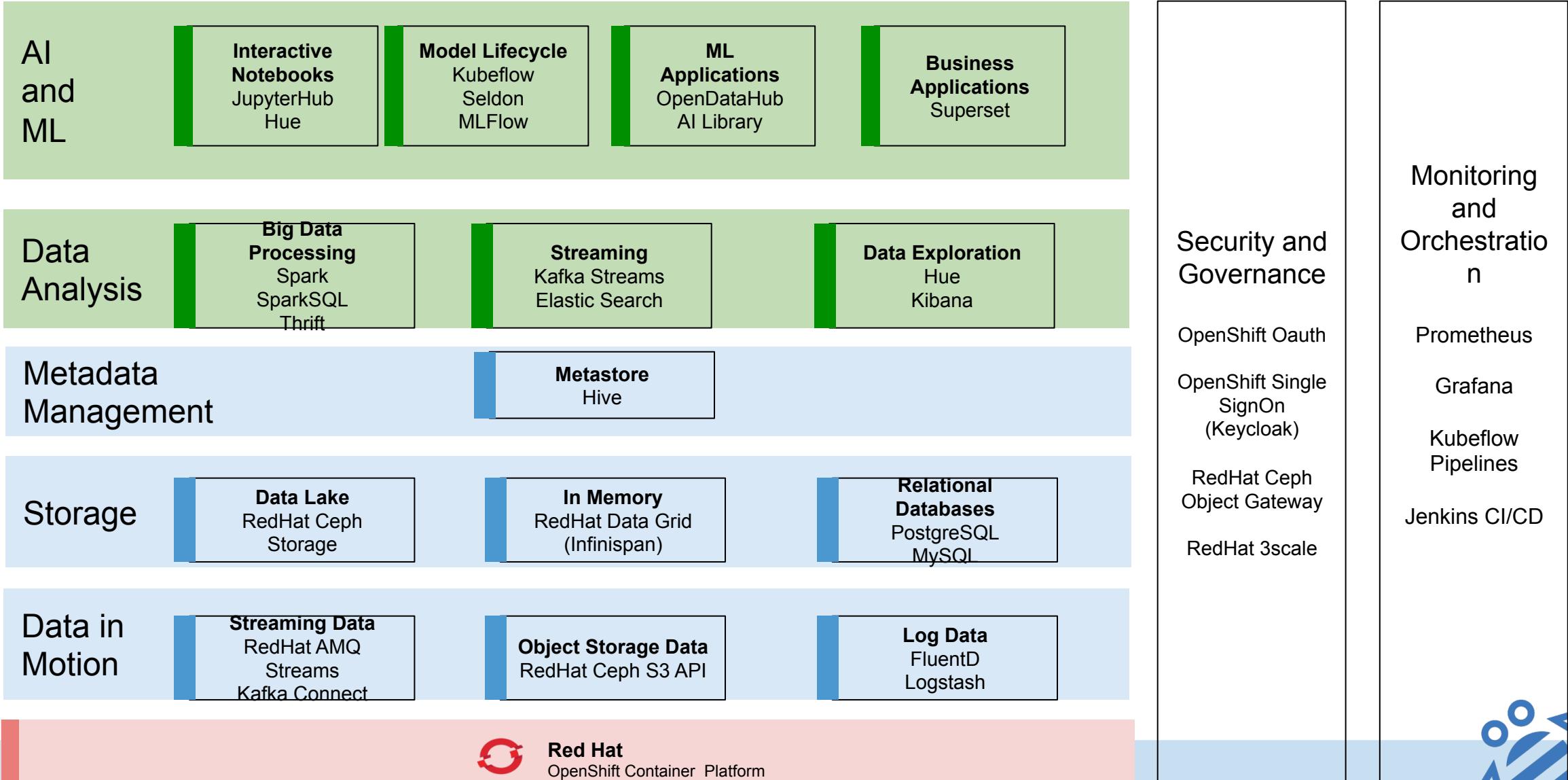


**Data Platform**

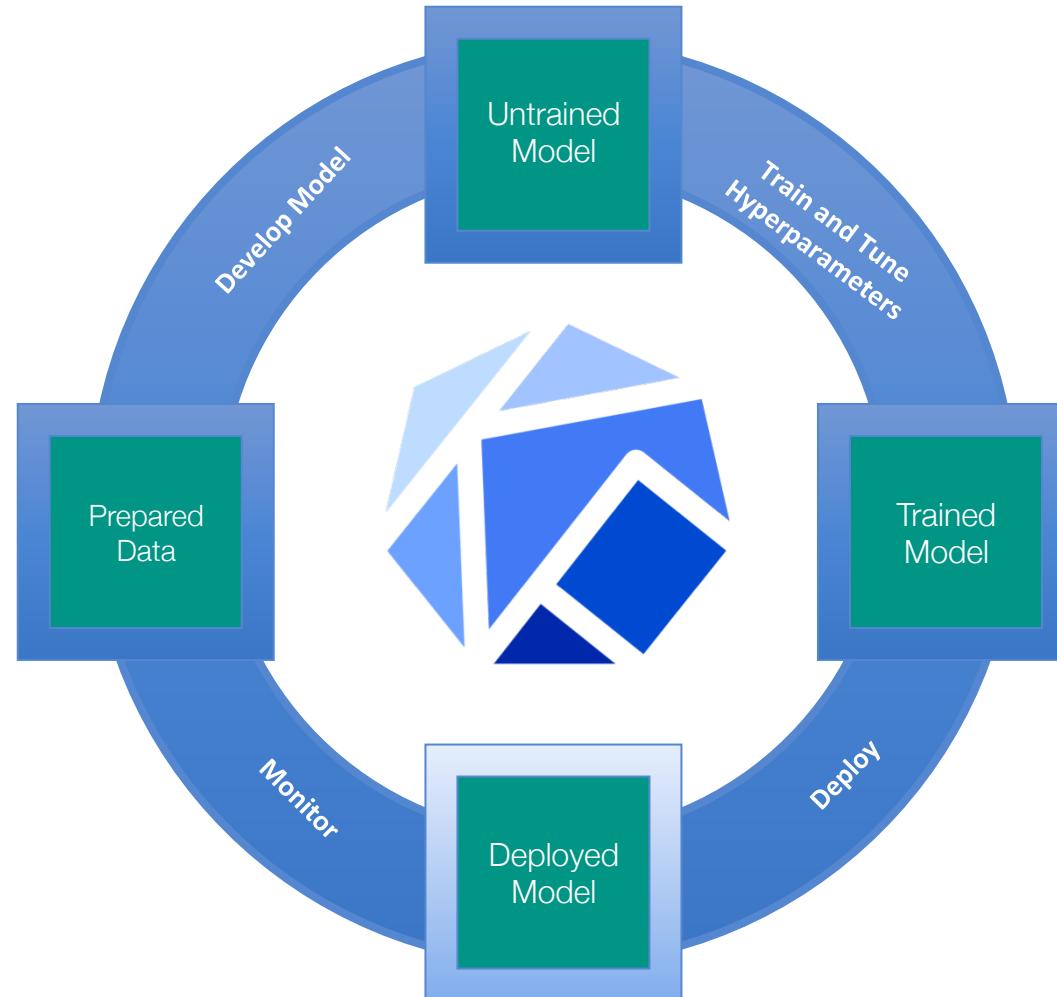
Operator Hub - [operatorhub.io](https://operatorhub.io)





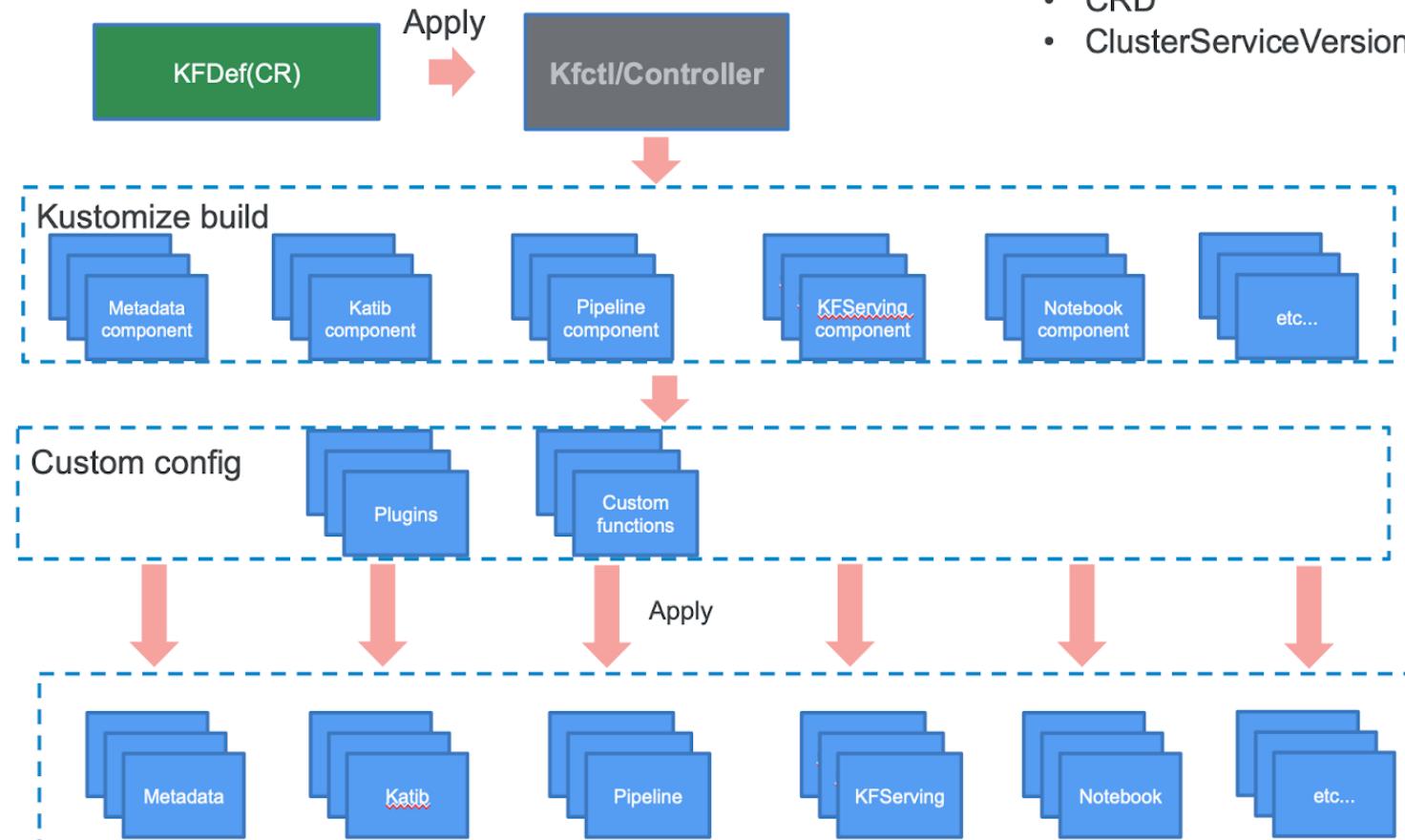


# OpenDataHub and Kubeflow: Relationship



- Deploy, manage and monitor Kubeflow
- On various environments
  - IBM Cloud
  - GCP
  - AWS
  - Azure
  - OpenShift
  - Other K8S

## KFCTL CONTROLLER - Initial deployment



- An version of the Operator based on Kubeflow Architecture released-  
[https://developers.redhat.com/blog/2020/05/07/open-data-hub-0-6-brings-component-updates-and-kubeflow-architecture/?sc\\_cid=7013a000002DTqEAAW](https://developers.redhat.com/blog/2020/05/07/open-data-hub-0-6-brings-component-updates-and-kubeflow-architecture/?sc_cid=7013a000002DTqEAAW)
- Most of the components converted  
<https://github.com/opendatahub-io/odh-manifests>
- Still a separate deployment – needs to do both ODH and Kubeflow in one go.

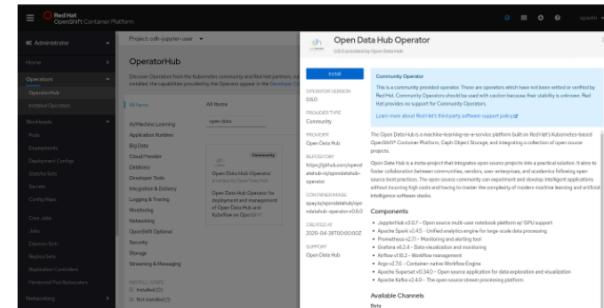
## Future

- KF 1.0 on OpenShift
- Disconnected deployment
- Open Data Hub CI/CD
- Kubeflow on OpenShift CI
- UBI based ODH & KF
- Multitenancy model
- Mixing KF & ODH



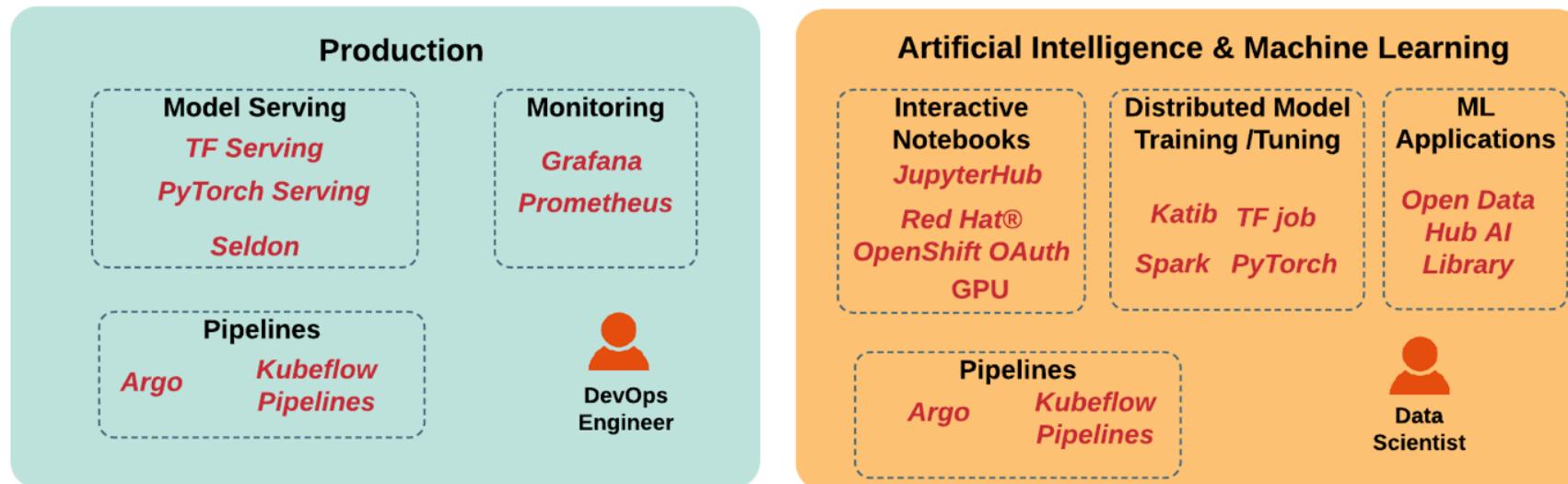
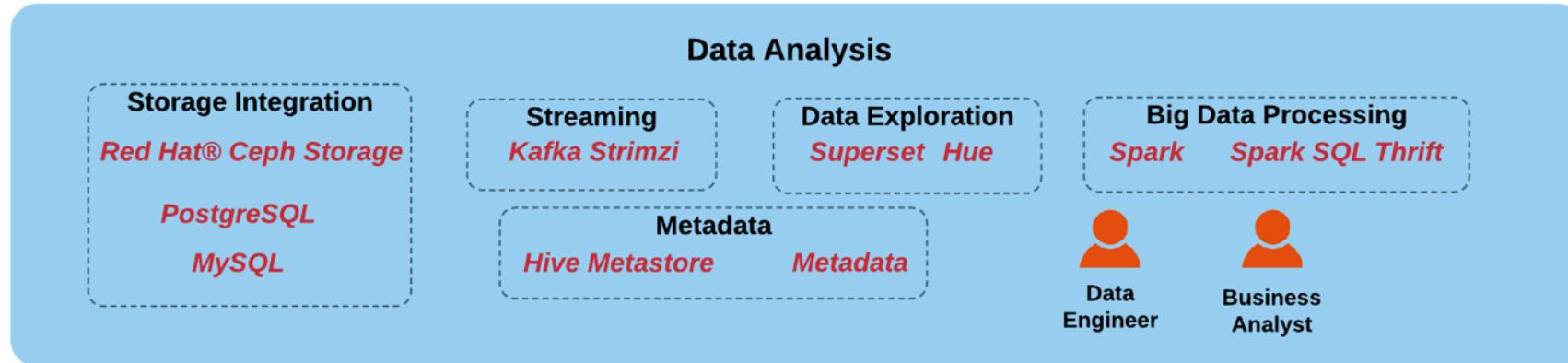
## Open Data Hub 0.6 brings component updates and Kubeflow architecture

By Václav Pavlin May 7, 2020



Open Data Hub (ODH) is a blueprint for building an AI-as-a-service platform on Red Hat's [Kubernetes](#)-based [OpenShift 4.x](#). Version 0.6 of Open Data Hub comes with significant changes to the overall architecture as well as component updates and additions. In this article, we explore these changes.





# Open Data Hub in OpenShift

Red Hat OpenShift Container Platform

Home Projects Status Search Events Catalog Workloads Networking Storage

You are logged in as a temporary administrative user. Update your password.

| NAME ↑                            | STATUS |
|-----------------------------------|--------|
| PR airflow-on-k8s-operator-system | Active |
| PR anonymous                      | Active |
| PR default                        | Active |
| PR kube-public                    | Active |
| PR kube-system                    | Active |
| PR opendatahub                    | Active |

Red Hat OpenShift Container Platform

Home Projects Status Search Events Catalog Workloads Networking Storage Builds Monitoring Compute Administration other resources

Project: opendatahub ▾

P jupyterhub-nb-kube-3aadmin

spark-operator

D spark-operator, #1

strimzi

D odh-message-bus-entity-operator, #1

SS odh-message-bus-kafka

SS odh-message-bus-zookeeper

superset

DC superset, #1

other resources

D ailibrary-operator, #1

D airflow-on-k8s-operator-controller-manager, #1

D argo-server, #1



# Apache Superset

**Superset** Security Manage Sources Charts Dashboards SQL Lab

Growth Analysis Scratchpad

Database: main Schema: superset Add a table (43)

**slices**

- created\_on
- changed\_on
- id**
- slice\_name
- datasource\_type
- datasource\_name
- viz\_type
- params
- created\_by\_fk**
- changed\_by\_fk**
- description
- cache\_timeout
- perm
- datasource\_id

**dashboards**

- created\_on
- changed\_on
- id**
- dashboard\_title
- position\_json
- created\_by\_fk**
- changed\_by\_fk**
- CSS

Database: main Schema: superset Add a table (43)

SELECT b.dashboard\_id, a.dashboard\_title, b.slice\_id, c  
**JOIN** dashboards a  
**JOIN** dashboard\_slices b ON a.id = b.dashboard\_id  
**JOIN** slices c ON c.id = b.slice\_id

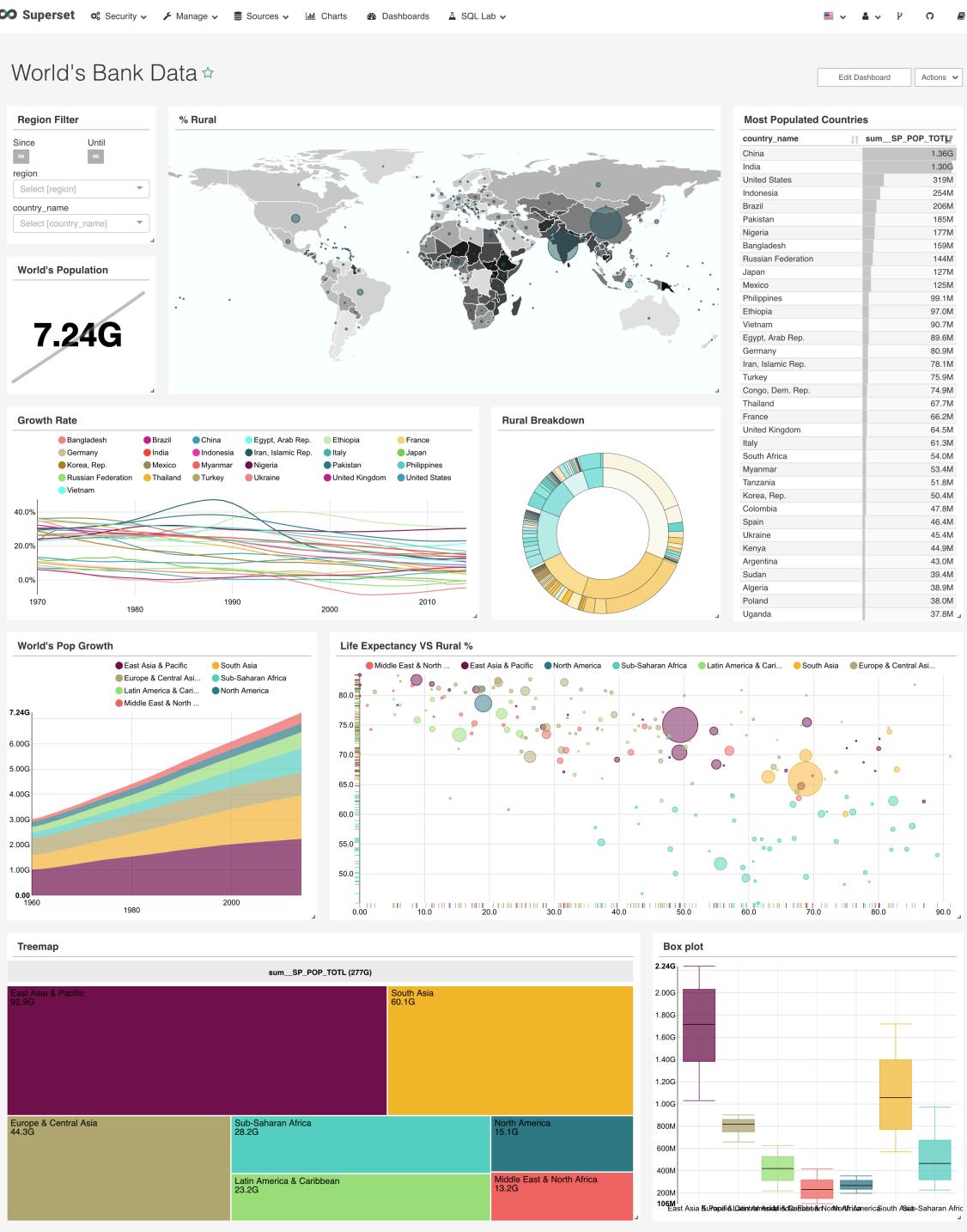
dashboards  
dashboard\_title  
DATABASE  
datasource\_type  
datasource\_name  
datasource\_id

Run Query Save Query Share Query

Results Query History Preview for slices Preview for d

Visualize .CSV

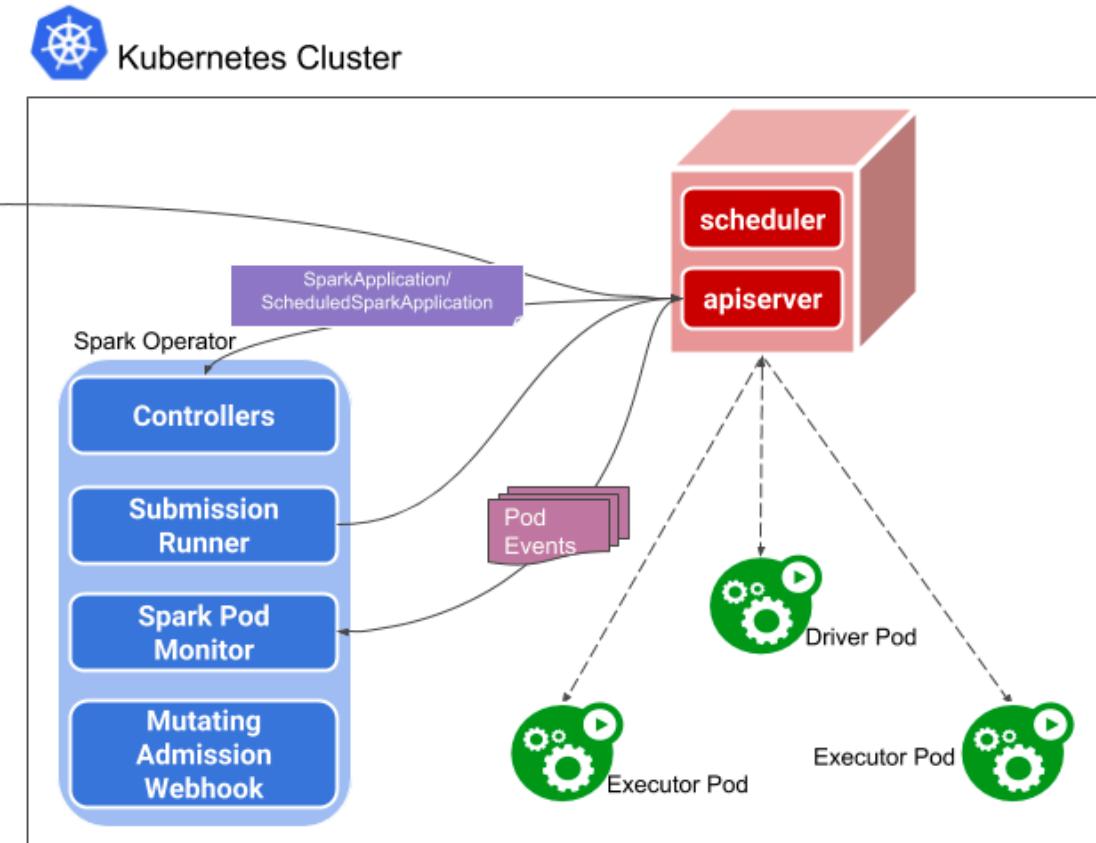
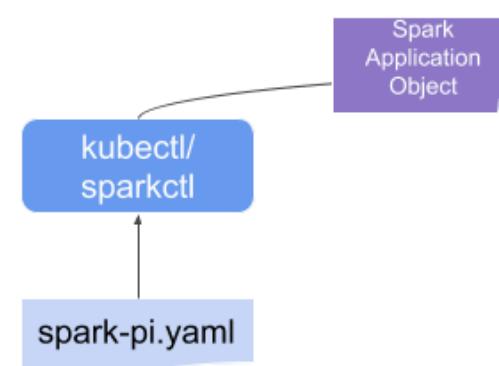
| dashboard_id | dashboard_title | slice_id | slice_name              |
|--------------|-----------------|----------|-------------------------|
| 2            | Births          | 882      | Girls                   |
| 2            | Births          | 883      | Boys                    |
| 2            | Births          | 884      | Participants            |
| 2            | Births          | 885      | Genders                 |
| 2            | Births          | 886      | Genders by State        |
| 2            | Births          | 887      | Trends                  |
| 2            | Births          | 888      | Average and Sum Treemap |
| 2            | Births          | 889      | Title                   |
| 2            | Births          | 890      | Name Cloud              |



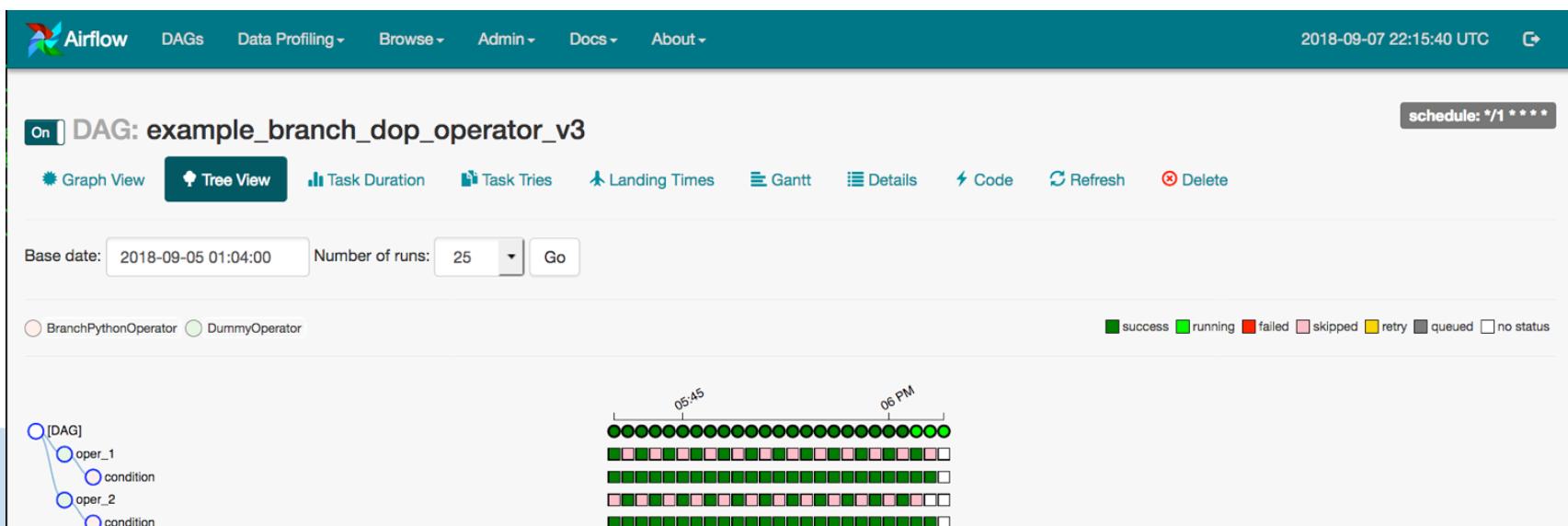
# IBM Spark with Open Data Hub



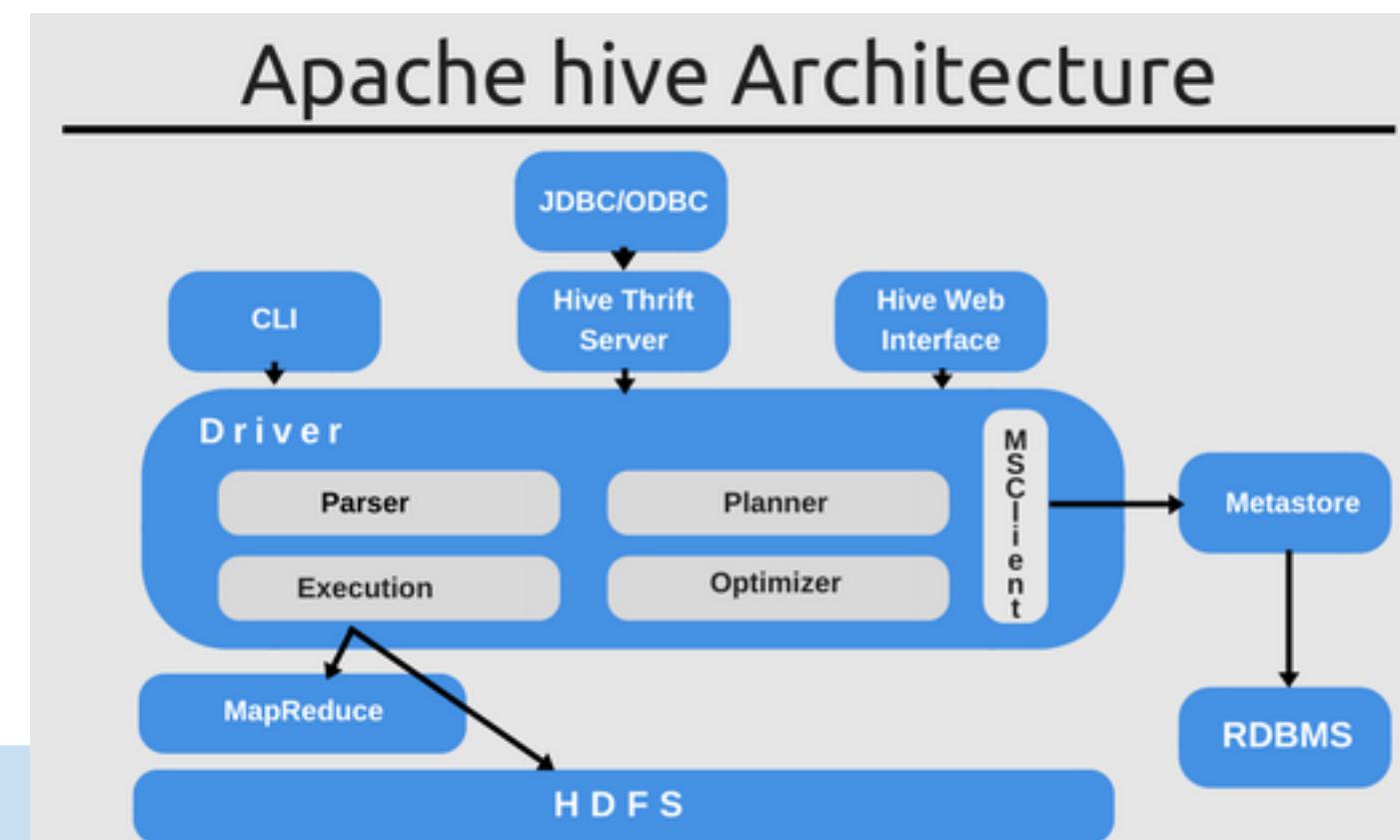
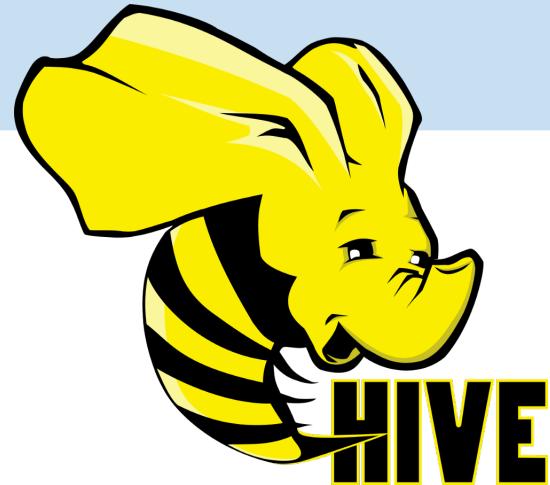
- Open Data Hub will also deploy the Spark Operator to manage Spark as an application.
- Two versions of Spark – Spark in dedicated mode and Spark on K8s
- Currently moving towards Spark on K8s Operator from Google for serverless Spark. IBM Hummingbird team investigating this



- Open Data Hub will also deploy the Airflow Operator to manage Airflow as an application.
- Using the Airflow Operator originally developed in the GoogleCloudPlatform repository and later donated to Apache.
- The Operator creates a controller-manager pod which will be created as a part of the Open Data Hub deployment.
- Users can then install the Airflow components they need from the available options (eg: CeleryExecutor or KubernetesExecutor, Postgres deployment or MySQL deployment etc. )



- Hive was one of the first abstraction engines to be built on top of MapReduce.
- Started at Facebook to enable data analysts to analyse data in Hadoop by using familiar SQL syntax without having to learn how to write MapReduce.
- Hive an essential tool in the Hadoop ecosystem that provides an SQL dialect for querying data stored in HDFS, other file systems that integrate with Hadoop such as MapR-FS and Amazon's S3 and databases like HBase(the Hadoop database) and Cassandra.
- Hive is a Hadoop based system for querying and analysing large volumes of structured data which is stored on HDFS.
- Hive is a query engine built to work on top of Hadoop that can compile queries into MapReduce jobs and run them on the cluster.



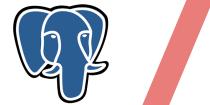


Superset



Apache  
Airflow

jupyterhub



PostgreSQL

OpenShift  
Ready



Data Platform

Operator Hub - [operatorhub.io](https://operatorhub.io)



PYTORCH



XGBoost



**seldon**

Spark

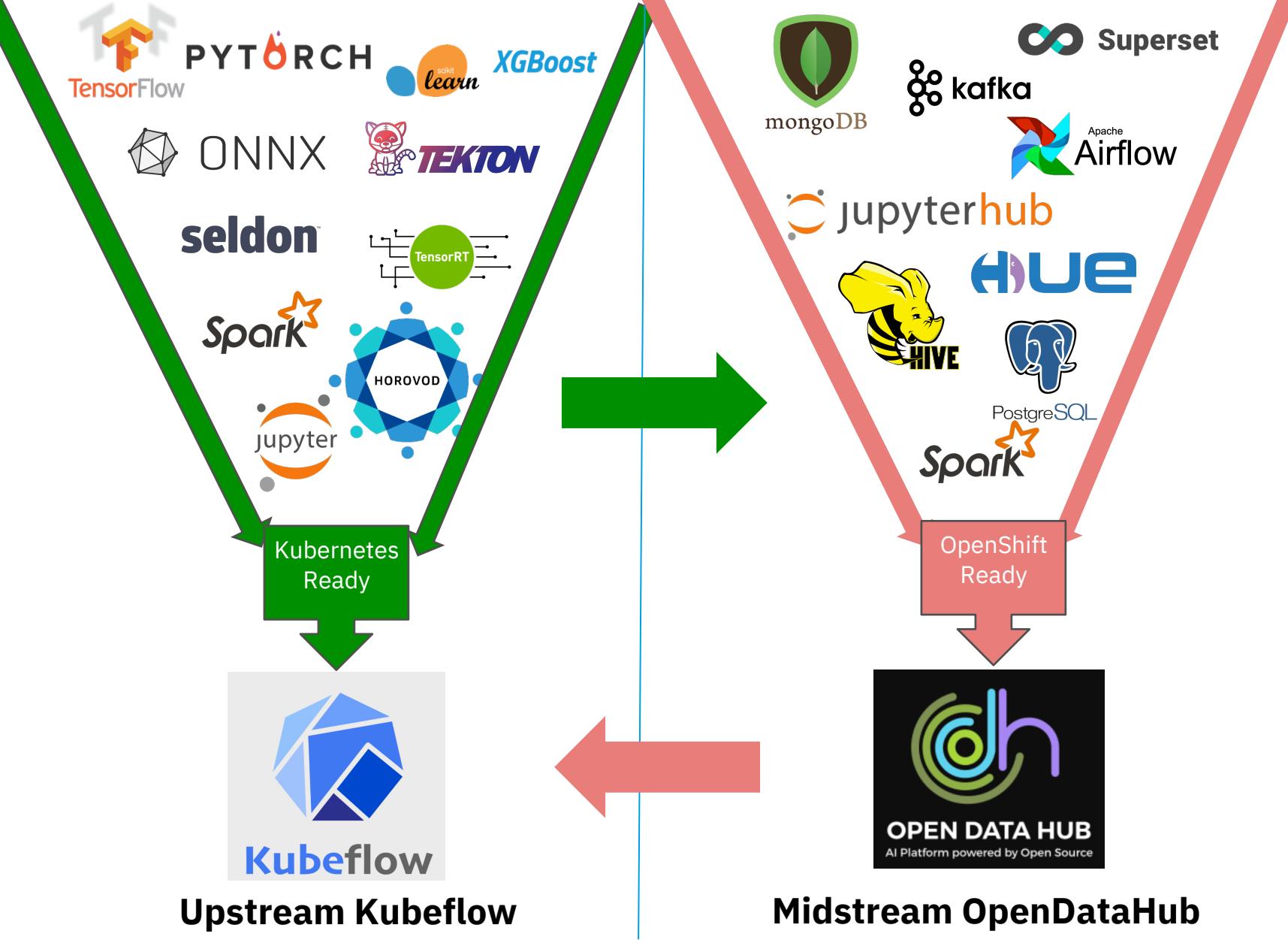
jupyter

Kubernetes  
Ready



**ML and AI Platform**

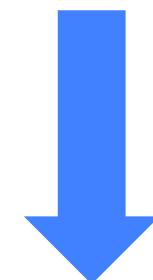
Operator Hub - [operatorhub.io](https://operatorhub.io)



# OpenDataHub



## Kubeflow



**Open Source End To End  
Data and AI Platform**

**Upstream Kubeflow**

**Midstream OpenDataHub**

Operator Hub - [operatorhub.io](https://operatorhub.io)

RedHat MarketPlace <https://marketplace.redhat.com/en-us>

Date: Wed July 15, 2020

| Time               | Topic                                           | Presenter    |
|--------------------|-------------------------------------------------|--------------|
| 8:00am - 8:30 am   | Data and AI Open Source at CODAIT               | Animesh      |
| 8:30am - 9:30 am   | Kubeflow Overview - End to end ML on Kubernetes | Animesh      |
| 9:30am - 9:45am    | Break                                           |              |
| 9:45am - 10:45am   | Git and Github                                  | Tom & Morgan |
| 10:45am - 11:00am  | Break                                           |              |
| 11:00am - 11:30am  | Kubeflow development environment                | Weiqiang     |
| 11:30am - 12:00 pm | Control plane deep dive                         | Weiqiang     |
| 12:00pm - 1:00pm   | Lunch break                                     |              |
| 1:00pm - 2:00pm    | Kubeflow deployment handson                     |              |
| 2:00pm - 3:00pm    | Tryout Kubeflow Components                      | Tommy        |
| 3:00pm - 4:00pm    | Q&A                                             |              |

[https://github.com/IBM/  
KubeflowDojo](https://github.com/IBM/KubeflowDojo)



<https://github.com/kubeflow>

<https://github.com/opendatahub-io>

Date: Thu July 16, 2020

| Time              | Topic                                  | Presenter                  |
|-------------------|----------------------------------------|----------------------------|
| 8:00am - 8:30am   | Overview of Kubeflow repos             | Tommy                      |
| 8:30 am - 9:30am  | Kubeflow Pipelines deep dive           | Animesh, Tommy, Christian  |
| 9:30am - 9:45am   | Break                                  |                            |
| 9:45 am - 10:45am | Kubeflow Pipelines-Tekton hands on     | Christian Kadner, Tommy Li |
| 10:45am - 11 am   | Break                                  |                            |
| 11:00am - 12 am   | KFServing deep dive                    | Animesh, Tommy             |
| 12:00pm - 1:00pm  | Lunch break                            |                            |
| 1:00pm - 2:00pm   | Distributed Training and HPO Deep Dive | Andrew, Kevin, Animesh     |
| 2:00pm - 2:15pm   | Break                                  |                            |
| 2:15pm - 2:30pm   | Kubeflow PR workflow                   | Weiqiang                   |
| 2:30pm - 3:30pm   | PR workflow handson                    |                            |
| 3:30pm - 4:00pm   | Wrap up and final Q&A                  | Animesh                    |