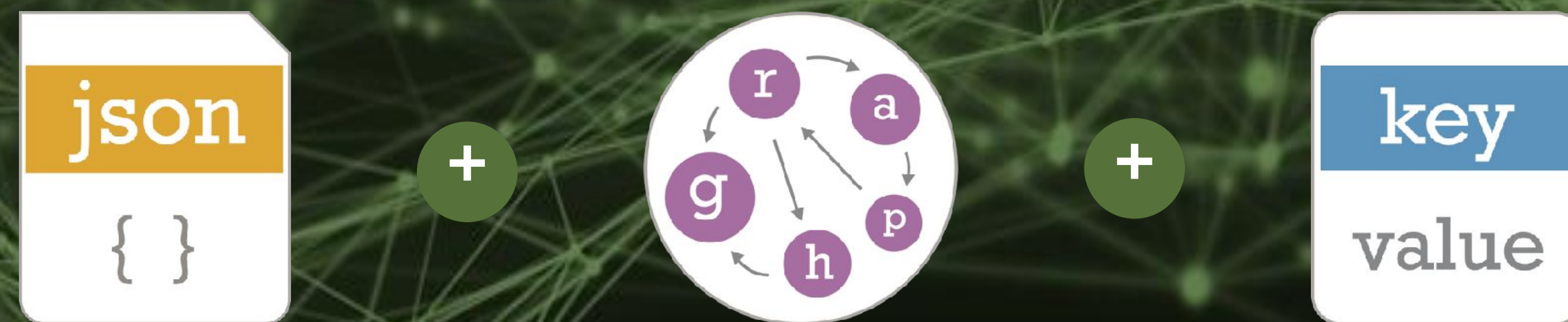# The case for a common metadata layer for machine learning platforms
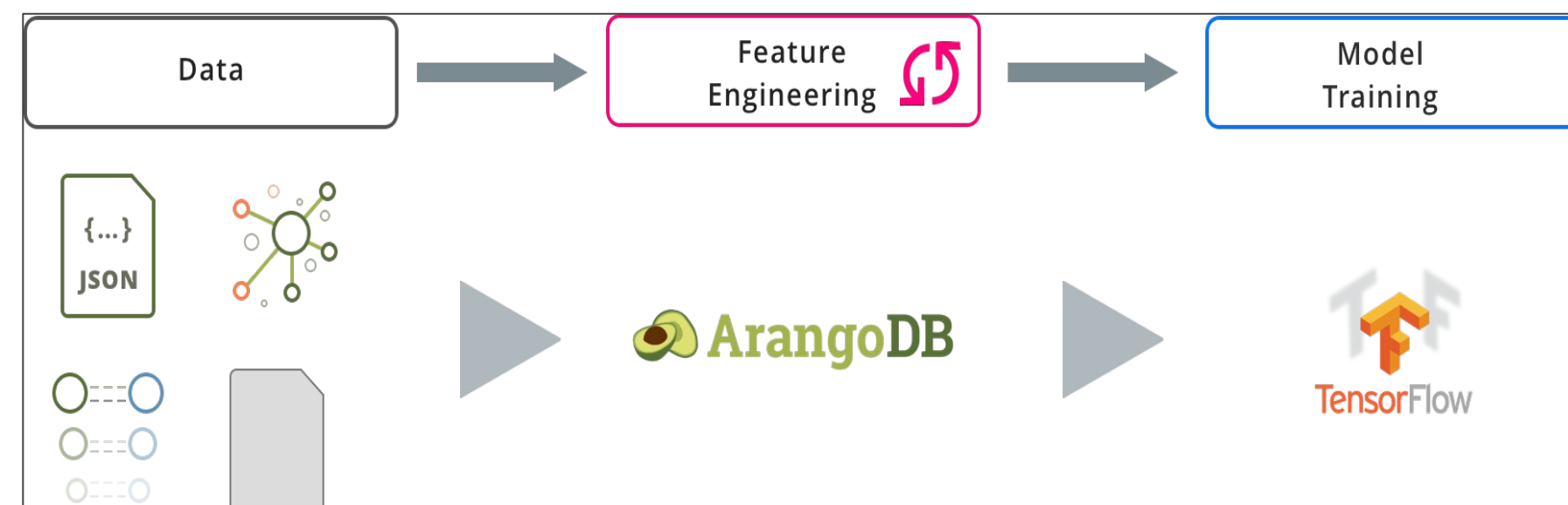## *From Data to Metadata*
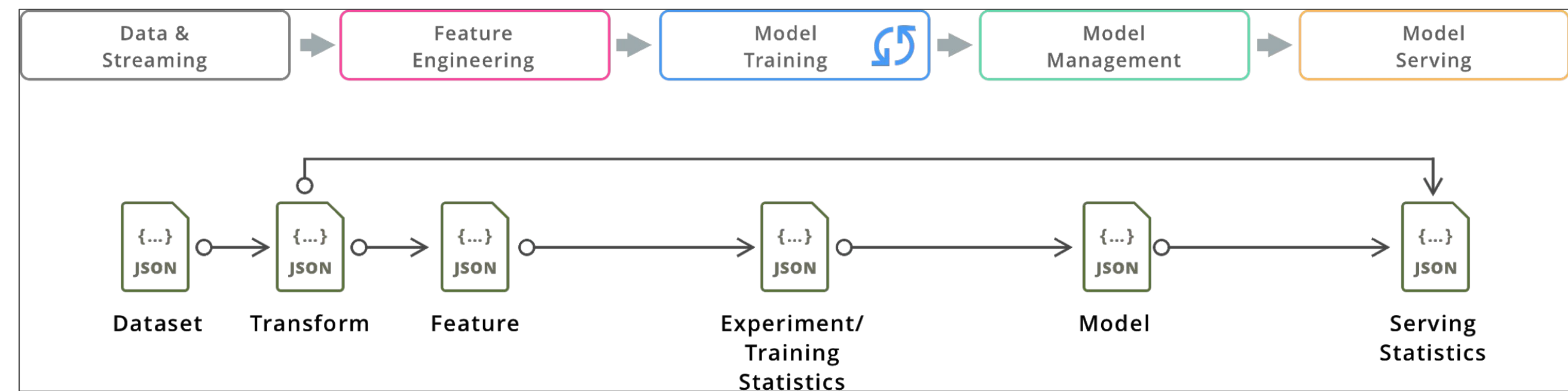
# Databases and Machine Learning

## Multi-Model-Powered Machine Learning

- Feature and Model Engineering

## Databases for Machine Learning Infrastructure

- Utilize Multi-Model for managing heterogeneous metadata across Machine Learning Pipelines

# ArangoML

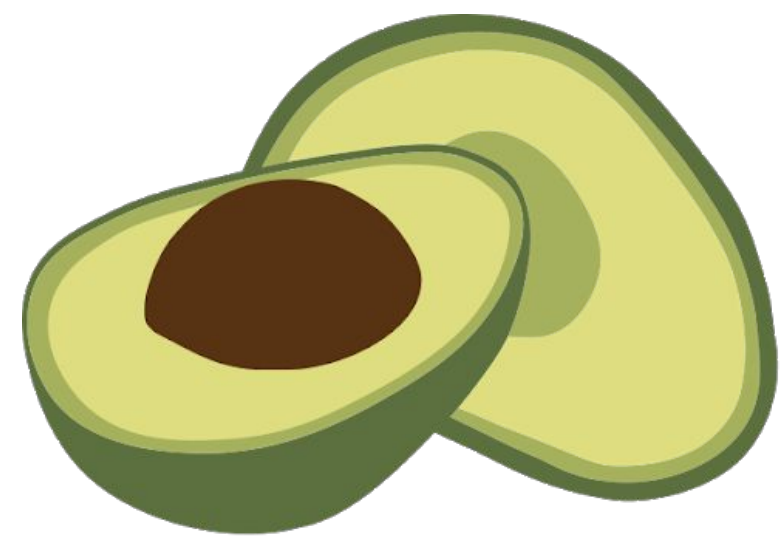## Multi-Model-Powered Machine Learning

- Feature and Model Engineering

## Databases for Machine Learning Infrastructure

- Utilize Multi-Model for managing heterogeneous metadata across Machine Learning Pipelines
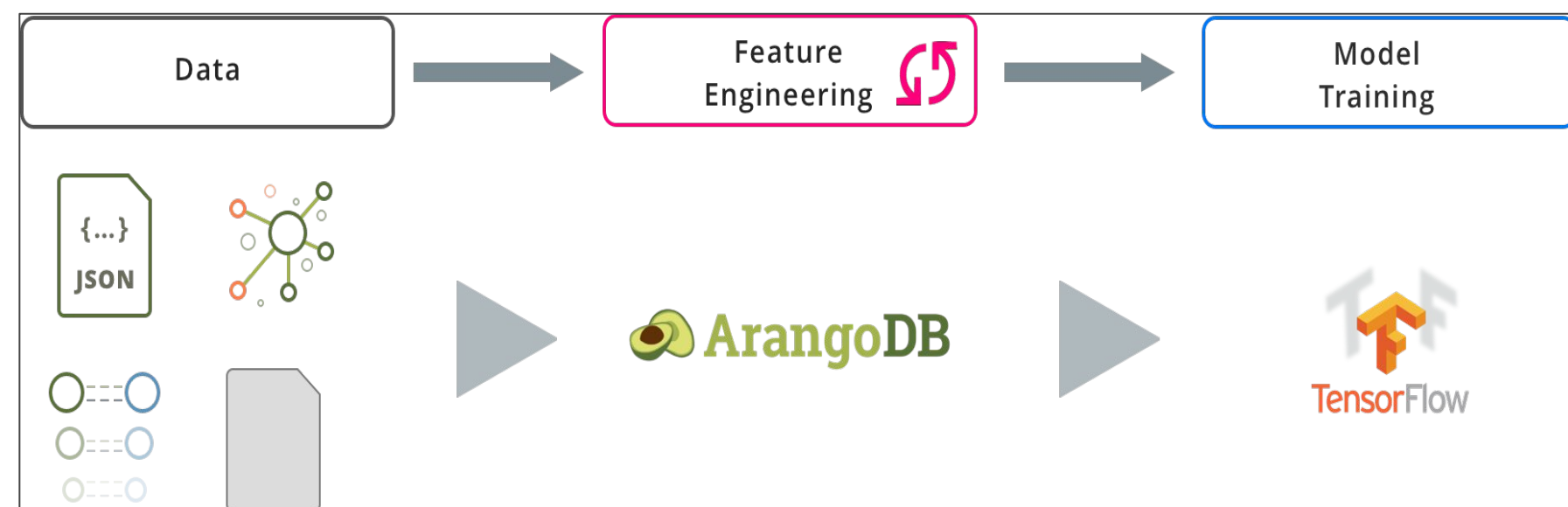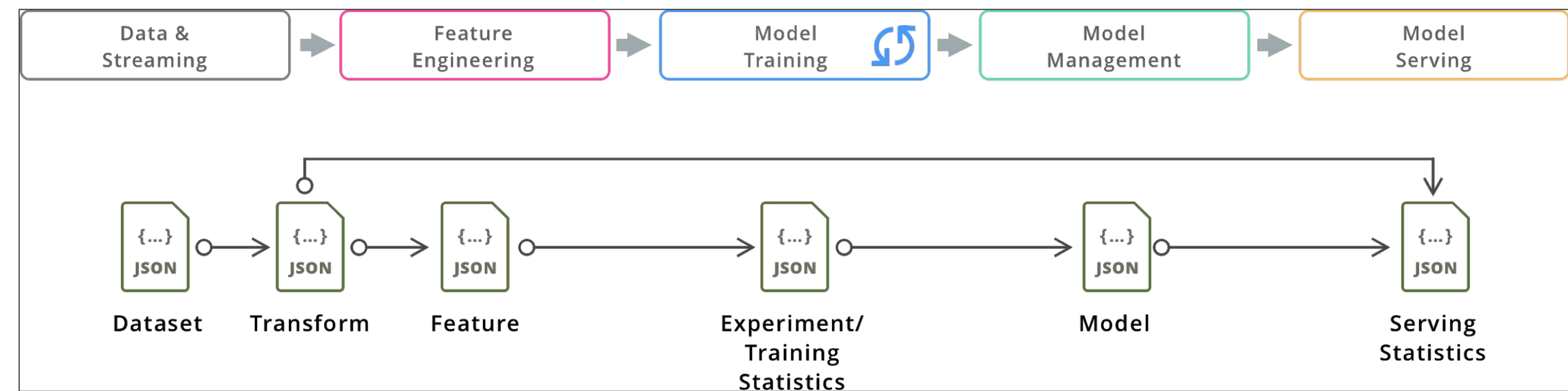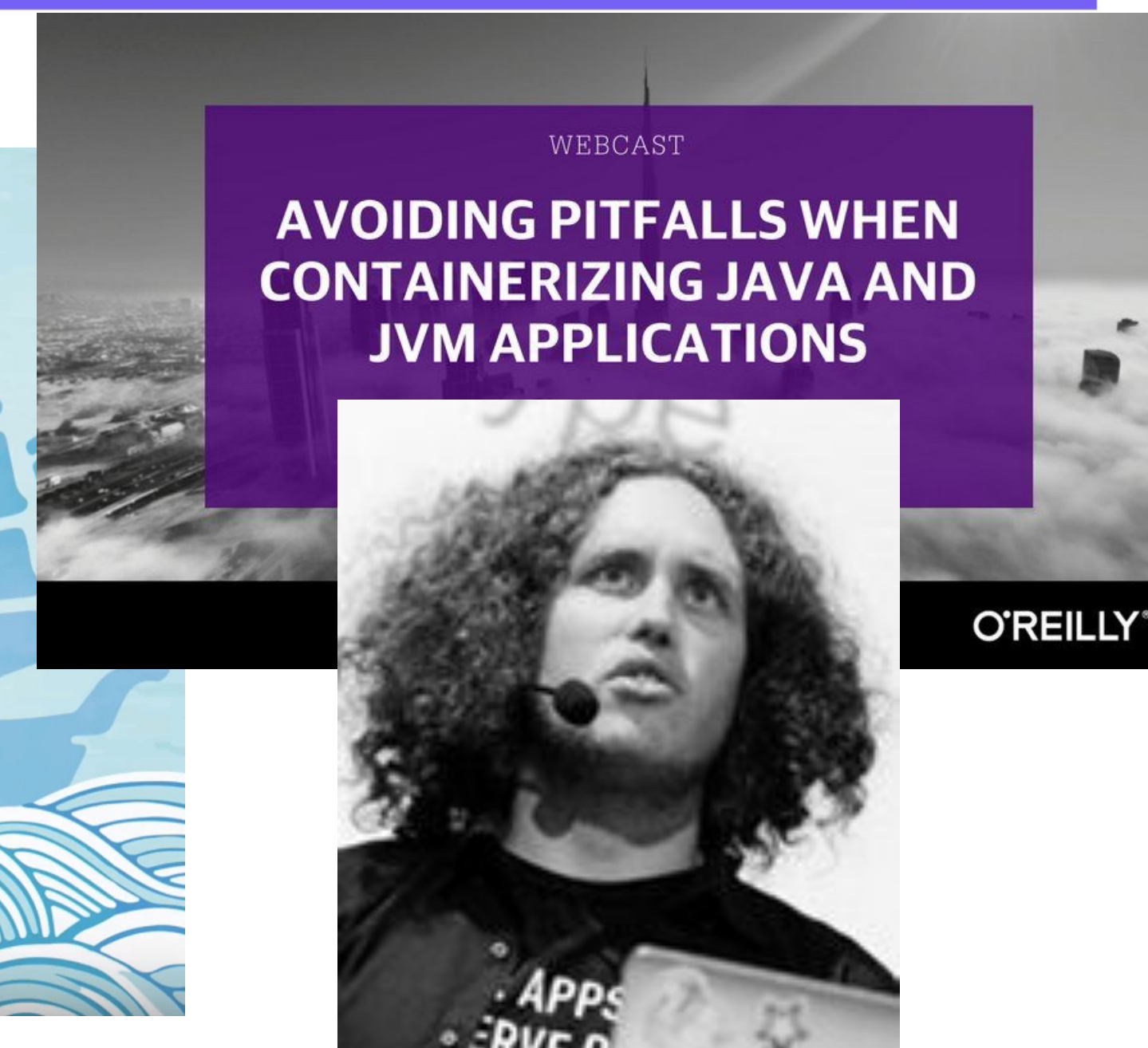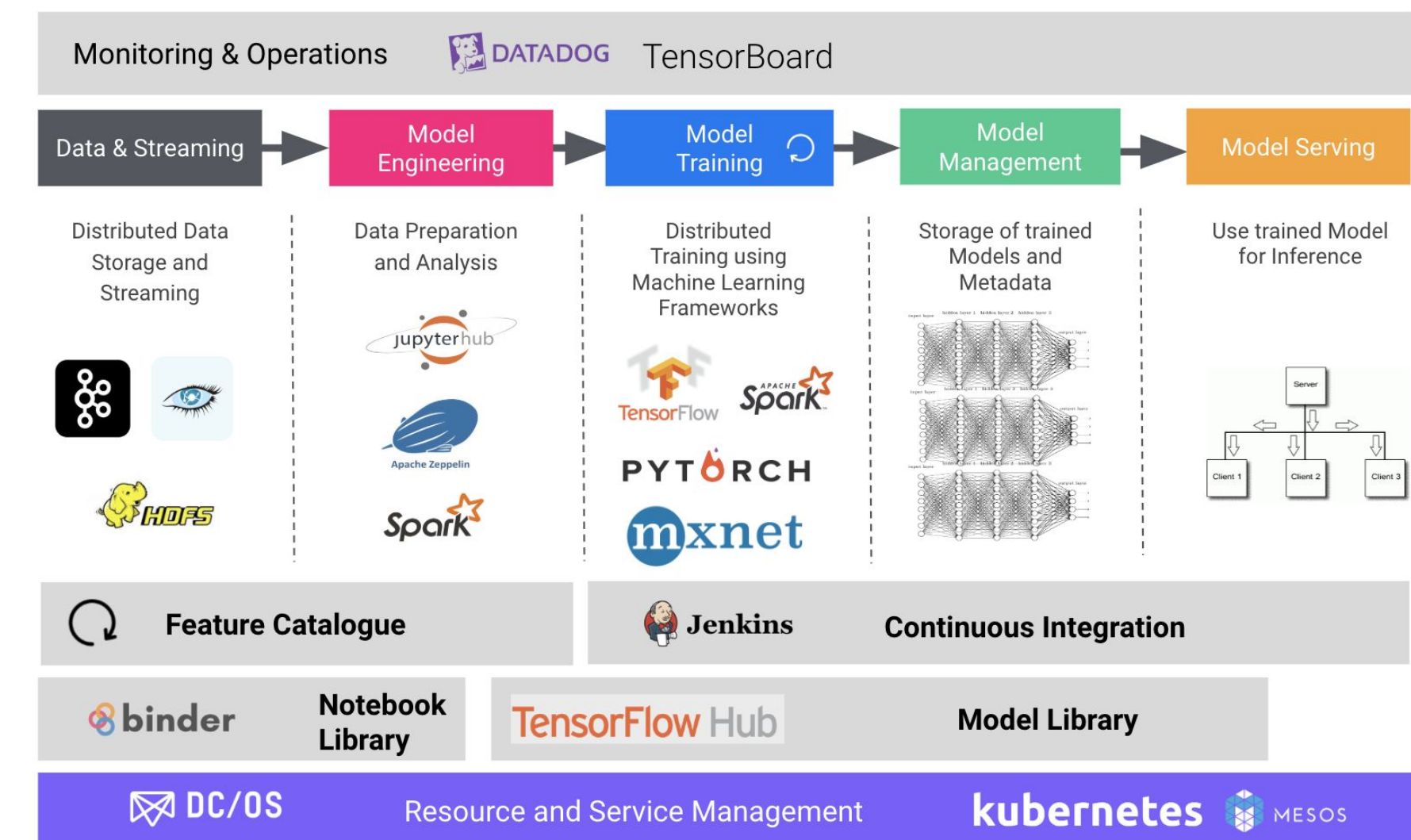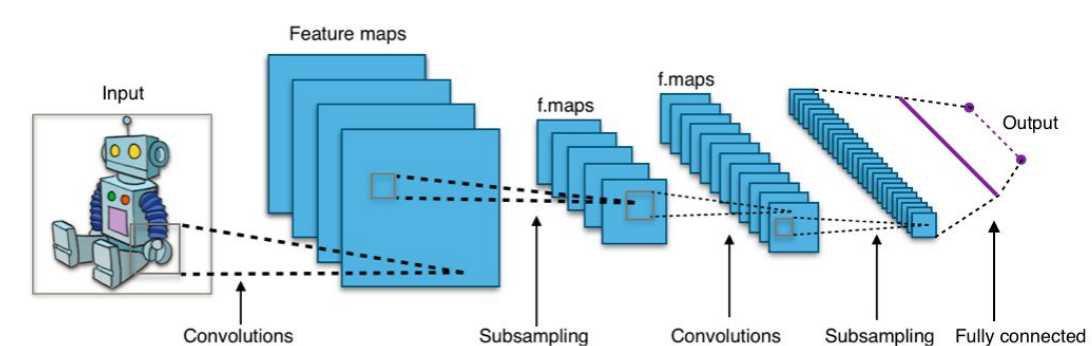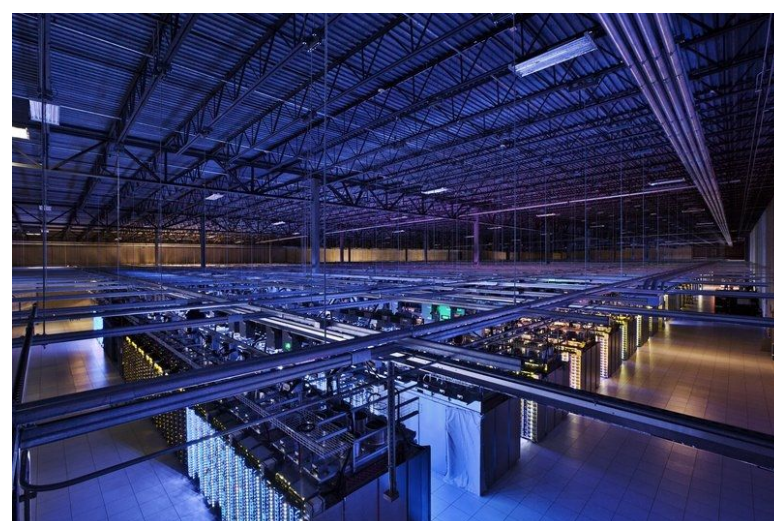
# Jörg Schad, PhD

**Head of Engineering and ML @ArangoDB**

- **Suki.ai**
- **Mesosphere**
- **Architect @SAP Hana**
- **PhD Distributed DB Systems**

- **Twitter: @joerg_schad**





Operating Deep Learning Pipelines Anywhere Using Kubeflow
Jörg Schad & Gilbert Song, Mesosphere



WEBCAST
AVOIDING PITFALLS WHEN CONTAINERIZING JAVA AND JVM APPLICATIONS

# Why is machine learning taking off?

Input    Feature maps    f.maps    f.maps    Output

Convolutions    Subsampling    Convolutions    Subsampling    Fully connected

Bach OR computer?

# DEEPBACH: A STEERABLE MODEL FOR BACH CHORALES GENERATION

# What Data Scientist should be doing…



Get Data → **Write intelligent machine learning code** → Train Model → Run Model

**Repeat**

# What Data Scientist are doing…



*Sculley, D., Holt, G., Golovin, D. et al. Hidden Technical Debt in Machine Learning Systems*

# Challenge: Persona(s)

# The Rise of the *DataOps Engineer*

Combines two key skills:

- Data science
- Distributed systems engineering

The equivalent of *DevOps* for *Data Science*

- **Build** automation software to run machine learning systems
- **Operate** systems so they're available, scalable, and performant
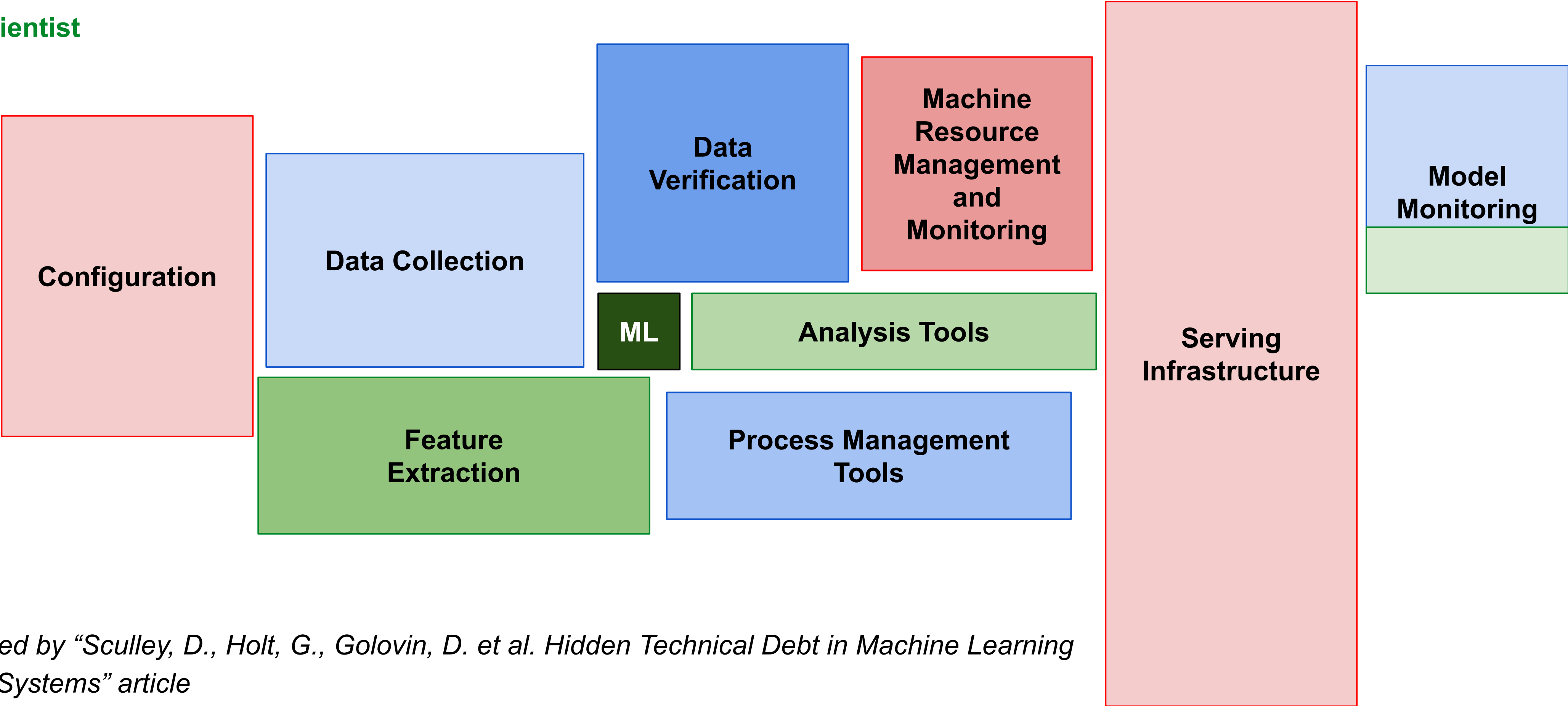- **Evangelize** tools and best practices among data scientists
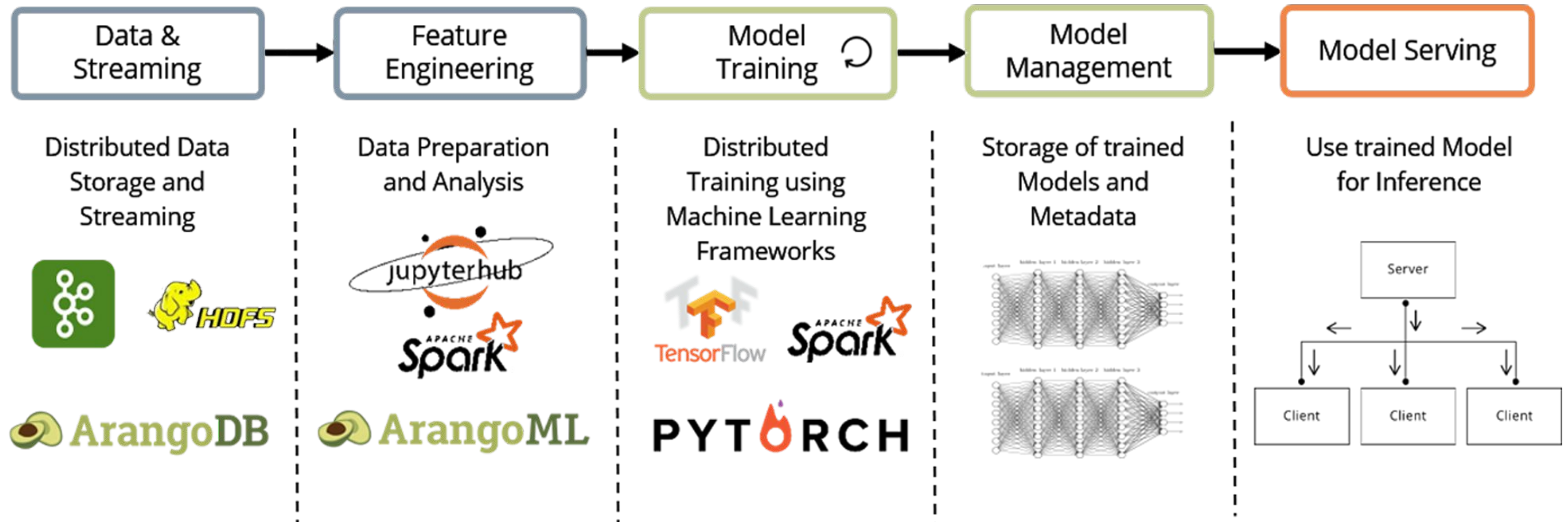
# Division of Labor

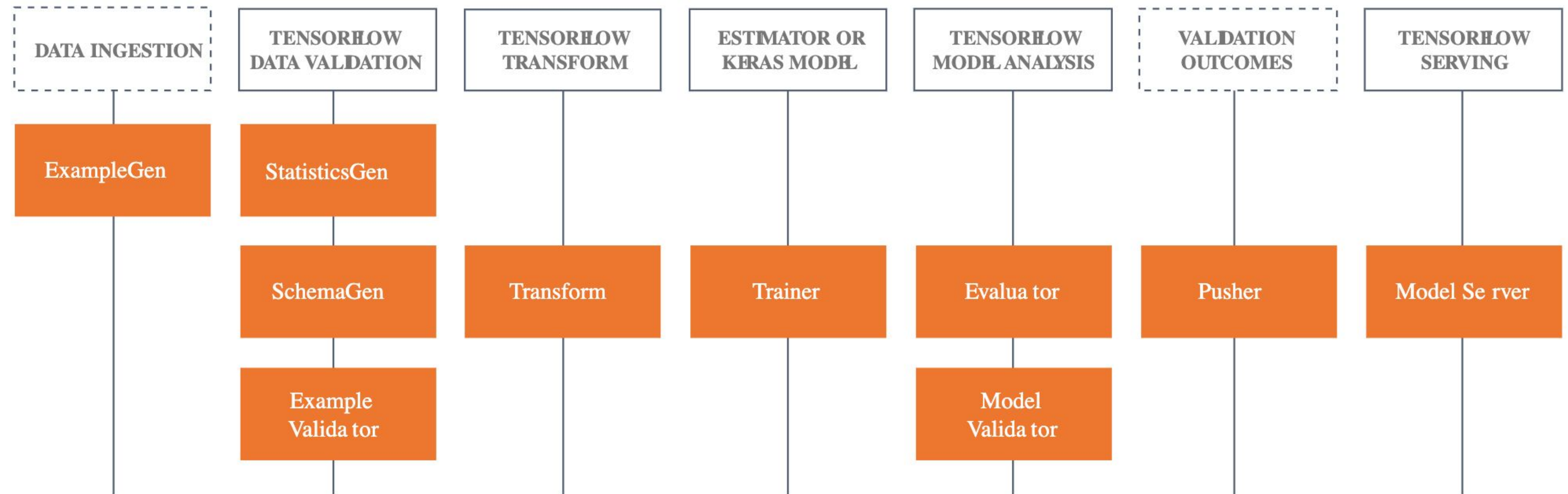System Admin/ DevOps

Data Engineer/DataOps

Data Scientist

Configuration

Data Collection

Data Verification

Machine Resource Management and Monitoring

Model Monitoring

ML

Analysis Tools

Serving Infrastructure

Feature Extraction

Process Management Tools

*Inspired by "Sculley, D., Holt, G., Golovin, D. et al. Hidden Technical Debt in Machine Learning Systems" article*
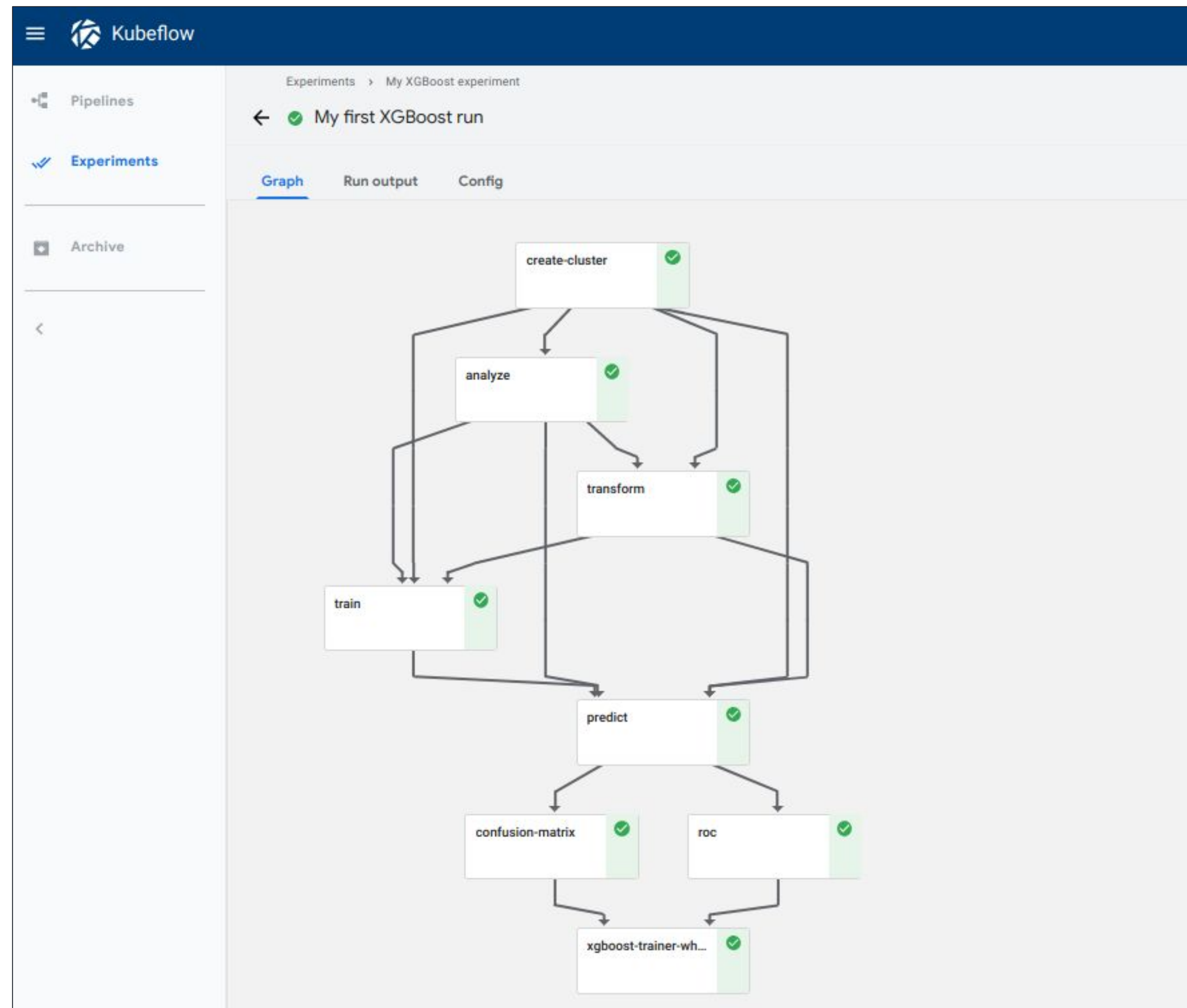
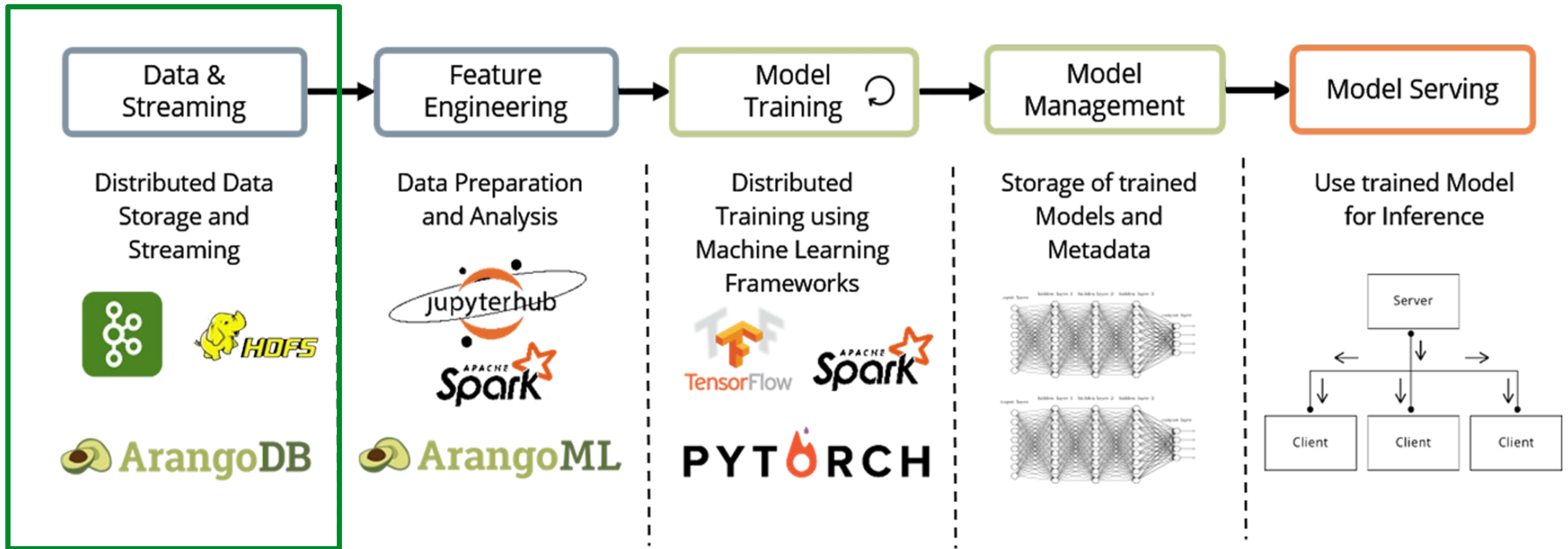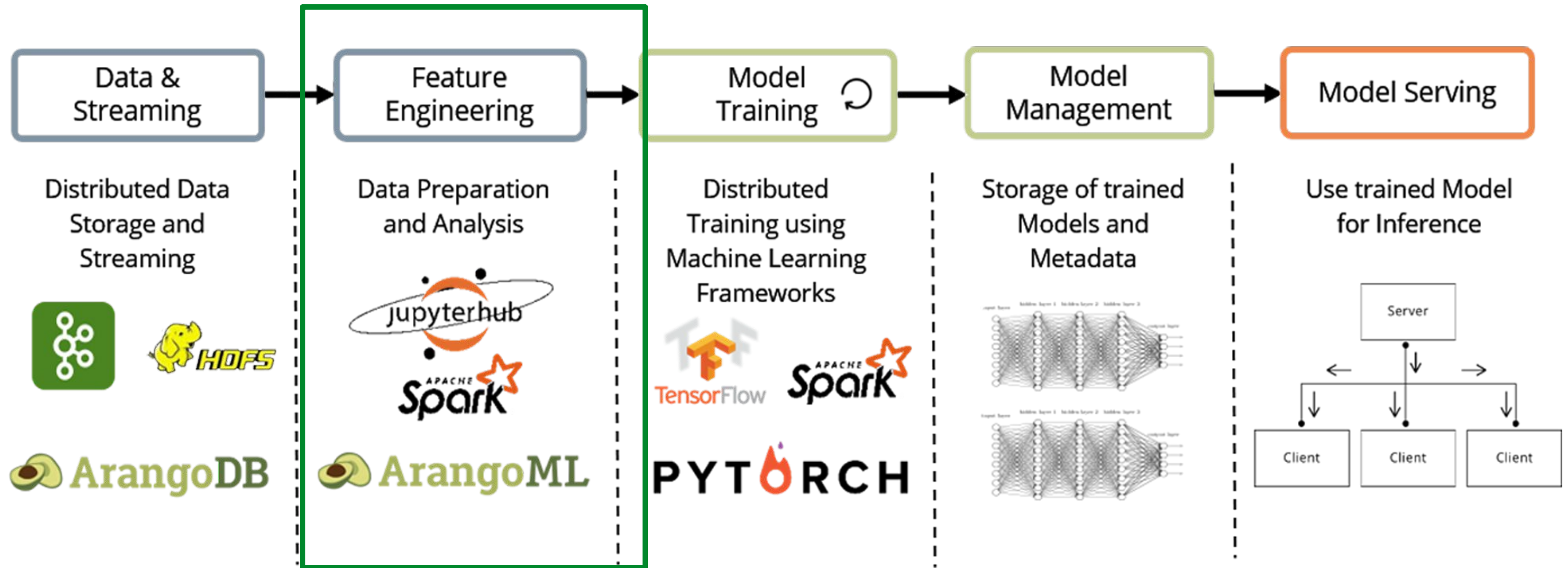# Machine Learning Pipeline

# TensorFlow Extended



| DATA INGESTION | TENSORFLOW DATA VALIDATION | TENSORFLOW TRANSFORM | ESTIMATOR OR KERAS MODEL | TENSORFLOW MODEL ANALYSIS | VALIDATION OUTCOMES | TENSORFLOW SERVING |
|---|---|---|---|---|---|---|
| ExampleGen | StatisticsGen | | | | | |
| | SchemaGen | Transform | Trainer | Evaluator | Pusher | Model Server |
| | Example Validator | | | Model Validator | | |

https://www.tensorflow.org/tfx/guide

# Kubeflow Pipelines



https://www.kubeflow.org/docs/pipelines/

# Databases I

# Databases II

# Graphs and Machine Learning

# Feature Engineering



| Director | Number Movies |
|----------|---------------|
| George_S_Fleming | 10 |
| …. | ….. |
| | |

# Feature Engineering

# ArangoDB

- Native Multi Model Database
  - Stores, K/V, Documents & Graphs

- Distributed
  - Graphs can span multiple nodes

- AQL - SQL-like multi-model query language

- ACID Transactions including Multi Collection Transactions

# Multi-Model?

# Feature Catalogue



- Feature Catalogue ≈ Preprocessing
  Cache + Discovery
  - Uber Michelangelo
  - Logical Clocks
  - Kubeflow FEAST

# Uber Michelangelo

"..there were no systems in place to build reliable, uniform, and reproducible pipelines for creating and managing training and prediction data at scale."

- **Feature store**



https://eng.uber.com/michelangelo/

# Feature Store

**ML Feature Data Warehouse**
Reduce the cost of generating and storing the feature data

**Just-in-time Feature Transforms**
Allow the research teams to experiment with new features and new feature engineering techniques

# Multi-Model ML Demo

# What is next?

**Abstract**

Document Sections

I. Introduction

II. The Graph Neural
Network Model

III. Computational
Complexity Issues

IV. Experimental Results

**Abstract:**
Many underlying relationships among data in several areas of science and engineering, e.g., computer vision, molecular chemistry, molecular biology, pattern recognition, and da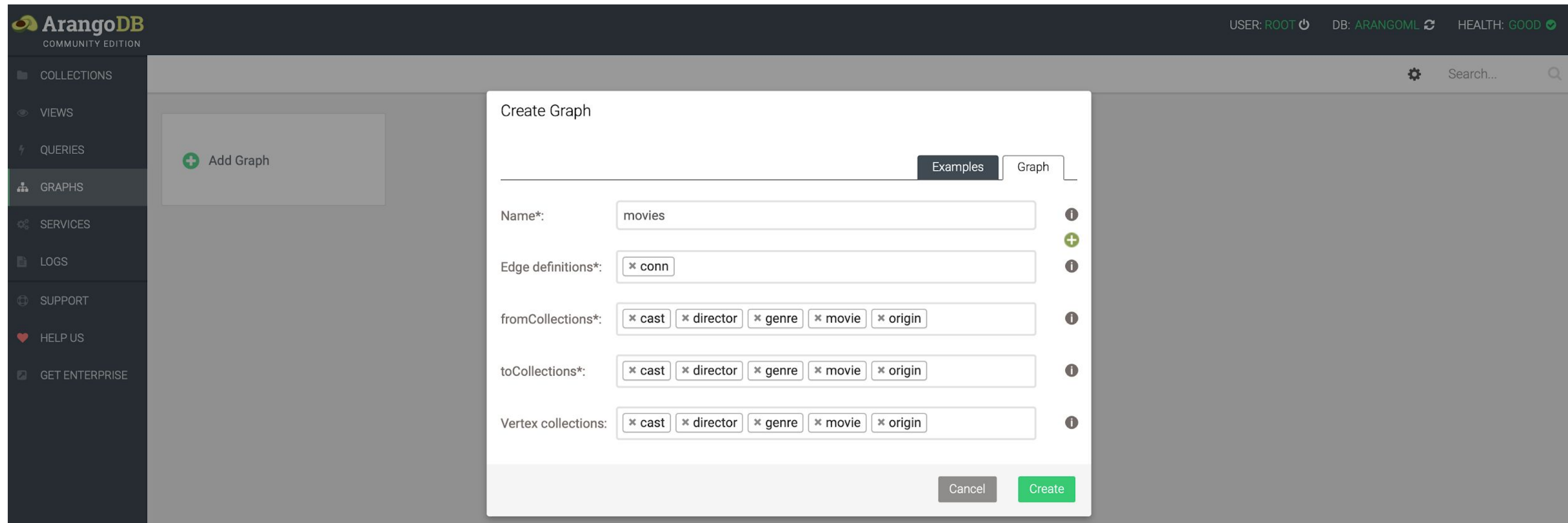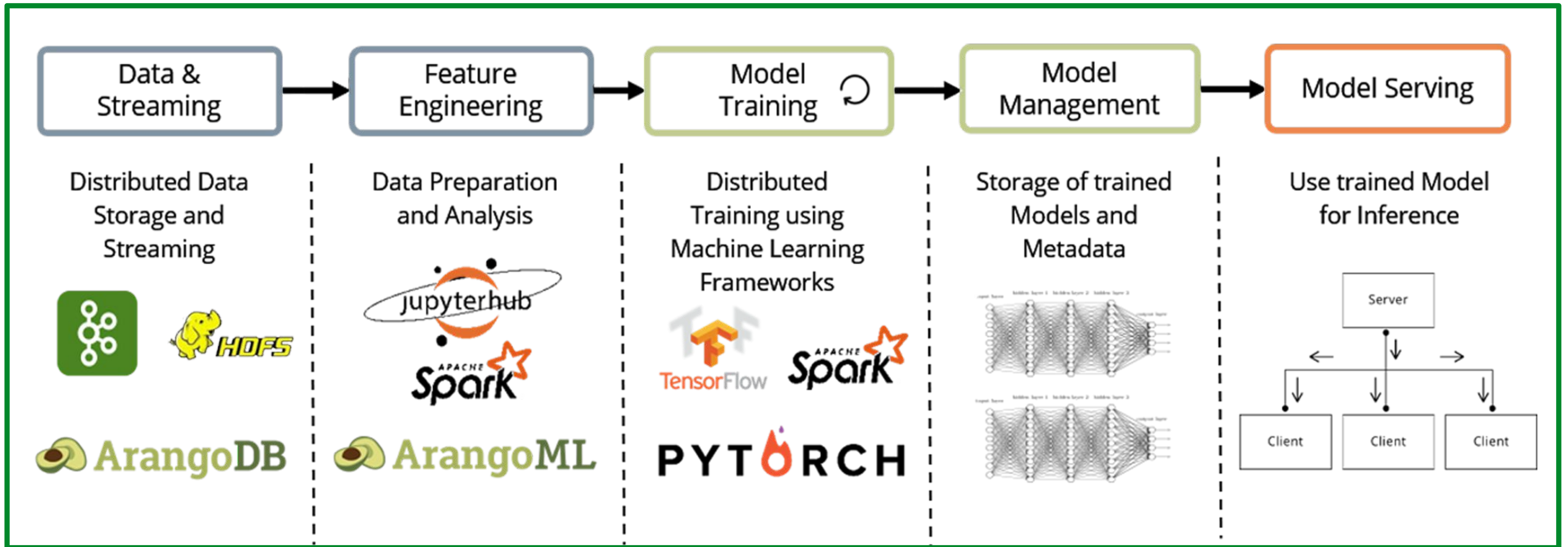ta mining, can be represented in terms of graphs. In this paper, we propose a new neural network model, called graph neural network (GNN) model, that extends existing neural network methods for processing the data represented in graph domains. This GNN model, which can directly process most of the practically useful types of graphs, e.g., acyclic, cyclic, directed, and undirected, implements a function tau(G,n) isin $IR^m$ that maps a graph G and one of its nodes $n$ into an $m$-dimensional Euclidean space. A supervised learning algorithm is derived to estimate the parameters of the proposed GNN model. The computational cost of the proposed algorithm is also considered. Some experimental results are shown to validate the proposed learning algorithm, and to demonstrate its generalization capabilities.

https://ieeexplore.ieee.org/abstract/document/4700287

# Databases III

# Challenges



The Secret Sharer: Evaluating and Testing
Unintended Memorization in Neural Networks

Nicholas Carlini[1,2]    Chang Liu[2]    Úlfar Erlingsson[1]    Jernej Kos[3]    Dawn Song[2]

[1]*Google Brain*    [2]*University of California, Berkeley*    [3]*National University of Singapore*

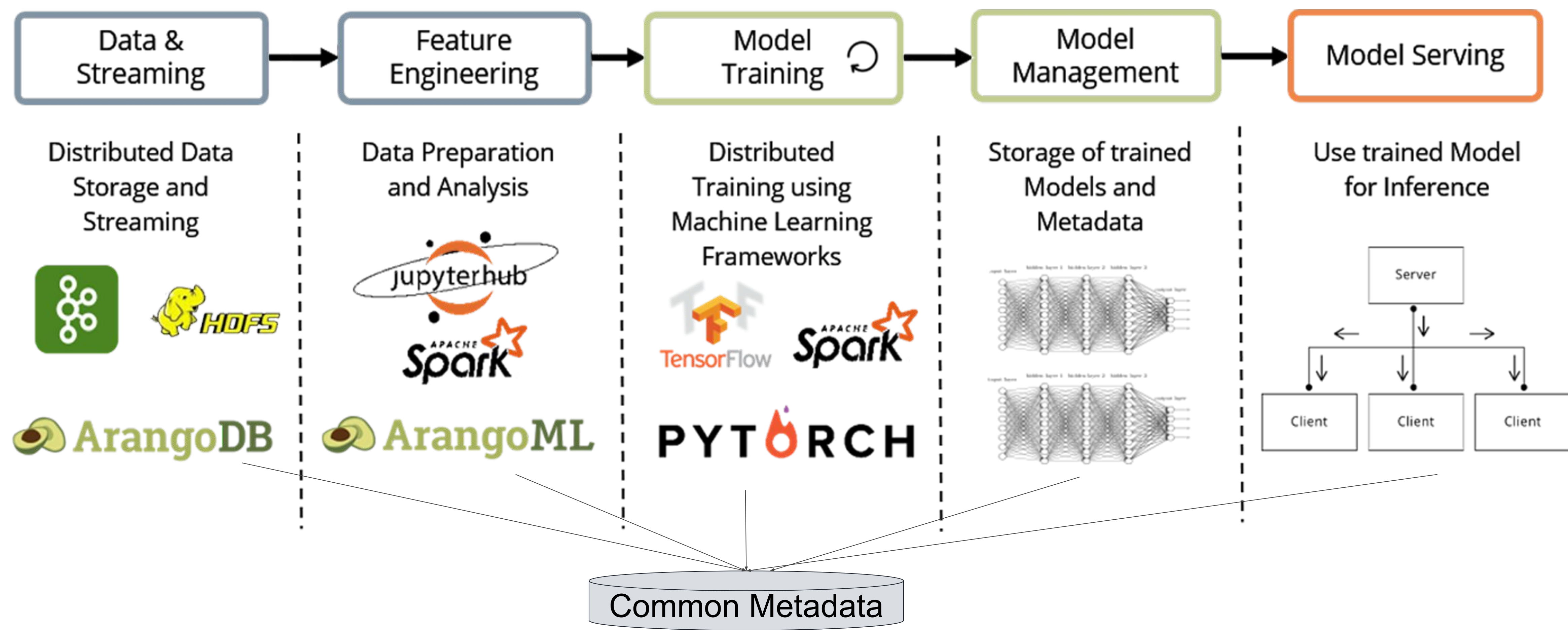https://blog.acolyer.org/2019/09/23/the-secret-sharer/

# Challenges



- **Understand complete provenance of Model**
  - a. Understand Provenance
  - b. Complete version history
  - c. Audit
- **Find all Models in production derived from dataset x**
- **Compare performance of different model performance**
- **Identify reusable steps**
- **Is my serving data distribution the same as for training data**
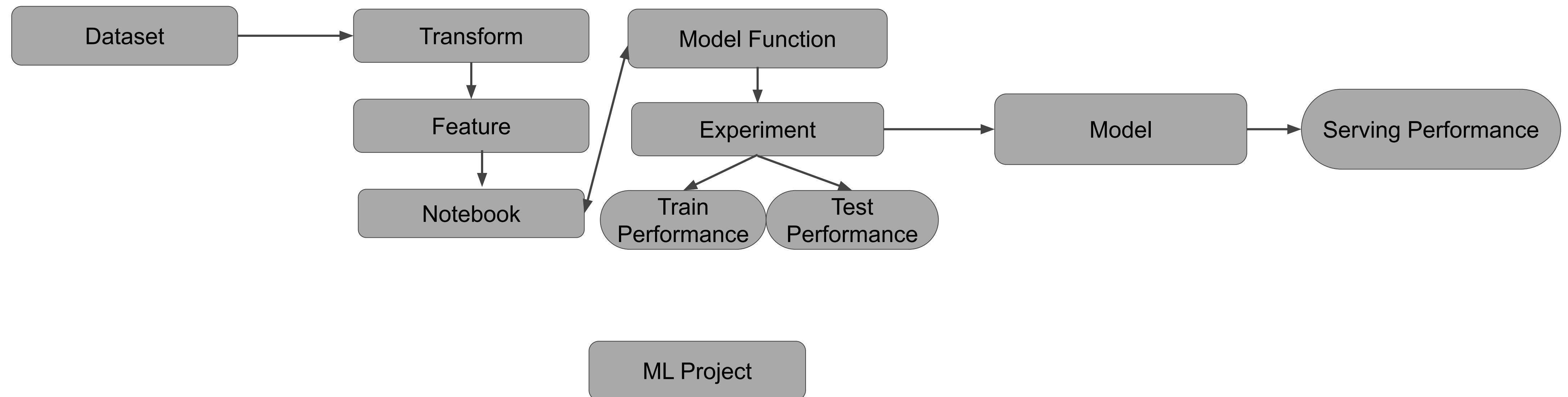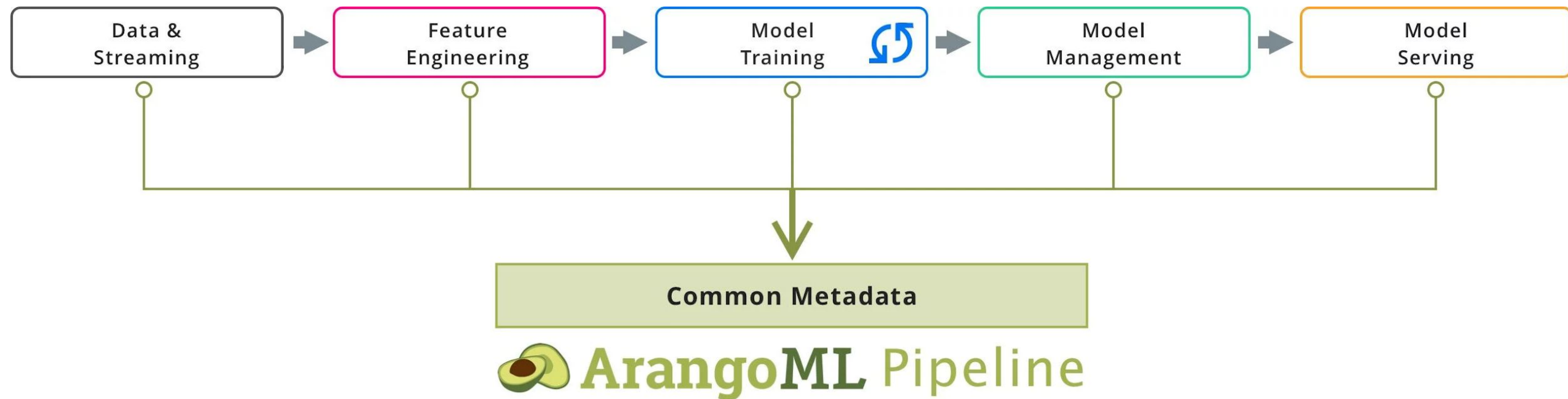- ...

# From Data to Metadata….

# Metadata?

In this context, *metadata* means information about executions (runs), models, datasets, and other artifacts.
*Artifacts* are the files and objects that form the inputs and outputs of the components in your ML workflow.
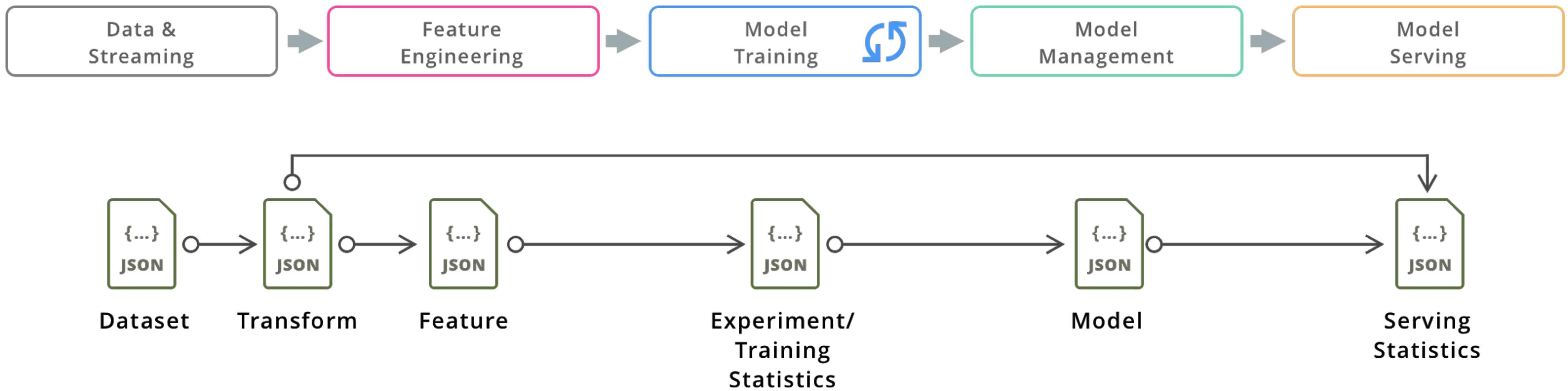
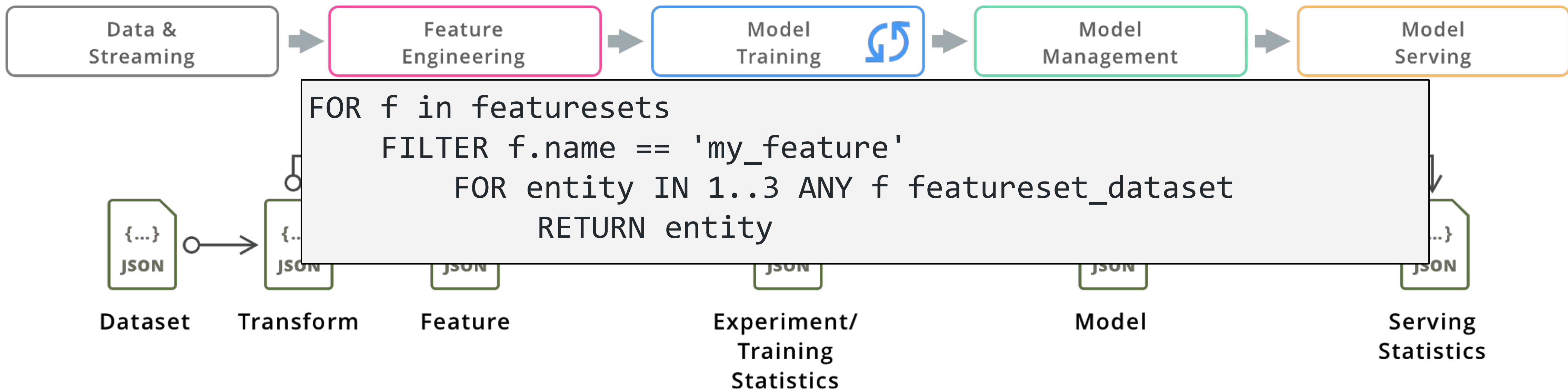https://www.kubeflow.org/docs/components/misc/metadata/

# ArangoML Pipeline

"A common extensible metadata layer for ML pipelines which
allows Data Scientists and DataOps to manage all information
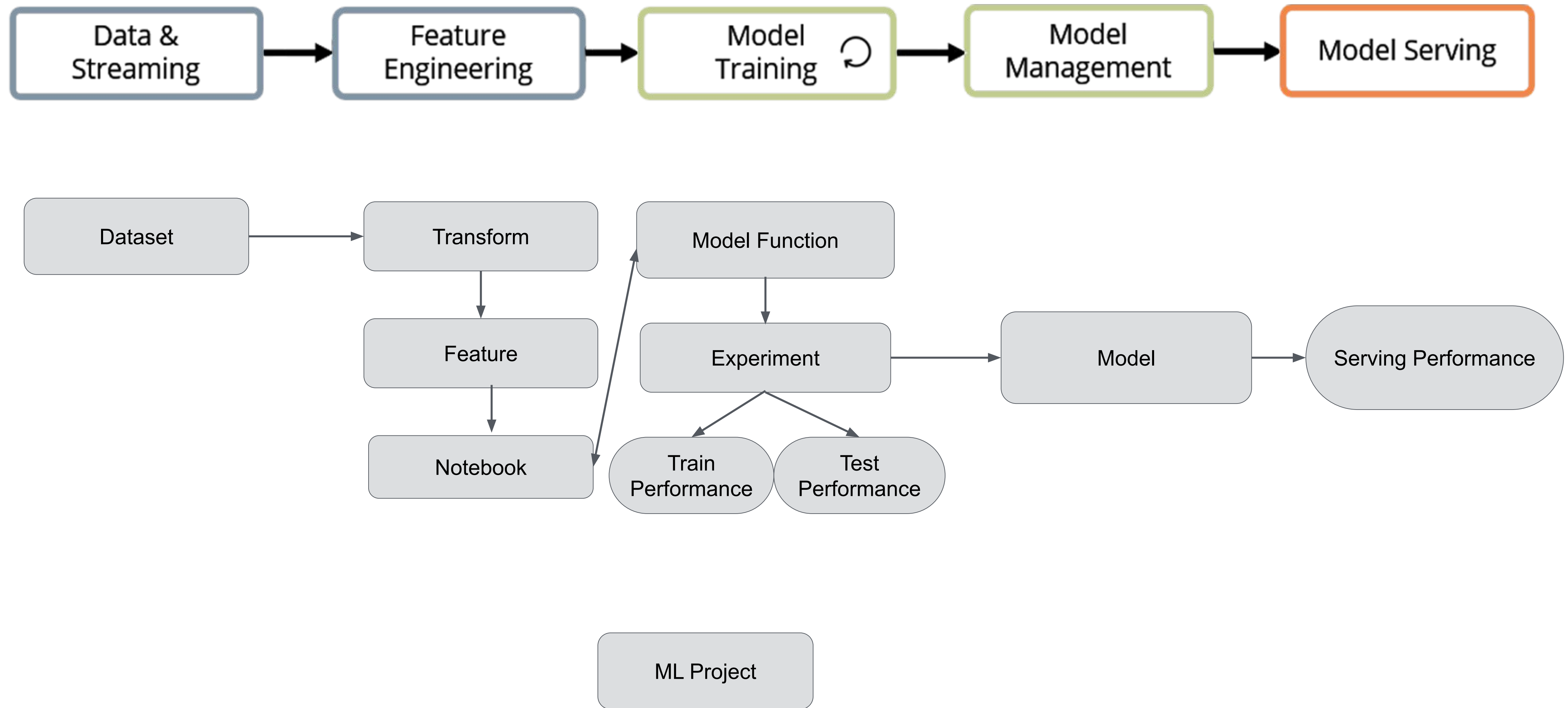related to their ML pipelines in one place."

# Multi-Model Metadata

# Multi-Model Metadata



```
FOR f in featuresets
    FILTER f.name == 'my_feature'
        FOR entity IN 1..3 ANY f featureset_dataset
            RETURN entity
```

# ArangoML "Schema"

# Discover



https://github.com/arangoml/arangopipe

# Visualization



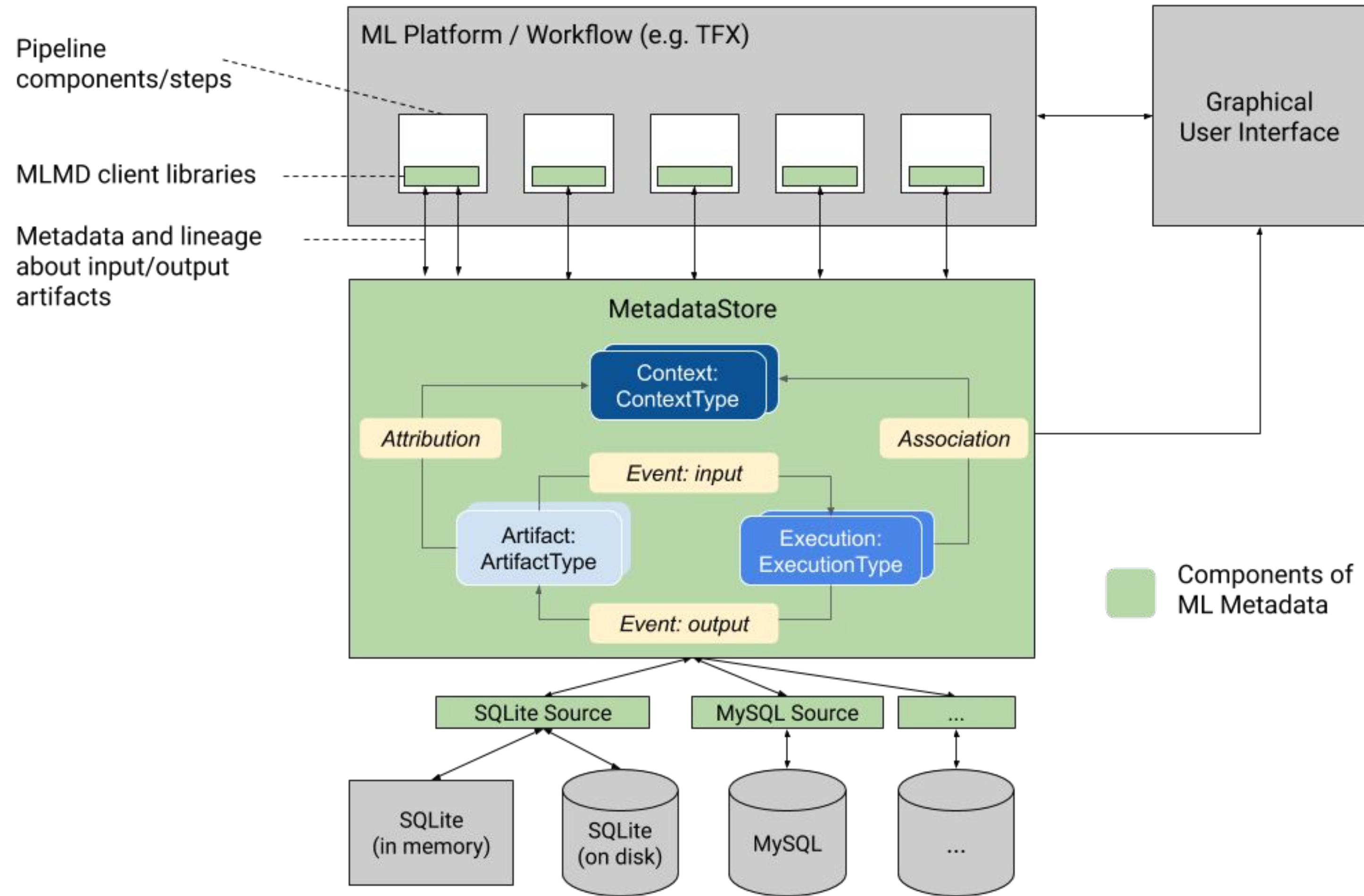https://github.com/arangoml/arangopipe

# **Arango**ML

- Python package
- HTTP API
- TFX Integration [coming shortly]

```
from arangopipe.arangopipe_api import ArangoPipe

ap = ArangoPipe(conn_config)
model_info = {"name": "hyper-param-optimization",  "type": "hyper-opt-experiment"}
model_reg = ap.register_model(model_info, project = "Housing_Price_Estimation_Project")
```

https://github.com/arangoml/arangopipe

# TFX MLMD



https://www.tensorflow.org/tfx/guide/mlmd

# Kubeflow Metadata



https://www.kubeflow.org/docs/components/misc/metadata/

# Thanks for listening!





- @arangoml
- https://github.com/arangoml/arangopipe
- Demo

  https://github.com/arangoml/arangopipe/blob/master/arangopipe/arangopipe_examples.ipynb

- @arangodb
- https://www.arangodb.com/
- Demo

  https://github.com/arangoml/knowlegegraph-demo