# Brief History of Hops

**World's fastest HDFS**
16X increase in throughput on Spotify workload (USENIX FAST)

**World's First GPUs-as-a-Resource**
support in a Hadoop platform

**World's First Open Source Feature Store**
for Machine Learning

2017 — 2018 — 2019 →

**Winner of IEEE Scale Challenge 2017**
with HopsFS - 1.2m ops/sec

**World's First**
Hierarchical File System to store small files in metadata on NVMe disks

**World's First**
Hierarchical Filesystem with Multi Data Center High Availability

"If you're working with big data and Hadoop, **this one paper could repay your investment** in the Morning Paper many times over…. **HopsFS is a huge win.**"
- *Adrian Colyer, The Morning Paper*

LOGICAL CLOCKS

# Quick overview of Hops/Hopsworks

The **only** open-source data platform to support:

- Project-based multi-tenancy
- On-premise resource management of GPUs (>1 server)
- Per-Project Python Dependencies with Conda
- Feature Store
- Jupyter notebooks as Jobs (Airflow)
- Free-text search for files/dirs in the filesystem
- NVMe to store small files in filesystem metadata

# Example workflow in Hopsworks at Scale

1.  Insert 1m images (<100kb) in seconds
2.  Train a DNN classifier using 100s of GPUs
3.  Run a Spark job to identify all objects in the 1m images  and add the image annotations (JSON) as extended metadata to HopsFS
4.  "show me the images with >3 bicycles" and get a sub-second response.

Ops folks: Remove the image directory, and elasticsearch is auto-cleaned up!
Data scientists: Do it all in Jupyter notebooks and Python (if you want)!

# The Future is Cloud-Native...but what about the FS?

**Kubernetes**

Does it have to be S3?

What will the Cloud-Native Filesystem be?
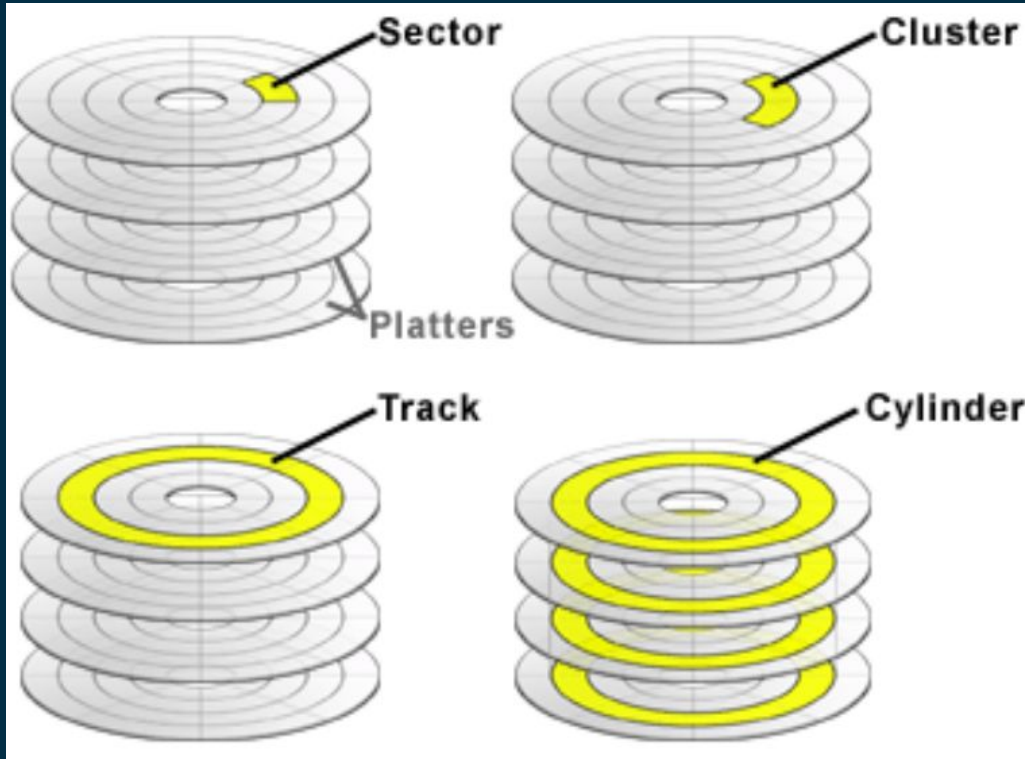
# A Brief History of Data

# MapReduce v0.01-alpha



IBM 082 Punch Card Sorter

Scan -> Sort -> Scan -> Sort ….



Not Fault Tolerant!

LOGICAL CLOCKS
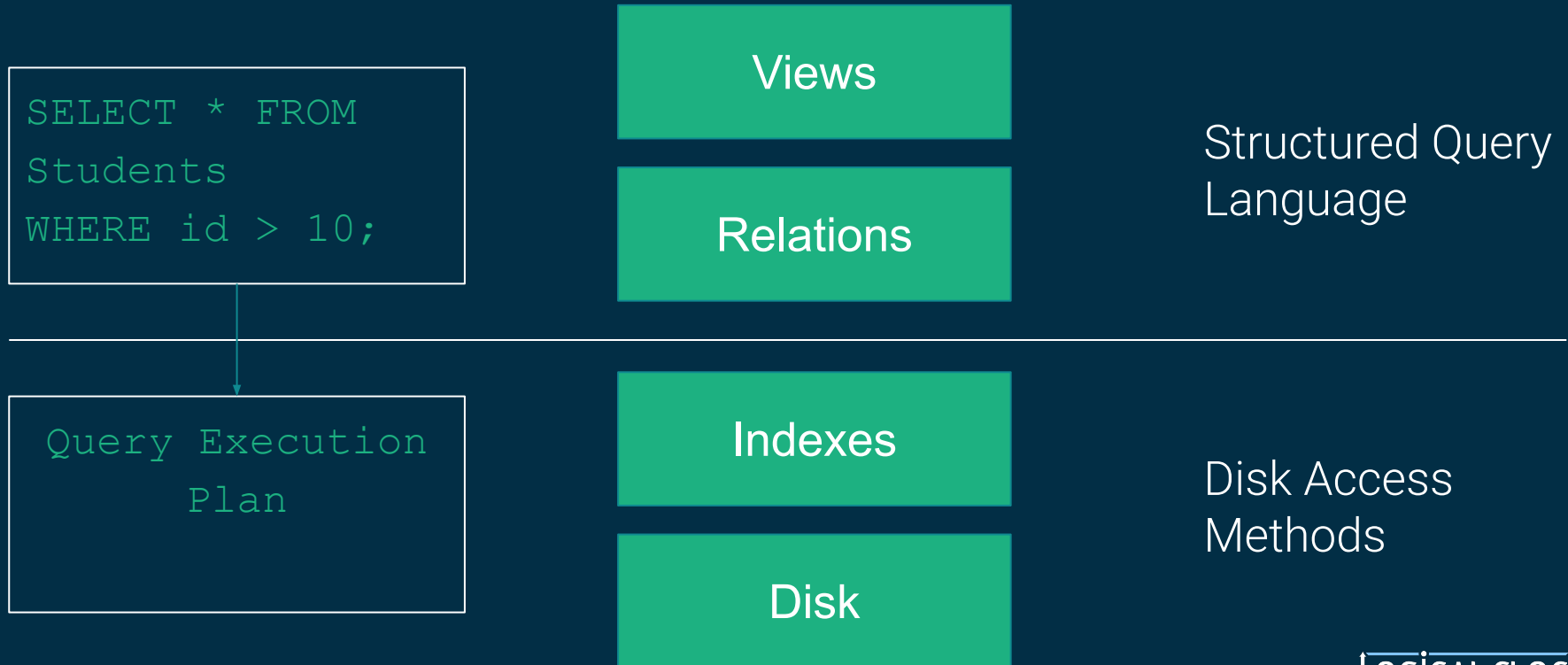
# First DBMS' and Filesystems were Disk Aware



Early Filesystems' block size was tightly coupled to the sector size of a disk

LOGICAL CLOCKS

# Hierarchical and Network DB Systems



You had to know what you want, and how to find it on disk.

LOGICAL CLOCKS

# Codd's Relational Model and SystemR

```
SELECT * FROM
Students
WHERE id > 10;
```

Query Execution
Plan

Views

Relations

Indexes

Disk

Structured Query
Language

Disk Access
Methods

LOGICAL CLOCKS

# +30 years..Data Volumes outgrew Relational DBs



Data volumes got too large for single-server SQL DBs
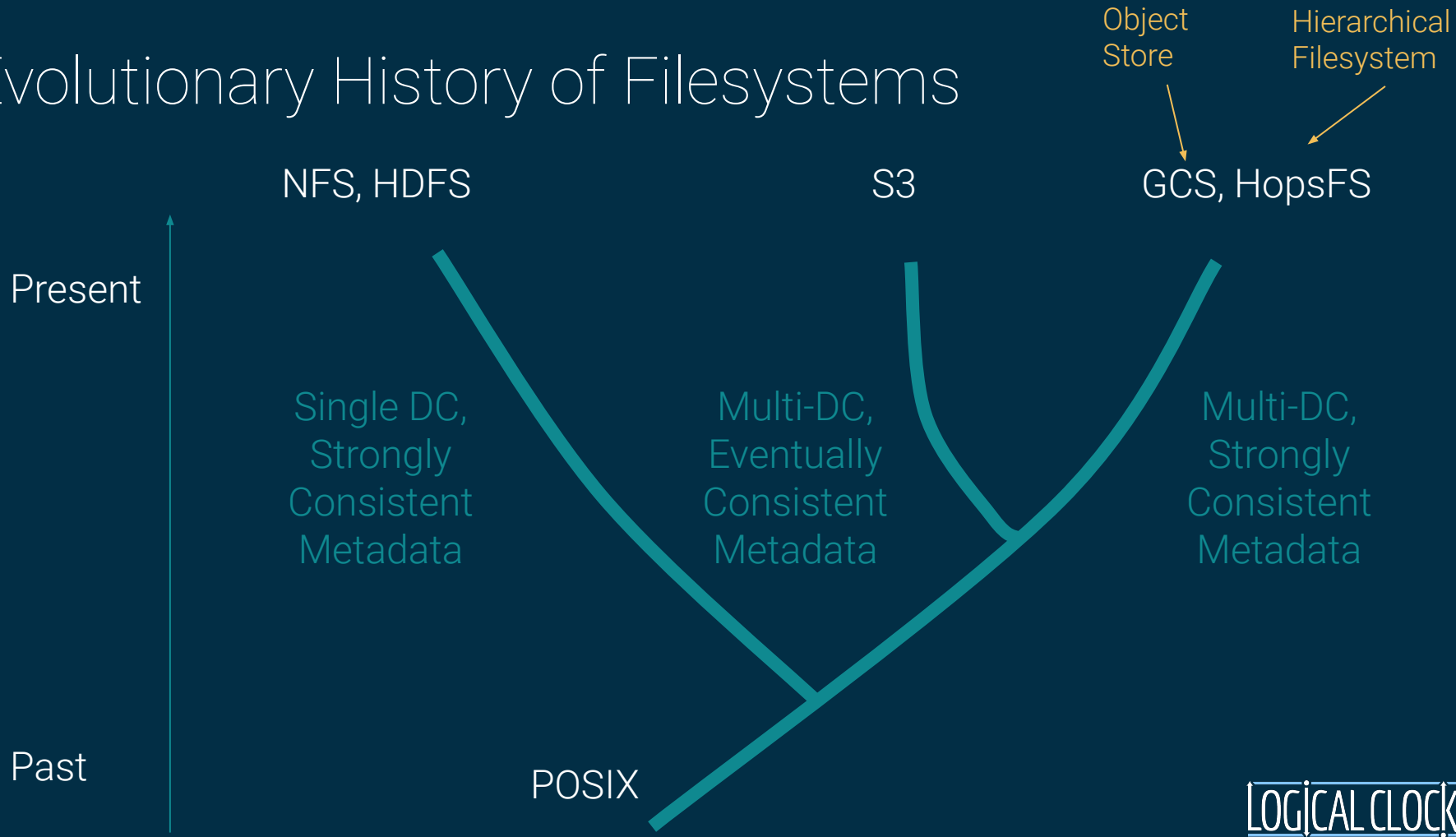
And thus, the NoSQL movement was born…

…..only to be quickly out-evolved

LOGICAL CLOCKS

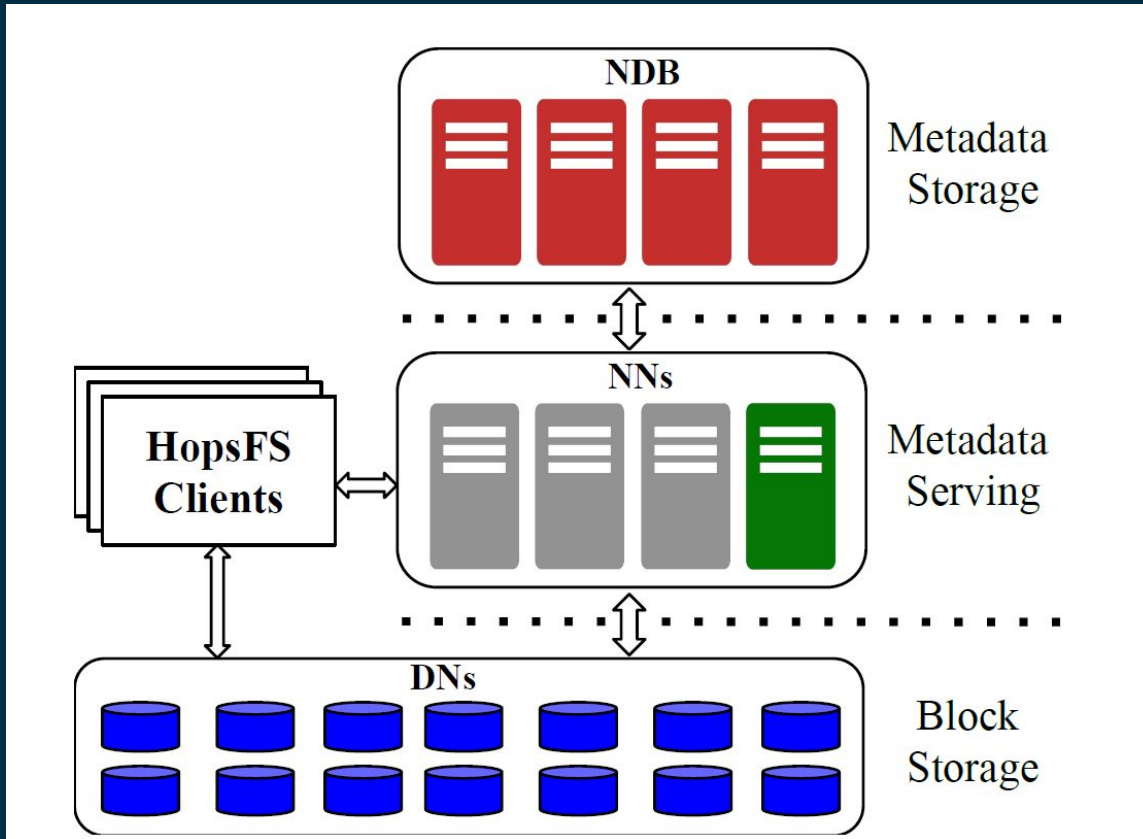# Evolutionary History of SQL Datastores

# What about Filesystems?

# Evolutionary History of Filesystems

NFS, HDFS                    S3                    GCS, HopsFS

Object Store          Hierarchical Filesystem

Present

Single DC,            Multi-DC,             Multi-DC,
Strongly              Eventually            Strongly
Consistent            Consistent            Consistent
Metadata              Metadata              Metadata

Past

POSIX

LOGICAL CLOCKS

# Why is Strongly Consistent Metadata important?

- ● POSIX-like semantics
  - ○ Insert a file in a dir, and yes, it will be there!
- ● Atomic rename
  - ○ Building block for scalable SQL systems
- ● Consistent change data capture (changelog)
  - ○ Data provenance
  - ○ Search/Index/tag the filesystem namespace

# HopsFS uses NDB for Strongly Consistent Metadata



Make these
Layers
Data-Center HA

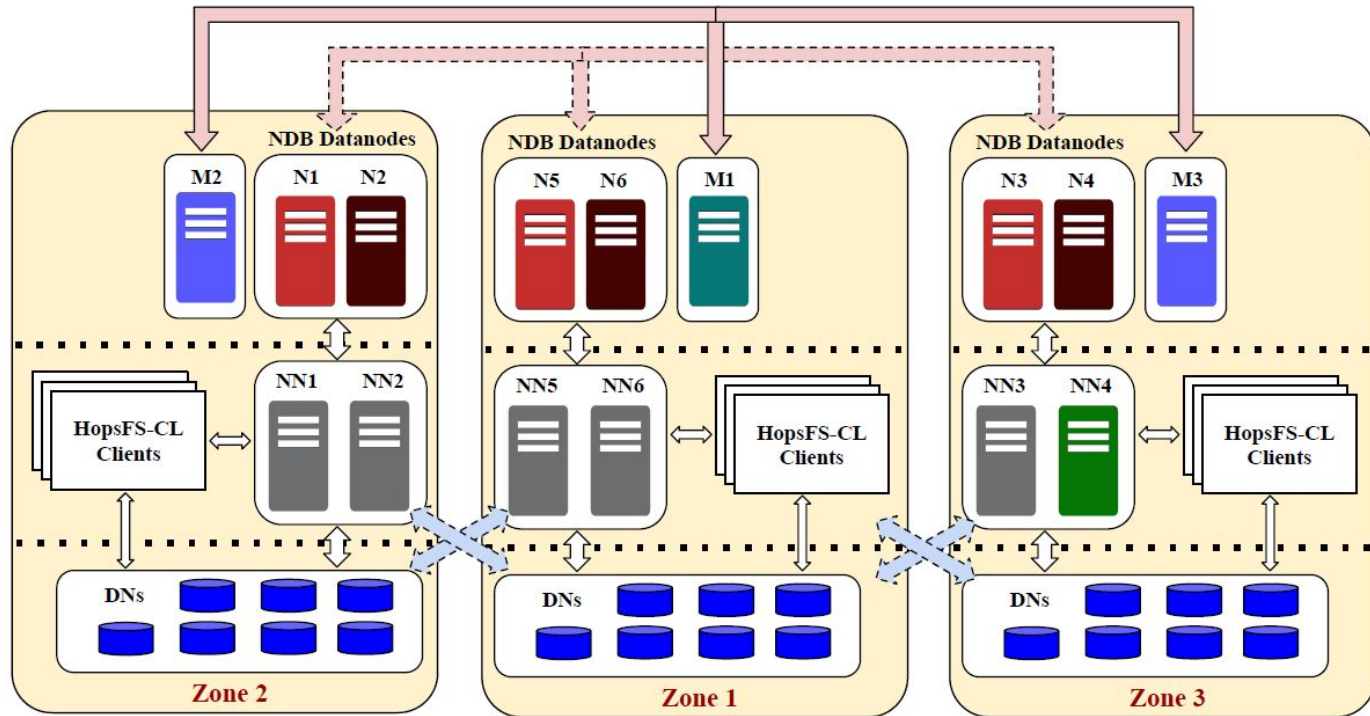LOGICAL CLOCKS

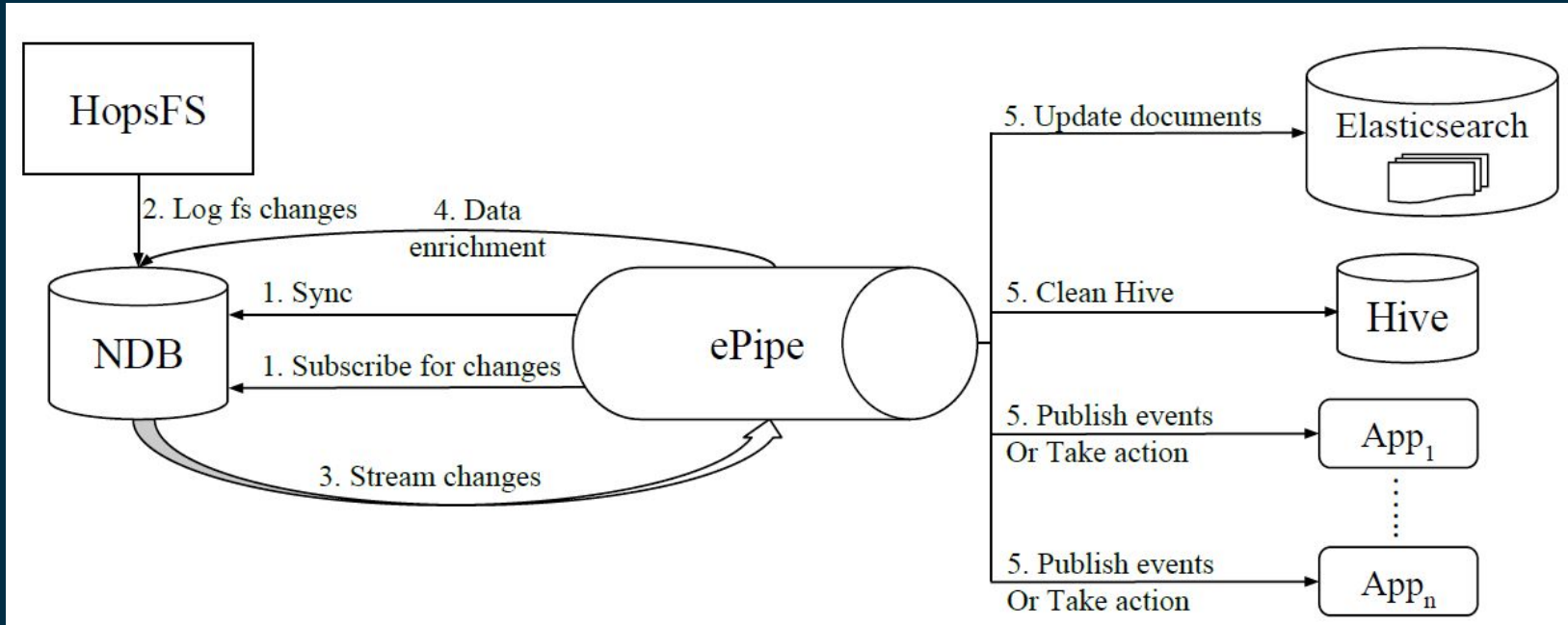# Multi-DC HopsFS affects every layer of the stack



Database nodes DC-aware

Namenodes DC-aware

36% performance improvements by optimizing for  DC local operations

LOGICAL CLOCKS

# Triple replication also possible with HopsFS

# Change Data Capture for HopsFS with Epipe



Overhead of running ePipe on the Spotify Hadoop workload: 4.77%

ePipe: Near Real-Time Polyglot Persistence of HopsFS Metadata, Ismail et al, CCGrid, 2019.

# HOPS STRIKES BACK IN THE CLOUD

**Availability:**   Highly available across Data Centers (AZ)

**Performance:**  >1.6m Ops/Second on Spotify workload (GCE, 3 AZs)
                  NVMe disks used to store small files in metadata layer

**Security:**     TLS-based security

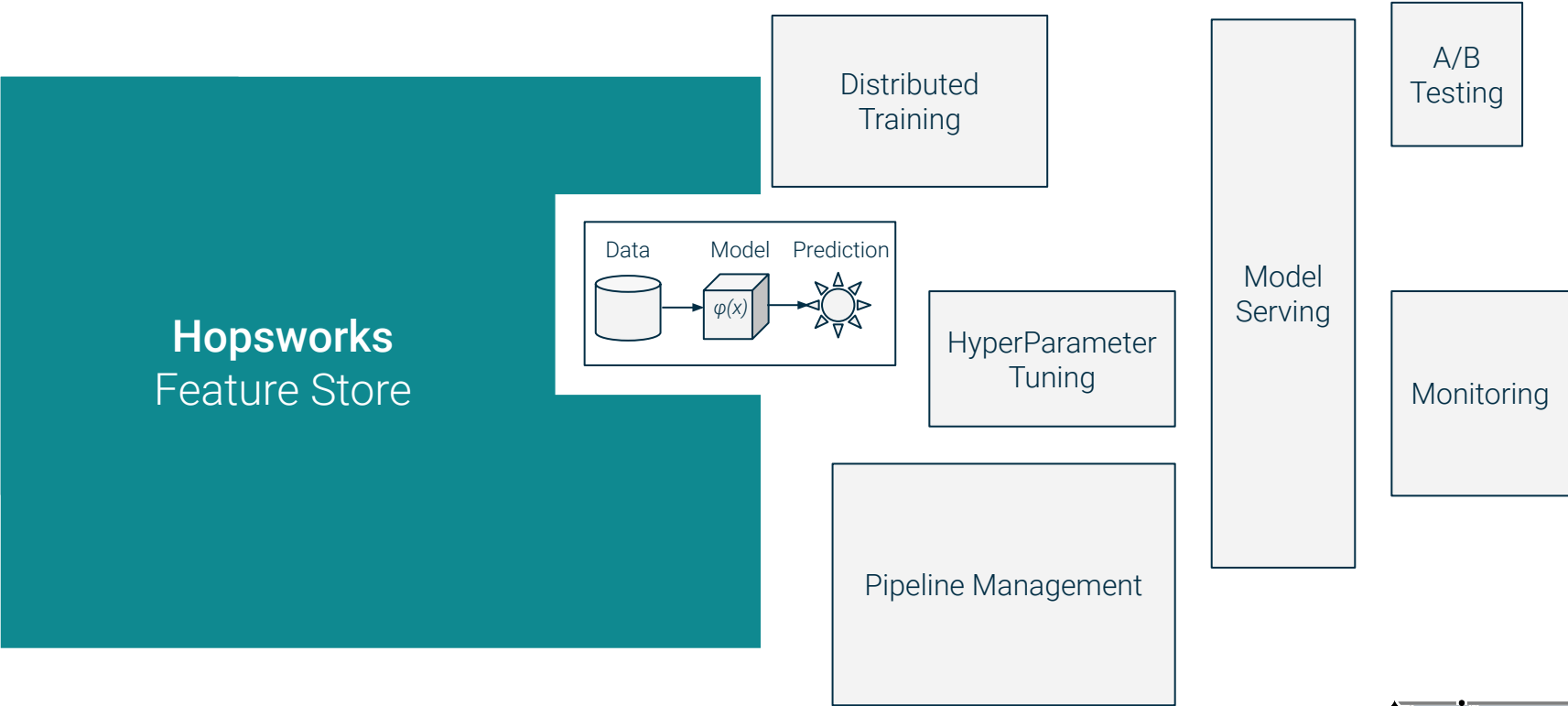**HDFS API**:     Native support in Spark, Flink, TensorFlow, etc.

# Hopsworks - a platform for Data Intensive AI built on Hops

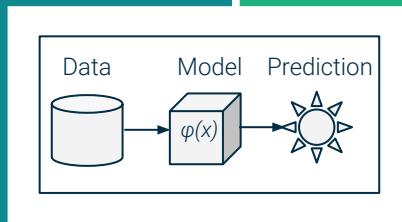# Hopsworks hides the Complexity of Deep Learning

Data validation

Distributed Training

A/B Testing

Data Collection



Data    Model    Prediction
φ(x)

HyperParameter Tuning

Model Serving

Hardware Management

Feature Engineering

Pipeline Management

Monitoring

*Figure from "Technical Debt in Machine Learning Systems", Google research paper*

LOGICAL CLOCKS

# Hopsworks hides the Complexity of Deep Learning

**Hopsworks**
Feature Store

Data    Model    Prediction

$\varphi(x)$

Distributed Training

HyperParameter Tuning

Model Serving

A/B Testing

Monitoring

Pipeline Management

LOGICAL CLOCKS

# Hopsworks hides the Complexity of Deep Learning



**Hopsworks**
Feature Store

Data    Model    Prediction

φ(x)

**Hopsworks**
REST API

**Datasources**

# Hopsworks

The Platform for Data Intensive AI
-
Machine Learning, Deep Learning &
Model serving

Applications

**API**

Dashboards

LOGICAL CLOCKS

# What is Hopsworks?

## Efficiency & Performance

**Feature Store**
Data warehouse for ML

**Distributed Deep Learning**
Faster with more GPUs

**HopsFS**
NVMe speed with Big Data

**Horizontally Scalable**
Ingestion, DataPrep,
Training, Serving

## Usability & Process

**Jupyter/Python Development**
Notebooks in pipelines

**Version Everything**
Code, Infrastructure, Data

**Model Serving on Kubernetes**
TF Serving, MLeap, SkLearn

**End-to-End ML Pipelines**
Orchestrated by Airflow

## Security & Governance

**Secure Multi-Tenancy**
Project-based restricted access

**Encryption At-Rest, In-Motion**
TLS/SSL everywhere

**AI-Asset Governance**
Models, experiments, data, GPUs

**Data/Model/Feature Lineage**
Discover/track dependencies

LOGICAL CLOCKS

# Which services require Distributed Metadata (HopsFS)?

## Efficiency & Performance

**Feature Store**
Data warehouse for ML

**Distributed Deep Learning**
Faster with more GPUs

**HopsFS**
NVMe speed with Big Data

**Horizontally Scalable**
Ingestion, DataPrep,
Training, Serving

## Usability & Process

**Jupyter/Python Development**
Notebooks in pipelines

**Version Everything**
Code, Infrastructure, Data

**Model Serving on Kubernetes**
TF Serving, MLeap, SkLearn

**End-to-End ML Pipelines**
Orchestrated by Airflow

## Security & Governance

**Secure Multi-Tenancy**
Project-based restricted access

**Encryption At-Rest, In-Motion**
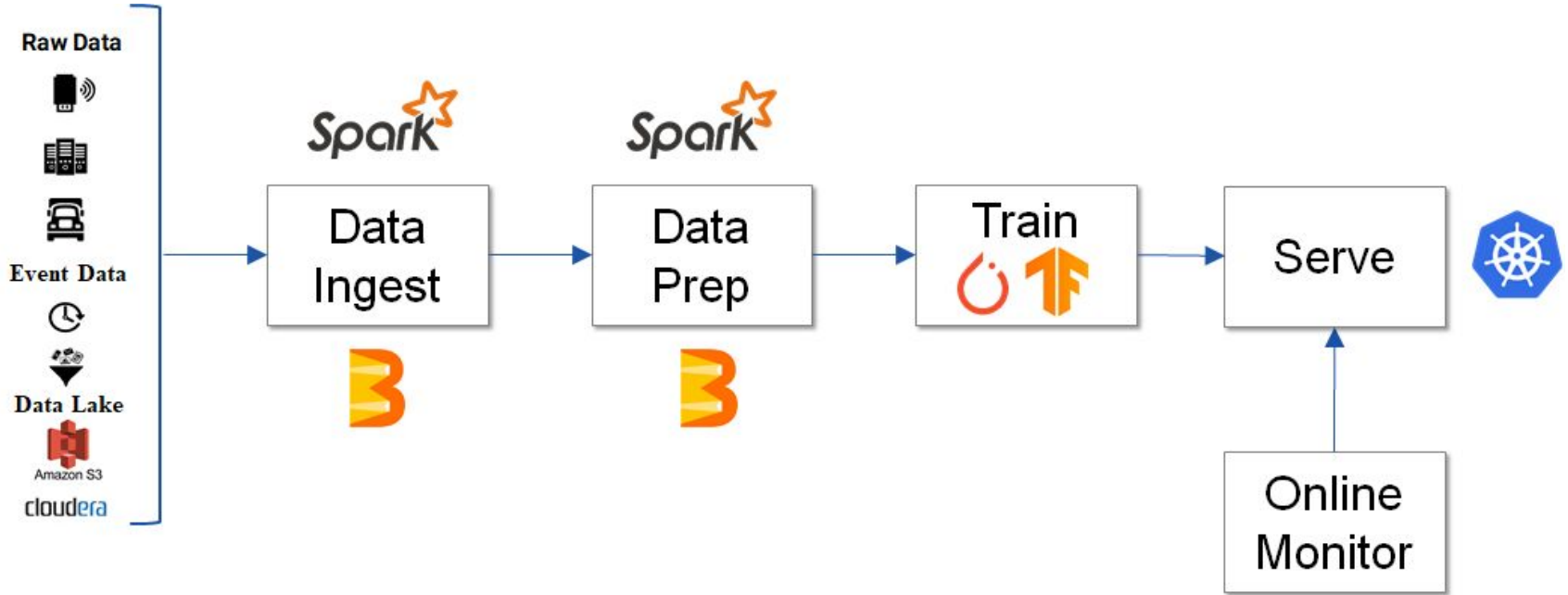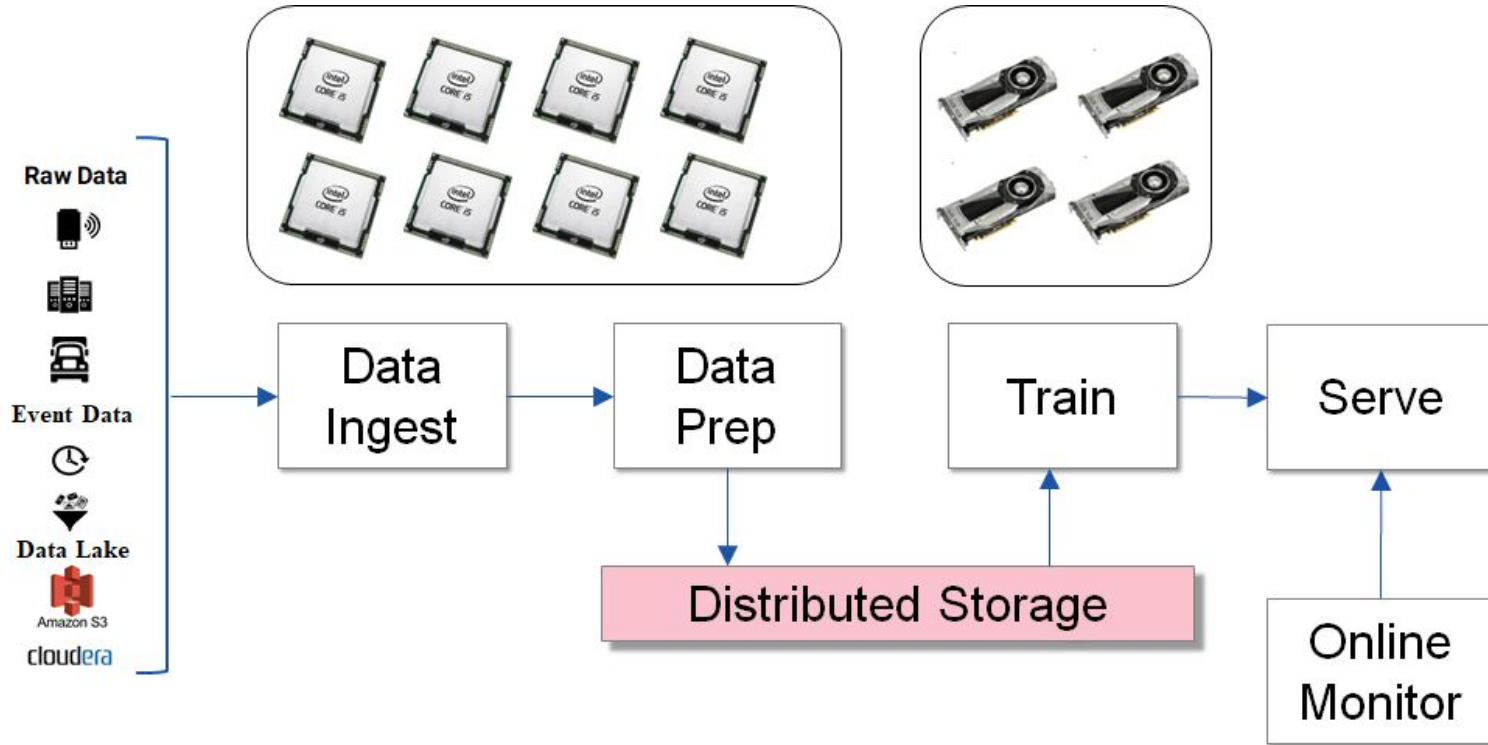TLS/SSL everywhere

**AI-Asset Governance**
Models, experiments, data, GPUs

**Data/Model/Feature Lineage**
Discover/track dependencies
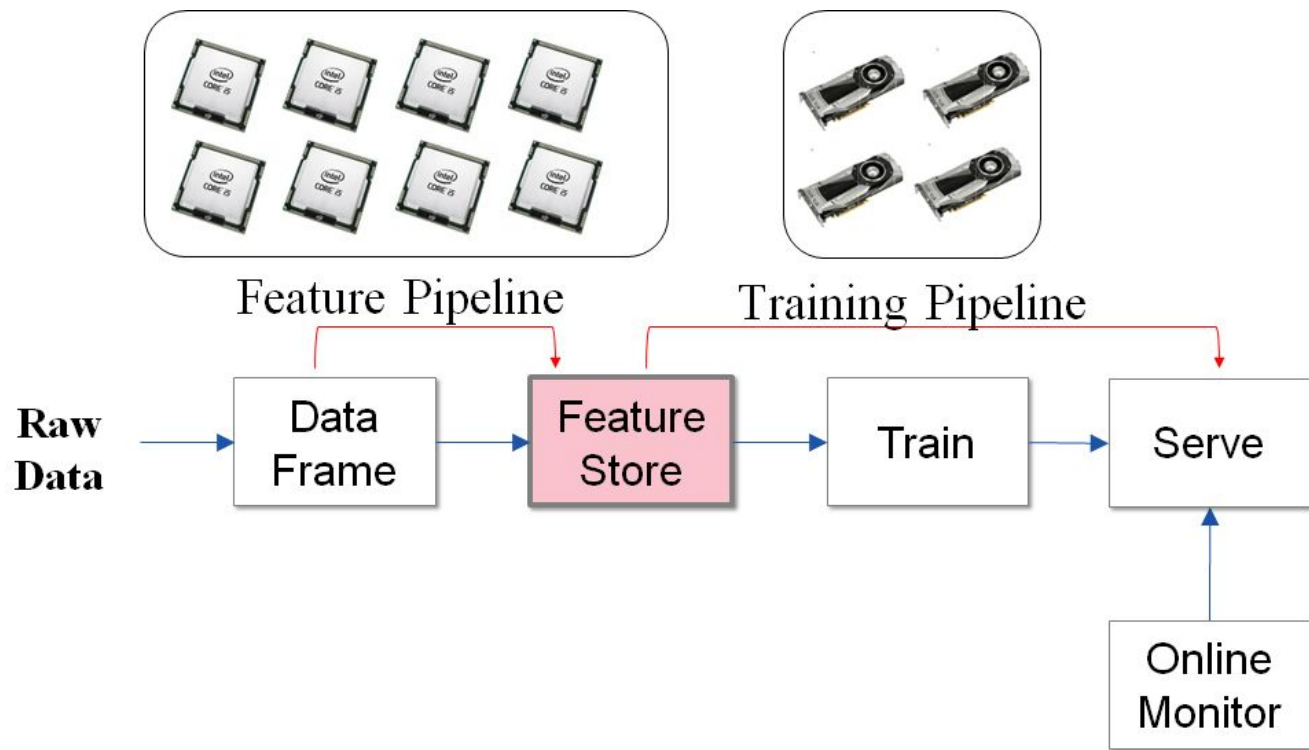
LOGICAL CLOCKS

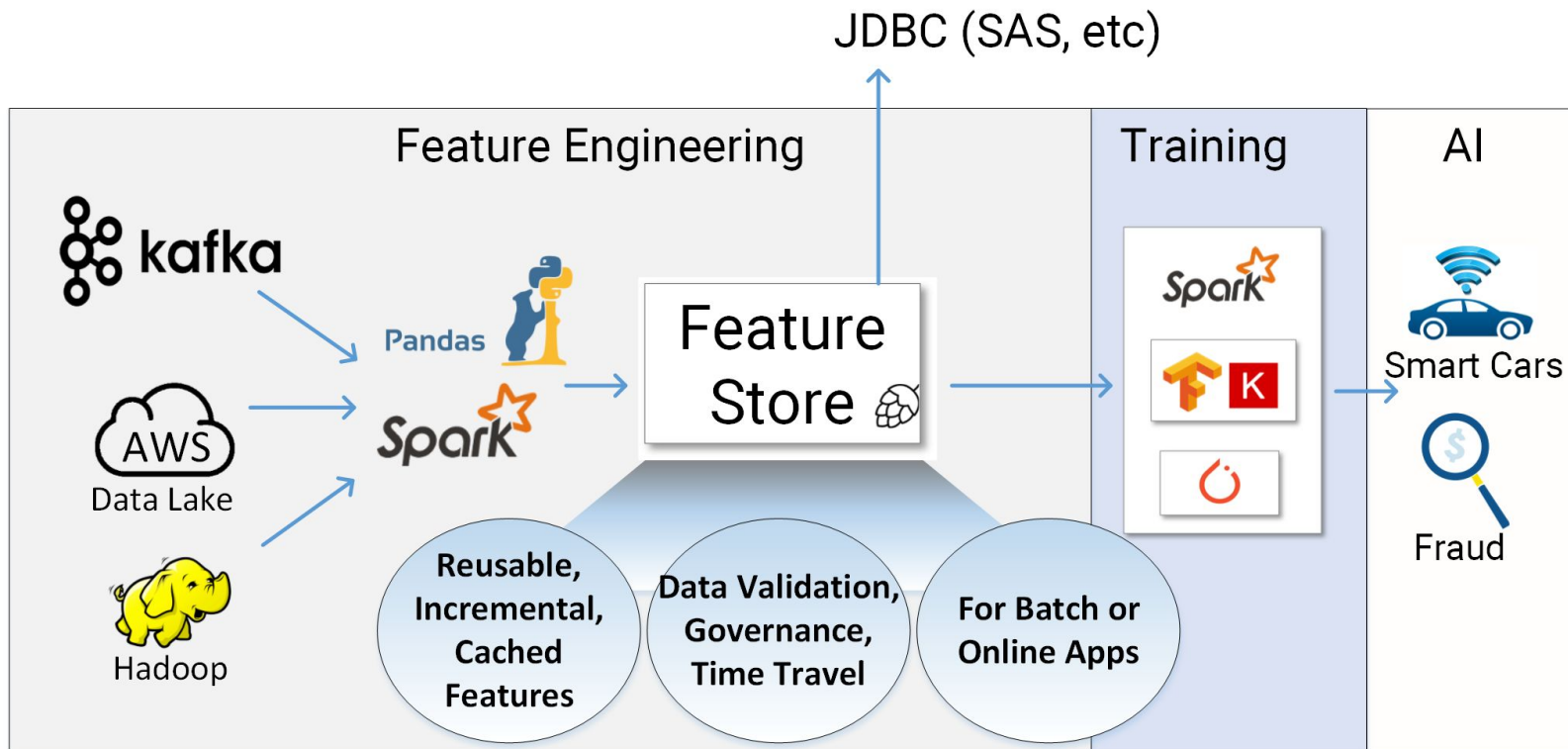# End-to-End ML Pipelines in Hopsworks

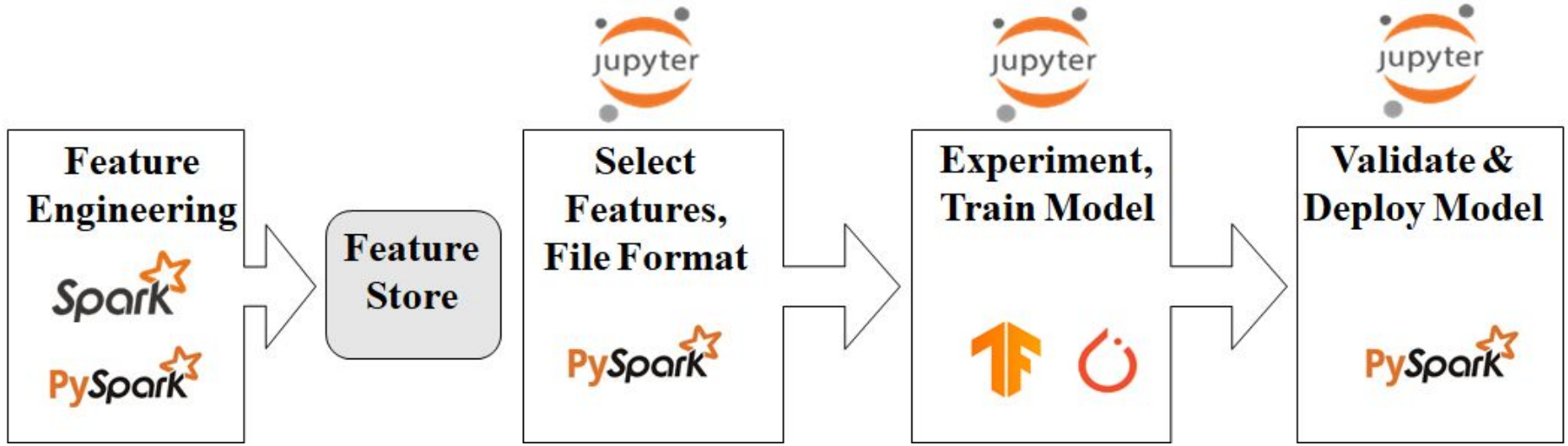# End-to-End Pipelines can be factored into stages

# Typical Feature Store Pipelines

# Hopsworks' Feature Store

Dev View: Pipelines of Jupyter Notebooks in Airflow

# How to get started with Hopsworks?

Register for a free account at: www.hops.site
Images available for AWS, GCE, Virtualbox.

We need your support. Star us, tweet about us!

https://github.com/logicalclocks/hopsworks

@hopsworks