

# Feature Store: the missing data layer in ML pipelines?<sup>1</sup>

*Spotify ML Guild Fika*

Kim Hammar

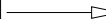
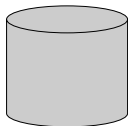
*kim@logicalclocks.com*

February 26, 2019

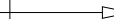
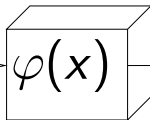


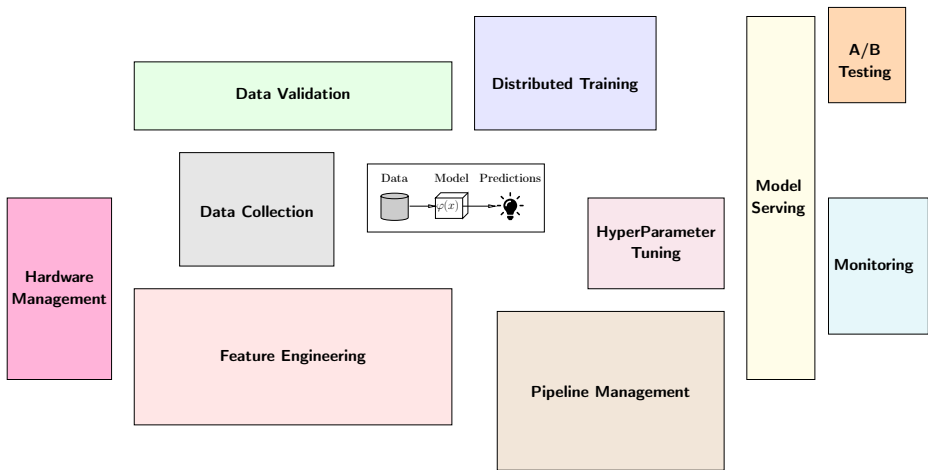
<sup>1</sup>Kim Hammar and Jim Dowling. *Feature Store: the missing data layer in ML pipelines?*  
<https://www.logicalclocks.com/feature-store/>. 2018.

Data



Model Predictions





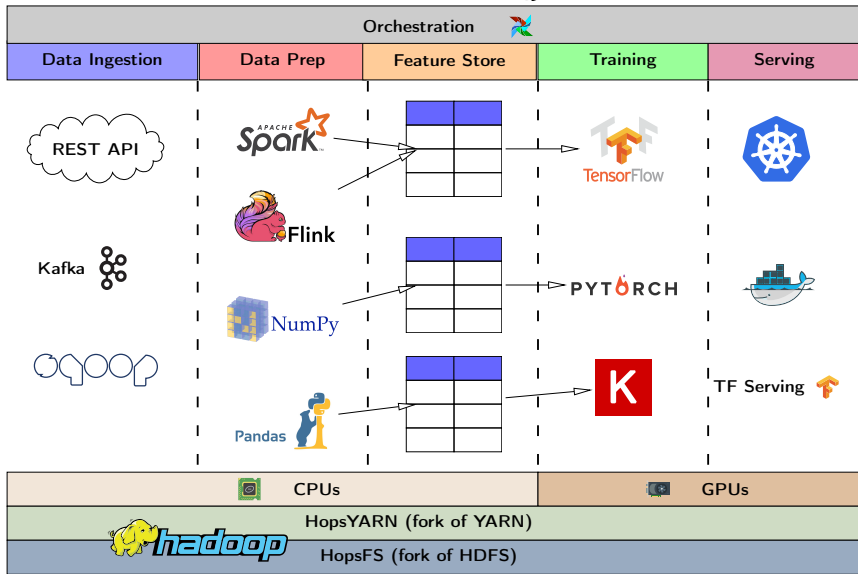
2

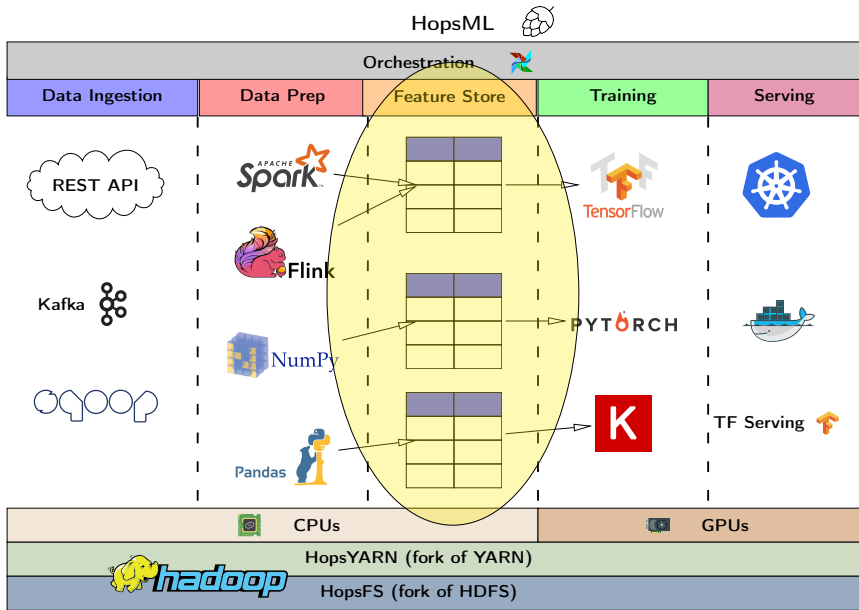
<sup>2</sup>Image inspired from Sculley et al. (Google) *Hidden Technical Debt in Machine Learning Systems* ▶

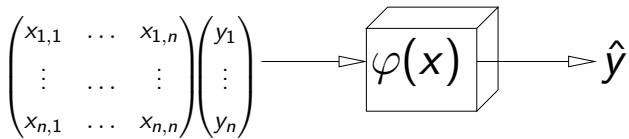
- 1 **Hopsworks**: Quick background of the platform
- 2 **What** is a Feature Store
- 3 **Why** You Need a Feature Store, Things to Consider:
  - How to encourage feature reuse?
  - How to store large-scale datasets for deep learning?
  - How to serve features for inference?
- 4 **How** to Build a Feature Store (Hopsworks Feature Store Case Study)
- 5 **Demo**

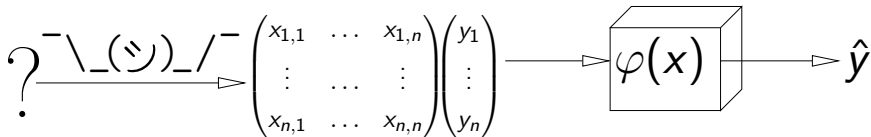


## Orchestration



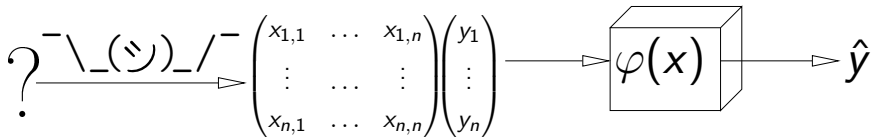






<sup>3</sup>Jeremy Hermann and Mike Del Balso. *Scaling Machine Learning at Uber with Michelangelo*. <https://eng.uber.com/scaling-michelangelo/>. 2018.



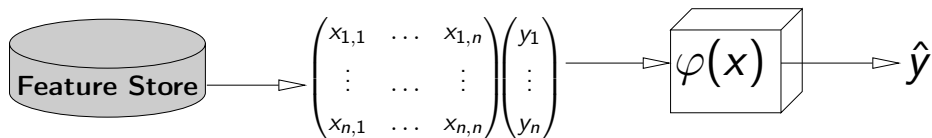


*"Data is the hardest part of ML and the most important piece to get right."*

*Modelers spend most of their time selecting and transforming features at training time and then building the pipelines to deliver those features to production models."*

- Uber<sup>3</sup>

<sup>3</sup>Jeremy Hermann and Mike Del Balso. *Scaling Machine Learning at Uber with Michelangelo*. <https://eng.uber.com/scaling-michelangelo/>. 2018.



*"Data is the hardest part of ML and the most important piece to get right."*

*Modelers spend most of their time selecting and transforming features at training time and then building the pipelines to deliver those features to production models."*

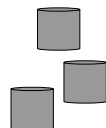
- Uber<sup>4</sup>

---

<sup>4</sup>Jeremy Hermann and Mike Del Balso. *Scaling Machine Learning at Uber with Michelangelo*. <https://eng.uber.com/scaling-michelangelo/>. 2018.

# Disentangle ML Pipelines with a Feature Store

## Raw/Structured Data



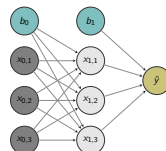
*Feature Engineering*



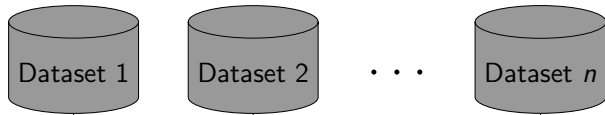
*Training*



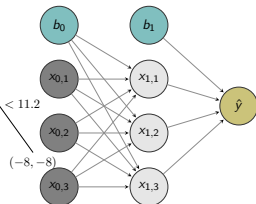
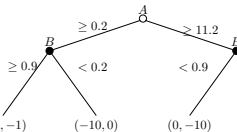
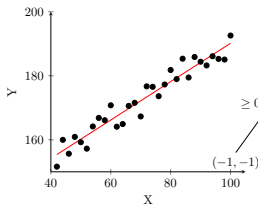
## Models

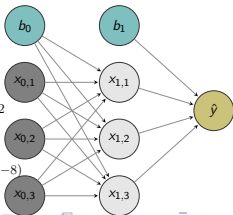
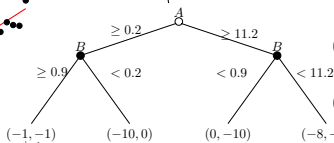
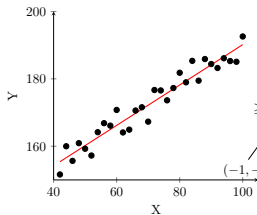
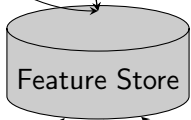
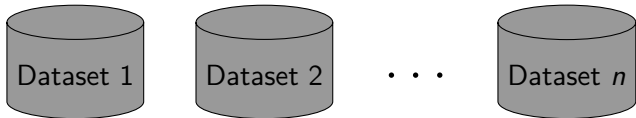


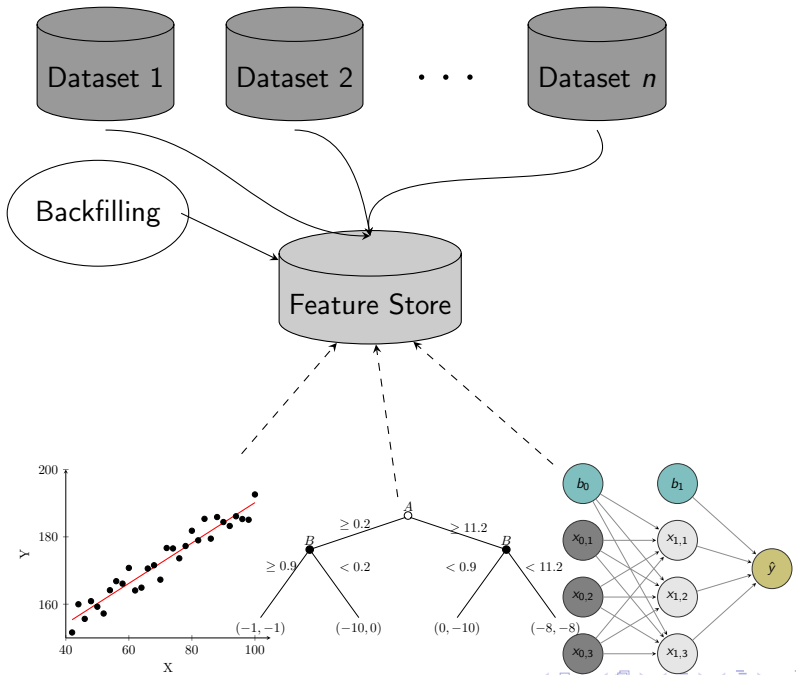
- A feature store is a central vault for storing documented, curated, and access-controlled features.
- The feature store is the interface between data engineering and data model development

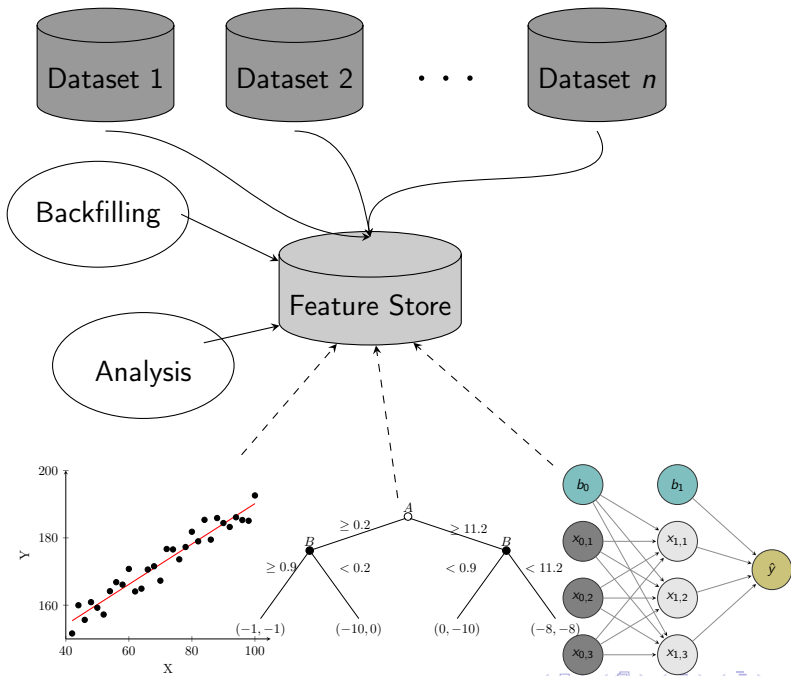


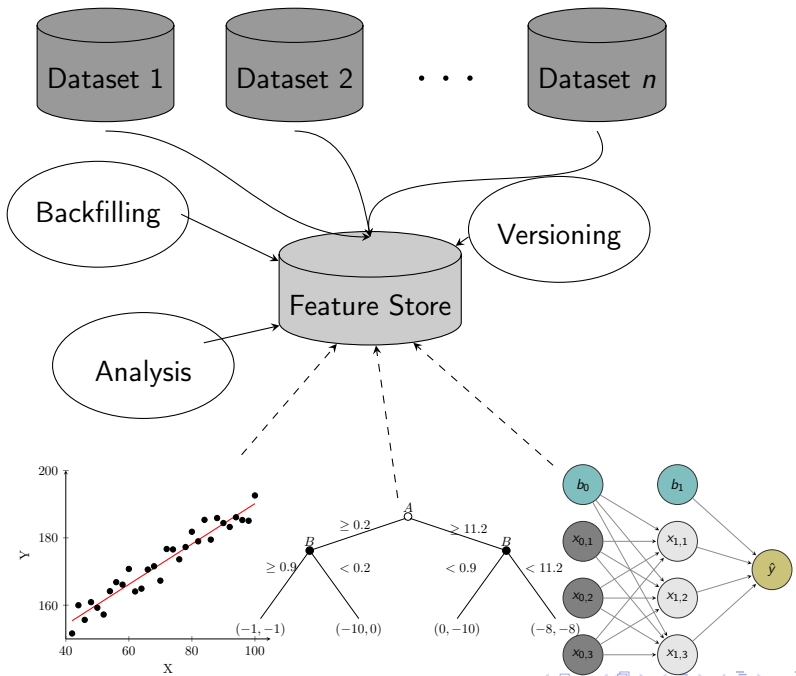
Feature Engineering



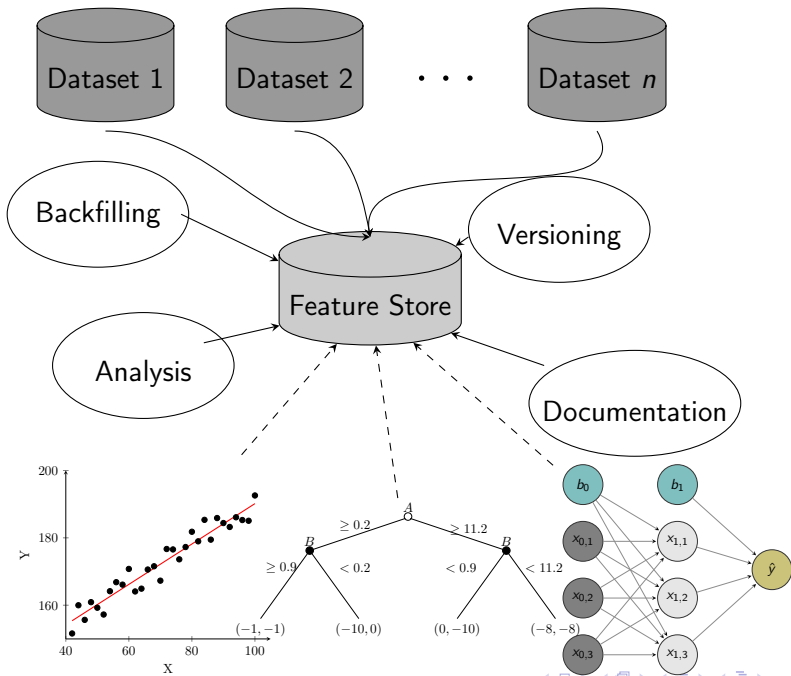












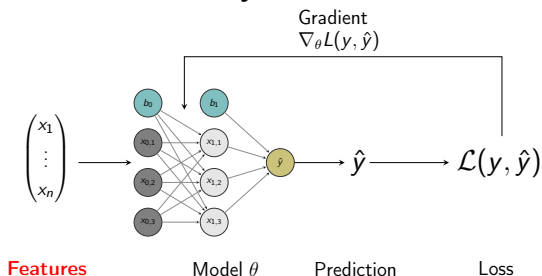
# What is a Feature?

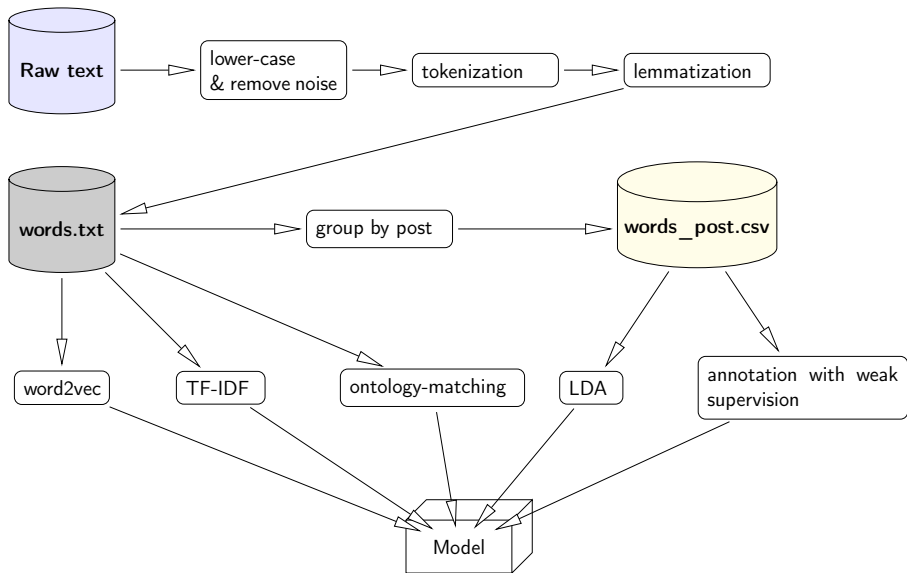
*A feature is a measurable property of some data-sample*

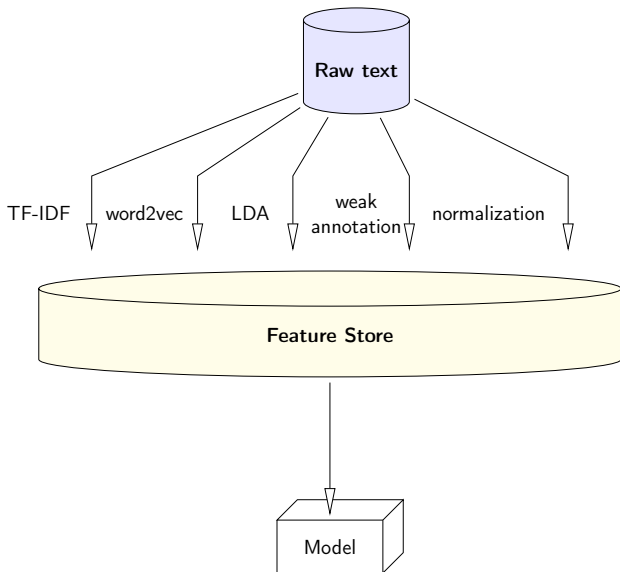
A feature could be..

- An aggregate value (min, max, mean, sum)
- A raw value (a pixel, a word from a piece of text)
- A value from a database table (the age of a customer)
- A derived representation: e.g an embedding or a cluster

Features are the fuel for AI systems:

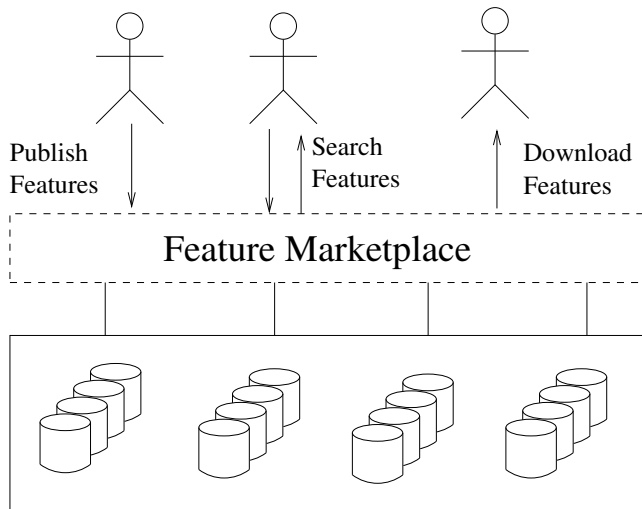






# How to Encourage Feature Reusage?

# Feature Marketplace



# Feature Store API Service

```
from hops import featurestore
features_df =
featurestore.get_features([
    "average_attendance",
    "average_player_age"
])
```

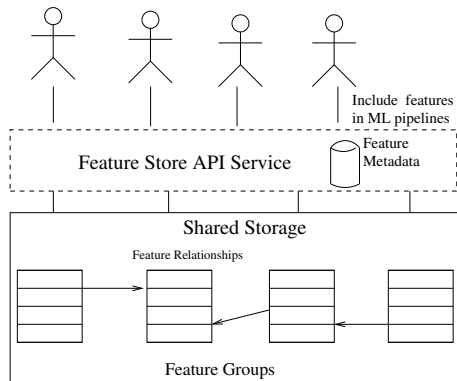


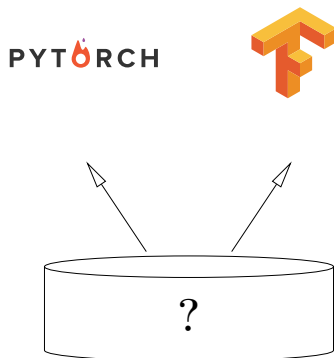
Figure: Feature Store API Service

# How to Store Datasets for Deep Learning?

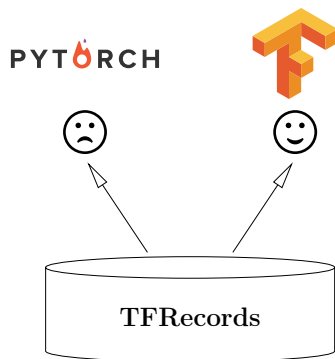
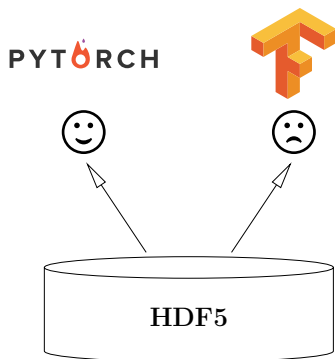


# How to Store Datasets for Deep Learning?

- Should be framework agnostic
- Need to be able to store tensor datasets
- Should support sharding for distributed training
- Advanced features: row-predicate filtering, SQL interface, columnar selection.

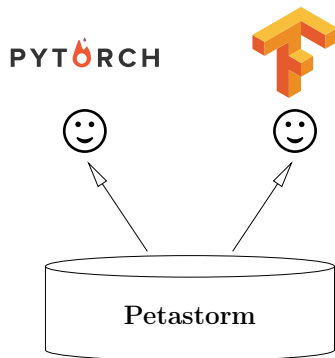


# How to Store Datasets for Deep Learning?



# How to Store Datasets for Deep Learning?

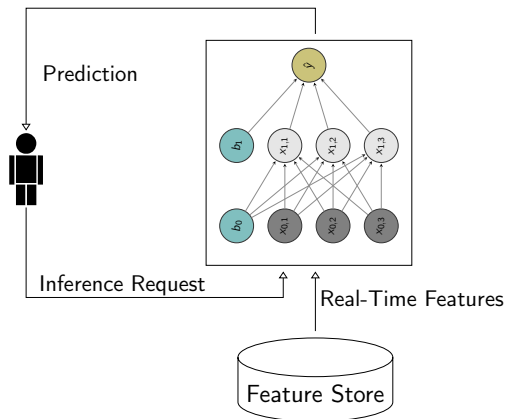
- Petastorm is a dataset format designed for deep learning
- Petastorm stores data as parquet files with extra metadata to handle multi-dimensional tensors
- Petastorm contains readers for the popular machine learning frameworks such as SparkML, Tensorflow, PyTorch



# How to Serve Features for Inference?

# Delivering Features for Training and Serving is Different

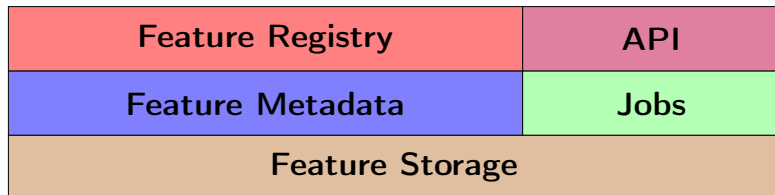
- Serving can require real-time features
- Ideally we want consistency between real-time features and batch features used for training
- Complex engineering problem



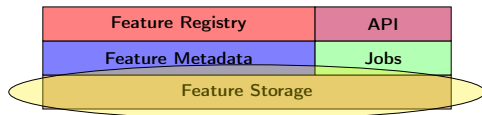
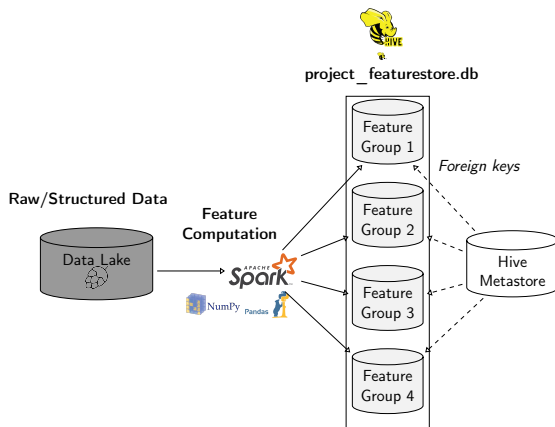
# How to Implement a (batch) Feature Store?

# The Components of a Feature Store

- **The Storage Layer:** For storing feature data in the feature store
- **The Metadata Layer:** For storing feature metadata (versioning, feature analysis, documentation, jobs)
- **The Feature Engineering Jobs:** For computing features
- **The Feature Registry:** A user interface to share and discover features
- **The Feature Store API:** For writing/reading to/from the feature store

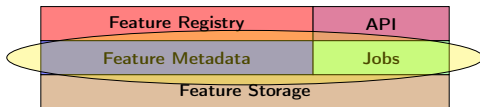
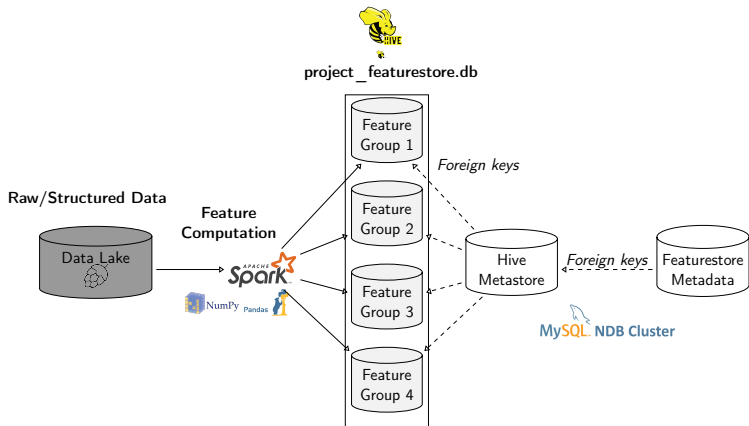


# Feature Storage

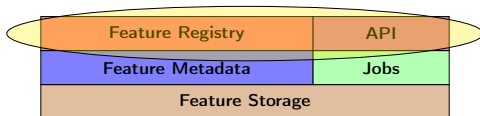
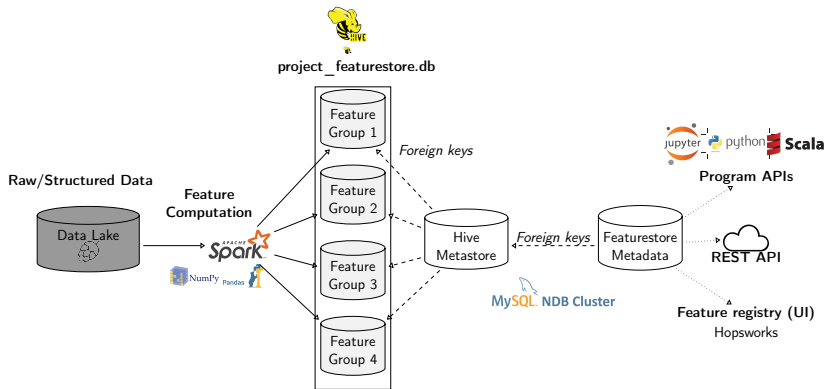




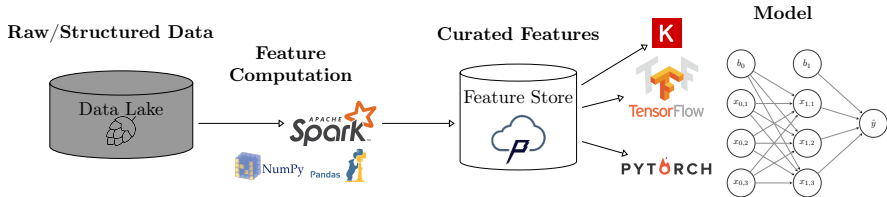
# Feature Metadata



# Feature Registry and API



# Demo-Setting



# Summary

- Machine learning comes with a high technical cost
- Machine learning pipelines needs proper data management
- A **feature store** is a place to store curated and documented features
- The feature store serves as an interface between feature engineering and model development, it can help disentangle complex ML pipelines
- *Hopsworks*<sup>6</sup> provides the world's first open-source feature store



@hopshadoop

[www.hops.io](http://www.hops.io)

@logicalclocks

[www.logicalclocks.com](http://www.logicalclocks.com)

LOGICAL CLOCKS

We are open source:

<https://github.com/logicalclocks/hopsworks>

<https://github.com/hopshadoop/hops>

7

<sup>6</sup>Jim Dowling. *Introducing Hopsworks*. <https://www.logicalclocks.com/introducing-hopsworks/>. 2018.

<sup>7</sup>Thanks to Logical Clocks Team: Jim Dowling, Seif Haridi, Theo Kakantousis, Fabio Buso, Gautier Berthou, Ermias Gebremeskel, Mahmoud Ismail, Salman Niazi, Antonios Kouzoupis, Robin Andersson, and Alex Ormensean

# References

- Hopsworks' feature store<sup>8</sup> (the only open-source one!)
- Uber's feature store<sup>9</sup>
- Airbnb's feature store<sup>10</sup>
- Comcast's feature store<sup>11</sup>
- GO-JEK's feature store<sup>12</sup>
- HopsML<sup>13</sup>
- Hopsworks<sup>14</sup>

---

<sup>8</sup>Kim Hammar and Jim Dowling. *Feature Store: the missing data layer in ML pipelines?* <https://www.logicalclocks.com/feature-store/>. 2018.

<sup>9</sup>Li Erran Li et al. "Scaling Machine Learning as a Service". In: *Proceedings of The 3rd International Conference on Predictive Applications and APIs*. Ed. by Claire Hardgrove et al. Vol. 67. Proceedings of Machine Learning Research. Microsoft NERD, Boston, USA: PMLR, 2017, pp. 14–29. URL: <http://proceedings.mlr.press/v67/li17a.html>.

<sup>10</sup>Nikhil Simha and Varant Zanoian. *Zipline: Airbnb's Machine Learning Data Management Platform*. <https://databricks.com/session/zipline-airbnbs-machine-learning-data-management-platform>. 2018.

<sup>11</sup>Nabeel Sarwar. *Operationalizing Machine Learning—Managing Provenance from Raw Data to Predictions*. <https://databricks.com/session/operationalizing-machine-learning-managing-provenance-from-raw-data-to-predictions>. 2018.

<sup>12</sup>Willem Pienaar. *Building a Feature Platform to Scale Machine Learning | DataEngConf BCN '18*. <https://www.youtube.com/watch?v=0iCXY6VnpCc>. 2018.

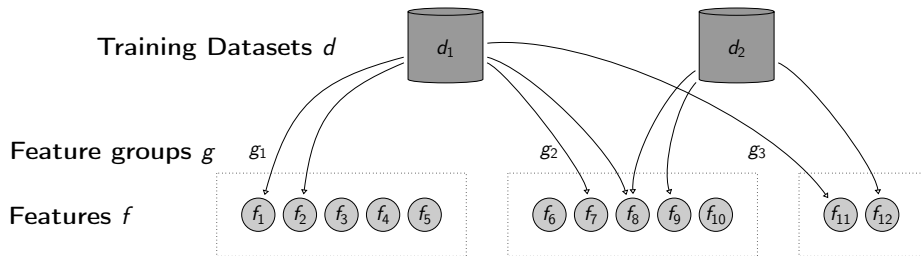
<sup>13</sup>Logical Clocks AB. *HopsML: Python-First ML Pipelines*. <https://hops.readthedocs.io/en/latest/hopsml/hopsML.html>. 2018.

<sup>14</sup>Jim Dowling. *Introducing Hopsworks*. <https://www.logicalclocks.com/introducing-hopsworks/>. 2018.

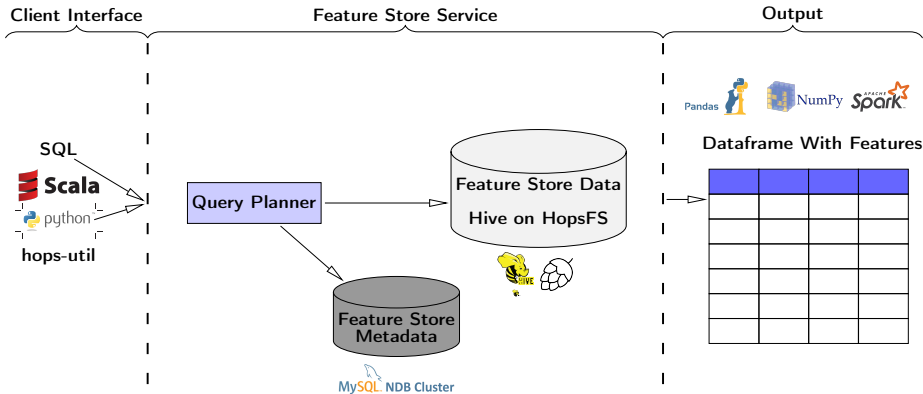
# Backup Slides

# Modeling Data in the Feature Store

- A **feature group** is a logical grouping of **features**
  - Typically from the same input dataset and computed with the same job
- A **training dataset** is a set of features suitable for a prediction task
  - Features in a training dataset are often from several feature groups
  - E.g features on customers, features on user activities, etc.



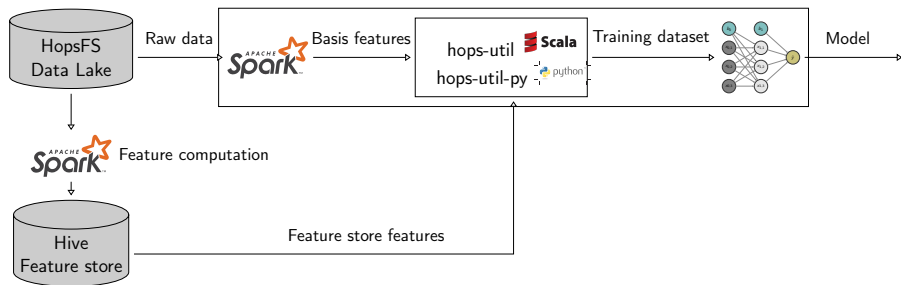
# Hopworks Feature Store API Service





# Training Pipeline in HopsML

- 1 Create job/notebook to **compute features** and publish to the feature store
- 2 Create job/notebook to read features/labels and **save to a training dataset**
- 3 **Read the training dataset into your model** for training



## Reading from the Feature Store:

---

```
from hops import featurestore
features_df = featurestore.get_features([
    "average_attendance",
    "average_player_age"
])
```

---

## Writing to the Feature Store:

---

```
from hops import featurestore
raw_data = spark.read.parquet(filename)
pol_features = raw_data.map(lambda x: x^2)
featurestore.insert_into_featuregroup(pol_features, "pol_featuregroup")
```

---