

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

*FOSDEM 2019 Brussels | HPC, Big Data and Data Science DevRoom*

Kim Hammar

*kim@logicalclocks.com*

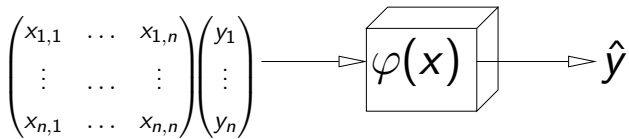
@KimHammar1 

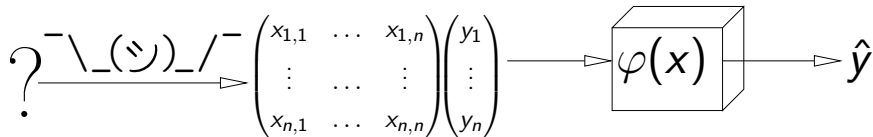
February 1, 2019



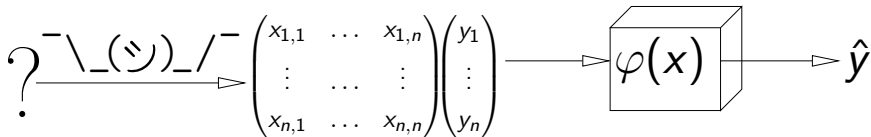
LOGICAL CLOCKS

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





<sup>2</sup>Jeremy Hermann and Mike Del Balso. *Scaling Machine Learning at Uber with Michelangelo*.  
<https://eng.uber.com/scaling-michelangelo/>. 2018.

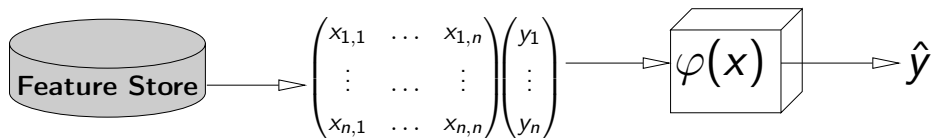


*"Data is the hardest part of ML and the most important piece to get right."*

*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>2</sup>

<sup>2</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."*

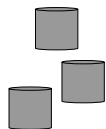
*Modelers spend most of their time selecting and transforming features at training time and then building the pipelines to deliver those features to production models."*

- 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.

# Solution: Disentangle ML Pipelines with a Feature Store

Raw/Structured Data



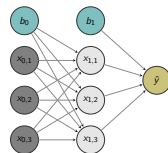
*Feature Engineering*



*Training*



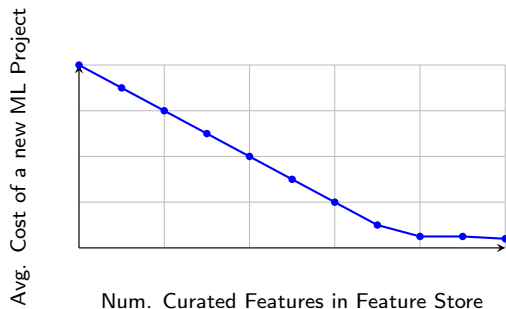
Models



- 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

## The feature store enables:

- Reusability of features between models and teams
- Automatic backfilling of features
- Automatic feature documentation and analysis
- Feature versioning
- Standardized access of features between training and serving
- Feature discovery



# Reusing Features Without a Feature Store is Complex



**Siloed Feature Sets**  
Without a feature store it is typical to have feature sets stored in isolation from each other.

$$\begin{pmatrix} w_{1,1} & \dots & w_{1,n} \\ \vdots & \ddots & \vdots \\ w_{n,1} & \dots & w_{n,n} \end{pmatrix}$$

Features

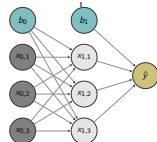
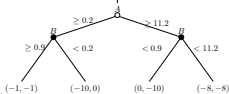
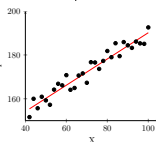
$$\begin{pmatrix} w_{1,1} & \dots & w_{1,n} \\ \vdots & \ddots & \vdots \\ w_{n,1} & \dots & w_{n,n} \end{pmatrix}$$

Features

$$\begin{pmatrix} w_{1,1} & \dots & w_{1,n} \\ \vdots & \ddots & \vdots \\ w_{n,1} & \dots & w_{n,n} \end{pmatrix}$$

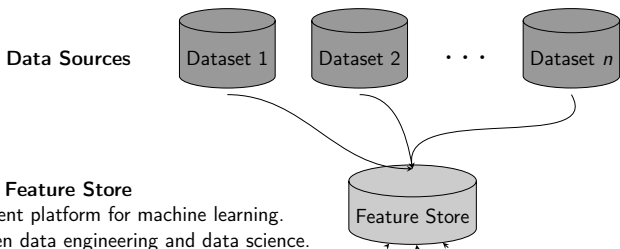
Features

**Models**  
Models are trained using sets of features. Without a feature store each model typically defines its own feature definitions, without feature sharing across models.



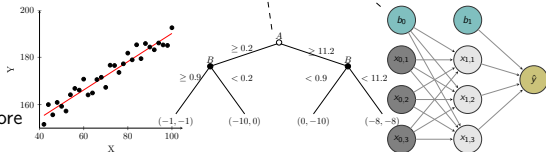


# Reusing Features With a Feature Store is Simple

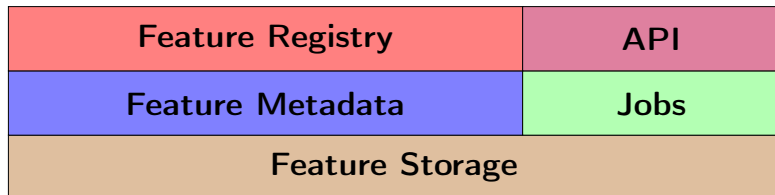


**Models**

Models are trained using sets of features.  
The features are fetched from the feature store  
and can overlap between models.



- **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



## Reading from the Feature Store:

---

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

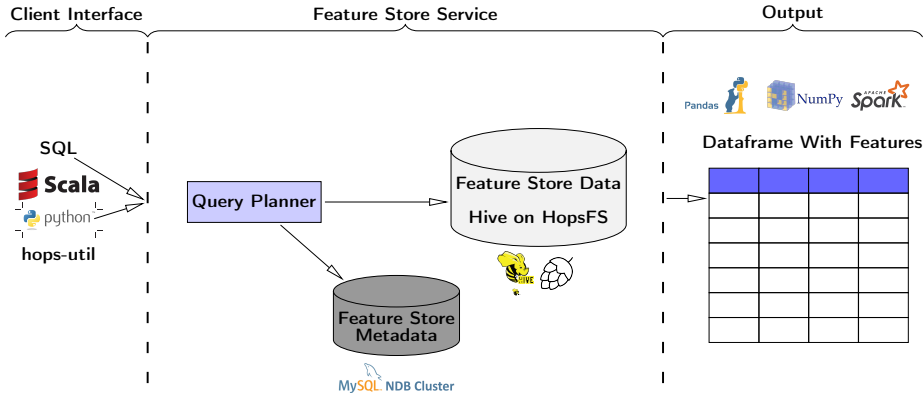
---

## 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")
```

---



- 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>4</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>

<sup>4</sup>Jim Dowling. *Introducing Hopsworks*. <https://www.logicalclocks.com/introducing-hopsworks/>. 2018.

<sup>5</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 Ormenisan

- Hopsworks' feature store<sup>6</sup> (the only open-source one!)
- Uber's feature store<sup>7</sup>
- Airbnb's feature store<sup>8</sup>
- Comcast's feature store<sup>9</sup>
- GO-JEK's feature store<sup>10</sup>
- HopsML<sup>11</sup>
- Hopsworks<sup>12</sup>

---

<sup>6</sup>Kim Hammar and Jim Dowling. *Feature Store: the missing data layer in ML pipelines?* <https://www.logicalclocks.com/feature-store/>. 2018.

<sup>7</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>8</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>9</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>10</sup>Willem Pienaar. *Building a Feature Platform to Scale Machine Learning | DataEngConf BCN '18*. <https://www.youtube.com/watch?v=0iCXY6VnpCc>. 2018.

<sup>11</sup>Logical Clocks AB. *HopsML: Python-First ML Pipelines*. <https://hops.readthedocs.io/en/latest/hopsml/hopsML.html>. 2018.

<sup>12</sup>Jim Dowling. *Introducing Hopsworks*. <https://www.logicalclocks.com/introducing-hopsworks/>. 2018.