

# Hopsworks Feature Store after 4 years: Lessons learned and what's next

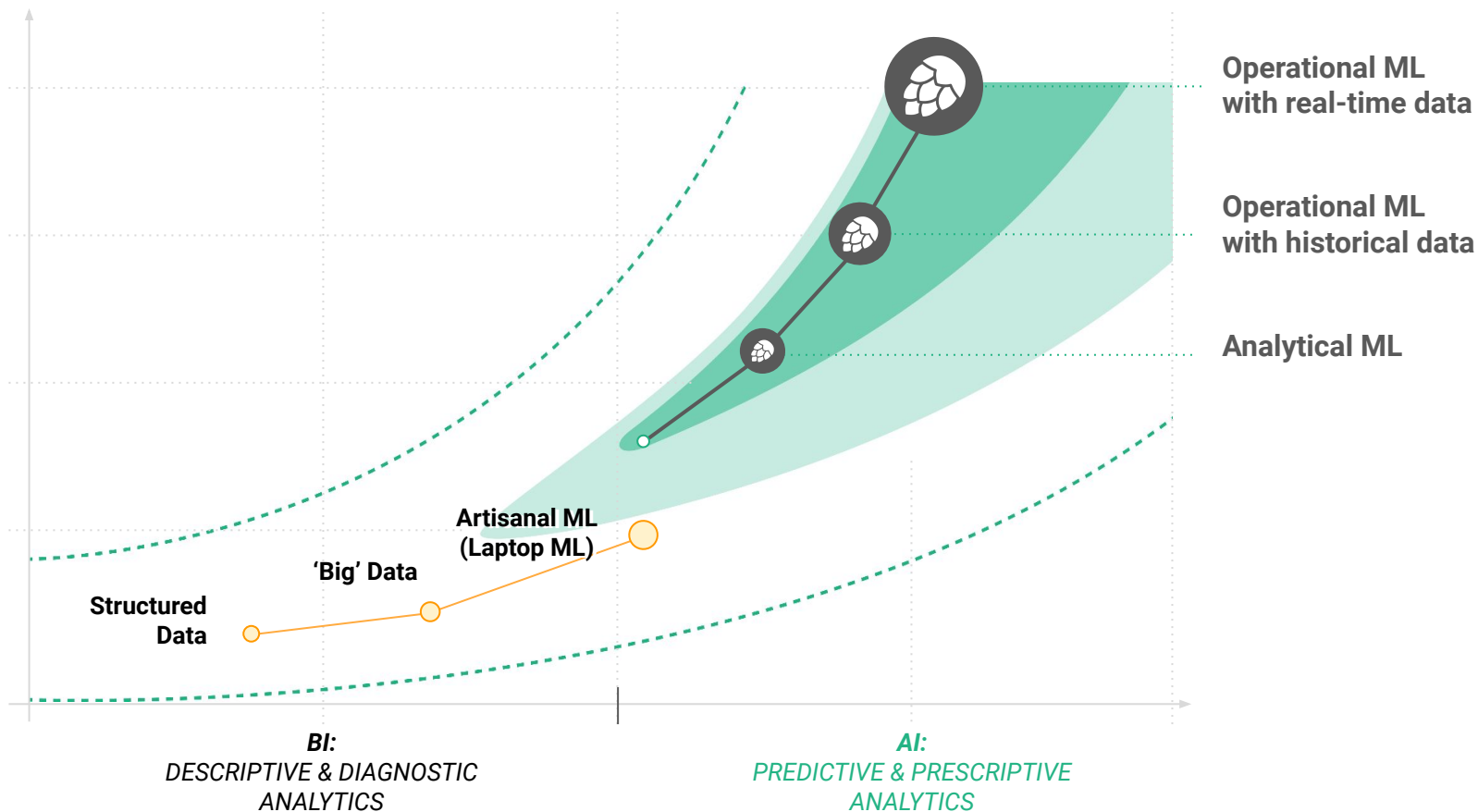
Fabio Buso, VP of Engineering, Hopsworks

Moritz Meister, Head of Feature Store Engineering, Hopsworks

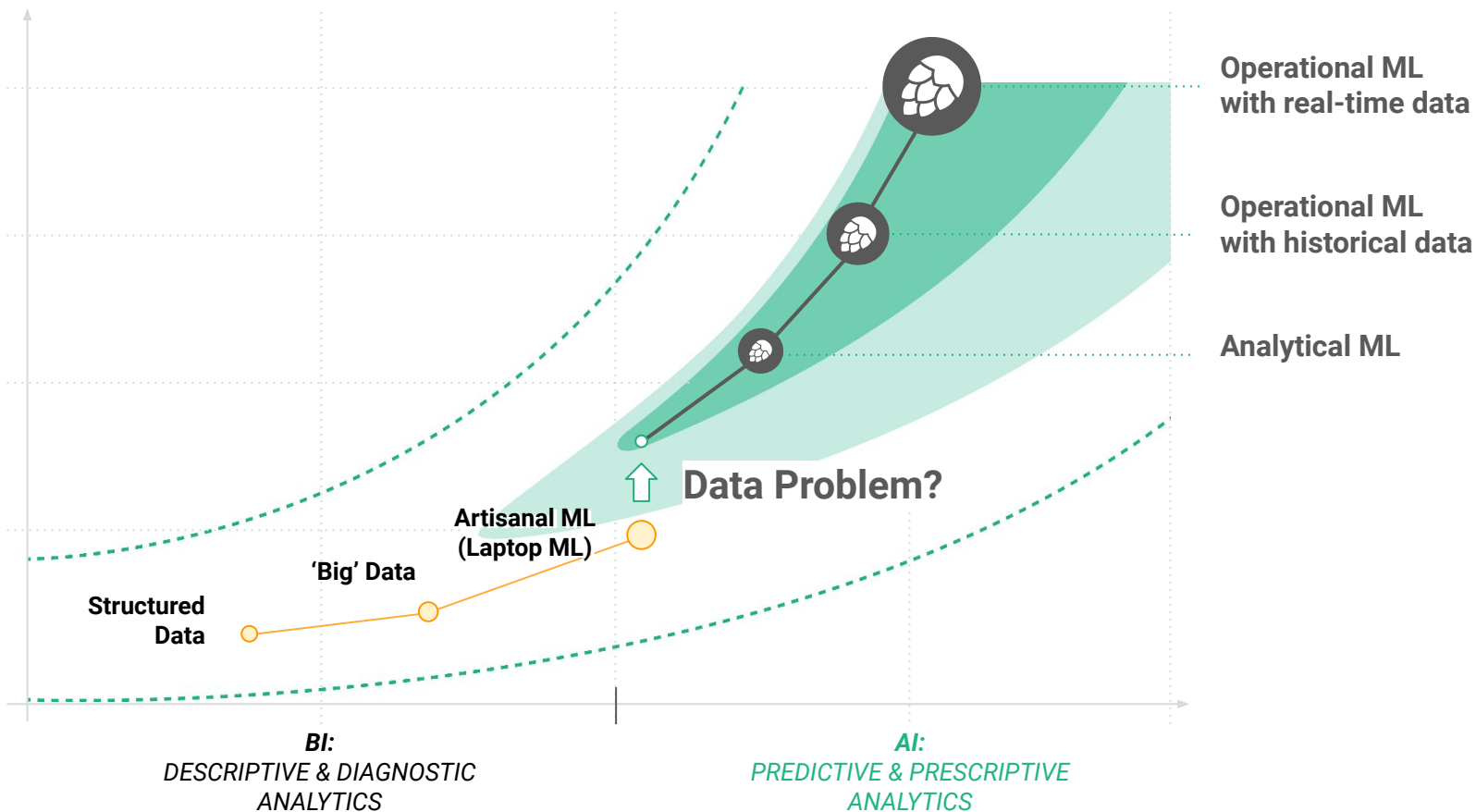
# Feature Stores and the Data Problem



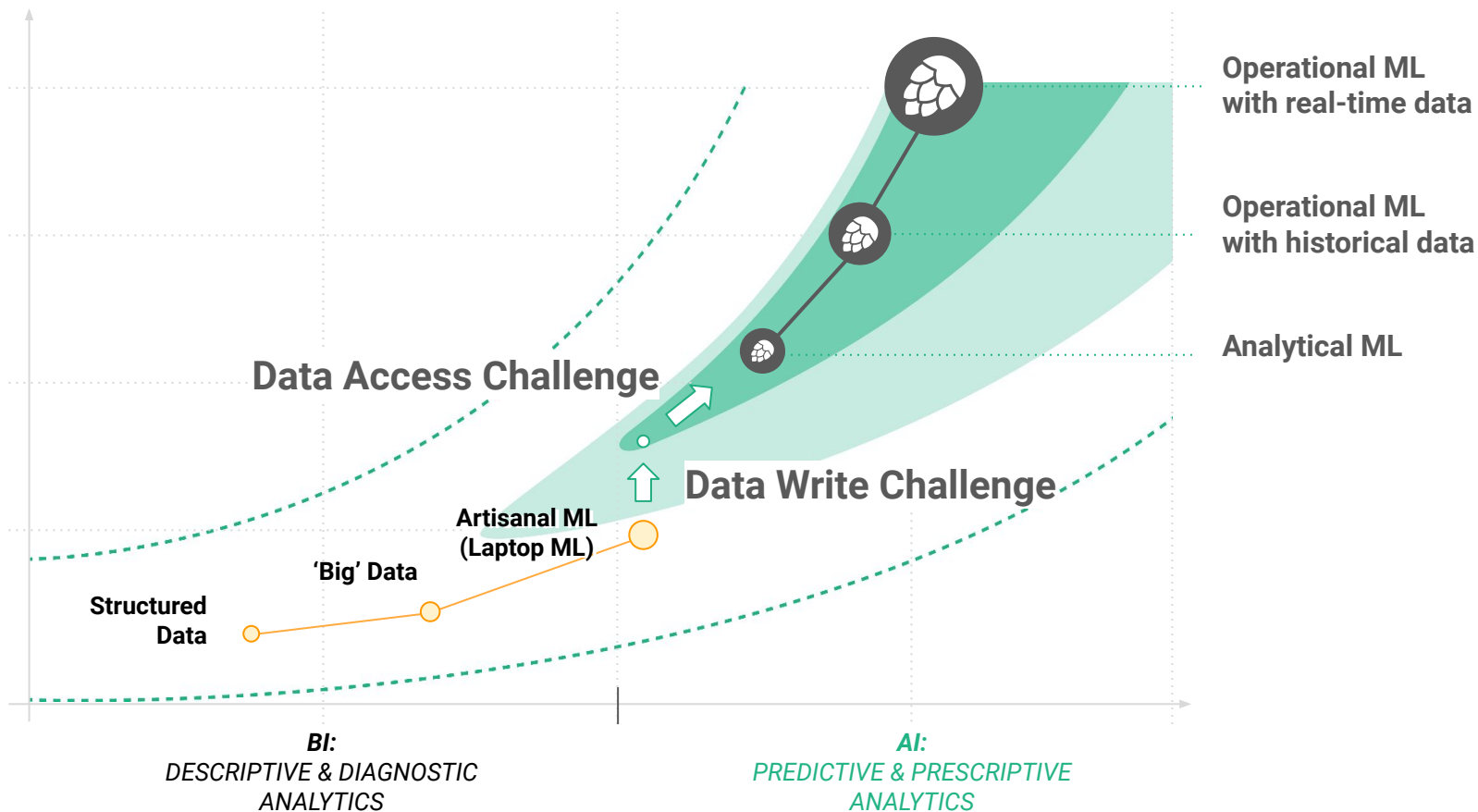
# Business Value



# Business Value



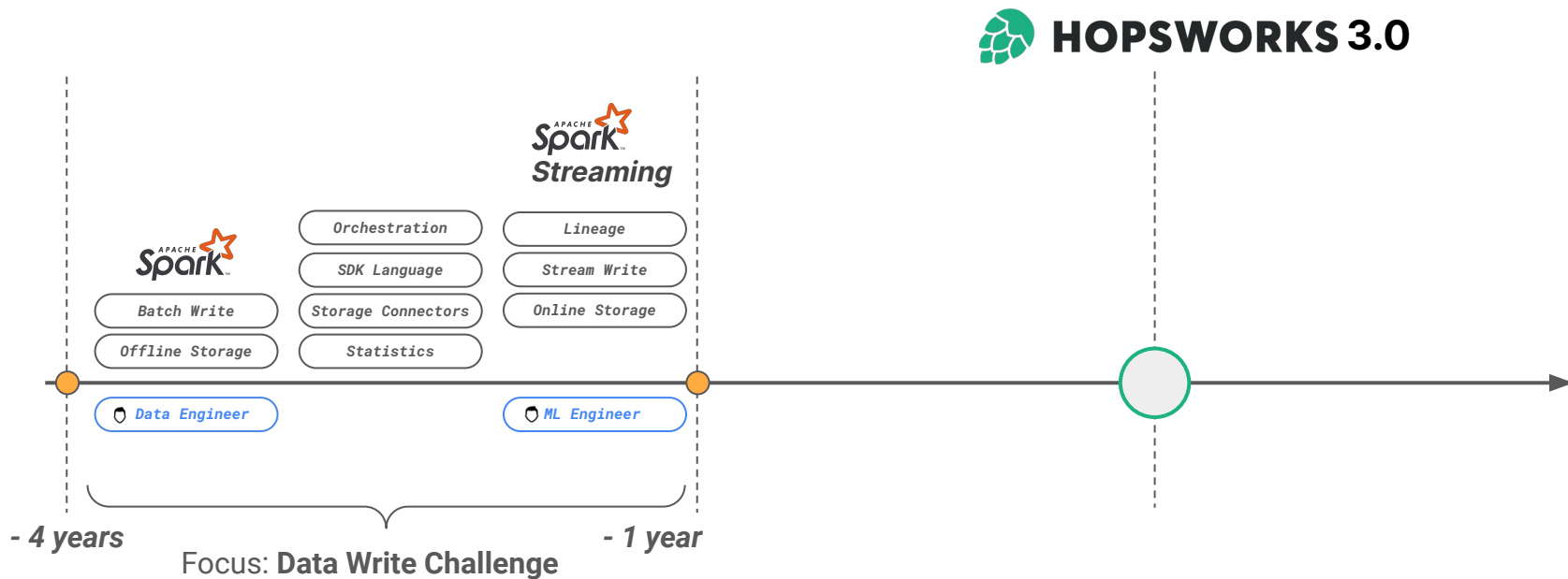
Business Value





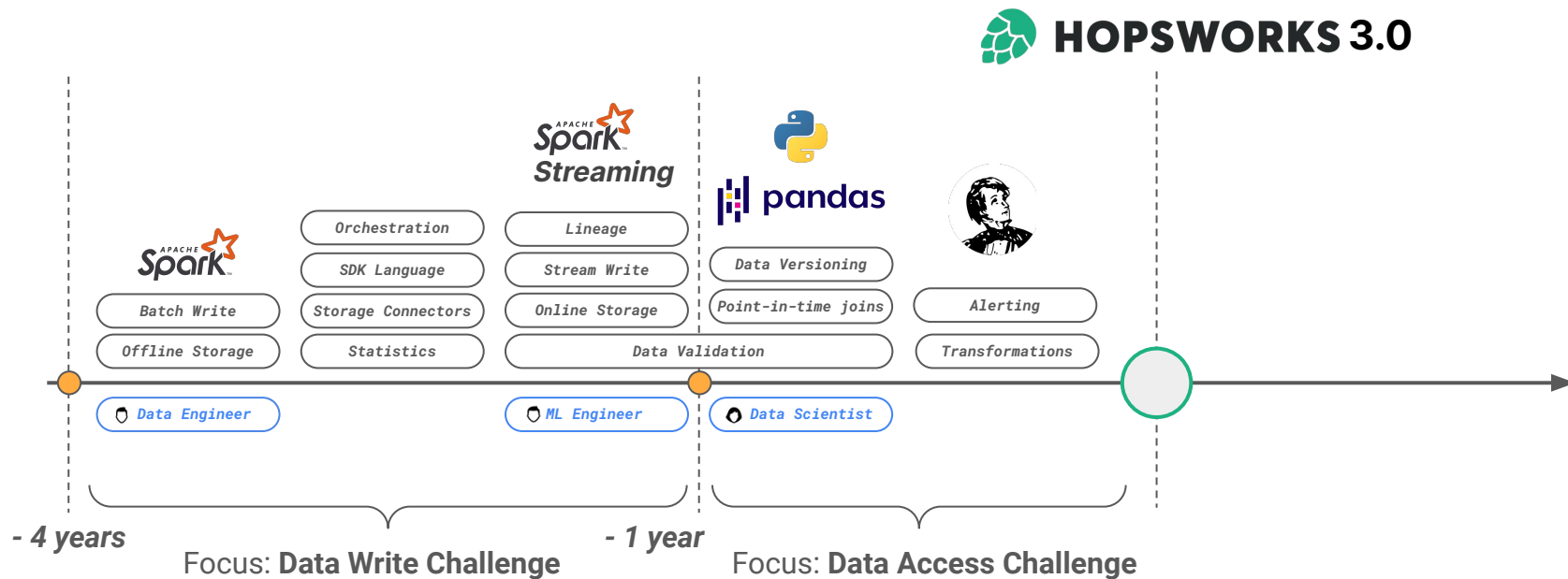
# A (short) history of Feature Store capabilities

From batch to prediction service



# A (short) history of Feature Store capabilities

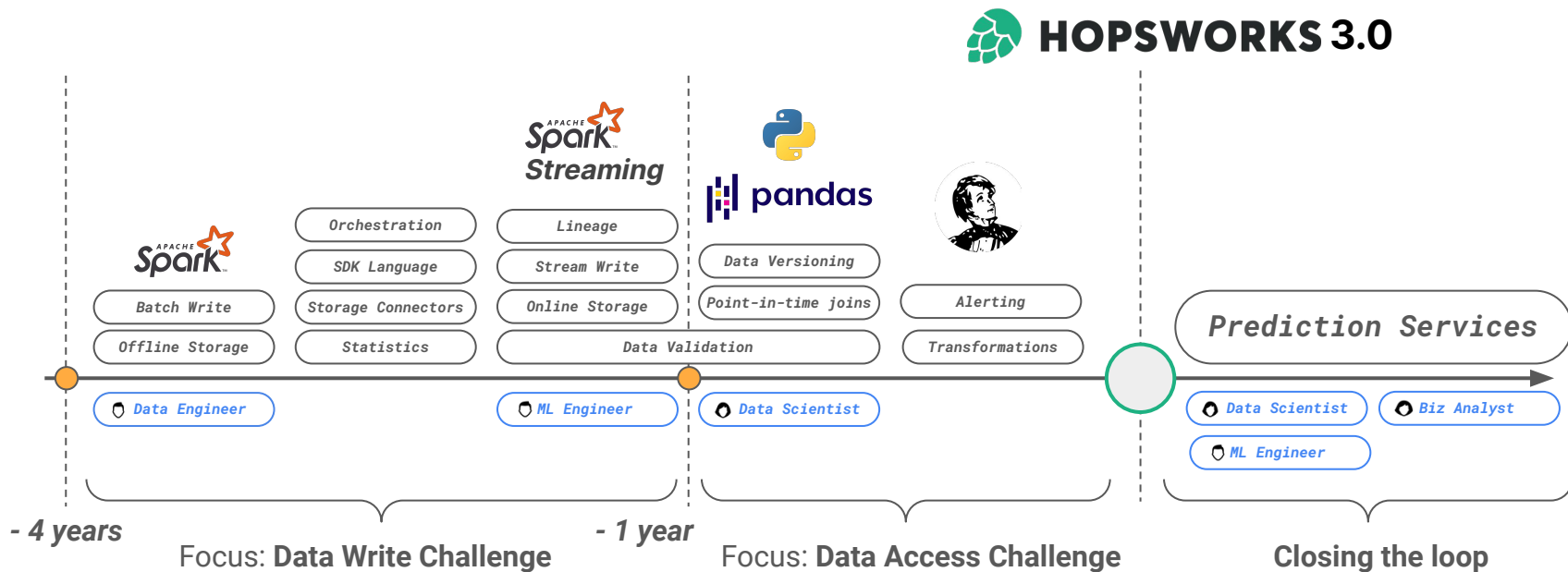
From batch to prediction service





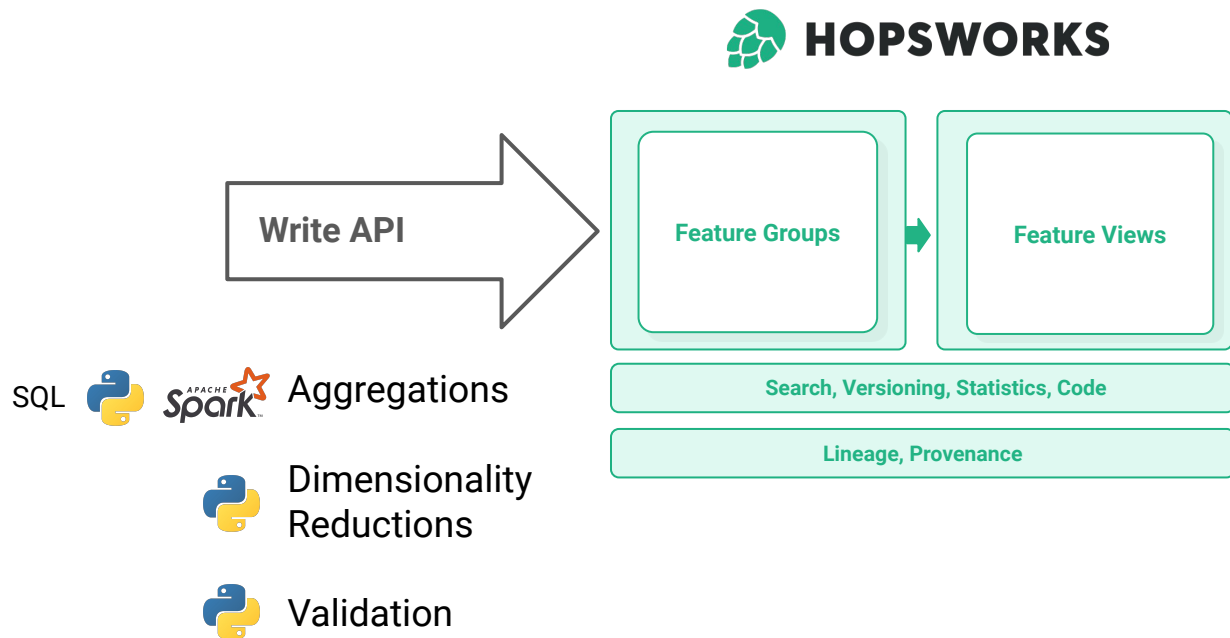
# A (short) history of Feature Store capabilities

From batch to prediction service



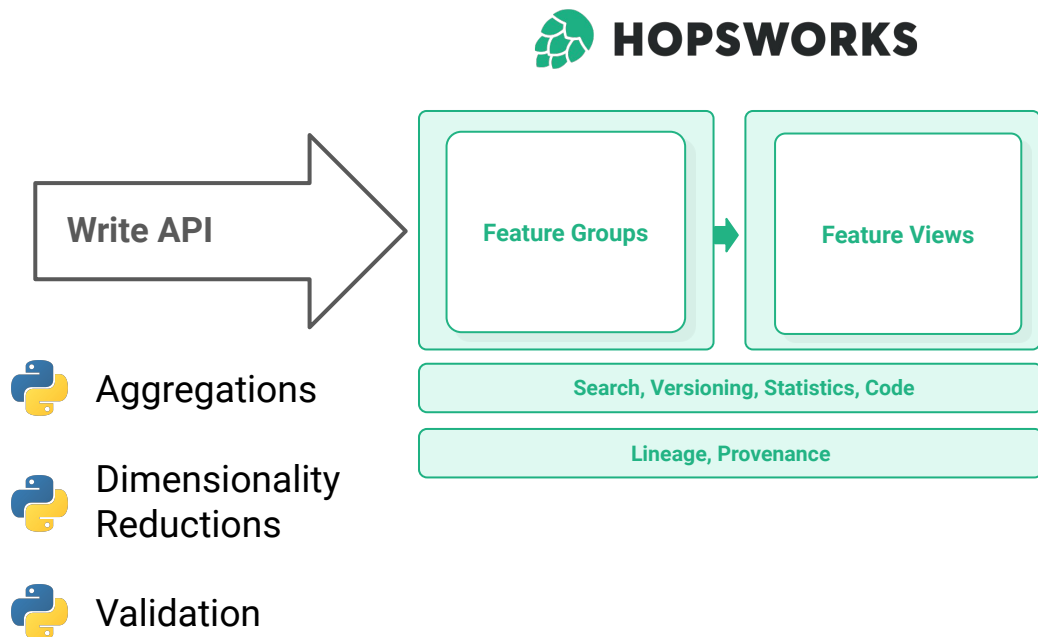
# Write to Feature Groups, read from Feature Views

Effort 1. Making Feature Pipelines accessible to Data Scientists



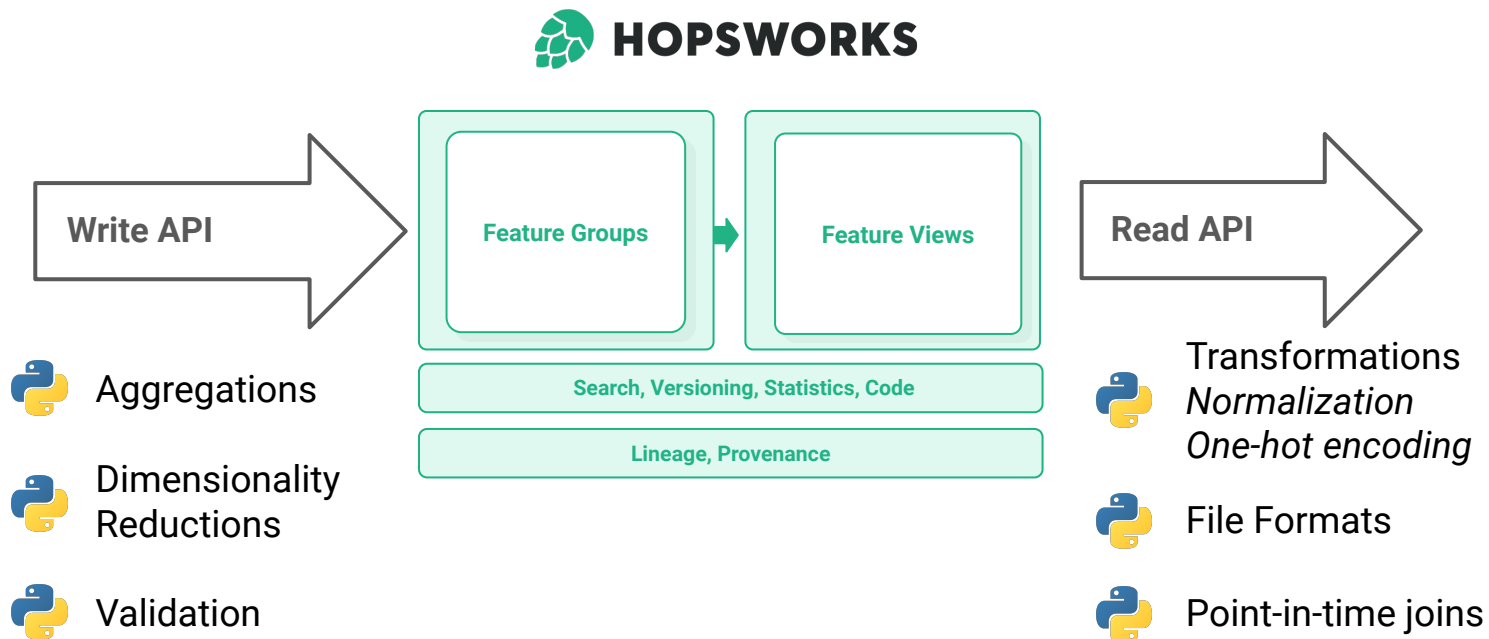
# Write to Feature Groups, read from Feature Views

Effort 1. Making Feature Pipelines accessible to Data Scientists



# Write to Feature Groups, read from Feature Views

Effort 2. Making Data accessible to Data Scientists





# Feature Groups and the Dataframe API

*Data Write Challenge*



## Why DSLs don't make the cut

### Advantages of Dataframes

1. **Flexibility** - *DSLs are use case specific*
2. **User Experience** - *No additional learning curve and no lock in*
3. **Bring your own pipeline** - *Keep existing libraries and pipelines*

# Why Spark is not suitable for Data Scientists

1. **Resource Estimation**
  - *Distributed environment resources don't map 1:1 to local environments*
2. **Debugging**
  - *Debugging in distributed systems quickly becomes a complex task*
3. **Models are rarely trained on Spark Dataframes**
  - *Modeling Frameworks are often not distributed tools*
4. **Wrapping SQL in Python functions is not pythonic**
  - *At serving time there is often no Spark context available*

Feature engineering

Initialize Feature Group metadata

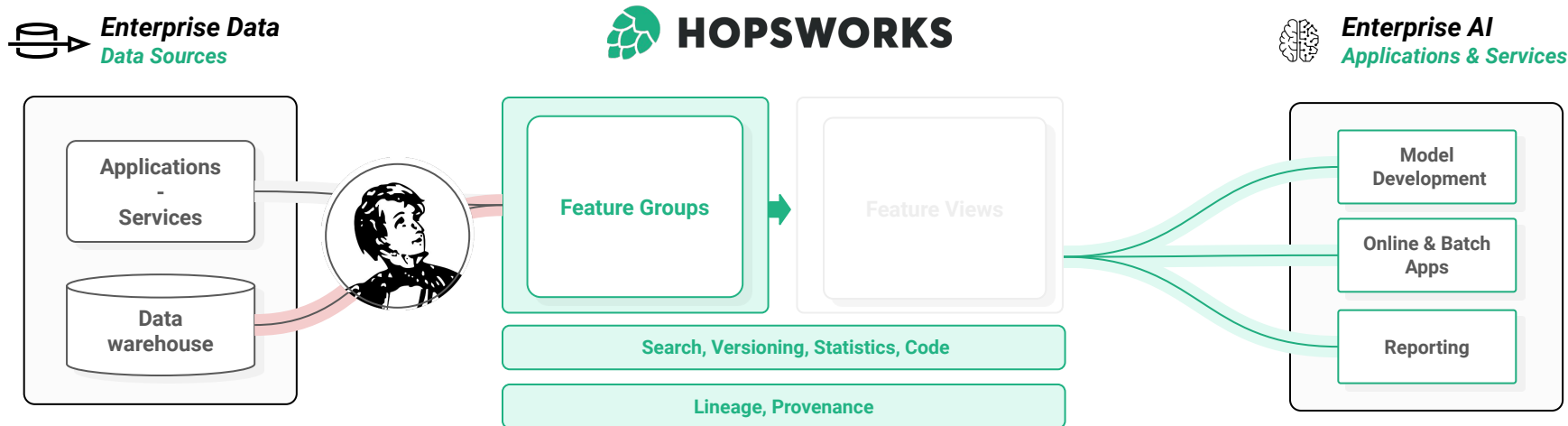
Write feature data

```
1 window_len = "4h"
2 cc_group = trans_df[["cc_num", "amount", "datetime"]] \
3     .groupby("cc_num").rolling(window_len, on="datetime")
4
5 # Compute aggregations
6 df_4h_count = pd.DataFrame(cc_group.mean())
7 df_4h_count.columns = ["trans_freq", "datetime"]
8 df_4h_count = df_4h_count.reset_index(level=["cc_num"])
9 df_4h_count = df_4h_count.drop(columns=["cc_num", "datetime"])
10 df_4h_count = df_4h_count.sort_index()
11 window_aggs_df = window_aggs_df.merge(df_4h_count, left_index=True, right_index=True)
12
13 # Initialize metadata of feature group
14 window_aggs_fg = fs.get_or_create_feature_group(
15     name="transactions_4h_aggs_fraud_batch_fg",
16     version=1,
17     description="Aggregate transaction data over 4h windows.",
18     primary_key=["cc_num"],
19     event_time="datetime"
20 )
21
22 # Write Dataframe to the feature group
23 window_aggs_fg.insert(window_aggs_df)
```



# Where does Great Expectations fit in?

Validate data before it is made available to data scientists



# Where does Great Expectations fit

Validate data before it is made available to data scientists

*Pandas Dataframe*



*Expectation*


```
1 ExpectationConfiguration(
2     expectation_type="expect_column_min_to_be_between",
3     kwargs={
4         "column": "age_at_transaction",
5         "min_value": 18,
6         "max_value": 130
7     }
8 )
```


*Result*


```
1 validation_result = {
2     "success": false,
3     "result": {
4         "observed_value": 17.738671,
5         "element_count": 106020,
6         "missing_count": null,
7         "missing_percent": null
8     },
9     ...
10 }
```



# Where does Great Expectations fit



Validate data before it is made available to data scientists


 **HOPSWORKS**


dataval 


 Search for


 + 


 


 [← Back](#)

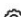
 **Overview**


 Features


 Provenance


 **Expectations**

 Tags

 Alerts









 API

 Data preview

 Feature statistics

## Expectations ⓘ

7 expectations — configured in **STRICT** Gatekeeper mode (data is ingested if and only if all individual validations are successful) [Edit Expectation Suite](#)

Validation Reports	Validation Results				
validation date ↕	success percent ↕	evaluated expectations ↕	successful expectations ↕	unsuccessful expectations ↕	result ↕
25-08-2022 22:25	57.14285714285714%	7	4	3	rejected  
25-08-2022 22:20	57.14285714285714%	7	4	3	ingested  
25-08-2022 22:17	57.14285714285714%	7	4	3	ingested  
25-08-2022 22:09	57.14285714285714%	7	4	3	ingested  

# Python, Pandas and Feature Views

*Data Access Problem*



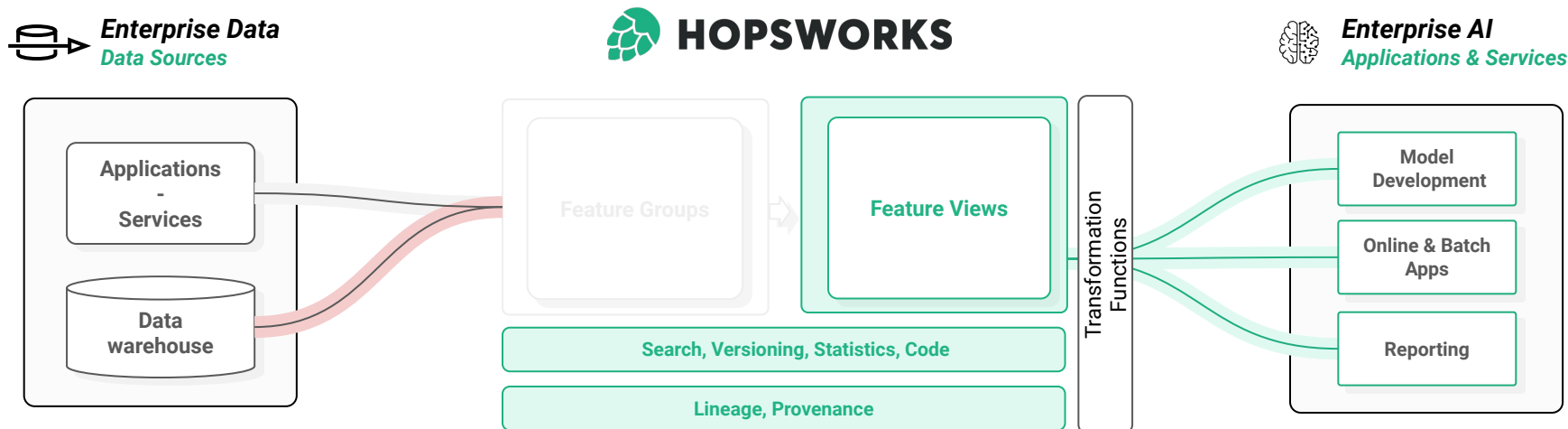
## Feature Reusability requires Views - Feature Views

Why Data Scientists want Python-centric APIs

1. **Models are trained with Python libraries** - *DSLs are use case specific*
2. **Training frameworks might require file formats** - Scikit-learn & Pandas, Tensorflow & Tfrecored, Pytorch & Numpy
3. **Point-in-time joins are hard** - *Writing PIT joins in SQL, Spark or Python is error prone*
4. **Transformations happen at feature retrieval time** - *Reusing a feature often means limiting the set of entities a model is trained on, which again leads to different statistics*
5. **Spark not available at feature retrieval time** - *At serving time there is often no Spark context available*

# Feature Views for training data, batch scoring and low latency serving

The API for Data Scientists



# Feature Reusability with Feature Views

Pandas like join making PIT joins transparent

Get Feature Group metadata

Select features

Create Feature View along with transformations

```
1 # Select feature groups
2 trans_fg = fs.get_feature_group("transactions_fraud_batch_fg",
3     version=1)
4 window_aggs_fg = fs.get_feature_group("transactions_4h_aggs_fraud_batch_fg",
5     version=1)
6
7 # Join the feature groups and select the features
8 ds_query = trans_fg.select(["fraud_label", "category", "amount",
9     "age_at_transaction", "days_until_card_expires", "loc_delta"]) \
10     .join(window_aggs_fg.select_except(["cc_num"]))
11
12 feature_view = fs.create_feature_view(
13     name="transactions_view",
14     query=ds_query,
15     labels=["fraud_label"],
16     transformation_functions={
17         "category": label_encoder,
18         "amount": min_max_scaler,
19         "days_until_card_expires": min_max_scaler,
20     }
21 )
```

# Reading from the Feature View

One API to rule them all

Training Data as files

Training Data in-memory  
as Pandas Dataframe

New batches of data  
for scoring

Real-time feature lookup

```

1  # Create training data in a specific file format
2  td_version, td_job = feature_view.create_train_validation_test_split(
3      description = 'transactions fraud batch training dataset',
4      data_format = 'csv',
5      validation_size = 0.2,
6      test_size = 0.1
7  )
8
9  # Retrieve previously materialised training data
10 X_train, y_train, X_val, y_val, X_test, y_test = feature_view.get_train_validation_test_split(1)
11
12 # Generate batches of data from a certain time range
13 start_time = int(float(datetime.strptime("2022-01-03 00:00:01",
14     date_format).timestamp()) * 1000)
15 end_time = int(float(datetime.strptime("2022-03-31 23:59:59",
16     date_format).timestamp()) * 1000)
17 march_transactions = feature_view.get_batch_data(
18     start_time = start_time, end_time = end_time)
19
20 # Initialize with training data version for transformations and get single feature vectors
21 feature_view.init_serving(1)
22 feature_array = feature_view.get_feature_vector(
23     {"primary_key": 1})
  
```

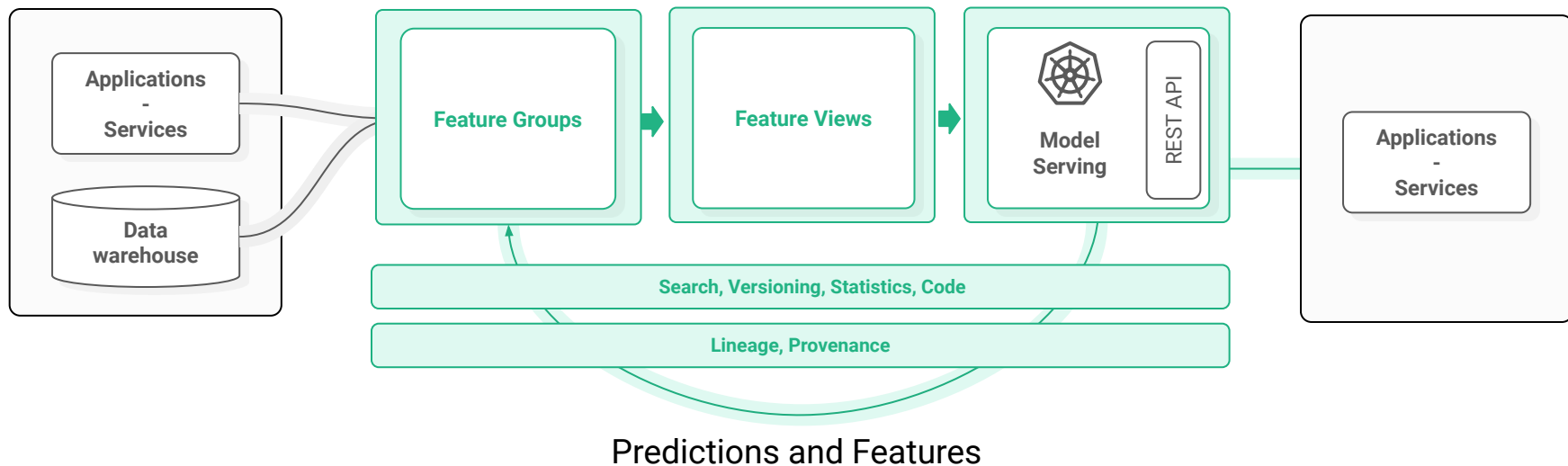


# Endboss: Prediction Service

Integrating Data APIs with Model APIs

# Building Prediction Services with KServe Integration

Closing the loop



# Tying it all together

## The predictor interface

Establish connection to Feature Store

Initialize Feature View

Loading the model from a model registry

Implementing the predict interface using the Feature View

```

1  import os
2  import numpy as np
3  import hsfs
4  import joblib
5
6  class Predict(object):
7
8      def __init__(self):
9          """ Initializes the serving state, reads a trained model"""
10         # get feature store handle
11         fs_conn = hsfs.connection()
12         self.fs = fs_conn.get_feature_store()
13
14         # get feature views
15         self.fv = self.fs.get_feature_view("transactions_view", 1)
16
17         # initialise serving
18         self.fv.init_serving(1)
19
20         # load the trained model
21         self.model = joblib.load(
22             os.environ["ARTIFACT_FILES_PATH"] + "/fraud.model.pkl")
23         print("Initialization Complete")
24
25     def predict(self, inputs):
26         """ Serves a prediction request using a trained model"""
27         # Numpy Arrays are not JSON serializable
28         return self.model.predict(
29             np.asarray(
30                 self.fv.get_feature_vector({"cc_num": inputs[0]})
31                 .reshape(1, -1)
32                 .tolist()
  
```

**app.hopsworks.ai**

Try it out! Serverless and free!