

The logo for CONF42, featuring the text "CONF42" in white on a green rectangular background. The letter "O" is stylized with a leaf-like shape inside it.

CONF42

CONF42 Python, March 2023

The mythical machine learning pipeline - use feature/training/inference pipelines

Jim Dowling (@jim_dowling) - CEO at Hopsworks



HOPSWORKS



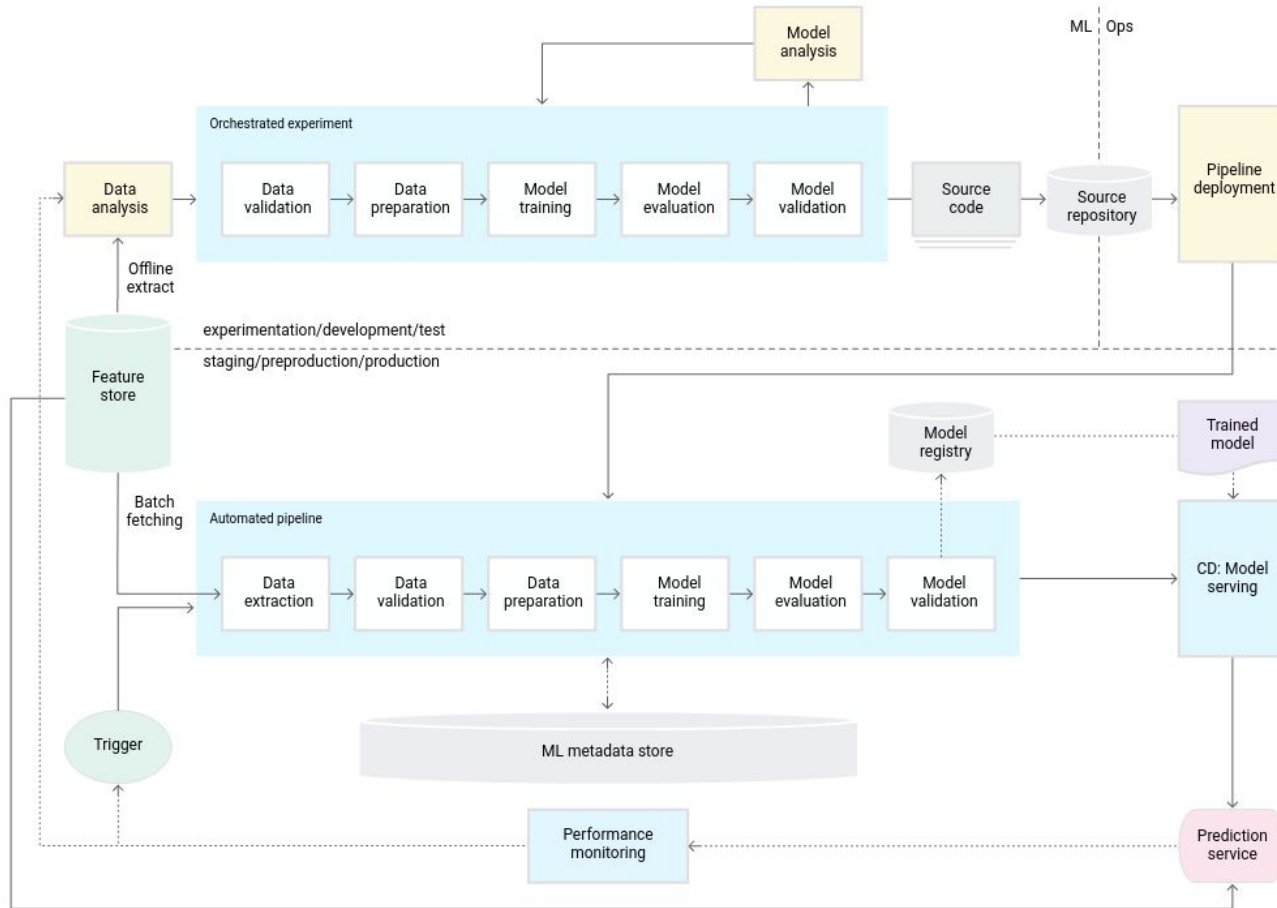
MLOps Today

00

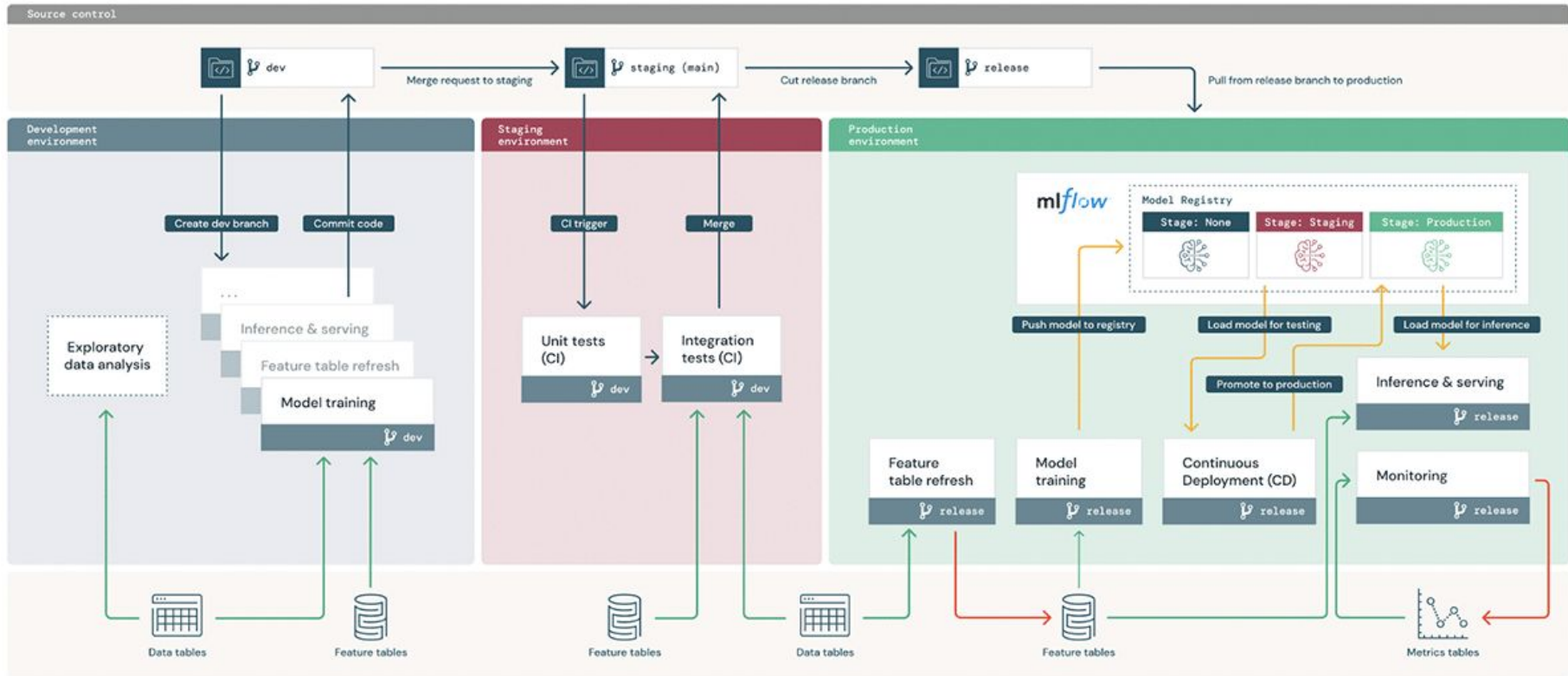
**MLOps as it exists today
is too hard**



// MLOps according to Google

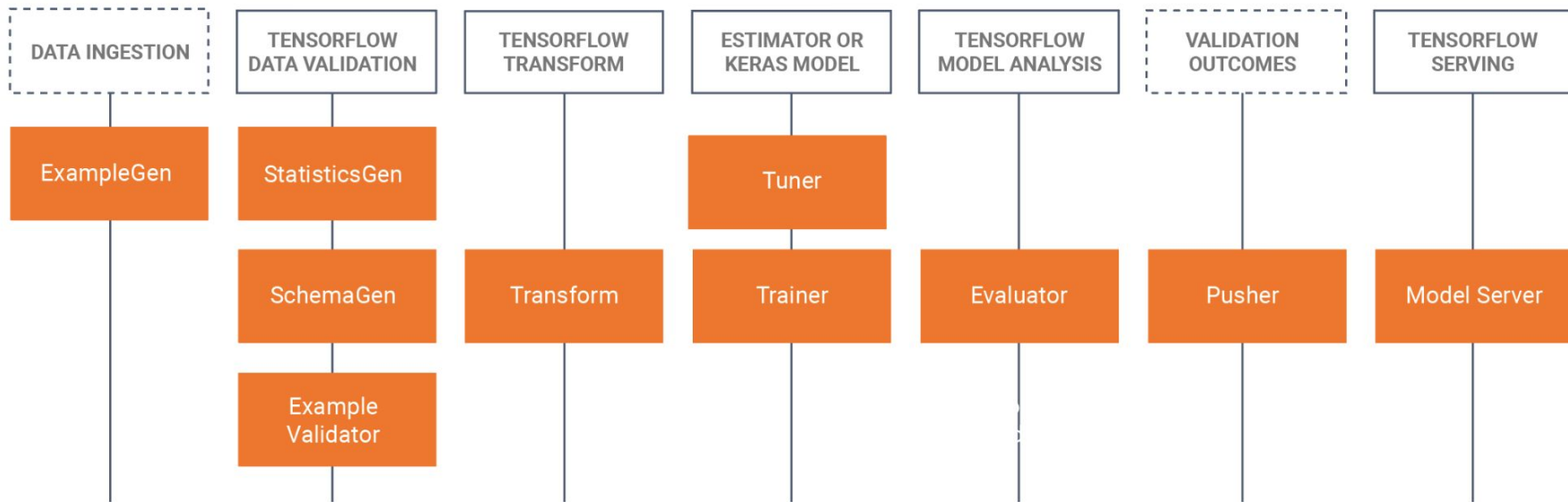


// MLOps according to Databricks





// The Mythical End-to-End Machine Learning Pipeline



[Image from <https://www.tensorflow.org/tfx/guide>]

// Problems with the Monolithic ML Pipelines

“Kubeflow Pipelines is a platform for building and deploying portable, scalable machine learning (ML) workflows ...End-to-end orchestration...**enabling you to re-use components and pipelines** to quickly create end-to-end solutions without having to rebuild each time” [[KubeFlow Website](#)] (Feb '23)]

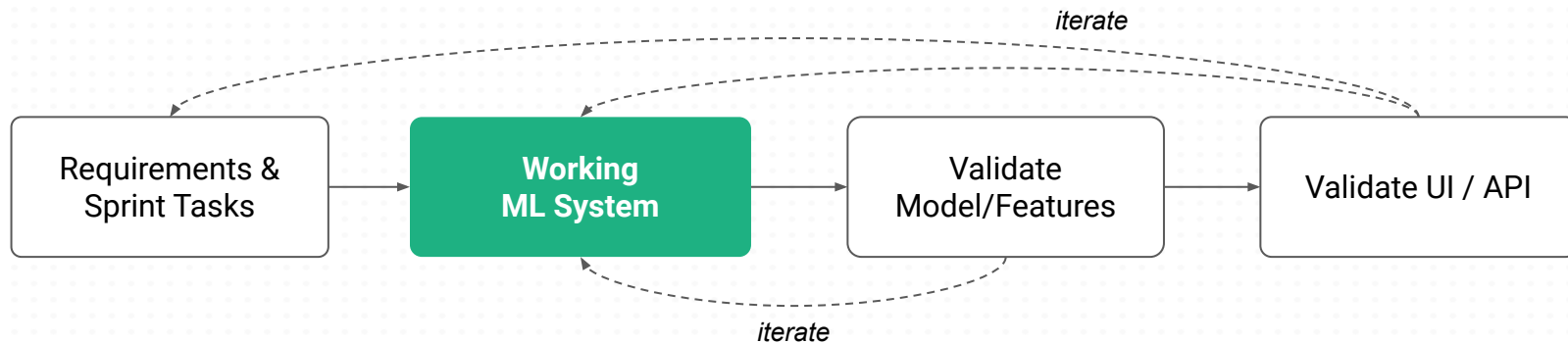
- Coupling feature engineering, model training, and inference adds development and operational complexity
- No reuse of precomputed features across multiple models (feature store)
- Too hard to get to a working MVP (minimal viable ML product)





// For developers, MLOps is about Testing, Versioning, and incremental changes

- Get to a **working ML system** with a baseline, then iteratively improve it.
- Automated **testing and versioning** of features and models
- Improve both **iteration speed** and **quality**

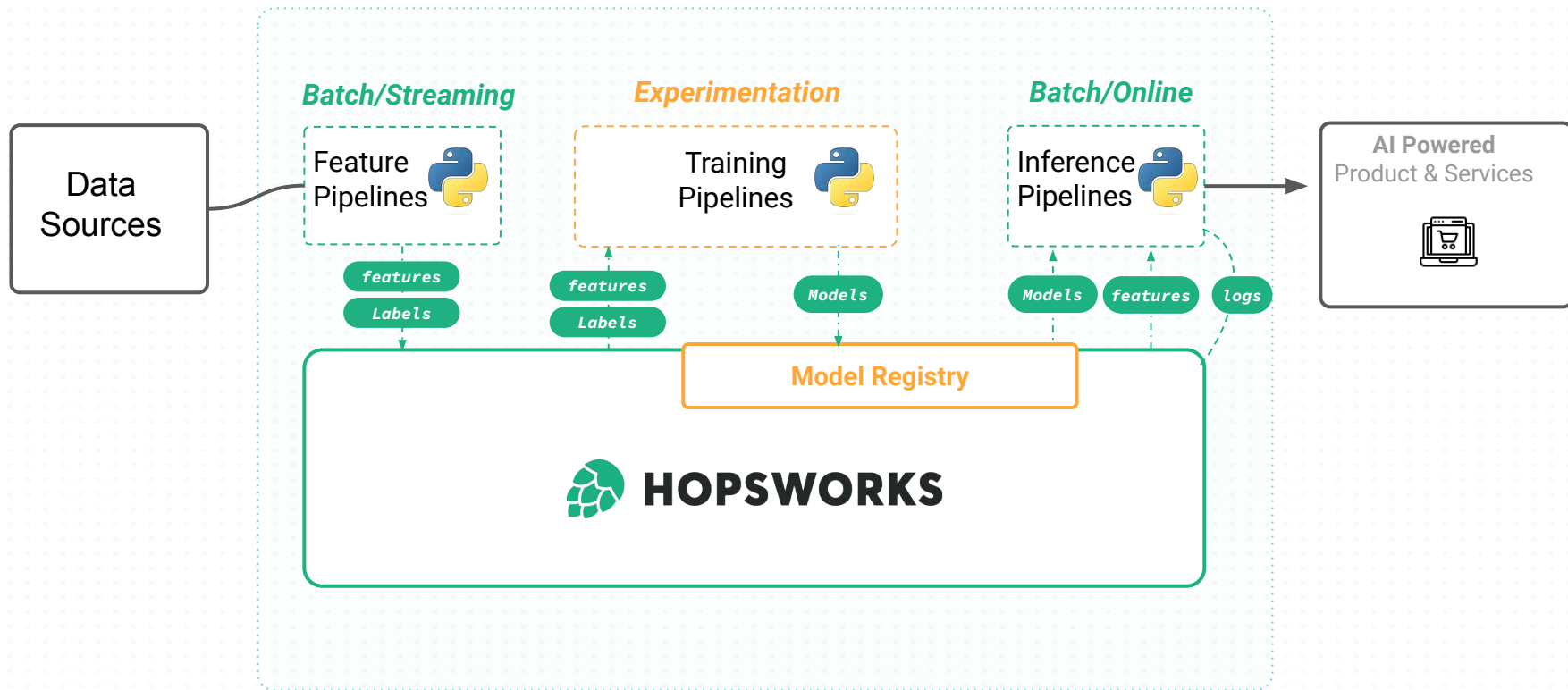


01

MLOps in Python without the infrastructure



// Feature Pipelines, Training Pipelines, and Inference Pipelines



02 Feature Pipelines



// Feature Pipelines - transform raw data into DataFrames

Feature engineering

Initialize Feature Group metadata

Validate and 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)
```



// Feature Pipelines: Python or Spark or SQL or Flink?

Spark

Good: scale, testing, complex features.

Bad: resource estimation, debugging, operations.

SQL

Good: low operational overhead, easy to implement aggregations

Bad: Complex features (e.g., embeddings) become very complex with UDFs

Python

Good: low operational overhead, rich library support, complex features possible. Pandas or Polars.

Bad: does not scale (~10 GBs with Pandas, Polars ~100 GBs),

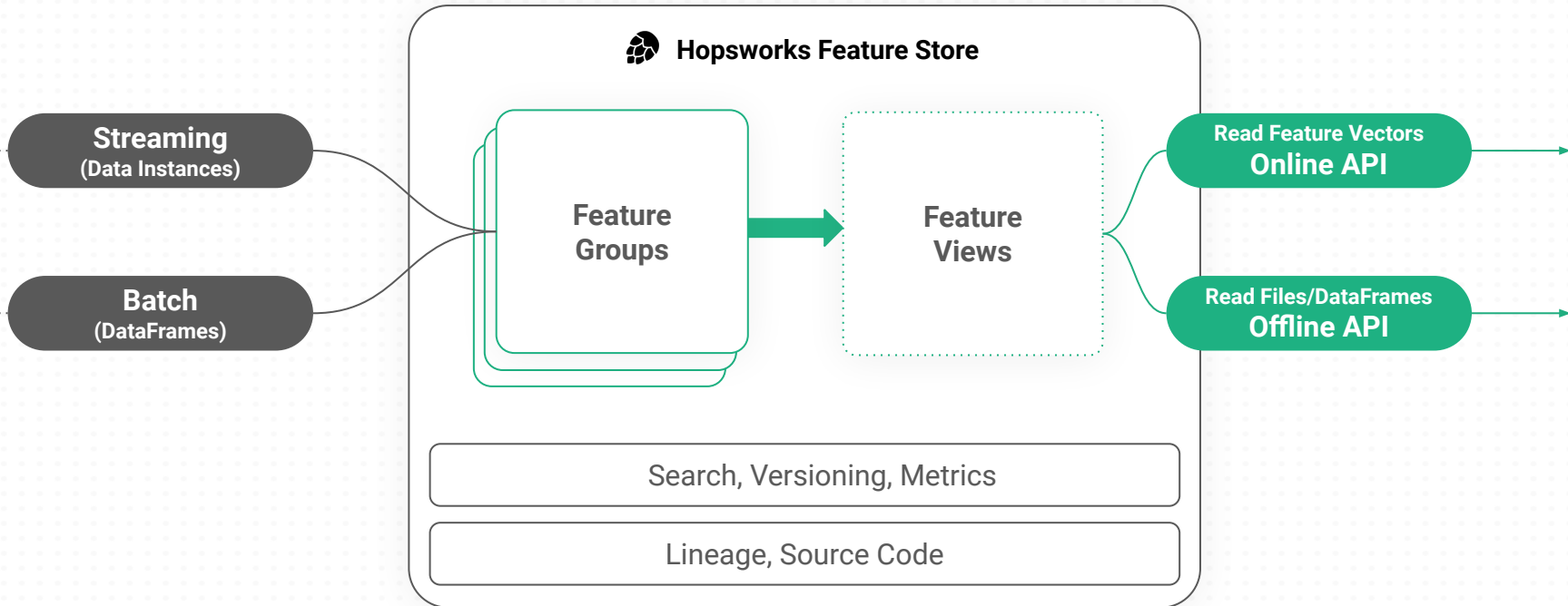
Flink Challenges

Good: low-latency, real-time feature computation that scales to huge clusters (10k+ nodes)

Bad: complex, Java-centric, high operational overhead

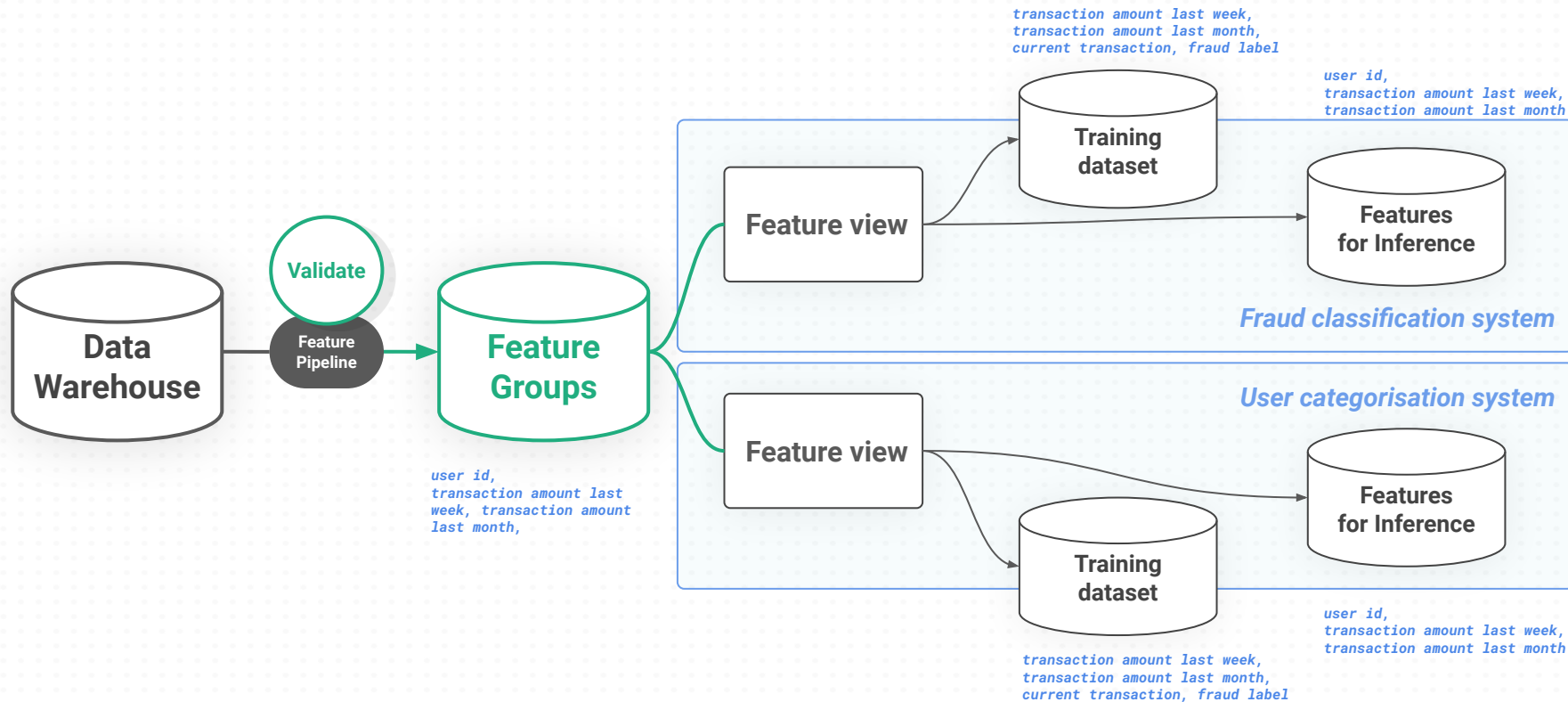


// Feature Store: Write to Feature Groups, Read from Feature Views



// One Feature Pipeline - Many Models

A **single feature pipeline** to ingest features into feature group and the feature view can reuse the features for training and inference on **many models**.



03 Training Pipelines

// Feature Store: Read from Feature Views

What does a Feature View do?

1. **API** for ML model development and operations
2. **Define** features, labels, and model-specific transformations
3. **Generate** training dataset, batch inference data, and feature vectors



// Feature View: Point-in-Time Correct JOINS

Label Feature Group

cust_id	event_time	category
1	2021-01-01	short_term
2	2021-01-02	short_term
3	2021-01-16	long_term

Transactions Feature Group (weekly)

cust_id	event_time	amount_last_week	amount_last_month
1	2021-01-01	14.00	87.00
1	2021-01-07	44.00	80.00
1	2021-01-14	78.00	82.00

Point-in-time Correct JOIN
(no data leakage)

trans.event_time => time	category	amount_last_week	amount_last_month
2022-01-01	short_term	14.00	87.00
2022-01-02	short_term	44.00	80.00
2022-01-16	long_term	78.00	82.00

Training Data

// Feature View: Point-in-Time Correct JOINS

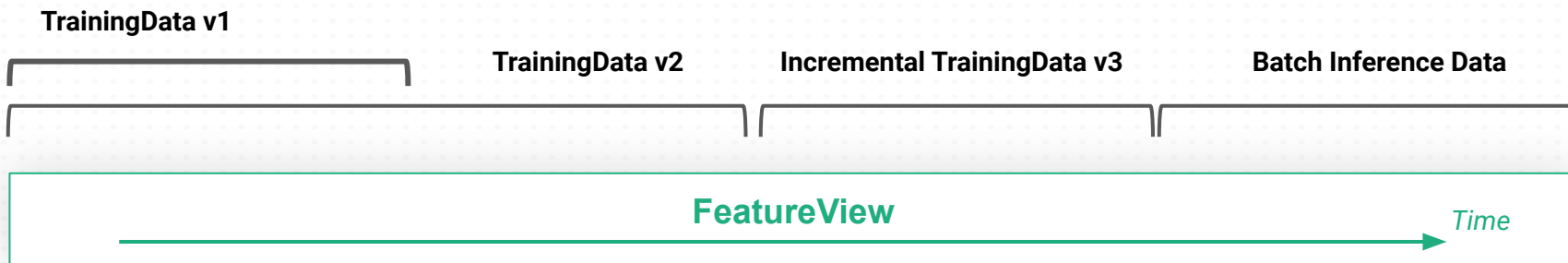
```
1 WITH right_fg0 AS (  
2 SELECT *  
3 FROM (  
4     SELECT `fg1`.`category`  
5           , `fg0`.`amount_last_seek`  
6           , `fg0`.`amount_last_month`  
7           , RANK() OVER (PARTITION BY `fg0`.`cust_id`, `fg1`.`event_time` ORDER BY `fg0`.`event_time` DESC) pit_rank_hopsworks  
8     FROM `label_fg` `fg1`  
9     INNER JOIN `transactions_fg` `fg0`  
10    ON `fg1`.`cust_id` = `fg0`.`cust_id` AND `fg1`.`event_time` >= `fg0`.`event_time`  
11 )  
12 WHERE `pit_rank_hopsworks` = 1)  
13 (SELECT `right_fg0`.`category`, `right_fg0`.`mean_cloud_perc_se3`, `right_fg0`.`mean_cloud_perc_se3`  
14 FROM right_fg0)
```



```
query = label_fg.select(["category"]).join(  
    transactions_fg.select(["amount_last_week", "amount_last_month"])  
)
```



// Feature View: Training Data and Batch Inference Data for Models



```
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 = datetime.strptime("2022-01-03 00:00:01", date_format)
14 end_time = datetime.strptime("2022-03-31 23:59:59", date_format)
15 march_transactions = feature_view.get_batch_data(
16     start_time = start_time, end_time = end_time)
```



// Feature View: Model-Specific Transformations

Training Pipeline

```
1 standard_scaler = fs.get_transformation_function(name="standard_scaler")
2 label_encoder = fs.get_transformation_function(name="label_encoder")
3
4 transformation_functions = {
5     "amount_last_month": standard_scaler,
6     "amount_last_week": standard_scaler,
7     "category": label_encoder,
8 }
9
10 fv = fs.create_feature_view(name="category_predictin",
11                             ...
12                             transformation_functions=transformation_functions)
13 X_train,X_test,y_train, y_test = fv.train_test_split(0.1)
```

Online Inference Pipeline

```
1 keys= {cust_id: 1}
2 feature_vector = fv.get_feature_vector(keys)
```

Model-Specific
Transformation
functions (UDFs)
applied
transparently

04 Inference Pipelines



// Tying it all together

The predictor interface

Establish connection to Feature Store

Initialize Feature View

Loading the model from a model registry

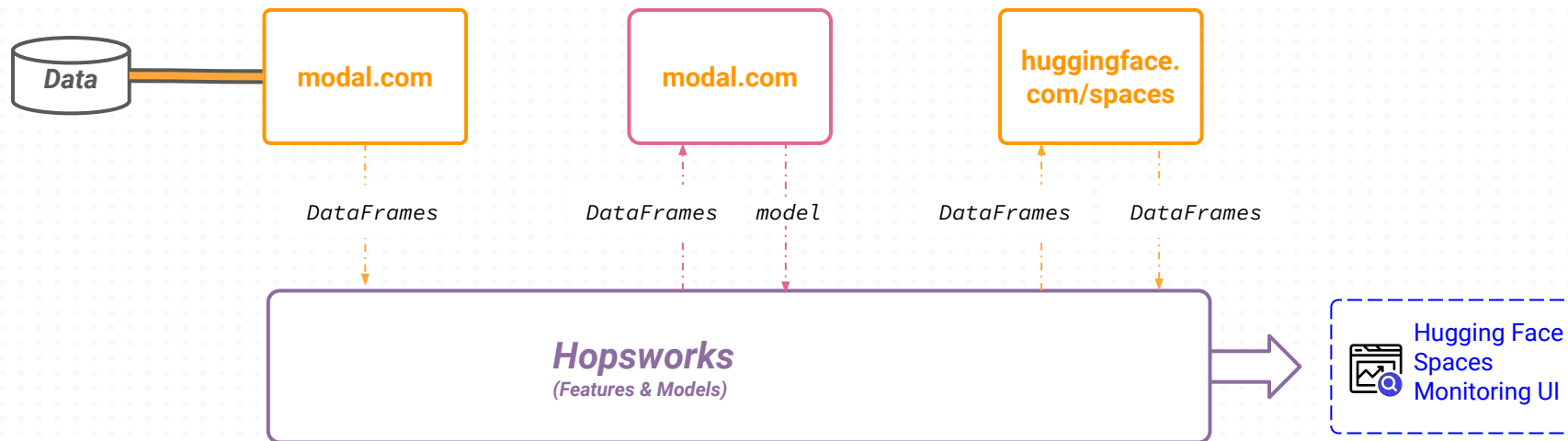
Implementing the predict interface using the Feature View

```
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()
```

05 Demo



// Serverless Machine Learning with Hopsworks, Modal, and HF Spaces

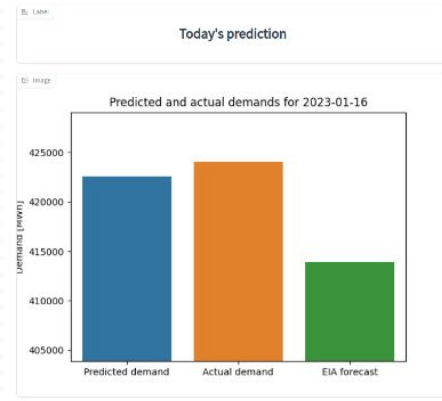




// Example Serverless Machine Learning Systems



Air Quality Predictions for Poland



New York Electricity Price Demand

Like It Or Not

A machine learning service that predicts the number of likes and upvote ratio of your reddit post before you submit it. The initial computation may take a few seconds, as the model must be downloaded. Please be patient.

Output

Likes: **580** Upvote Ratio: **61.16%**

It seems like your post will get a lot of likes, you should try to make it even more interesting! You can expect an upvote ratio of 61.3% which means that 61.3% of the people who see your post will upvote it (and 38.7% will downvote it). This means that people will have mixed feelings about your post.

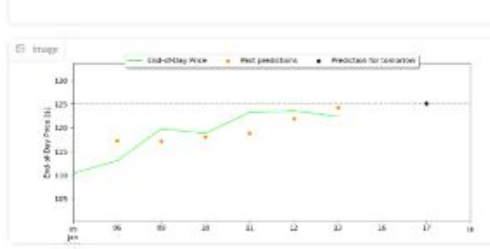
Explanation

Below you can see how different features of your post affected the final prediction. The diagram shows the default value that would have been predicted in case no features about your post were known. In addition, every feature is associated with a bar the color and length of which indicate the magnitude and type of impact it had on the prediction. A long bar with red color states that the feature increased the prediction value by a large amount. The exact meaning of the feature names and their values can be found at the [page of the model's metadata](#).

Recent Predictions

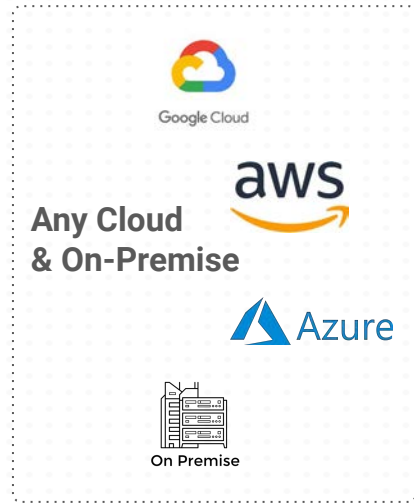
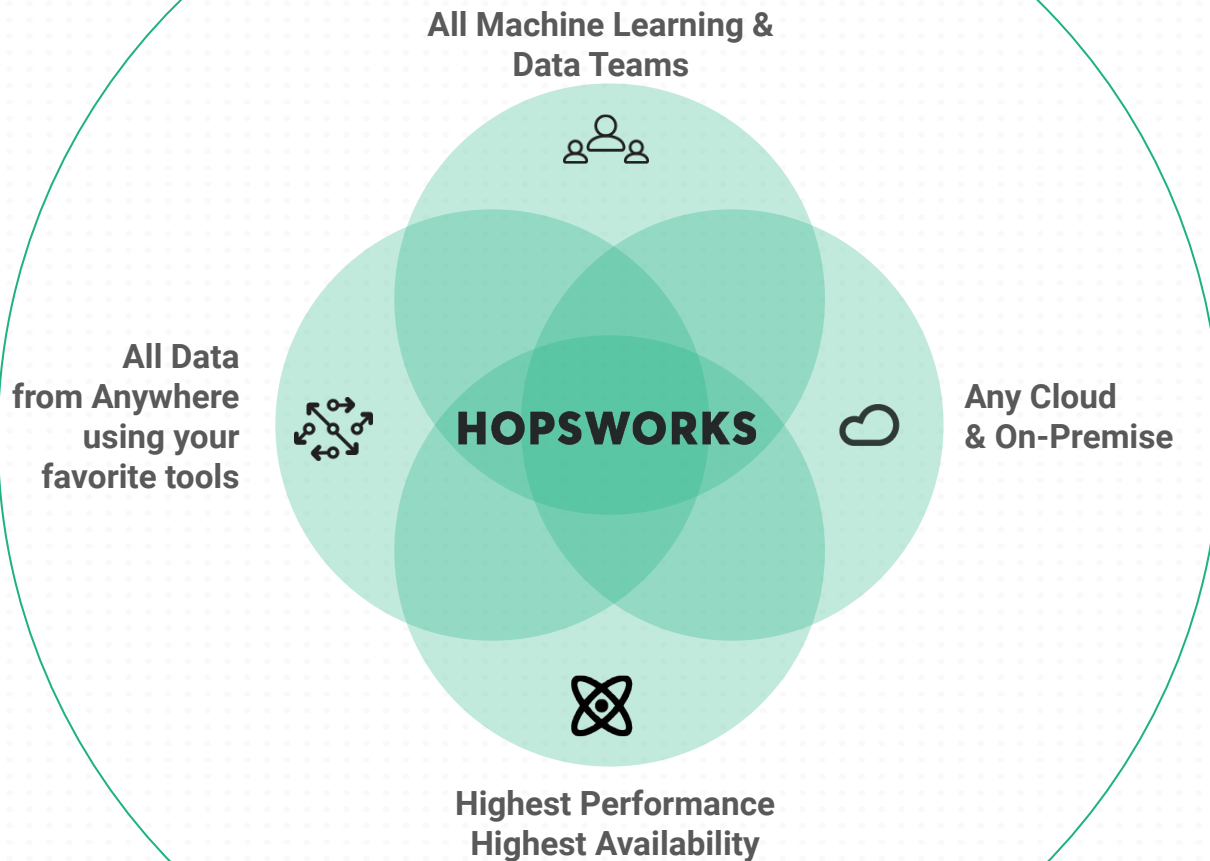
Index	Date	Pred [\$]	Pred [%]	True [\$]	True [%]
0	2023-01-06	117.289150	0.062979	113.059998	0.024651
1	2023-01-09	117.124921	0.035954	119.769997	0.059349
2	2023-01-10	118.120352	-0.013773	118.849998	-0.007681
3	2023-01-11	118.876772	0.000225	123.220001	0.036769
4	2023-01-12	121.877310	-0.010897	123.559998	0.002759
5	2023-01-13	124.272293	0.005765	122.400002	-0.009388
6	2023-01-17	125.150664	0.022473	NaN	NaN

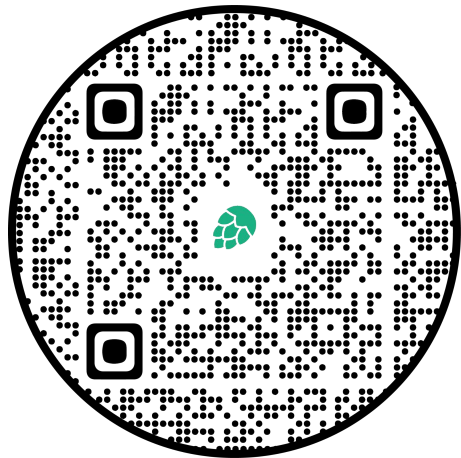
Recent Stock Prices



Will your reddit post be liked or not?

Tesla Stock Price Prediction





public-hopsworks.slack.com

- # Join our slack community
- # Explore our latest tutorials
- # Ask us any questions



SERVERLESS ML

www.serverless-ml.org