



An open source, enterprise-grade,
high-performance feature store

<https://github.com/feathr-ai/feathr>

- built at LinkedIn & Microsoft
- a Linux Foundation AI & Data Sandbox project
- Now natively integrated with Azure/AWS, and Databricks



Agenda

1 The Problem

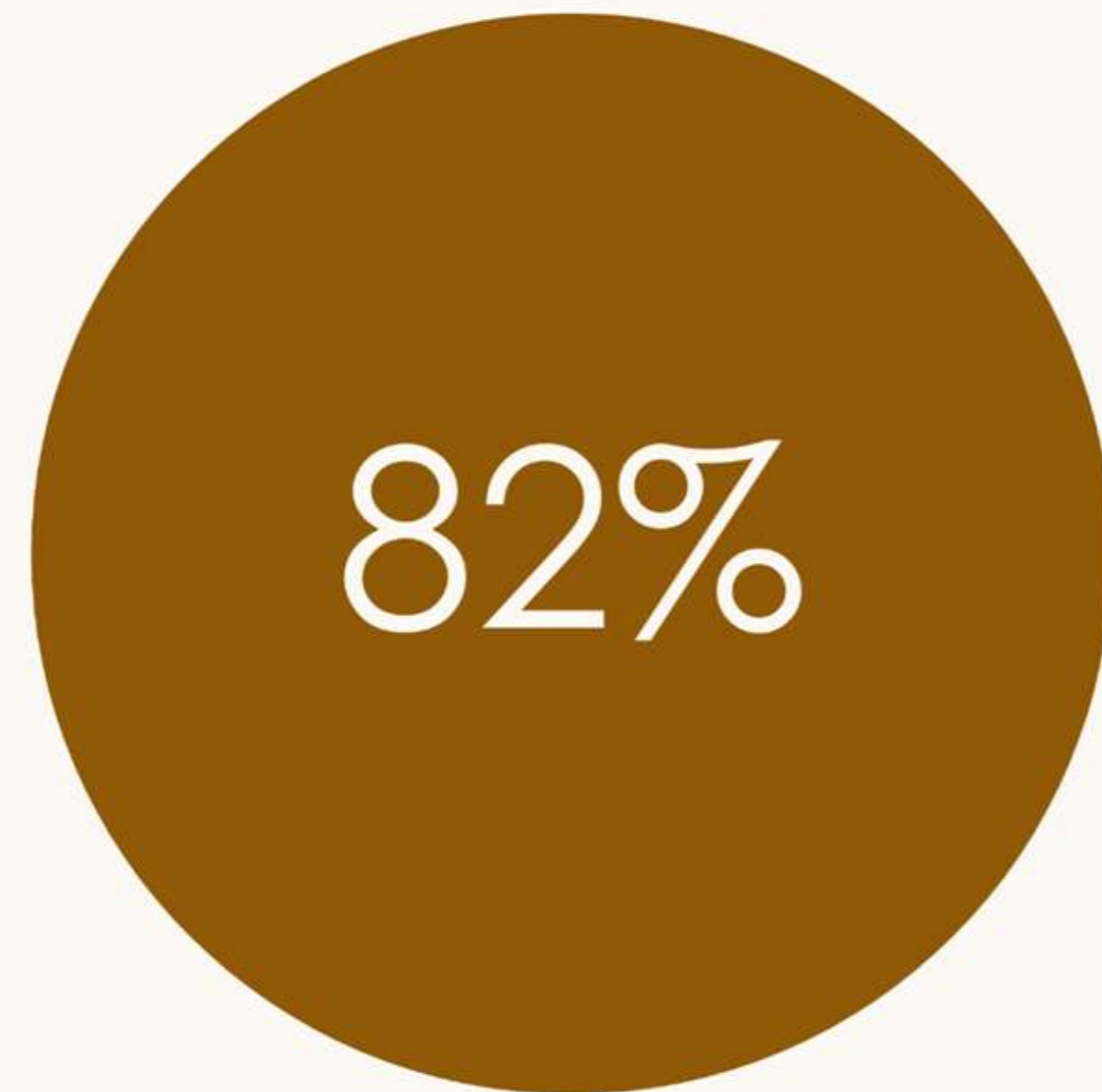
2 The Solution

3 The Use Case

4 Summary

Survey in Forbes: Big data engineering for AI

Re: Building training sets + Cleaning and organizing data + Collecting datasets



Time spent on data preparation



Respondents said data preparation 'least enjoyable' part of data science

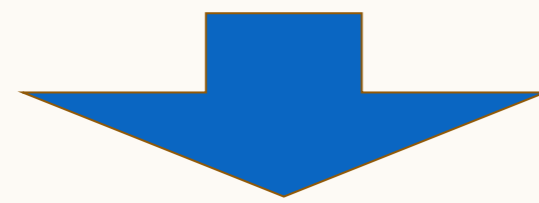
Like a music streaming app ...

for feature engineering

Music "Workflow"

Old
Way

- Manually get music files from **various sources**
- Convert them to a **format** my device can play
- **Load** onto my device (different for home/car)
- Worry about **storage**, bitrate, **compatibility**



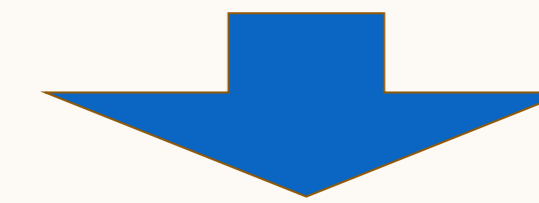
New
Way

- **Just ask virtual assistant to play the song.**



ML Feature Workflow

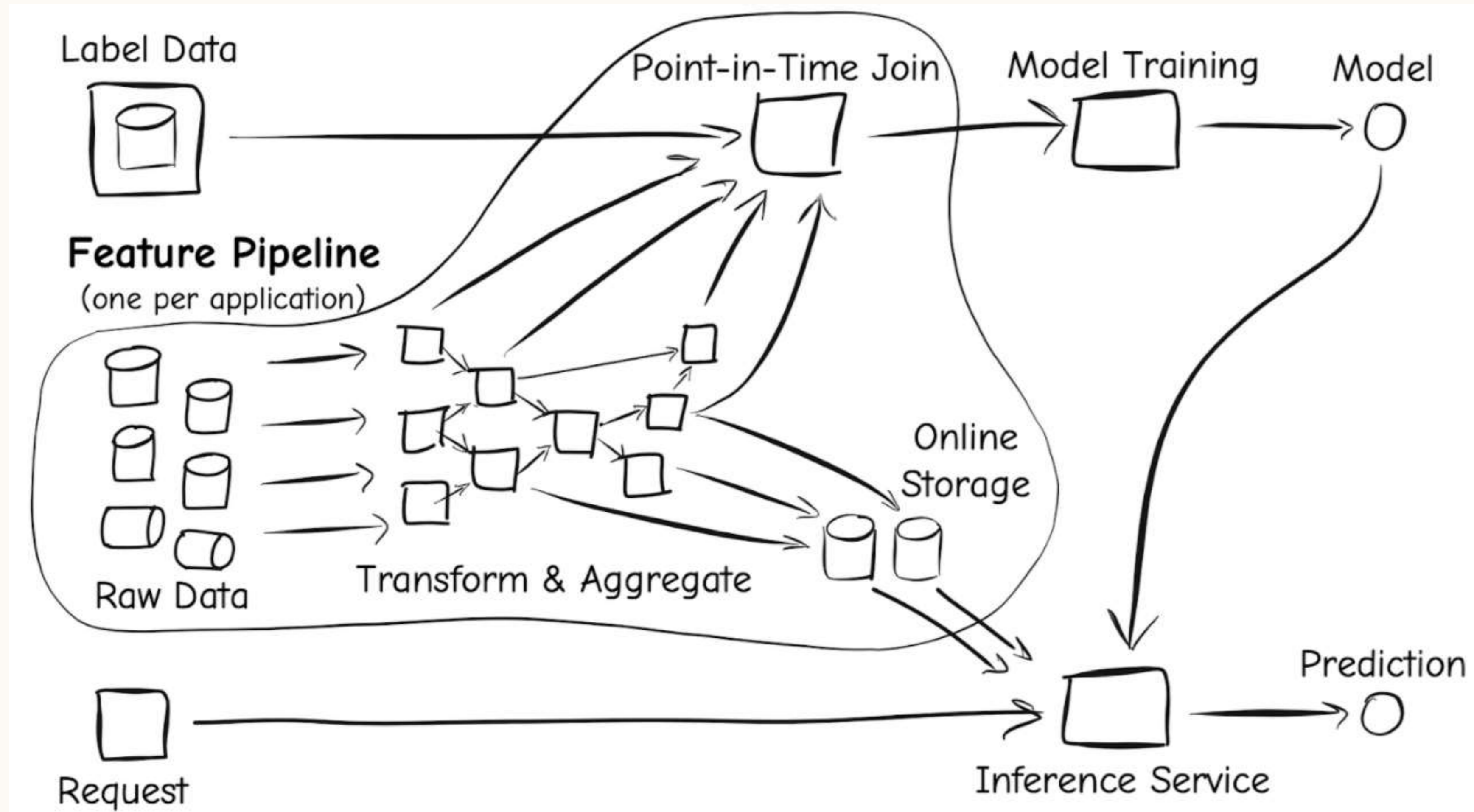
- Write jobs to get entity data from various sources
- Extract, aggregate, join, convert into proper format
- Load into model framework (different for train/serving)
- Worry about scale, perf, leakage, train/serve skew



- **Just import the feature by name into model code.**
- If feature doesn't exist, define and register it via simple APIs.

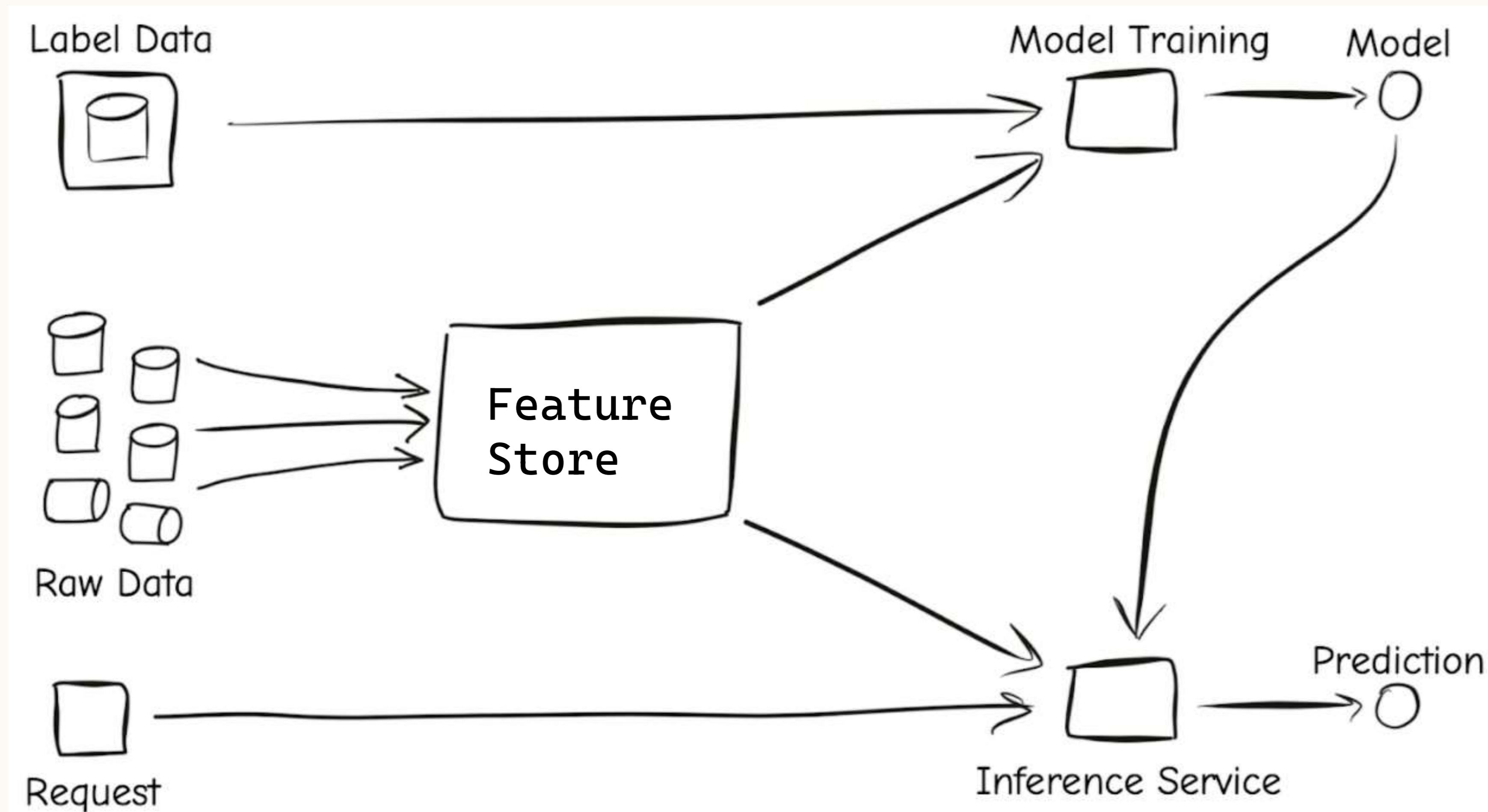
Why Feathr

Handle all the complex weight-lifting and “boring” work automatically



Why Feathr

Handle all the complex weight-lifting and “boring” work automatically



Like a package manager for feature engineering

Code

```
import module1  
import module2  
import module3  
import module4
```



Features

```
query = FeatureQuery(  
    feature_list=[  
        "feature_1",  
        "feature_2",  
        "feature_3",  
        "feature_4"  
    ],  
    key=item_id)
```

Problem: The complexity of feature preparation pipelines

1. Load & Transform

- Different programming APIs for different environments, e.g. online, offline, nearline, etc.
- Boilerplate and repetitive work
- Hard to test and debug

2. Train/Inference Skew

- Offline training and online inference usually require different data serving pipelines.
- Ensuring features are generated consistently is time intensive and error prone.
- Teams are deterred from using real time data for inferencing due to the difficulty in serving the right data.

3. Re-use and share

- The cost of building and maintaining feature pipelines was borne redundantly across many teams.
- Team-specific pipelines also made it impractical to reuse features across projects. e.g. no common type system, no common feature namespace

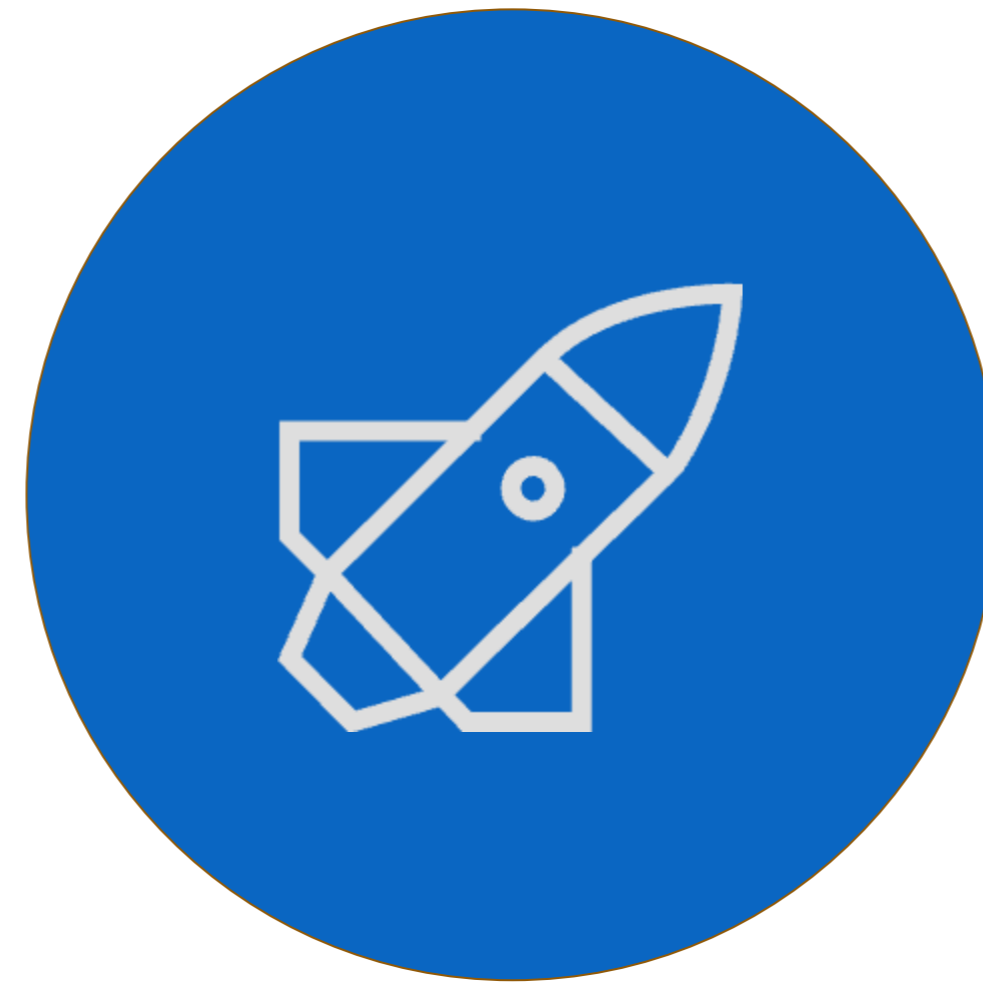
What a feature store
should be

Feature store principal use cases



Develop Features

Based on raw data,
using simple APIs



Deploy Features

For training and online
model inferencing

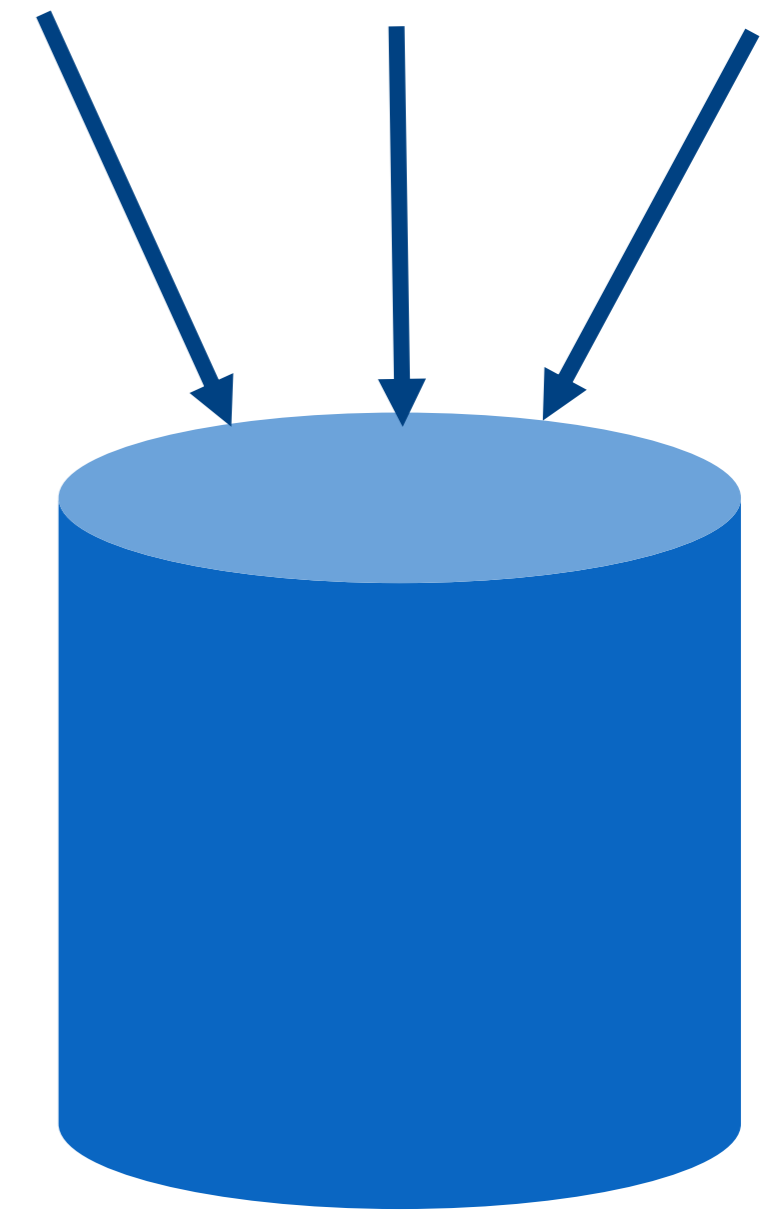


Manage Features

Monitor feature health
and share across teams

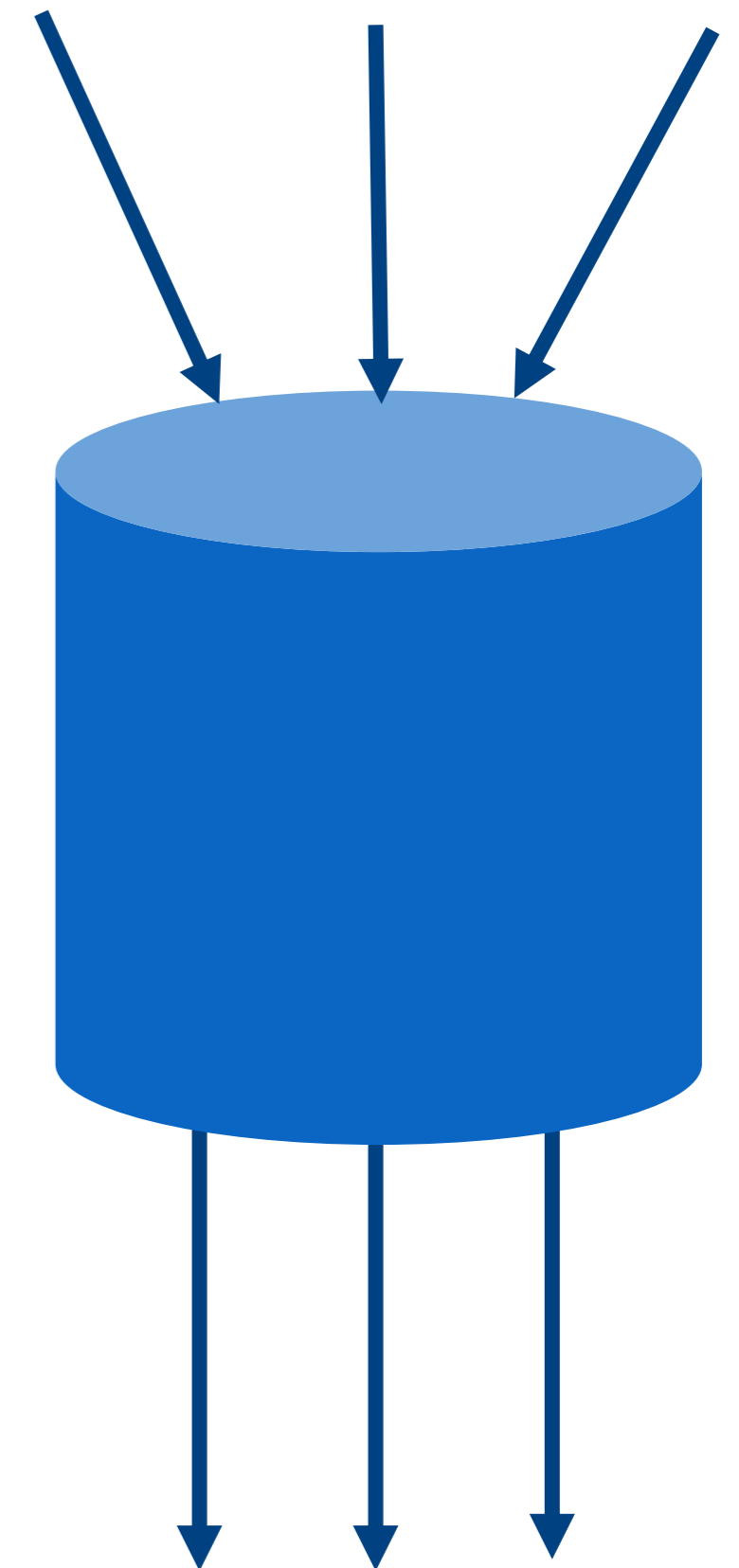
The “feature store” abstraction

- “Put a feature in” (Producer)
 - Develop a feature based on **raw data sets**
 - Sliding time windows
 - Aggregations
 - Transformations
 - Lookups/joins
 - Develop a feature based on other feature(s)



The "feature store" abstraction

- "Put a feature in" (Producer)
 - Develop a feature based on **raw data sets**
 - Sliding time windows
 - Aggregations
 - Transformations
 - Lookups/joins
 - Develop a feature based on other feature(s)
- "Get some features out" (Consumer)
 - Join features to training labels
 - Backfill historical values of features
(**point-in-time correctness**)
 - Efficiently compute, store, and serve features for **online inference**





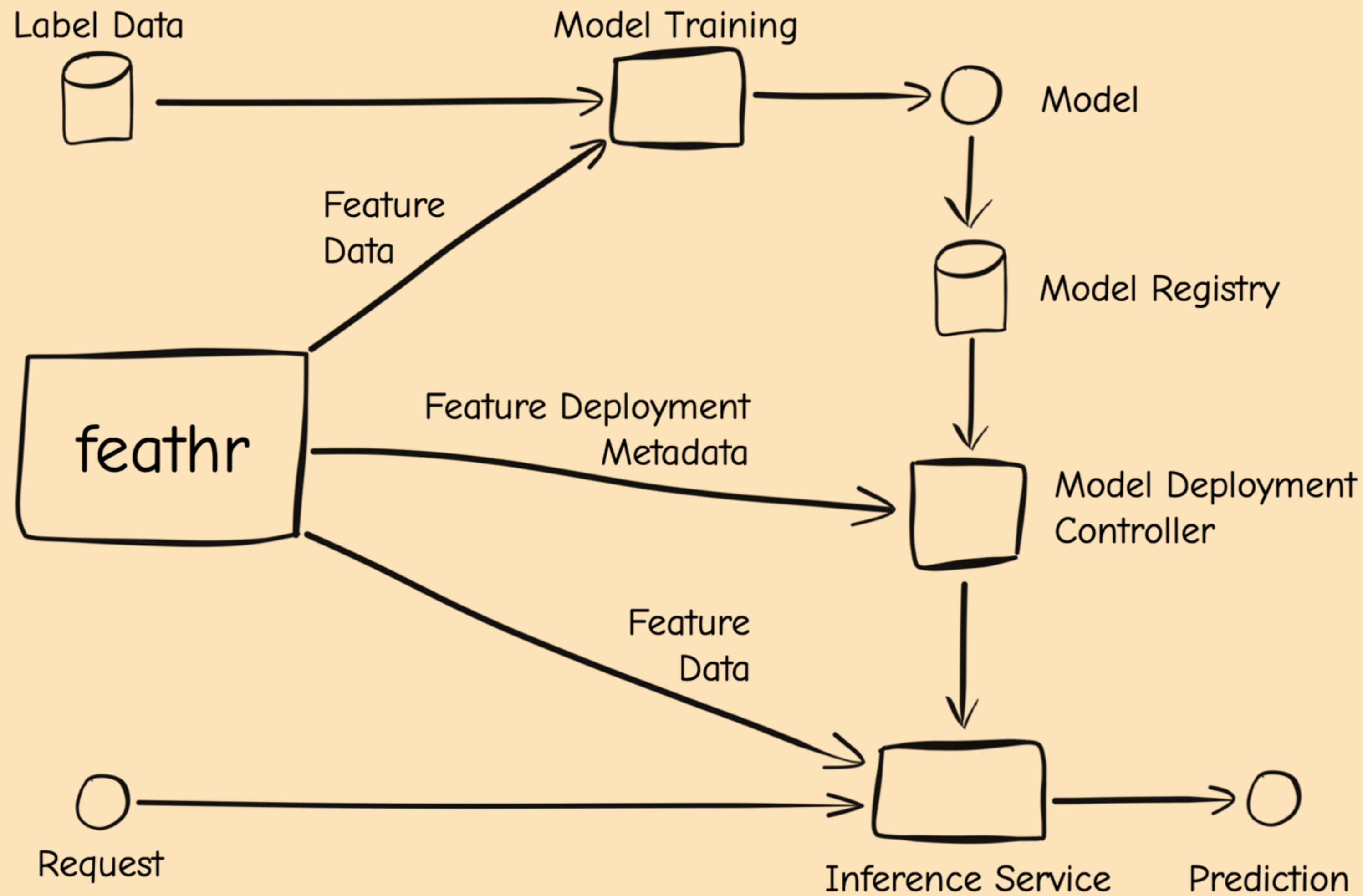
Feathr at LinkedIn

Introducing Feathr, a battle tested feature store built by LinkedIn

The screenshot shows the LinkedIn Jobs interface. At the top, there are navigation options: JOBS, MANAGE JOBS, and POST A JOB. A search bar is present with the text "Search for something or someone". Below this, a job listing for "Mechanical Engineers" is shown, with details like "San Francisco Bay Area", "Opened 28d ago", "Daily budget: \$15.00", and "Total spend: \$30.00". There are buttons for "Edit job & budget" and "Close job". Below the job listing, there are tabs for "Candidates" and "Pipeline (0)", and a "13 Applicants" button. The main content area shows a candidate profile for Mae Norris, a User Experience Designer at LinkedIn. Her profile includes a summary, a "Summary" section with a quote, and an "Experience" section listing her role as "Account Manager" at "Freshing" from Jan 2014 to Present. There are buttons for "Skip", "Not interested", "1-click message", and "Preview message".

The screenshot shows a mobile LinkedIn feed. At the top, there is a search bar with the text "Search for people, jobs, posts and m...". Below this, there are tabs for "Feed" and "Featured". A post by Tomer Cohen, Head of Content, Search & Discovery Products, is shown. The post text says "Congrats Lior Ron on the acquisition! Drinks are on you next time we meet." Below the text is a photo of a white brain. The post has 1.2k Likes and 117 Comments. Below the post, there are buttons for "Like", "Comment", and "Share". At the bottom, there is a navigation bar with icons for Home, My Network, Messaging, Notifications, and Jobs. To the right of the feed, there is a job listing for "Principal Data Scientist" at Microsoft in Bangalore, IN. The job listing includes a "Save" button and an "Apply" button. Below the job listing, there is a section for "32 connections can refer you" with a button for "AS".

[Learning Hiring Preferences: The AI Behind LinkedIn Jobs](#)
[Personalized Recommendations in LinkedIn Learning](#)
[Helping members connect to opportunity through AI](#)
[Near real-time features for near real-time personalization](#)



Feathr is a pillar of LinkedIn's ML platform

Model deployment service uses Feathr to ensure a model's feature dependencies are deployed, before deploying the model.

Feathr at LinkedIn

- hundreds of models
- thousands of features
- many kinds of entities (economic graph)
- petabyte scale

Timeline

- | | |
|------|---|
| 2017 | Initial development and launch |
| 2018 | Broad adoption within LinkedIn |
| 2020 | Majority of LinkedIn ML applications onboarded |
| 2022 | Open source, Azure Integration, joined Linux Foundation AI & Data |

Impact at LinkedIn

Majority of ML applications at LinkedIn have adopted Feathr



Improved Productivity

Faster experimentation with new features, from weeks to days



Improved Performance

Running time improved over custom pipelines, as much as 50%



Improved Collaboration

Applications can share features, which was difficult previously

What is Feathr

An abstraction layer between raw data and model

- **Define** features based on raw data sources using simple APIs.
- **Get** those features by their names during model training and model inferencing.
- **Share** features across your team and organizations.

Highlights



Cloud-native

Native integration with
Azure and AWS



Rich Transformation

Python built-in
transformations and
PySpark UDF, on-demand
evaluation



Scalable & High Performance

Highly optimized feature
compute engine

Use case: Create Feature Definition

Load raw source data, and define transformation

```
batch_source = HdfsSource(  
    name="nycTaxiBatchSource", # Source name to enrich your metadata  
    path="abfss://green_tripdata_2020-04.csv", # Path to your data  
    event_timestamp_column="lpep_dropoff_datetime", # Event timestamp for point-in-time correctness  
    timestamp_format="yyyy-MM-dd HH:mm:ss") # Supports various formats including epoch  
  
trip_id = TypedKey(key_column="trip_id",  
                  key_column_type=ValueTypes.INT64,  
                  description="trip id")  
  
features = [  
    Feature(name="f_trip_distance", # Ingest feature data as-is  
           feature_type=FLOAT,  
           key=trip_id),  
    Feature(name="f_is_long_trip_distance",  
           feature_type=BOOLEAN,  
           transform="cast_float(trip_distance)>30",  
           key=trip_id) # SQL-like syntax to transform raw data into feature  
]  
  
anchor = FeatureAnchor(name="anchor_features", # Features anchored on same source  
                       source=batch_source,  
                       features=features)
```

Use Case - Streaming Feature

Create features from streaming source

```
stream_source = KafkaSource(name="kafkaStreamingSource",
                             kafkaConfig=KafkaConfig(brokers=["feathrazureci.servicebus.windows.net:"],
                                                       topics=["feathrcieventhub"],
                                                       schema=schema)
                             )

driver_id = TypedKey(key_column="driver_id",
                    key_column_type=ValueTypes.INT64,
                    description="driver id",
                    full_name="nyc driver id")

kafkaAnchor = FeatureAnchor(name="kafkaAnchor",
                             source=stream_source,
                             features=[Feature(name="f_modified_streaming_count",
                                              feature_type=INT32,
                                              transform="trips_today + 1",
                                              key=driver_id),
                                      Feature(name="f_modified_streaming_count2",
                                              feature_type=INT32,
                                              transform="trips_today + 2",
                                              key=driver_id)]
                             )
```


Use Case - Feature Materialization

Materialize feature values to online storage for realtime access

```
client = FeathrClient()
redisSink = RedisSink(table_name="nycTaxiDemoFeature")
# Materialize two features into a redis table.
settings = MaterializationSettings("nycTaxiMaterializationJob",
sinks=[redisSink],
feature_names=["f_location_avg_fare", "f_location_max_fare"])
client.materialize_features(settings)
```

Use Case - Feature Sharing and Discovery

Share features and discover features

The screenshot displays the Feathr web interface. On the left is a dark sidebar with navigation options: Homepage, Data Sources, Features (highlighted), Jobs, and Monitoring. The main content area is titled "Lineage feathr_ci_registry_39_6_728496" and includes filter tabs for "All Features", "Source", "Anchor", "Anchor Feature", and "Derived Feature".

The feature lineage diagram shows two source features: "nycTaxiBatchSource" (feathr_source_v1) and "PASSTHROUGH" (feathr_source_v1). "nycTaxiBatchSource" is linked to "aggregationFeatures" (feathr_anchor_v1) and "f_location_avg_fare" (feathr_anchor_feature_v1). "PASSTHROUGH" is linked to "f_is_long_trip_distance" (feathr_anchor_feature_v1), "request_features" (feathr_anchor_v1), "f_day_of_week" (feathr_anchor_feature_v1), "f_trip_time_duration" (feathr_anchor_feature_v1), and "f_trip_distance" (feathr_anchor_feature_v1). "f_trip_time_duration" and "f_trip_distance" are further linked to "f_trip_time_rounded" (feathr_derived_feature_v1) and "f_trip_time_distance" (feathr_derived_feature_v1). "f_trip_time_rounded" is linked to "f_trip_time_rounded_plus" (feathr_derived_feature_v1).

On the right, there are two metrics charts under the "Metrics" tab. The first chart, titled "avg", shows a line graph with data points fluctuating around a value of 6.5 from 2022-01-02 to 2022-01-11. The second chart, titled "max", shows a line graph with data points consistently around a value of 9.0 over the same period.

Use Case - Derived Feature

Define features on top of other features

```
# Compute a new feature(a.k.a. derived feature) on top of an existing feature
derived_feature = DerivedFeature(name="f_trip_time_distance",
                                feature_type=FLOAT,
                                key=trip_key,
                                input_features=[f_trip_distance, f_trip_time_duration],
                                transform="f_trip_distance * f_trip_time_duration")

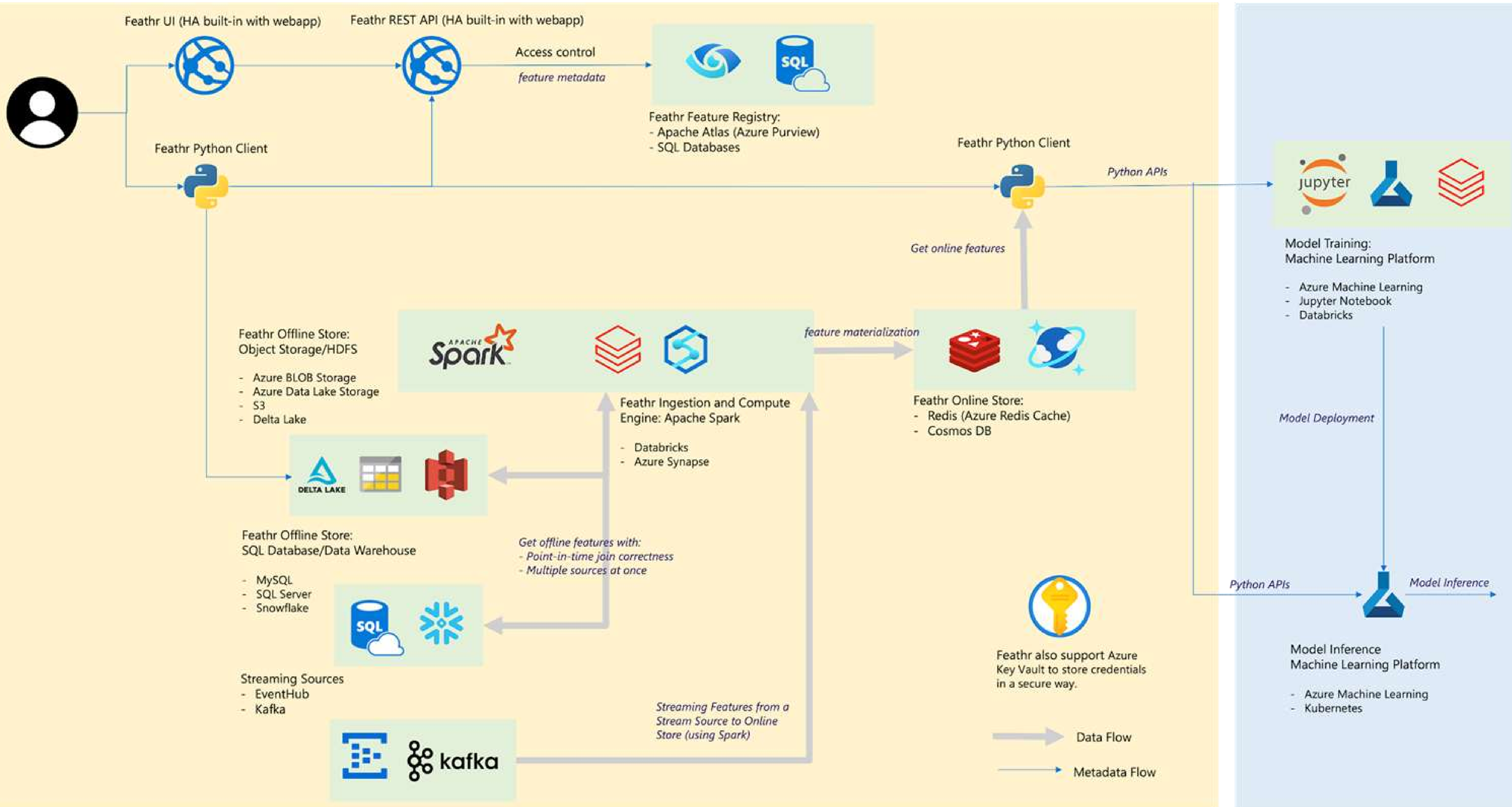
# Another example to compute embedding similarity
user_embedding = Feature(name="user_embedding", feature_type=DENSE_VECTOR, key=user_key)
item_embedding = Feature(name="item_embedding", feature_type=DENSE_VECTOR, key=item_key)

user_item_similarity = DerivedFeature(name="user_item_similarity",
                                     feature_type=FLOAT,
                                     key=[user_key, item_key],
                                     input_features=[user_embedding, item_embedding],
                                     transform="cosine_similarity(user_embedding, item_embedding)")
```

Feathr Highlights – Scalability

- Capable of processing tens of billions of rows and PB scale data
- Native optimizations like bloom filters, join plan optimizer, salted join
- Incremental joins for large dataset

Feathr Architecture



Demo and Q&A

More Resources

Source code – welcome to start & fork!

<https://github.com/feathr-ai/feathr>

Tutorials:

[Introduction to Feathr - Beginner's guide](#)

[Notebook tutorial: Build a Product Recommendation Machine Learning Model with Feathr Feature Store](#)

Slack invitation:

https://join.slack.com/t/feathrai/shared_invite/zt-1ffva5u6v-voq0Us7bbKAw873cEzHOSg

Summary

- **Feathr** is an open-source feature store which can be seen as an abstraction layer between raw data and model.
- **Feathr** allows users to define features with transformation on top of raw data source and get feature values by feature name during both training and inferencing.
- **Feathr** simplifies feature preparation workflows and enables feature sharing across teams and company.

Thank you

(Check out our GitHub: <https://github.com/feathr-ai/feathr>)