

A Query-Based Feature Store at OLX



OLX Brasil

OLX Brasil

The logo for OLX Brasil, featuring the word "OLX" in a bold, purple, sans-serif font, followed by the word "Brasil" in a smaller, purple, sans-serif font.

Biggest
C2C
Marketplace in
Brasil

Tudo pra sua casa! 🏠 Use CASA10 e garanta 10% OFF com a Compra Segura. Parcela em até 12x sem juros! [Vem!](#)



Buscar



Plano Profissional



Meus Anúncios



Chat



Notificações



Entrar

Anunciar

Estou procurando por...



Imóveis



Autos e peças



Para a sua casa



Eletrônicos e celulares



Vagas de emprego



Serviços



Músicas e hobbies



Esportes e lazer



Moda e beleza



Agro e indústria



Todas as Categorias

Anúncios recentes



Caixa de som automotiva
R\$ 170



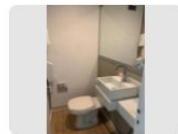
Projeto Benq 4.000 mil lumens + Telão 02 x 02...
R\$ 3.000



Compressor de ar honda fit 2005
R\$ 800



Vendo Pia de banheiro com coluna bege 69 rea...
R\$ 69



Sala Comercial Centro da Barra
R\$ 220.000



Câmera de ré automotiva ?
R\$ 70



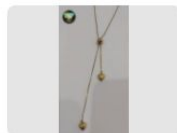
Divisor para Antena Digital e Tv a Cabo 1/ 2...
R\$ 14



PS4 1tb
R\$ 1.350



Ótima casa com 2 quartos, na rua Cristóvão...
R\$ 160.000



Colar ajustável de ouro 18k/750.
R\$ 840



Astra 2006/2006
R\$ 22.000



Hb20 1.6 automatico 2015
R\$ 48.000

Feature Stores for ML

6M

Daily Active Users

500K

New Ads Per Day

5M

Chat Messages Per Day

Problems



How to
minimize
scam
impact?

How to
minimize
scam
impact?

We use
Machine
Learning!!

Chat Moderation

User Moderation

Ad Moderation



Chat Moderation

User Moderation

Ad Moderation



create

ML Model

Chat Moderation

User Moderation

Ad Moderation



Create

ML Model

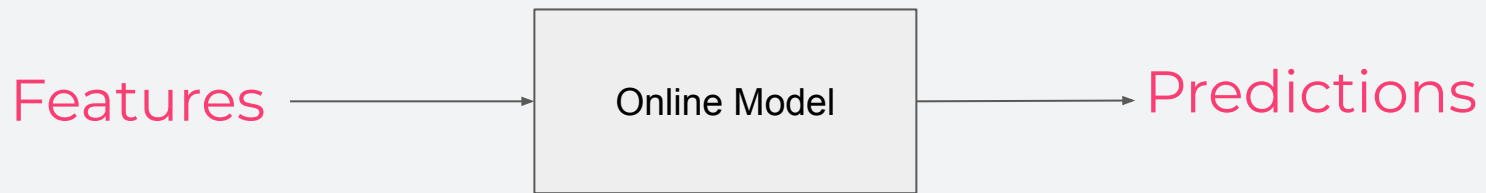
Take decision

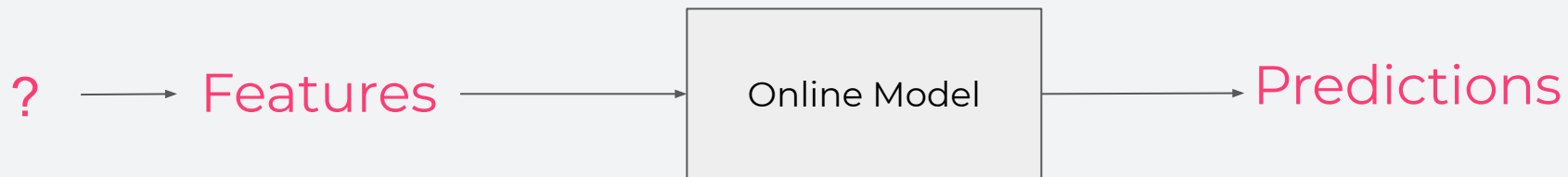


Chat Moderation

Chat Moderation

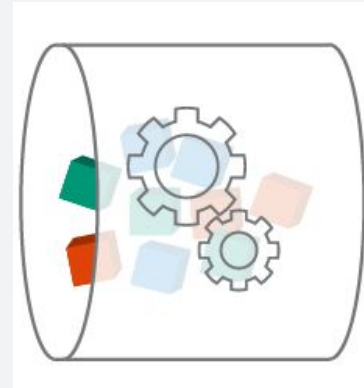
Need to ban users ASAP,
in minutes, based on
chat behavior



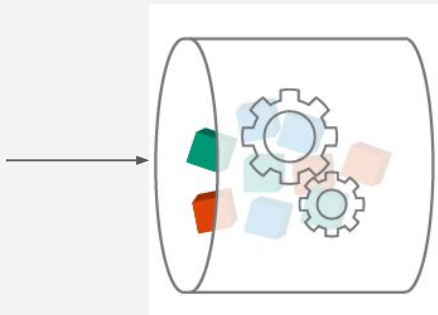


Chat Moderation

They needed a
Transformation
Engine to serve
features in Real
Time



Raw
Data

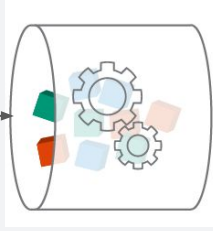


Online Model

Predictions

Imagine if each
model in
production
should replicate
this infra?

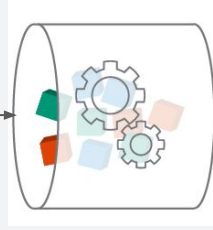
Raw Data



Online Model

Predictions

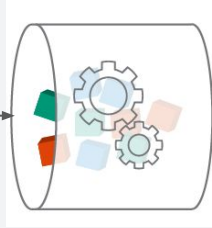
Raw Data



Online Model

Predictions

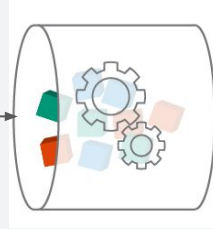
Raw Data



Online Model

Predictions

Raw Data

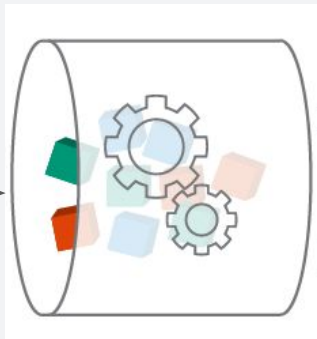


Online Model

Predictions

Feature Stores for ML

Raw Data



Online Model

Predictions

Online Model

Predictions

Online Model

Predictions

Who you gonna
call?

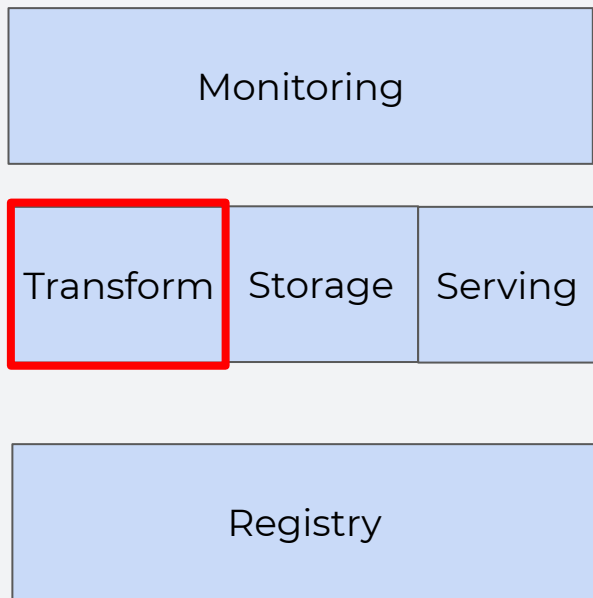
Who you gonna
call?

The MLOps
Platforms Team

This would be the
start of our
Feature Store

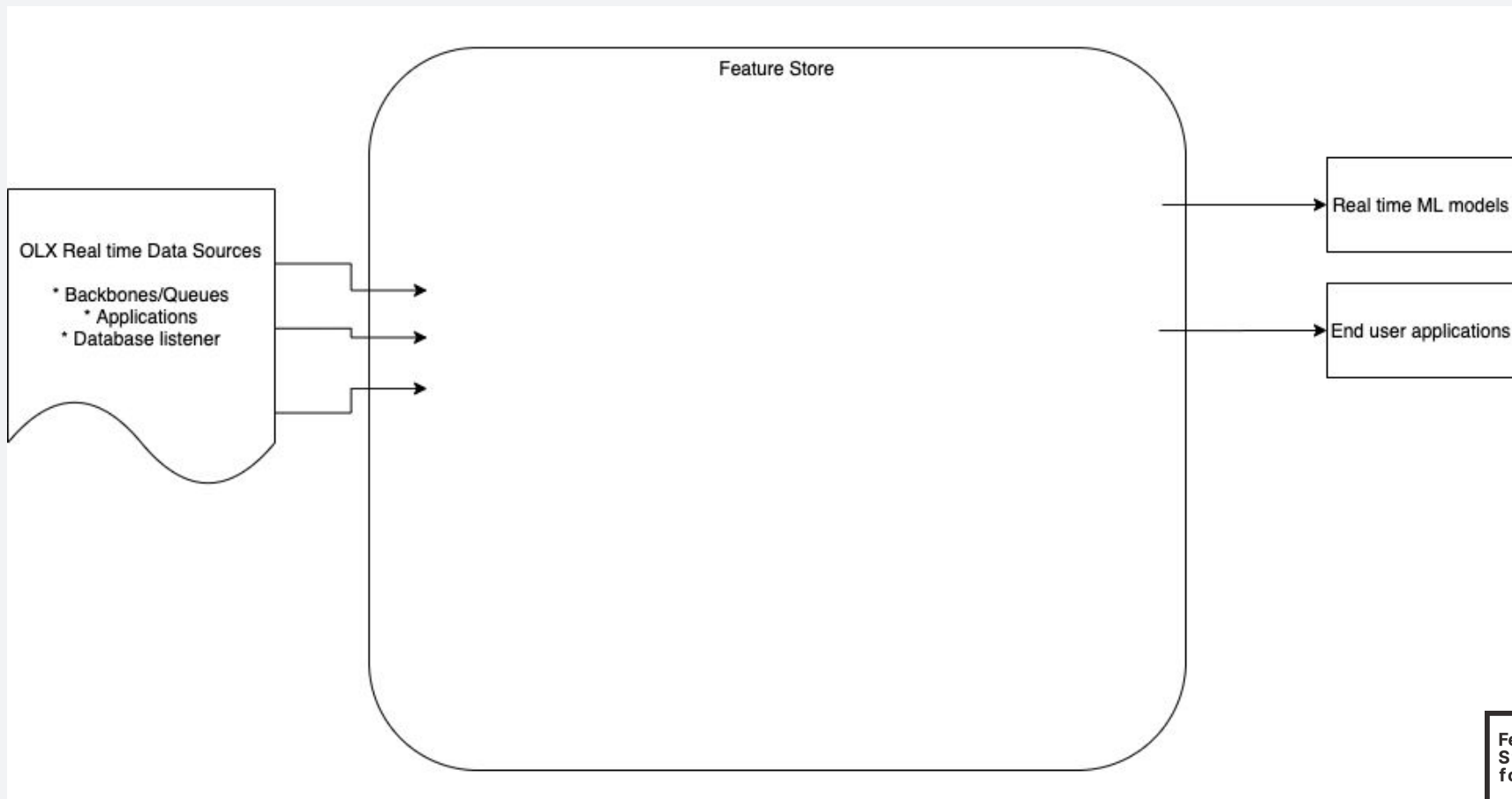
This would be the
start of our
Feature Store

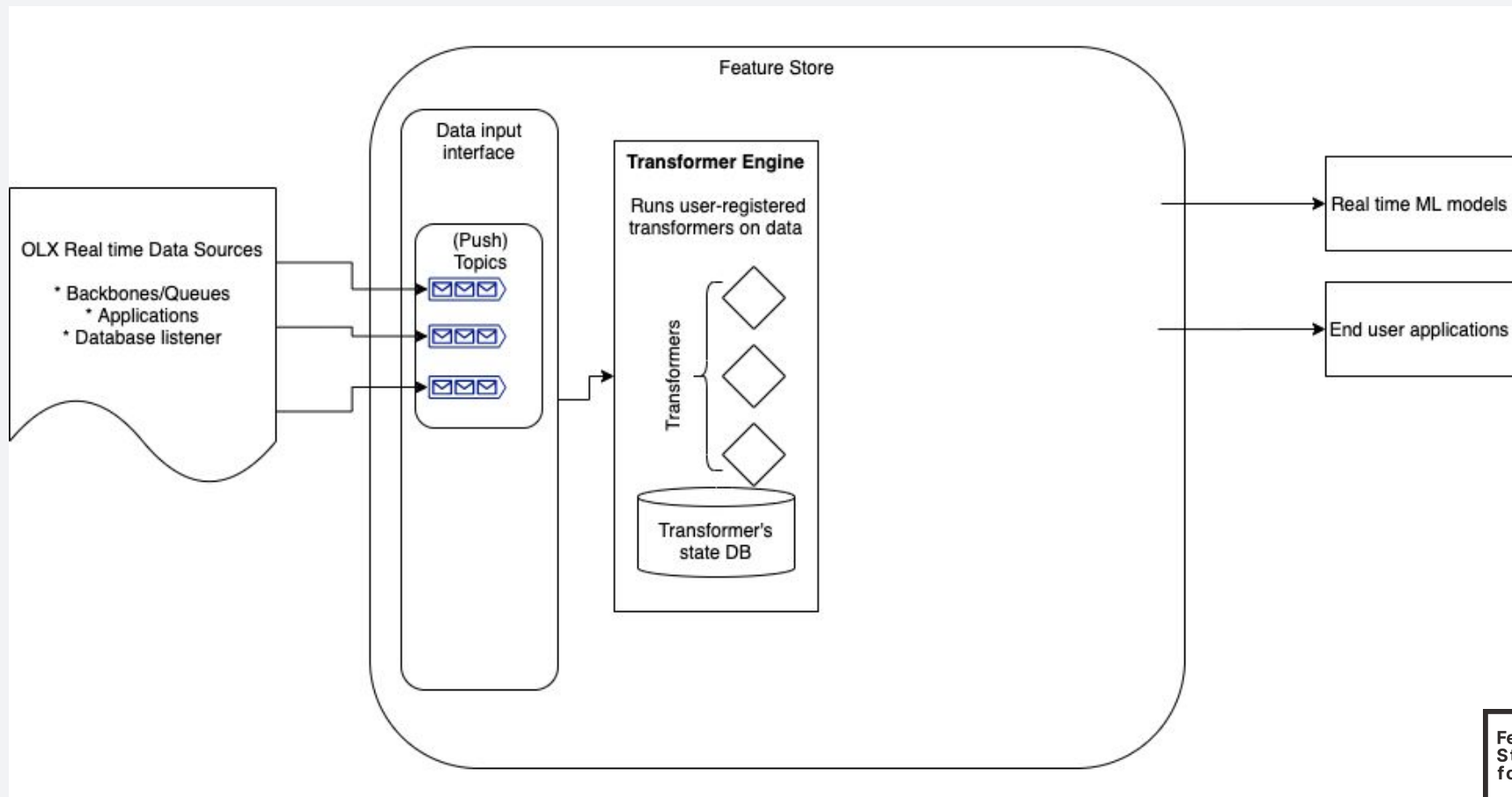
Feature Store

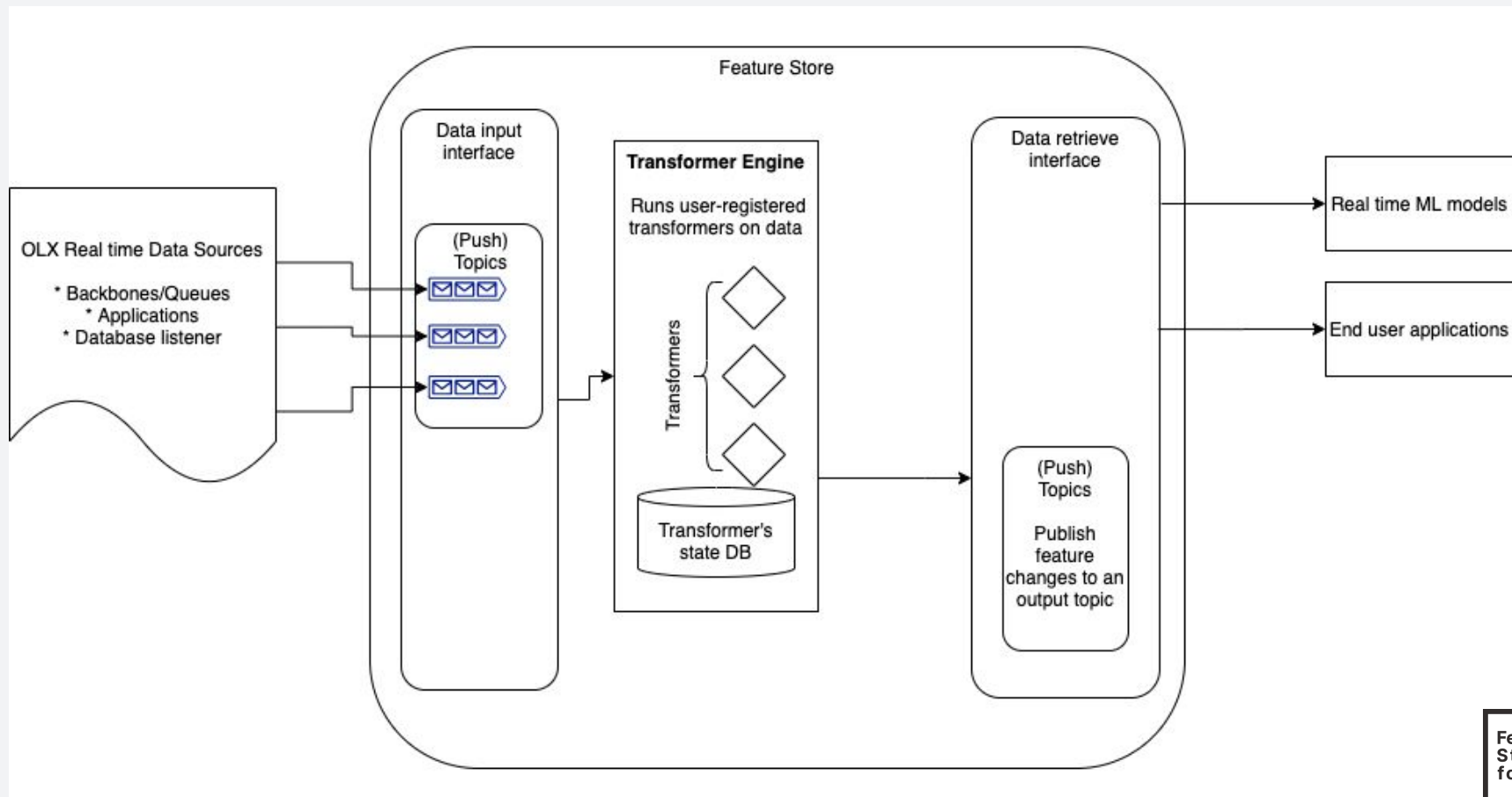


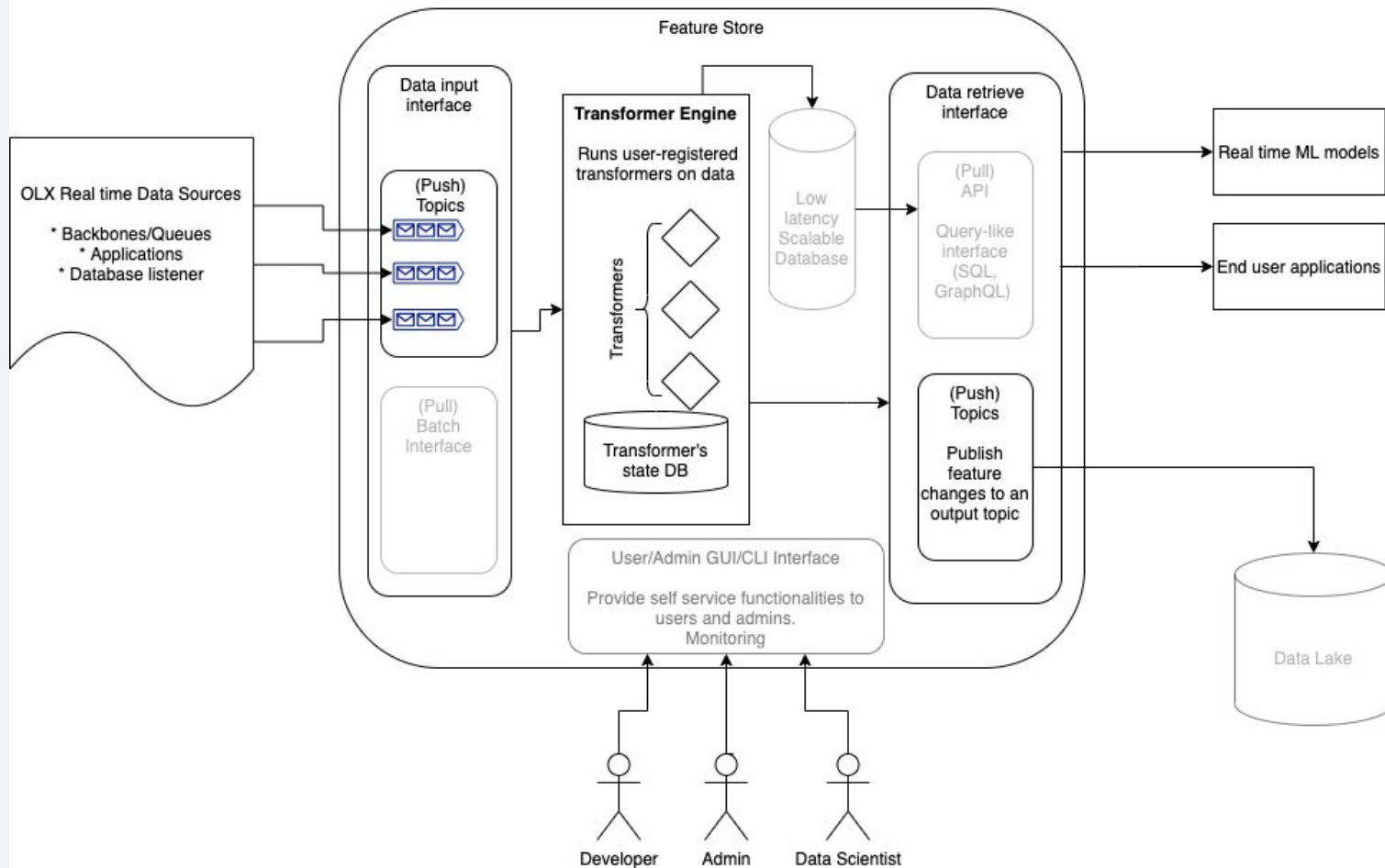
FS Conception

Platform that
centralizes
processment, offers
easy to create
features and
promotes reusability









FS Implementation

Main Concerns:

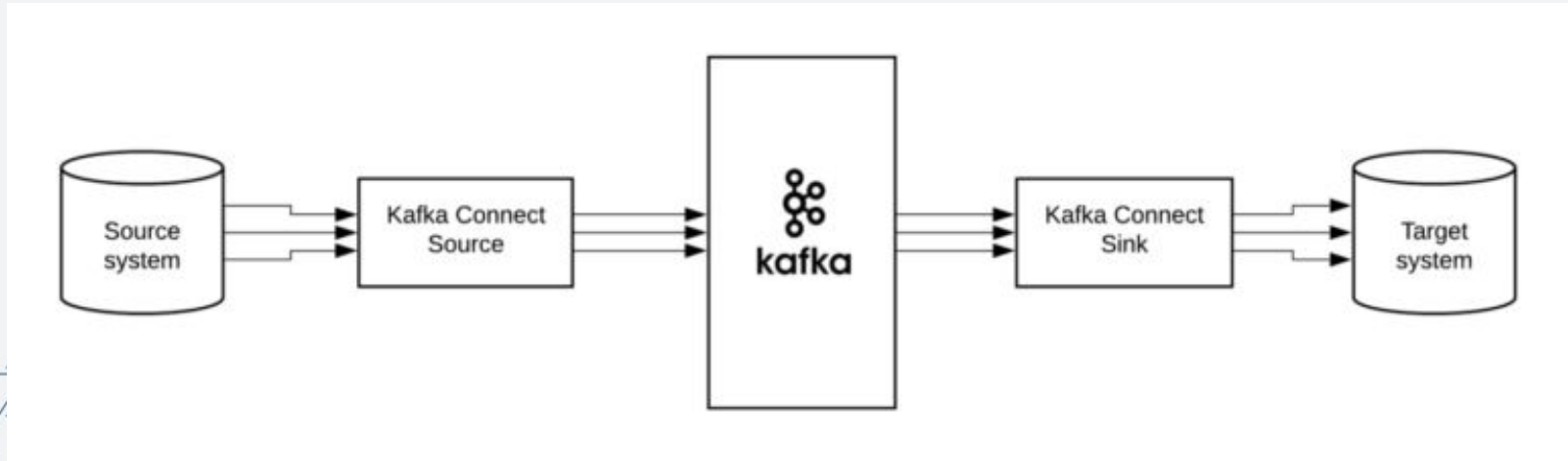
- Must be easy to create features
- Pricing
- Infrastructure Complexity

Messaging System



Kafka Managed
by AWS

Integration with external components

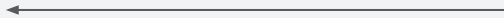
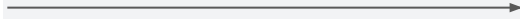


Kafka Connect

What about the transformations?



Raw Data



Aggregations

Discovery with some Transformation Engines

	Flink
Language	Java/SQL
Maturity	Good
Documentation	Bad
Feature Development	Regular
Beginner-friendly	Bad
Scalability	Good
Cost	Bad
Flexibility	Good

	Flink	PyFlink
Language	Java/SQL	Python
Maturity	Good	Regular
Documentation	Bad	Bad
Feature Development	Regular	Regular
Beginner-friendly	Bad	Bad
Scalability	Good	Good
Cost	Bad	Bad
Flexibility	Good	Good

	Flink	PyFlink	Kafka Streams
Language	Java/SQL	Python	Java
Maturity	Good	Regular	Good
Documentation	Bad	Bad	Good
Feature Development	Regular	Regular	Regular
Beginner-friendly	Bad	Bad	Bad
Scalability	Good	Good	Good
Cost	Bad	Bad	Regular
Flexibility	Good	Good	Good

	Flink	PyFlink	Kafka Streams	Faust
Language	Java/SQL	Python	Java	Python
Maturity	Good	Regular	Good	Bad
Documentation	Bad	Bad	Good	Bad
Feature Development	Regular	Regular	Regular	Good
Beginner-friendly	Bad	Bad	Bad	Regular
Scalability	Good	Good	Good	Bad
Cost	Bad	Bad	Regular	Regular
Flexibility	Good	Good	Good	Regular

	Flink	PyFlink	Kafka Streams	Faust	KSQLdb
Language	Java/SQL	Python	Java	Python	SQL
Maturity	Good	Regular	Good	Bad	Good
Documentation	Bad	Bad	Good	Bad	Good
Feature Development	Regular	Regular	Regular	Good	Good
Beginner-friendly	Bad	Bad	Bad	Regular	Good
Scalability	Good	Good	Good	Bad	Good
Cost	Bad	Bad	Regular	Regular	Regular
Flexibility	Good	Good	Good	Regular	Regular

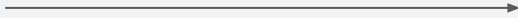
Language
Maturity
Documentation
Feature Development
Beginner-friendly
Scalability
Cost
Flexibility

KSQLdb
SQL
Good
Good
Good
Good
Good
Regular
Regular

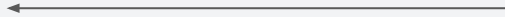
```
'''
CREATE OR REPLACE STREAM test_stateless (
  adId VARCHAR,
  chatId VARCHAR
) WITH (
  KAFKA_TOPIC='test_topic',
  VALUE_FORMAT='JSON',
  PARTITIONS=1
)
''',
'''
CREATE OR REPLACE TABLE test_stateful as
SELECT adId, COUNT(*)
FROM test_stateless
GROUP BY adId
EMIT CHANGES
'''
```

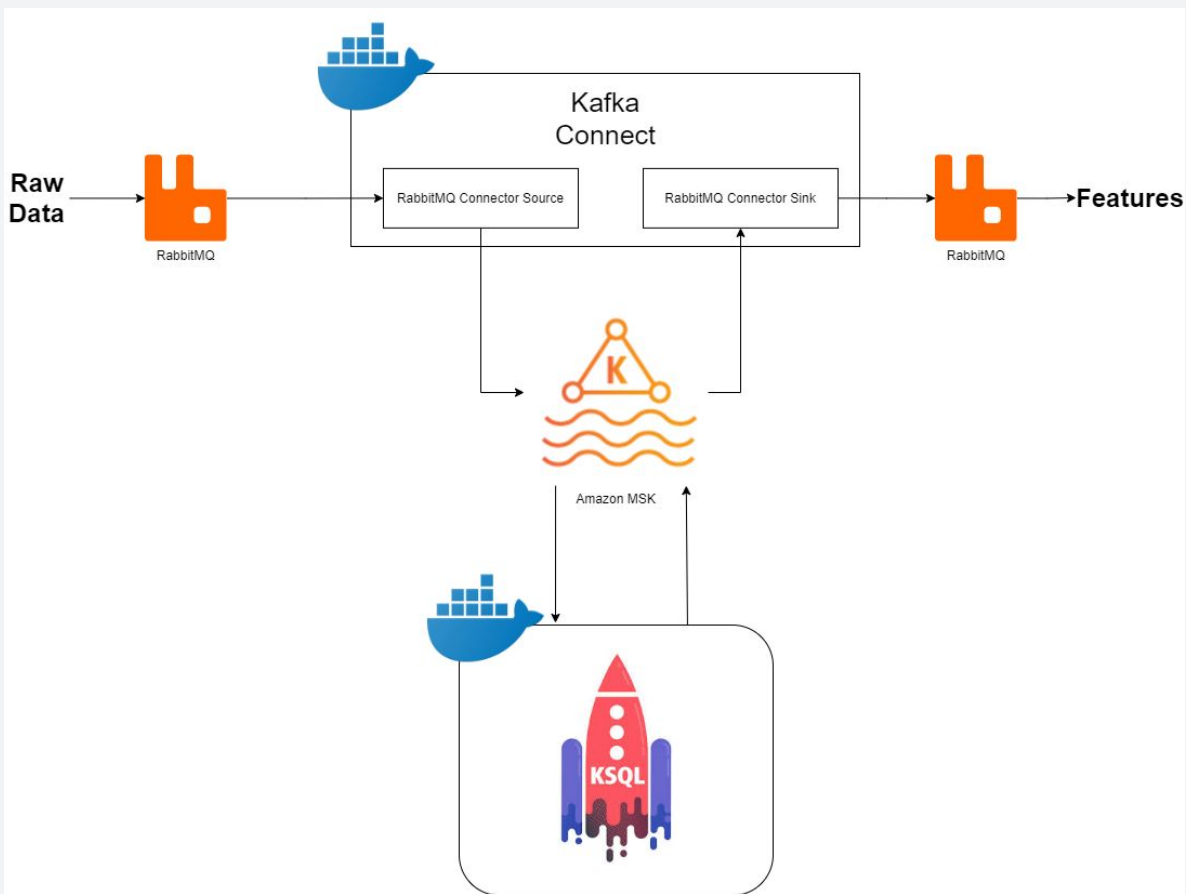


Raw Data



Aggregations





Example of usage

We can create a stream to get the Raw Data

```
1 CREATE STREAM chatmod_message ( --noqa
2     userIp VARCHAR,
3     tipVisibleTo VARCHAR,
4     tipValue VARCHAR,
5     tipType VARCHAR,
6     textMessage VARCHAR,
7     senderAccountId VARCHAR,
8     receiverAccountId VARCHAR,
9     publicId VARCHAR,
10    platform VARCHAR,
11    messageType VARCHAR,
12    messageTimestamp VARCHAR,
13    messageId VARCHAR,
14    listId VARCHAR,
15    chatId VARCHAR,
16    adId VARCHAR,
17    senderId VARCHAR,
18    senderType VARCHAR,
19    isFirstMessage BOOLEAN
20 ) WITH (
21     KAFKA_TOPIC = 'chatmod_message',
22     VALUE_FORMAT = 'json',
23     PARTITIONS = 1,
24     TIMESTAMP = 'messageTimestamp',
25     TIMESTAMP_FORMAT = 'yyyy-MM-dd HH:mm:ss.SSS'
26 );
```

```
1 CREATE TABLE chatmod_first_message_similarity_feature WITH (KAFKA_TOPIC = 'chatmod_first_message_similarity_feature') AS --noqa
2 SELECT
3     buyer_id AS buyer_id_key,
4     AS_VALUE(buyer_id) AS buyer_id,
5     MIN_LEVENSHTEIN_DISTANCE(textMessage) AS first_message_similarity
6 FROM chatmod_buyer_message WINDOW SESSION (2 HOURS) --noqa
7 WHERE isFirstMessage = true
8 GROUP BY buyer_id EMIT CHANGES; --noqa
```

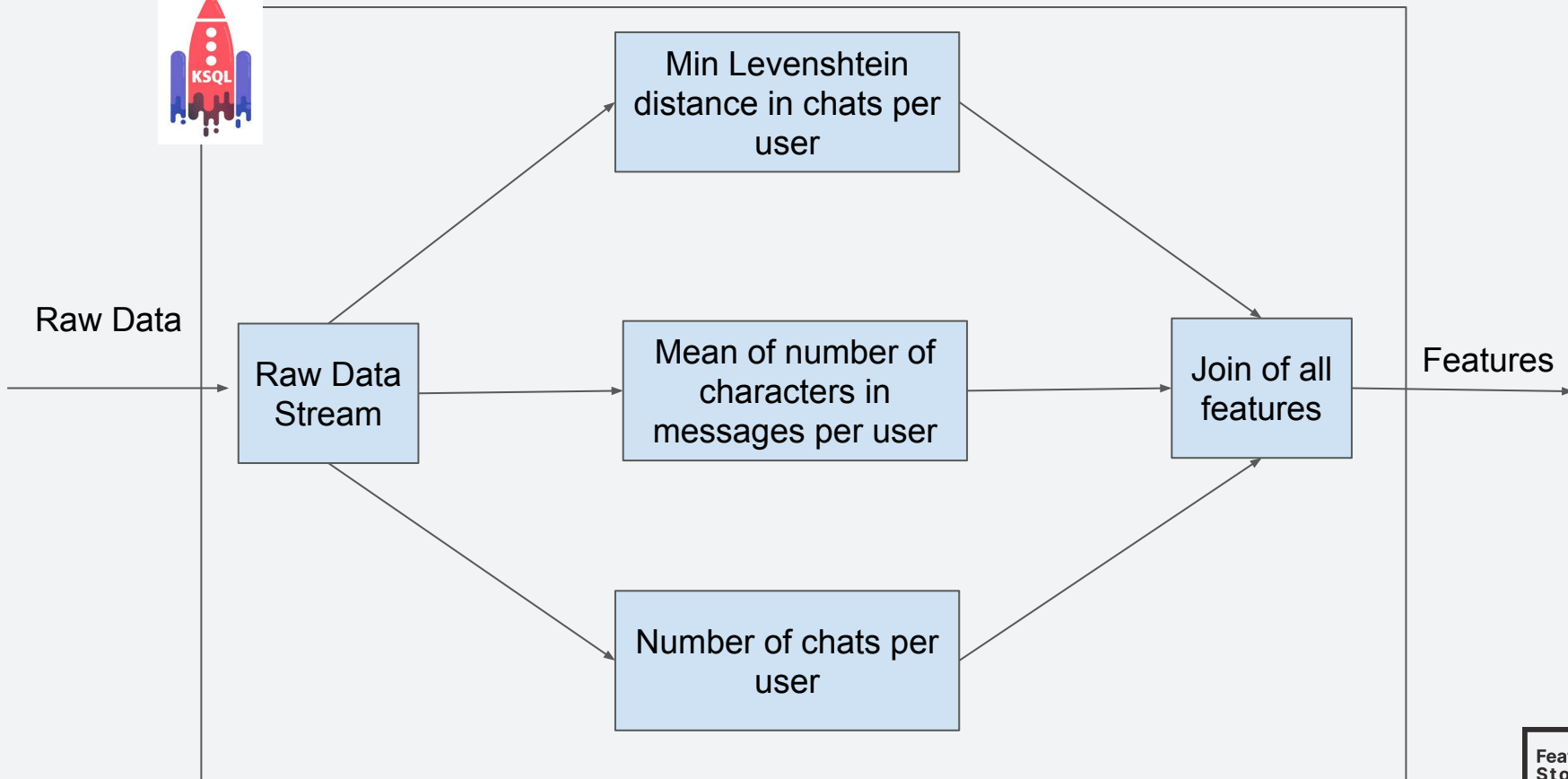
We can create a TABLE
that aggregates data

User Defined Function

```
1 CREATE TABLE chatmod_first_message_similarity_feature WITH (KAFKA_TOPIC = 'chatmod_first_message_similarity_feature') AS --noqa
2 SELECT
3     buyer_id AS buyer_id_key,
4     AS_VALUE(buyer_id) AS buyer_id,
5     MIN_LEVENSHTAIN_DISTANCE(textMessage) AS first_message_similarity
6 FROM chatmod_buyer_message_window_session (2 HOURS) --noqa
7 WHERE isFirstMessage = true
8 GROUP BY buyer_id EMIT CHANGES; --noqa
9
```

We can generate a stream
that reads from this Table
and saves into a topic

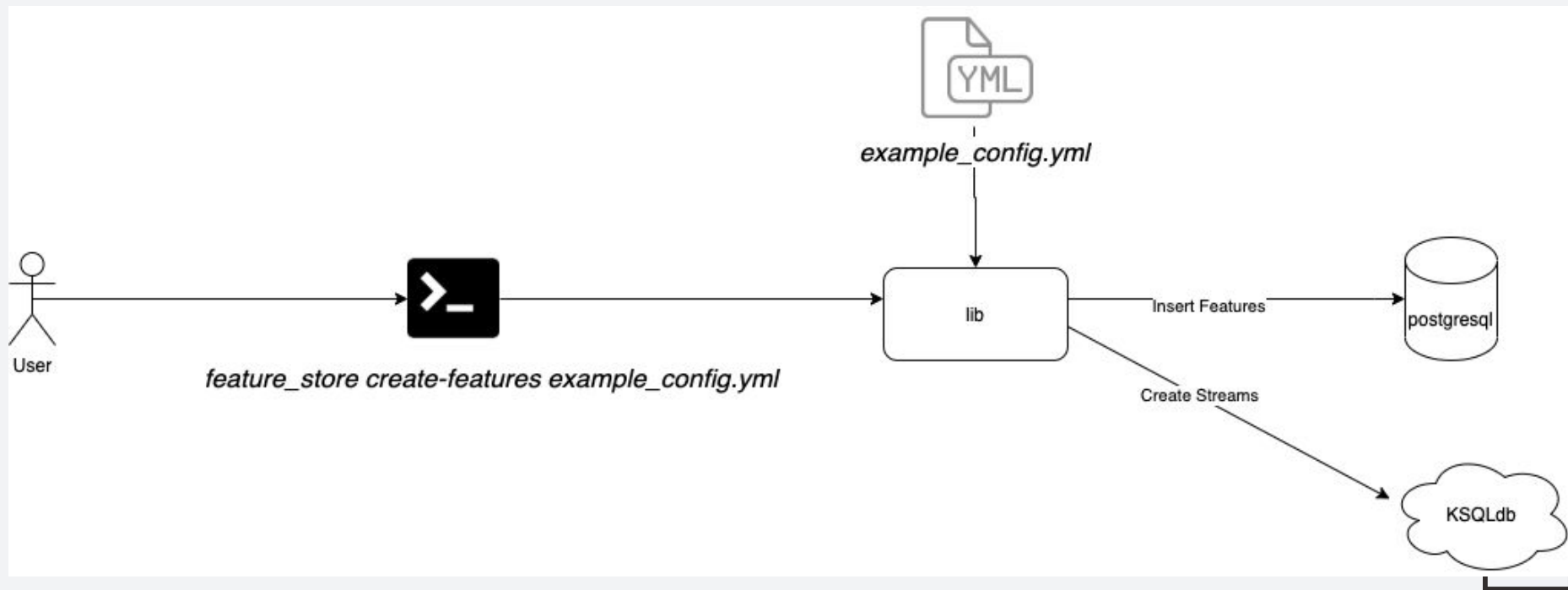
```
9
10 CREATE OR REPLACE STREAM chatmod_first_message_similarity_feature_windowed ( --noqa
11     buyer_id VARCHAR,
12     buyer_id_key VARCHAR KEY,
13     first_message_similarity INTEGER
14 ) WITH (
15     KAFKA_TOPIC='chatmod_first_message_similarity_feature',
16     VALUE_FORMAT='JSON',
17     WINDOW_TYPE='SESSION'
18 );
```



Declaring Features

```
1 product:
2   name: test_name_6
3   team: test_team
4   description: This is the product description
5
6 connectors:
7   connector_1:
8     name: source_connector
9     connector_type: source
10    topic_name: topic_source_connector
11   connector_2:
12     name: sink_connector
13     connector_type: sink
14     topic_name: topic_sink_connector
15
16 features:
17   feature_1:
18     name: first_feature_6
19     description: This is my first feature
20     feature_streams:
21     - stream_1
22     - stream_2
23     sql_file: example_sql_scripts/file_1.sql
24     input_topic: topic_1
25     output_topic: topic_2
26     is_external: True
27   feature_2:
28     name: second_feature_6
29     description: This is my second feature
30     feature_streams:
31     - stream_3
32     - stream_4
33     sql_file: example_sql_scripts/file_2.sql
34     input_topic: topic_2
35     output_topic: topic_3
36     is_external: True
37     dependencies:
38     - feature_1
39
```

Feature Declaration Interface



Results

**Messages Input per
day**

5 M

Features

8 created

**Feature Output per
day**

1.5 M

**Models using the
Feature Store**

2 models

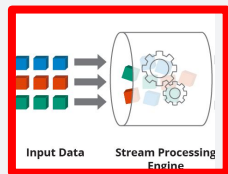
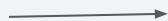
**Fraudulent messages
detected per day**

300-400

**Complaint reduced
by**

22-26%

Roadmap



Monitoring

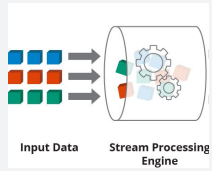
Storage

Serving

Registry

Creation of features in batch,
from datalake

Monitoring



Storage

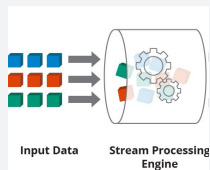
Serving

Registry

Feature Versioning

Feature Governance

Monitoring



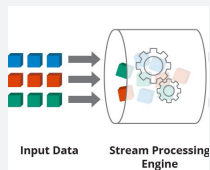
Storage

Serving

Registry

Export features into Datalake,
for training

Monitoring



Storage

Serving

Registry

Feature serving with database
and API



Thank you!

Do you have any questions?

Augusto Acioli Pinho Vanderley - OLX Brasil



https://www.linkedin.com/in/augusto-vanderley/?locale=en_US

