

Real-time machine learning apps with NoSQL

Attila Toth
Developer Advocate at ScyllaDB

Attila Toth



- + Developer advocate at ScyllaDB
- + Working as a software engineer / dev advocate in the data space for 6+ years
- + Lives in Budapest, Hungary

Agenda

- What is a feature in ML?
- About feature stores
- NoSQL as a Feature Store
- DEMO

What is a feature?

Post

Post Features

post_likes_10m
post_shares_7d

post_topic

Cross Features

user_creator_likes_30d
user_creator_dwell_time_7d



Creator

Creator Features

creator_comments_30d
creator_views_1h

creator_state

User

User Features

user_time_spent_7d
user_video_plays_30m

user_dob

What is a feature?

Individual measurable property

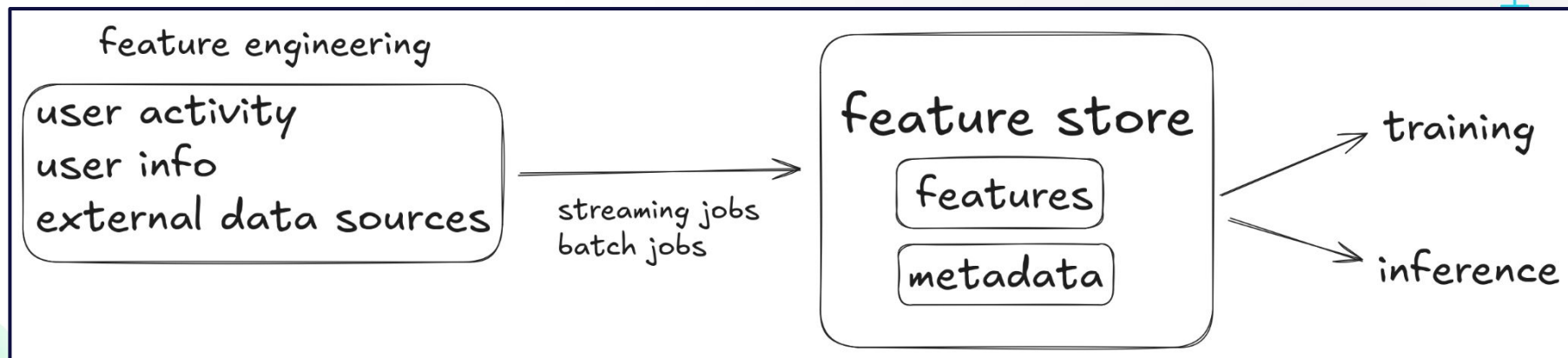
Any data point that is used to train models

Example feature vector (collection of related features):

zipcode	person_age	person_income	loan_amount	loan_int_rate
94109	25	120000	10000	12

What is a feature store?

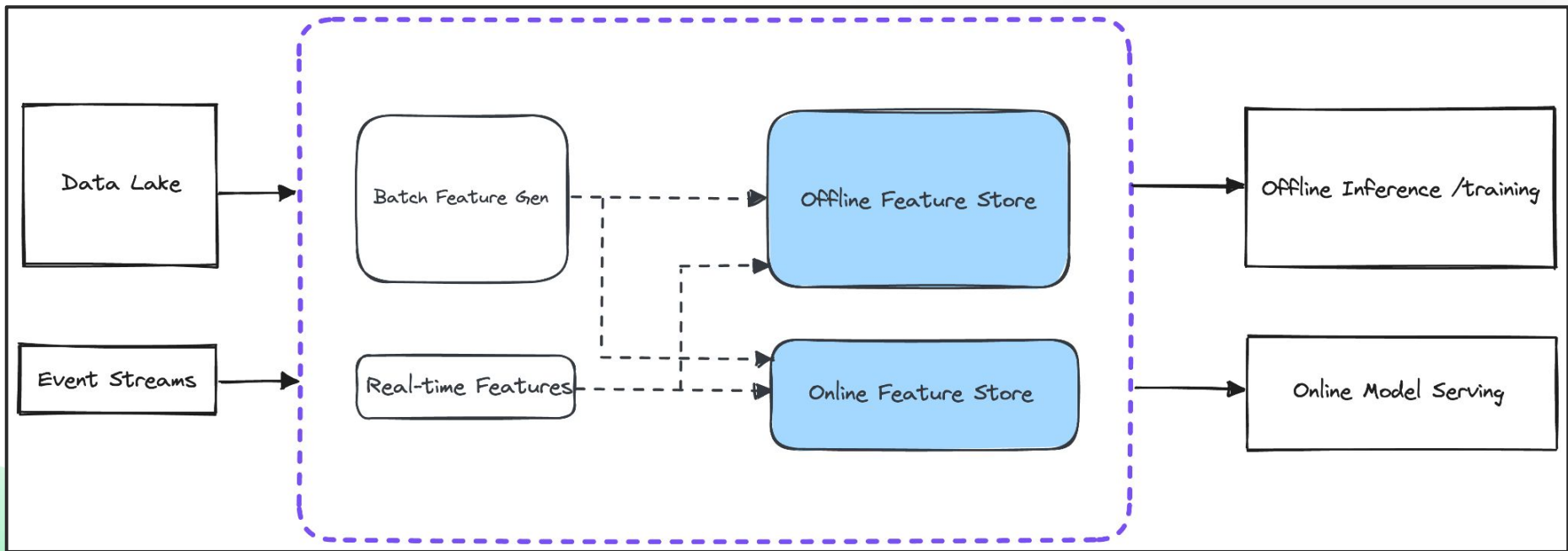
Feature store is a central database to **create, store, discover, and access features** for model training and inference.



Benefits of a feature store

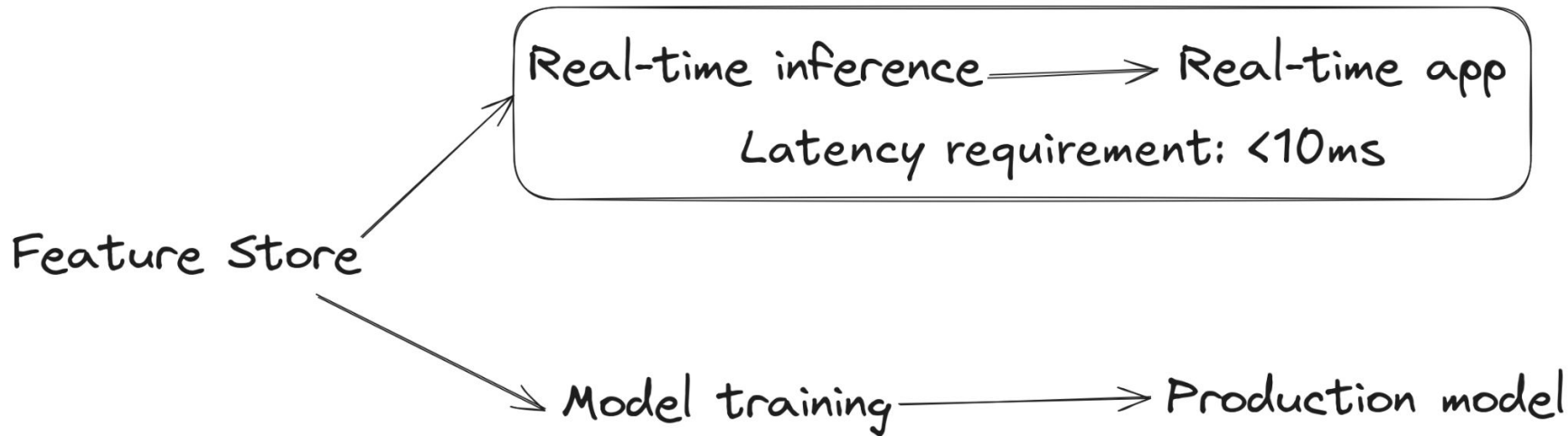
- Reusing features (avoiding duplicated effort)
- Sharing features between teams
- Data scientists can create and access features without sync with engineering teams
- Unifying the training and online inference feature pipelines
- Consolidate different ML workloads into one system
- etc...

Feature store architecture (example)





Different needs & workloads



Different needs & workloads

Online store (real-time)

- Focused on now/today
- Short-term retention with TTL
- For fast rapid retrieval of specific data
- Focus on real-time live site requirements

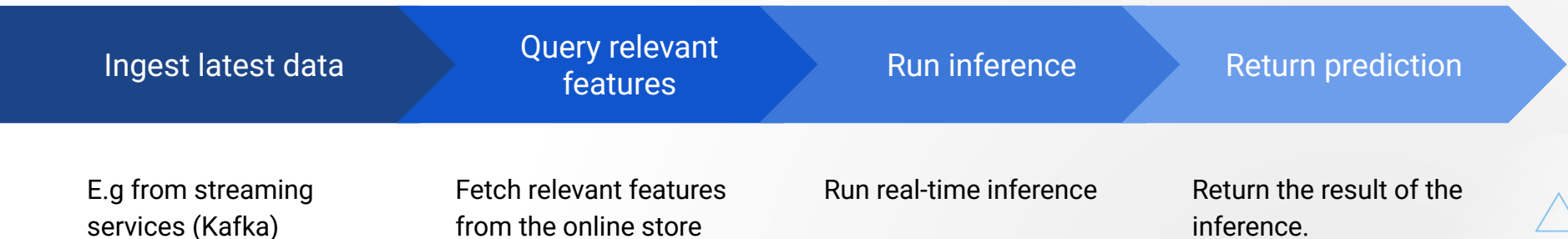
Offline store

- Huge date range (historical data)
- For ML model training
- Long-term storage
- Less latency-sensitive
- Multipurpose

Latency in real-time ML apps

Database latency → The time it takes from receiving a query to returning the result

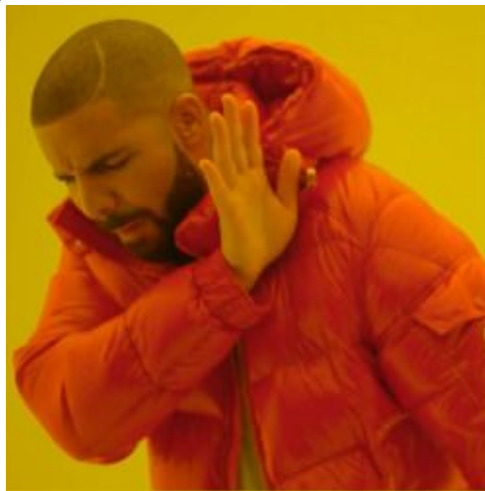
- Online stores power user-facing, time-sensitive decisions
 - E.g. recommendation engine, fraud detection, pricing etc



latency requirement (real-time apps):
~~seconds~~ milliseconds

Feature store solutions

- **Open Source frameworks:**
 - Hopsworks
 - Feast
- **Major Cloud providers:**
 - AWS SageMaker, GCP Vertex
- **Other commercial solutions:**
 - Databricks
 - Tecton
 - Etc



use the
feature store's
built-in database
as online store



use whatever
database
I want as
online store

Others

Platform	Open-Source	Offline	Online	Real-Time Ingestion
Hopsworks	AGPL-V3	Hudi/Hive and pluggable	RonDB	Flink, Spark Streaming
Iguazio Data Science	No	Parquet	V3IO, proprietary DB	Nuclio
Databricks	No	Delta Lake	Mysql or Aurora	None
Amazon SageMaker	No	S3, Parquet	DynamoDB	None
Featureform	Mozilla	Pluggable	Pluggable	None
Feathr	Yes	Pluggable	Pluggable	None
Feast	Yes	Pluggable	Pluggable	Sync'd from Offline Store
Qwak	No	S3	Redis	Computed at Request time
Azure	No	Delta Lake?	CosmosDB?	??
Chalk	No	Snowflake, BigQuery, R...	Redis, DynamoDB, Bigta...	Strreaming, Batch

<https://airtable.com/appgjFbGAL0btUNmQ/shr6oH1eXkGfOeghl/tblgEP2mglcsW6tZ7?viewControls=on>
<https://www.featurestore.org/>

NoSQL as an Online Feature Store

About ScyllaDB

- NoSQL database, drop-in replacement for Cassandra and DynamoDB
- Superior performance over legacy NoSQL
 - >5x higher throughput at same cost
 - >20x lower latency
 - >75% TCO savings
- Avoids vendor lock-in
 - ScyllaDB can run anywhere
 - Also available as DBaaS on AWS and GCP
 - Open-Source version available
- ScyllaDB is a great choice if you require...
 - High availability
 - Low latency
 - High throughput
 - "Big data" scale



ScyllaDB is Key-Value



NoSQL as an online feature store: Benefits

- Low latency (<5ms P99)
- High throughput (millions of operations per second)
- Large scale (petabytes of data)
- High availability (the system remains operational even if a node goes down)
- Easy migration from legacy systems (Cassandra, DynamoDB, MongoDB, Redis, etc)
- Integrations (Cassandra and DynamoDB compatibility, ingesting from Kafka, CDC, Feast integration etc)
- Self-hosted (anywhere) or managed service with ScyllaDB Cloud (AWS or GCP)

Feature Store

DEMO

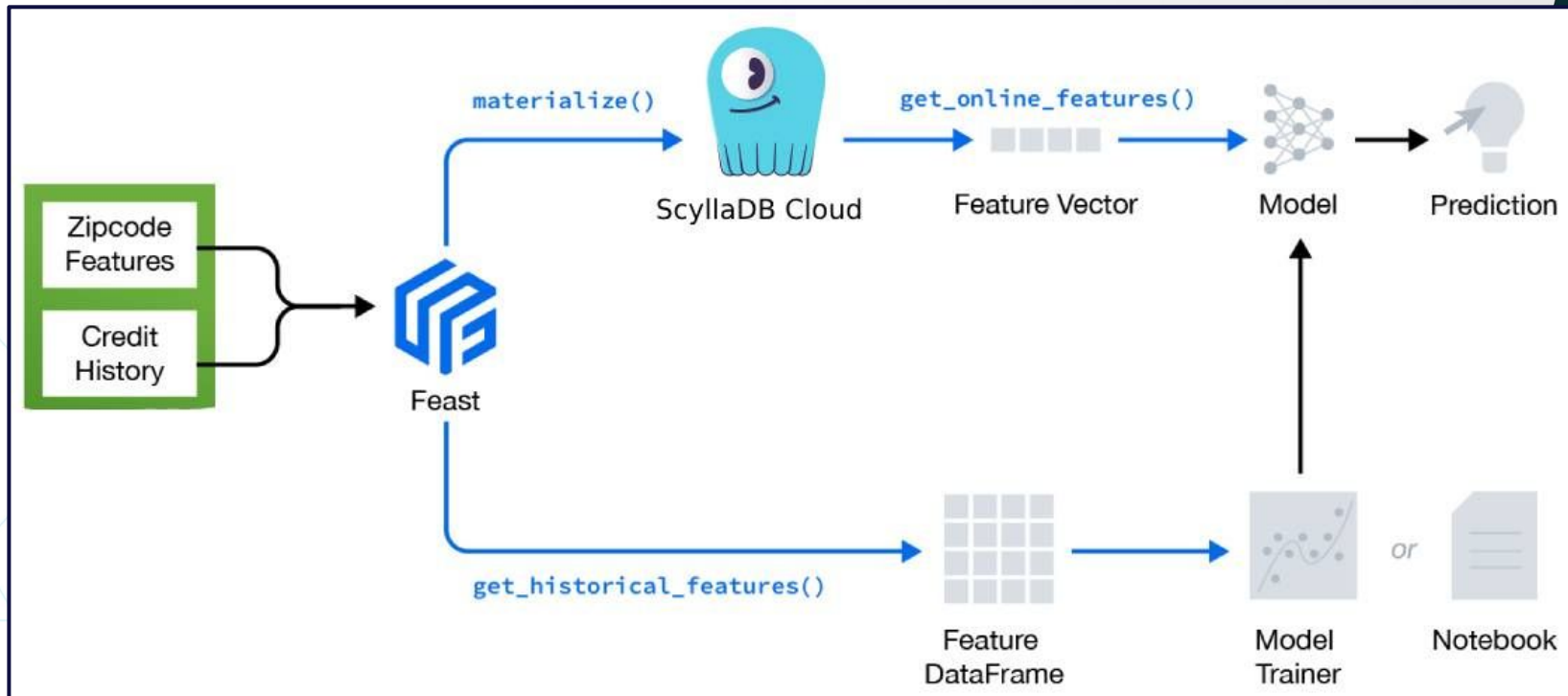
DEMO God

AI Prompt:

“Create an image of a God who governs the success or failure of live software demos”



Feast example



<https://github.com/scylladb/scylladb-feature-store>

Feast + ScyllaDB configuration

```
project: repo
# By default, the registry is a file (but can be turned into a more scalable SQL-backed registry)
registry: data/registry.db
# The provider primarily specifies default offline / online stores & storing the registry in a given cloud
provider: local
online_store:
  type: cassandra
  hosts:
    - 172.17.0.2
  username: scylla
  password:
  keyspace: feast
entity_key_serialization_version: 2
```

Thank you!

Email:
ATTILA.TOTH@SCYLLADB.COM

Slides (and more):
AttilaToth.dev/Budapest-Data-2025

