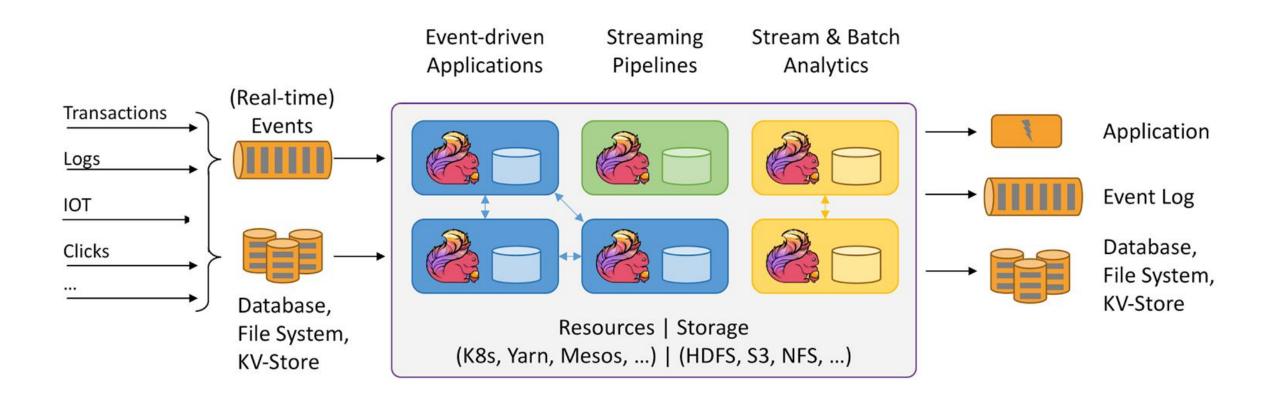# Introduction to SQL on Apache Flink®

Flink SQL Training

https://github.com/ververica/sql-training

# Apache Flink is a Distributed Data Processing System

# Scalable and Consistent Data Processing

- Flexible and expressive APIs

- Guaranteed correctness
  - Exactly-once state consistency
  - Event-time semantics

- In-memory processing at massive scale
  - Runs on 10000s of cores
  - Manages 10s TBs of state

# Powered By Apache Flink



Details about their use cases and more users are listed on Flink's website at https://flink.apache.org/poweredby.html
Also check out the Flink Forward YouTube channel with more than 350 recorded talks at https://www.youtube.com/channel/UCY8_lgiZLZErZPF47a2hXMA

# Why SQL for Stream Processing?

- Implementing Flink stream processing apps requires special skills
  - Java/Scala experience
  - In-depth knowledge of streaming concepts like time and state
  - Knowledge of distributed data processing

- Everybody knows and uses SQL

- SQL queries are optimized and efficiently executed

- Unified syntax and semantics for batch & streaming data

# Flink SQL in a Nutshell

*A standard-compliant SQL service*

*to query static and streaming data alike*

*that leverages the performance, scalability, and consistency of*

*Apache Flink.*

# How is streaming SQL different from traditional SQL?

- Basically all tables that are processed with SQL queries change over time
  - Transactions from applications
  - Bulk inserts from ETL processes

- Traditional processors run SQL queries on static snapshots of the tables
  - The query input is finite
  - The query result is final and finite

- Stream SQL processors run continuous queries on changing (dynamic) tables
  - The query input is unbounded
  - The query result is never final, continuously updated, and potentially unbounded

- The semantics of a query are the same regardless whether it is executed one-time on a table snapshot or continuously on a changing table

# Running a One-time Query on a Changing Table

Take a snapshot when the query starts

A final result is produced

| user | cTime | url |
|------|-------|-----|
| Mary | 12:00:00 | https://... |
| Bob | 12:00:00 | https://... |
| Mary | 12:00:02 | https://... |
| Liz | 12:00:03 | https://... |

```
SELECT
    user,
    COUNT(url) as cnt
FROM clicks
GROUP BY user
```

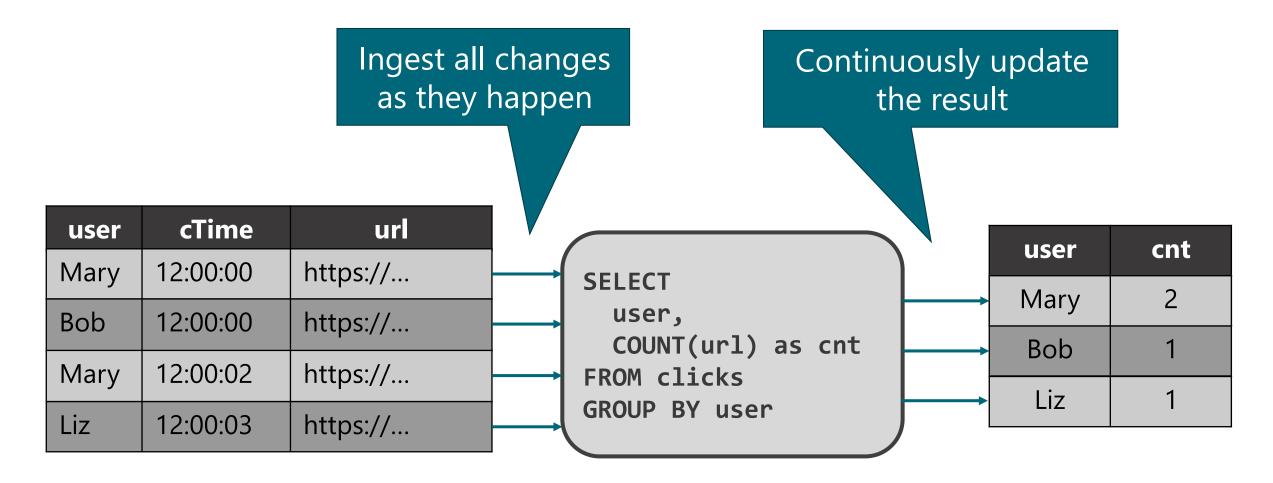| user | cnt |
|------|-----|
| Mary | 2 |
| Bob | 1 |

The query terminates

A row that was added after the query was started is not considered

# Running a Continuous Query on a Changing Table

Ingest all changes as they happen

Continuously update the result

| user | cTime | url |
|------|-------|-----|
| Mary | 12:00:00 | https://... |
| Bob | 12:00:00 | https://... |
| Mary | 12:00:02 | https://... |
| Liz | 12:00:03 | https://... |

```
SELECT
    user,
    COUNT(url) as cnt
FROM clicks
GROUP BY user
```

| user | cnt |
|------|-----|
| Mary | 2 |
| Bob | 1 |
| Liz | 1 |

The result is identical to the one-time query (at this point)

# SQL Feature Set in Flink 1.10

## STREAMING & BATCH

- SELECT FROM WHERE
- GROUP BY [HAVING]
  - Non-windowed
  - TUMBLE, HOP, SESSION windows
- JOIN
  - Time-Windowed INNER + OUTER JOIN
  - Non-windowed INNER + OUTER JOIN
- User-Defined Functions
  - Scalar
  - Aggregation
  - Table-valued

## STREAMING ONLY

- OVER / WINDOW
  - UNBOUNDED + BOUNDED PRECEDING
- INNER JOIN with
  - Time-versioned table
  - External lookup table
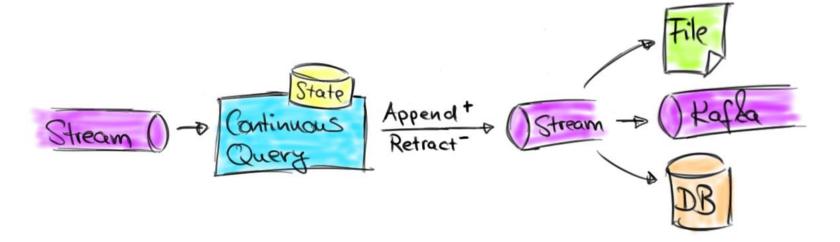- MATCH_RECOGNIZE
  - Pattern Matching/CEP (SQL:2016)

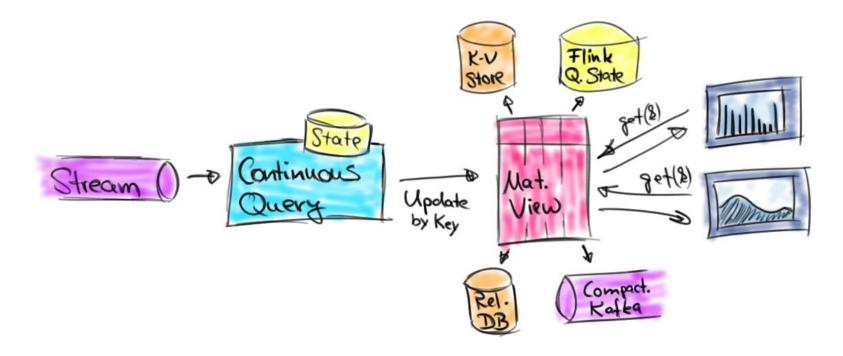## BATCH ONLY

- Full TPC-DS support

# Data Pipelines

- Transform, aggregate, and move events in real-time

- Low-latency ETL
    - Convert and write streams to file systems, DBMS, K-V stores, indexes, ...
    - Ingest appearing files to produce streams

# Stream & Batch Analytics

- Stream & Batch Analytics
  - o Run analytical queries over bounded and unbounded data
  - o Query and compare historic and real-time data
  - o Compute and update data to visualize in real-time

# Training Environment

https://github.com/ververica/sql-training/

# What You Will Learn in This Training?

- Querying streaming data with SQL

- Expressing common stream processing operations with SQL
  - Window aggregations, stream joins, and pattern matching

- Piping the results of continuous queries into Kafka

- Materializing the results of continuous queries in MySQL

- Using Flink's SQL CLI client

# Training Scenario: Taxi Ride Data

- We are working with data about taxi rides in New York

- Three tables
  - `Rides`           One start and one end event for each ride
  - `Fares`           One payment event for each ride
  - `DriverChanges`   One event for each driver change of a taxi

  - All tables are registered and available in the environment

  - Each tables is backed by a Kafka topic

# Training Scenario: Taxi Ride Data

```
Flink SQL> SELECT * FROM Rides;
```

| rideId | taxiId | isStart | lon | lat | rideTime | psgCnt |
|---|---|---|---|---|---|---|
| 1 | 2013000001 | true | -73.99078 | 40.76088 | 2013-01-01T00:00 | 1 |
| 2 | 2013000002 | true | -73.978325 | 40.77809 | 2013-01-01T00:00 | 5 |
| 3 | 2013000003 | true | -73.98962 | 40.72999 | 2013-01-01T00:00 | 1 |

```
Flink SQL> SELECT * FROM Fares;
```

| rideId | payTime | payMethod | tip | toll | fare |
|---|---|---|---|---|---|
| 65 | 2013-01-01T00:00:36 | CSH | 0.0 | 0.0 | 3.5 |
| 137 | 2013-01-01T00:01 | CSH | 0.0 | 0.0 | 3.5 |
| 77 | 2013-01-01T00:01:22 | CSH | 0.0 | 0.0 | 4.0 |

```
Flink SQL> SELECT * FROM DriverChanges;
```

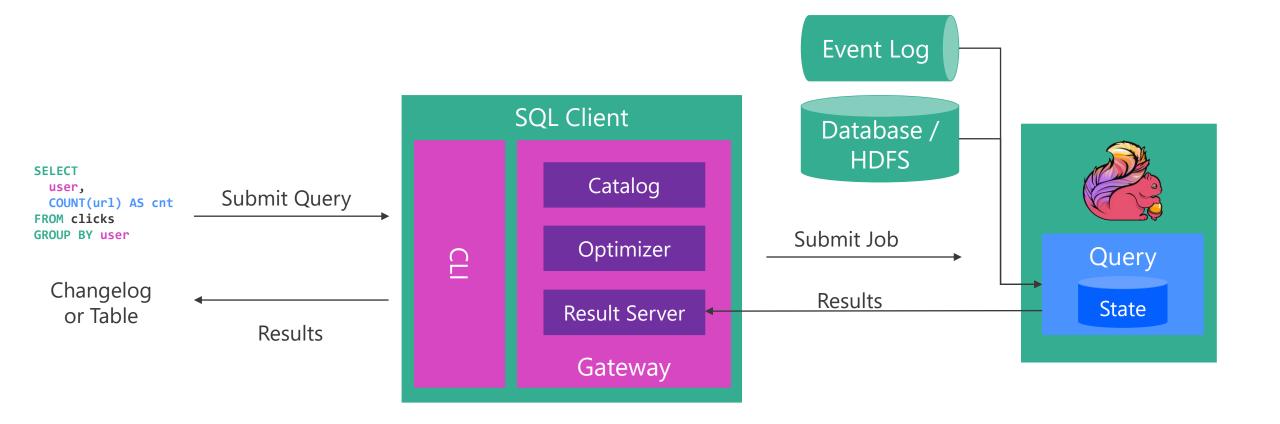| taxiId | driverId | usageStartTime |
|---|---|---|
| 2013000061 | 2013000061 | 2013-01-01T00:00:02 |
| 2013000062 | 2013000062 | 2013-01-01T00:00:03 |
| 2013000063 | 2013000063 | 2013-01-01T00:00:08 |

# Our Training Environment



SQL Client — **Submit** query → JobManager ⟷ **Assign & monitor query tasks** → TaskManager

**Execute query tasks**

**WebUI: http://localhost:8081**

**Push events at 10x speed** → kafka

**Read & write data** → MySQL

kafka ⟷ **Coordinate** → APACHE ZooKeeper

© 2019 Ververica

# Introduction to SQL Client

# Interactive Query Submission via SQL Client



```
SELECT
  user,
  COUNT(url) AS cnt
FROM clicks
GROUP BY user
```

# Detached Query Submission via SQL Client



Event Log

Database / HDFS

## SQL Client

CLI

Catalog

Optimizer

Result Server

Gateway

Query

State

```
INSERT INTO dashboard
SELECT
  user,
  COUNT(url) AS cnt
FROM clicks
GROUP BY user
```

Submit Query

Cluster ID & Job ID

Submit Job

# Hands On Exercises

# Introduction to SQL on Flink

Continue with the "Introduction to the Training Environment"
in "Introduction to SQL on Flink"

https://github.com/ververica/sql-training/wiki/Introduction-to-SQL-on-Flink

We are here to help!

www.ververica.com                    @VerericaData