# Query Live Data with SQL

Why, how, and what's next?
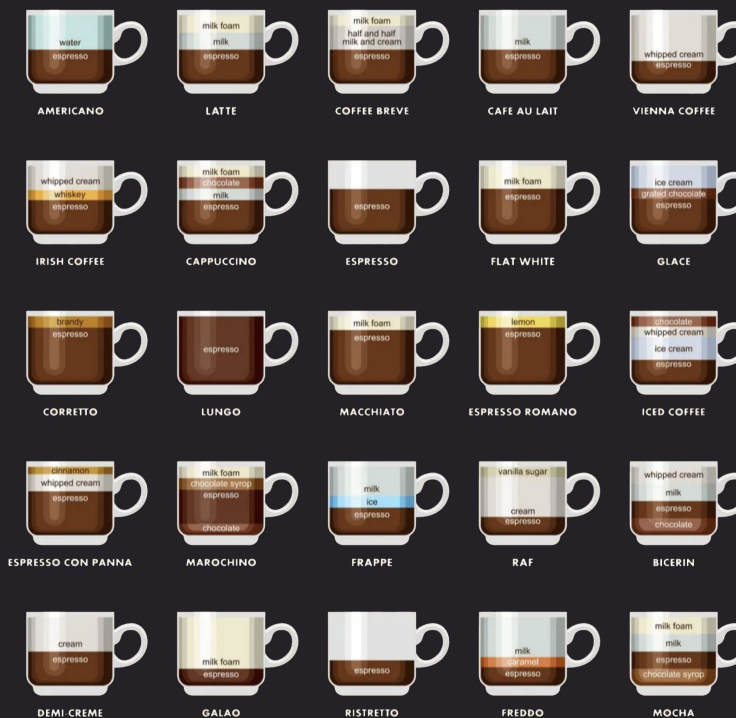
Dec 12, 2023
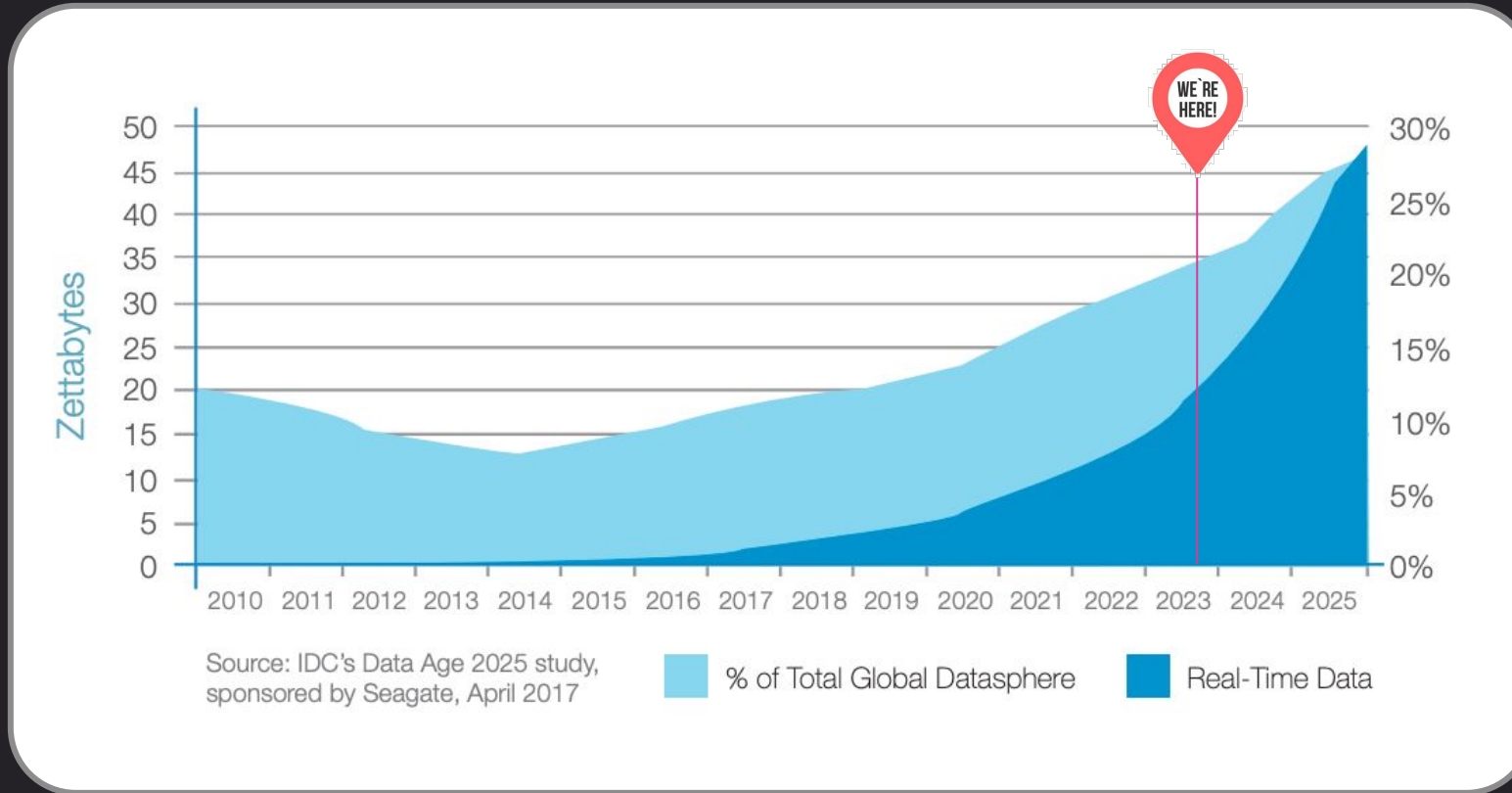
**Jove Zhong**
Co-Founder and Head of Product, Timeplus

**Gang Tao**
Co-Founder and CTO, Timeplus

# Live data is everywhere, at the edge and cloud



**46 ZB**
of data created by billions of IoT by 2025

**30%**
of data generated will be real-time by 2025

Only **1%**
of data is analyzed and streaming data is primarily untapped

# Why SQL For ~~Big~~ *(Live)* Data?



**Reliable**

**Fast**

**Easy**

**Powerful**

**Descriptive**

# Sample Use Cases

## FinTech

- Real-time post-trade analytics
- Real-time pricing

## DevOps

- Real-time Github insights
- Real-time o11y and usage based pricing

## Security Compliance

- SOC2 compliance
- Container vulnerability monitoring
- Monitor Superblocks user activities
- Protect sensitive info in Slack

## IoT

- Real-time fleet monitoring
- Oil well edge monitoring

## Customer 360

- Auth0 notifications for new signups
- HubSpot custom dashboards/alerts
- Jitsu clickstream analytics
- Real-time Twitter marketing

## Misc

- Wildfire monitoring and alerting
- Data-driven parent

Learn more: https://docs.timeplus.com/showcases

timeplus    https://github.com/timeplus-io/proton                    OSA CON 23

# How do you like your coffee?



| | Streaming Processor | Streaming Database | Real-Time Database |
|---|---|---|---|
| Java | Flink | ksqlDB, Hazelcast | Druid, Pinot, Trino |
| C++ | | proton | ClickHouse, StarRocks |
| Rust | Arroyo | RisingWave | Databend |

FlinkSQL
since 2016

espresso

OSA CON 23

**FlinkSQL**
since 2016

**Coffee Tasting Notes** ☕

Community
☕☕☕☕
Real-time     ☕☕☕
Streaming     ☕☕☕
Historical    ☕
JOIN
☕☕☕☕
Largescale
☕☕☕☕
Lightweight ☕☕
Easy to use ☕☕

| | |
|---|---|
| SQL | High-level Language |
| Table API | Declarative DSL |
| DataStream / DataSet API | Core APIs |
| Stateful Stream Processing | Low-level building block (streams, state, [event] time) |

```sql
CREATE TABLE kafka (
  `timestamp` BIGINT,
  `user_id`   STRING,
  `page_id`   STRING,
  `action`    STRING,
  `ts`        TIMESTAMP(3) METADATA FROM 'timestamp'
) WITH (
  'connector'                   = 'kafka',
  'topic'                       = 'demo-stream',
  'properties.bootstrap.servers' = 'localhost:9092',
  'properties.group.id'         = 'testGroup',
  'properties.auto.offset.reset' = 'earliest',
  'scan.startup.mode'           = 'earliest-offset',
  'format'                      = 'json'
);
SELECT * FROM kafka JOIN lookup USING (user_id);
```

**ksqlDB**
since 2019

**Coffee Tasting Notes**

```
Community    ☕☕☕
Real-time    ☕☕☕
Streaming    ☕☕☕
Historical   ☕☕
JOIN         ☕☕☕
Largescale   ☕☕
Lightweight  ☕☕
Easy to use  ☕☕☕
```

**ksqlDB** ⇄ **Kafka**

*Compute*          *Storage*

```
CREATE STREAM githubEvents (
  id VARCHAR,
  created_at VARCHAR,
  actor VARCHAR,
  type VARCHAR,
  repo VARCHAR,
  payload VARCHAR
)
WITH (kafka_topic='github_events', value_format='json');
SELECT * FROM githubEvents WHERE type='CreateEvent';
```

# HAZELCAST

Distributed computation and storage platform

No dependency on disk storage, it keeps all its operational state in the RAM of the cluster.

| Streaming Processor | Streaming Database | Real-time Database |
|---|---|---|
| Flink | ksqlDB / Hazelcast | Druid / Pinot / Trino |

```
CREATE OR REPLACE MAPPING trades (
  id BIGINT,
  ticker VARCHAR,
  price_usd DECIMAL,
  amount BIGINT)
TYPE Kafka
OPTIONS (
  'valueFormat' = 'json-flat',
  'bootstrap.servers'='..:9092',
  'security.protocol'='SASL_SSL',
  'sasl.jaas.config'='org.apache.kafka.common.security.plain.PlainLoginModule
  required username=".."
  password="..";',
  'sasl.mechanism'='PLAIN'
);
SELECT ticker, price_usd, amount FROM trades WHERE price_usd * amount > 100;
```

# pinot

1. create a schema json (columns, PKs)
2. create a table configuration json (streamType=Kafka)
3. docker run .. apachepinot/pinot:latest AddTable \
   -schemaFile /tmp/transcript-schema.json \
   -tableConfigFile /tmp/transcript-table-realtime.json \
      ..
   -exec

# druid

1. load the druid-kafka-indexing-service extension on both the Overlord and the MiddleManagers
2. Create a supervisor-spec.json containing the Kafka supervisor spec file.
3. curl -X POST -H 'Content-Type: application/json' -d @supervisor-spec.json http://localhost:8090/druid/indexer/v1/supervisor

Add a catalog properties file etc/catalog/kafka.properties
for the Kafka connector.

```
connector.name=kafka
kafka.nodes=localhost:9092
kafka.table-names=aSchema.table_name
kafka.hide-internal-columns=false
```

```
$ ./trino --catalog kafka --schema aSchema

trino:aSchema> SELECT count(*) FROM customer;
```

**Streaming Processor**

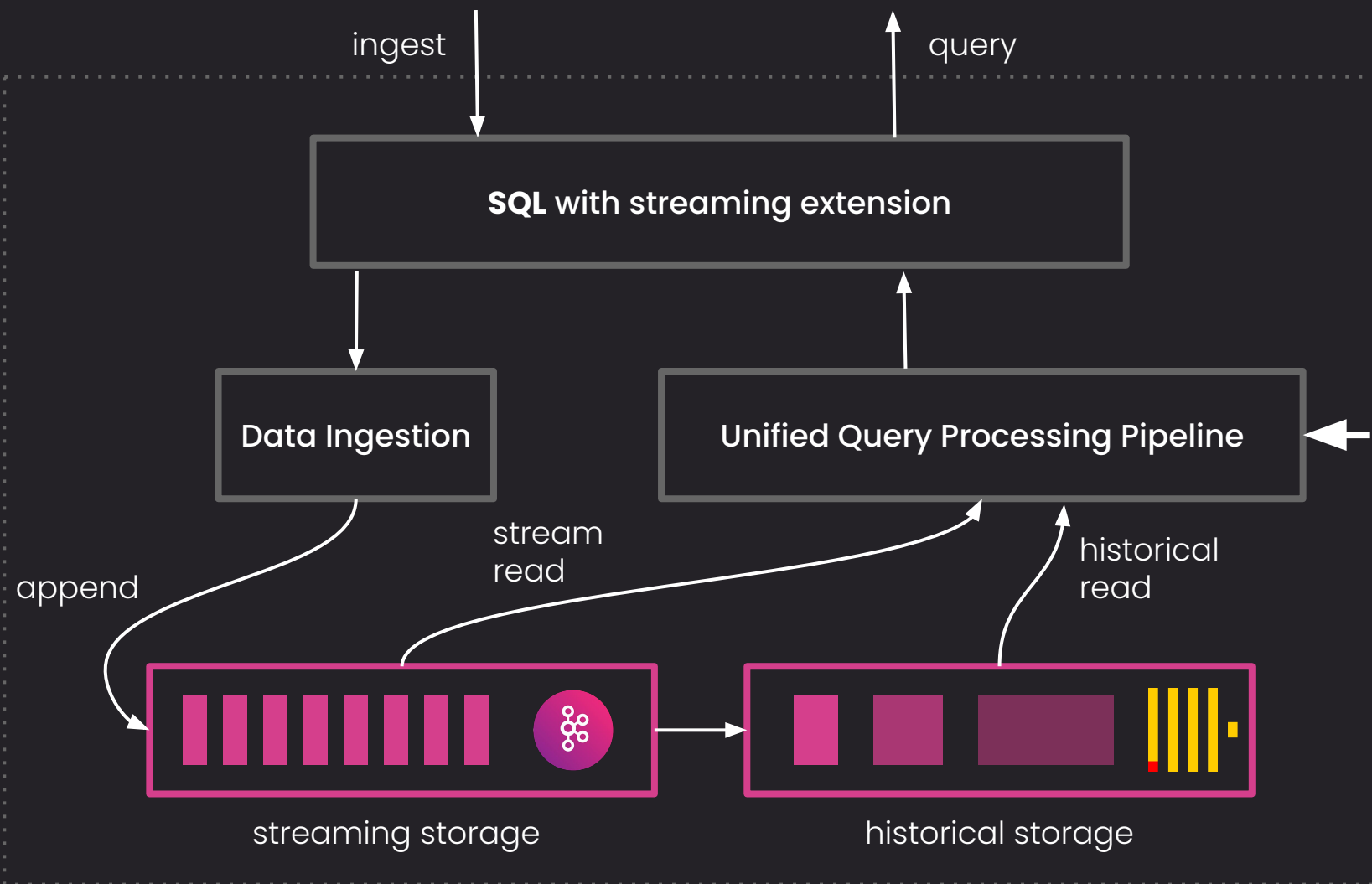**Streaming Database**

**Real-time Database**

proton

ClickHouse    StarRocks

timeplus    https://github.com/timeplus-io/proton

OSA CON 23

# proton

**Stream tail**

```sql
SELECT * FROM car_live_data
```

**Historical query**

```sql
SELECT * FROM table(car_live_data)
```

**Global aggregation**

```sql
SELECT count(*) FROM car_live_data
```

**Window aggregation**

```sql
SELECT window_start, count(*)
FROM tumble(car_live_data, 1m)
GROUP BY window_start
```

**Sub streams**

```sql
SELECT cid,
       speed_kmh,
       lag(speed_kmh) OVER
       (PARTITION BY cid) AS last_spd
FROM car_live_data
```

**Late event**

```sql
SELECT window_start, count(*)
FROM tumble(car_live_data, 5s)
GROUP BY window_start
EMIT AFTER WATERMARK AND DELAY 2s
```

**Time travel**

```sql
SELECT *
FROM car_live_data
WHERE
  _tp_time > now() - 1d
```

**Stream join**

```sql
SELECT
  device, cpu_usage, timestamp
FROM
  device_utils
INNER JOIN
table(device_products_info) AS dim
ON device_utils.product_id = dim.id
```

timeplus   https://github.com/timeplus-io/proton

OSA CON 23

**proton**

Since 2021

Mocha

Coffee Tasting Notes

| | |
|---|---|
| Community | ☕☕ |
| Real-time | ☕☕☕☕ |
| Streaming | ☕☕☕ |
| Historical | ☕☕☕☕ |
| JOIN | ☕☕☕ |
| Largescale | ☕☕ |
| Lightweight | ☕☕☕☕ |
| Easy to use | ☕☕☕ |

# proton → timeplus

## Add data to your workspace

We support many systems and methods to pull or push data into Timeplus.

Still exploring? Try o

Apache K

Redpand

Import f
Choose a

**Query**

SQL Helper

● Tab 1 ×  +

1  select * from car_live_data

Streaming Query  Total Results

_tp_time (Event Time

From:
2023-09-14T23:15:25.86Z
To:
2023-09-14T23:15:37.202Z

2023-09-14T23:15:37.1
2023-09-14T23:15:37.1
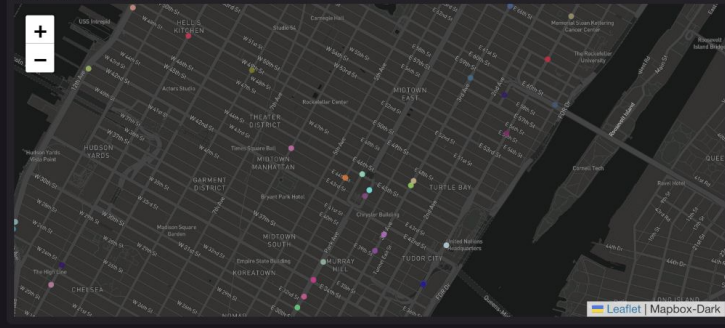2023-09-14T23:15:37.1
2023-09-14T23:15:37.1
2023-09-14T23:15:37.1
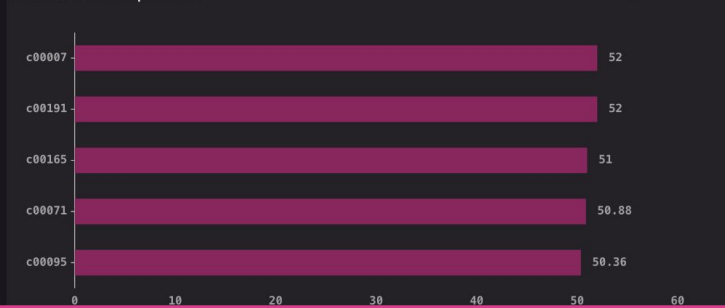2023-09-14T23:15:37.1
2023-09-14T23:15:37.1
2023-09-14T23:15:37.1
2023-09-14T23:15:37.1
2023-09-14T23:15:37.2
2023-09-14T23:15:37.2

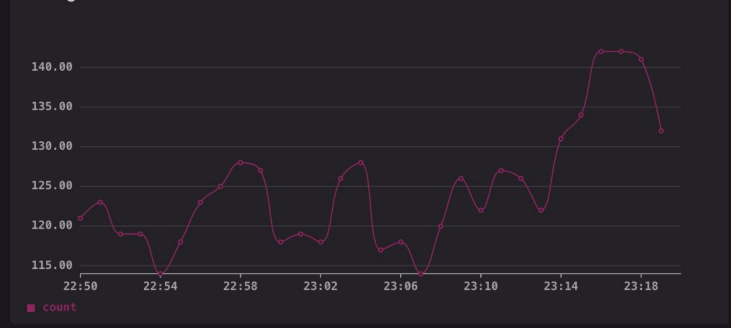Monitor car in real-time

**Real-time car loaction**

Leaflet | Mapbox-Dark

**Running Car Per Minute**                Last updated: 18s ago

140.00
135.00
130.00
125.00
120.00
115.00
    22:50  22:54  22:58  23:02  23:06  23:10  23:14  23:18
■ count

**Fastest 5 cars in past min**          Last updated: Just Now

c00007                                          52
c00191                                          52
c00165                                          51
c00071                                          50.88
c00095                                          50.36
     0    10    20    30    40    50    60

**Revenue for last 5 minutes**          Last updated: Just Now

# $1,570.00

−3.04

lo_u...

View
stock_portfolio_u...

itie...

Stream
account

View
car_info

timeplus   https://github.com/timeplus-io/proton

OSA CON 23

**ClickHouse**

```
CREATE TABLE queue2 (
    timestamp UInt64,
    level String,
    message String
 )
ENGINE = Kafka
SETTINGS
  kafka_broker_list = 'localhost:9092',
  kafka_topic_list = 'topic',
  kafka_group_name = 'group1',
  kafka_format = 'JSONEachRow',
  kafka_num_consumers = 4;
```

**StarRocks**

```
CREATE ROUTINE LOAD test_db.table102
ON table1
COLUMNS TERMINATED BY ",",
COLUMNS (user_id, user_gender, event_date, event_type)
WHERE event_type = 1
FROM KAFKA
(
    "kafka_broker_list" = "broker:port",
    "kafka_topic" = "topic1",
    "property.kafka_default_offsets" = "OFFSET_BEGINNING"
);
```

**ClickHouse** features highlights
- Table engine and table function
- Rich functions 1500+
- Rich data types - Array, Map etc

**StarRocks** features highlights
- More capable of joins
- High concurrency
- High frequency changes

OSA CON 23

**RisingWave**

**Cappuccino**

frothed milk

milk

espresso

**Coffee Tasting Notes** ☕

Community   ☕☕☕
Real-time   ☕☕☕
Streaming   ☕☕☕
Historical  ☕☕
JOIN        ☕☕☕
Largescale  ☕☕☕
Lightweight ☕☕☕☕
Easy to use ☕☕☕

```
docker run -it --pull=always -p 4566:4566 -p 5691:5691
ghcr.io/risingwavelabs/risingwave:latest playground

psql -h localhost -p 4566 -d dev -U root
```
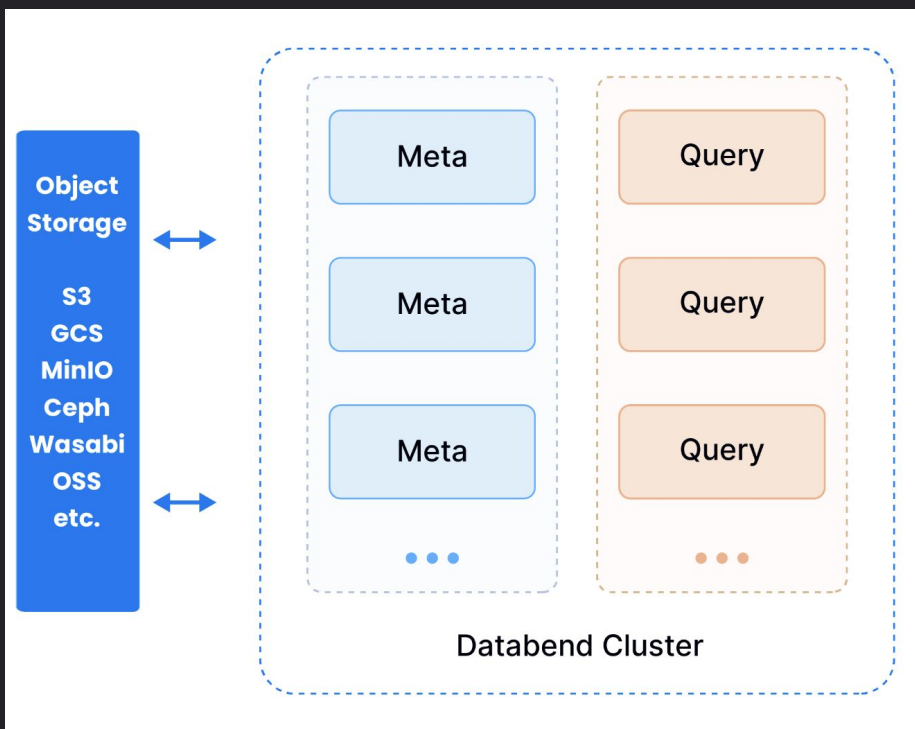
```sql
CREATE TABLE github_events (
  id varchar,
  created_at timestamp,
  actor varchar,
  type varchar,
  repo varchar,
  payload jsonb
) WITH (
  connector = 'kafka',
  topic = 'github_events',
  properties.bootstrap.server = 'xyz.aws.confluent.cloud:9092',
  scan.startup.mode = 'earliest',
  properties.security.protocol = 'SASL_SSL',
  properties.sasl.mechanism = 'PLAIN',
  properties.sasl.username = 'username',
  properties.sasl.password = 'password'
) FORMAT PLAIN ENCODE JSON;
```

```sql
SELECT window_start, window_end, count(*) as events
FROM HOP (github_events, created_at,
          INTERVAL '1 MINUTES', INTERVAL '2 MINUTES')
GROUP BY window_start, window_end
ORDER BY window_start ASC;
```

# Query Kafka with SQL: Open Source + Cloud + Closed Source

**Streaming Processor**

Flink
decodable
databricks

Arroyo
DOZER
Feldera

**Streaming Database**

ksqlDB
Hazelcast

DeltaStream

proton

RisingWave
Materialize

**Real-Time Database**

Druid
Pinot
Trino

ClickHouse
StarRocks
ROCKSET
tinybird

Databend
pathway

OSA CON 23

# Q+A / Thank you!

Try Timeplus Proton (Open Source)
Or sign up for a free cloud account

timeplus.com

**Jove Zhong**          **Gang Tao**