



Continuous SQL with SQL Stream Builder

Kenny Gorman - Product Owner

Timothy Spann - Principal DataFlow Field Engineer

John Kuchmek - Senior Solutions Engineer

06-May-2021

<https://www.meetup.com/futureofdata-newyork/>

@PaasDev

Welcome to Future of Data - Virtual

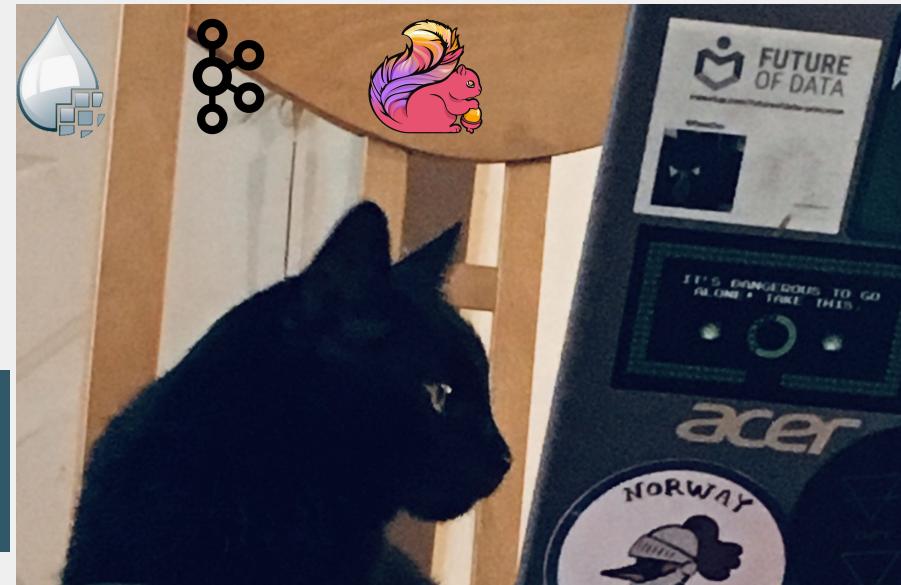


[Princeton Future of Data Meetup](#)

[New York Future of Data Meetup](#)

[Philadelphia Future of Data Meetup](#)

**From Big Data to AI to Streaming to Containers to
Cloud to Analytics to Cloud Storage to Fast Data to
Machine Learning to Microservices to ...**



AGENDA

- Introductions with Kenny, John and Tim
- Flink Quick Overview
- SQL Stream Builder Overview
- Q&A
- Demos
- Q&A - Interactive Panel Session
- Next Meetups
- Raffle

Cloudera DataFlow Use Cases

Data Movement

Optimize resource utilization by moving data between data centers or between on-premises and cloud infrastructures

e.g. *intercontinental data exchange*

Logging Modernization

Optimize log analytics solutions by with CDF in simplifying log ingestion from the edge, reducing costs and gaining key analytics

e.g. *Splunk / Logstash offload*

Streaming analytics insights

Make key business decisions by analyzing streaming data for complex patterns, gaining actionable intelligence etc.

e.g. *Fraud detection, Network threat analysis, app monitoring, Clickstream analysis*

360° view of customer

Ingest, transform and combine customer data from multiple sources into a single data view / lake

e.g. *Real-time customer offers, Loan approvals*

IoT & Edge use cases

e.g. *Predictive Maintenance, Asset Tracking / Monitoring, Patient Monitoring, Quality Processes, Fleet Management, Connected Cars and more*

Enterprise data management

Managing massive volumes of high-velocity data to/from legacy systems, ETL tools and other data stores

e.g. *Flume offload, ETL replacement, payment data processing, integration with Oracle*

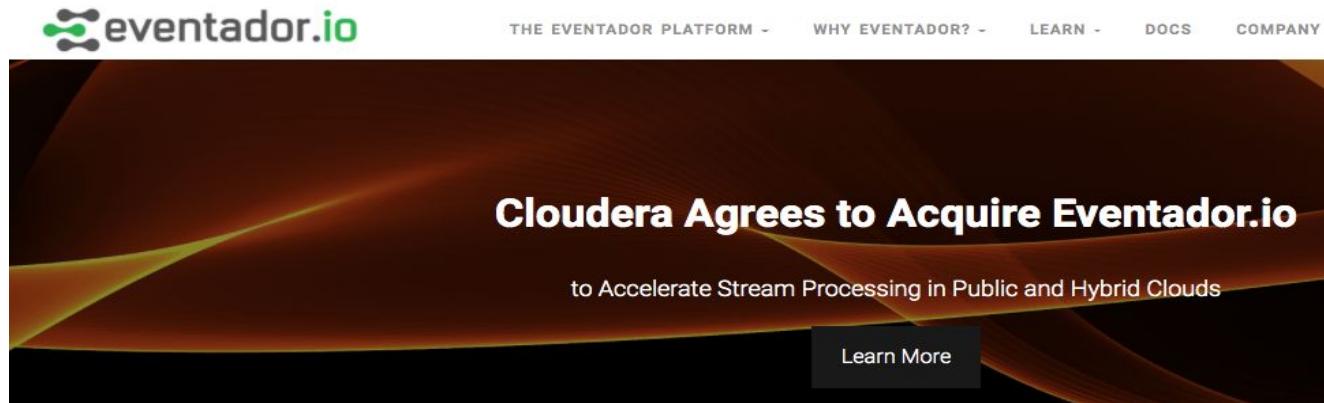


Simplifying the User Experience



October 12, 2020

Cloudera acquires Eventador to accelerate Stream Processing in Public & Hybrid Clouds



While the acquisition may be seen as incremental or small in nature as Eventador isn't yet a household name, the capabilities and improvements it has the potential to make for Cloudera shouldn't be understated. A solid move by Cloudera to continue its momentum forward.

APACHE FLINK

3B+

3B+ data points daily streaming in from 25 million customers running real time machine learning prediction



Flink

USE CASE

Streaming real-time data pipelines that need to handle complex stream or batch data event processing, analytics, and/or support event-driven applications

TECHNOLOGY

Flink performs compute at in-memory speed at any scale

Flink parses SQL using Apache Calcite, which supports standard ANSI SQL

Flink runs standalone, on YARN, and has a K8s Operator

APPLICATION

Comcast a global media uses Flink for operationalizing machine learning models and near-real-time event stream processing

Flink helps deliver a personalized, contextual interaction reducing time to support resolutions saving millions of dollars per year

CONSIDERATION

Data Freshness SLAs

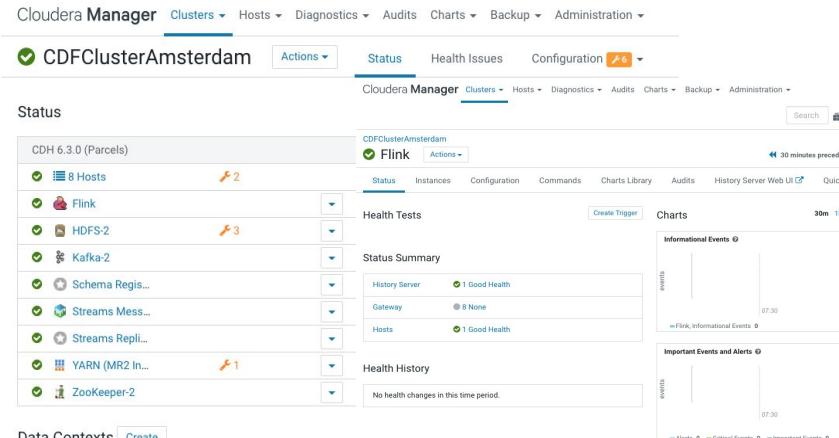
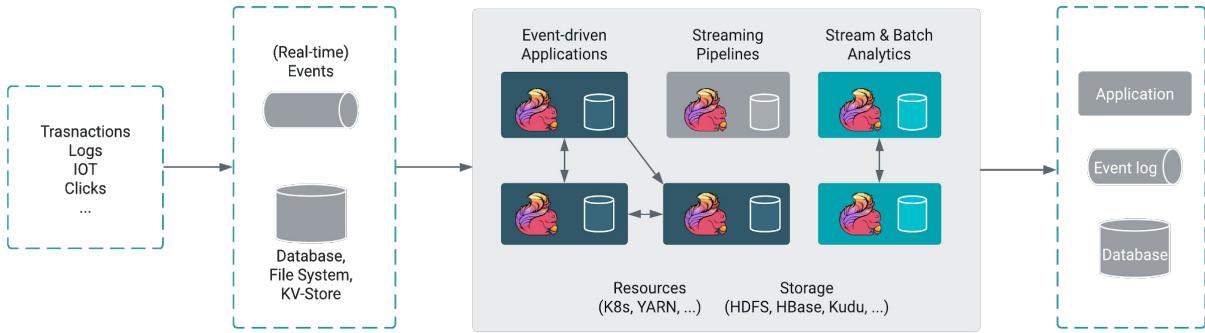
Flink can read and write from Hive data

Review requirements for fault tolerance, resilience, and HA

Other technologies play in this space like Hive storage handler to connect to Kafka

FLINK FEATURES

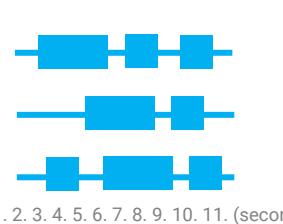
- Distributed processing engine for stateful computations
- 10s TBs of managed state
- Flexible and expressive APIs
- Guaranteed correctness & Exactly-once state consistency
- Event-time semantics
- Flexible deployment & large ecosystem (K8s, YARN, S3, HDFS..)
- Support for Flink SQL API



DELIVERING STREAMING ANALYTICS

Capture Events that Matter

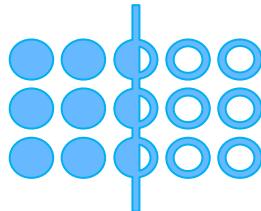
Low-latency analytics use cases



Events Processing

Parsing and Blending Data

Both offline and streaming data



Data Analysts Can Write Queries

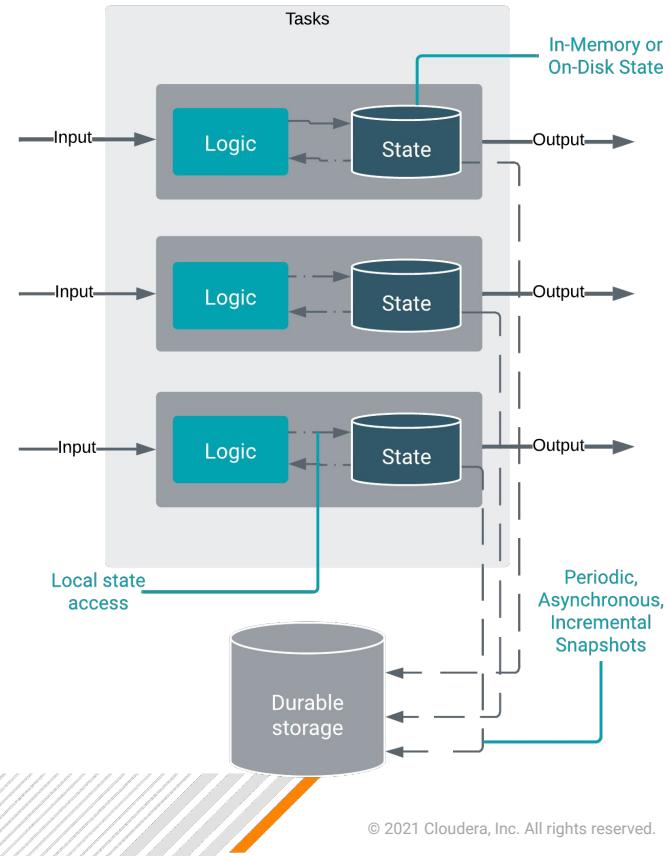
Across the Line of Businesses



Streaming Analytics

MAINTAINS & CHECKPOINTS STATE

- Flink maintains state locally per task (in-mem / on-disk)
 - Fast access!
- State is periodically checkpointed to durable storage
 - A checkpoint is a **consistent** snapshot of the state of all tasks



Integrated Governance

Unified Governance & Lineage

Apache nifi Flow Management

Reports Entity and Lineage information about NiFi Flows

Connects with existing Lineage information



Streams Messaging

Topic access centrally managed supporting granular CRUD operations

Manage permissions on dedicated clusters or manage multiple clusters at once

Manage schemas centrally and make them available to consumers/producers



Stream Processing

Reports Flink Apps as an operation

Lineage through integration with existing Lineage information like Kafka topics, HBase tables etc.

Integrated SQL and materialized view engine via SQL Stream Builder.

SQL Stream Builder

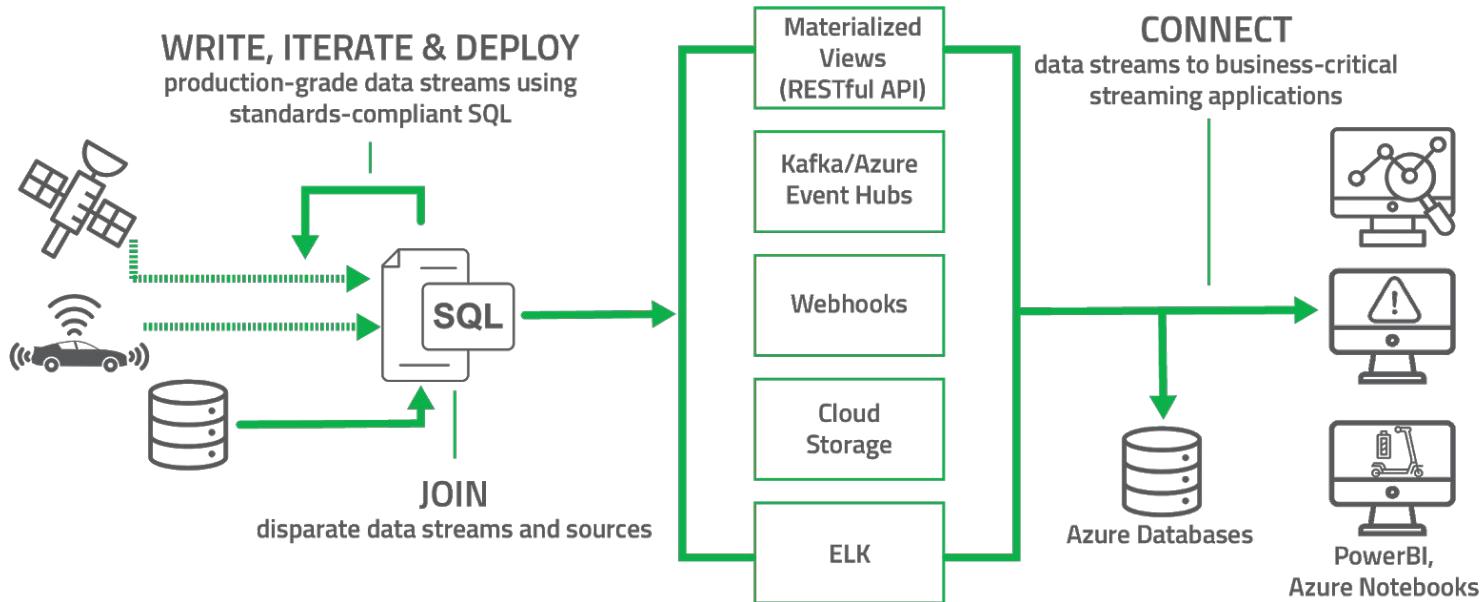
- **Democratize** data access across enterprise - anyone who knows SQL can create powerful stream processors.
- **Iterative interface** - Just like SQL on databases, run queries and reason about the data with an interactive UI.
- Leverages **Apache Flink** for running of SQL jobs - production grade, scalable and high performance
- Deep integration and **features** above and beyond just UI features - UDF's, input transforms, Kafka key and time integration, CEP framework and more.
- Create **Materialized Views** to integrate with downstream components like notebooks, visualizations and applications.

The screenshot shows the Stream Builder interface. At the top, there's a code editor window titled "SQL" containing a query:1 SELECT TUMBLE_END(geo_events.eventTimestamp, INTERVAL '3' MINUTE) AS windowEnd,
2 geo_events.driverId, geo_events.driverName, geo_events.route,
3 AVG(speed_events.speed) AS driverAvgSpeed
FROM
4 geo_events,
5 speed_events
WHERE
6 geo_events.driverId = speed_events.driverId AND
7 geo_events.eventTimestamp BETWEEN
8 speed_events.eventTimestamp - INTERVAL '1' SECOND AND
9 speed_events.eventTimestamp + INTERVAL '1' SECOND
10 GROUP BY
11 TUMBLE(geo_events.eventTimestamp, INTERVAL '3' MINUTE),
12 geo_events.driverId,A toolbar below the editor includes buttons for "default mode", "solarized dark", "Sample", "Stop", and "Restart". Below the editor is a log viewer with tabs for "Logs", "Results", and "Help". The log pane displays the following output:

```
[9/4/2020, 3:35:12 PM][INFO] No persistent sink specified, using ephemeral sink.  
[9/4/2020, 3:35:12 PM][INFO] StreamBuilder job Speeding Drivers Over 3 Minute Window is starting.  
[9/4/2020, 3:35:25 PM][INFO] SSB version 8.0.4 selected for job.  
[9/4/2020, 3:35:25 PM][INFO] Streaming job is now running in the background. You can safely navigate to other pages now, and re-visit the running job by clicking the SQL Jobs tab.  
[9/4/2020, 3:35:25 PM][INFO] Stream sample is running, and will display the next 100 messages matching your query.  
[9/4/2020, 3:35:25 PM][INFO] ⏺ Wait for results from stream...  
[9/4/2020, 3:36:58 PM][INFO] Stoped job Speeding Drivers Over 3 Minute Window with job ID 4582  
[9/4/2020, 3:37:00 PM][INFO] Job Speeding Drivers Over 3 Minute Window is stopped.  
[9/4/2020, 3:37:17 PM][INFO] StreamBuilder job Speeding Drivers Over 3 Minute Window is starting.
```

Streaming SQL

Democratizing access to streams of data via structured query language



Download these assets today

CLOUDERA

WHITE PAPER

CLOUDERA DELIVERS THE BEST KAFKA ECOSYSTEM TODAY

Serving Hundreds of Customers Globally




CLOUDERA

SOLUTION BRIEF

DATA-IN-MOTION PHILOSOPHY

A Blueprint for Enterprise-wide Streaming Data Architecture




CLOUDERA

WHITE PAPER

CHOOSE THE RIGHT STREAM PROCESSING ENGINE FOR YOUR DATA NEEDS

Technical and Operational Factors that are Crucial to the Decision Making Process




The Complete Edge-to-Cloud Streaming Data Platform

This paper has focused primarily on the streams messaging aspects of the Kafka ecosystem with regard to how to secure, monitor, balance, and replicate large scale Kafka environments across on-premises, hybrid, private, and public cloud environments. The Data-in-Motion reference architecture diagram in Figure 7 below, puts this all in perspective.

A DATA-IN-MOTION REFERENCE ARCHITECTURE

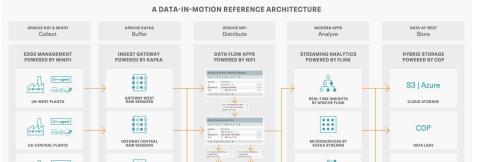


Figure 7

As a result of supporting our customers through their data journeys, we learned that it is not enough to have the best messaging solution at the heart of your end-to-end streaming architecture. As represented in the diagram above, flow management, along with stream processing and analytics, are two additional tenets that need to be unified with streams messaging capabilities. These three tenets, if properly integrated, will ensure a sustainable, scalable and adaptable end-to-end streaming architecture and is the basis to our data-in-motion philosophy.

Real-time Data-in-Motion

Below describes how the data-in-motion capabilities described in this solution brief have been successfully applied by a Cloudera customer, end-to-end.

A global medical device manufacturer successfully modernized their messaging architecture to support a highly implementable real-time data pipeline that generates more data, more often, and at a higher resolution than ever before.

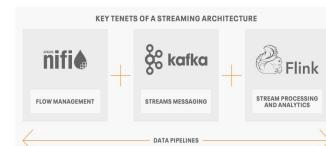
- Flow management**—Due to the physical nature of medical data, the data flow was complex, requiring a motion and at rest state. The company made user interface and analysis available to 100% business driven, only engaging the technology teams as needed.
- Streams messaging**—Messaging volume justifies the quarterly reporting of data states to enable continuous monitoring. Kafka enables the business to scale that volume across multiple on-premises and cloud environments.
- Stream processing and analytics**—How you generate real-time analytical insights from the data streaming between producers and consumers.

Streaming Architecture in Context

Below we have listed the three tenets that together provide a unified end-to-end streaming architecture.

- Flow management**, broadly speaking, refers to the collection, distribution, and transformation of data across multiple points of producers and consumers.
- Streams messaging** is the provisioning and distribution of messages between producers and consumers.
- Stream processing and analytics** is how you generate real-time analytical insights from the data streaming between producers and consumers.

KEY TENETS OF A STREAMING ARCHITECTURE



Cloudera's data-in-motion philosophy is rooted in the complementary powers that are brought to the table by Apache Nifi for flow management, Apache Kafka for streams messaging, and Apache Flink for stream processing and analytics.

Operational Features Table

The table below gives an operational comparison across four modern stream processing engines. Refer to it when evaluating the nonfunctional aspects of your project.

	Flink 1.10	Kafka Streams 2.4	Spark Structured Streaming 2.4	Storm 2.0 and Trident
Deployment model	<ul style="list-style-type: none"> Clustered Kubernetes Microservices Kafka Docker IoT Microservices 	<ul style="list-style-type: none"> Not clustered Kubernetes Microservices 	<ul style="list-style-type: none"> Clustered Kubernetes 	<ul style="list-style-type: none"> Clustered Microservices
Documentation	<ul style="list-style-type: none"> Good technical documentation Growth in popularity Good developer coverage Stack Overflow coverage 	<ul style="list-style-type: none"> Extensive documentation Extensive examples Stack Overflow coverage 	<ul style="list-style-type: none"> Extensive documentation Extensive examples Stack Overflow coverage 	<ul style="list-style-type: none"> Good documentation for 1.x
Maturity/community	<ul style="list-style-type: none"> Similar but fastest growing community with strong research and production development 	<ul style="list-style-type: none"> Newest, strong community with strong growth 	<ul style="list-style-type: none"> Spark Structured Streaming community is strong, though streaming is a small, quiet community 	<ul style="list-style-type: none"> Oldest framework, community edged by newer engines
Use cases	<ul style="list-style-type: none"> Unbounded and bounded streams Batch Complex event processing IoT Microservices Others 	<ul style="list-style-type: none"> Microservices/event driven, embedded in another application 	<ul style="list-style-type: none"> Unified ETL, semi-RT processing 	<ul style="list-style-type: none"> IoT, complex event processing
Enterprise management	<ul style="list-style-type: none"> Rich OSS Enhanced vendor offerings 	<ul style="list-style-type: none"> Minimal OSS Some via vendor offerings 	<ul style="list-style-type: none"> Rich OSS Some via vendor offerings 	<ul style="list-style-type: none"> Some integrations
Push button security	<ul style="list-style-type: none"> Complex Some OSS support Cloud native integration Cloud vendor offerings 	<ul style="list-style-type: none"> Simple, some OSS support, good vendor offerings 	<ul style="list-style-type: none"> Complex Cloud OSS support Cloud vendor offerings 	<ul style="list-style-type: none"> Complex Good OSS support Cloud vendor offerings
Logging/metrics	<ul style="list-style-type: none"> Cloud native integration Cloud vendor offerings 	<ul style="list-style-type: none"> BTO/microservices 	<ul style="list-style-type: none"> Cloud logging integration 	<ul style="list-style-type: none"> Cloud logging integration
Scaling up/down	<ul style="list-style-type: none"> Not yet auto-scaling, but all requirements available 	<ul style="list-style-type: none"> BTO/microservices, scaling limits (e.g. shuffle sort) 	<ul style="list-style-type: none"> Not yet auto-scaling, but all requirements available 	<ul style="list-style-type: none"> Management tools help but tuning is challenging
Overall fit	<ul style="list-style-type: none"> Dread fit for purpose 	<ul style="list-style-type: none"> Fits with some work 	<ul style="list-style-type: none"> Fits with a lot of work 	<ul style="list-style-type: none"> Not fit for purpose

CLOUDERA

© 2021 Cloudera, Inc. All rights reserved.

14

