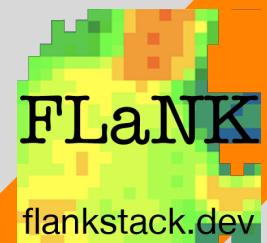




# Real-Time Streaming Pipelines With FLaNK (Apache Flink, Apache NiFi & Apache Kafka)

Timothy Spann - Principal DataFlow Field Engineer

15-April-2021



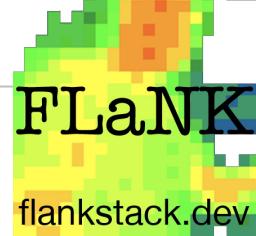
# Tim SPANN

<https://github.com/tspannhw>

<https://www.datainmotion.dev/>

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





# FLaNK Stack for Cloud Data Engineers

Multiple users, frameworks, languages, clouds, data sources & clusters



CLOUD DATA ENGINEER

- Experience in ETL/ELT
- Coding skills in Python or Java
- Knowledge of database query languages such as SQL
- Experience with Streaming
- Knowledge of Cloud Tools



CAT

- Expert in ETL (Eating, Ties and Laziness)
- Edge Camera Interaction
- Typical User
- No Coding Skills
- Can use NiFi
- Questions your cloud spend



AI / Deep Learning / ML / DS

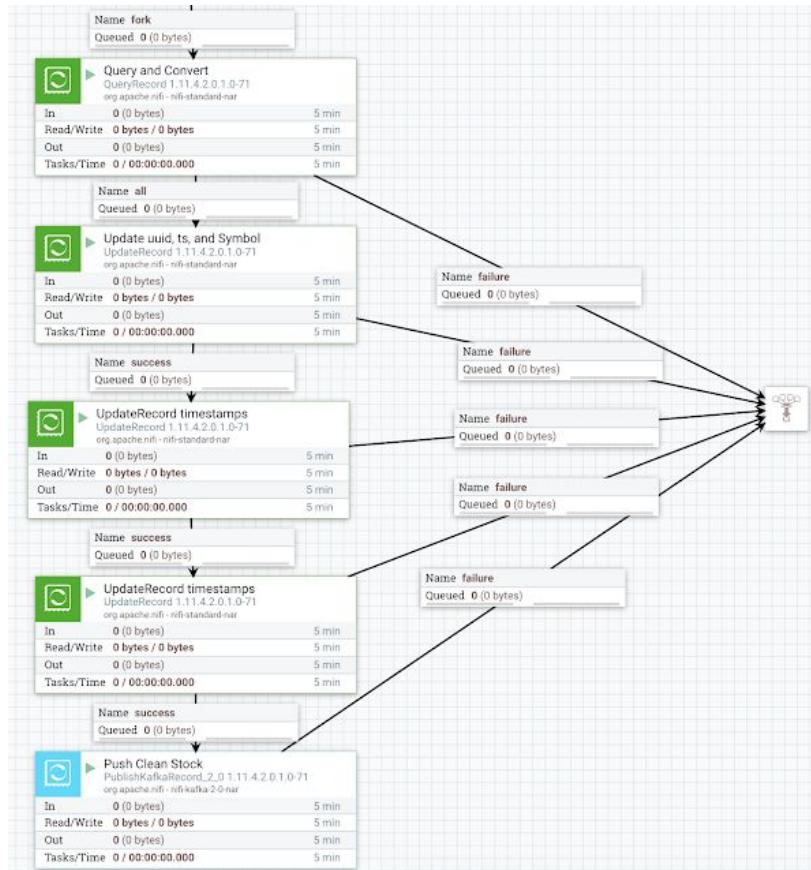
- Can run in Apache NiFi
- Can run in Kafka Streams
- Can run in Apache Flink
- Can run in MiNiFi Agents

# Today's Data. REST and Websocket JSON “stonks”

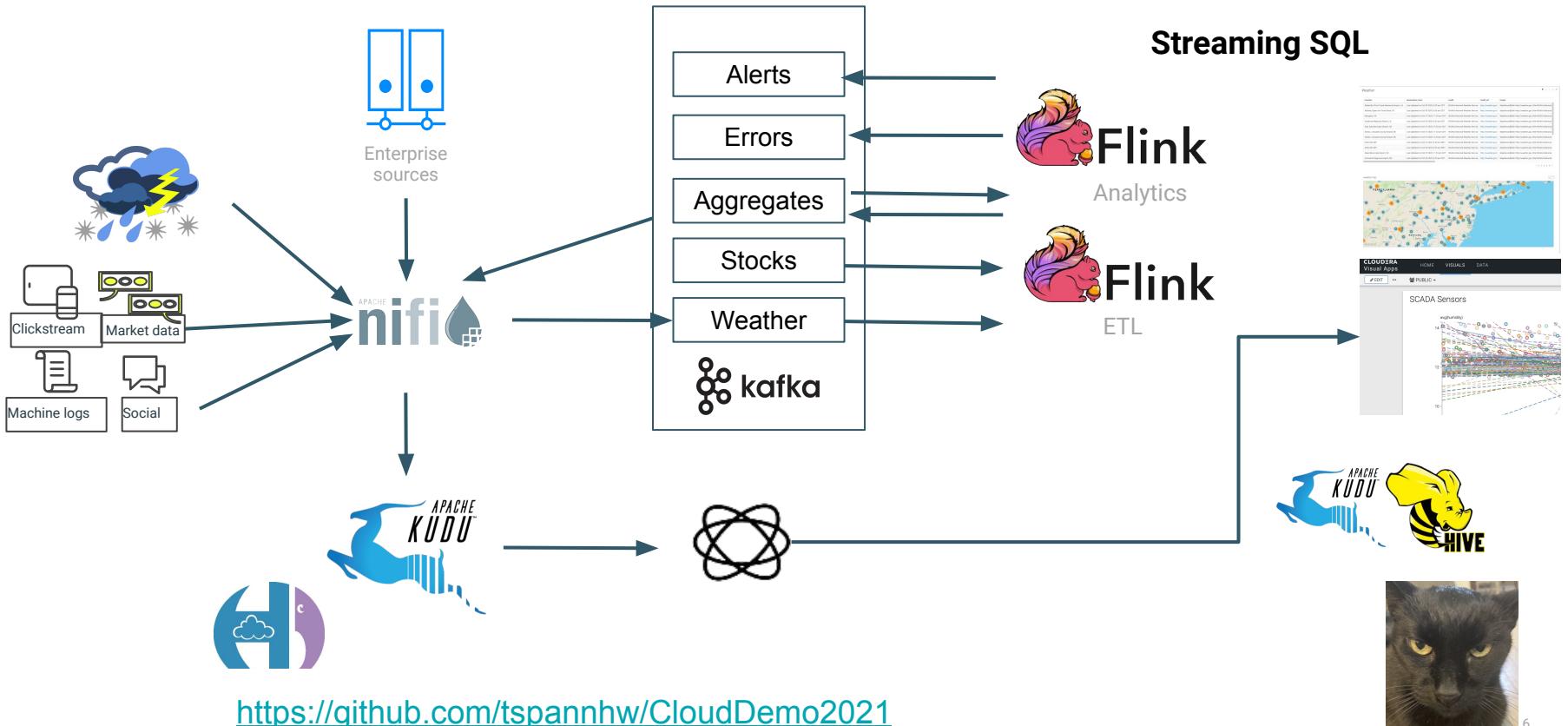


```
{"symbol":"CLDR",
"uuid":"10640832-f139-4b82-8780-e3ad37b3d0
ce",
"ts":1618529574078,
"dt":1612098900000,
"datetime":"2021/01/31 08:15:00",
"open":12.24500,
"close":12.25500,
"high":12.25500,
"volume":12353,
"low":12.24500}
```





# End to End Streaming Demo Pipeline



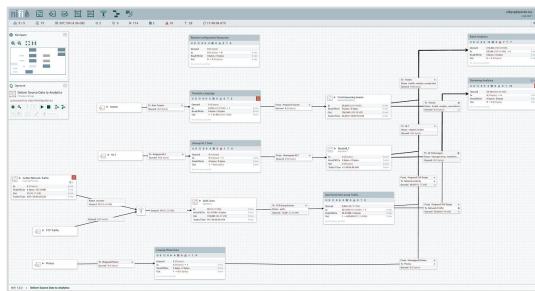


# FLINK

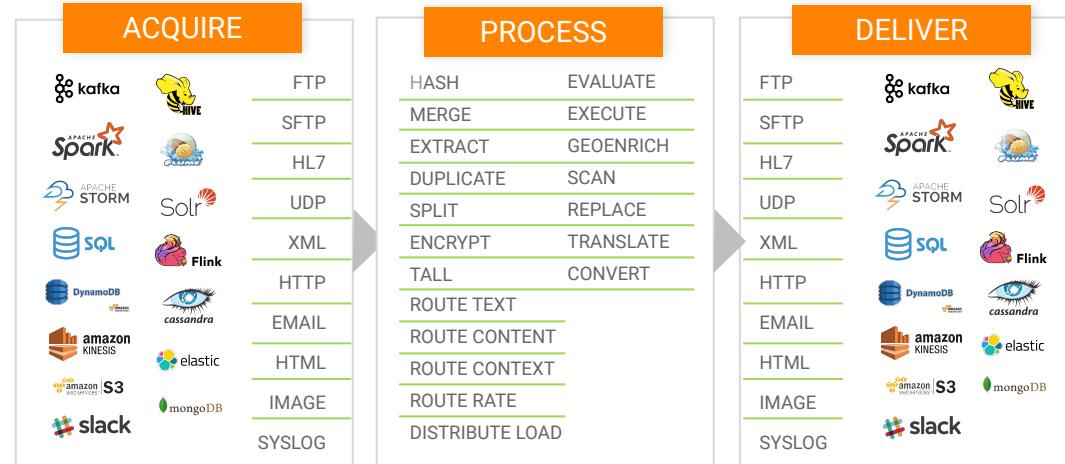
## NIFI - KAFKA - FLINK

# Apache NiFi in a Nutshell

Enable easy ingestion, routing, management and delivery of any data anywhere (*Edge, cloud, data center*) to any downstream system with built in end-to-end security and provenance



Advanced tooling to industrialize flow development (*Flow Development Life Cycle*)



- Over 340 Prebuilt Processors
- Easy to build your own
- Parse, Enrich & Apply Schema
- Filter, Split, Merger & Route
- Throttle & Backpressure

- Guaranteed Delivery
- Full data provenance from acquisition to delivery
- Diverse, Non-Traditional Sources
- Eco-system integration

# ParquetReader / ParquetWriter Records

- Native Record Processors for Apache Parquet Files!
- CVS <-> Parquet
- XML <-> Parquet
- AVRO <-> Parquet
- JSON <-> Parquet
- More...

<https://www.datainmotion.dev/2019/10/migrating-apache-flume-flows-to-apache-7.html>

<https://www.datainmotion.dev/2019/11/exploring-apache-nifi-110-parameters.html>

Property	Value
Record Reader	JsonTreeReader
Record Writer	ParquetRecordSetWriter
Merge Strategy	Bin-Packing Algorithm
Correlation Attribute Name	No value set
Attribute Strategy	Keep Only Common Attributes
Minimum Number of Records	
Maximum Number of Records	
Minimum Bin Size	
Maximum Bin Size	Requires Controller Service RecordReaderFactory 1.10.0.2.0.0.0-35 from org.apache.nifi - nifi-standard-services-api-nar
Max Bin Age	Compatible Controller Services ParquetReader 1.10.0.2.0.0.0-35
Maximum Number of Bins	
Add Controller Service	
Controller Service Name	ParquetReader
Bundle	org.apache.nifi - nifi-parquet-nar
Tags	reader, record, parse, row, parquet
Description	

CANCEL CREATE

# UpdateRecord

- Use with LookupRecord
- ELT
- Works on CSV, XML, JSON, AVRO, ...
- RecordPath or Literals
- Use Schemas and Schema Registry

The screenshot shows two instances of the 'UpdateRecord timestamps' component from org.apache.nifi - nifi-standard-nar. The top instance is associated with a 'success' flowfile and has the following properties:

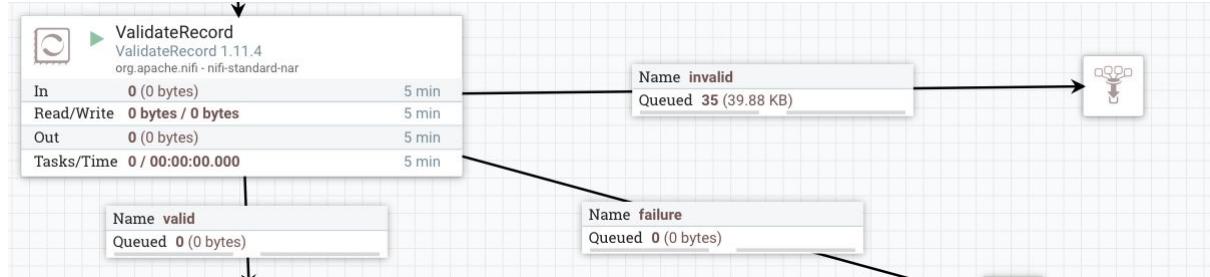
Property	Value
Record Reader	SchemaAwareJsonTreeReader
Record Writer	Standard Inherit JsonRecordSetWriter
Replacement Value Strategy	Record Path Value
/datetime	format(/datetime,"yyyy/MM/dd HH:mm:SS")

The bottom instance is associated with a 'success' flowfile and has the following properties:

Property	Value
Record Reader	SchemaAwareJsonTreeReader
Record Writer	Standard Inherit JsonRecordSetWriter
Replacement Value Strategy	Literal Value
/symbol	\${field.value:toDate("yyyy-mm-dd HH:mm:ss"):toNumber()}
/ts	\${symbol}
/uuid	\${now():toNumber()}

# ValidateRecord

- Works on CSV, XML, JSON, AVRO, ...
- RecordPath or Literals
- Use Schemas and Schema Registry
- Checks fields, types, nullable



<b>Record Reader</b>	<b>SchemaAwareJsonTreeReader</b>
<b>Record Writer</b>	<b>Standard Inherit JsonRecordSetWriter</b>
Record Writer for Invalid Records	No value set
<b>Schema Access Strategy</b>	<b>Use Schema Name Property</b>
Schema Registry	Edge2ai Cloudera SchemaRegistry CSchema
Schema Name	<code>\$(schema.name)</code>
Schema Text	<code>\$(avro.schema)</code>
<b>Allow Extra Fields</b>	<code>true</code>
<b>Strict Type Checking</b>	<code>false</code>

# RestLookupService

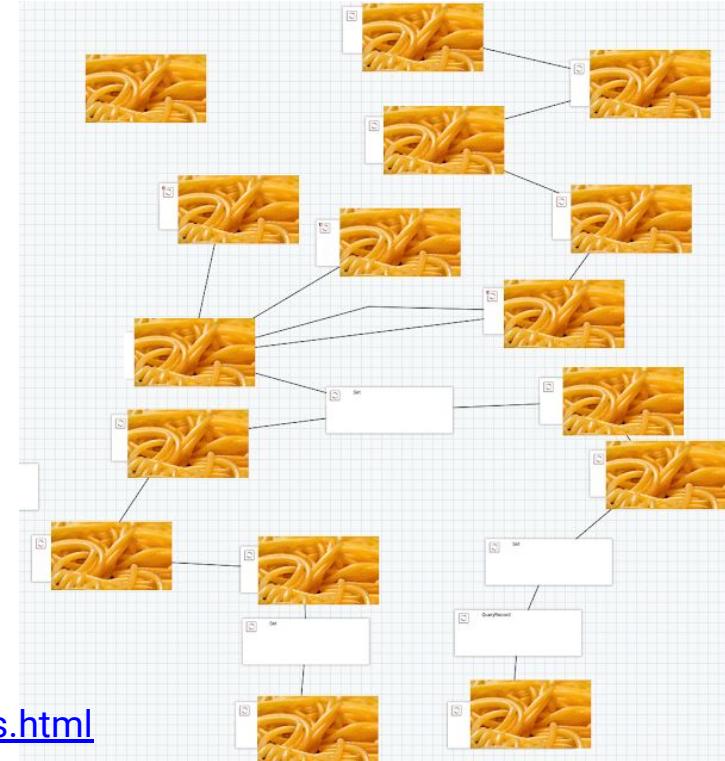
- Works on CSV, XML, JSON, AVRO, ...
- Use Schemas and Schema Registry
- Can call Cloudera ML Models
- SSL and Proxy enabled

Property	Value
URL	#{model_url}
Record Reader	JsonTreeReader
Record Path	#{sentiment.record.path}
SSL Context Service	No value set
Proxy Configuration Service	No value set
Basic Authentication Username	No value set
Basic Authentication Password	No value set
Use Digest Authentication	false
Connection Timeout	15 secs
Read Timeout	45 secs
Authorization	Bearer #{model.api.key}

# No More Spaghetti Flows

- Reduce, Reuse, Recycle. Use Parameters to reuse common modules.
- Put flows, reusable chunks into separate Process Groups.
- Write custom processors if you need new or specialized features
- Use Cloudera supported NiFi Processors
- Use Record Processors everywhere

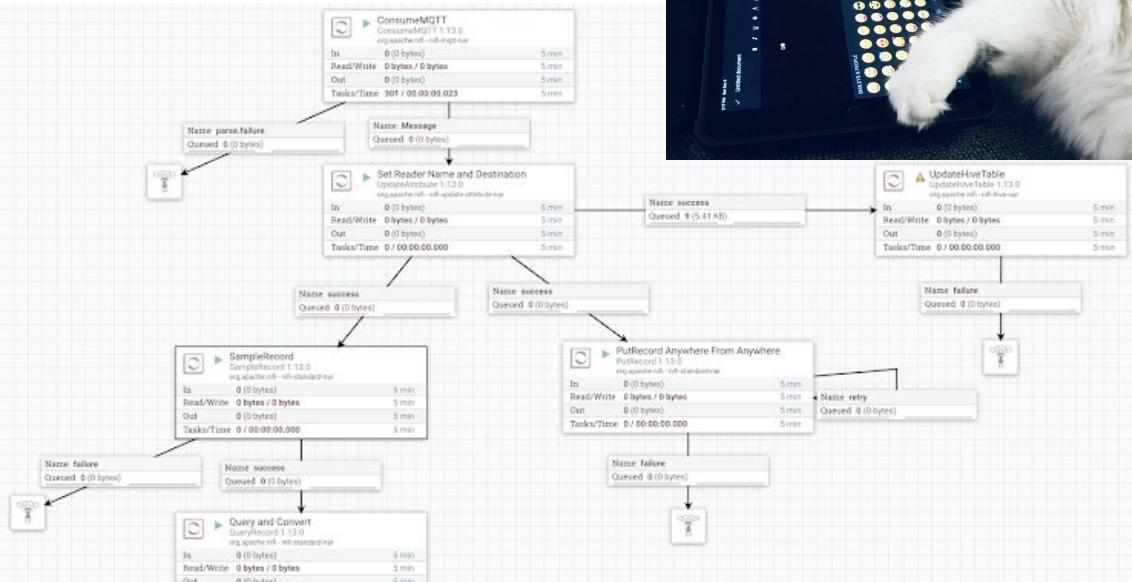
<https://www.datainmotion.dev/2020/06/no-more-spaghetti-flows.html>



# New Features

... based on Apache NiFi 1.13.2

- ListenFTP
- Data Drift
- SampleRecord
- Generic Record Sink
- Generic Record Reader
- PutRecord
- WindowsEventLogReader



<https://www.datainmotion.dev/2021/02/new-features-of-apache-nifi-1130.html>



Yes, Franz, It's Kafka

Let's do a metamorphosis on your data. Don't fear changing data.

## You don't need to be a brilliant writer to stream data.



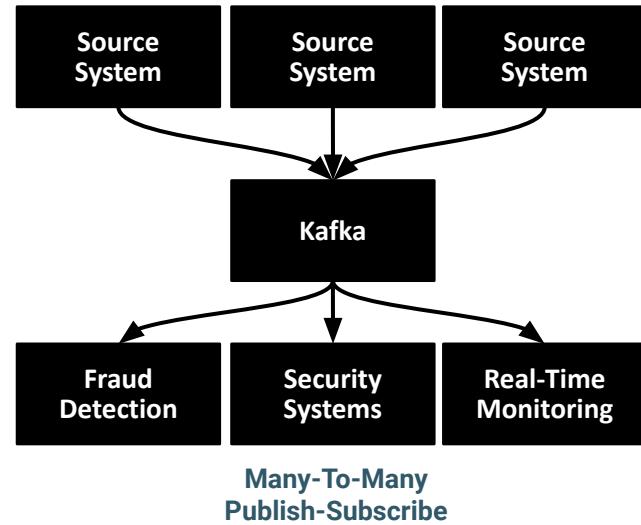
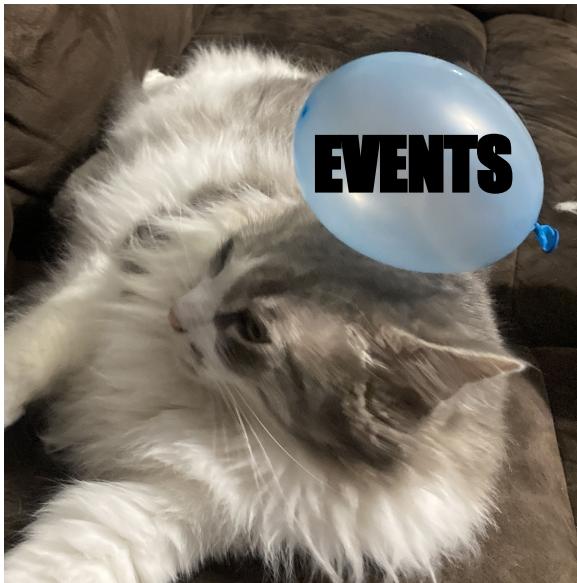
Franz Kafka was a German-speaking Bohemian novelist and short-story writer, widely regarded as one of the major figures of 20th-century literature. His work fuses elements of realism and the **fantastic**.

[Wikipedia](#)



# Apache Kafka

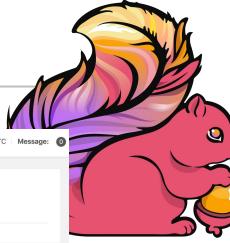
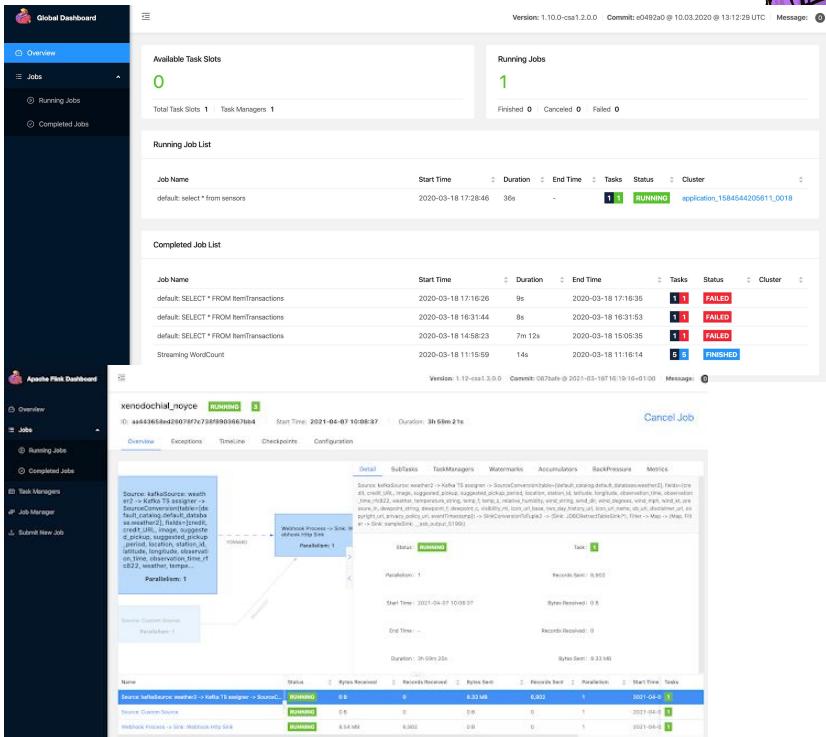
- Highly reliable distributed messaging system
- Decouple applications, enables many-to-many patterns
- Publish-Subscribe semantics
- Horizontal scalability
- Efficient implementation to operate at speed with big data volumes
- Organized by topic to support several use cases



# Flink SQL

... based on Apache Flink 1.12

- Streaming Analytics
  - Continuous SQL
  - Continuous ETL
  - Complex Event Processing
  - Standard SQL Powered by Apache Calcite
  - Deployed Apache Flink Apps on YARN
  - Scalable Stream Processing



<https://www.datainmotion.dev/2021/04/cloudera-sql-stream-builder-ssb-updated.html>

# Flink SQL

## Key Takeaway: Rich SQL grammar with advanced time and aggregation tools

```
-- specify Kafka partition key on output
SELECT foo AS _eventKey FROM sensors

-- use event time timestamp from kafka
-- exactly once compatible
SELECT eventTimestamp FROM sensors

-- nested structures access
SELECT foo.'bar' FROM table; -- must quote nested
column

-- timestamps
SELECT * FROM payments
WHERE eventTimestamp > CURRENT_TIMESTAMP-interval
'10' second;

-- unnest
SELECT b.* , u.*
FROM bgp_avro b,
UNNEST(b.path) AS u(pathitem)

-- aggregations and windows
SELECT card,
MAX(amount) as theamount,
TUMBLE_END(eventTimestamp, interval '5' minute) as
ts
FROM payments
WHERE lat IS NOT NULL
AND lon IS NOT NULL
GROUP BY card,
TUMBLE(eventTimestamp, interval '5' minute)
HAVING COUNT(*) > 4 -- >4==fraud

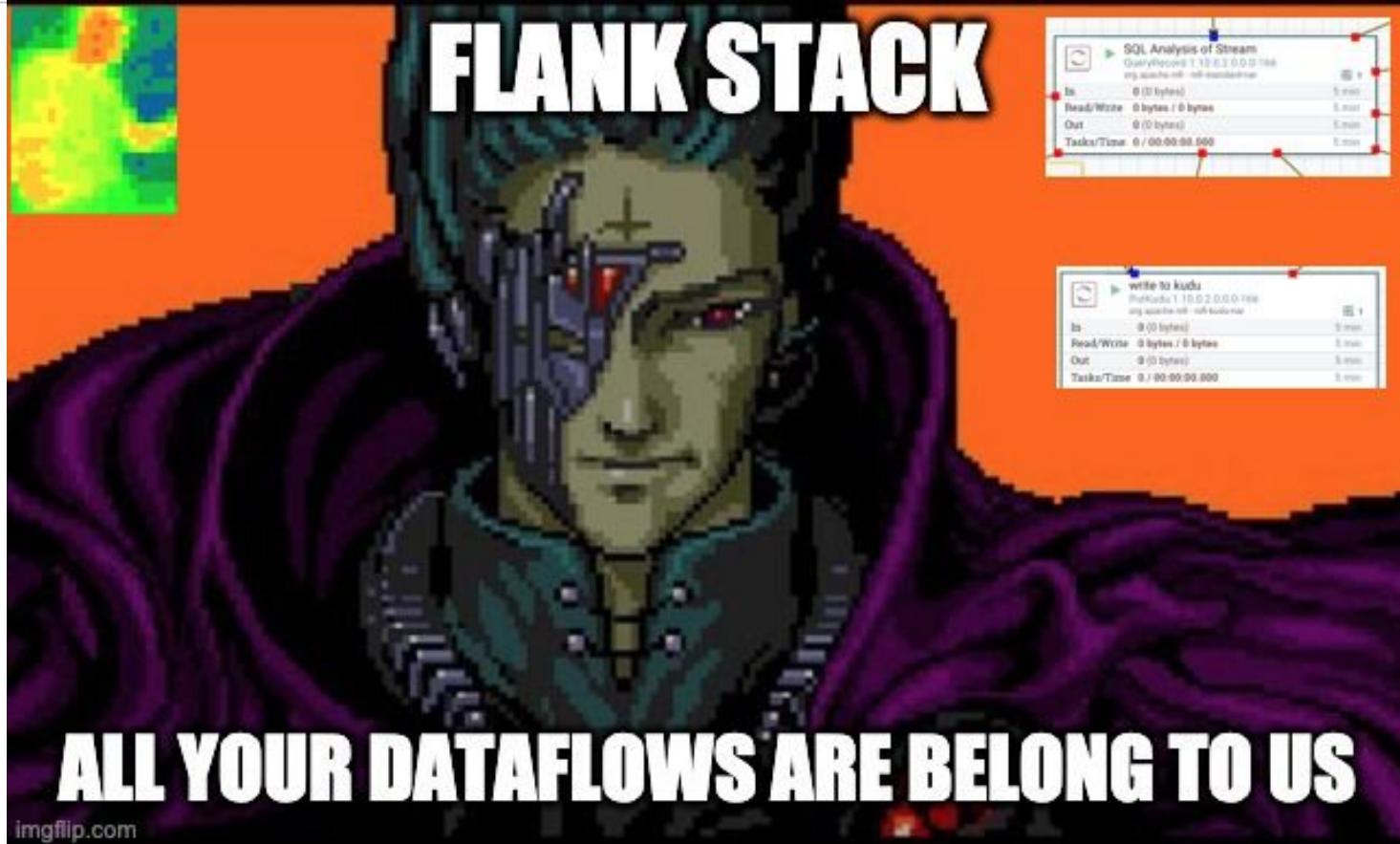
-- try to do this ksql!
SELECT us_west.user_score+ap_south.user_score
FROM kafka_in_zone_us_west us_west
FULL OUTER JOIN kafka_in_zone_ap_south ap_south
ON us_west.user_id = ap_south.user_id;
```

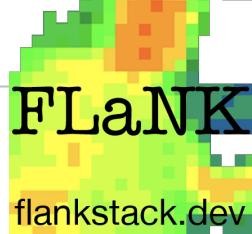
# Flink SQL

```
SELECT location, station_id, latitude, longitude, observation_time, weather, temperature_string,  
relative_humidity, wind_string, wind_dir, wind_degrees, wind_mph, pressure_in, dewpoint_string,  
dewpoint_f, dewpoint_c FROM weather2 WHERE location is not null and location <> 'null' and  
trim(location) <> "" and location like '%NJ'
```

```
SELECT HOP_END(eventTimestamp, INTERVAL '1' SECOND, INTERVAL '30' SECOND) as  
windowEnd, count("close") as closeCount, sum(cast("close" as float)) as closeSum, avg(cast("close" as  
float)) as closeAverage, min("close") as closeMin, max("close") as closeMax, sum(case when "close" >  
14 then 1 else 0 end) as stockGreaterThan14 FROM stocksraw GROUP BY HOP(eventTimestamp,  
INTERVAL '1' SECOND, INTERVAL '30' SECOND)
```





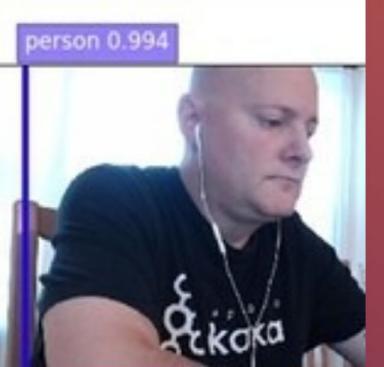


# Upcoming Events

April 27

DEVELOPERWEEK  
EUROPE

**Timothy Spann**  
Cloudera  
Principal DataFlow Field Engineer



person 0.994

May 19

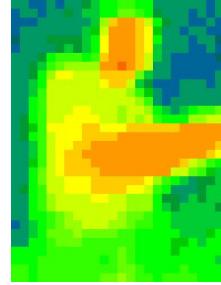
Live Demo Jam Expands

## THE LEADING-EDGE STREAMING DATA PLATFORM WITH NIFI, KAFKA, AND FLINK

May 19, 2021  
10:00 AM PT | 1:00 PM ET

Tim Spann  
Principal DataFlow Field Engineer, Cloudera





TH<sup>DATA</sup>N<sup>FLINK</sup> Y<sup>HIVE</sup> U<sup>SPARK</sup>

