

Self-Service Data Ingestion Using NiFi, StreamSets & Kafka

Guido Schmutz – 6.12.2017



@gschmutz



guidoschmutz.wordpress.com

BASEL • BERN • BRUGG • DÜSSELDORF • FRANKFURT A.M. • FREIBURG I.BR. • GENF
HAMBURG • KOPENHAGEN • LAUSANNE • MÜNCHEN • STUTTGART • WIEN • ZÜRICH

trivadis
makes **IT** easier. 

■ Guido Schmutz

Working at Trivadis for more than 21 years

Oracle ACE Director for Fusion Middleware and SOA



Consultant, Trainer Software Architect for Java, Oracle, SOA and Big Data / Fast Data



Head of Trivadis Architecture Board

Technology Manager @ Trivadis

More than 30 years of software development experience

Contact: guido.schmutz@trivadis.com

Blog: <http://guidoschmutz.wordpress.com>

Slideshare: <http://www.slideshare.net/gschmutz>

Twitter: [@gschmutz](https://twitter.com/gschmutz)

gschmutz 21:08 on April 18, 2017

Tags: flink (1), kafka (59), kafka-connect (4), kafka-streams (17), spark-streaming (53), storm (39), streams (4)

Last week in Stream Processing & Analytics – 18.4.2017

This is the 62nd edition of my blog series blog series around Stream Processing and Analytics!

Every week I'm also updating the following two lists with the presentations/videos of the current week:

- [Presentations from SlideShare](#)
- [Videos from YouTube](#)

As usual, find below the new blog articles, presentations, videos and software releases from last week:

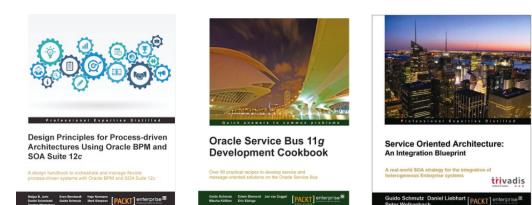
News and Blog Posts

General

- [Multi-Master Replication For Geo-Distributed Data: It's more than you think by Ellen Friedman](#)
- [Understanding Indicators of Attack \(IOAs\): The Power of Event Stream Processing in CrossStrike Falcon by Dan Brown](#)
- [Stream processing and messaging systems for the IoT age by Ben Lorica](#)

Apache Kafka / Kafka Streams / Confluent Platform

- [Creating a Data Pipeline with Kafka Connect API – from Architecture to Operations by Alexandra Wang](#)
- [Streaming Spring Boot Application Logs to ELK Stack—Part 1 by ksadayamuthu](#)
- [Streaming Spring Boot Application Logs to Apache Kafka—ELK/Kafka Stack—Part 2 by ksadayamuthu](#)

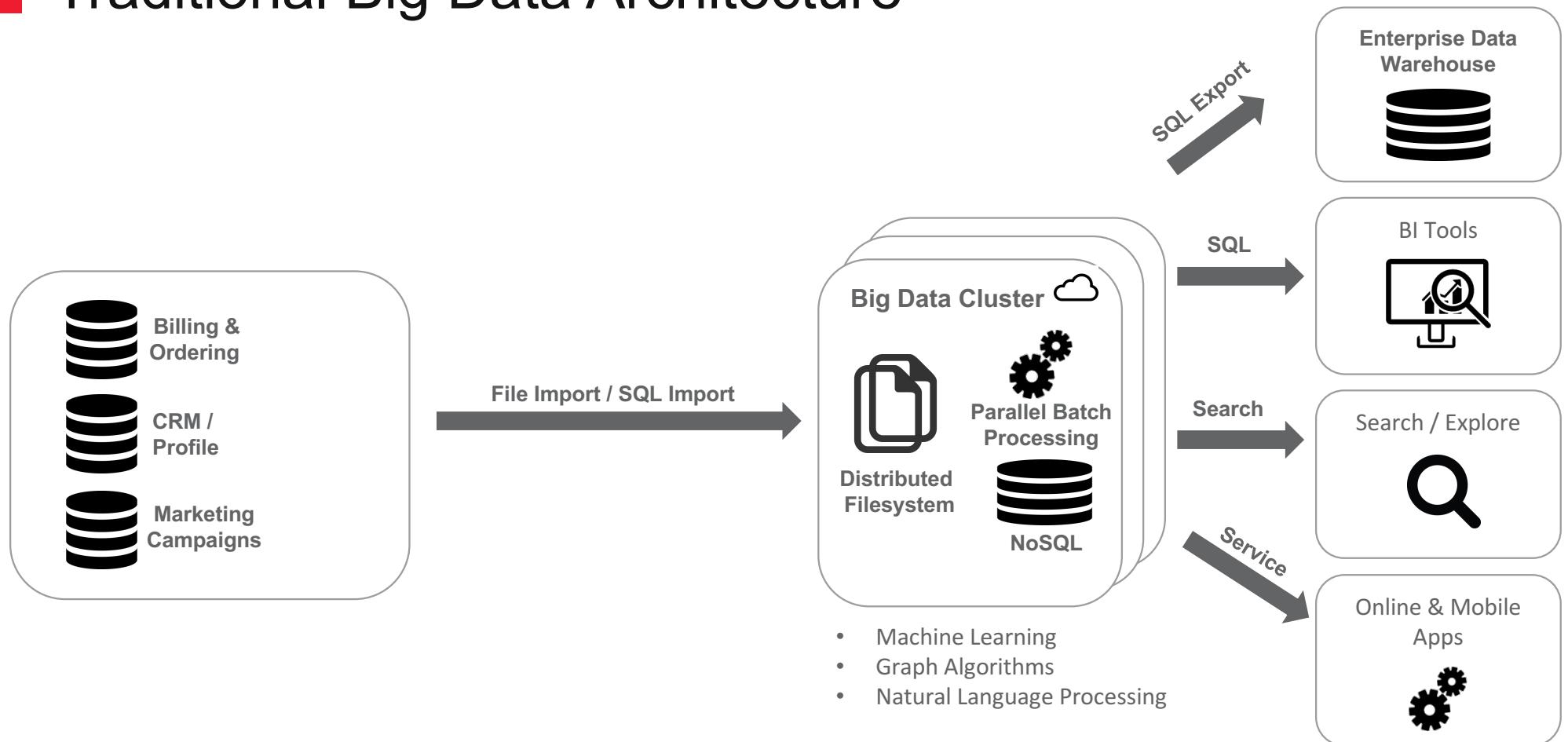


■ Agenda

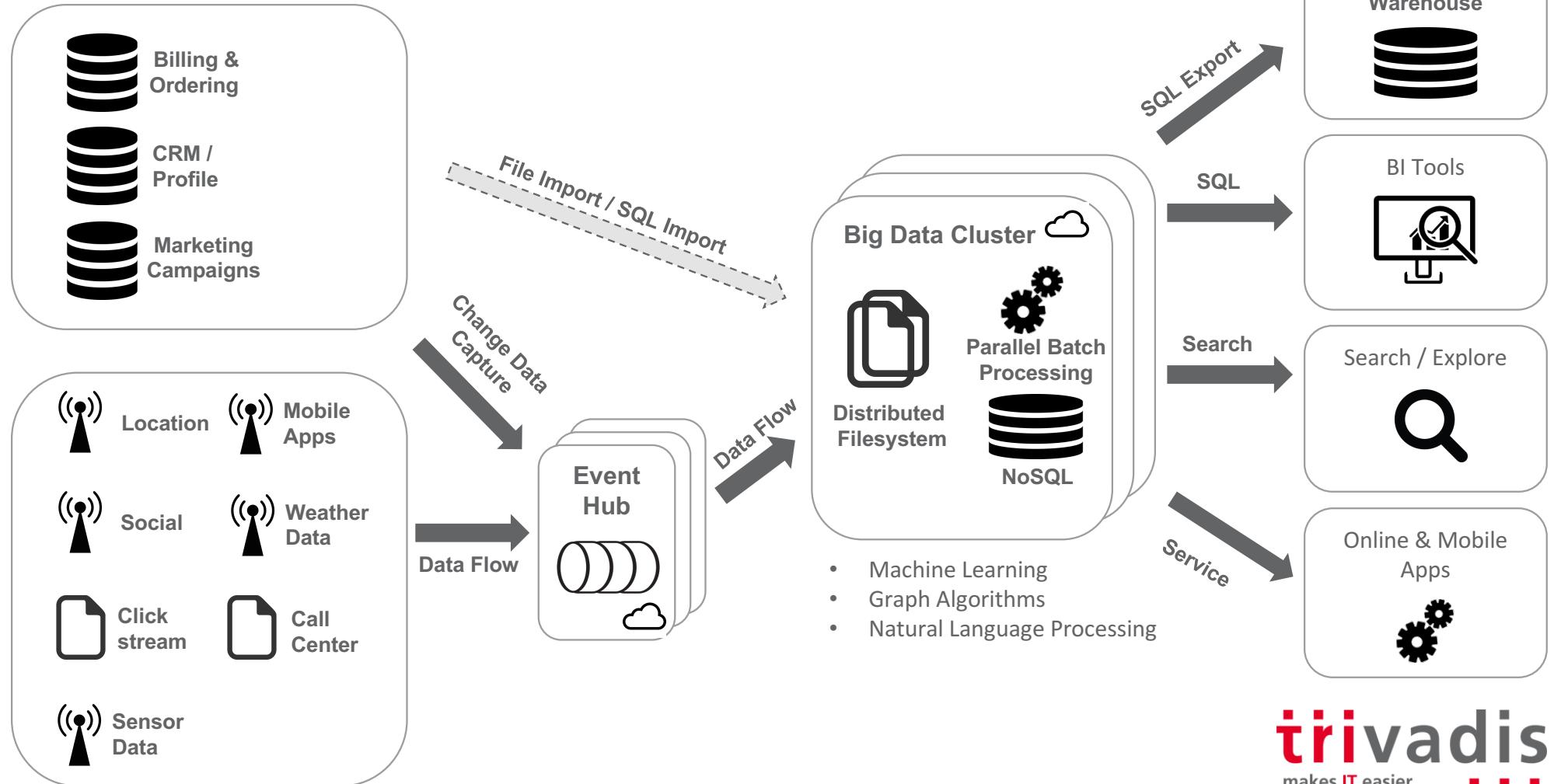
1. Data Flow Processing
2. Apache NiFi
3. StreamSets Data Collector
4. Kafka Connect
5. Summary

Data Flow Processing

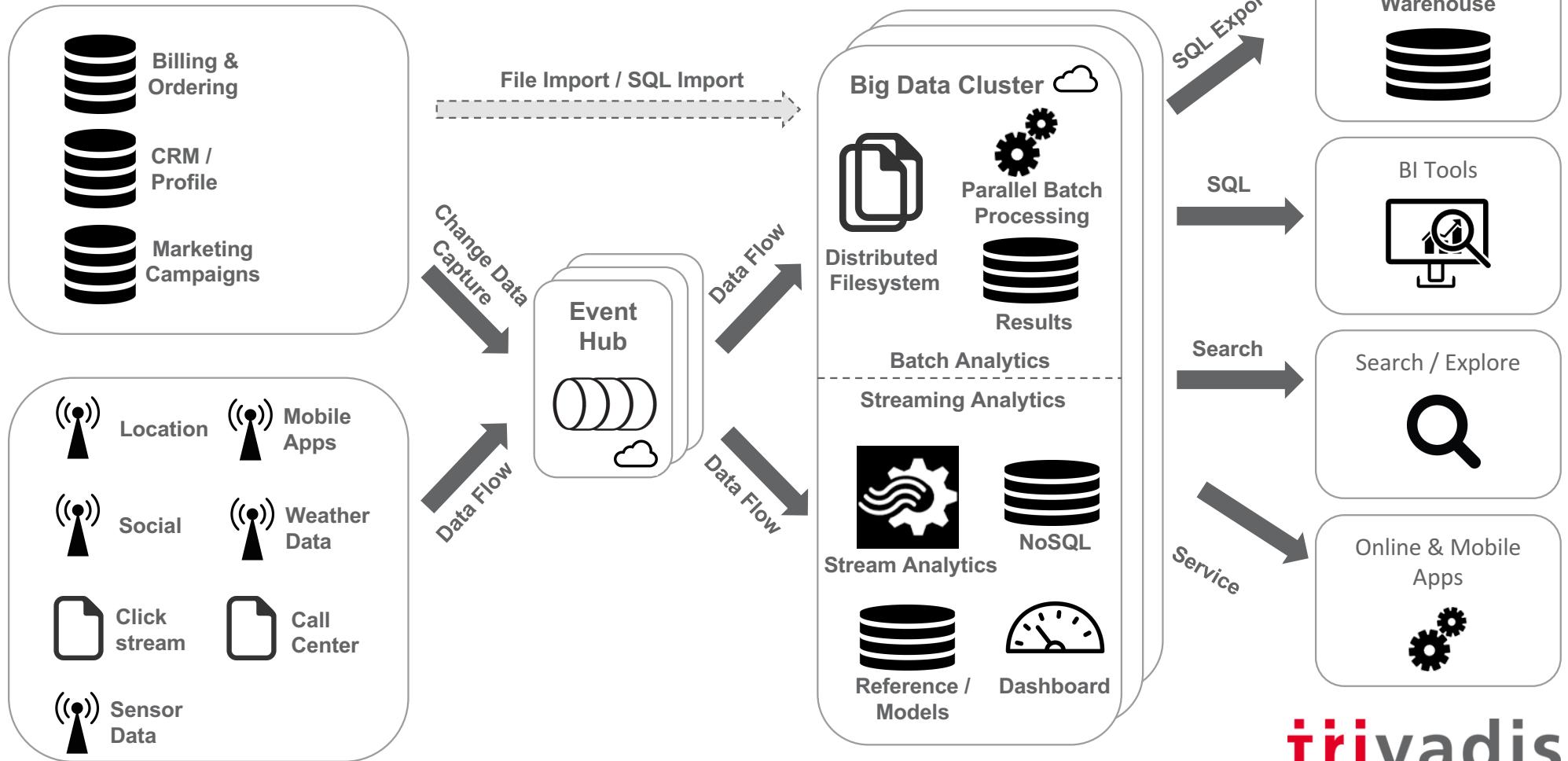
■ Traditional Big Data Architecture



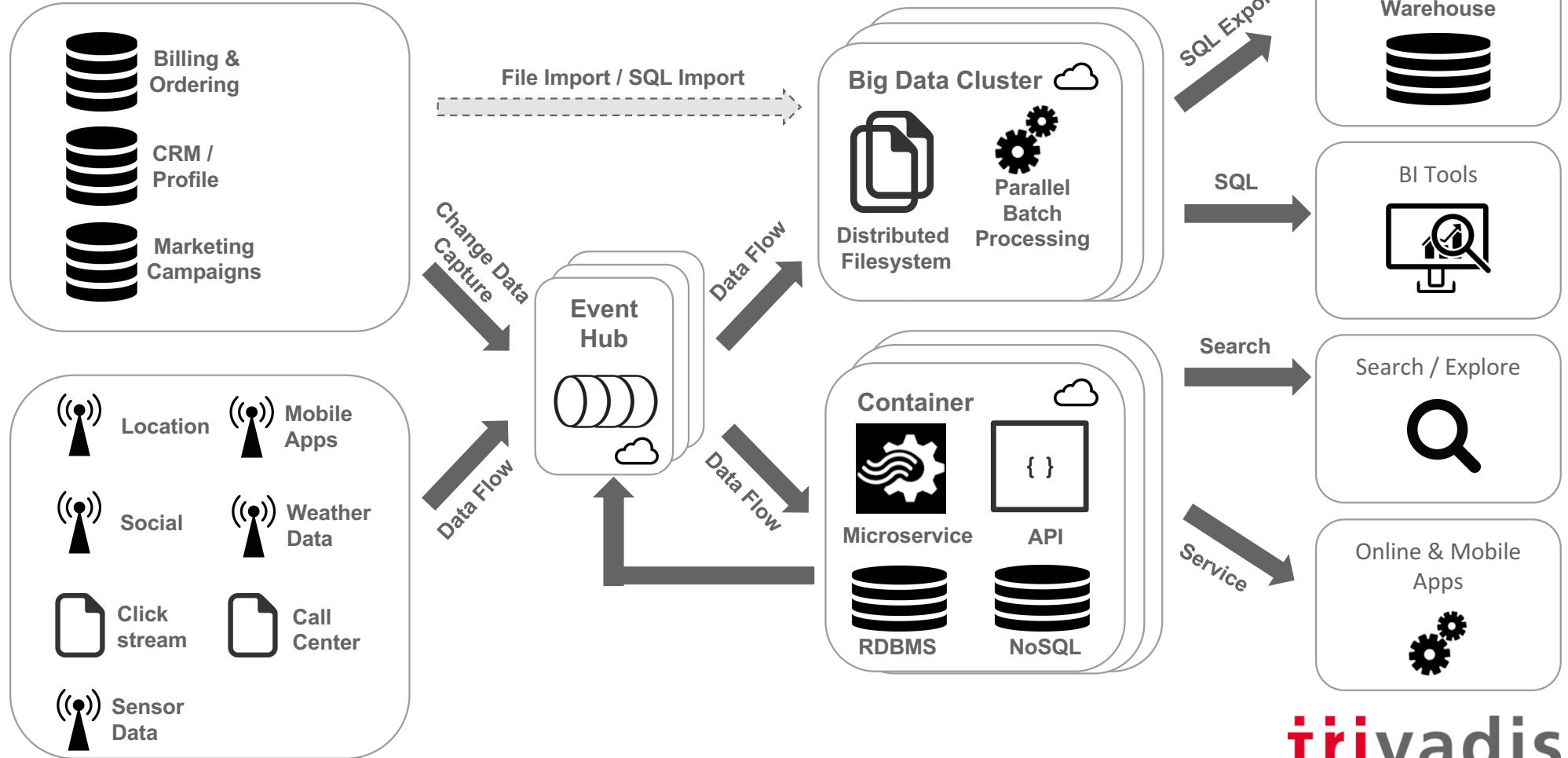
Event Hub – handle event stream data

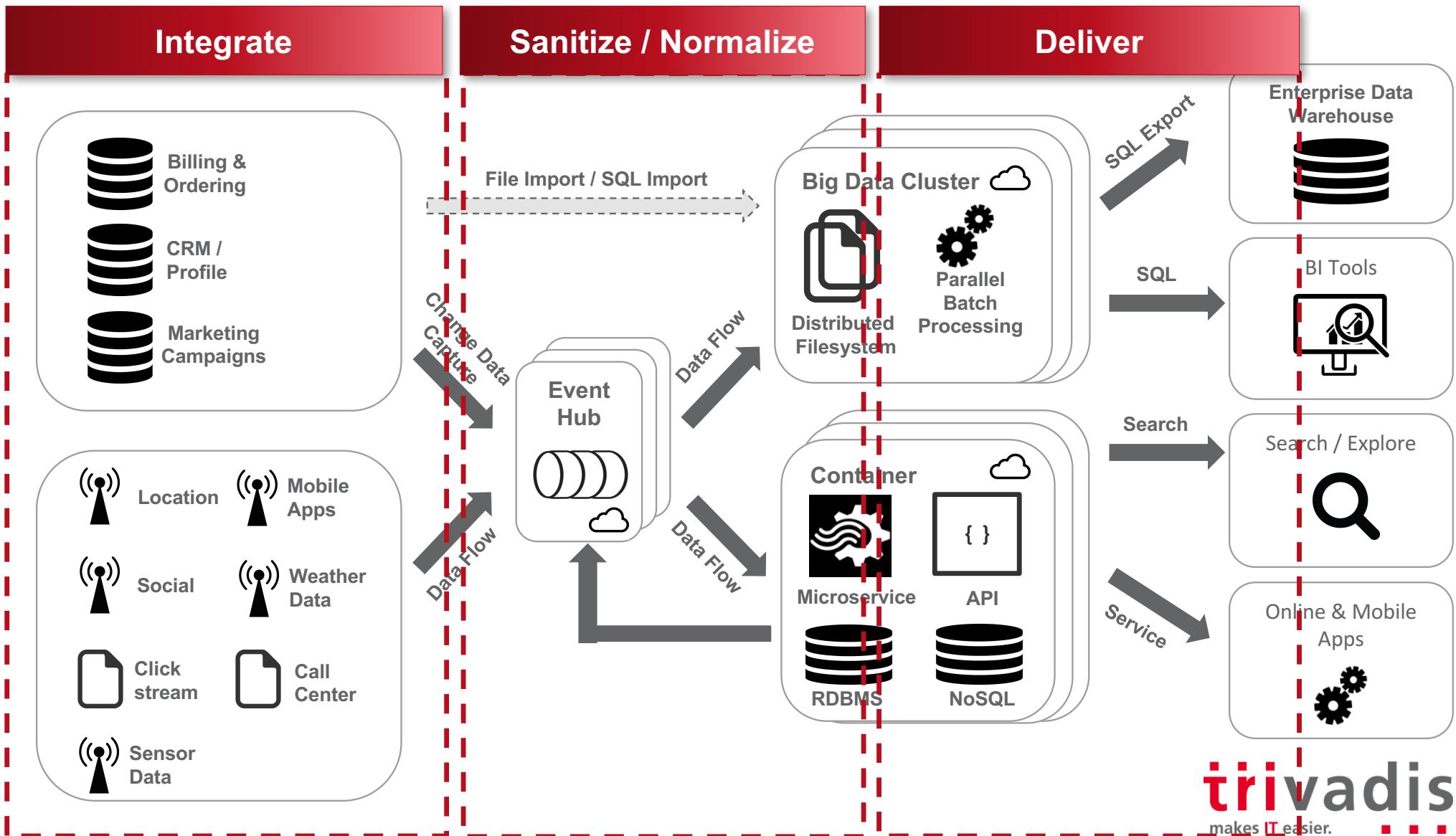


■ Event Hub – taking Velocity into account



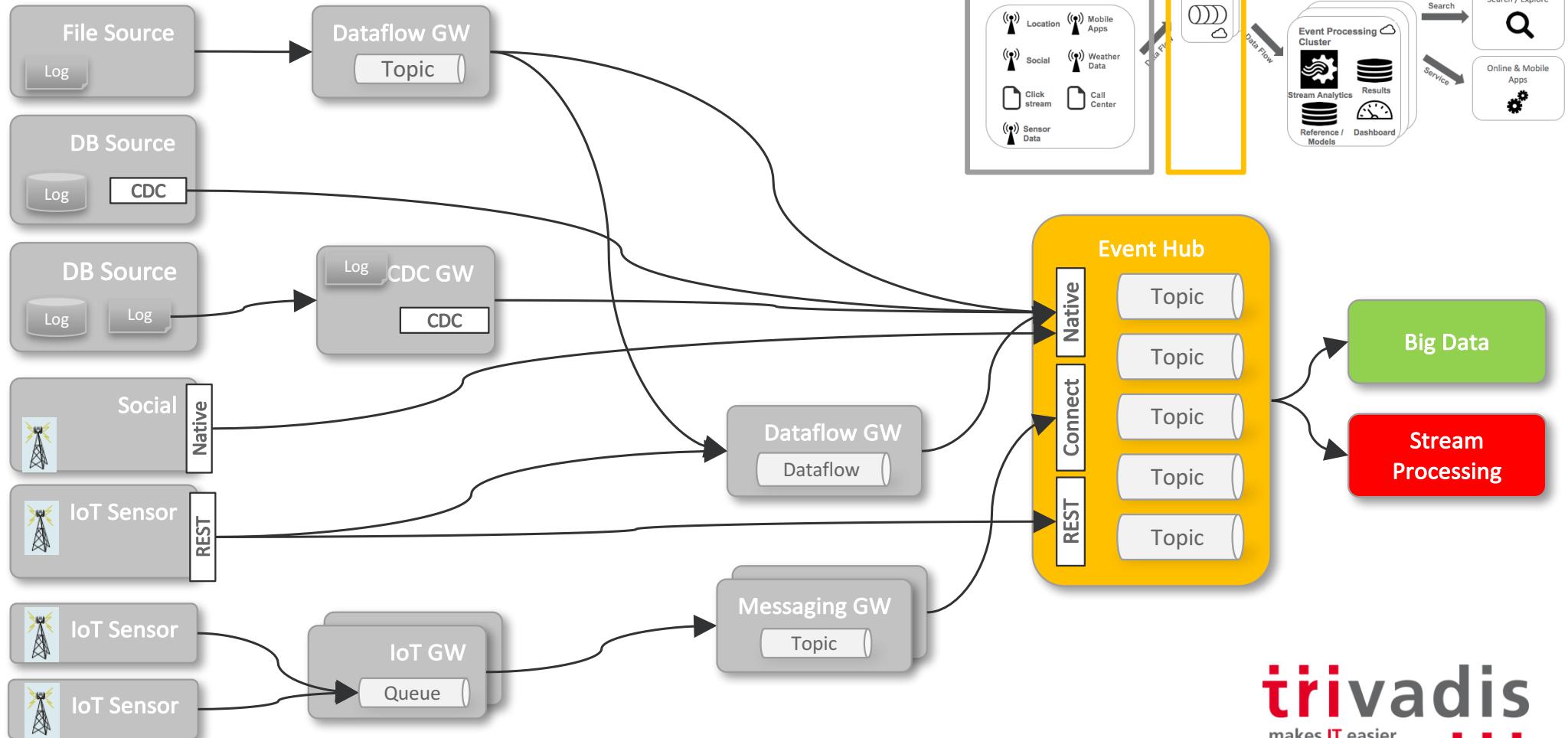
Event Hub – Asynchronous Microservice Architecture





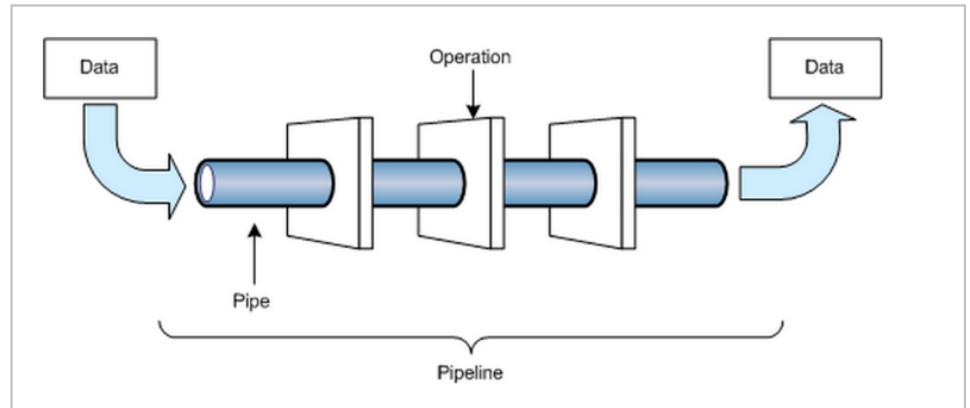
trivadis
makes IT easier.

Continuous Ingestion - DataFlow Pipelines

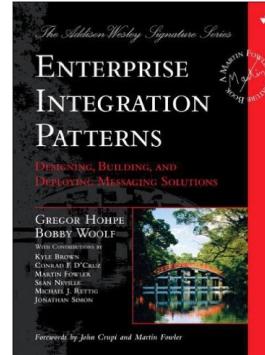
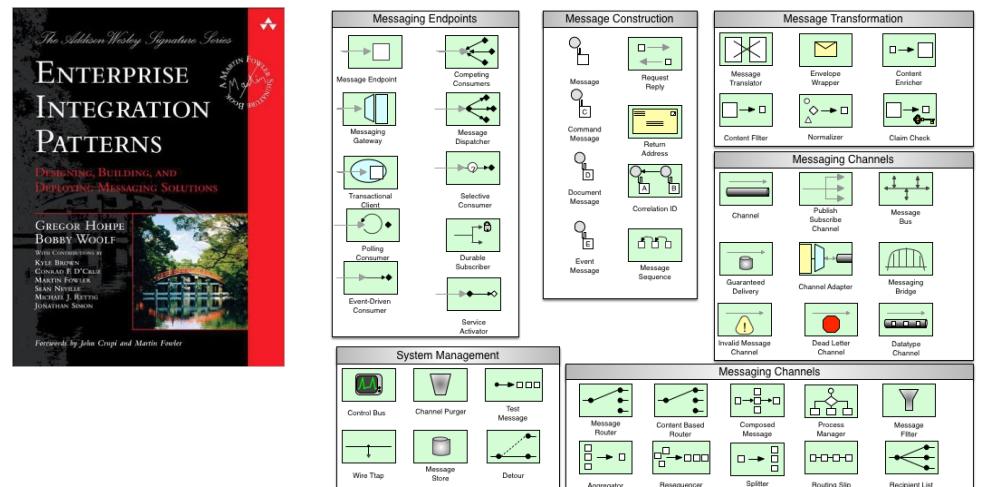


DataFlow Pipeline

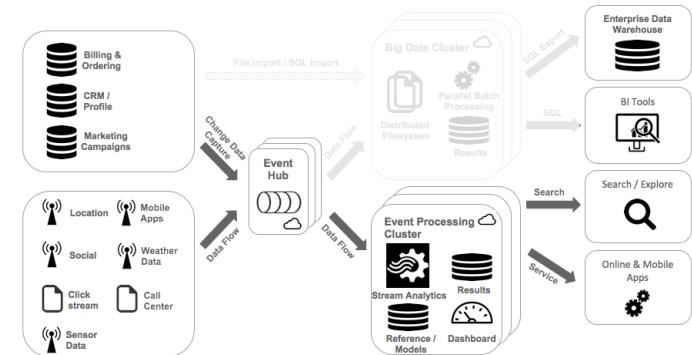
- Flow-based "programming"
- Ingest Data from various sources
- Extract – Transform – Load
- High-Throughput, straight-through data flows
- Data Lineage
- Batch- or Stream-Processing
- Visual coding with flow editor
- Event Stream Processing (ESP) **but not** Complex Event Processing (CEP)



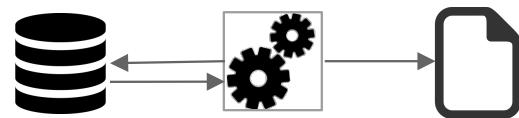
Source: Confluent



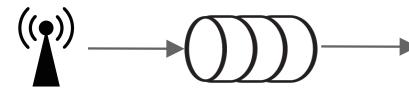
Continuous Ingestion – Integrating data sources



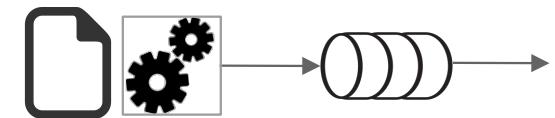
SQL Polling



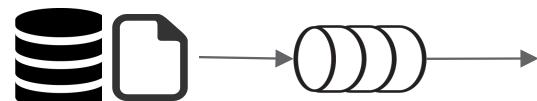
Sensor Stream



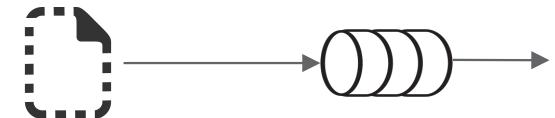
File Polling (File Tailing)



Change Data Capture (CDC)



File Stream (Appender)



Ingestion with/without Transformation?

Zero Transformation

- No transformation, plain ingest, no schema validation
- Keep the original format – Text, CSV, ...
- Allows to store data that may have errors in the schema

Format Transformation

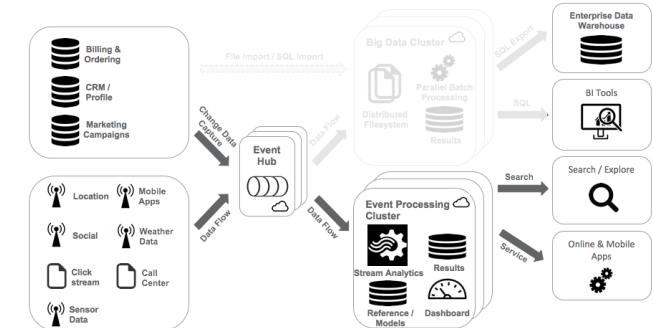
- Prefer name of Format Translation
- Simply change the format
- Change format from Text to Avro
- Does schema validation

Enrichment Transformation

- Add new data to the message
- Do not change existing values
- Convert a value from one system to another and add it to the message

Value Transformation

- Replaces values in the message
- Convert a value from one system to another and change the value in-place
- Destroys the raw data!



■ Why is Data Ingestion Difficult?

Infrastructure Drift

Physical and Logical Infrastructure changes rapidly

Key Challenges:

Infrastructure Automation
Edge Deployment

Structure Drift

Data Structures and formats evolve and change unexpectedly

Key Challenges:

Consumption Readiness
Corruption and Loss

Semantic Drift

Data semantics change with evolving applications

Key Challenges

Timely Intervention
System Consistency

Source: Streamsets

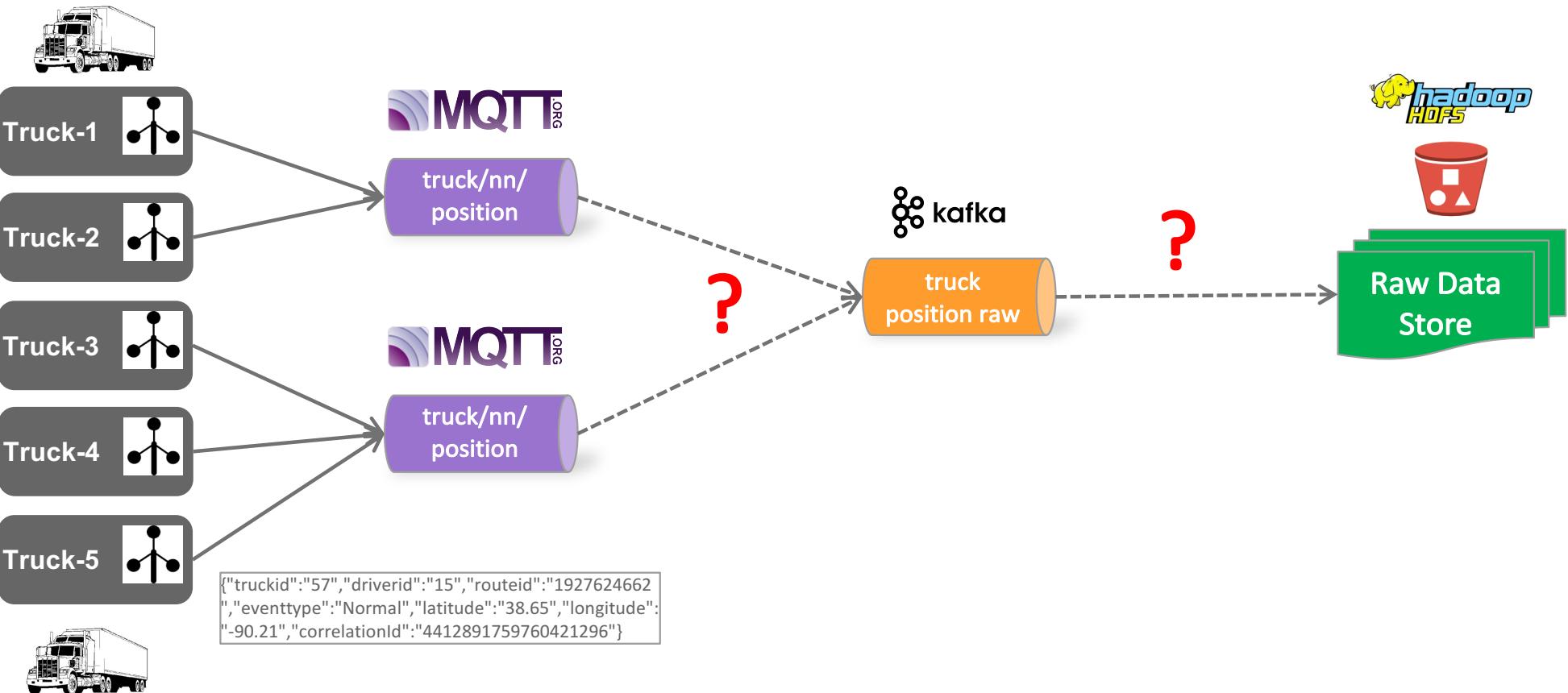
■ Challenges for Ingesting Sensor Data

- Multitude of sensors
- Real-Time Streaming
- Multiple Firmware versions
- Bad Data from damaged sensors
- Regulatory Constraints
- Data Quality

```
{  
  "firmware": "1.0",  
  "device_id": "938AC7",  
  "container_id": "AD34BA42",  
  "reading_date": "1463847641",  
  "temperature": "1.5",  
  "temp_unit": "C",  
  "relative_humidity": "56"  
}  
  
{  
  "firmware": "2.0",  
  "device_id": "934DC7",  
  "container_id": "AD34BA42",  
  "reading_date": "1463867641",  
  "temperature": "1.5",  
  "temp_unit": "C",  
  "relative_humidity": "56",  
  "orientation": {  
    "roll": "-1",  
    "pitch": "3",  
    "yaw": "160"  
  }  
}  
  
{  
  "firmware": "3.0",  
  "device_id": "824DC7",  
  "container_id": "AD34BA42",  
  "reading_date": "1463967641",  
  "temperature": "1.5",  
  "temp_unit": "C",  
  "relative_humidity": "56",  
  "orientation": {  
    "roll": "-1",  
    "pitch": "3",  
    "yaw": "160"  
  },  
  "location": {  
    "lat": "37.354108",  
    "lon": "-121.955236"  
  }  
}
```

Source: Cloudera

Demo Case

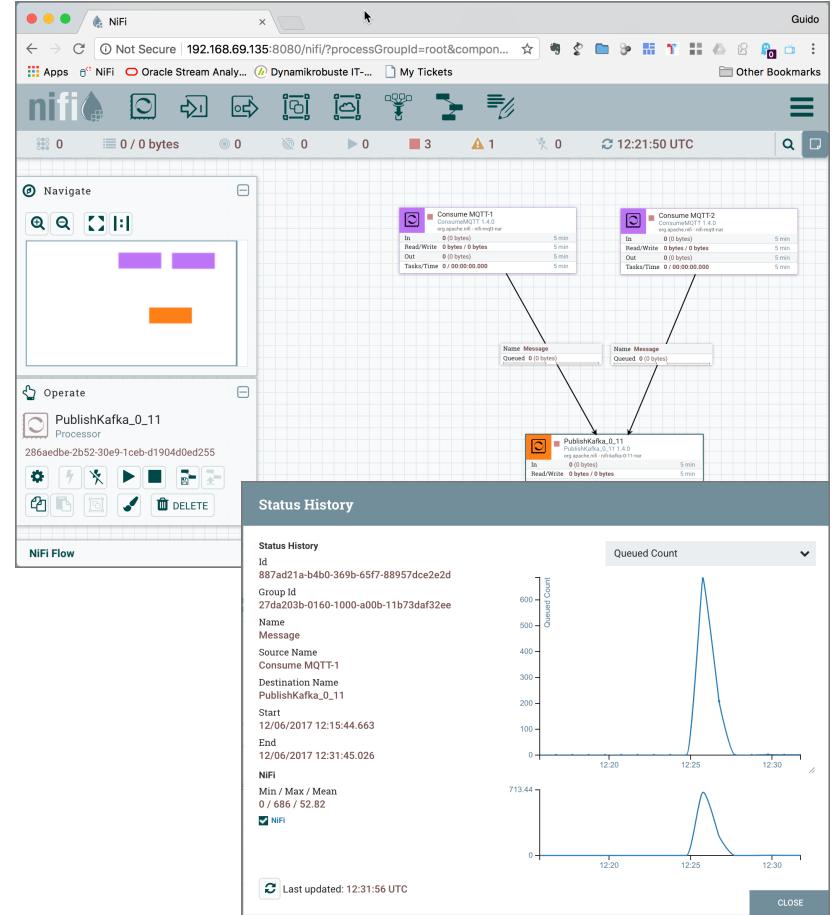


Apache NiFi



Apache NiFi

- Originated at NSA as Niagarafiles – developed behind closed doors for 8 years
- Open sourced December 2014, Apache Top Level Project July 2015
- Look-and-Feel modernized in 2016
- Opaque, “file-oriented” payload
- Distributed system of processors with centralized control
- Based on *flow-based programming* concepts
- Data Provenance and Data Lineage
- Web-based user interface



trivadis
makes IT easier.



■ Processors for Source and Sink

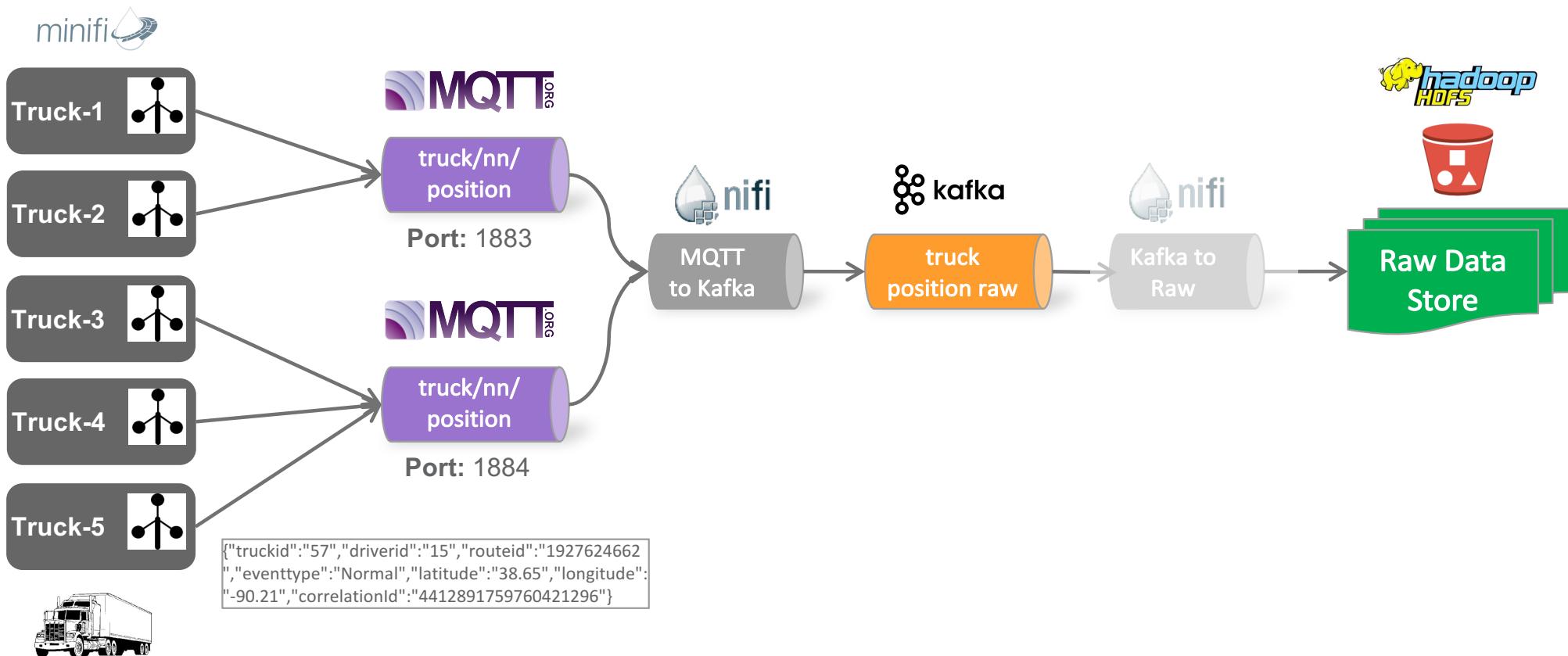
- **ConsumeXXXX** (AMQP, EWS, IMAP, JMS, Kafka, MQTT, POP3, ...)
- **DeleteXXXX** (DynamoDB, Elasticsearch, HDFS, RethinkDB, S3, SQS, ...)
- **FetchXXXX** (AzureBlobStorage, ElasticSearch, File, FTP, HBase, HDFS, S3 ...)
- **ExecuteXXXX** (FlumeSink, FlumeSource, Script, SQL, ...)
- **GetXXXX** (AzureEventHub, Couchbase, DynamoDB, File, FTP, HBase, HDFS, HTTP, Ignite, JMSQueue, JMSTopic, Kafka, Mongo, Solr, Splunk, SQS, TCP, ...)
- **ListenXXXX** (HTTP, RELP, SMTP, Syslog, TCP, UDP, WebSocket, ...)
- **PublishXXXX** (Kafka, MQTT)
- **PutXXXX** (AzureBlobStorage, AzureEventHub, CassandraQL, CloudWatchMetric, Couchbase, DynamoDB, Elasticsearch, Email, FTP, File, Hbase, HDFS, HiveQL, Kudu, Lambda, Mongo, Parquet, Slack, SQL, TCP,)
- **QueryXXXX** (Cassandra, DatabaseTable, DNS, Elasticserach)



■ Processors for Processing

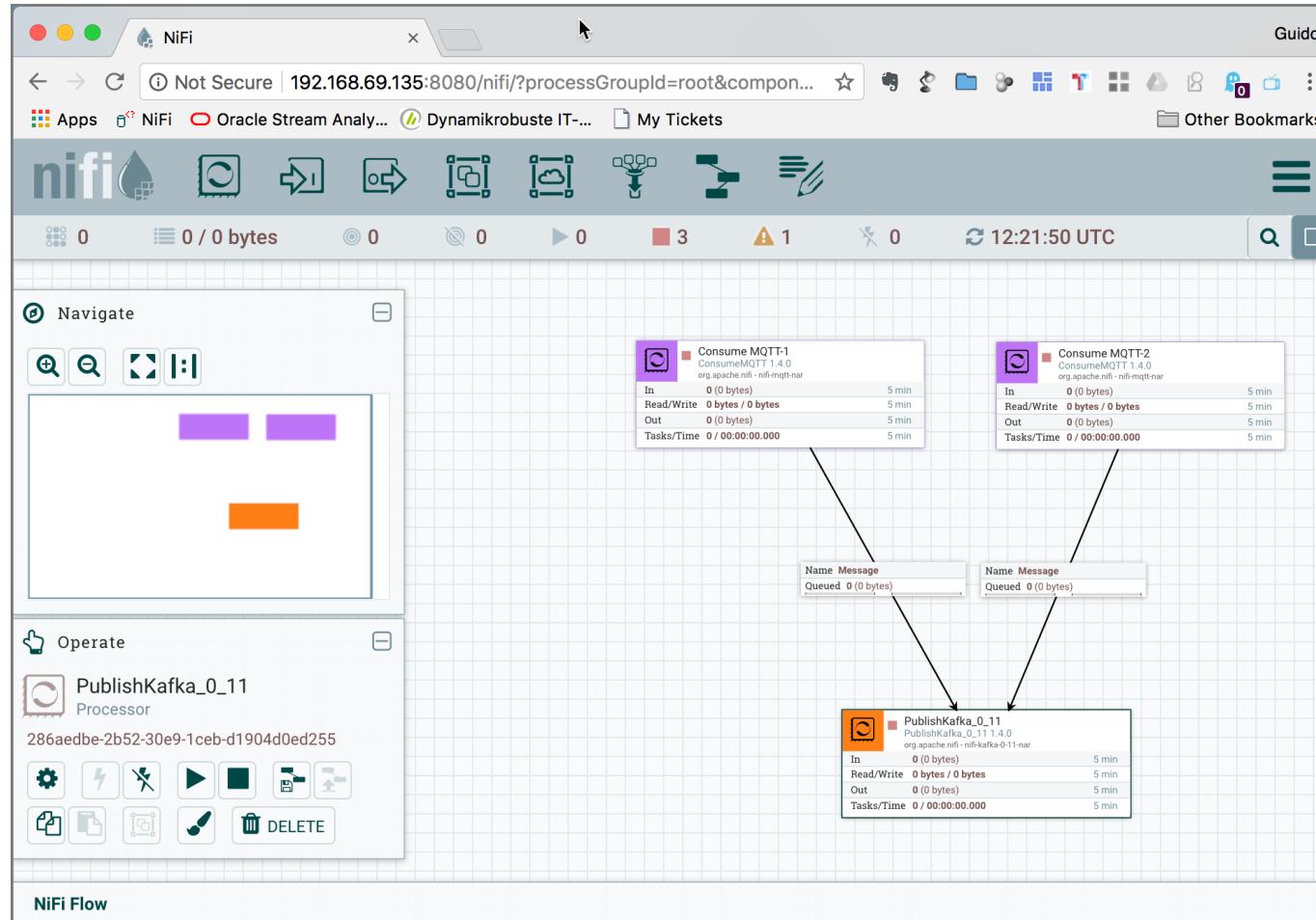
- **ConvertXxxxToYyyy**
- **ConvertRecord**
- **EnforceOrder**
- **EncryptContent**
- **ExtractXXXX** (AvroMetadata,
EmailAttachments, Grok,
HL7Attributes, ImageMetadata, ...)
- **GeoEnrichIP**
- **JoltTransformJSON**
- **MergeContent**
- **ReplaceText**
- **ResizeImage**
- **SplitXXXX** (Avro, Content, JSON,
Record, Xml, ...)
- **TailFile**
- **TransformXML**
- **UpdateAttribute**

Demo Case





Demo: Dataflow for MQTT to Kafka





Demo: MQTT Processor

Configure Processor

SETTINGS SCHEDULING PROPERTIES COMMENTS

Required field

Property	Value
Broker URI	tcp://192.168.69.135:1883
Client ID	consumer1
Username	No value set
Password	No value set
SSL Context Service	No value set
Last Will Topic	No value set
Last Will Message	No value set
Last Will Retain	No value set
Last Will QoS Level	No value set
Session state	Clean Session
MQTT Specification Version	AUTO
Connection Timeout (seconds)	30
Keep Alive Interval (seconds)	60
Topic Filter	truck/+/position

CANCEL APPLY



Demo: Kafka Processor

Configure Processor

SETTINGS SCHEDULING PROPERTIES COMMENTS

Required field

Property	Value
Kafka Brokers	192.168.69.135:9092
Security Protocol	PLAINTEXT
Kerberos Service Name	No value set
Kerberos Principal	No value set
Kerberos Keytab	No value set
SSL Context Service	No value set
Topic Name	truck_position
Delivery Guarantee	Guarantee Replicated Delivery
Use Transactions	false
Attributes to Send as Headers (Regex)	No value set
Message Header Encoding	UTF-8
Kafka Key	No value set
Key Attribute Encoding	UTF-8 Encoded
Message Demarcator	No value set

CANCEL APPLY



Demo: Masking Field with ReplaceText Processor

Configure Processor

SETTINGS SCHEDULING PROPERTIES COMMENTS

Required field

Property	Value
Search Value	"driverid":"(.*?)"
Replacement Value	"driverid":xxxxx"
Character Set	UTF-8
Maximum Buffer Size	1 MB
Replacement Strategy	Regex Replace
Evaluation Mode	Entire text

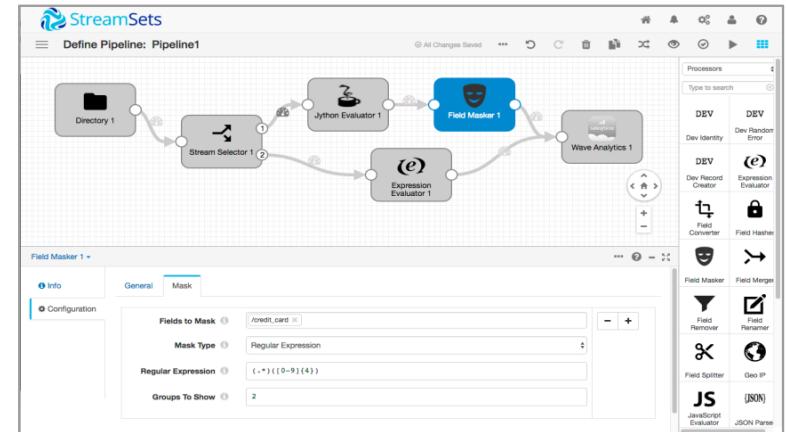
CANCEL APPLY

StreamSets

StreamSets Data Collector



- Founded by ex-Cloudera, Informatica employees
- Continuous open source, intent-driven, big data ingest
- Visible, record-oriented approach fixes combinatorial explosion
- Batch or stream processing
 - Standalone, Spark cluster, MapReduce cluster
- IDE for pipeline development by ‘civilians’
- Relatively new - first public release September 2015
- So far, vast majority of commits are from StreamSets staff



■ StreamSets Origins

An origin stage represents the source for the pipeline. You can use a single origin stage in a pipeline

Origins on the right are available out of the box

API for writing custom origins



Source: <https://streamsets.com/connectors>

■ StreamSets Processors

A processor stage represents a type of data processing that you want to perform

use as many processors in a pipeline as you need

Programming languages supported

- Java
- JavaScript
- Jython
- Groovy
- [Java Expression Language \(EL\)](#)
- [Spark](#)

Some of processors available out-of-the-box:

- Expression Evaluator
- Field Flattener
- Field Hasher
- Field Masker
- Field Merger
- Field Order
- Field Splitter
- Field Zip
- Groovy Evaluator
- JDBC Lookup
- JSON Parser
- Spark Evaluator
- ...

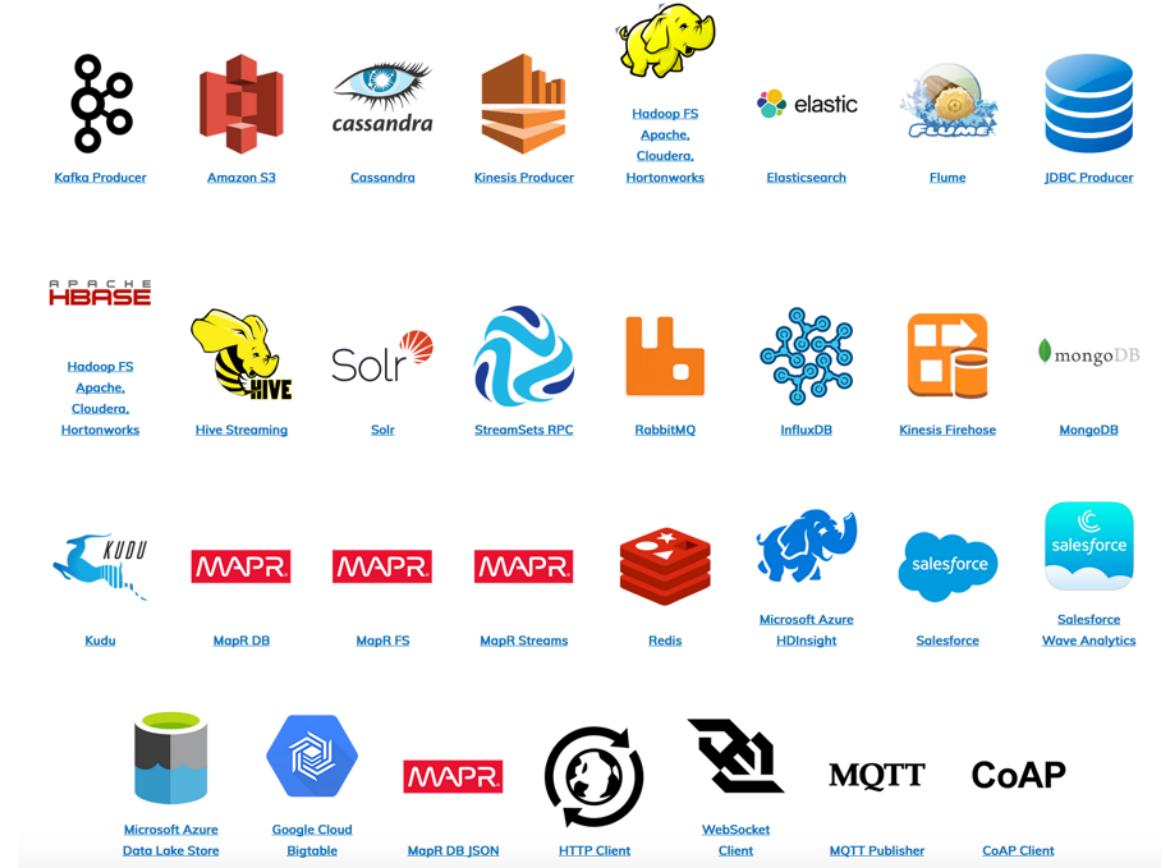
■ StreamSets Destinations



A destination stage represents the target for a pipeline. You can use one or more destinations in a pipeline

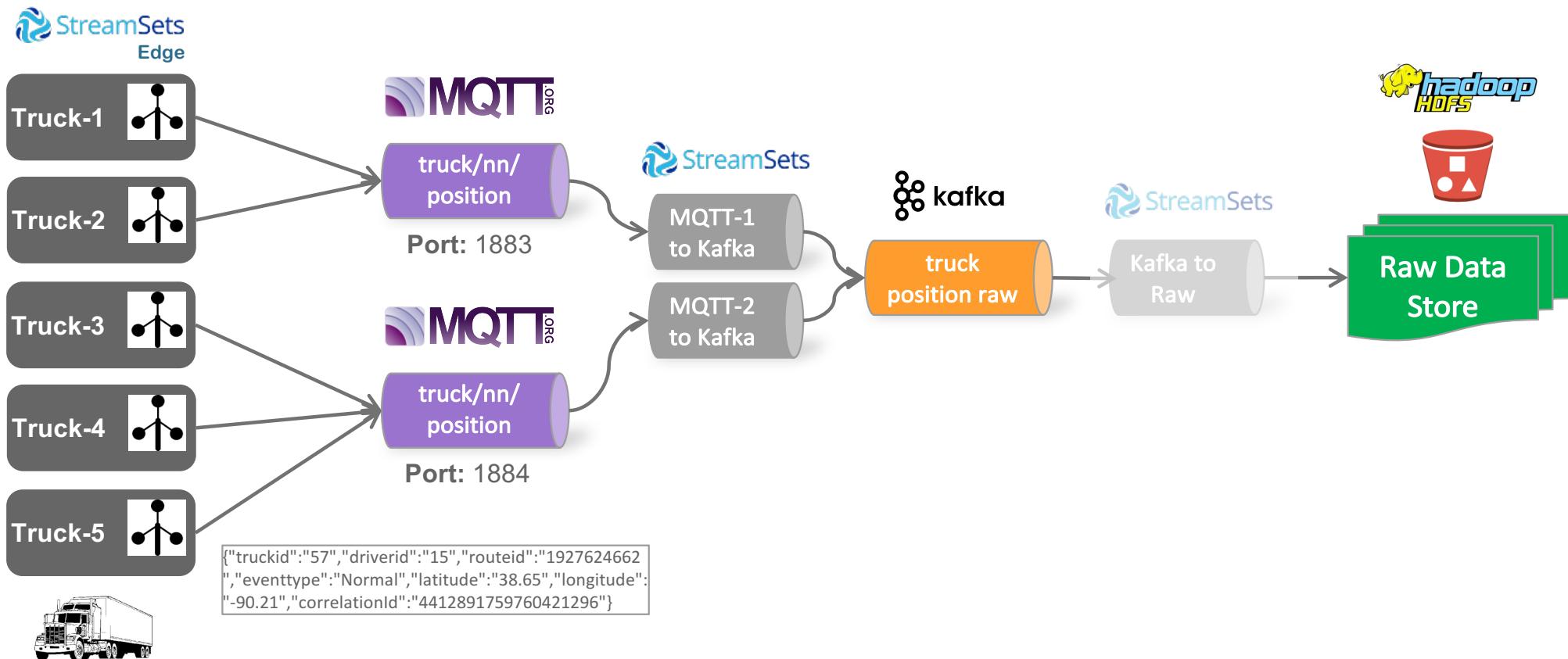
Destinations on the right are available out of the box

API for writing custom origins

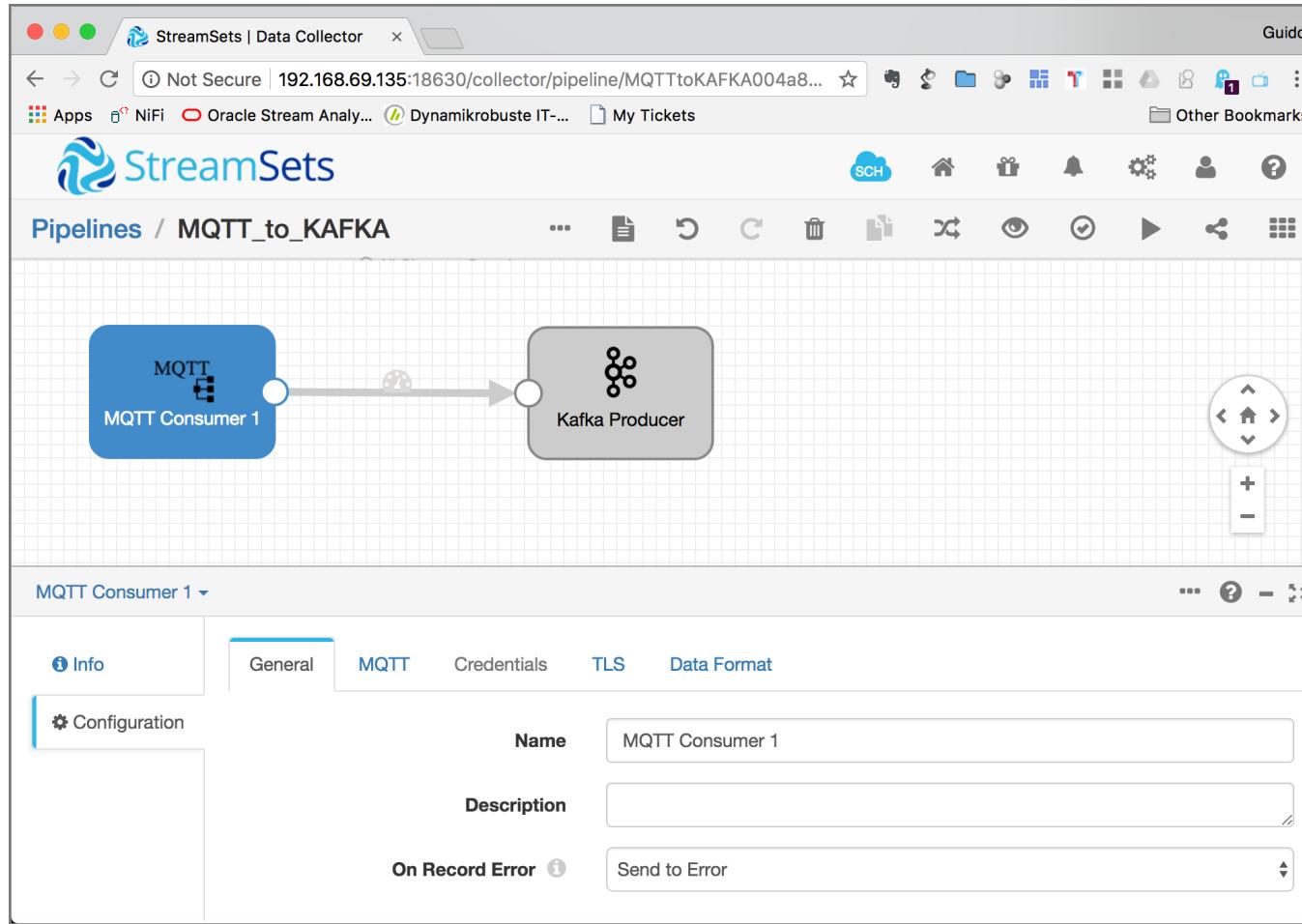


Source: <https://streamsets.com/connectors>

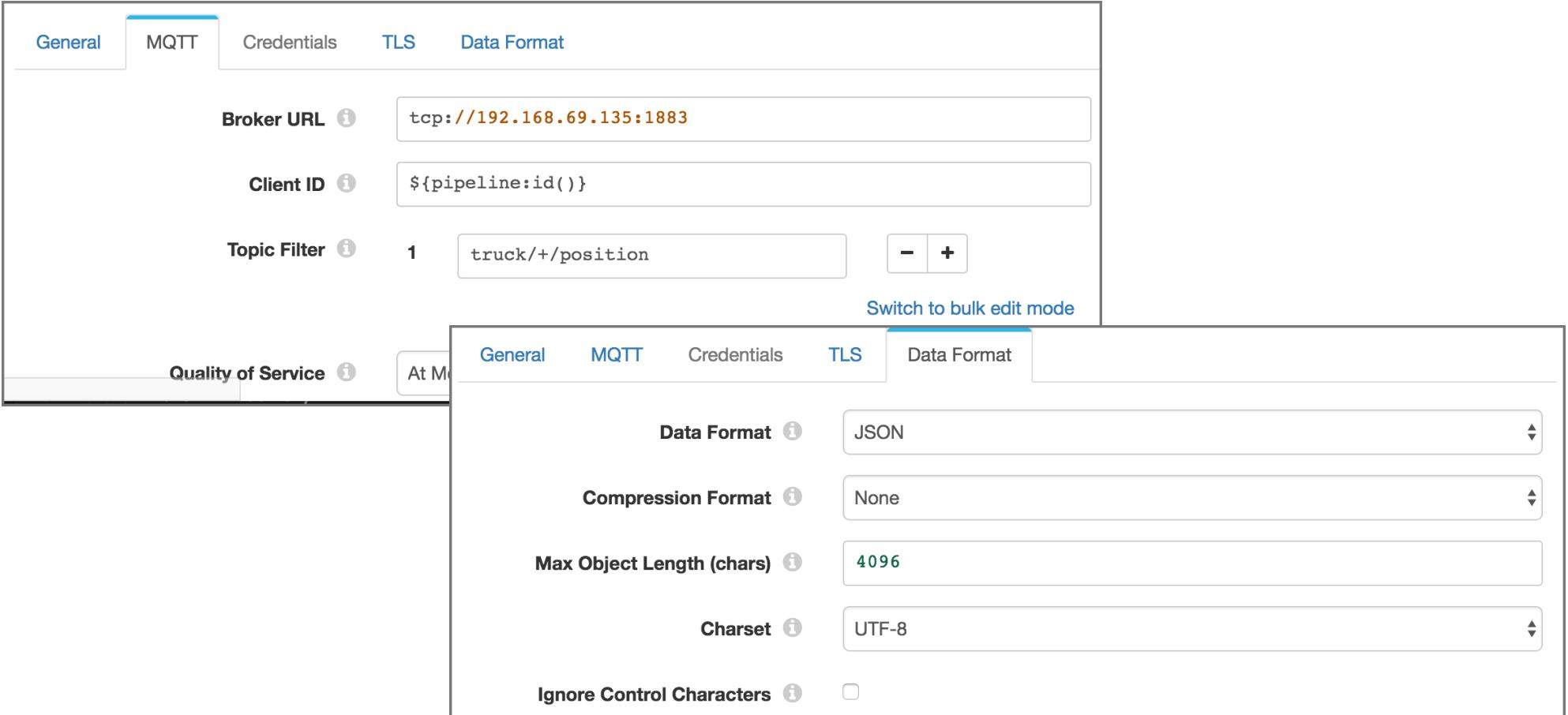
Demo Case



Demo: Dataflow for MQTT to Kafka



Demo: MQTT Source



The screenshot shows the configuration interface for an MQTT Source component in StreamSets Data Pipeline. The main configuration pane has tabs for General, MQTT (selected), Credentials, TLS, and Data Format. The MQTT tab contains fields for Broker URL (tcp://192.168.69.135:1883), Client ID (\${pipeline:id()}), and Topic Filter (1 truck/+/position). A 'Switch to bulk edit mode' link is visible. A modal dialog is open over the main pane, showing the Data Format tab with options for Data Format (JSON), Compression Format (None), Max Object Length (chars) (4096), Charset (UTF-8), and Ignore Control Characters (unchecked).

General MQTT Credentials TLS Data Format

Broker URL [i](#) tcp://192.168.69.135:1883

Client ID [i](#) \${pipeline:id()}

Topic Filter [i](#) 1 truck/+/position [-](#) [+](#)

Switch to bulk edit mode

Quality of Service [i](#) At Most Once

General MQTT Credentials TLS Data Format

Data Format [i](#) JSON

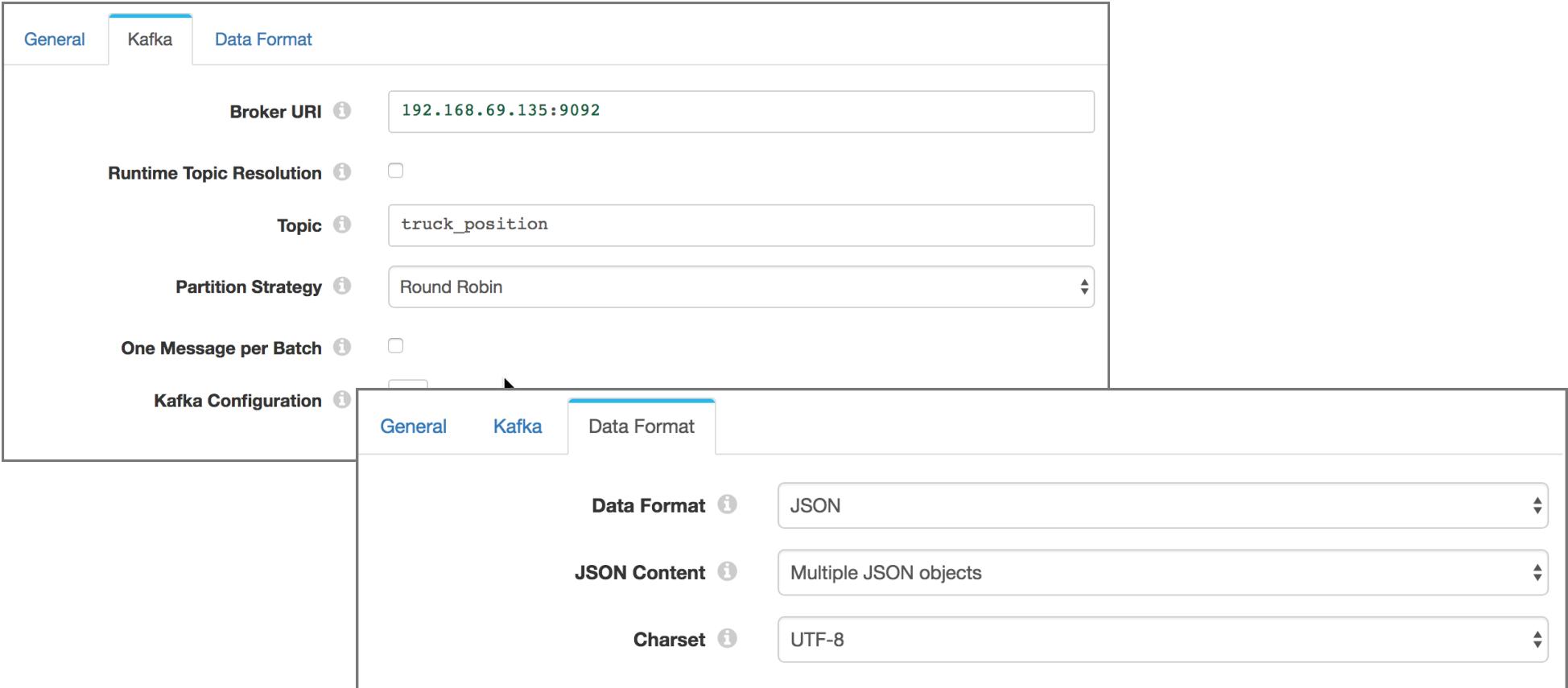
Compression Format [i](#) None

Max Object Length (chars) [i](#) 4096

Charset [i](#) UTF-8

Ignore Control Characters [i](#)

Demo: Kafka Sink



The screenshot shows the configuration interface for a Kafka Sink component in StreamSets Data Pipeline. The main window has tabs for General, Kafka, and Data Format, with the Kafka tab selected. The Data Format tab is expanded, showing sub-options for Data Format, JSON Content, and Charset.

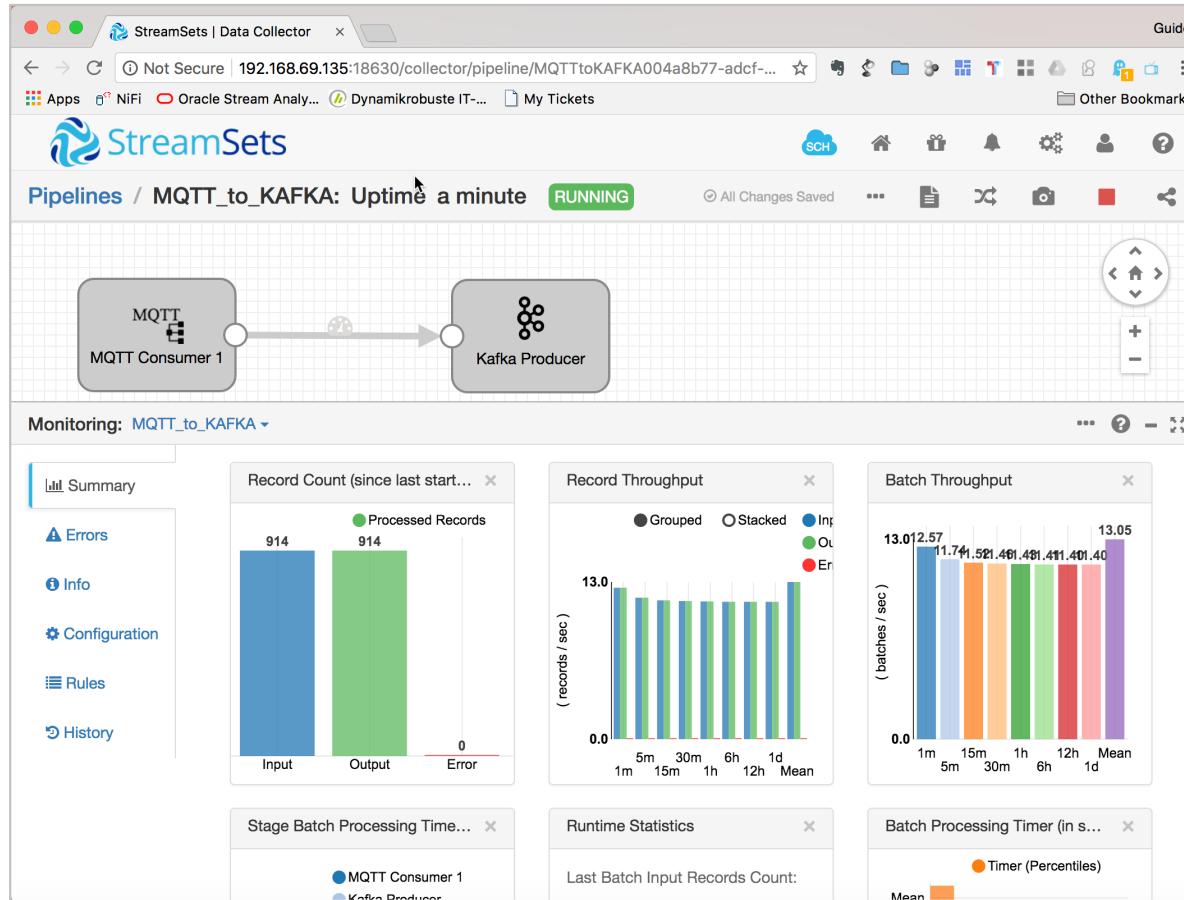
Kafka Configuration:

- Broker URI:** 192.168.69.135:9092
- Runtime Topic Resolution:** (checkbox)
- Topic:** truck_position
- Partition Strategy:** Round Robin
- One Message per Batch:** (checkbox)

Data Format Configuration:

- Data Format:** JSON
- JSON Content:** Multiple JSON objects
- Charset:** UTF-8

Demo: Dataflow for MQTT to Kafka



Demo: Masking fields



General Mask

Fields to Mask i /driverid x - +

Mask Type i Variable length ▼

[Switch to bulk edit mode](#)



Demo: Sending Message to Kafka in Avro



General Kafka Data Format

Data Format ⓘ Avro

Avro Schema Location ⓘ Confluent Schema Registry

Schema Registry URLs ⓘ 1 - +

[Switch to bulk edit mode](#)

Lookup Schema By ⓘ Subject

Schema Subject ⓘ truck-movement-v1.0

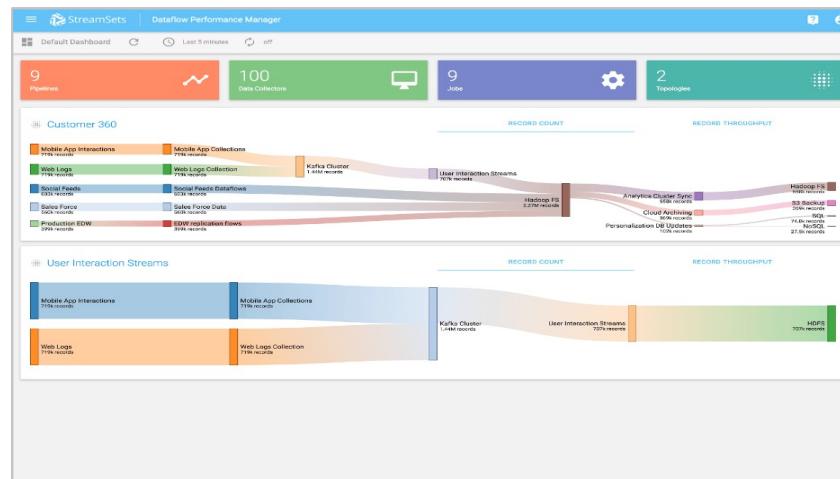
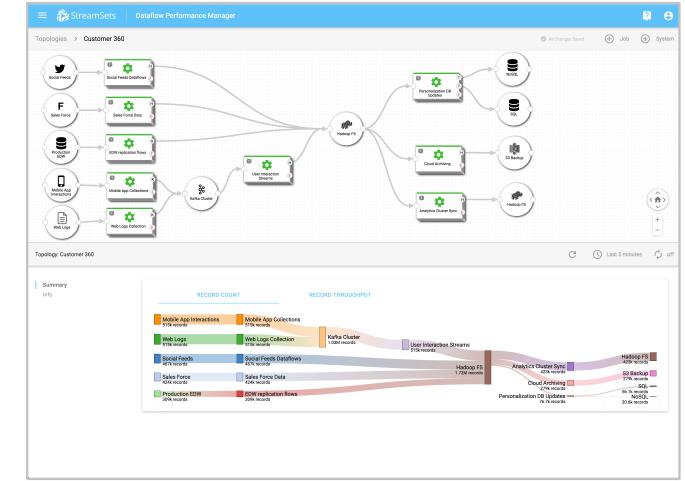
Include Schema ⓘ

Avro Compression Codec ⓘ Snappy

StreamSets Dataflow Performance Manager



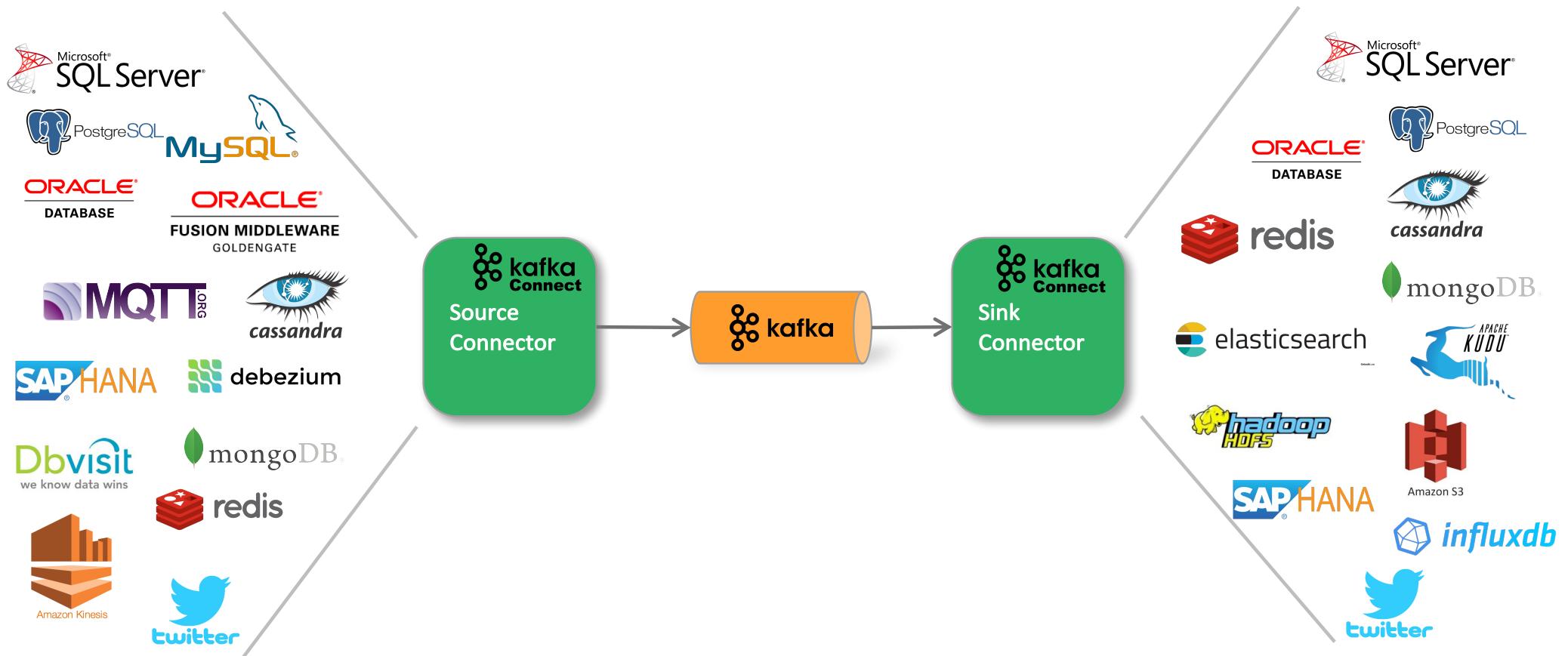
- Map dataflows to topologies, manage releases & track changes
- Measure KPIs and establish baselines for data availability and accuracy
- Master dataflow operations through Data SLAs



Source: <https://streamsets.com/connectors>

Kafka Connect

Kafka Connect - Overview



■ Kafka Connect – Single Message Transforms (SMT)

Simple Transformations for a single message

Defined as part of Kafka Connect

- some useful transforms provided out-of-the-box
- Easily implement your own

Optionally deploy 1+ transforms with each connector

- Modify messages produced by source connector
- Modify messages sent to sink connectors

Makes it much easier to mix and match connectors

Some of currently available transforms:

- InsertField
- ReplaceField
- MaskField
- ValueToKey
- ExtractField
- TimestampRouter
- RegexRouter
- SetSchemaMetaData
- Flatten
- TimestampConverter

Kafka Connect – Many Connectors

60+ since first release (0.9+)

20+ from Confluent and Partners

Confluent supported Connectors

CONNECTOR	TAGS	DEVELOPER/SUPPORT	DOWNLOAD
Amazon S3 (Sink)	S3	Confluent	Confluent
Elasticsearch (Sink)	search, Elastic, log, analytics	Confluent	Confluent
HDFS (Sink)	HDFS, Hadoop, Hive	Confluent	Confluent
JDBC (Source)	JDBC, MySQL	Confluent	Confluent
JDBC (Sink)	JDBC, MySQL	Confluent	Confluent

Certified Connectors

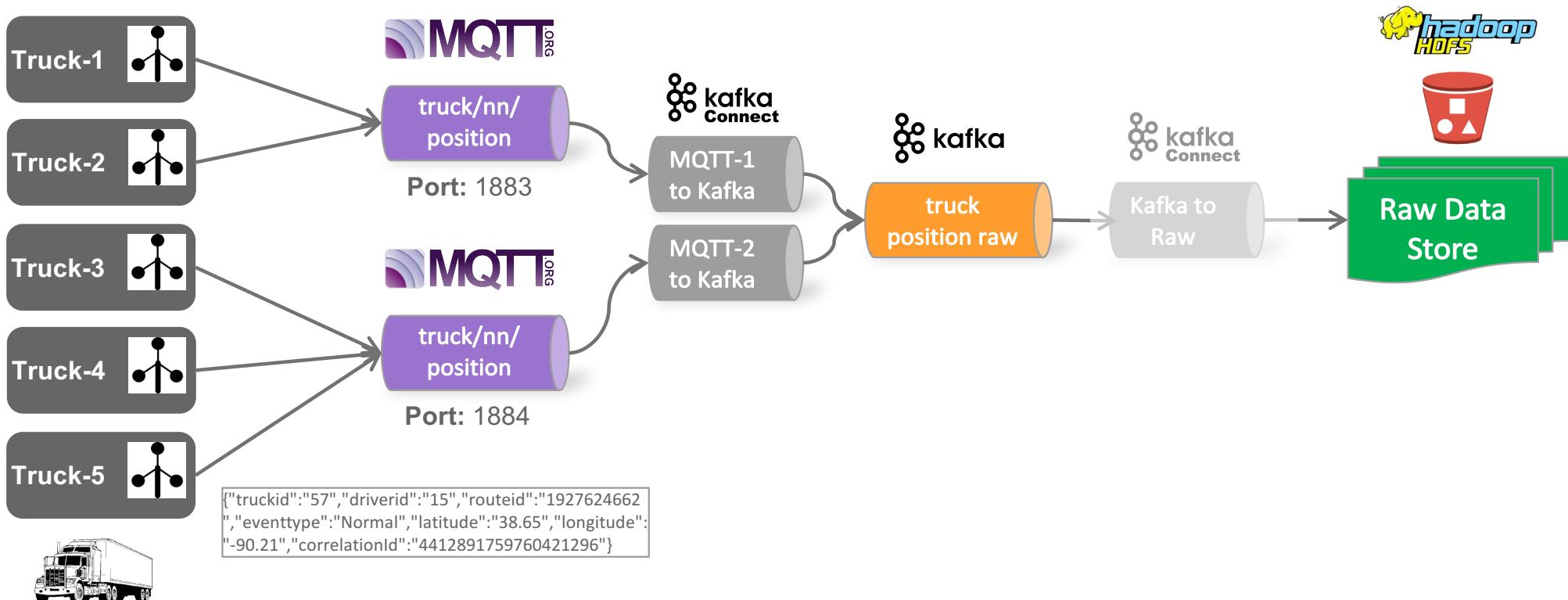
CONNECTOR	TAGS	DEVELOPER/SUPPORT	DOWNLOAD
Attunity (Source)	CDC	Attunity	Attunity
Azure IoTHub (Source)	IoT, messaging	Microsoft	Azure
Couchbase (Source)	Couchbase, NoSQL	Couchbase	Couchbase
Couchbase (Sink)	Couchbase, NoSQL	Couchbase	Couchbase
DataStax (Sink)	Cassandra, DataStax	Data Mountaineer	Data Mountaineer
Dbvisit Replicate (Source)	CDC, Oracle	Dbvisit	Dbvisit
IBM Data Replication (Source)	CDC	IBM	IBM
JustOne (Sink)	Postgres	JustOne	JustOne
Kinetica (Source)	GPU, RDBMS	Kinetica	Kinetica
Kinetica (Sink)	GPU, RDBMS	Kinetica	Kinetica
SAP HANA (Sink)	HANA, RDBMS	SAP	Community
SAP HANA (Source)	HANA, RDBMS	SAP	Community
Striim (Source)	CDC, MS SQLServer, Oracle, MySQL	Striim	Striim
Syncsort DMX (Source)	DB2, IMS, VSAM, CICS	Syncsort	Syncsort
Syncsort DMX (Sink)	DB2, IMS, VSAM, CICS	Syncsort	Syncsort
Vertica (Source)	Vertica	HP Enterprise	HP Enterprise
Vertica (Sink)	Vertica	HP Enterprise	HP Enterprise
VoltDB (Sink)	VoltDB, NewSQL	VoltDB	VoltDB

Community Connectors

CONNECTOR	TAGS	DEVELOPER/SUPPORT	DOWNLOAD
Apache Ignite (Source)	File System	Community	Community
Apache Ignite (Sink)	File System	Community	Community
Azure DocumentDB (Sink)	DocumentDB, Azure, NoSQL	Community	Community
Blockchain (Source)	Bitcoin, Blockchain	Community	Community
Bloomberg Ticker (Source)	Application feed	Community	Community
Cassandra (Source)	Cassandra	Community	Community 1
Cassandra (Sink)	Cassandra	Community	Community
CoAP (Source)	Constrained Application Protocol, CoAP	Community	Community
DynamoDB	Dynamo, NoSQL	Community	Community
Elasticsearch (Sink)	Elastic, search, log, analytics	Community	Community 1 Community 2 Community 3
Files/Directories (Source)	File System, Directories, Logs	Community	Community 1 Community 2
FileSystem Connector (Source)	File System, S3, HDFS	Community	Community
FTP (Source)	File System	Community	Community 1 Community 2
Github (Source)	Application Feed	Community	Community
Google BigQuery (Sink)	Analytics, Data Warehouse	Community	Community

Source: <http://www.confluent.io/product/connectors>

Demo Case



Demo: Dataflow for MQTT to Kafka

```
#!/bin/bash
curl -X "POST" "http://192.168.69.138:8083/connectors" \
  -H "Content-Type: application/json" \
  -d $'{
"name": "mqtt-source",
"config": {
  "connector.class": "com.datamountaineer.streamreactor.connect.mqtt.source.MqttSourceConnector",
  "connect.mqtt.connection.timeout": "1000",
  "tasks.max": "1",
  "connect.mqtt.kcql": "INSERT INTO truck_position SELECT * FROM truck/+position",
  "name": "MqttSourceConnector",
  "connect.mqtt.service.quality": "0",
  "connect.mqtt.client.id": "tm-mqtt-connect-01",
  "connect.mqtt.converter.throw.on.error": "true",
  "connect.mqtt.hosts": "tcp://mosquitto-1:1883"
}
}'
```

Summary

■ Summary

Apache NiFi

- visual dataflow modelling
- very powerful – “with power comes responsibility”
- special package for Edge computing
- data lineage and data provenance
- supports for backpressure
- no transport mechanism (DEV/TST/PROD)
- custom processors
- supported by Hortonworks

StreamSets

- visual dataflow modelling
- very powerful – “with power comes responsibility”
- special package for Edge computing
- data lineage and data provenance
- no transport mechanism
- custom sources, sinks, processors
- supported by StreamSets

Kafka Connect

- declarative style data flows
- simplicity - “simple things done simple”
- very well integrated with Kafka – comes with Kafka
- Single Message Transforms (SMT)
- use Kafka Streams for complex data flows
- custom connectors
- supported by Confluent

Technology on its own won't help you. You need to know how to use it properly.



trivadis
makes **IT** easier. 