



Powered by **Lucidworks**

Real-Time Cloud Native Open Source Streaming Of Any Data to Apache Solr

Timothy Spann, Developer Advocate



Speaker Bio

Developer Advocate

DZone Zone Leader and Big Data MVB;
@PaasDev

<https://github.com/tspannhw> <https://www.datainmotion.dev/>

<https://github.com/tspannhw/SpeakerProfile>

<https://dev.to/tspannhw>

<https://sessionize.com/tspann/>

<https://www.slideshare.net/bunkertor>



Agenda

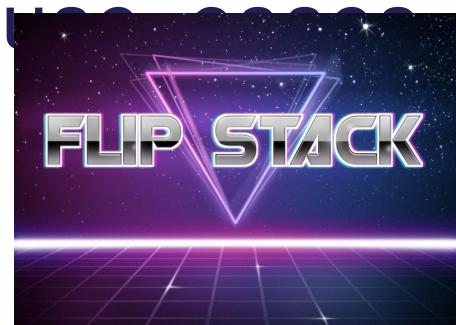
Utilizing Apache Pulsar and Apache NiFi we can parse any document in real-time at scale. We receive a lot of documents via cloud storage, email, social channels and internal document stores. We want to make all the content and metadata to Apache Solr for categorization, full text search, optimization and combination with other datastores. We will not only stream documents, but all REST feeds, logs and IoT data. Once data is produced to Pulsar topics it can instantly be ingested to Solr through Pulsar Solr Sink.

Utilizing a number of open source tools, we have created a real-time scalable any document parsing data flow. We use Apache Tika for Document Processing with real-time language detection, natural language processing with Apache OpenNLP, Sentiment Analysis with Stanford CoreNLP, Spacy and TextBlob. We will walk everyone through creating an open source flow of documents utilizing Apache NiFi as our integration engine. We can convert PDF, Excel and Word to HTML and/or text. We can also extract the text to apply sentiment analysis and NLP categorization to generate additional metadata about our documents. We also will extract and parse images that if they contain text we can extract with TensorFlow and Tesseract.

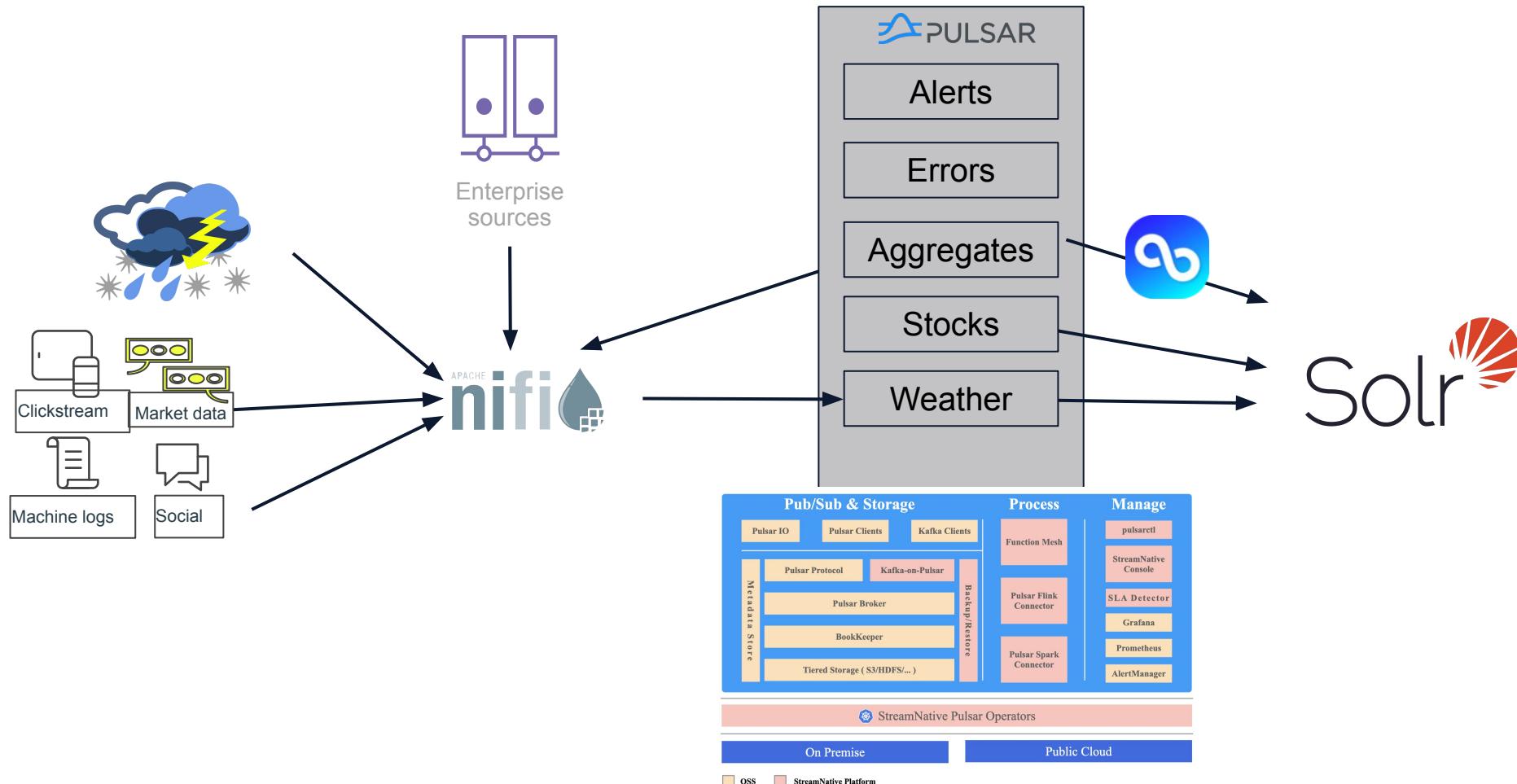
FLIP Stack

- Apache Flink
- Apache Pulsar
- StreamNative's Flink Connector for Pulsar
- Apache +++

Apache projects are the way for all streaming

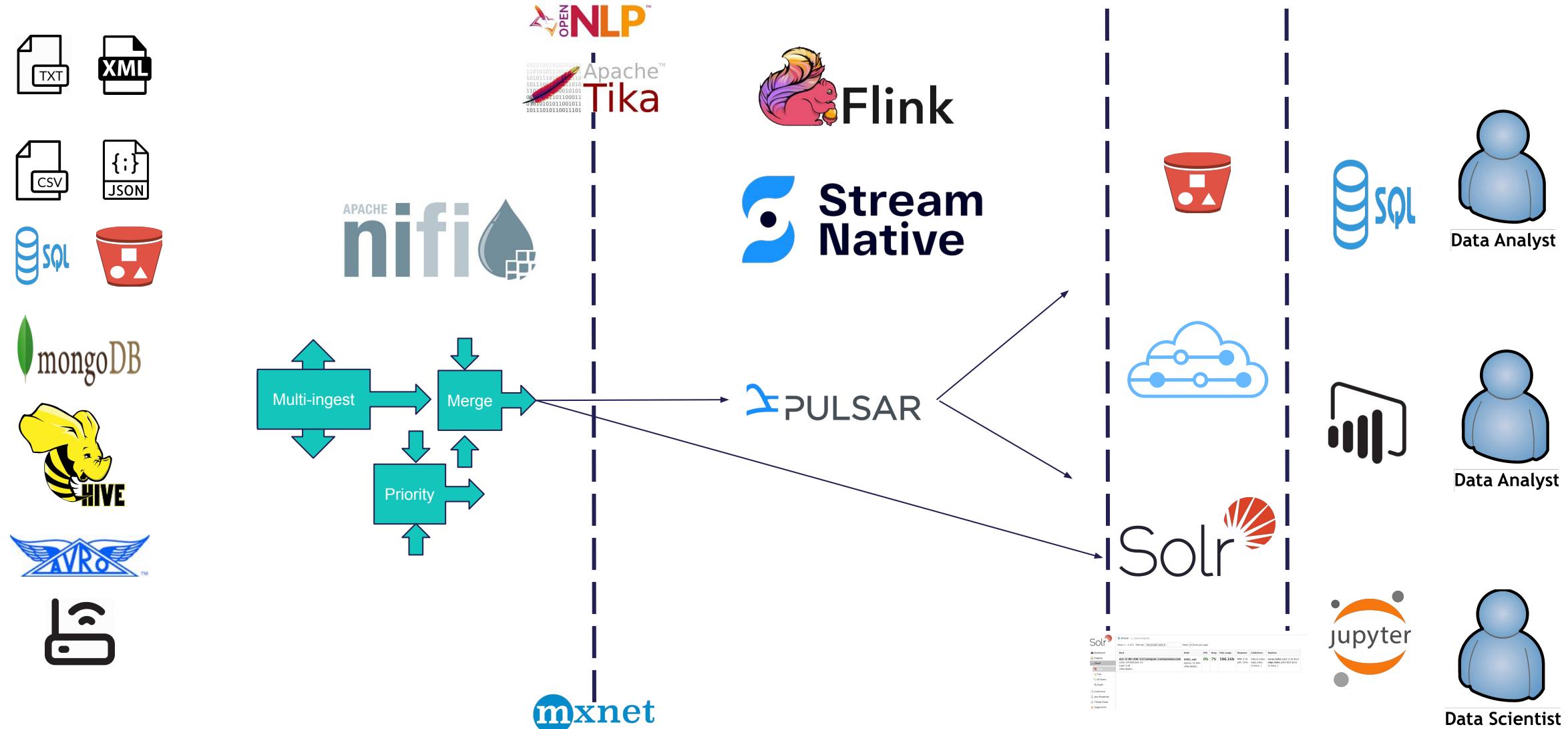


End to End Streaming Demo Pipeline



<https://hub.streamnative.io/connectors/solr-sink/2.5.1>

All Data - Anytime - Anywhere - Multi-Cloud - Multi-Protocol



StreamNative Cloud

Powered by Apache Pulsar, StreamNative provides a cloud-native, real-time messaging and streaming platform to support multi-cloud and hybrid cloud strategies.



Cloud Native



kubernetes

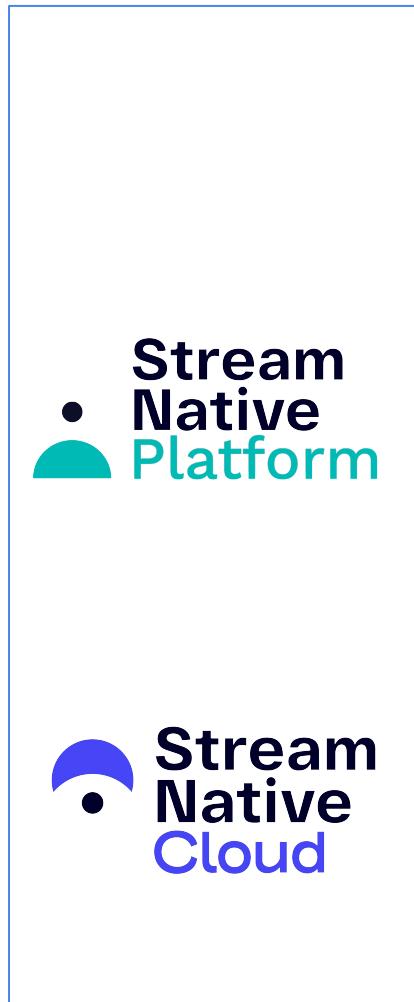
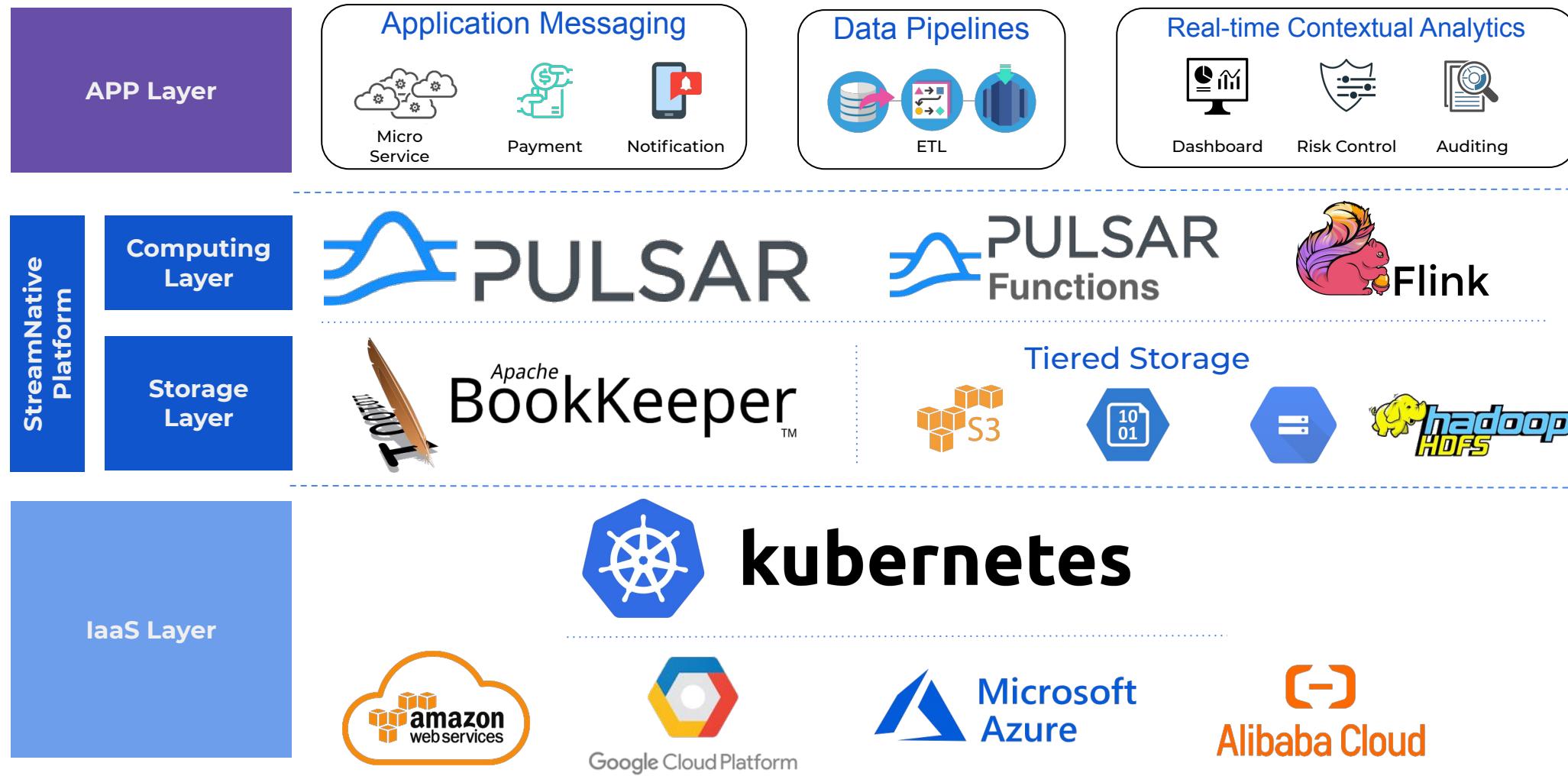
Built for Containers



Flink

Flink SQL

StreamNative Solution



Apache Pulsar



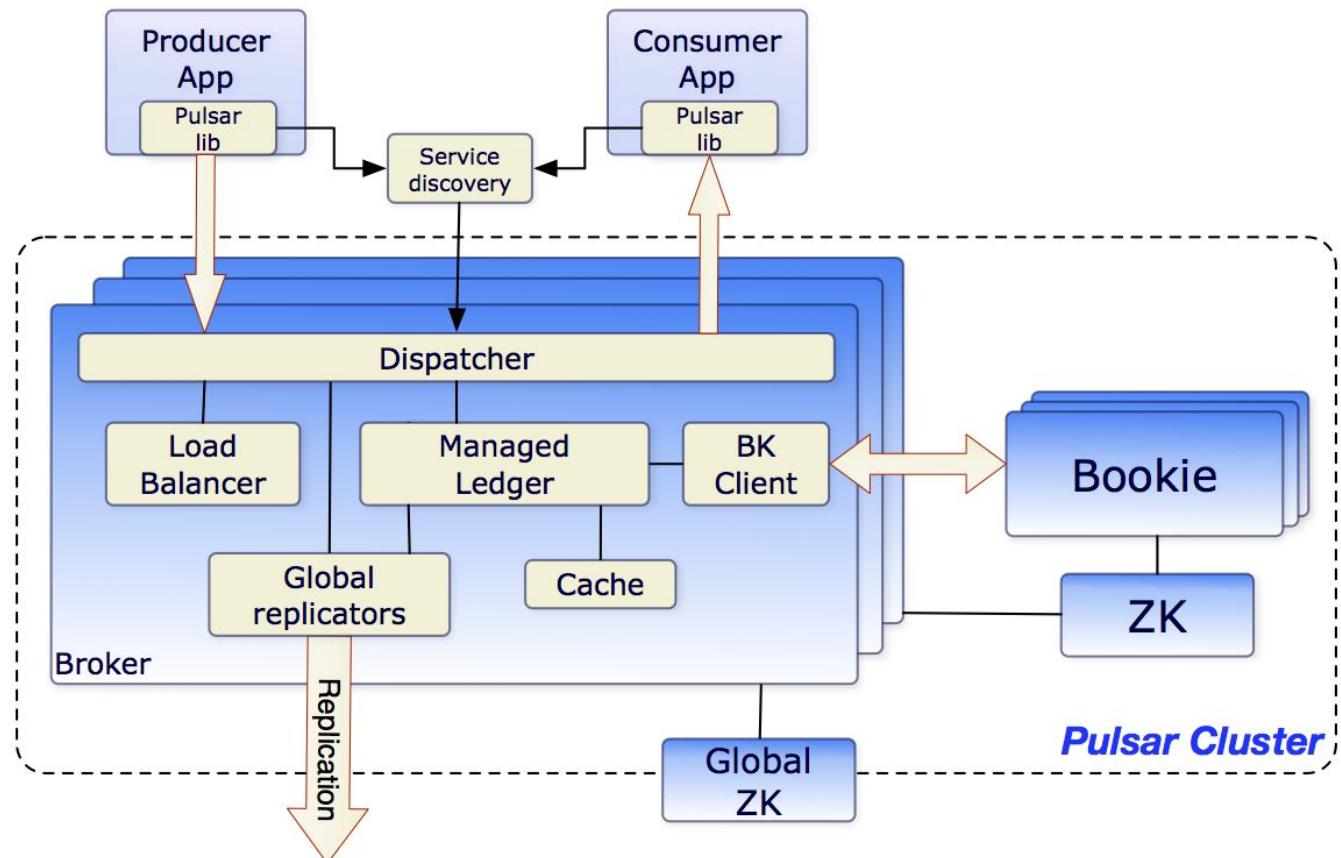
Apache Pulsar is Cloud-Native Messaging and Event-Streaming Platform



Apache Pulsar Overview

Enable Geo-Replicated Messaging

- Pub-Sub
- Geo-Replication
- Pulsar Functions
- Horizontal Scalability
- Multi-tenancy
- Tiered Persistent Storage
- Pulsar Connectors
- REST API
- CLI
- Many clients available
- Four Different Subscription Types
- Multi-Protocol Support
 - MQTT
 - AMQP
 - JMS
 - Kafka
 - ...



What are the Benefits of Pulsar?



Multi-Tenancy

Scalability

Geo-Replication

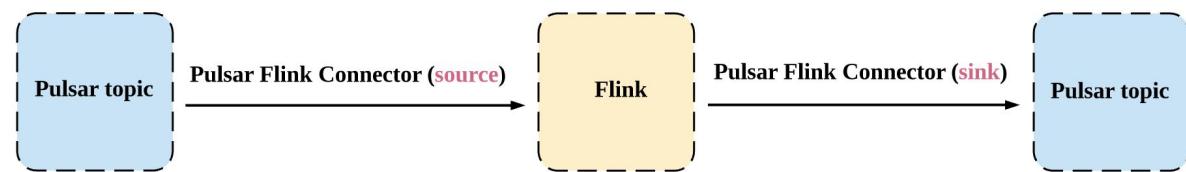
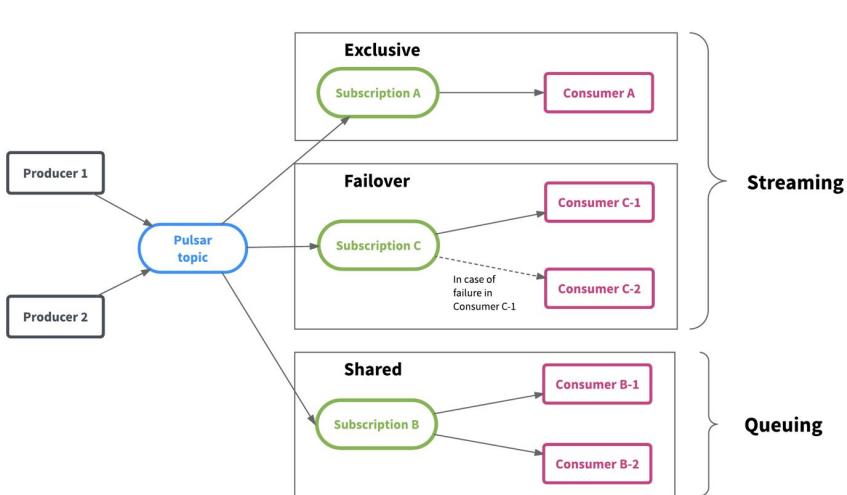
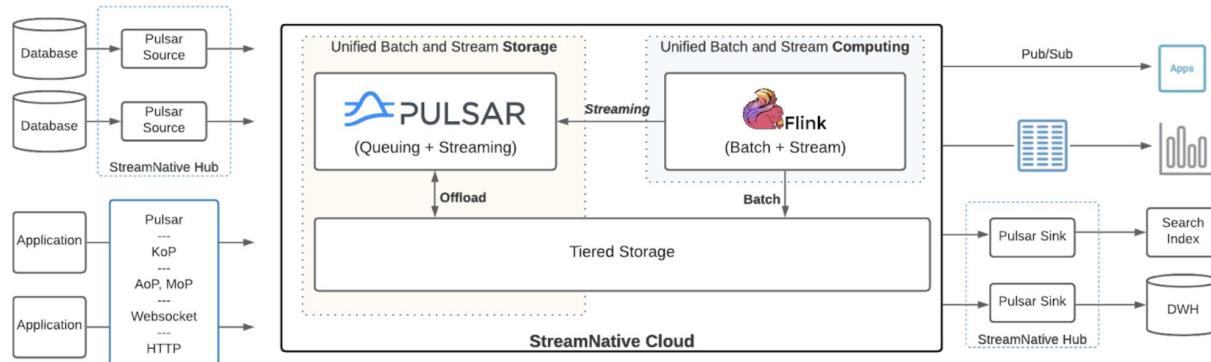
Unified Messaging
Model

Data Durability

A Unified Messaging Platform



Flink + Pulsar (FLiP)



<https://flink.apache.org/2019/05/03/pulsar-flink.html>

<https://github.com/streamnative/pulsar-flink>

<https://streamnative.io/en/blog/release/2021-04-20-flink-sql-on-streamnative-cloud>

Apache Solr



Apache Solr As a Destination

The screenshot shows the Apache Solr Admin interface with the 'Query' tab selected. The left sidebar contains links for Logging, Core Admin, Java Properties, Thread Dump, energy (selected), Overview, Analysis, Dataimport, Documents, Files, Ping, Plugins / Stats, Query (selected), Replication, Schema, and Segments info. The main area has fields for q.op (set to OR), fq (empty), sort (empty), start, rows (0, 10), fl (empty), df (empty), wt (set to -----), and checkboxes for indent off, debugQuery, defType (set to lucene), hl, facet, spatial, and spellcheck. Below these are Raw Query Parameters (key1=val1&key2=val2) and an Execute Query button.

```
params :{
  "q": "*:*",
  "q.op": "OR",
  "_": "1631799072231"}},
"response": {"numFound": 24, "start": 0, "numFoundExact": true, "docs": [
  {
    "current": [0.077649],
    "voltage": [117.896479],
    "power": [3.848307],
    "total": [1.251],
    "dev_name": ["Wi-Fi Smart Plug With Energy Monitoring"],
    "day1": [0],
    "systemtime": ["09/16/2021 09:29:50"],
    "cpu": [3.3],
    "memory": [29.8],
    "diskusage": [32607.4],
    "uuid": ["20210916132950_ef6a416c-8675-4332-9f34-401794f08fle"],
    "macaddress": ["lc:bf:ce:la:7f:a0"],
    "id": "0bdcd768-1868-406d-ae2f-eec27e1ba31f",
    "_version_": 1711065259866849280},
  {
    "current": [0.077437],
    "voltage": [117.892421],
    "power": [3.828128],
    "total": [1.252],
    "dev_name": ["Wi-Fi Smart Plug With Energy Monitoring"],
    "day1": [0],
    "systemtime": ["09/16/2021 09:30:51"],
    "cpu": [3.0],
    "memory": [29.9],
    "diskusage": [32607.3],
    "uuid": ["20210916133050_82338f49-7d7f-4c37-b4b0-4ee9042dbc0c"],
    "macaddress": ["lc:bf:ce:la:7f:a0"],
    "id": "b0d0fcfb1-e5ed-4ac0-b0c2-499e2b4dab59",
    "_version_": 1711065323473469440},
  {
    "current": [0.0770391]
```

Apache NiFi

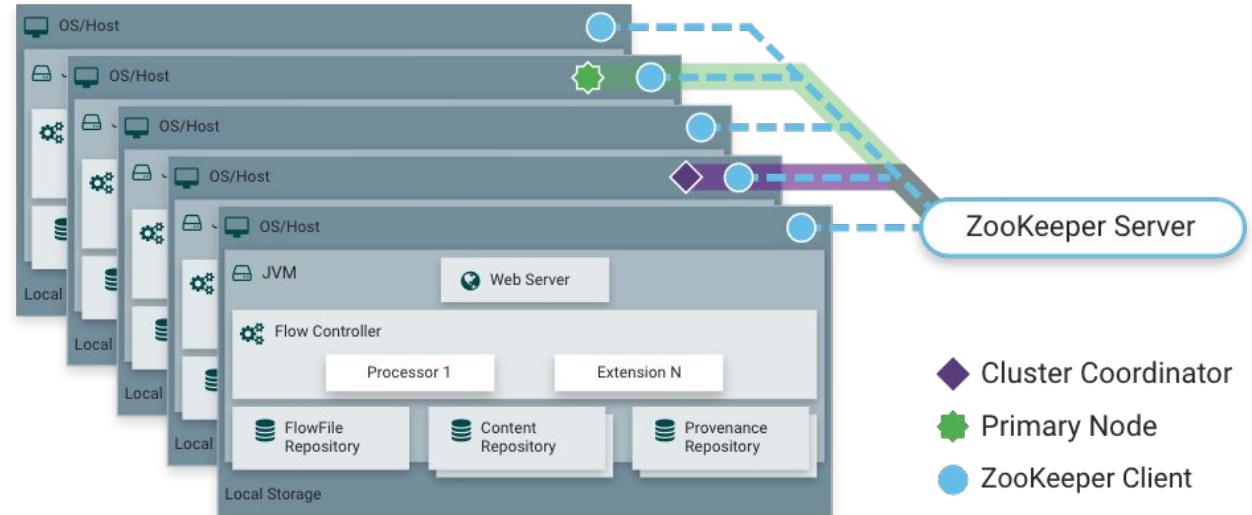
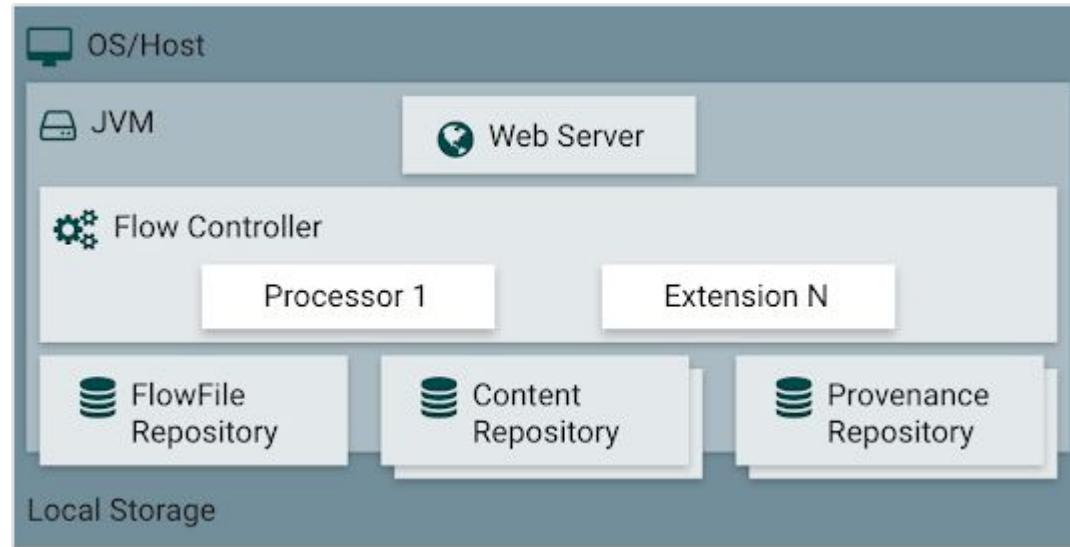


Why Apache NiFi?



- Guaranteed delivery
- Data buffering
 - Backpressure
 - Pressure release
- Prioritized queuing
- Flow specific QoS
 - Latency vs. throughput
 - Loss tolerance
- Data provenance
- Supports push and pull models
- Hundreds of processors
- Visual command and control
- Over a sixty sources
- Flow templates
- Pluggable/multi-role security
- Designed for extension
- Clustering
- Version Control

Architecture



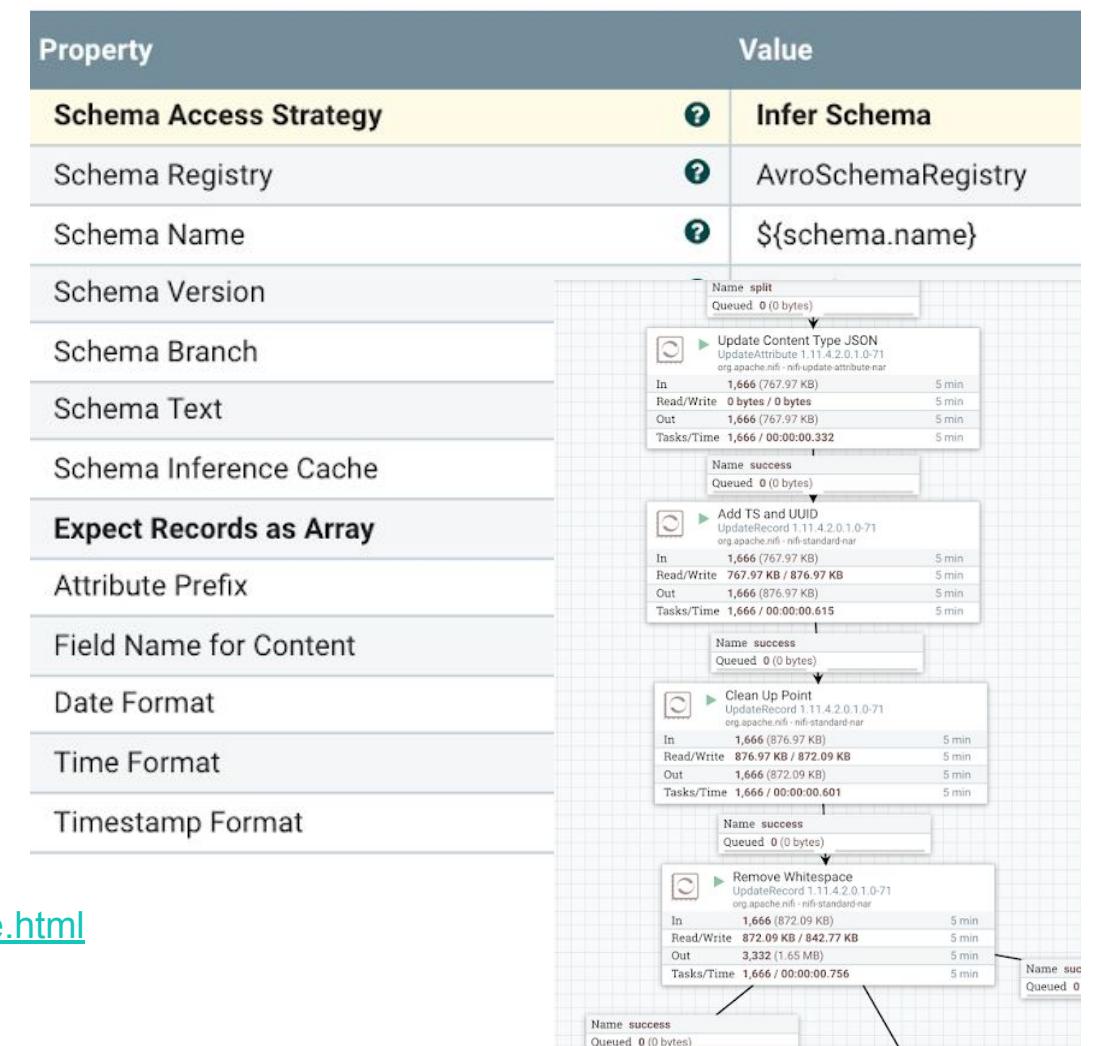
<https://nifi.apache.org/docs/nifi-docs/html/overview.html>

Record Processors

- XML, CSV, JSON, AVRO and more
- Schemas or Inferred Schemas
- Easily convert between them
- Support SQL with Apache Calcite

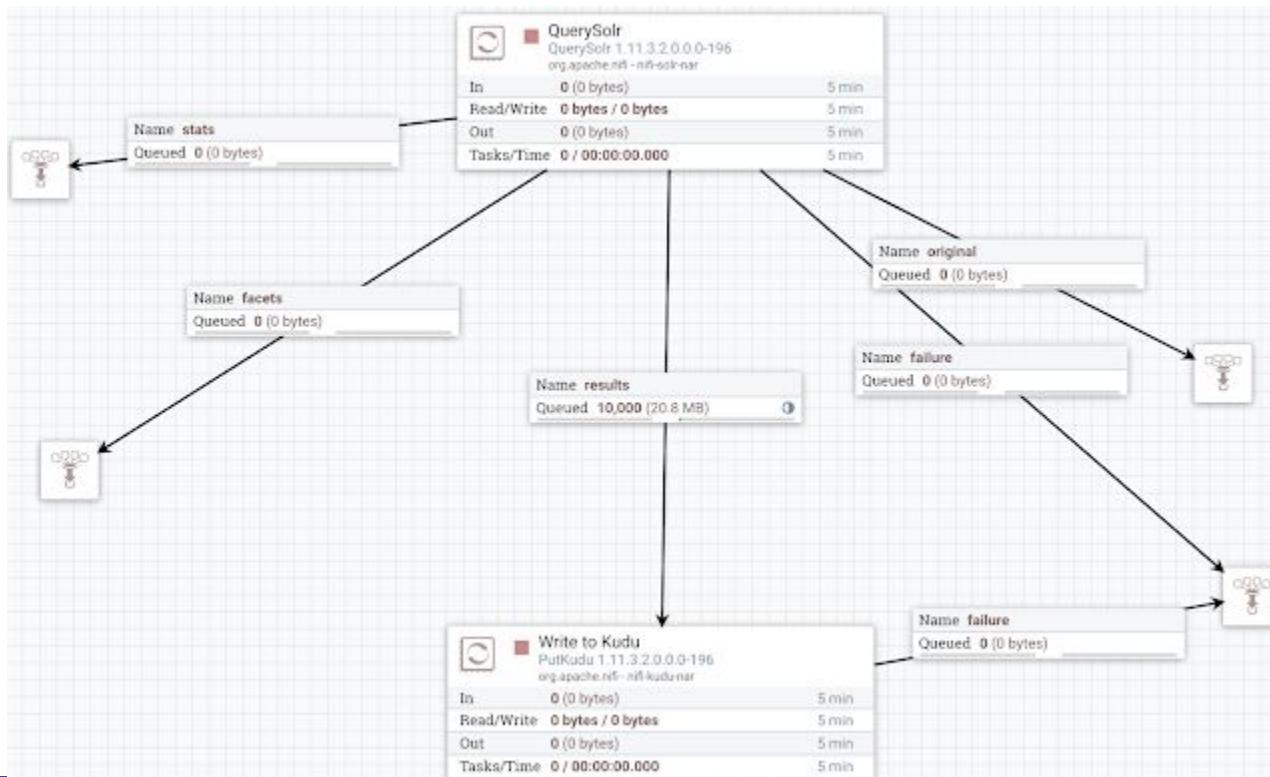
Property	Value
Record Reader	XMLReader
Record Writer	JsonRecordSetWriter
Include Zero Record FlowFiles	false
Cache Schema	true
query1	SELECT * FROM FLOWFILE

<https://www.datainmotion.dev/2019/03/advanced-xml-processing-with-apache.html>



SOLR Connectors

- XML, CSV, JSON, AVRO and more
- Schemas or Inferred Schemas
- Use Records or Raw Text
- Support SQL with Apache Calcite



Configure Processor

Stopped

SETTINGS SCHEDULING PROPERTIES COMMENTS

Required field

Property	Value
Solr Type	Cloud
Solr Location	.com:2181/solr
Collection	sensors
Solr Update Path	/update
Record Reader	JsonTreeReader
Fields To Index	sensor_id,sensor_ts
Commit Within	5000
Kerberos Credentials Service	No value set
Kerberos Principal	No value set
Kerberos Password	No value set
Basic Auth Username	No value set
Basic Auth Password	No value set

Apache OpenNLP with Apache NiFi

Apache OpenNLP for Entity Resolution Processor
<https://github.com/tspannhw/nifi-nlp-processor>

Requires installation of NAR and Apache OpenNLP Models
(<http://opennlp.sourceforge.net/models-1.5/>).

This is a non-supported processor that I wrote and put into the community. You can write one too!

FlowFile

DETAILS ATTRIBUTES

Attribute Values

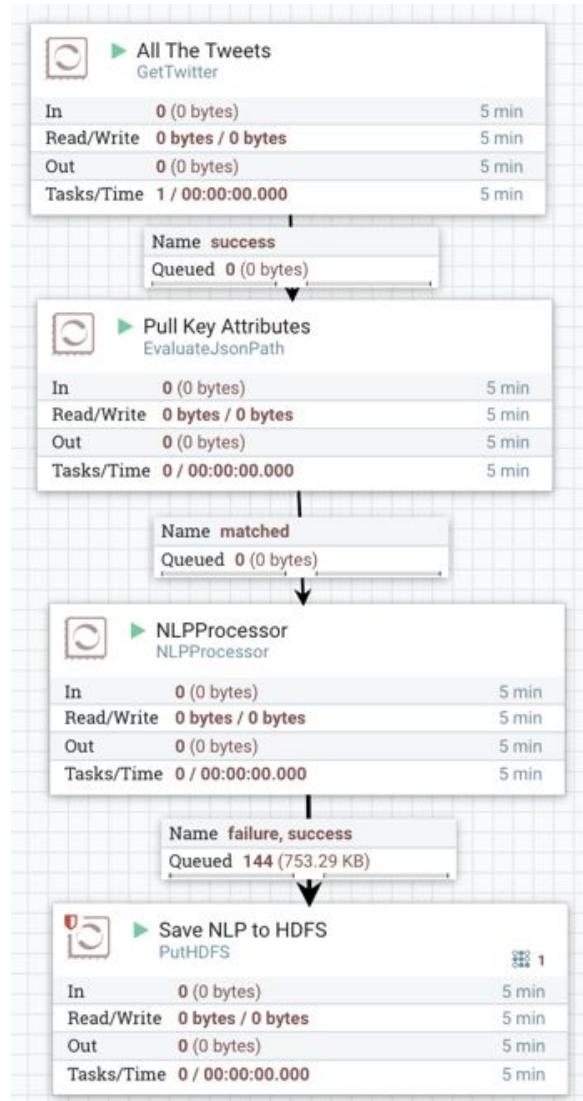
filename
2788601463132800.json

names
{"names": [{"name": "Tim Spann"}, {"name": "Peter Smith"}]}

followers_count
47

location
Columbus, Ohio

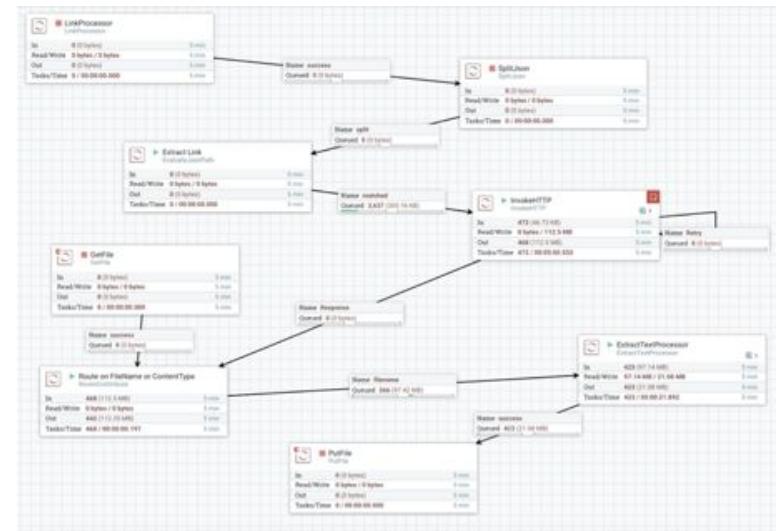
locations
{"locations": [{"location": "Sydney"}]}



Apache Tika with Apache NiFi

Displaying 1 of 195	
Type	Tags
ExtractTextProcessor	extracttextprocessor

Selected Processor:
ExtractTextProcessor
Run Tika Text Extraction from PDF, Word, Excel



View as: original ▾

streaming data. The combination of messaging and processing technologies enables stream processing at linear scale.

For example, Apache Storm ships with support for Kafka as a data source using Storm's core API or the higher-level, micro-batching Trident API. Storm's Kafka integration also includes support for writing data to Kafka, which enables complex data flows between components in a Hadoop-based architecture. For more information about Apache Storm, see the [Storm User Guide](#).

Legal notices

Contents Search

- 1. [Building a High-Throughput Messaging System with Apache Kafka](#)
- 2. [What's New](#)
- 3. [Apache Kafka Concepts](#)
- 4. [Installing Kafka 1. Prerequisites](#)
- 5. [Installing Kafka Using Ambari](#)
- 6. [Configuring Kafka for a Production Environment](#)
- 7. [Customizing Kafka Configuration Settings](#)
 - 1.1. Connection Settings
 - 1.2. Topic Settings
 - 1.3. Log Settings
 - 1.4. Compaction Settings
 - 1.5. Advanced kafka-env Settings
 - 1.6. Adding Configuration Properties
- 8. [Configuring ZooKeeper for Multiple Applications](#)
- 9. [Enabling Audit to HDFS for a Secure Cluster](#)
- 10. [Mirroring Data Between Clusters: Using the MirrorMaker Tool](#)
 - 1.1. Running MirrorMaker
 - 1.2. Checking Mirroring Progress
 - 1.3. Avoiding Data Loss
 - 1.4. Running MirrorMaker on Kerberos-Enabled Clusters
- 11. [Developing Kafka Producers and Consumers](#)

Search

<https://github.com/tspannhw/nifi-extracttext-processor>

<https://community.hortonworks.com/articles/76924/data-processing-pipeline-parsing-pdfs-and-identify.html>

<https://community.hortonworks.com/articles/81694/extracttext-nifi-custom-processor-powered-by-apach.html>

<https://community.hortonworks.com/articles/163776/parsing-any-document-with-apache-nifi-15-with-apac.html>

<https://community.hortonworks.com/content/kbentry/177370/extracting-html-from-pdf-excel-and-word-documents.html>

Final Thoughts



Build Your Own Pulsar - SOLR Integration

<https://github.com/tspannhw/FLiP-Energy>



```
bin/pulsar-admin sinks create --tenant public  
--namespace default  
--name solr-sink-energy  
--sink-type solr  
--sink-config-file conf/solr-sink-energy.yml  
--inputs energy
```

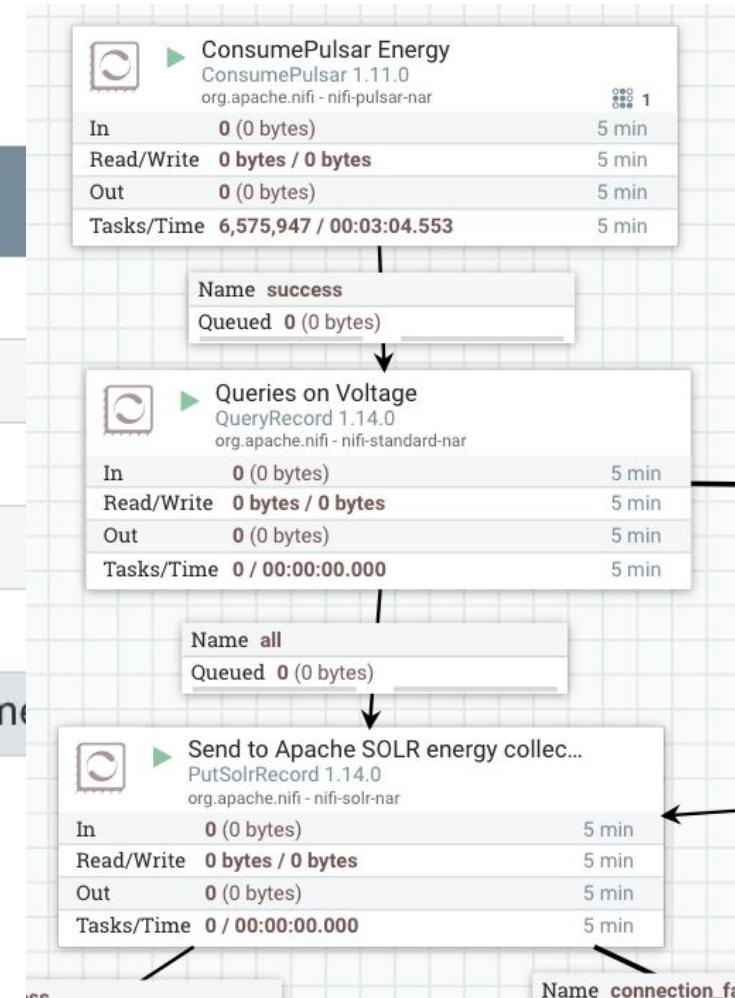
```
tspann@Timothys-mbp apache-pulsar-2.8.0 % cat conf/solr-sink.yml  
configs:  
  solrUrl: "http://localhost:8983/solr"  
  solrMode: "Standalone"  
  solrCollection: "iot"
```



Build Your Own Pulsar - NiFi Integration

PutSolrRecord

Property	Value
Solr Type	Standard
Solr Location	http://localhost:8983/solr/energy
Collection	energy
Solr Update Path	/update
Record Reader	CDP Infer JsonTreeReader
Fields To Index	current,voltage,power,total,diskusage,cpu,systemtime,uuid,message
Commit Within	5000



Connect with the Community & Stay Up-To-Date

- Join the Pulsar Slack channel - Apache-Pulsar.slack.com
- Follow [@streamnativeio](https://twitter.com/streamnativeio) and [@apache_pulsar](https://twitter.com/apache_pulsar) on Twitter
- [Subscribe](#) to Monthly Pulsar Newsletter for major news, events, project updates, and resources in the Pulsar community

Deeper Content

- <https://www.datainmotion.dev/2020/04/building-search-indexes-with-apache.html>
- <https://github.com/tspannhw/nifi-solr-example>
- <https://github.com/streamnative/pulsar-flink>
- <https://www.linkedin.com/pulse/2021-schedule-tim-spann/>
- https://github.com/tspannhw/SpeakerProfile/blob/main/2021/talks/20210729_HailHydrate!FromStreamtoLake_TimSpann.pdf
- <https://streamnative.io/en/blog/release/2021-04-20-flink-sql-on-streamnative-cloud>
- <https://docs.streamnative.io/cloud/stable/compute/flink-sql>
- <https://pulsar.apache.org/docs/en/client-libraries-websocket/>



@PaasDev



timothyspann

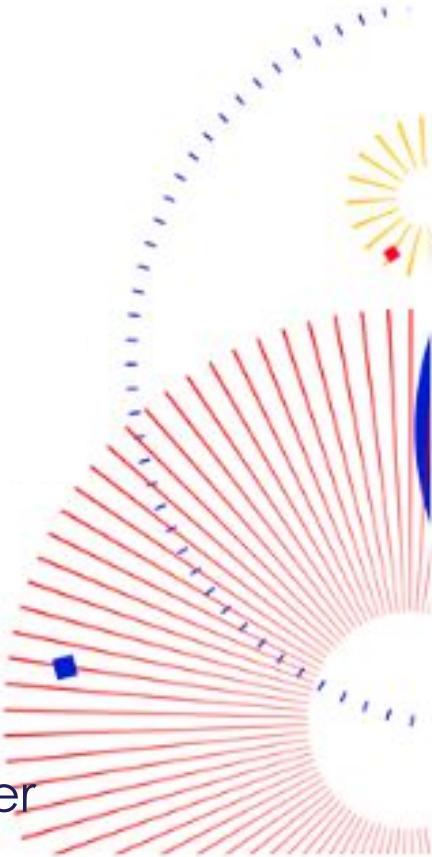
<https://www.pulsardeveloper.com/>



Pulsar Summit Asia

November 20-21, 2021

Contact us at partners@pulsar-summit.org to become a sponsor or partner



Announcing

Flink SQL on StreamNative Cloud



Thank You

ACTIVATE
VIRTUAL SEARCH & AI CONFERENCE
Powered by **Lucidworks**