



# Apache NiFi 101

**Timothy Spann** | Developer Advocate



**Tim Spann**

Developer Advocate

## DZone Zone Leader and Big Data MVB Data DJay

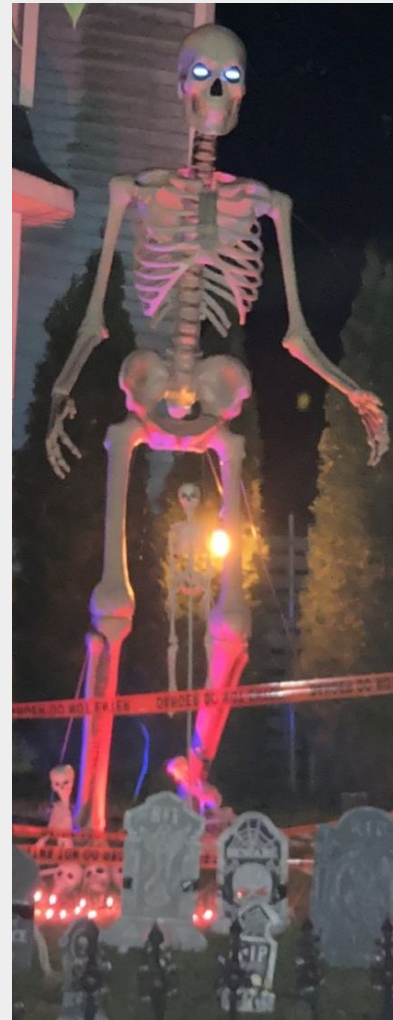
- <https://www.datainmotion.dev/>
- <https://github.com/tspannhw/SpeakerProfile>
- <https://dev.to/tspannhw>
- <https://sessionize.com/tspann/>
- <https://streamnative.io/>



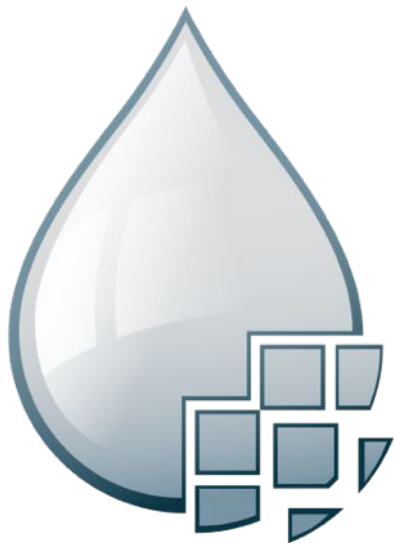
# Apache NiFi



# Don't Be Afraid of Open Source

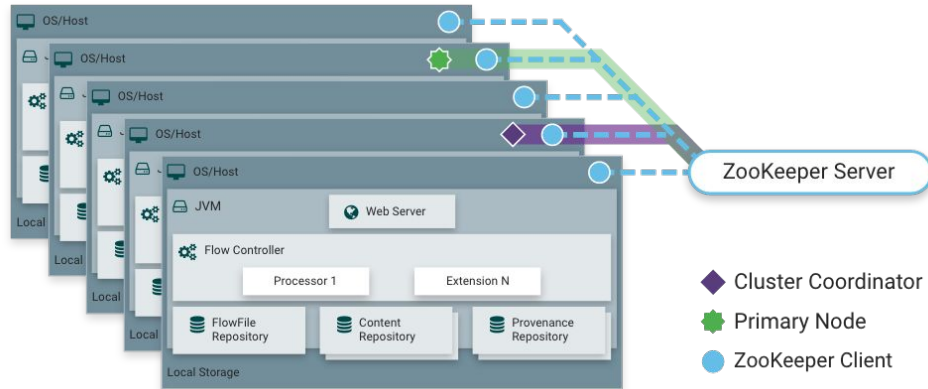
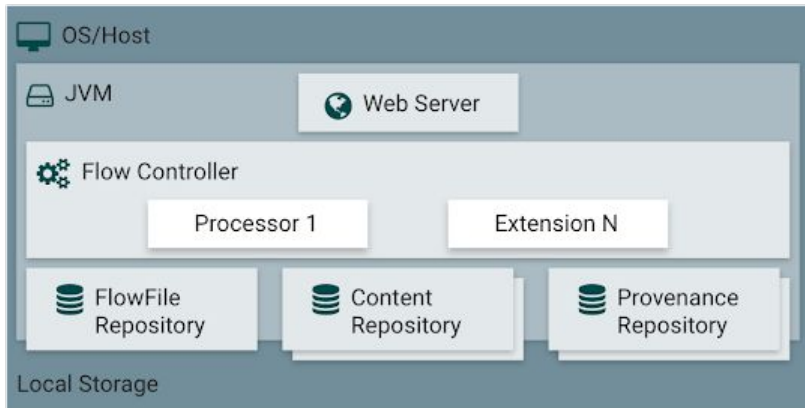


# Why Apache NiFi?



- 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>

# Provenance

## NiFi Data Provenance

Displaying 165 of 165

Oldest event available: 12/21/2020 16:55:33 UTC

Filter	by component name		
Date/Time	Type	FlowFile UUID	Size
12/22/2020 16:54:17 193 UTC	ATTRIBUTES_MODIFIED	6fdbae4f-6ba4-47c3-ba89-8830ec7cd3db	89 byt
12/22/2020 16:54:17 192 UTC	ATTRIBUTES_MODIFIED	1213e84d-e84d-4219-b3bd-c258be7f901e	81 byt
12/22/2020 16:54:14 194 UTC	ATTRIBUTES_MODIFIED	378c2d53-4195-48b0-a433-7aba74ed3718	81 byt
12/22/2020 16:54:03 297 UTC	ATTRIBUTES_MODIFIED	695dfce-9d71-4cd6-8f73-b5b46ed3622	83 byt
12/22/2020 16:53:59 296 UTC	ATTRIBUTES_MODIFIED	df43c05c-5aae-44c2-9eda-c620a8148604	84 byt
12/22/2020 16:53:59 295 UTC	ATTRIBUTES_MODIFIED	4b1dbb1e-1f83-4af9-b309-2da1277cd2e6	84 byt
12/22/2020 16:53:58 296 UTC	ATTRIBUTES_MODIFIED	45f8edd-ca55-431e-82b9-436ca4d092e	81 byt
12/22/2020 16:53:57 298 UTC	ATTRIBUTES_MODIFIED	bb07034b-63e1-d334-b	
12/22/2020 16:53:57 297 UTC	ATTRIBUTES_MODIFIED	d1a26b1e-a7af-4e16-b	
12/22/2020 16:53:57 297 UTC	ATTRIBUTES_MODIFIED	29966a0-4153-41bc-a	
12/22/2020 16:53:43 753 UTC	ATTRIBUTES_MODIFIED	1ca5c744-1cb4-4ff1-8b	
12/22/2020 16:53:37 747 UTC	ATTRIBUTES_MODIFIED	fafa47db-96db-48cc-a	
12/22/2020 16:53:21 646 UTC	ATTRIBUTES_MODIFIED	def1609-9446-460e-99	
12/22/2020 16:53:05 515 UTC	ATTRIBUTES_MODIFIED	964a95fc-df58-460c-a2	
12/22/2020 16:52:43 374 UTC	ATTRIBUTES_MODIFIED	79cfa9b-b16b-4dc4-8a	
12/22/2020 16:52:29 308 UTC	ATTRIBUTES_MODIFIED	3433eeb3-953c-4952-8	
12/22/2020 16:52:29 307 UTC	ATTRIBUTES_MODIFIED	a166e9a7-118a-42d2-9	
12/22/2020 16:52:29 307 UTC	ATTRIBUTES_MODIFIED	bd2946f6-5a89-40a7-b	
12/22/2020 16:52:29 307 UTC	ATTRIBUTES_MODIFIED	a16841bc-2505-4c8c-b	
12/22/2020 16:52:29 306 UTC	ATTRIBUTES_MODIFIED	578540f6-c446-47f1-85	
12/22/2020 16:52:29 306 UTC	ATTRIBUTES_MODIFIED	3d44c5f8-a737-4a8e-8c	
12/22/2020 16:52:29 306 UTC	ATTRIBUTES_MODIFIED	4dc93617-7059-424e-9c	
12/22/2020 16:52:29 306 UTC	ATTRIBUTES_MODIFIED	9f1b1cc1-f384-c111-93	
12/22/2020 16:52:29 306 UTC	ATTRIBUTES_MODIFIED	75a2b193-a45c-47b8-b	

### Provenance Event

DETAILS

ATTRIBUTES

#### Attribute Values

lastprice

123.66

No value set

symbol

IBM

No value set

timestamp

1608654962884

No value set

volume

100

No value set

<https://www.datainmotion.dev/2021/01/automating-starting-services-in-apache.html>

# Backpressure & Prioritizers

## Configure Connection

DETAILS

SETTINGS

Name

Id

3ca22430-cba4-3347-b45b-7bdc3530bd7e

FlowFile Expiration ?

0 sec

Back Pressure Object Threshold ?

10000

Size Threshold ?

1 GB

Load Balance Strategy ?

Do not load balance ▼

Available Prioritizers ?

FirstInFirstOutPrioritizer

NewestFlowFileFirstPrioritizer

OldestFlowFileFirstPrioritizer

PriorityAttributePrioritizer

Selected Prioritizers ?

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



# System Diagnostics

NiFi Summary

PROCESSORS

INPUT PORTS

OUTPUT PORTS

REMOTE PROCESS GROUPS

CONNECTIONS

PROCESS GROUPS

Displaying 1,338 of 1,338

Filter

by name

Name	Type	Process Group	Run Status	In (Size) 5 min	Read   Write 5 min	Out (Size) 5 min	Tasks   Time 5 min	
PublishPulsar	PublishPulsar	Status Pulsar	Running (1)	61 (50.13 KB)	50.13 KB   0 bytes	0 (0 bytes)	5,504,778   00:00:23.264	
Acquire Satellite Data	GenerateFlowFile	Satellite Data	Disabled	0 (0 bytes)	0 bytes   0 bytes	0 (0 bytes)	0   00:00:00.000	
Acquire Satellite Data	GenerateFlowFile	Satellite Data	Disabled	0 (0 bytes)	0 bytes   0 bytes	0 (0 bytes)	0   00:00:00.000	
Analyze Data in Stream	QueryRecord	Fresh Food	Disabled	0 (0 bytes)	0 bytes   0 bytes	0 (0 bytes)	0   00:00:00.000	
App Data	PublishKafkaRecord_2.0	Mobile Ingest	Disabled	0 (0 bytes)	0 bytes   0 bytes	0 (0 bytes)	0   00:00:00.000	
App Data	PublishKafkaRecord_2.0	Mobile Ingest	Disabled	0 (0 bytes)	0 bytes   0 bytes	0 (0 bytes)	0   00:00:00.000	
AttributeCleanerProcessor	AttributeCleanerProcessor	Mobile Ingest	Disabled	0 (0 bytes)	0 bytes   0 bytes	0 (0 bytes)	0   00:00:00.000	
AttributeCleanerProcessor	AttributeCleanerProcessor	Mobile Ingest	Disabled	0 (0 bytes)	0 bytes   0 bytes	0 (0 bytes)	0   00:00:00.000	
Attributes Grab	EvaluateJsonPath	Predict Temperature	Disabled	0 (0 bytes)	0 bytes   0 bytes	0 (0 bytes)	0   00:00:00.000	

# Flow File

Flow Files are content and key/value pairs for attributes that are each event/message/file that has been introduced into NiFi.



## FlowFile

DETAILS

ATTRIBUTES

### FlowFile Details

UUID  
29235844-0576-432b-9078-b477f9ff6f5

Filename  
35714351-11c8-48e0-b52c-6e295edbbb69

File Size  
1.15 KB

Queue Position  
No value set

Queued Duration  
1 days and 03:58:04.364

Lineage Duration  
1 days and 04:01:15.707

Penalized  
No

### Content Claim

Container  
default

Section  
2

Identifier  
1631110638665-2

Offset  
50840

Size  
1.15 KB

DOWNLOAD

VIEW

## FlowFile

DETAILS

ATTRIBUTES

### Attribute Values

invokehttp.tx.id  
4e97d934-a7f5-463b-ad4a-d74ae6da4325

link  
<https://status.cloudera.com/incidents/967wlx3x5rv7>

mime-type  
application/json

mime.type  
application/json

path  
./

pubdate  
Thu, 09 Jul 2020 00:14:58 +0000

record.count



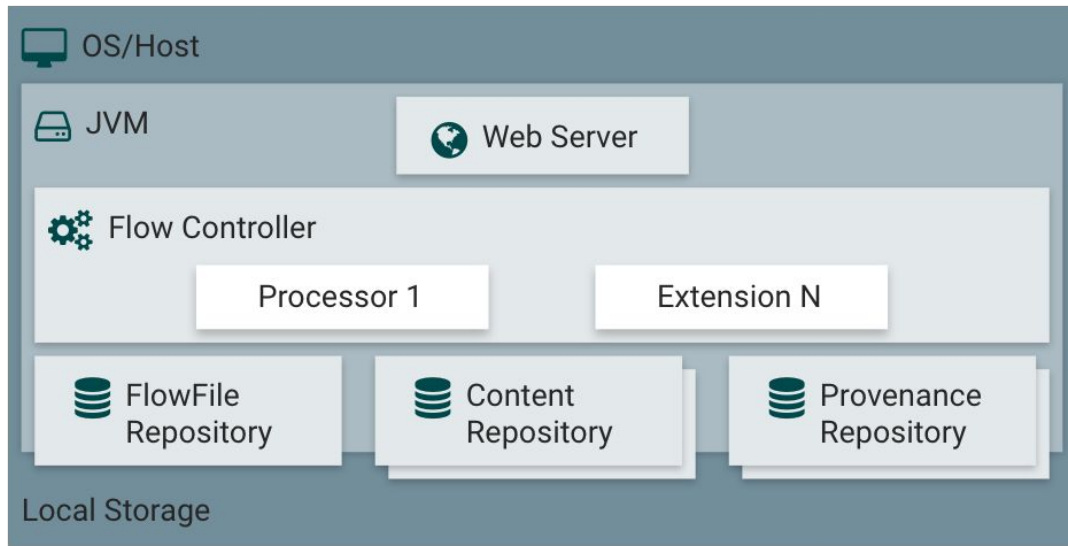
Streamlit <https://nifi.apache.org/docs/nifi-docs/html/overview.html>

# Version Control (Github and Beyond)

The screenshot shows the 'Tweets' group in the StreamNative interface. The group has a status bar with icons for target, no target, play, and a green checkmark. Below the status bar, there are four rows: 'Queued', 'In', 'Read/Write', and 'Out'. The 'Queued' row shows '0' for both 'In' and 'Out'. The 'In' row shows '0' for 'Read/Write' and 'Out'. The 'Read/Write' row shows '0' for 'In' and 'Out'. The 'Out' row shows '0' for 'Read/Write' and 'Out'. A context menu is open over the 'Queued' row, showing options: 'Configure', 'Parameters', 'Variables', 'Version', 'Enter group', 'Start', 'Stop', 'Enable', 'Disable', 'View status history', 'View connections', 'Manage access policies', 'Center in view', 'Group', 'Download flow', 'Create template', 'Copy', and 'Delete'. The 'Version' option is highlighted, and a sub-menu is open showing 'Commit local changes', 'Show local changes', 'Revert local changes', and 'Stop version control'.

The screenshot shows the 'Tweets' group in the StreamNative interface. The group has a status bar with icons for target, no target, play, a red square, a yellow triangle, and a green checkmark. Below the status bar, there are four rows: 'Queued', 'In', 'Read/Write', and 'Out'. The 'Queued' row shows '481 (66.59 MB)' for 'In' and 'Out'. The 'In' row shows '0 (0 bytes) → 0' for 'Read/Write' and 'Out'. The 'Read/Write' row shows '0 bytes / 0 bytes' for 'In' and 'Out'. The 'Out' row shows '0 → 0 (0 bytes)' for 'Read/Write' and 'Out'. A context menu is open over the 'Queued' row, showing options: 'Configure', 'Parameters', 'Variables', 'Version', 'Enter group', 'Start', 'Stop', 'Enable', 'Disable', 'View status history', 'View connections', 'Manage access policies', 'Center in view', 'Group', 'Download flow', 'Create template', 'Copy', and 'Delete'. The 'Download flow' option is highlighted.

# Repositories



<https://nifi.apache.org/docs/nifi-docs/html/nifi-in-depth.html#repositories>

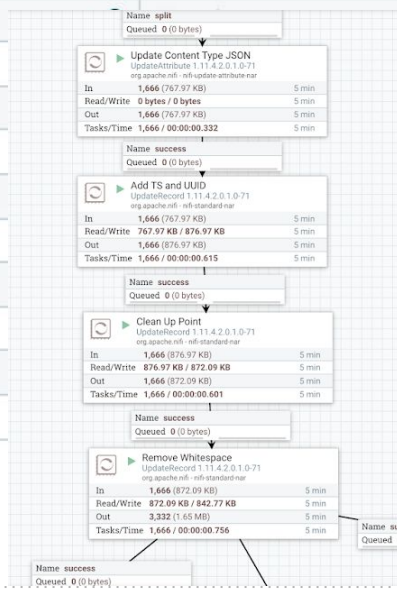
# Record Processors



- 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

Property		Value
Schema Access Strategy	?	Infer Schema
Schema Registry	?	AvroSchemaRegistry
Schema Name	?	\${schema.name}
Schema Version		
Schema Branch		
Schema Text		
Schema Inference Cache		
Expect Records as Array		
Attribute Prefix		
Field Name for Content		
Date Format		
Time Format		
Timestamp Format		



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

# Record Processors



Configure Processor

Invalid

SETTINGS SC

Required field

Property

Record Reader

Record Writer

Include Zero Record

**Add Controller Service**

Requires Controller Service  
RecordReaderFactory 1.13.0 from org.apache.nifi - nifi-standard-services-api-nar

Compatible Controller Services

AvroReader 1.13.0	✓
CSVReader 1.13.0	?
GrokReader 1.13.0	?
JsonPathReader 1.13.0	?
JsonTreeReader 1.13.0	?
ParquetReader 1.13.0	?
ReaderLookup 1.13.0	?
ScriptedReader 1.13.0	?
Syslog5424Reader 1.13.0	?
SyslogReader 1.13.0	?
WindowsEventLogReader 1.13.0	?
XMLReader 1.13.0	?

Required field

Property

Record Reader

Record Destination S

Include Zero Record

RecordSinkService 1.13.0 from org.apache.nifi - nifi-standard-services-api-nar

Compatible Controller Services

RecordSinkServiceLookup 1.13.0

Controller Service Name

RecordSinkServiceLookup

Bundle

org.apache.nifi - nifi-record-sink-service-nar

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

# Caching

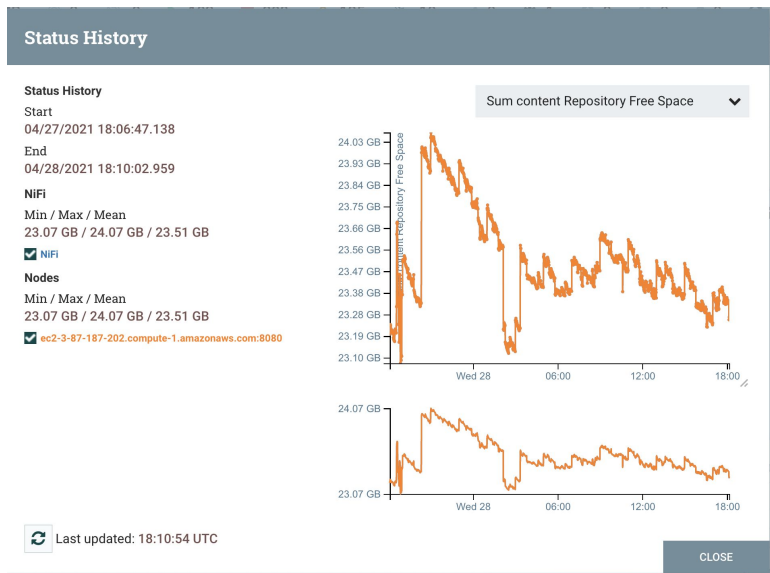


Property		Value
Record Reader	?	Infer JsonTreeReader
Record Writer	?	Standard Inherit JsonRecordSetWriter
Lookup Service	?	KuduLookupService
Result RecordPath	?	No value set
Routing Strategy	?	Route to 'matched' or 'unmatched'
Record Result Contents	?	Insert Entire Record
Record Update Strategy	?	Use Property
setid	?	/setid
version	?	/version

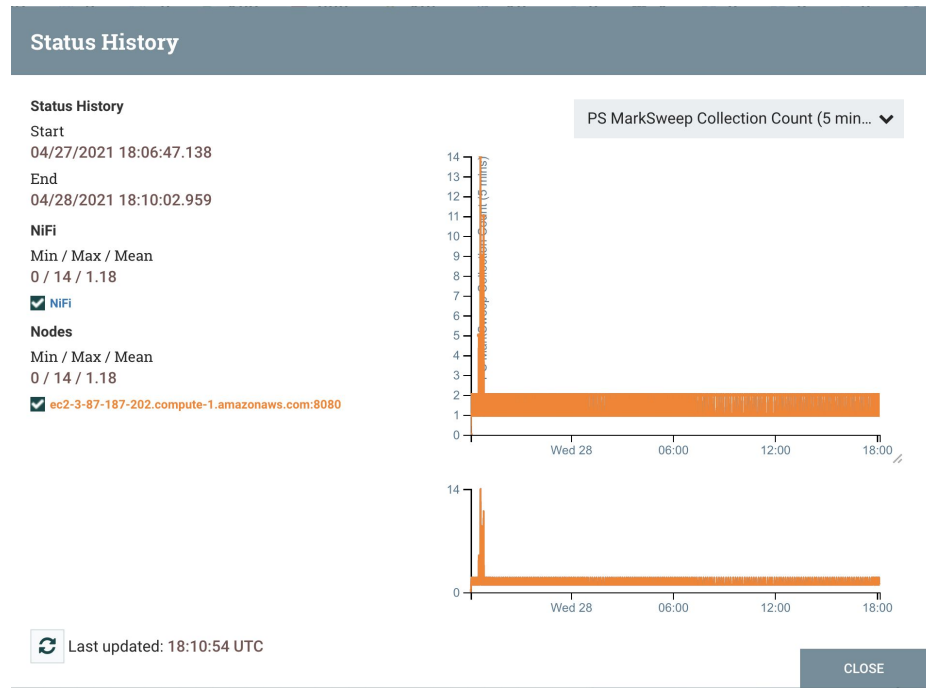
<https://dev.to/tspannhw/flank-using-apache-kudu-as-a-cache-for-fda-updates-4knj>

---

# Metrics, Status, Charts



<https://www.clouddataops.dev/data-flow-experience>





# DevOps



```
nifi-toolkit/bin/cli.sh nifi list-param-contexts -u http://edge2ai-1.dim.local:8080  
nifi-toolkit/bin/cli.sh nifi pg-list -u http://edge2ai-1.dim.local:8080  
nifi-toolkit/bin/cli.sh nifi pg-set-param-context -u http://edge2ai-1.dim.local:8080 ...
```

<https://www.datainmotion.dev/2021/01/automating-starting-services-in-apache.html>

<https://nipyapi.readthedocs.io/en/latest/>

---

# DevOps



```
nifi pg-list  
nifi pg-status  
nifi pg-get-services  
nifi pg-enable-services -u http://edge2ai-1.dim.local:8080 --processGroupId root  
nifi pg-start -u http://edge2ai-1.dim.local:8080 -pgid LOOKTHISUP  
nifi list-param-contexts -u http://edge2ai-1.dim.local:8080 -verbose  
nifi create-reporting-task -u http://edge2ai-1.dim.local:8080 -verbose -i
```

<https://dev.to/tspannhw/automating-starting-services-in-apache-nifi-and-applying-parameters-5h4n>

<https://github.com/tspannhw/ApacheConAtHome2020/blob/main/scripts/setupnifi.sh>

---

# Consume MQTT





This could read from Apache Pulsar - MoP (MQTT on Pulsar)

Property	ConsumeMQTT Processor	Value
<b>Session state</b>	?	<b>Clean Session</b>
<b>MQTT Specification Version</b>	?	<b>AUTO</b>
Connection Timeout (seconds)	?	30
Keep Alive Interval (seconds)	?	60
Group ID	?	No value set
<b>Topic Filter</b>	?	<b>No value set</b>
<b>Quality of Service(QoS)</b>	?	<b>0 - At most once</b>
<b>Max Queue Size</b>	?	<b>No value set</b>
Record Reader	?	No value set
Record Writer	?	No value set
<b>Add attributes as fields</b>	?	<b>true</b>
Message Demarcator	?	No value set

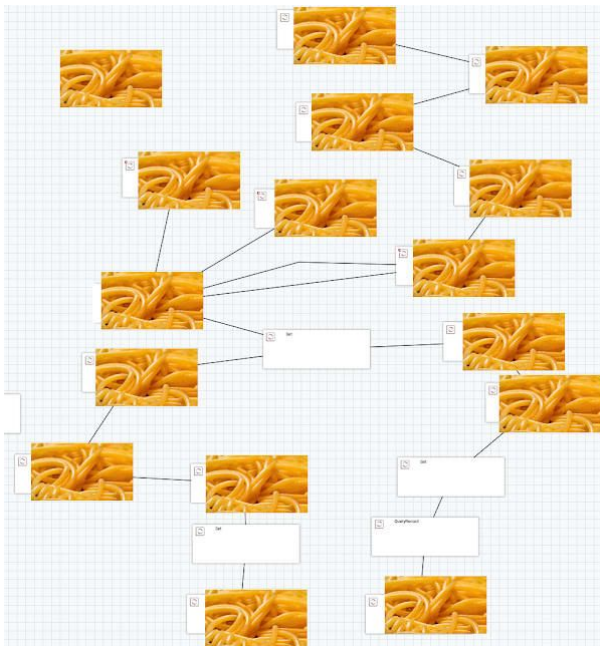
# Listen FTP



Let Apache NiFi be your FTP server

	 <b>ListenFTP</b> ListenFTP 1.13.0 org.apache.nifi - nifi-standard-nar	
In	0 (0 bytes)	5 min
Read/Write	0 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	0 / 00:00:00.000	5 min

# No More Spaghetti Flows - DO NOT

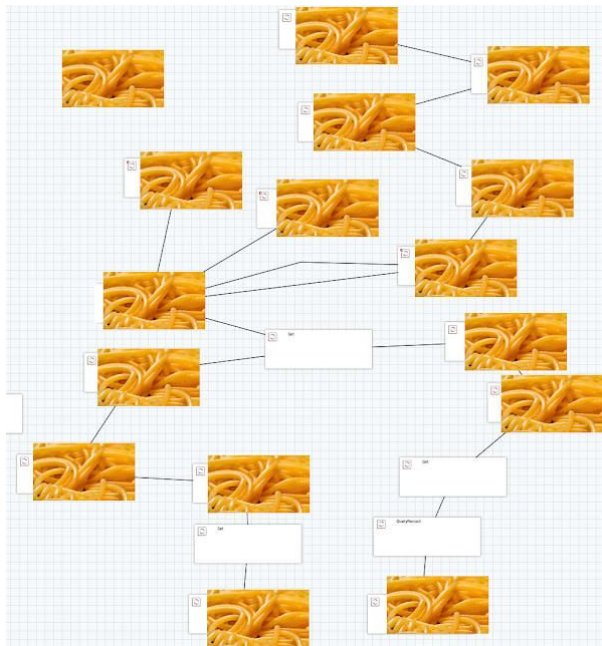


## Do Not

- Do not Put 1,000 Flows on one workspace.
- If your flow has hundreds of steps, this is a Flow Smell. Investigate why.
- Do not Use ExecuteProcess, ExecuteScripts or a lot of Groovy scripts as a default, look for existing processors
- Do not Use Random Custom Processors you find that have no documentation or are unknown.
- Do not forget to upgrade, if you are running anything before Apache NiFi 1.14, upgrade now!
- Do not run on default 512M RAM.
- Do not run one node and think you have a highly available cluster.
- Do not split a file with millions of records to individual records in one shot without checking available space/memory and back pressure.
- Use Split processors only as an absolute last resort. Many processors are designed to work on FlowFiles that contain many records or many lines of text. Keeping the FlowFiles together instead of splitting them apart can often yield performance that is improved by 1-2 orders of magnitude.

<https://dev.to/tspannhw/no-more-spaghetti-flows-2emd>

# No More Spaghetti Flows - DO

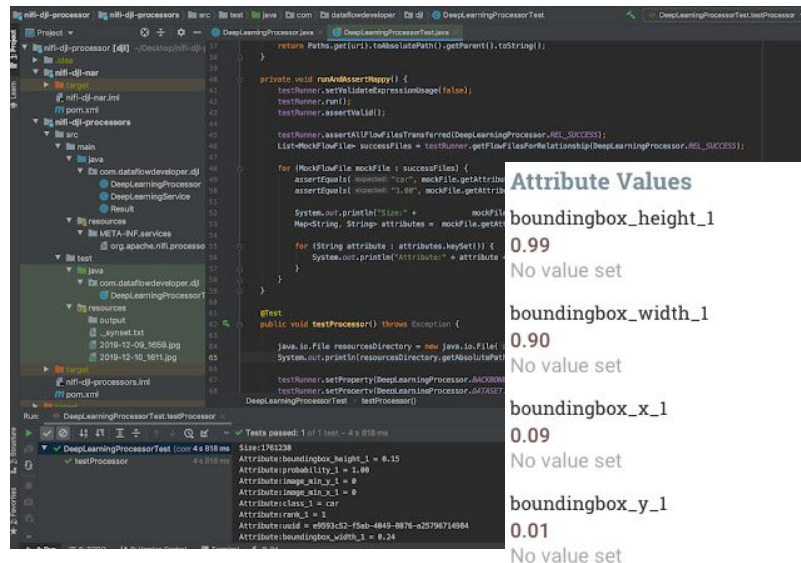
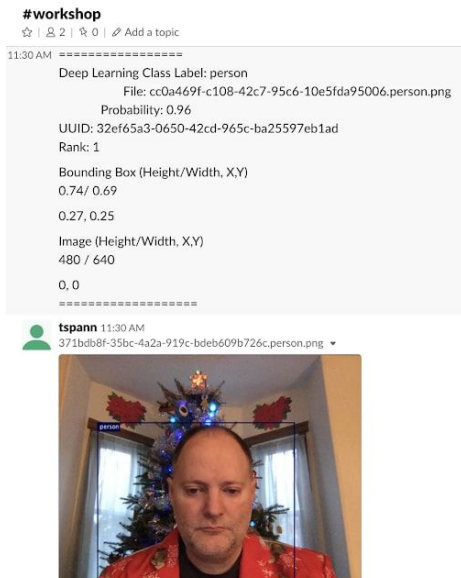


## Do

- Reduce, Reuse, Recycle. Use Parameters to reuse common modules.
- Put flows, reusable chunks (write to Slack, Database, Kafka) into separate Process Groups.
- Write custom processors if you need new or specialized features
- Use Record Processors everywhere
- Read the Docs!
- Use the NiFi Registry for version control.
- Use NiFi CLI and DevOps for Migrations.
- Walk through your flow and make sure you understand every step and it's easy to read and follow. Is every processor used? Are there dead ends?
- Do run Zookeeper on different nodes from Apache NiFi.
- Use routing based on content and attributes to allow one flow to handle multiple nearly identical flows is better than deploying the same flow many times with tweaks to parameters in same cluster.
- Use the correct driver for your database. There's usually a couple different JDBC drivers.

<https://dev.to/tspannhw/no-more-spaghetti-flows-2emd>

# Apache MXNet Native Processor through DJL.AI for Apache NiFi



## Attribute Values

boundingbox\_height\_1

0.99  
No value set

boundingbox\_width\_1

0.90  
No value set

boundingbox\_x\_1

0.09  
No value set

boundingbox\_y\_1

0.01  
No value set

class\_1

tvmonitor  
No value set

filename

2020-08-26\_1330.jpg.tvmonitor.png  
2020-08-26\_1330.jpg (previous)

This processor uses the DJL.AI Java Interface

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

<https://dev.to/tspannhw/easy-deep-learning-in-apache-nifi-with-djl-2d79>





# What are the Benefits of Pulsar?



Multi-Tenancy

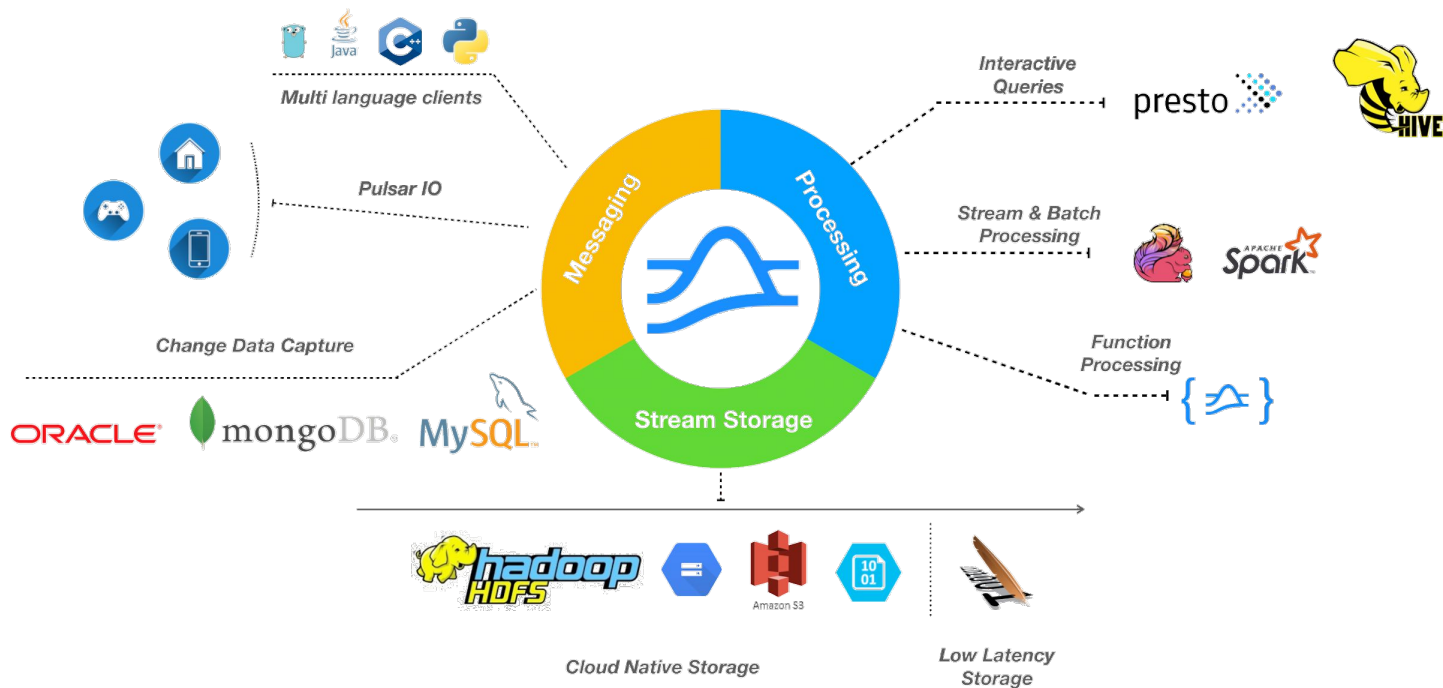
Scalability

Geo-Replication

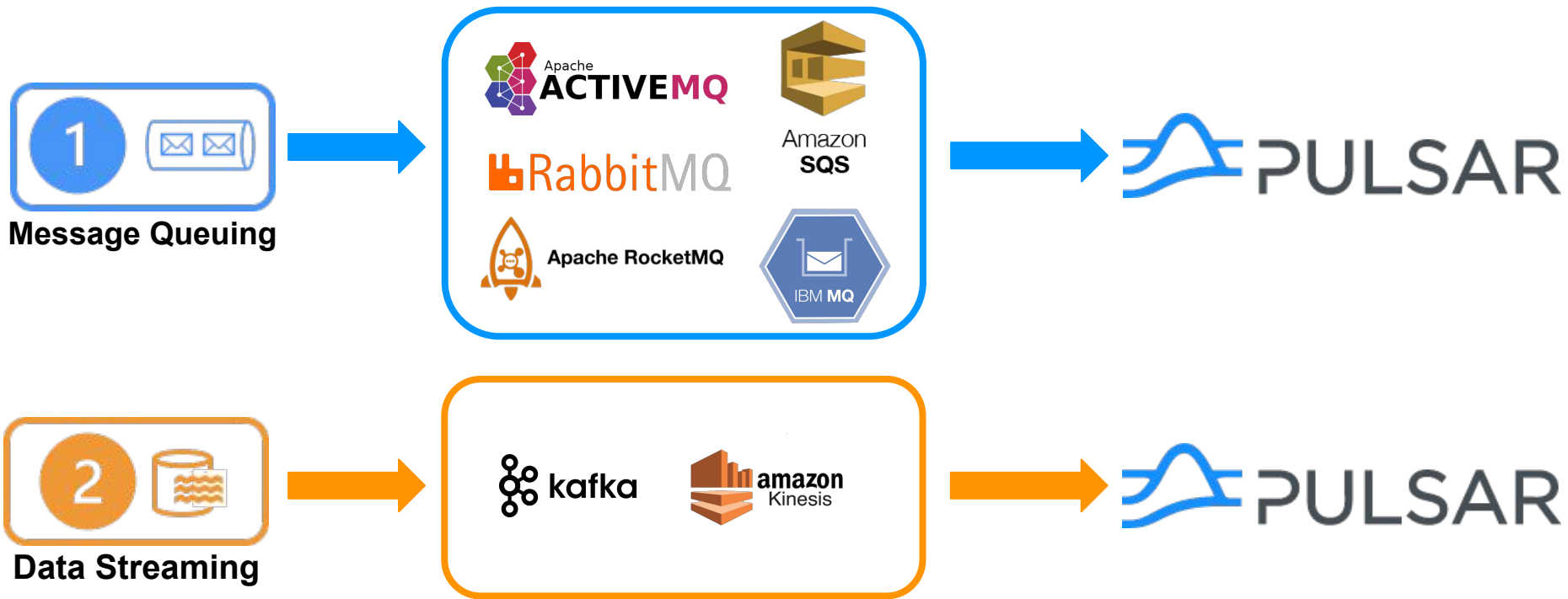
Unified Messaging  
Model

Data Durability

# Apache Pulsar



# A Unified Messaging Platform



# Demo



# Wrap-Up



# Stream Native

Founded by the original developers of Apache Pulsar and Apache BookKeeper, StreamNative builds a cloud-native event streaming platform that enables enterprises to easily access data as real-time event streams.

# Interested In Learning More?



## Resources

[Flink SQL Cookbook](#)

[The Github Source for Flink SQL Demo](#)

[The GitHub Source for Demo](#)



## Free eBooks

[Manning's Apache Pulsar in Action](#)

[O'Reilly Book](#)



## Upcoming Events

[11/8] [PASS Data Community](#)

[11/18] [Developer Week Austin](#)

[11/19] [Porto Tech Hub Con](#)

[12/3] [Data Science Camp](#)

Platform Engineer [Remote]

San Francisco

Platform Engineer (Flink/Spark) [Remote]

San Francisco

Product Engineer - Cloud [Remote]

San Francisco

Platform Engineer (Flink/Spark) [Remote]

San Francisco

Product Engineer - Cloud [Remote]

San Francisco

Sr. Product Manager [Remote]

San Francisco

# We're Hiring

[streamnative.io/careers/](https://streamnative.io/careers/)



# Let's Keep in Touch!



**Speaker Name**

Speaker title



@PassDev



<https://www.linkedin.com/in/timothyspann>



<https://github.com/tspannhw>