

# Apache NiFi and MiNiFi: Edge to Core

**Andy LoPresto - @yolopey**

Apache NiFi PMC

DataWorks Summit 2017 - Sydney

19 Sep 2017



# Agenda

*What is dataflow and what are the challenges?*

*Apache NiFi*

IoT Challenges

Apache MiNiFi

Exploration

Community

# Agenda



*What is dataflow and what are the challenges?*

*Apache NiFi*

*IoT Challenges*

*Apache MiNiFi*

*Exploration*

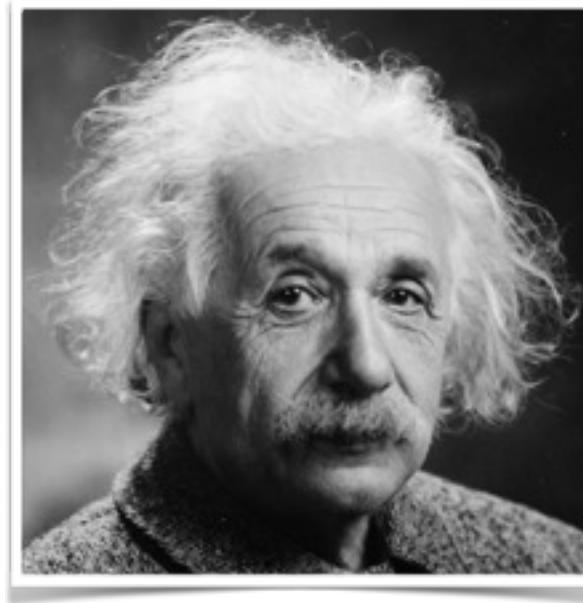
*Community*

# Gauging Audience Familiarity With NiFi



*"What's a NeeFee?"*

No experience with dataflow  
No experience with NiFi



*"I can pick this up pretty quickly"*

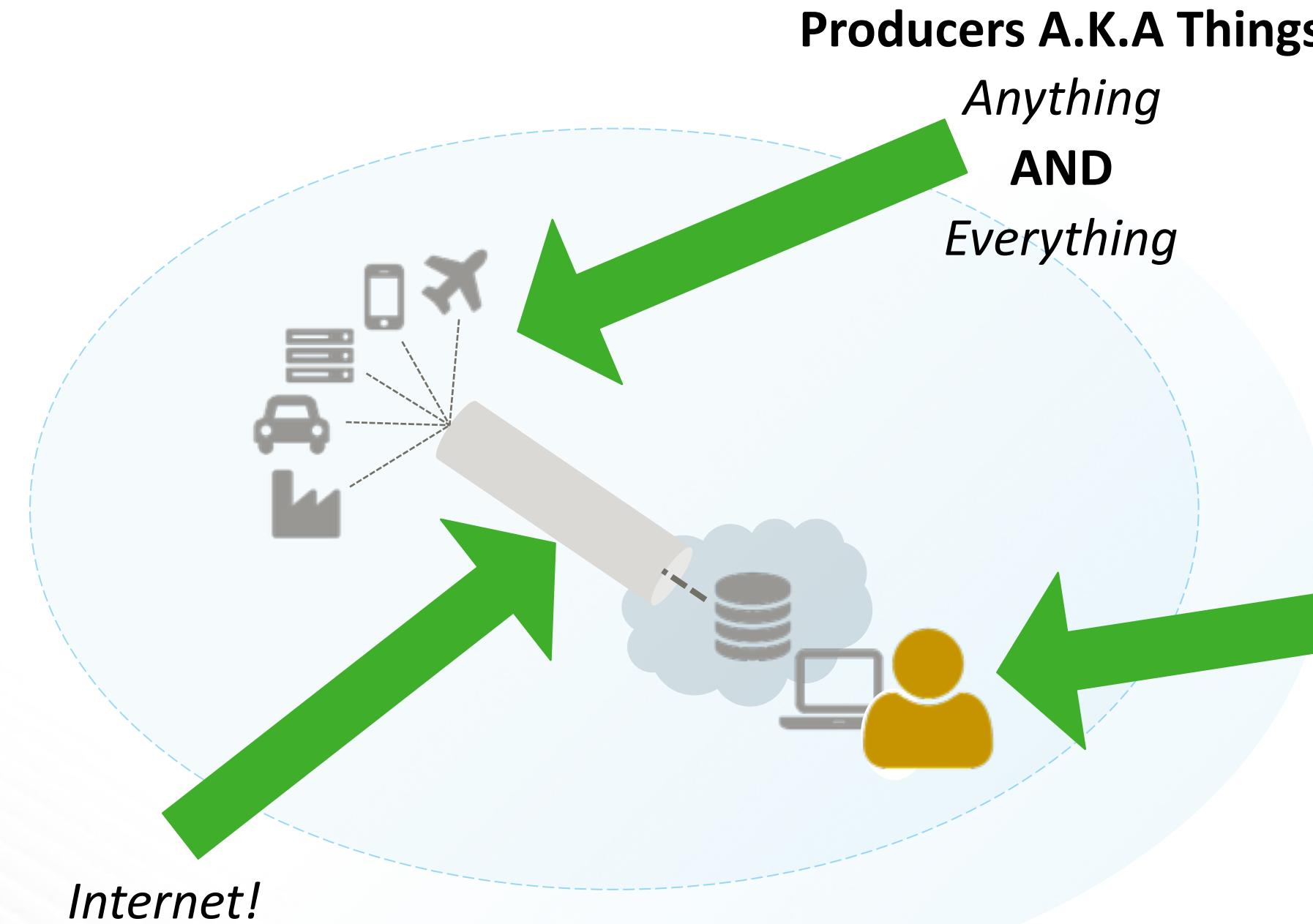
Some experience with dataflow  
Some experience with NiFi



*"I refactored the Ambari integration endpoint to allow for mutual authentication TLS during my coffee break"*

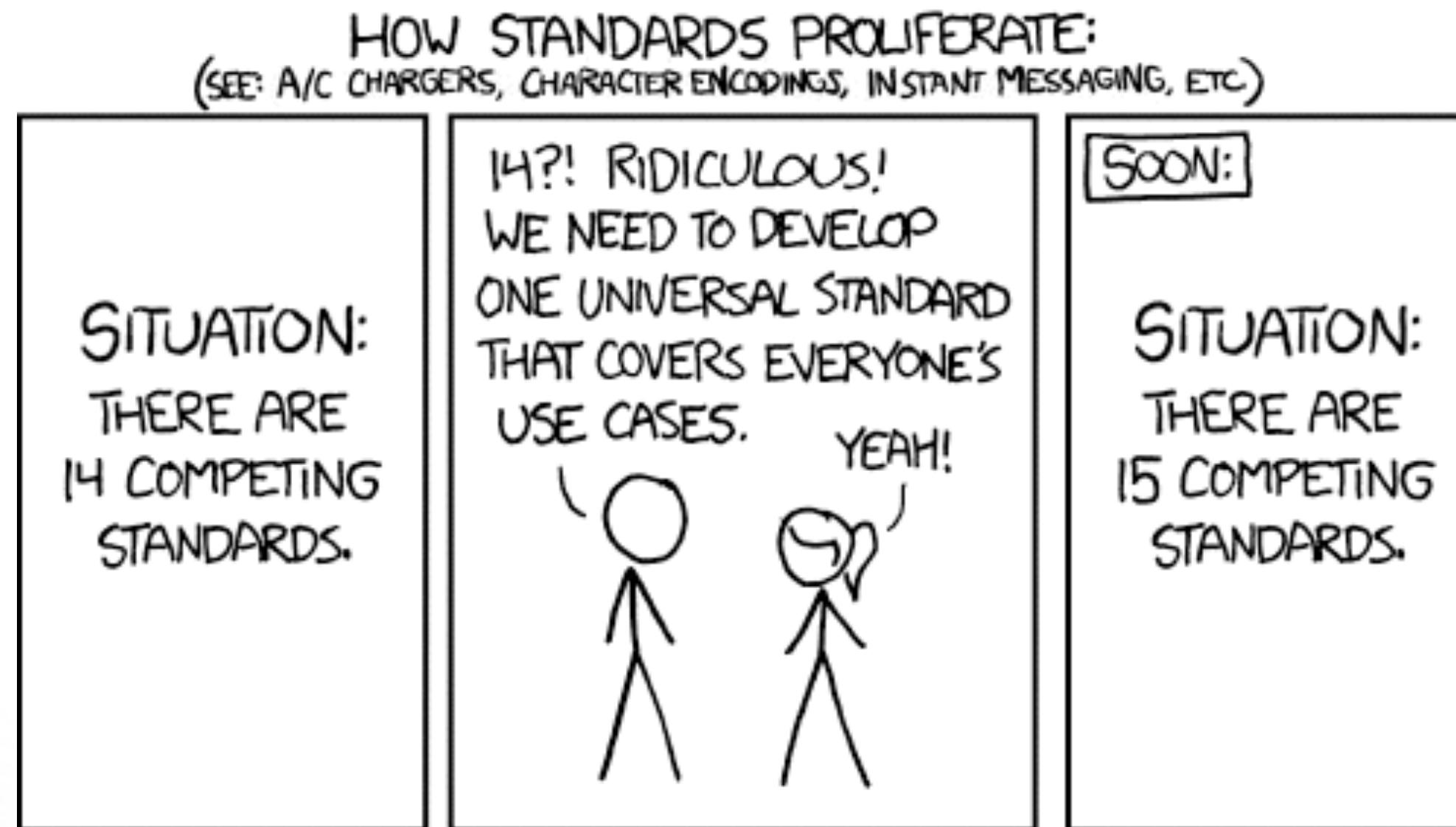
Forgotten more about NiFi  
than most of us will ever  
know

# Let's Connect A to B



- Consumers**
- User
  - Storage
  - System
  - ...More Things

# Moving data *effectively* is hard



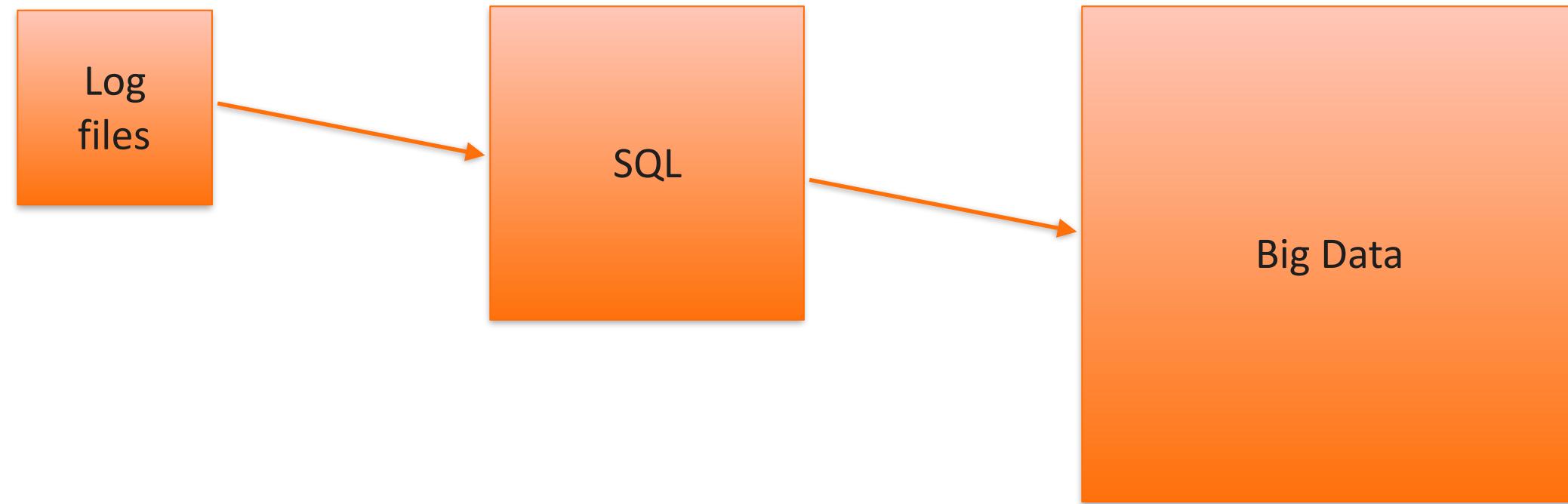
Standards: <http://xkcd.com/927/>

# Why is moving data *effectively* hard?

- ◆ Standards
- ◆ Formats
- ◆ “Exactly Once” Delivery
- ◆ Protocols
- ◆ Veracity of Information
- ◆ Validity of Information
- ◆ Ensuring Security
- ◆ Overcoming Security
- ◆ Compliance
- ◆ Schemas
- ◆ Consumers Change
- ◆ Credential Management
- ◆ “*That [person | team | group]*”
- ◆ Network\*
- ◆ “Exactly Once” Delivery

# Connecting A to B to C

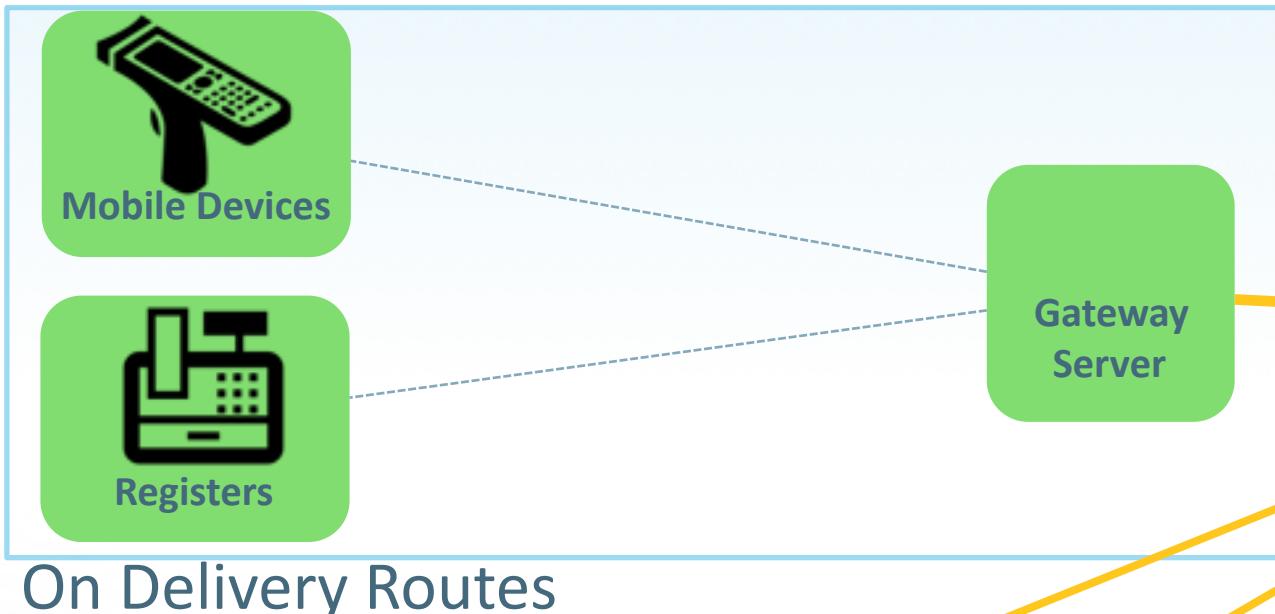
Easy enough with Bash scripts, Ruby/Python/Groovy, etc.



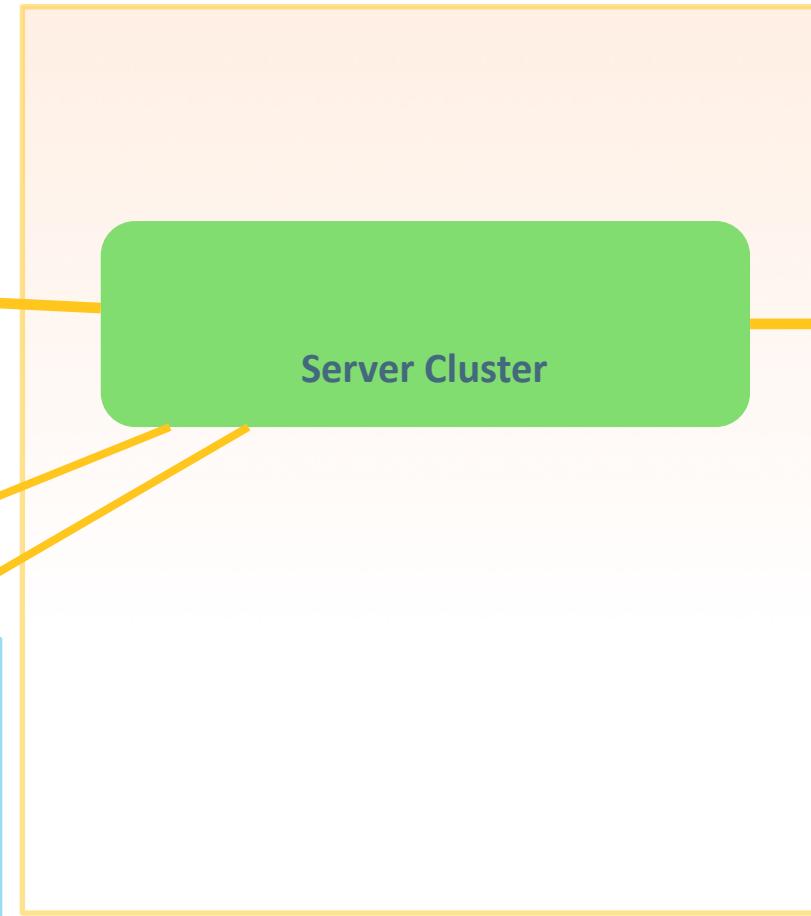
# Let's Connect Lots of As to Bs to As to Cs to Bs to Δs to Cs to ϕs

Let's consider the needs of a courier service

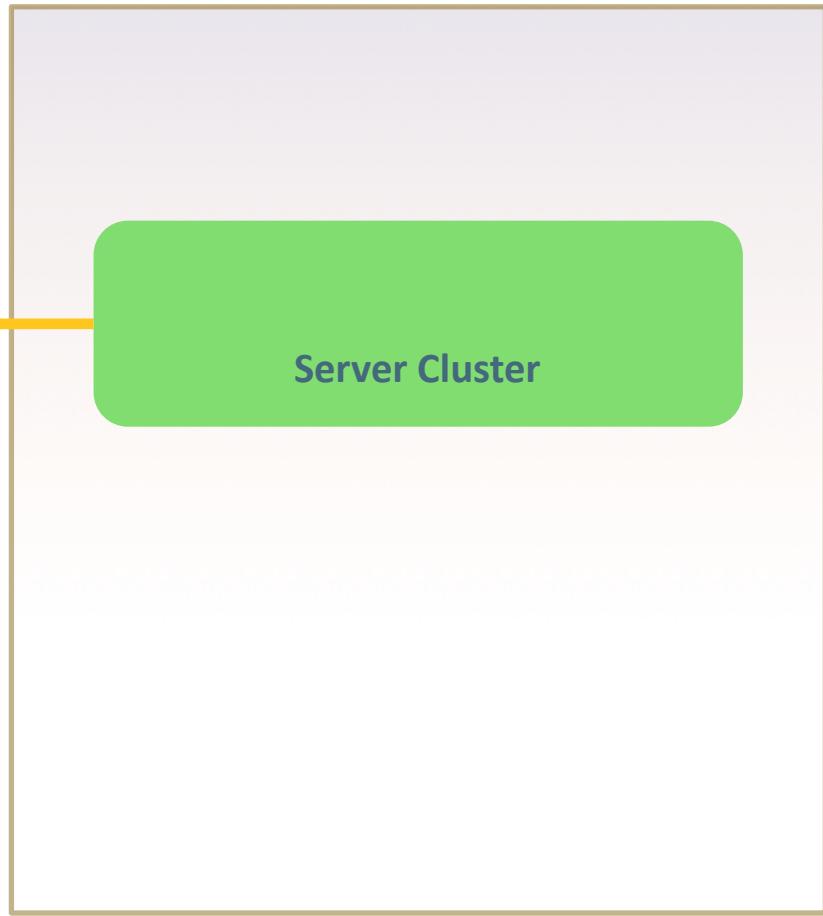
Physical Store



Distribution Center



Core Data Center at HQ



On Delivery Routes



Delivery Truck: Creative Stall, <https://thenounproject.com/creativestall/>

Deliverer: Rigo Peter, <https://thenounproject.com/rigo/>

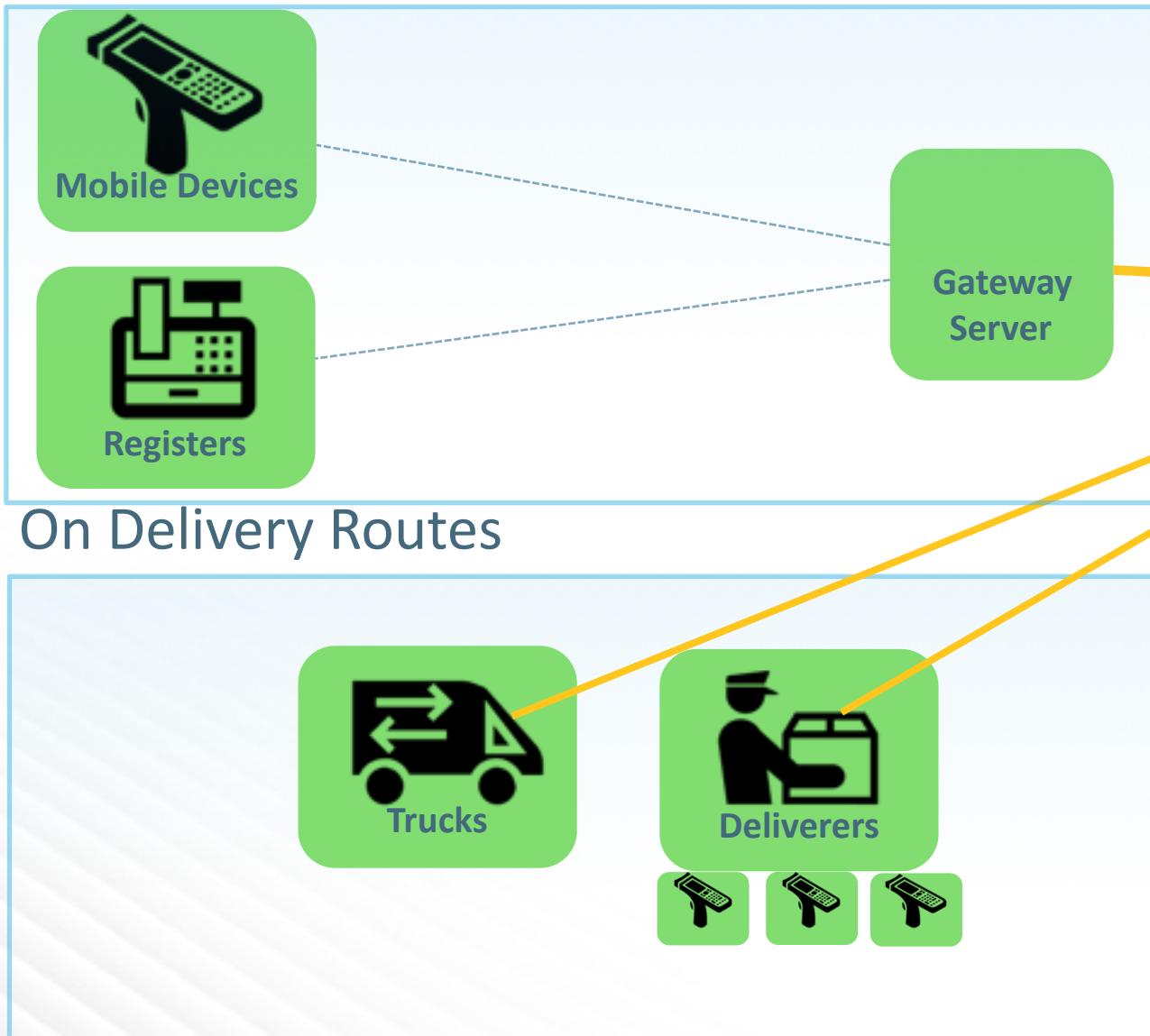
Cash Register: Sergey Patutin, <https://thenounproject.com/bdesign.by/>

Hand Scanner: Eric Pearson, <https://thenounproject.com/epearson001/>

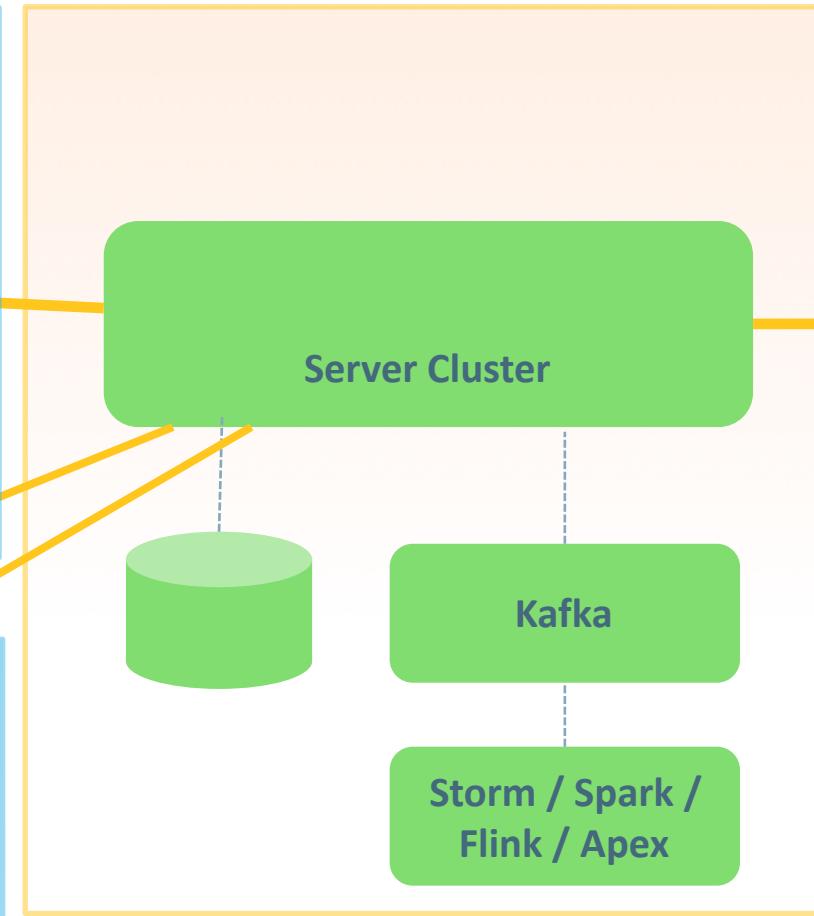
# Great! I am collecting all this data! Let's use it!

Finding our needles in the haystack

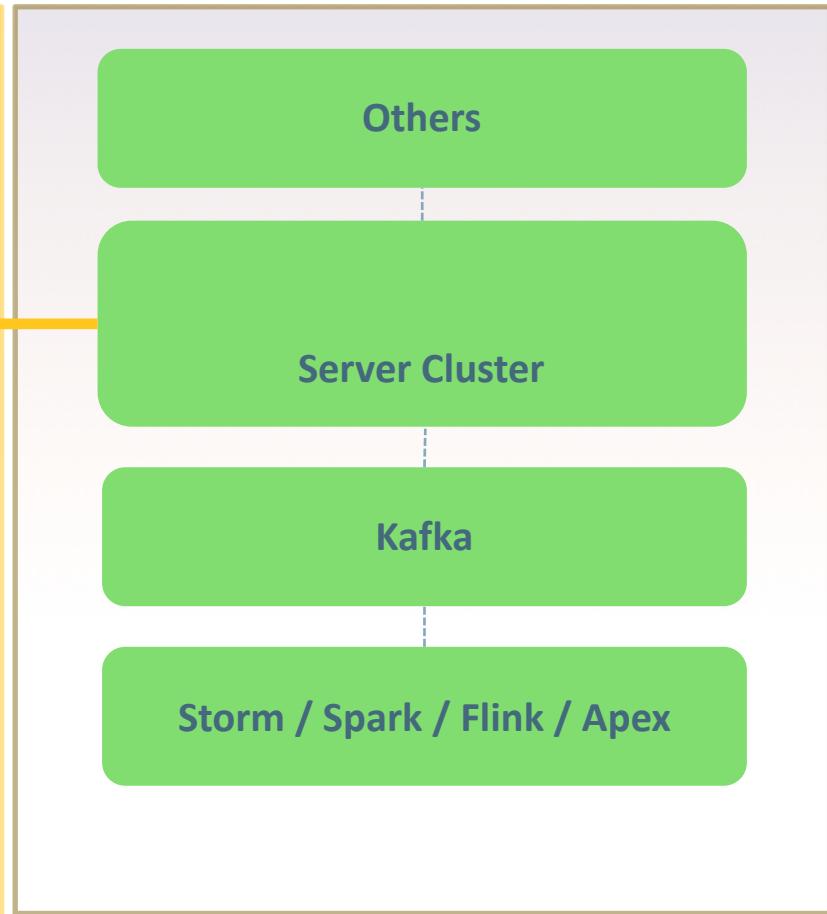
Physical Store



Distribution Center

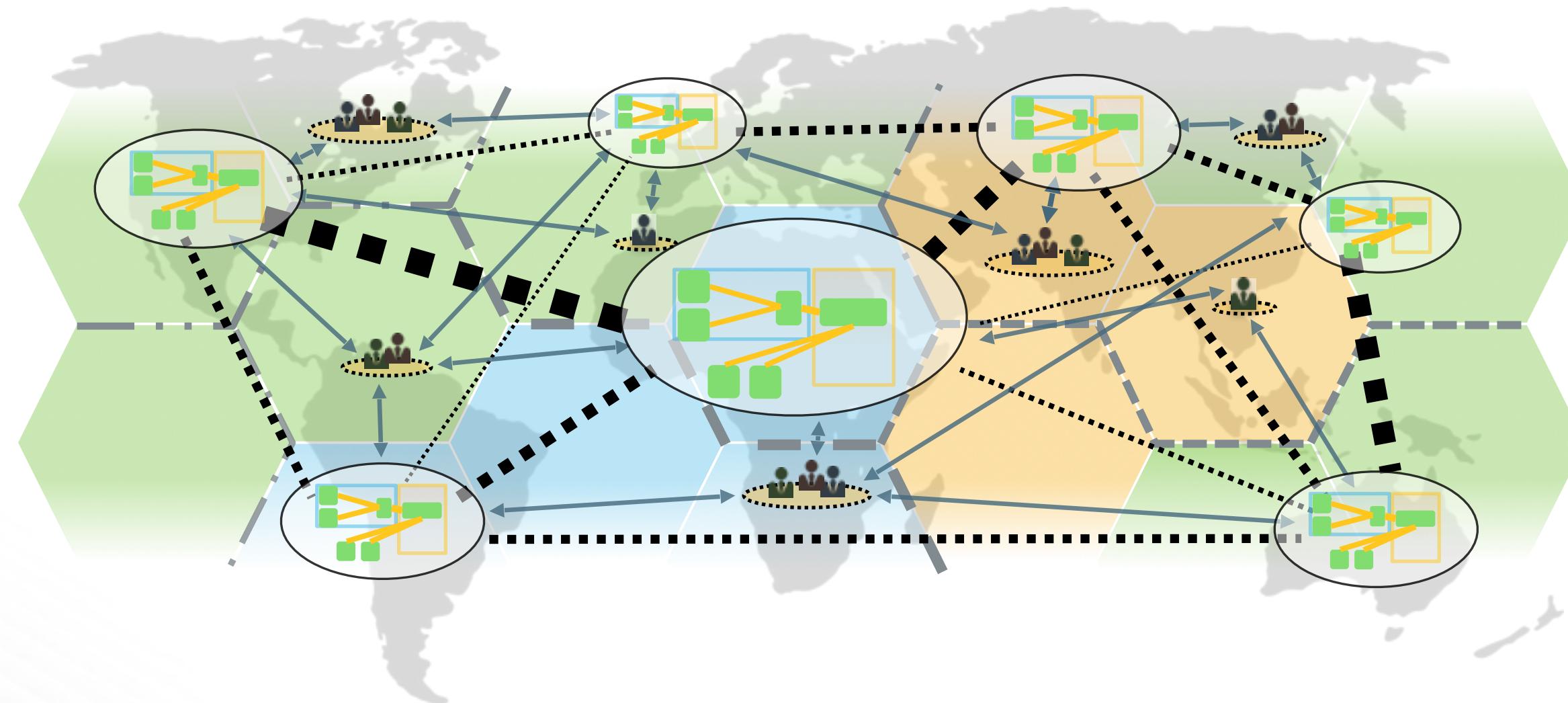


Core Data Center at HQ



# Let's Connect Lots of As to Bs to As to Cs to Bs to $\Delta$ s to Cs to $\varphi$ s

Raise your hand if you want to maintain Python scripts for the rest of your life



# Agenda



*What is dataflow and what are the challenges?*

*Apache NiFi*

IoT Challenges

Apache MiNiFi

Exploration

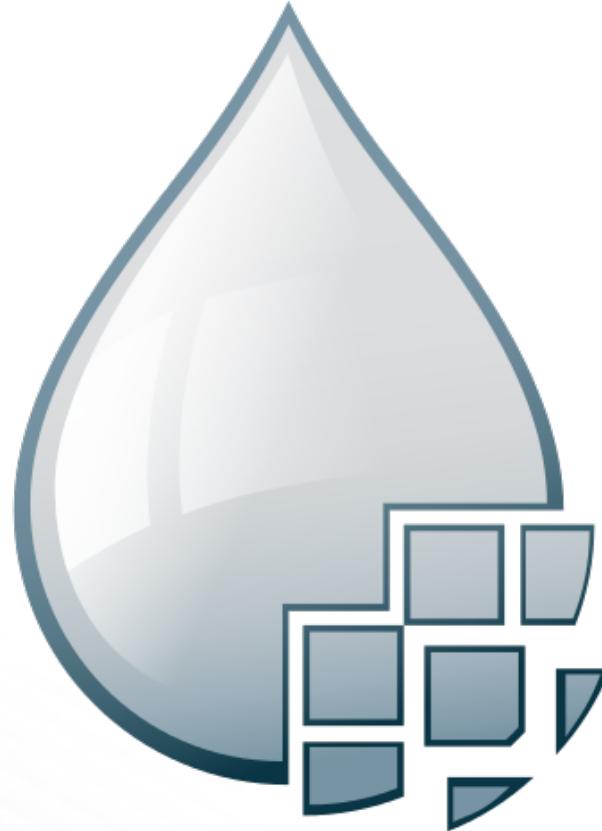
Community

# NiFi is based on Flow Based Programming (FBP)

FBP Term	NiFi Term	Description
Information Packet	FlowFile	Each object moving through the system.
Black Box	FlowFile Processor	Performs the work, doing some combination of data routing, transformation, or mediation between systems.
Bounded Buffer	Connection	The linkage between processors, acting as queues and allowing various processes to interact at differing rates.
Scheduler	Flow Controller	Maintains the knowledge of how processes are connected, and manages the threads and allocations thereof which all processes use.
Subnet	Process Group	A set of processes and their connections, which can receive and send data via ports. A process group allows creation of entirely new component simply by composition of its components.

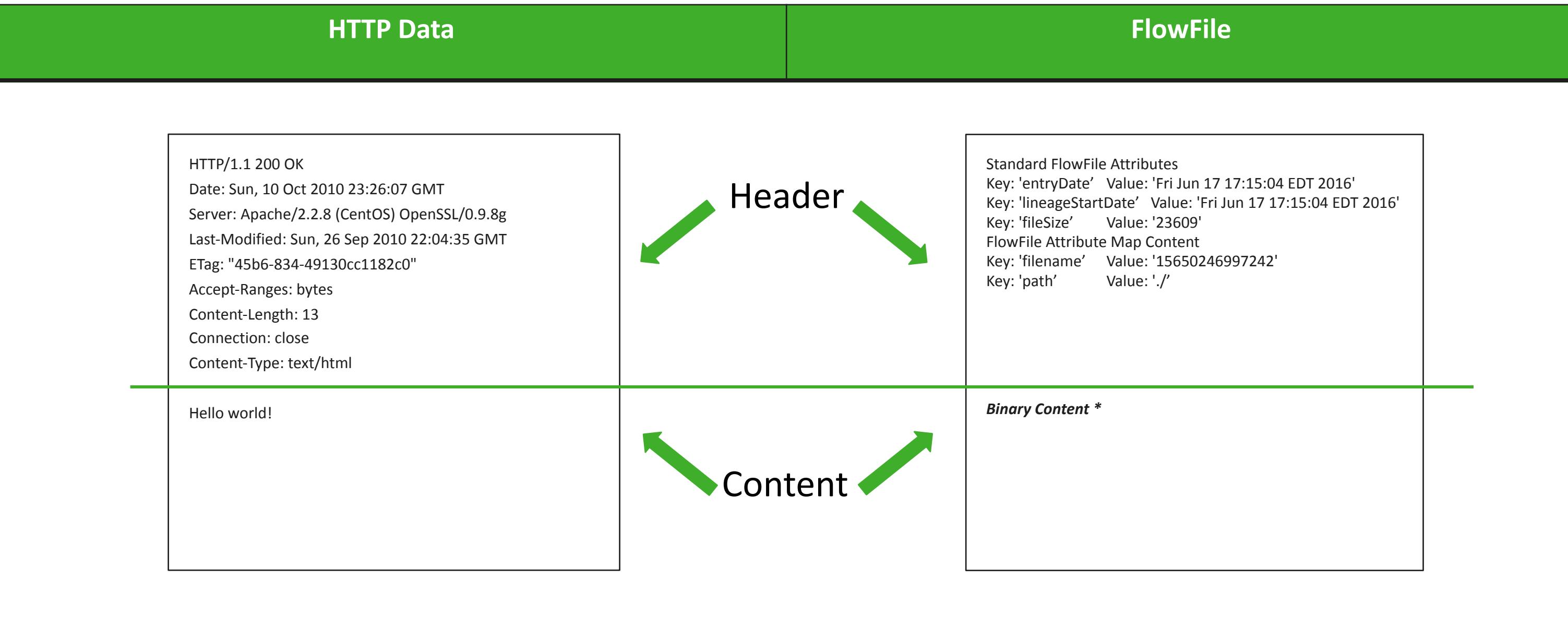
# Apache NiFi

## Key Features



- Guaranteed delivery
- Data buffering
  - Backpressure
  - Pressure release
- Prioritized queuing
- Flow specific QoS
  - Latency vs. throughput
  - Loss tolerance
- Data provenance
- Supports push and pull models
- Recovery/recording a rolling log of fine-grained history
- Visual command and control
- Flow templates
- Pluggable, multi-tenant security
- Designed for extension
- Clustering

# FlowFiles are like HTTP data

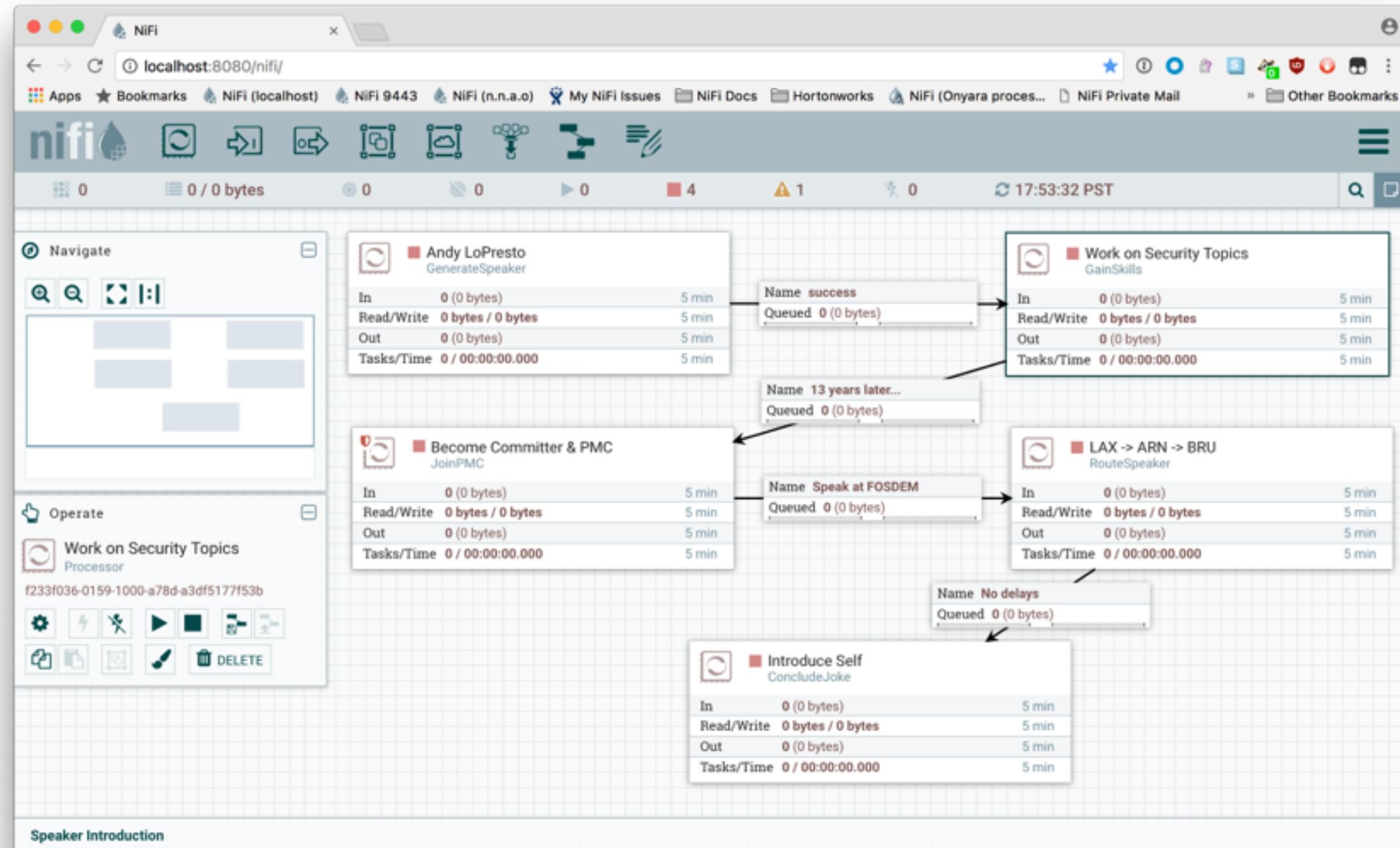


# User Interface

# Less of this... ... more of this



# User Interface



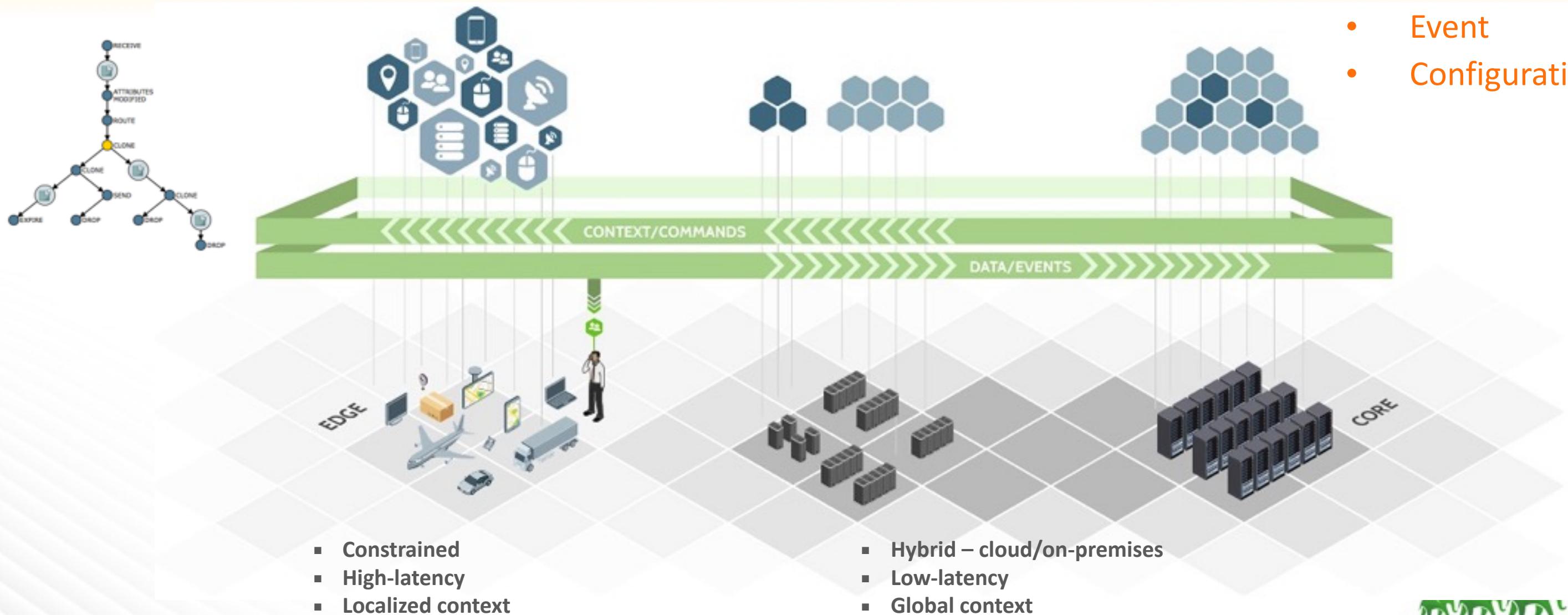
# Data Provenance

Origin – attribution  
Replay – recovery

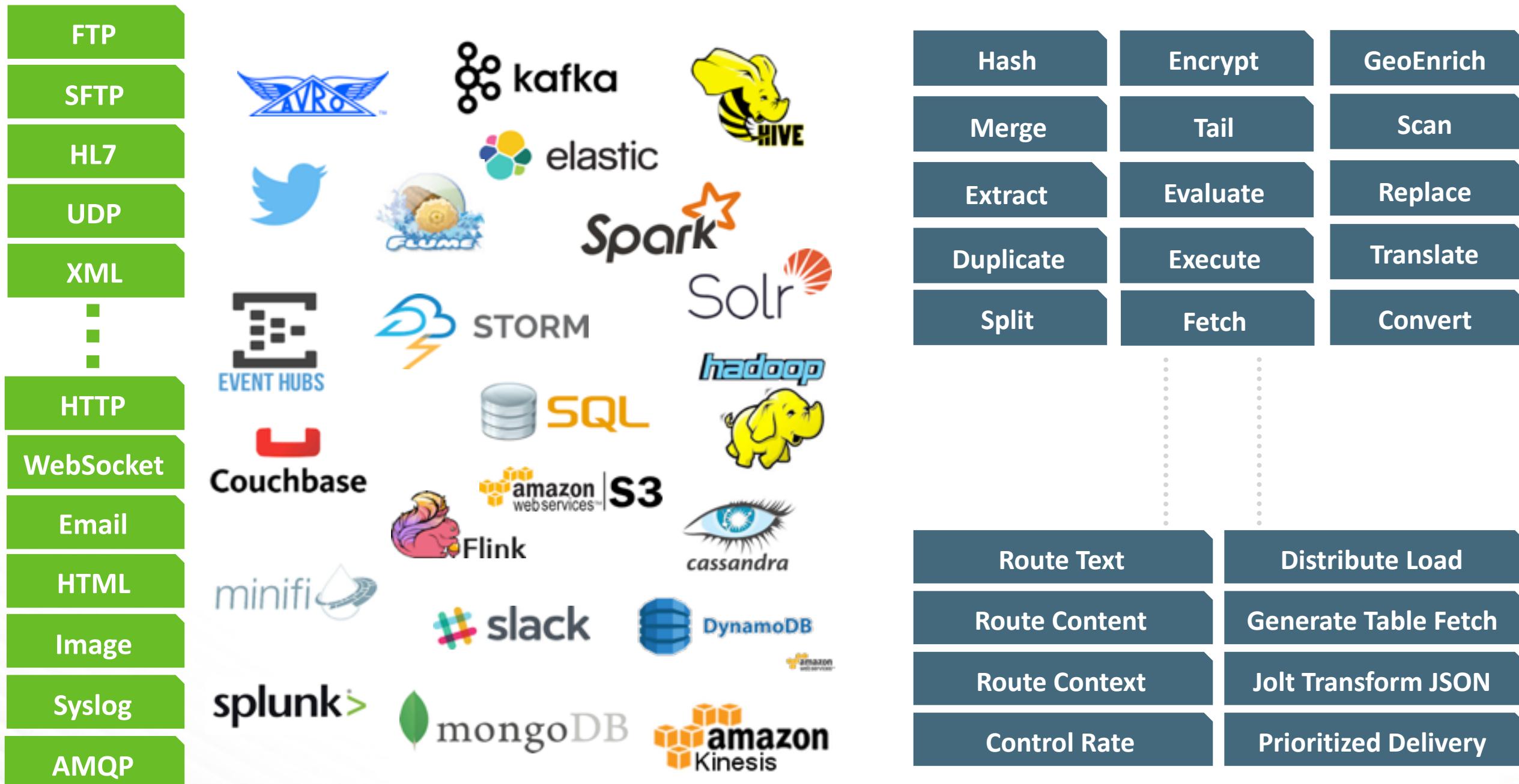
Evolution of topologies  
Long retention

## Types of Lineage

- Event
- Configuration



# Deeper Ecosystem Integration: 220+ Processors



All Apache project logos are trademarks of the ASF and the respective projects.

# Edge Challenges

- ◆ Limited computing capability
- ◆ Limited power/network
- ◆ Restricted software library/platform availability
- ◆ No UI
- ◆ Physically inaccessible
- ◆ Not frequently updated
- ◆ Competing standards/protocols
- ◆ Scalability
- ◆ Privacy & Security

# Recent Examples

- When the Mirai attack has its own Wikipedia page, that's not good

Dyn Analysis Summary Of Friday October 21 Attack  
Company News // Oct 26, 2016 // Scott Hilton

## IoTPOT: Analysing the Rise of IoT Compromises

Yin Minn Pa Pa<sup>†1</sup>, Shogo Suzuki<sup>†1</sup>, Katsunari Yoshioka<sup>†1</sup>, Tsutomu Matsumoto<sup>†1</sup>,  
Takahiro Kasama<sup>†2</sup>, Christian Rossow<sup>†3</sup>

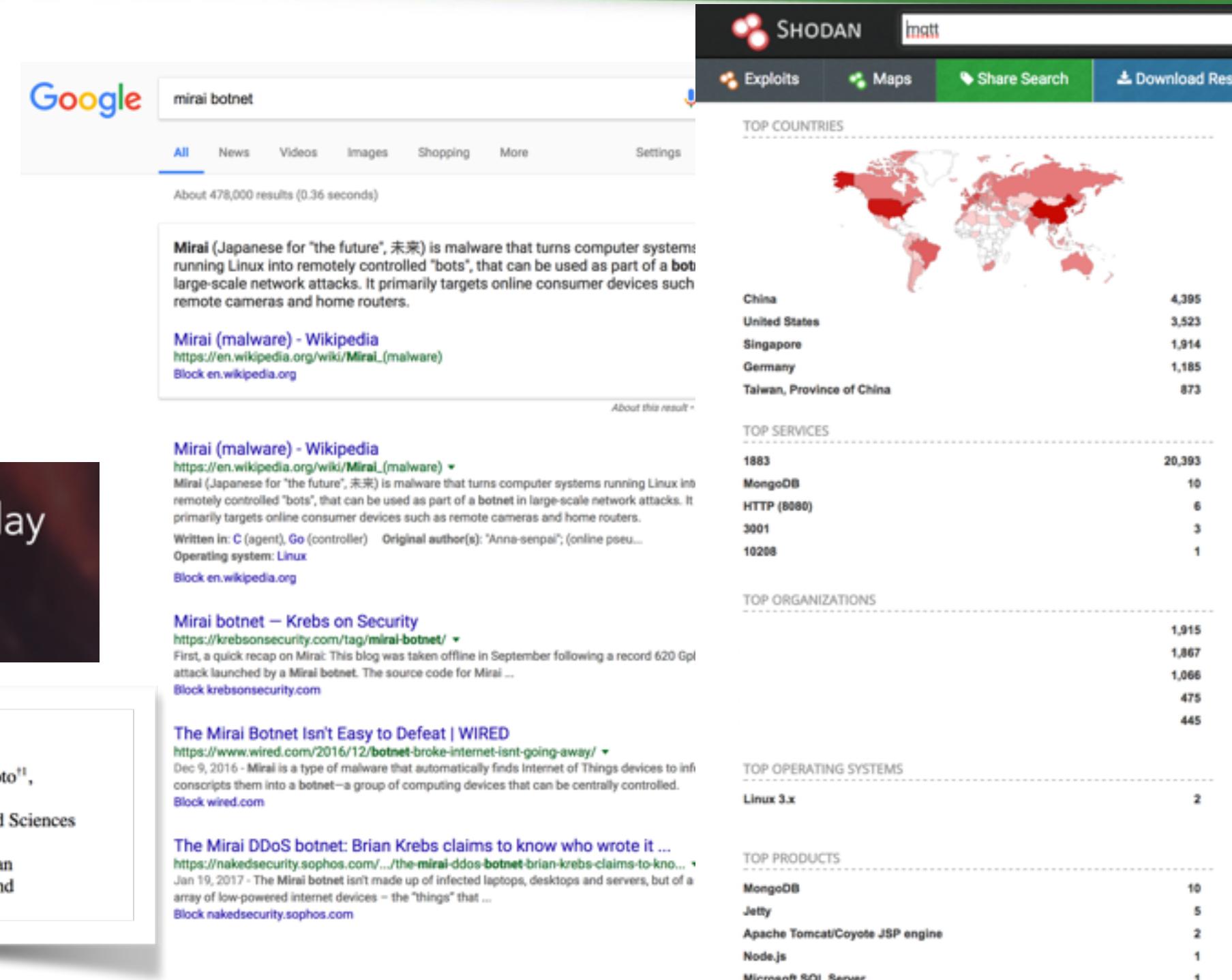
<sup>†1</sup>Graduate School of Environment and Information Sciences/Institute of Advanced Sciences

<sup>†1</sup> Yokohama National University, Japan

<sup>†2</sup>National Institute of Information and Communications Technology, Japan

<sup>†3</sup>Institute of Advanced Sciences, Yokohama National University, Japan and

<sup>†3</sup>Cluster of Excellence, MMCI, Saarland University, Germany



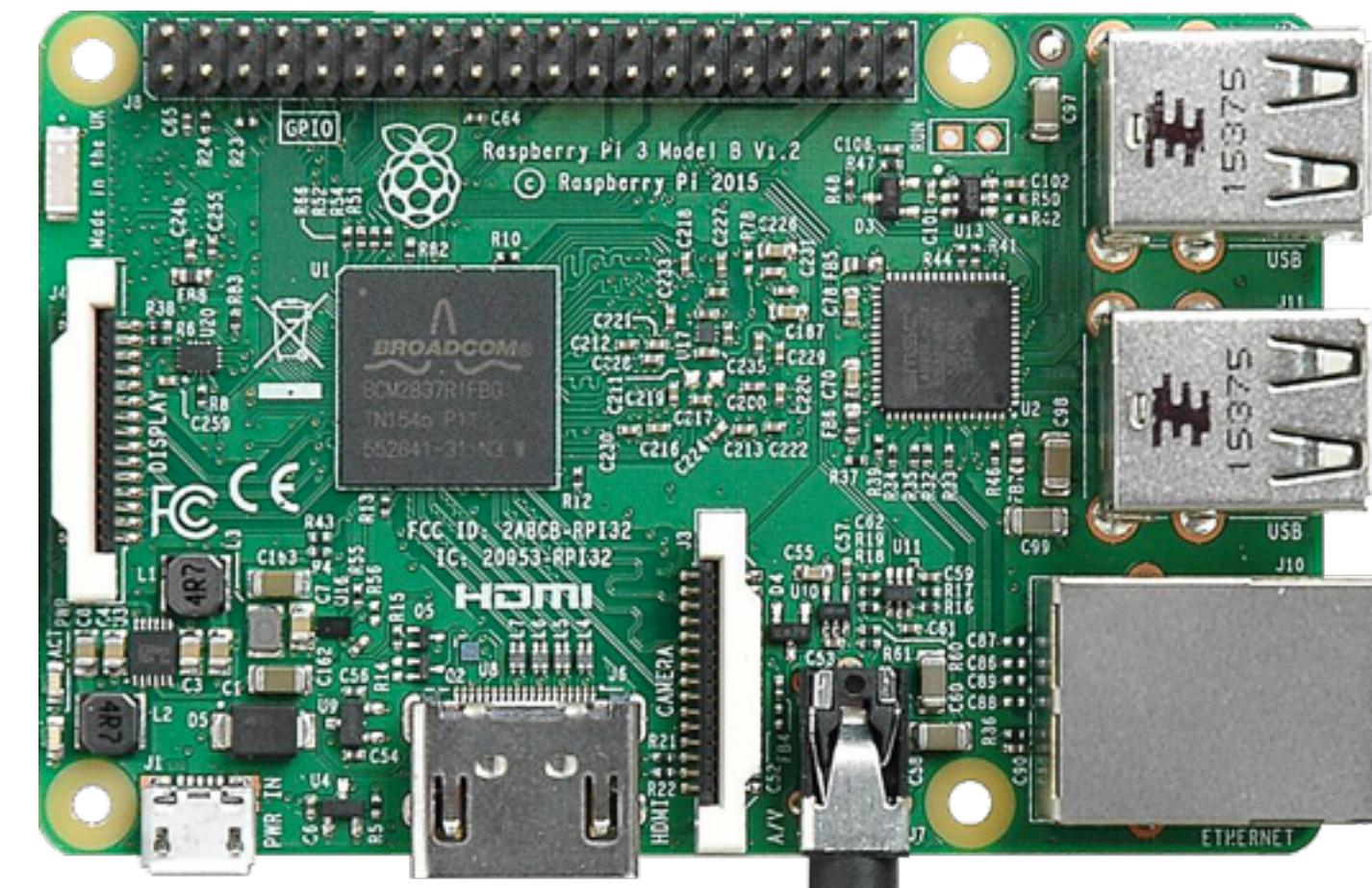
# NiFi Solves Everything\*

- ◆ Runs on JVM
- ◆ Provides UI for flow design & monitoring
- ◆ Security built-in
  - ◆ TLS, authn/authz, encrypted data
- ◆ Handles practically any format/protocol

# NiFi for IoT

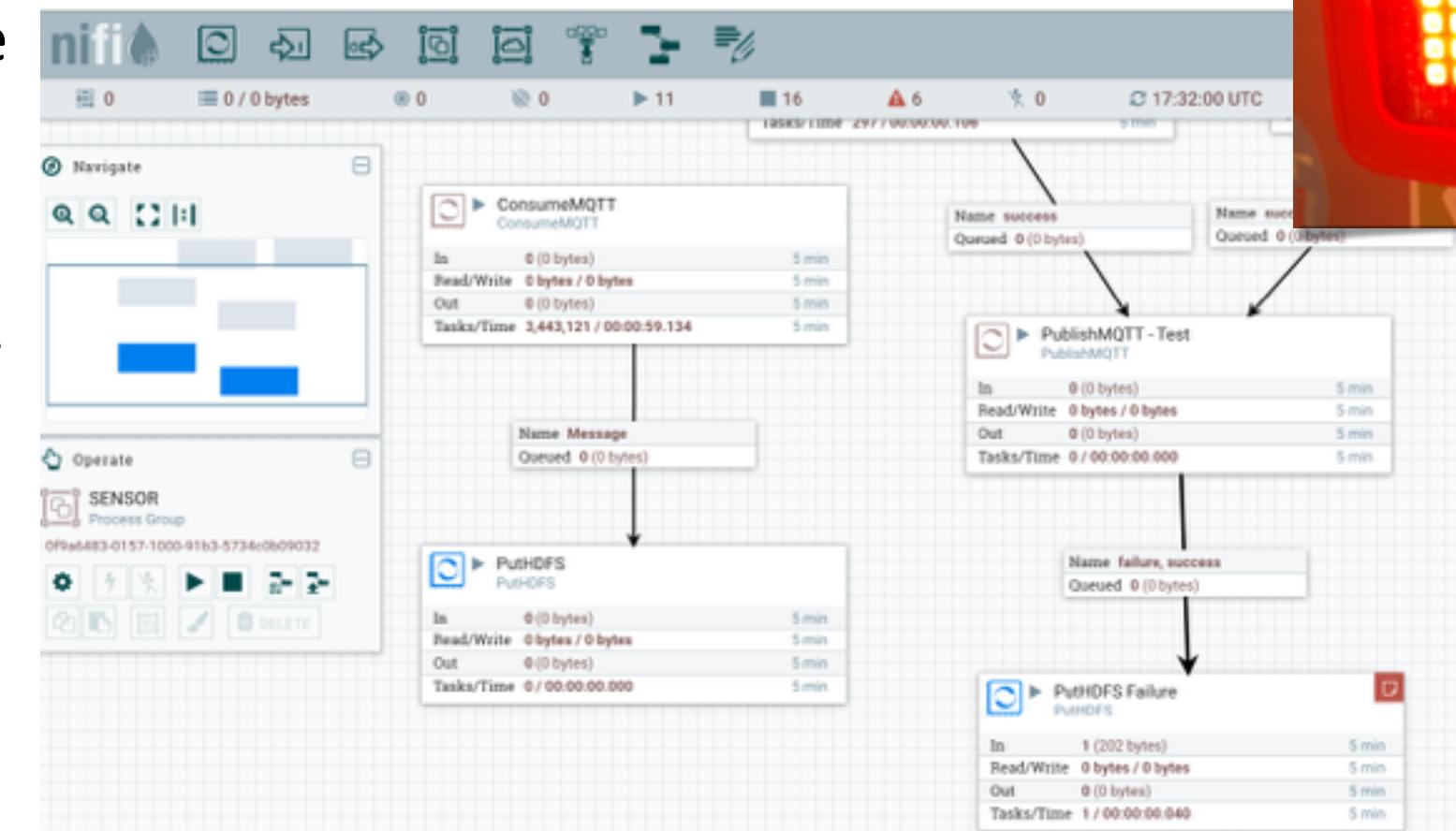
- ◆ NiFi supports AMQP, MQTT, UDP, TCP, HTTP(S), CEF, JMS, (S)FTP, *AWSIoT*
- ◆ With a little pruning, NiFi can run on a Raspberry Pi

```
├── bootstrap
├── jcl-over-slf4j-1.7.12.jar
├── jul-to-slf4j-1.7.12.jar
├── log4j-over-slf4j-1.7.12.jar
├── logback-classic-1.1.3.jar
├── logback-core-1.1.3.jar
├── nifi-api-0.6.1.jar
├── nifi-documentation-0.6.1.jar
├── nifi-framework-nar-0.6.1.nar
├── nifi-html-nar-0.6.1.nar
├── nifi-http-context-map-nar-0.6.1.nar
├── nifi-jetty-bundle-0.6.1.nar
├── nifi-kerberos-iaa-providers-nar-0.6.1.nar
├── nifi-ldap-iaa-providers-nar-0.6.1.nar
├── nifi-nar-utils-0.6.1.jar
├── nifi-properties-0.6.1.jar
├── nifi-provenance-repository-nar-0.6.1.nar
├── nifi-runtime-0.6.1.jar
├── nifi-scripting-nar-0.6.1.nar
├── nifi-ssl-context-service-nar-0.6.1.nar
├── nifi-standard-nar-0.6.1.nar
├── nifi-standard-services-api-nar-0.6.1.nar
├── nifi-update-attribute-nar-0.6.1.nar
└── slf4j-api-1.7.12.jar
```



# Example — Sensor Readings via RP3B

- ◆ Tim Spann
- ◆ Sense Hat sensor attachment
  - ◆ Temp, humidity, pressure
  - ◆ 8x8 LED display
- ◆ Python Flask server reading sensor and pushing to MQTT
- ◆ NiFi consuming MQTT



# So Why Do We Need A Different Solution?

- ◆ NiFi is designed to “own the box”
- ◆ NiFi 0.7.x started up in about 10-15 minutes on RP3 (593 MB)
- ◆ NiFi 1.x started up in about 30 minutes on RP3 (760 MB)
  - ◆ 33 new processors
  - ◆ Rewrite for multi tenant authorization
  - ◆ Complete UI overhaul

```
▶hw12203:/Users/alopresto/Workspace/scratch/rp3b-demo (master) alopresto
└─ 113s @ 17:09:05 $ ssh pi@my-raspberry-pi
^C
▶hw12203:/Users/alopresto/Workspace/scratch/rp3b-demo (master) alopresto
└─ 145s @ 17:09:37 $ █
```

# Apache NiFi Subproject: MiNiFi

- ◆ Get the key parts of NiFi close to where data begins and provide bidirectional communication
- ◆ NiFi lives in the data center — give it an enterprise server or a cluster of them
- ◆ MiNiFi lives as close to where data is born and is a guest on that device or system
  - ◆ IoT
  - ◆ Connected car
  - ◆ Legacy hardware

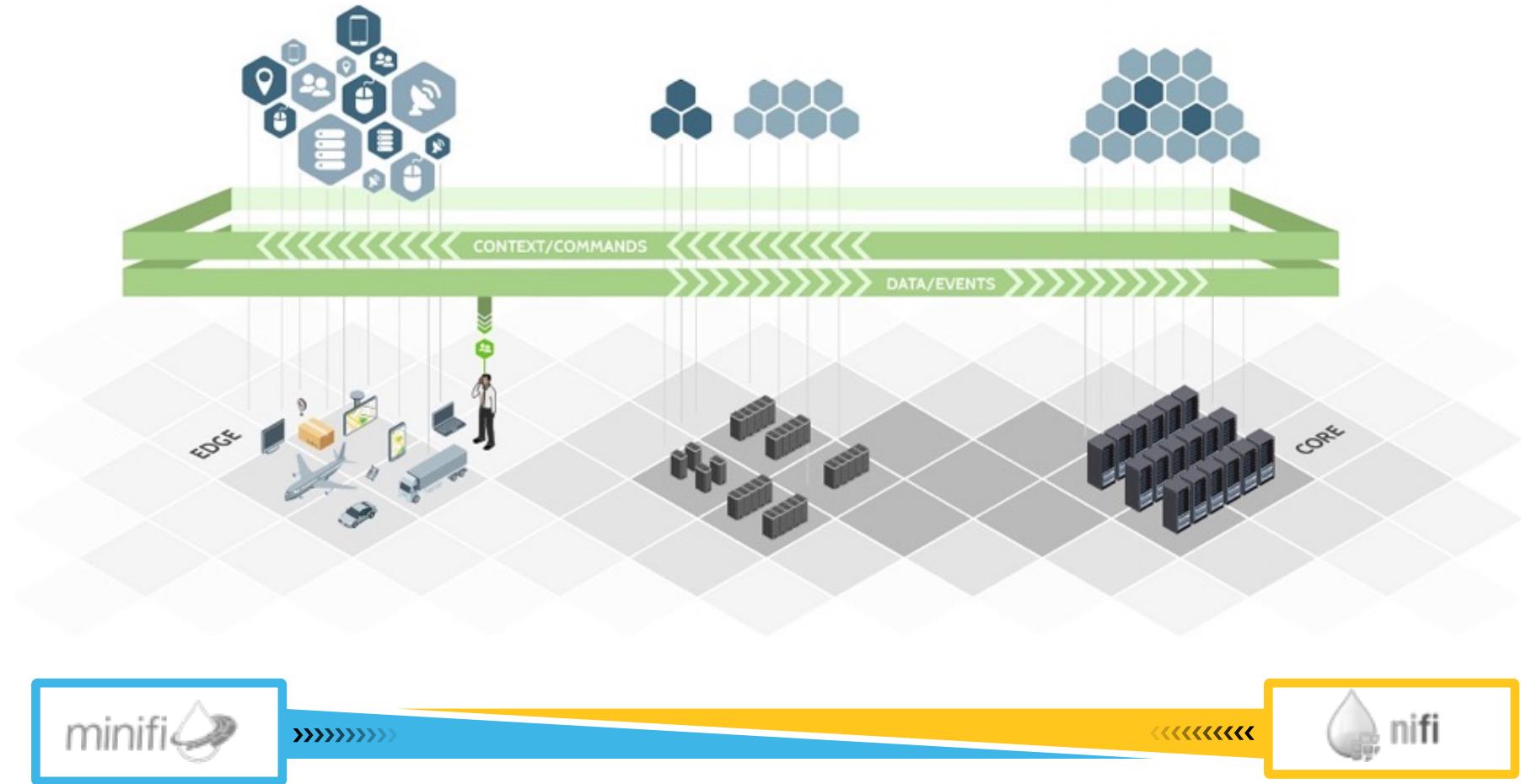


# Why build MiNiFi?

- ◆ NiFi is big
  - ◆ 1.3.0 release is 933 MB compressed
  - ◆ Can be modified to run in restricted environments, but requires manual surgery
  - ◆ Provides UI, provenance query, etc.
  - ◆ Runs on dedicated machines/clusters — “owns the box”
- ◆ MiNiFi lives at the edge
  - ◆ No UI
  - ◆ 0.1.0 Java binary is 45 MB, C++ binary is 746 KB
  - ◆ “Good guest”

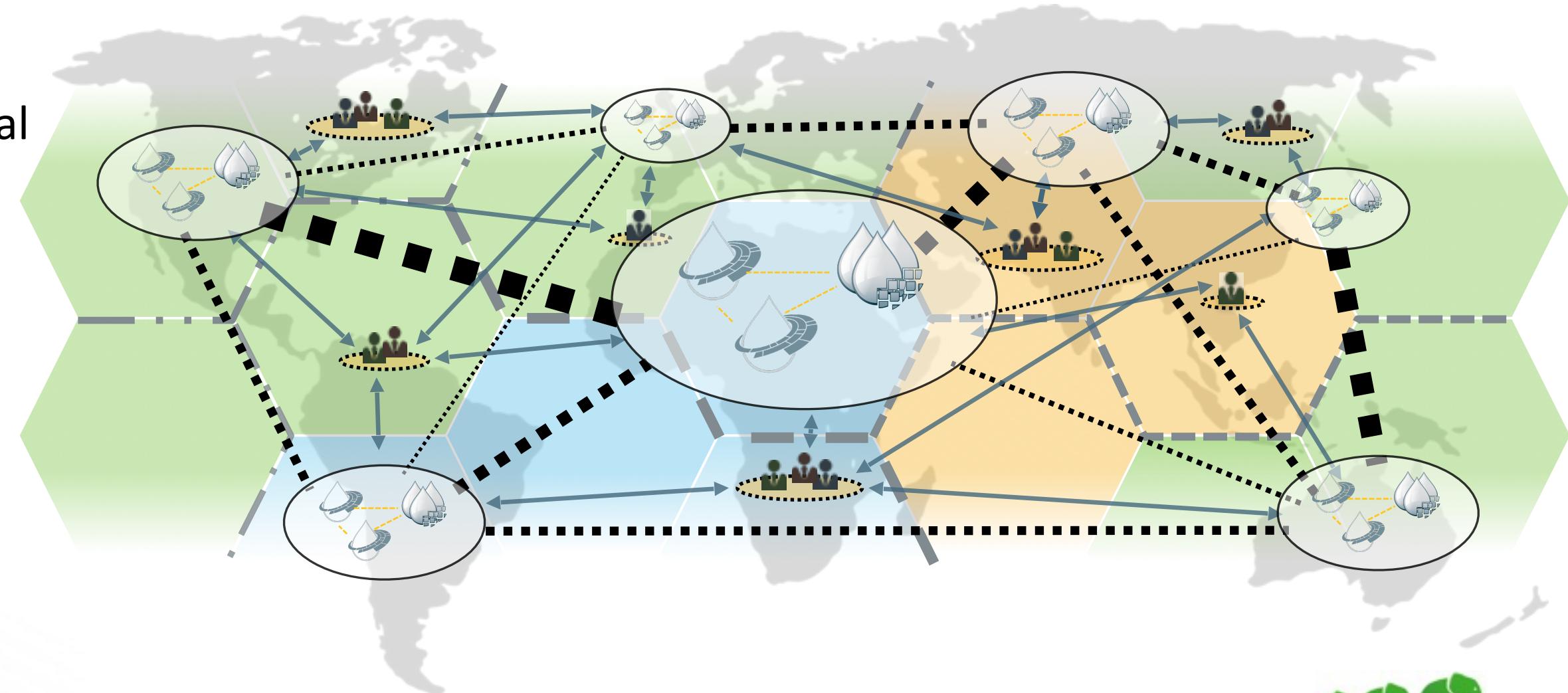
# How Does MiNiFi Interact With NiFi?

- ◆ NiFi
  - ◆ Design flows
  - ◆ Aggregate data from many sources
  - ◆ Perform routing/analysis/SEP
- ◆ MiNiFi
  - ◆ Receive flows
  - ◆ Collect data
  - ◆ Send for processing



# Let's Add Dimensionality

- ◆ We've been imagining EDGE to CORE as a bi-directional linear system
- ◆ Let's expand that to the real world



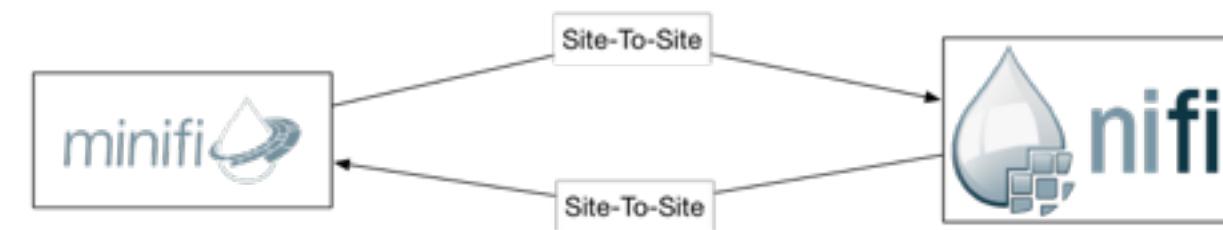
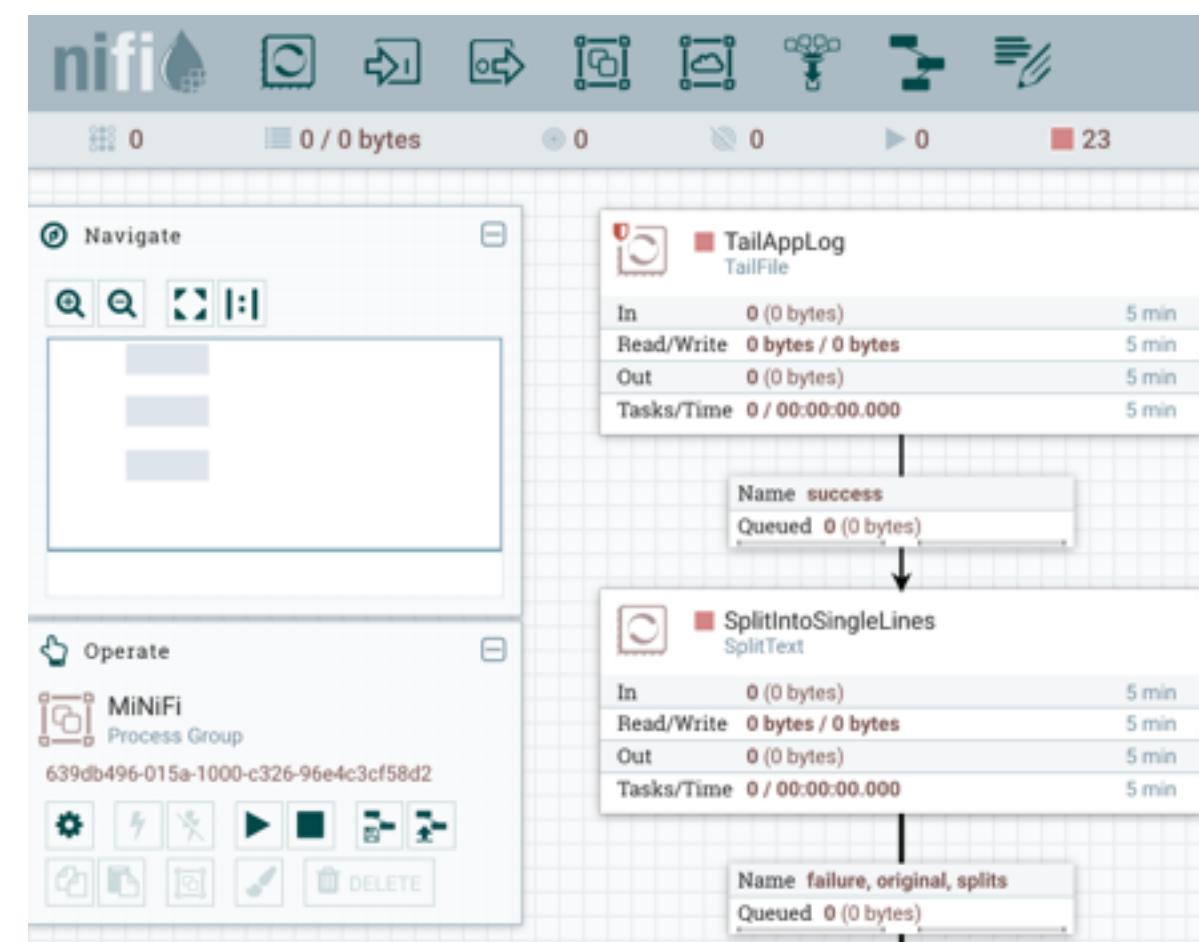
# Flavors of MiNiFi

## ◆ MiNiFi Java (v0.2.0)

- ◆ Modified version of NiFi
  - ◆ No UI
  - ◆ YAML configuration
- ◆ Reduced processor count
  - ◆ 110 by default, more available with additional NARs

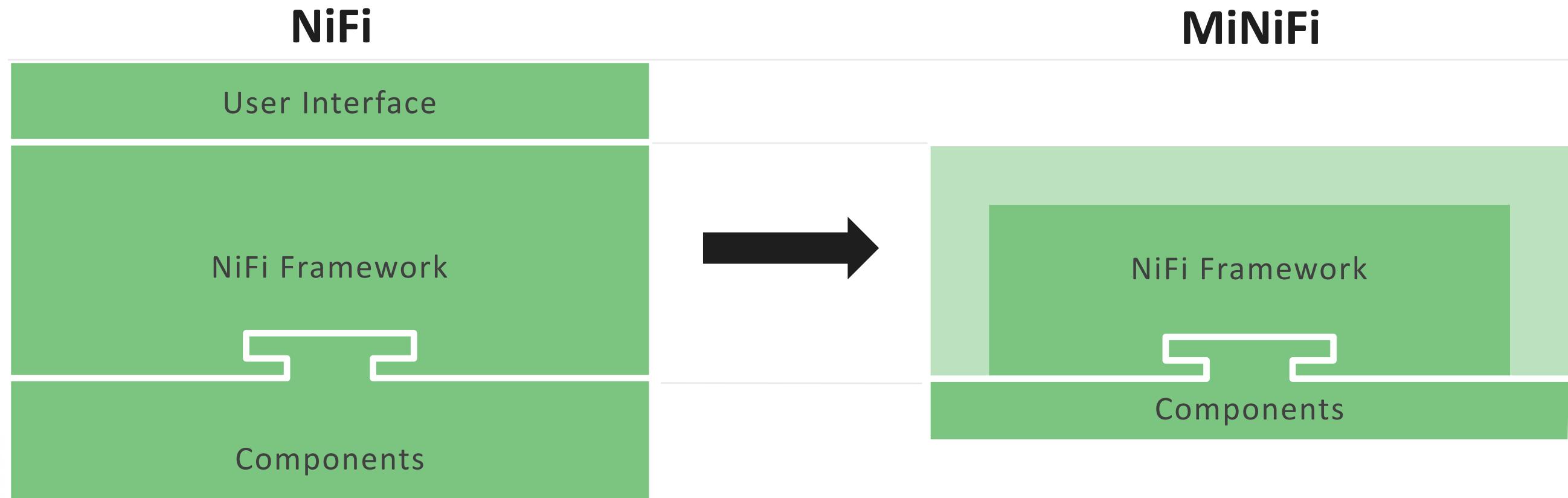
## ◆ MiNiFi C++ (v0.2.0)

- ◆ Written from scratch
- ◆ 10 processors by default
- ◆ Bi-directional site-to-site & provenance data

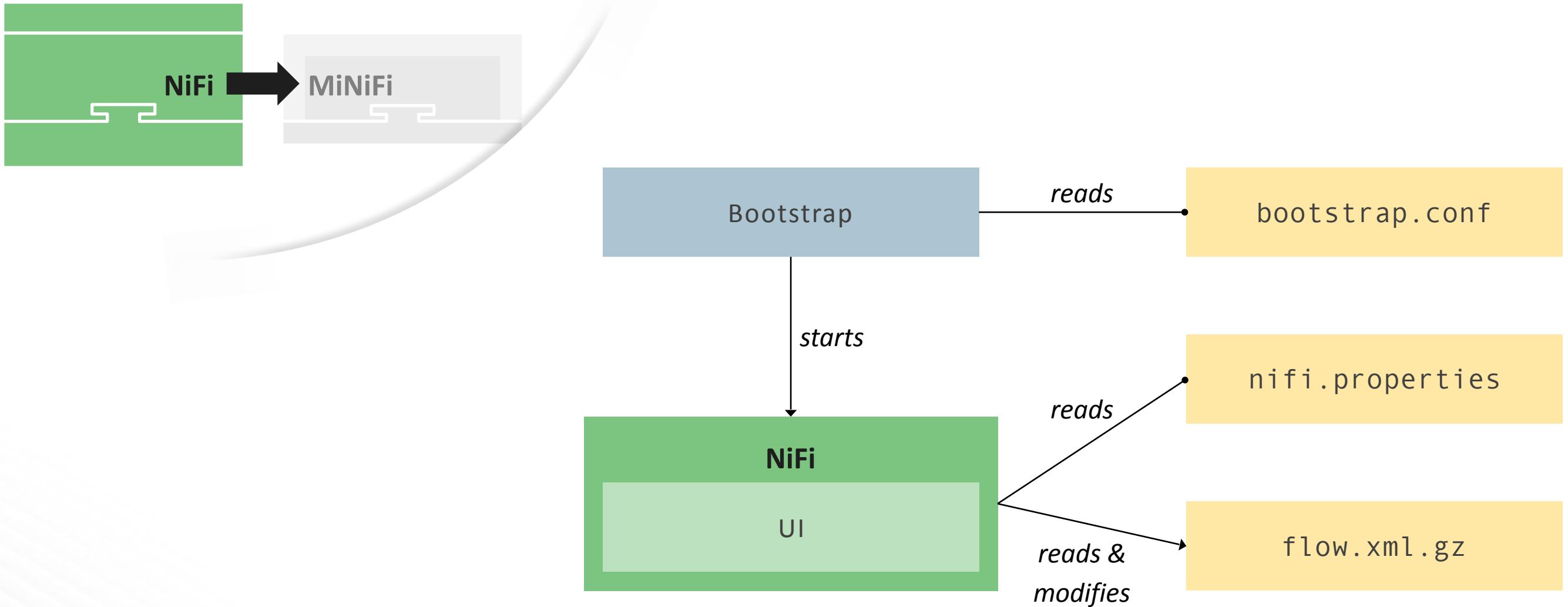


```
Security Properties:  
keystore: /tmp/ssl/localhost-ks.jks  
keystore type: JKS  
keystore password: localtest  
key password: localtest  
truststore: /tmp/ssl/localhost-ts.jks  
truststore type: JKS  
truststore password: localtest  
ssl protocol: TLS  
Sensitive Props:  
key:  
algorithm: PBWEWITHMD5AND256BITAES-CBC-OPENSSL  
provider: BC  
  
Processors:  
- name: TailAppLog  
  class: org.apache.nifi.processors.standard.TailFile  
  max concurrent tasks: 1  
  scheduling strategy: TIMER_DRIVEN  
  scheduling period: 10 sec  
  penalization period: 30 sec  
  yield period: 1 sec  
  run duration nanos: 0  
  auto-terminated relationships list:  
    File to Tail: logs/minifi-app.log  
    Rolling Filename Pattern: minifi-app*  
    Initial Start Position: Beginning of File  
- name: SplitIntoSingleLines  
  class: org.apache.nifi.processors.standard.SplitText  
  max concurrent tasks: 1  
  scheduling strategy: TIMER_DRIVEN  
  scheduling period: 0 sec  
  penalization period: 30 sec  
  yield period: 1 sec  
  run duration nanos: 0  
  auto-terminated relationships list:  
    - failure  
    - original  
Properties:  
  Line Split Count: 1  
  Header Line Count: 0  
  Remove Trailing Newlines: true
```

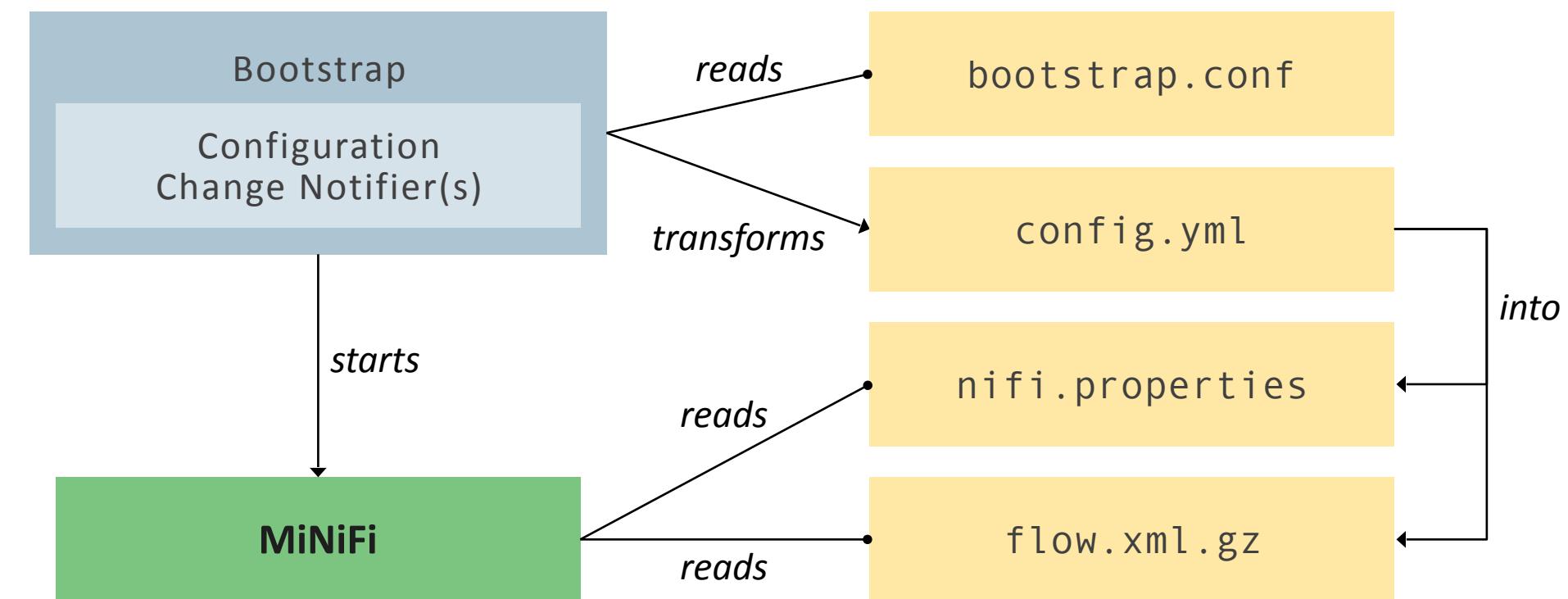
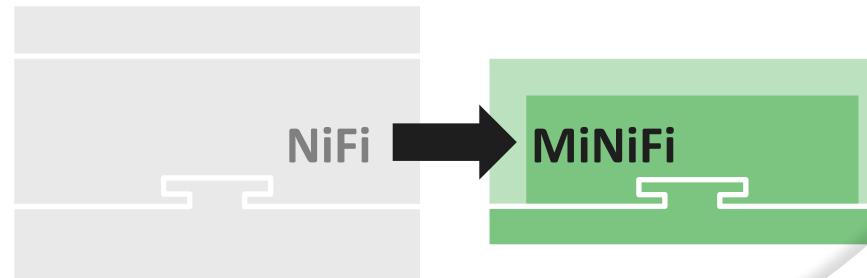
# NiFi vs MiNiFi Java Processes



# NiFi Java Processes



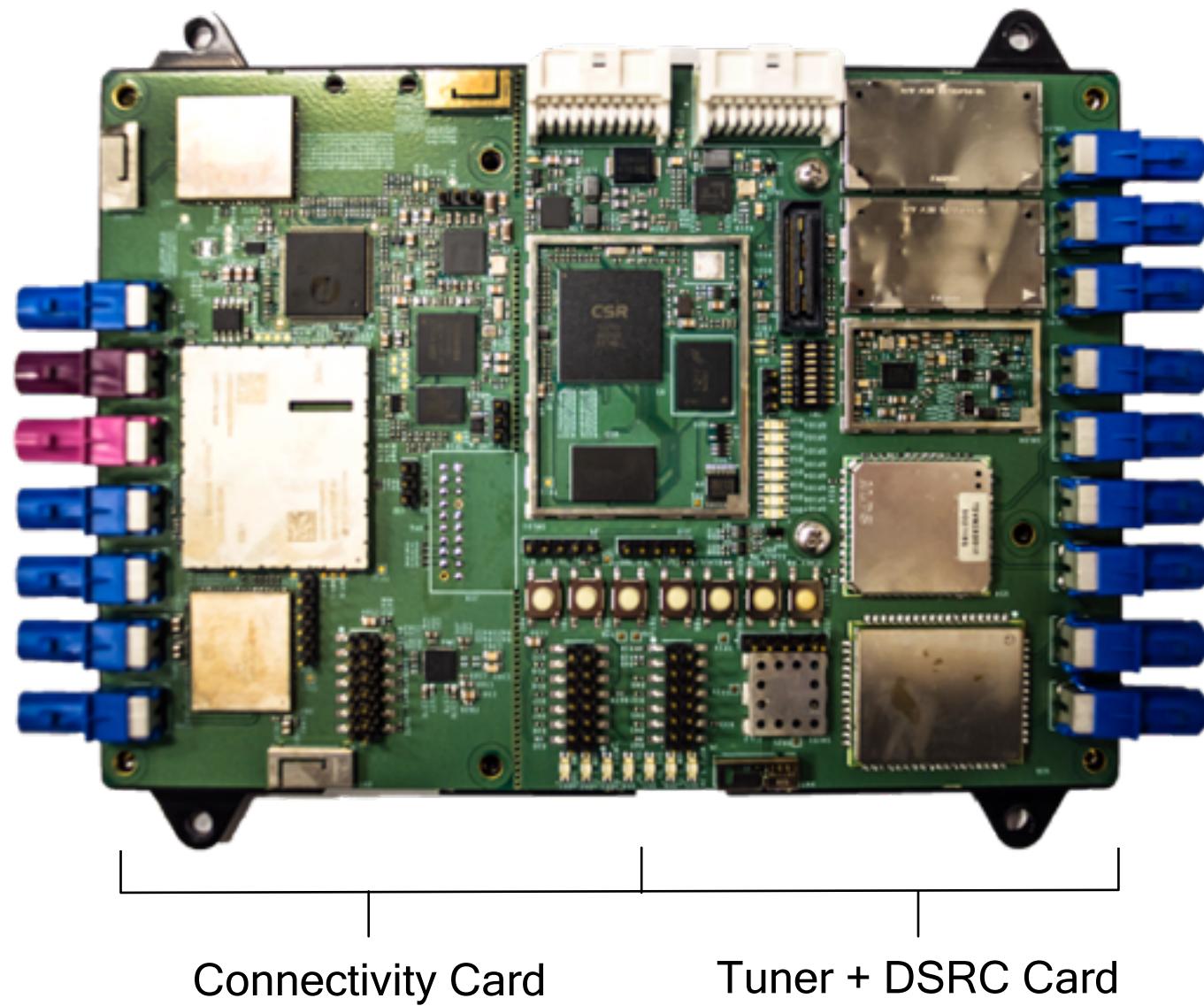
# MiNiFi Java Processes



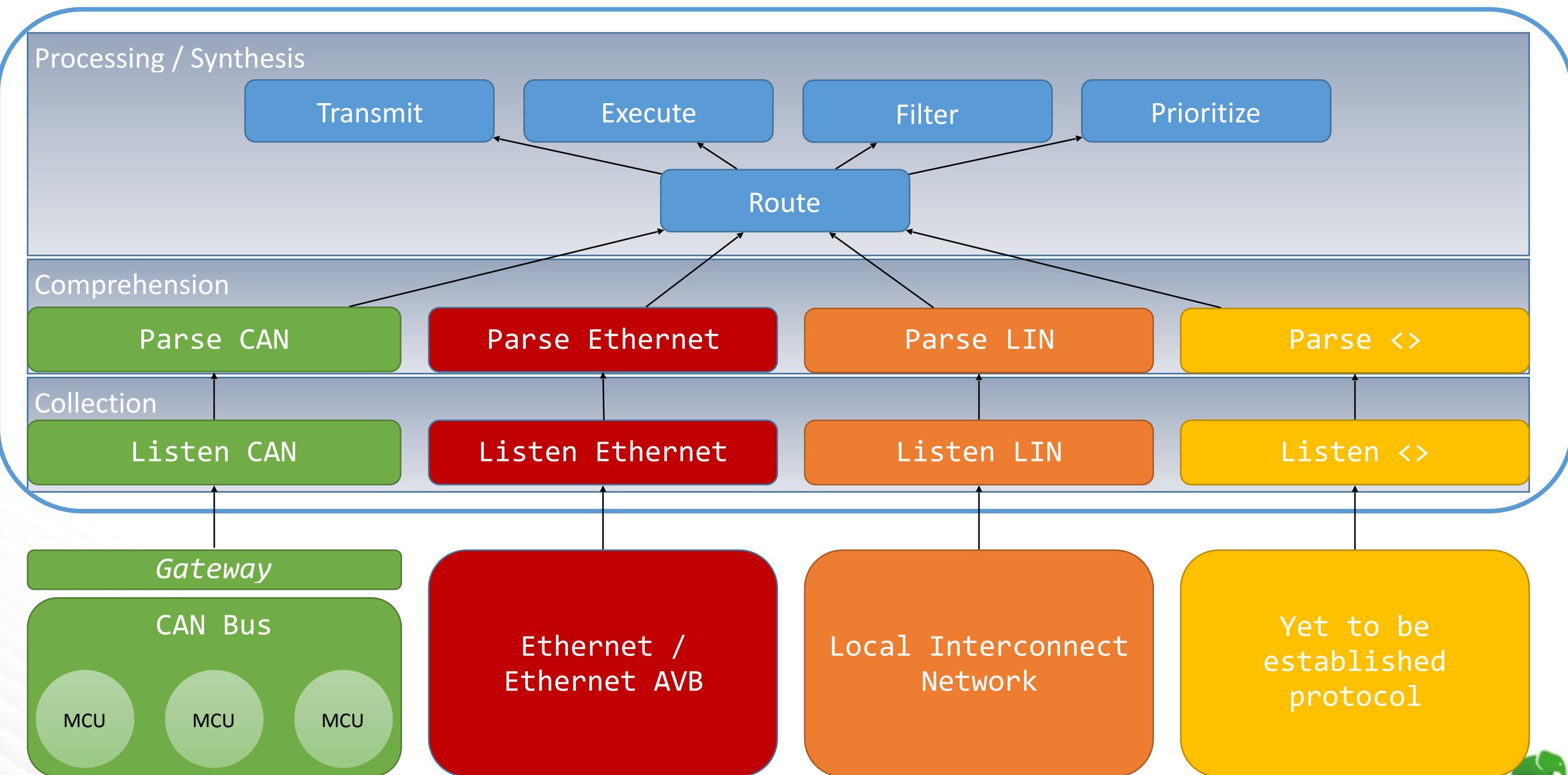
# What does MiNiFi provide?

- ◆ Data tagging/provenance
- ◆ Governance from edge (geopolitical restrictions)
- ◆ Security (encryption, certificate-based authentication)
- ◆ Low latency (immediate reactions & decision-making)

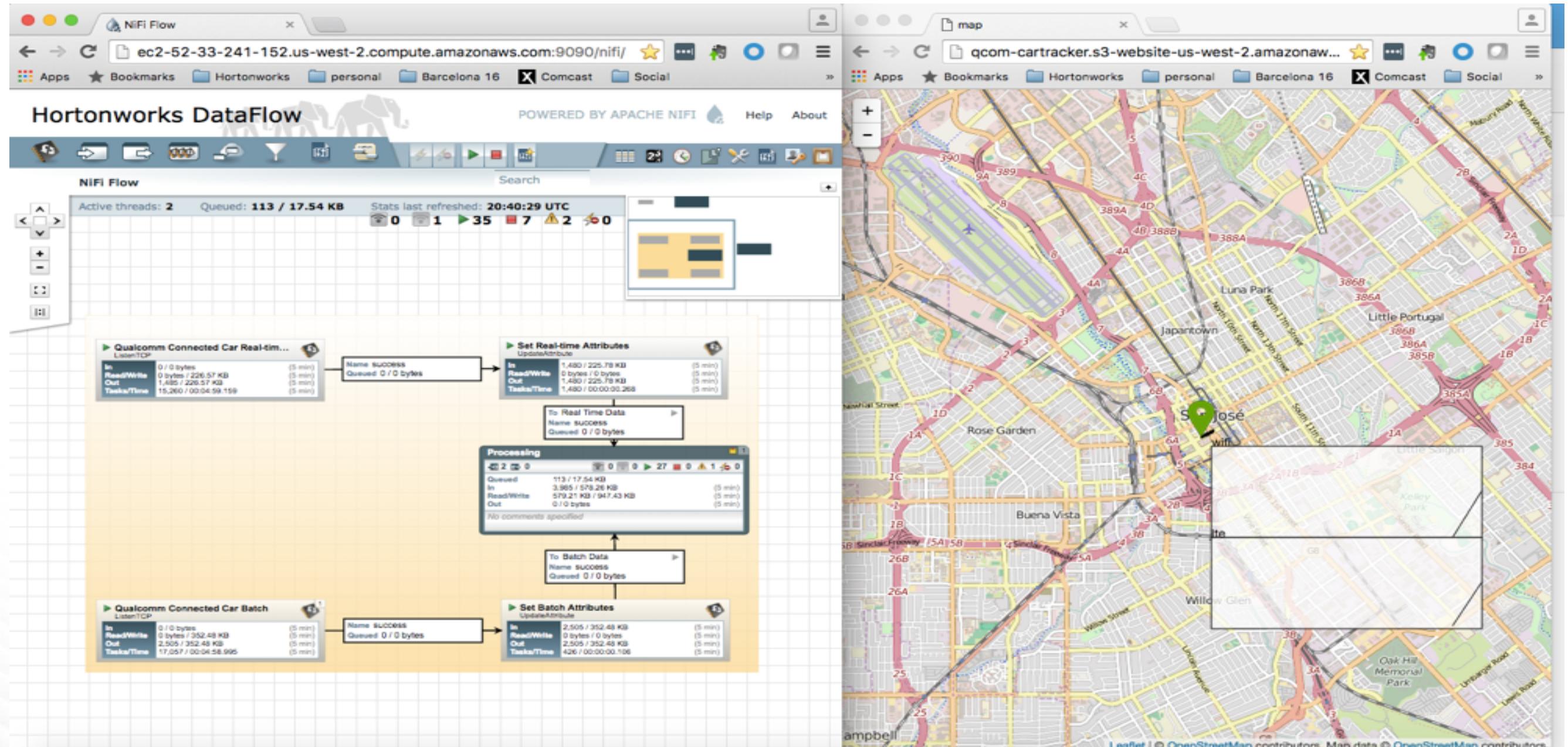
## Connected Car Reference Platform Box



# MiNiFi on a Connected Car



# MiNiFi on a Connected Car



# MiNiFi Exfil

- ◆ Site-to-Site
  - ◆ NiFi protocol
  - ◆ Two implementations
    - ◆ Raw socket
    - ◆ HTTP(S) (*Java only*)
  - ◆ Secured with mutual authentication TLS
- ◆ HTTP(S), (S)FTP, JMS, Syslog, File, Email, Process (*Java only*)

# Advanced Topics

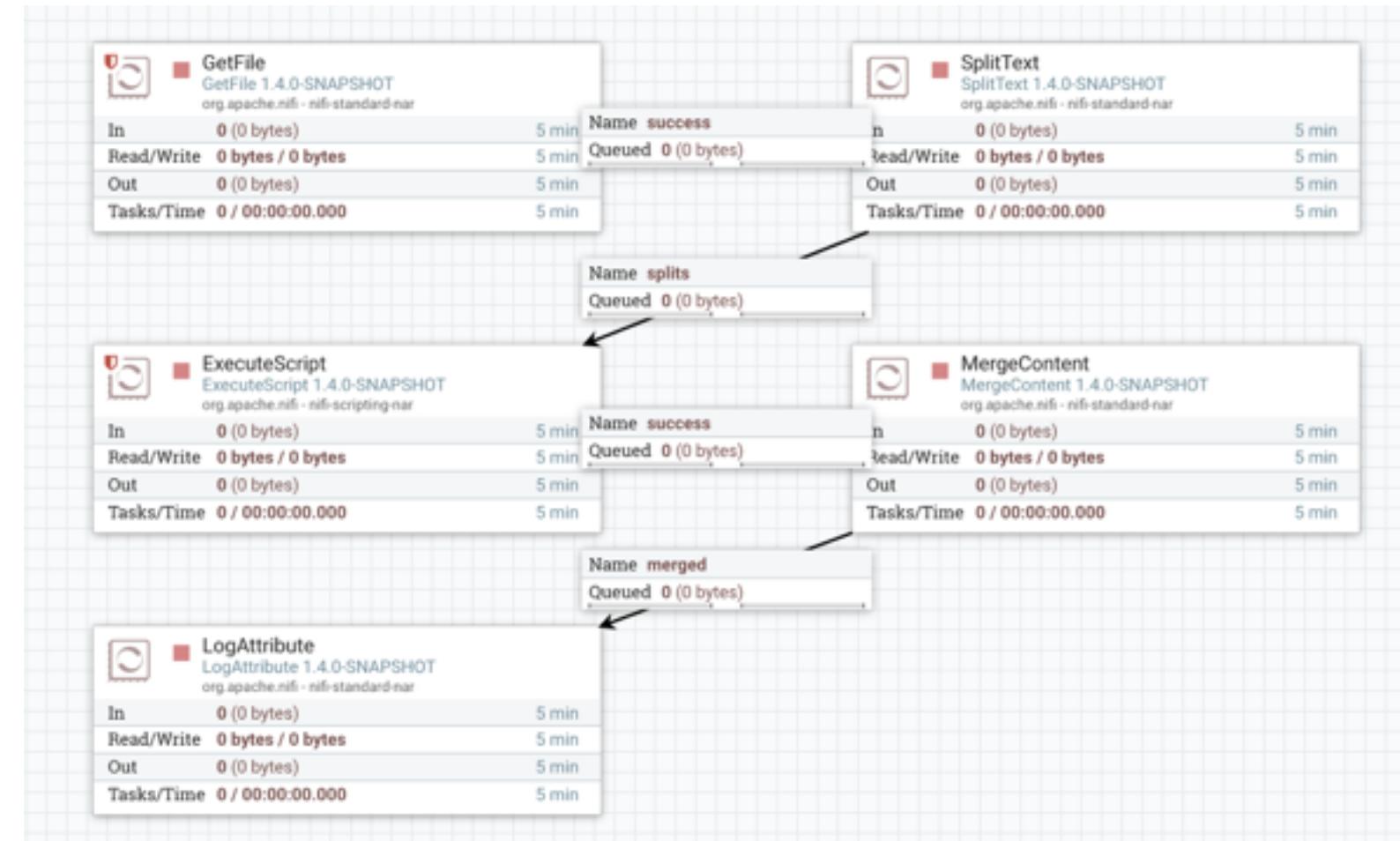
- ◆ New features in Apache NiFi 1.2.0 & 1.3.0
- ◆ New features in Apache MiNiFi Java 0.2.0 & C++ 0.2.0
- ◆ New subproject Apache NiFi Registry

## New in 1.2.0/1.3.0

- ◆ Record Parsing
- ◆ Encrypted Provenance Repository

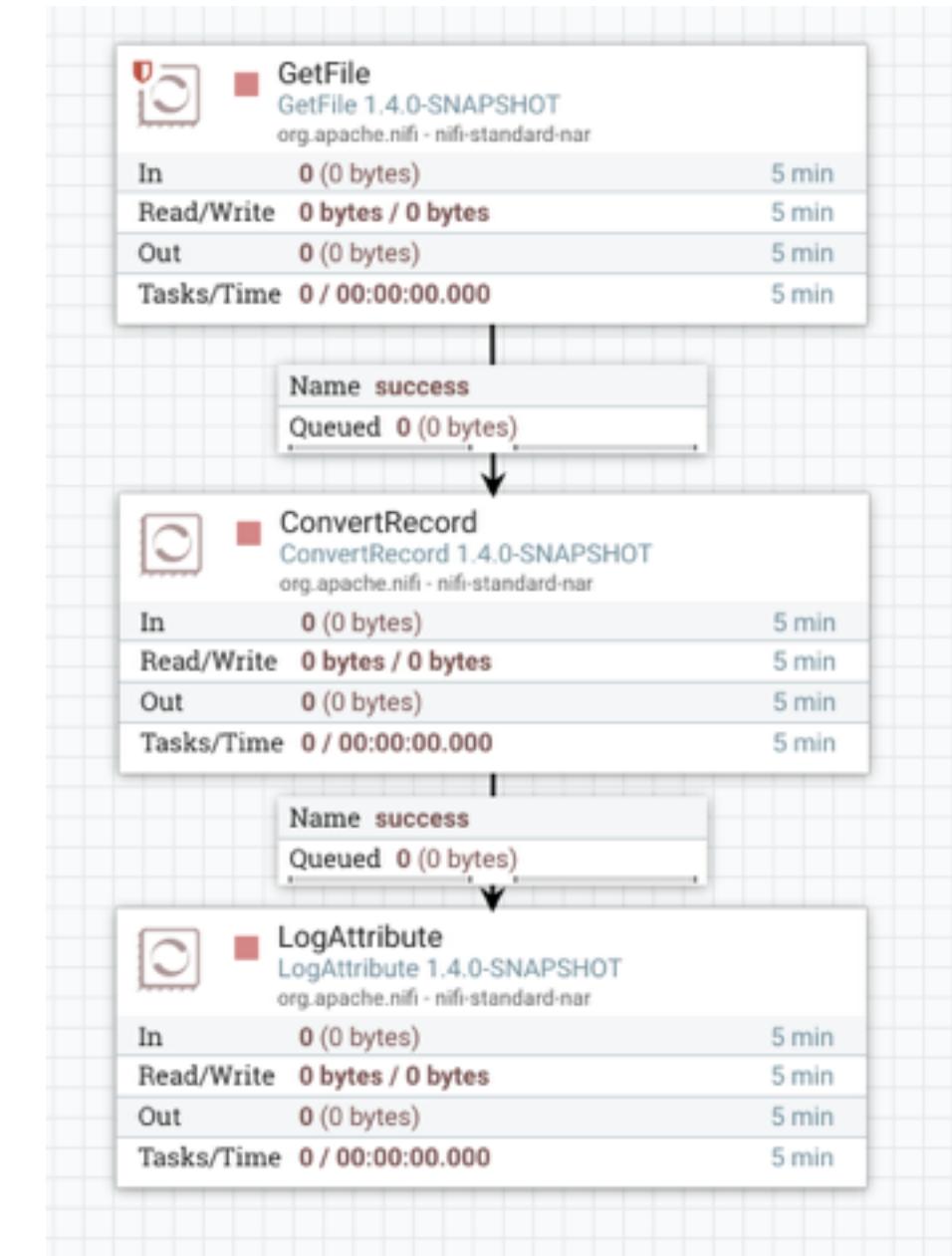
# Record Parsing

- Previously, data had to be divided into individual flowfiles to perform work
- CSV output with 50k lines would need to be split, operated on, remerged
- $1 + 50k + 50k + 1$  flowfiles = 100k flowfiles



# Record Parsing

- Now flowfile content can contain many “record” elements
- Read and write with *\*Reader* and *\*Writer* Controller Services
- Perform lookups, routing, conversion, SQL queries, validation, and more...
- 1 + 1 flowfiles = 2 flowfiles



# Encrypted Provenance Repository

- Every provenance event record is encrypted with AES GCM before being persisted to disk
  - Decrypted on deserialization for retrieval/query
  - Random access via offset seek
  - Handles key migration & rotation

Showing the events that match the specified query. <a href="#">Clear search</a>						
Filter	by component name					
	Date/Time	Type	Flowfile UUID	Size	Component Name	Component Type
1	06/05/2017 20:17:4...	CONTENT_MODIFIED	d602bdff-9d74-4c2e...	77 bytes	ConvertRecord	ConvertRecord
2	06/05/2017 20:17:4...	ROUTE	d602bdff-9d74-4c2e...	46 bytes	LookupRecord	LookupRecord
3	06/05/2017 20:17:4...	FORK	f5497cf-f1e41-4cb3...	40 bytes	LookupRecord	LookupRecord



# MiNiFi Java 0.2.0

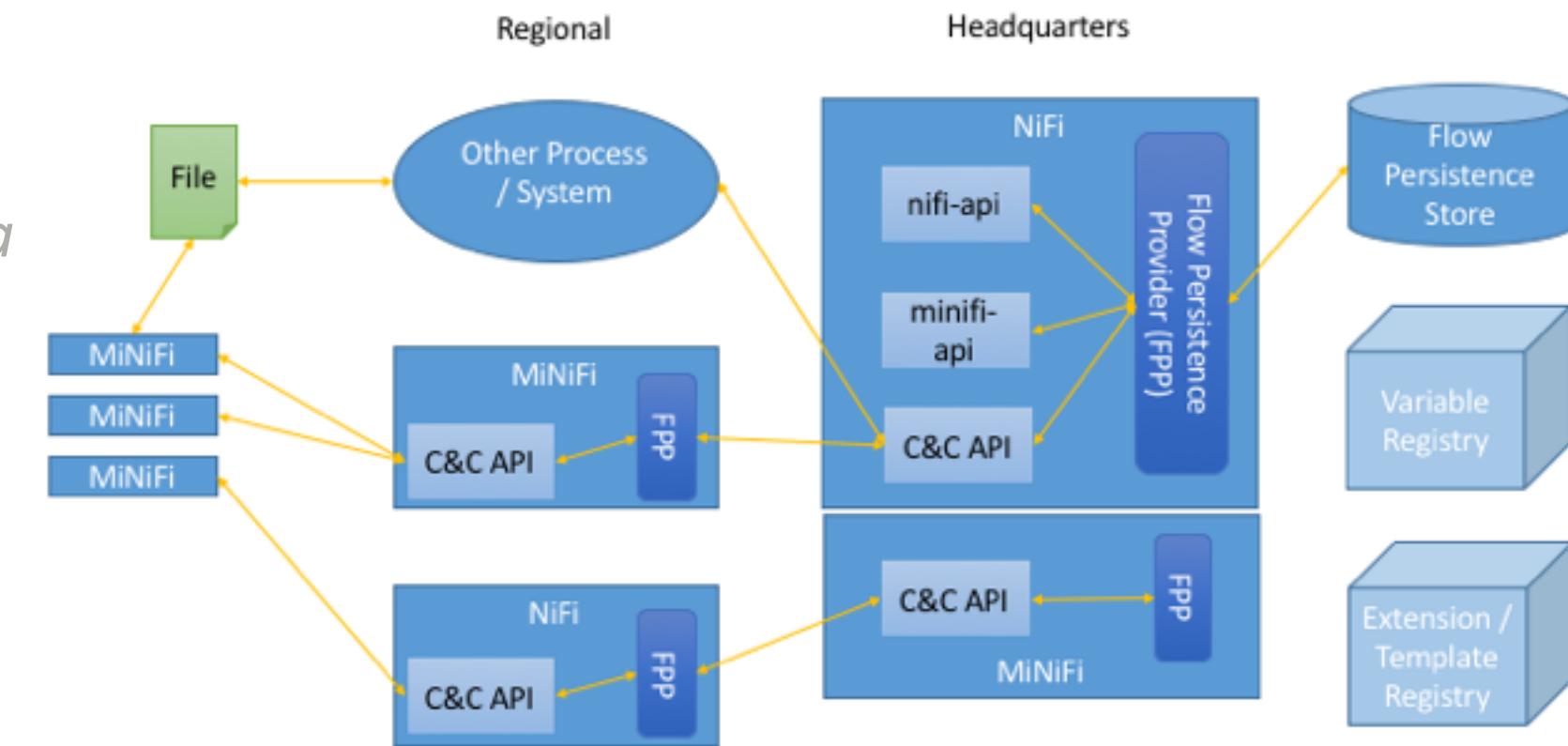
- ◆ Upgrading of core component dependencies to NiFi 1.2.0
- ◆ Initial command and control server capabilities
- ◆ Increased support for NiFi features in configuration YAML inclusive of:
- ◆ Support for HTTP Site to Site Proxy Properties
- ◆ Controller Services
- ◆ Binding site to site to a specific network interface

# MiNiFi C++ 0.2.0

- ◆ Incorporation of Catch testing framework and Google linting for code quality and enhanced test coverage
- ◆ Providing support for reporting tasks and an initial implementation of Site to Site Provenance reporting
- ◆ New Processors inclusive of PutFile, ListenHTTP

# MiNiFi Feature Proposals

- ◆ Flow Versioning
  - ◆ Develop flows for class of MiNiFi instances
- ◆ Command & Control (C2) API (*in Java master*)
  - ◆ FileChangelogestor
  - ◆ RestAPIIngestor
  - ◆ PullHTTPIngestor



# Apache NiFi Registry

- ◆ “...complementary application that provides a central location for storage and management of shared resources across one or more instances of NiFi and/or MiNiFi.”



# Apache NiFi Registry - Flow Registry

- ◆ Flow registry stores & manages versioned flow definitions
- ◆ Integrated with NiFi to allow save/retrieve/upgrade operations from canvas
- ◆ Admin of users, groups, and policies

# Apache NiFi Registry - Flow Registry

The screenshot shows the Apache NiFi Registry - Flow Registry interface. At the top, there are summary statistics: 70 Assets, 134 Extensions, 54 Flows, 43 Compliant, 5 Fleet, 23 Production Ready, and 47 Secure. Below this is a search bar and a sorting dropdown set to "Sort by: Last Update".

Three main items are listed:

- configuration-6eyq23u** Asset: Versions 2, Flows 2, Extensions 27, Assets 0, Deployments 0.
- nifi-email-bundle-lorem-ipsum-d...** Extension: Versions 5, Flows 45, Extensions 0, Assets 2, Deployments 0.
- Fraud Detection Flow** Data Flow: Versions 11, Flows 134, Extensions 134, Assets 14, Deployments 0.

A detailed view of the Fraud Detection Flow is shown on the right. It includes a description: "Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam." The flow diagram shows three main components: "Gather Network Traffic" (ExecuteProcess), "TCP Traffic" (TCPListener), and "Photos" (FileReader). Arrows indicate data flow from "Gather Network Traffic" to "TCP Traffic", and from "TCP Traffic" to "Photos". A tooltip for "Gather Network Traffic" provides performance metrics: In 0 (0 bytes), Read/Write 0 bytes / 45.14 MB, Out 93 (0.13 MB), Tasks/Time 420 / 00:06:20.325. The "CHANGE LOG" on the right lists recent changes:

- 11 3 days ago by Danyell Roten Production Ready
- 10 2 months ago by Marcelle Wisniesek
- 9 2 months ago by Danyell Roten
- 8 3 months ago by Marcelle Wisniesek
- 6 4 months ago by Marcelle Wisniesek
- 5 4 months ago by Marcelle Wisniesek
- 4 5 months ago

# Apache NiFi Registry - Flow Registry

The screenshot shows the Apache NiFi Registry - Administration interface. The top navigation bar includes links for NiFi Registry / Administration, General, Users, and Workflow. The Workflow tab is selected, indicated by an underline.

**Buckets (8)**

Create new bucket

Bucket Name	Action
Acme Fraud Team	<input type="button"/> <input type="button"/> <input type="button"/> <input type="button"/>
Acme Partners	<input type="button"/> <input type="button"/> <input type="button"/> <input type="button"/>
CS-prod	<input type="button"/> <input type="button"/> <input type="button"/> <input type="button"/>
CyberSec	<input type="button"/> <input type="button"/> <input type="button"/> <input type="button"/>
Flow dev	<input type="button"/> <input type="button"/> <input type="button"/> <input type="button"/>
FLP-org	<input type="button"/> <input type="button"/> <input type="button"/> <input type="button"/>
IOAT fraud team	<input type="button"/> <input type="button"/> <input type="button"/> <input type="button"/>
Security	<input type="button"/> <input type="button"/> <input type="button"/> <input type="button"/>

**Certifications (6)**

Create new certification

Label Name	Usage	Badge Design
Compliant	<input type="button"/> ON <input type="button"/> OFF	
Fleet	<input type="button"/> ON <input type="button"/> OFF	
Obsolete	<input type="button"/> ON <input type="button"/> OFF	
POC	<input type="button"/> ON <input type="button"/> OFF	
Production Ready	<input type="button"/> ON <input type="button"/> OFF	
Secure	<input type="button"/> ON <input type="button"/> OFF	

### Save Data Flow Version

Name  (2)

Location

Certifications

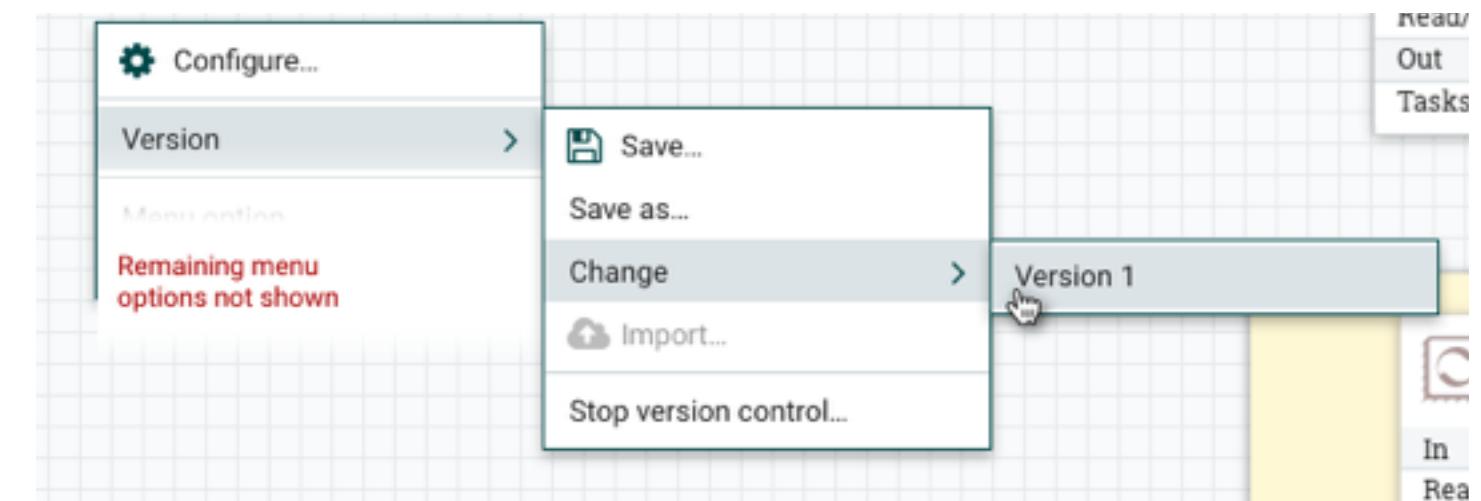
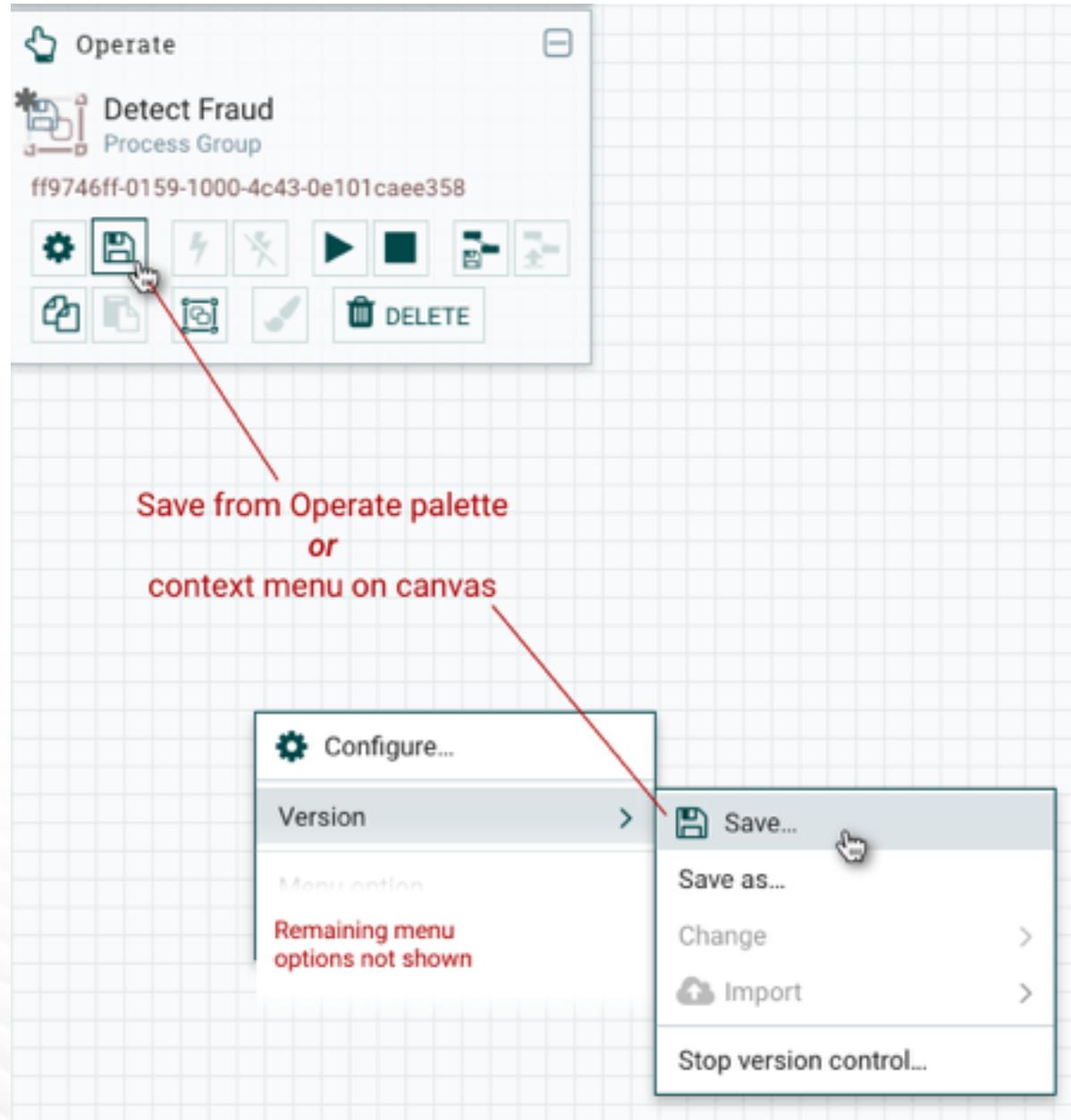
Compliant	Obsolete	POC	Production Ready
Secure			

Description

Change Comments

NOTICE: unsaved changes

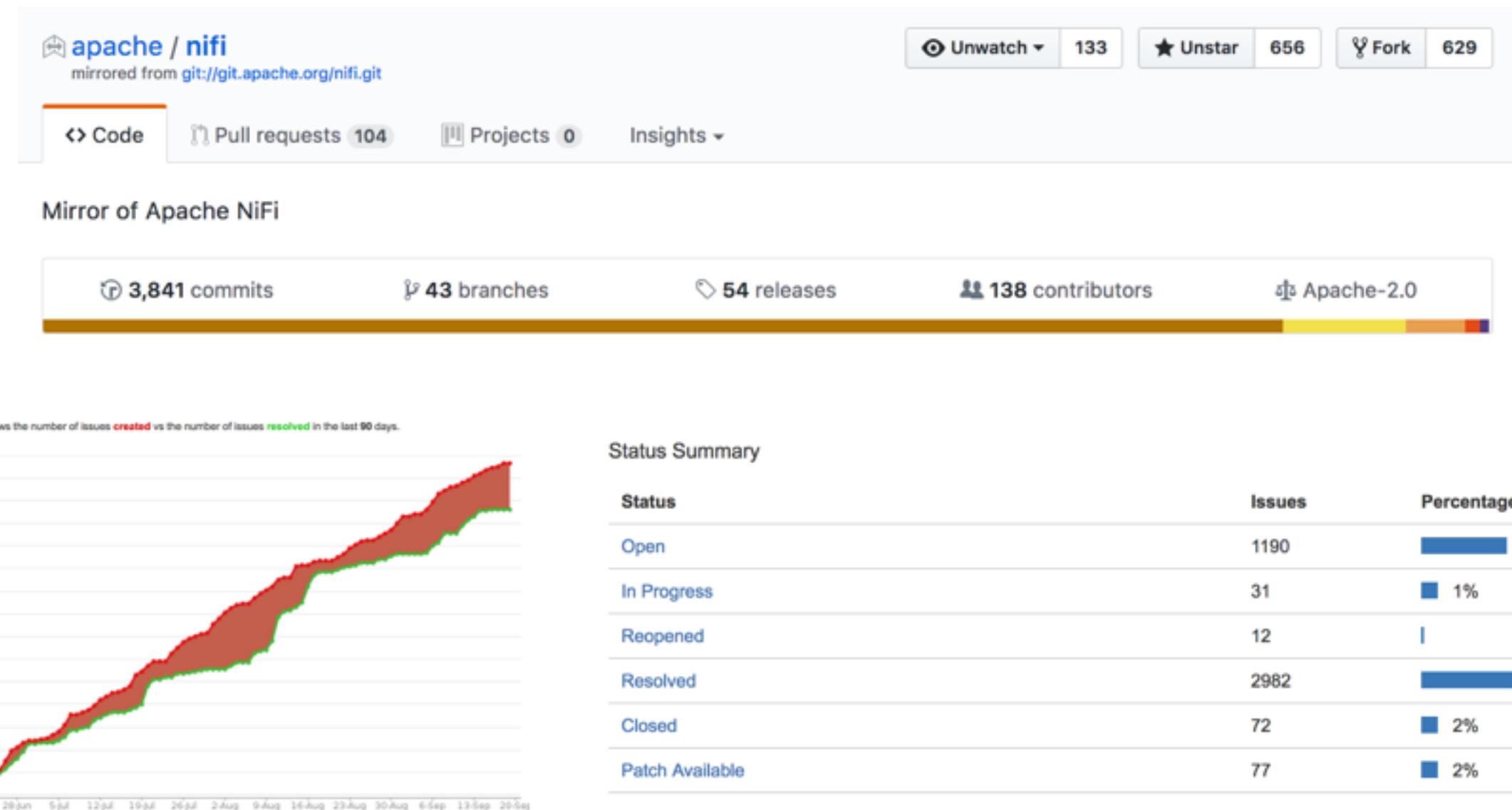
# Apache NiFi Registry - Flow Registry



# Why NiFi & MiNiFi?

- ◆ Moving data is multifaceted in its challenges and these are present in different contexts at varying scopes
  - Inter vs intra, domestically, internationally
- ◆ Provide common tooling and extensions that are needed but be flexible for extension
  - Leverage existing libraries and expansive Java ecosystem for functionality
  - Allow organizations to integrate with their existing infrastructure
- ◆ Empower folks managing your infrastructure to make changes and reason about issues that are occurring
  - Data Provenance to show context and data's journey
  - User Interface/Experience a key component

# Healthy Community



# Learn more and join us

**Apache NiFi site**

<https://nifi.apache.org>

**Subproject MiNiFi site**

<https://nifi.apache.org/minifi/>

**Subscribe to and collaborate at**

[dev@nifi.apache.org](mailto:dev@nifi.apache.org)

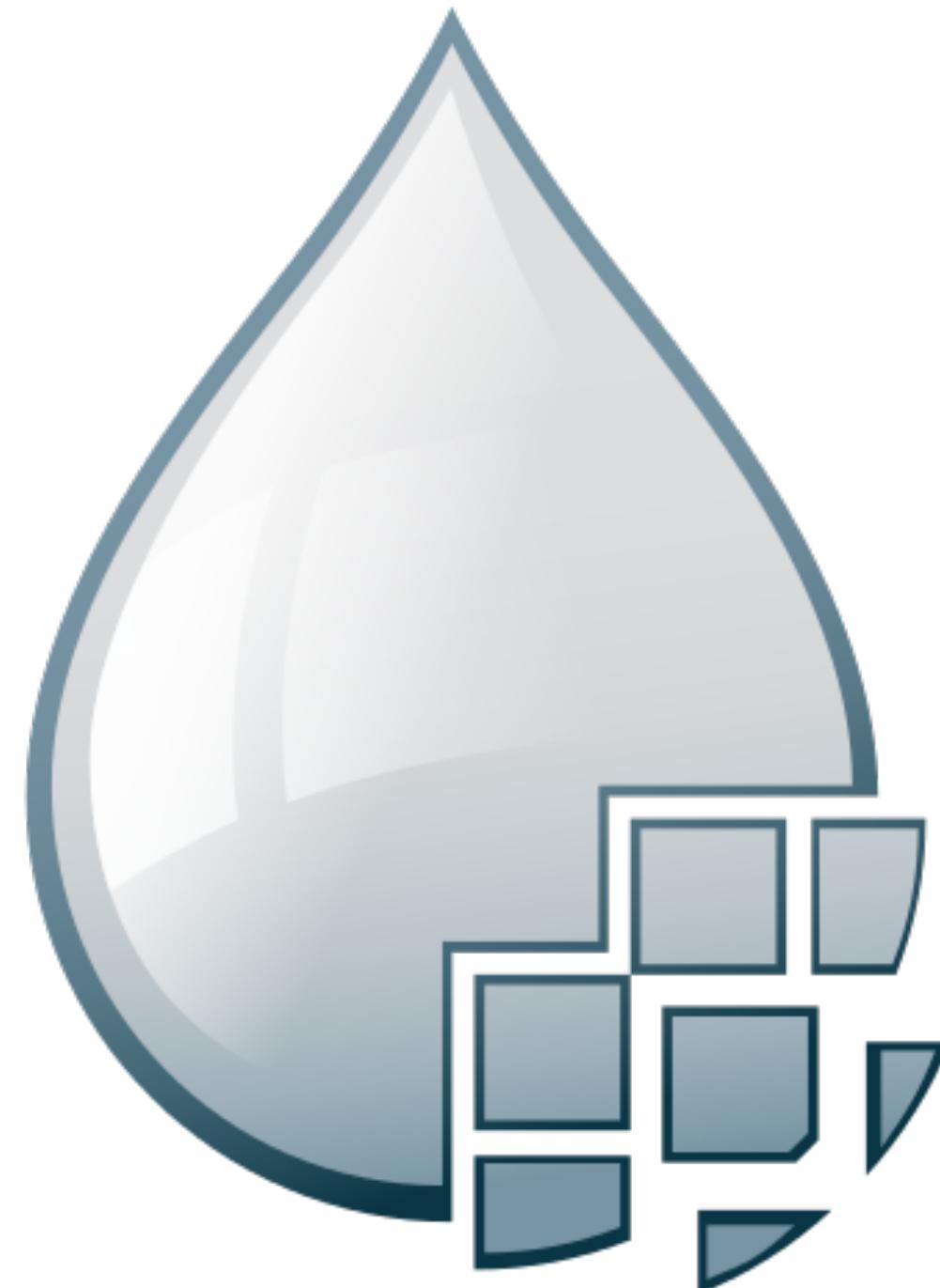
[users@nifi.apache.org](mailto:users@nifi.apache.org)

**Submit Ideas or Issues**

<https://issues.apache.org/jira/browse/NIFI>

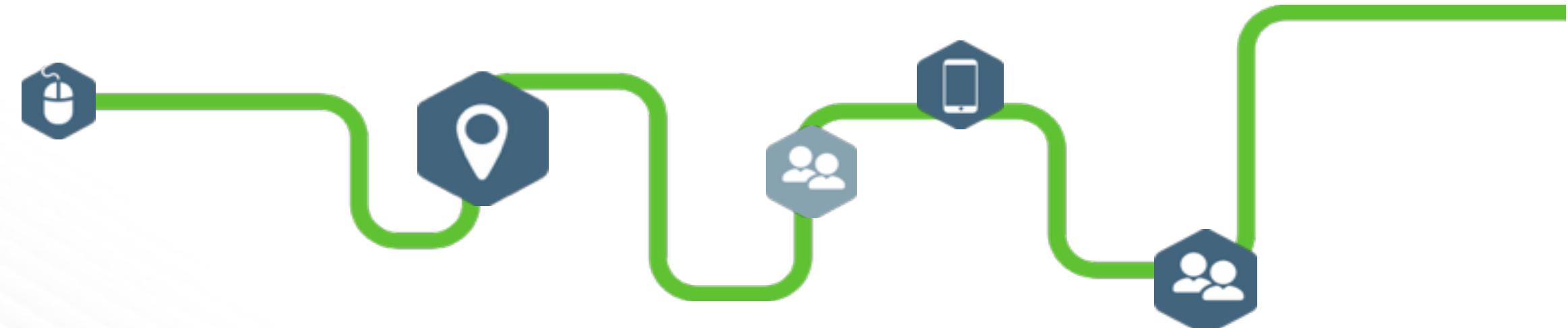
**Follow us on Twitter**

[@apachennifi](https://twitter.com/apachennifi)



# Learn and share at Birds of a Feather IOT, STREAMING & DATA FLOW

Thursday September 21  
6:00 pm, C4.6



# Thank You

I'm sticking around for discussions/questions



@yolopey / @apachennifi

[alopresto@apache.org](mailto:alopresto@apache.org)

PGP: 70EC B3E5 98A6 5A3F D3C4 BACE 3C6E F65B **2F7D EF69**