# Hands-on with Apache NiFi and MiNiFi
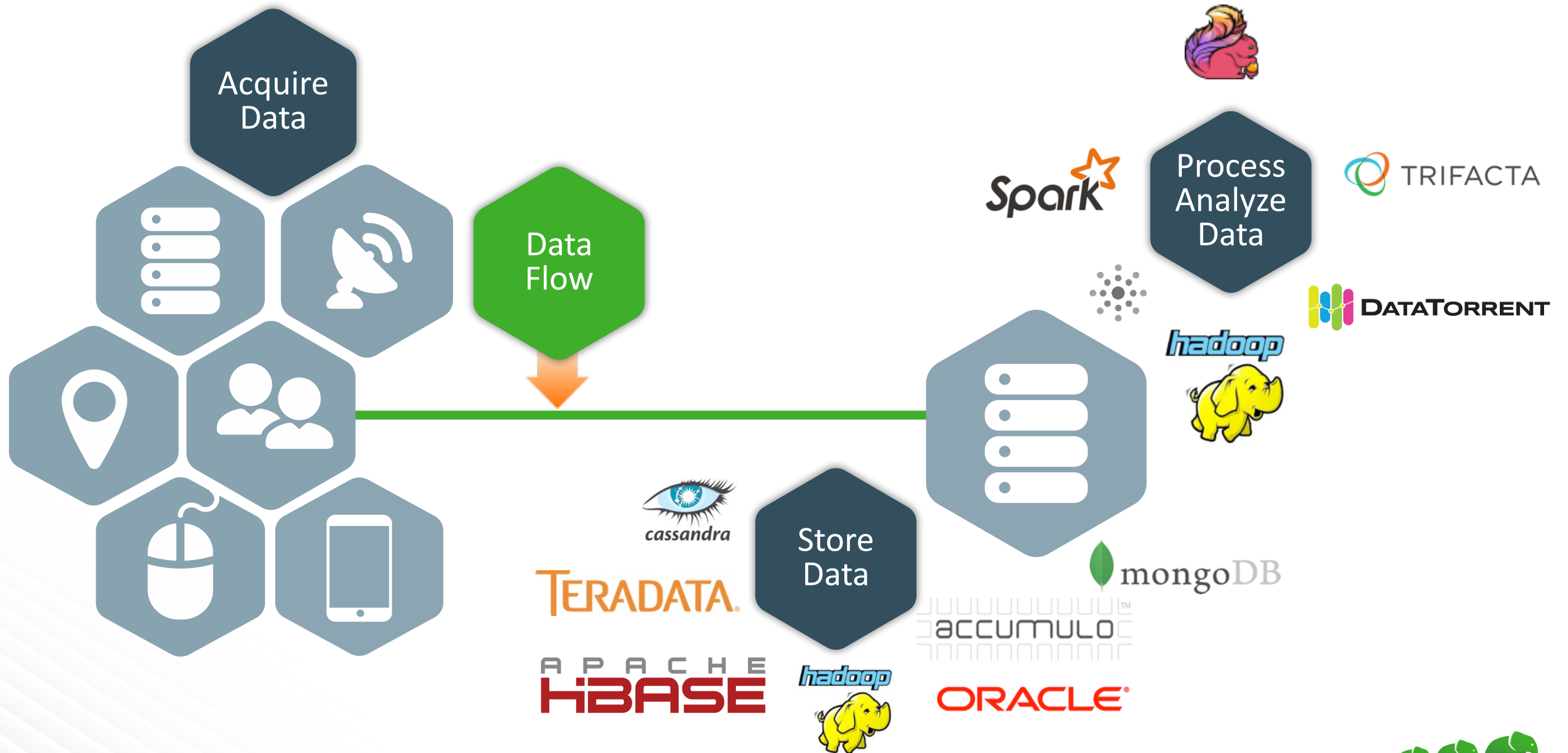
Andrew Psaltis - @itmdata

Berlin Buzzwords 2017

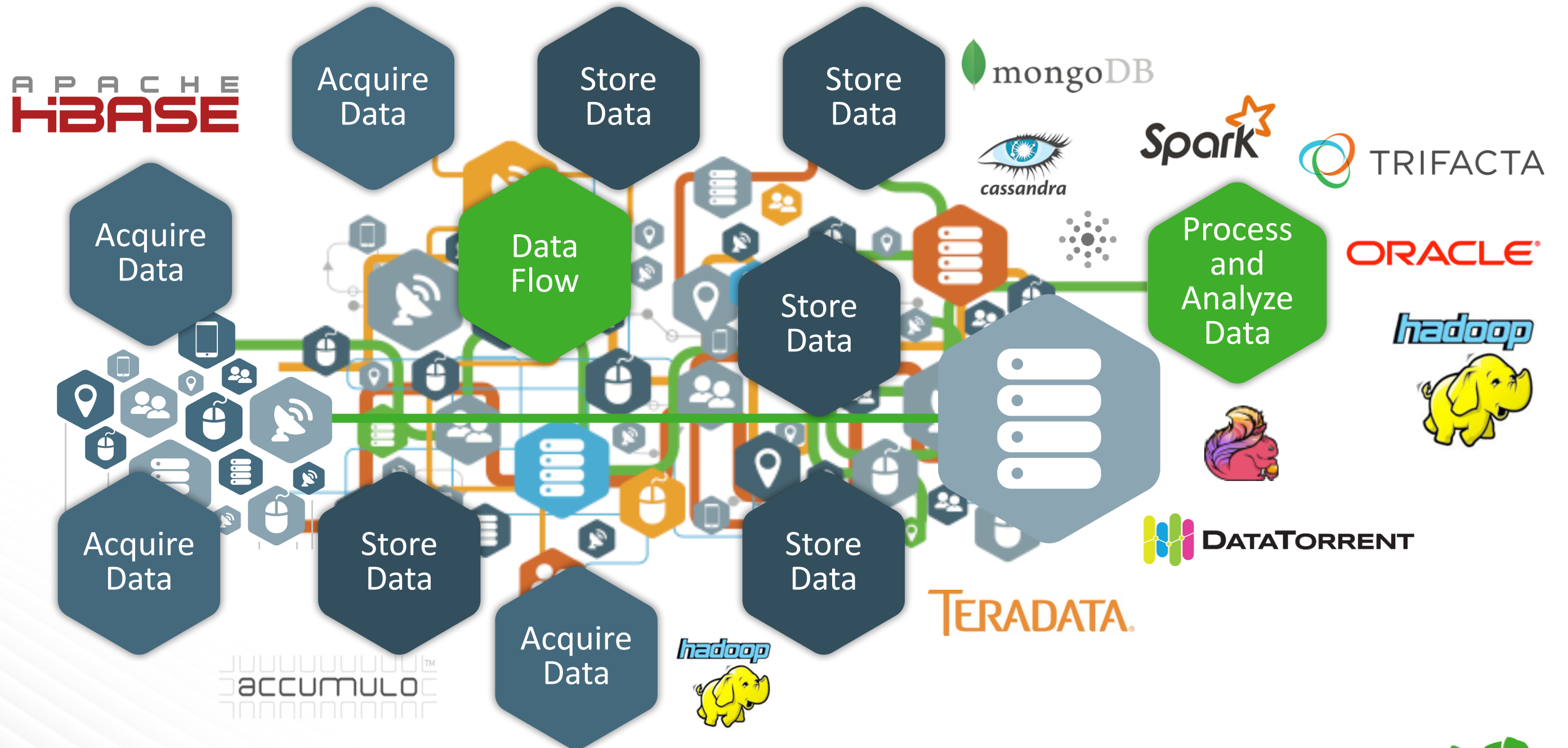HORTONWORKS®
POWERING THE FUTURE OF DATA™
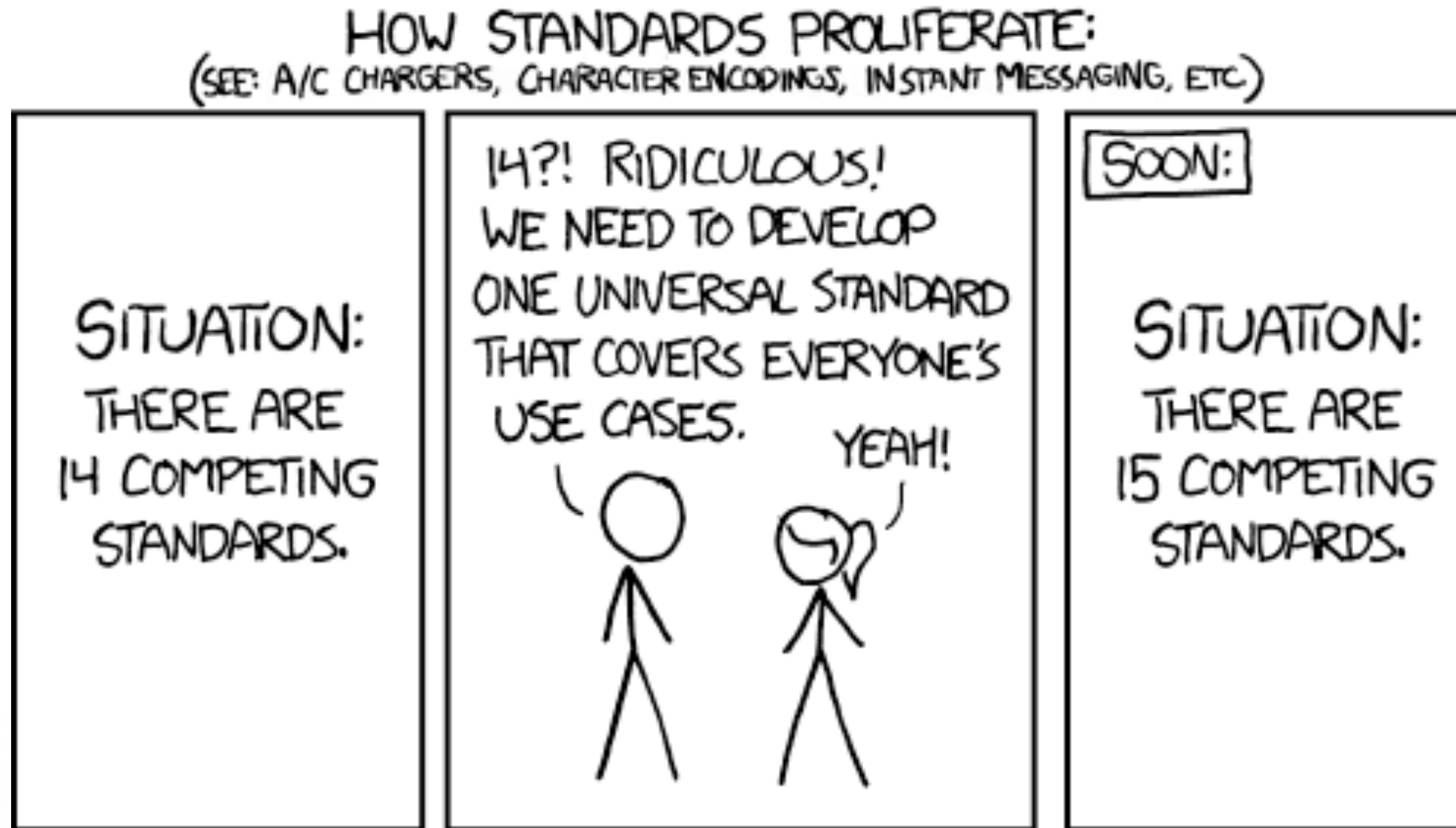
# Dataflow and the associated problems

# Simplistic View of Dataflows: Easy, Definitive

# Realistic View of Dataflows: Complex, Convoluted
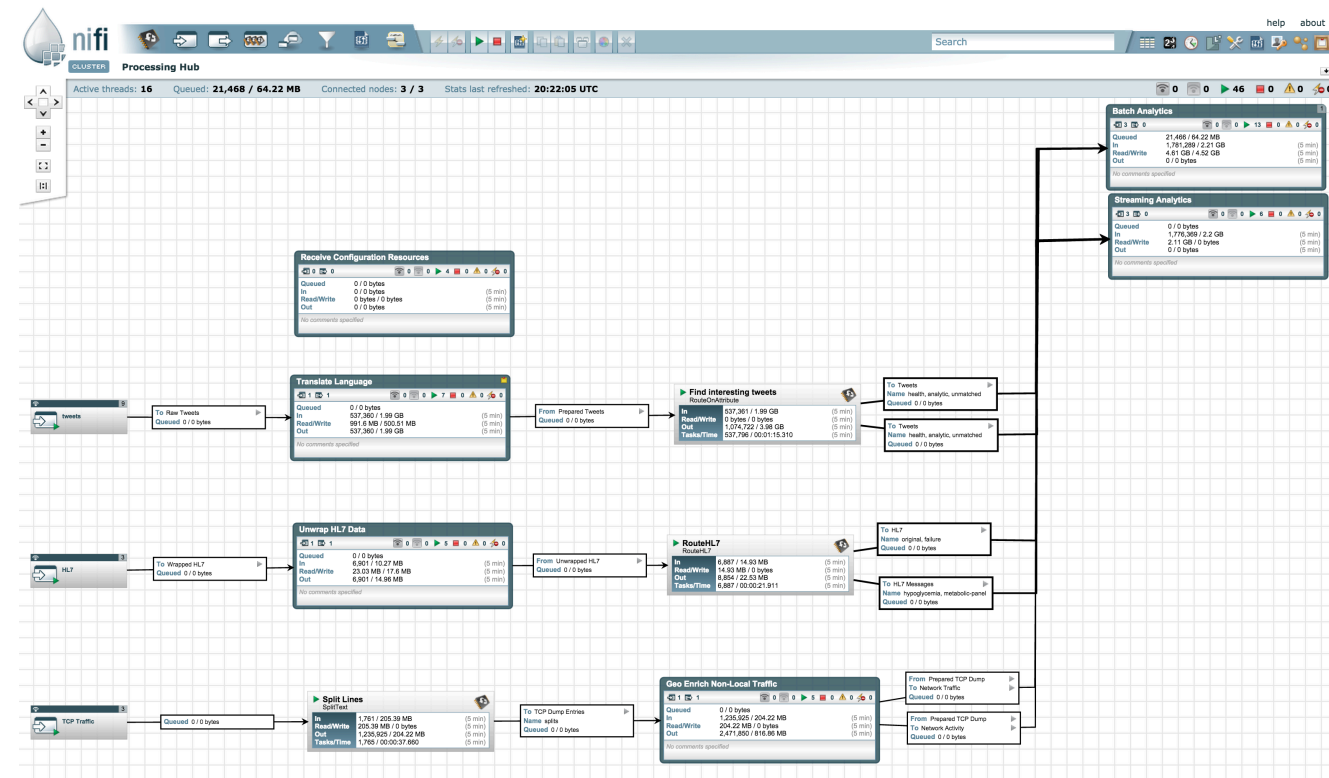
# Moving data *effectively* is hard



*Standards: http://xkcd.com/927/*

# Apache NiFi

# Apache NiFi

## Dataflow

- Web-based User Interface for creating, monitoring, & controlling data flows

- Directed graphs of data routing and transformation

- Highly configurable - modify data flow at runtime, dynamically prioritize data

- Easily extensible through development of custom components

- Data Provenance tracks data through entire system



[1] https://nifi.apache.org/

# NiFi - Terminology

| HTTP Data | FlowFile |
|---|---|

HTTP/1.1 200 OK

Date: Sun, 10 Oct 2010 23:26:07 GMT

Server: Apache/2.2.8 (CentOS) OpenSSL/0.9.8g

Last-Modified: Sun, 26 Sep 2010 22:04:35 GMT

ETag: "45b6-834-49130cc1182c0"

Accept-Ranges: bytes

Content-Length: 13

Connection: close

Content-Type: text/html

**Header**

Standard FlowFile Attributes
Key: 'entryDate'          Value: 'Fri Jun 17 17:15:04 EDT 2016'
Key: 'lineageStartDate'   Value: 'Fri Jun 17 17:15:04 EDT 2016'
Key: 'fileSize'   Value: '23609'
FlowFile Attribute Map Content
Key: 'filename'Value: '15650246997242'
Key: 'path'      Value: './'

Hello world!

**Content**

*Binary Content *

**HORTONWORKS**®

# FlowFiles are like HTTP data

| HTTP Data | FlowFile |
|---|---|

HTTP/1.1 200 OK

Date: Sun, 10 Oct 2010 23:26:07 GMT

Server: Apache/2.2.8 (CentOS) OpenSSL/0.9.8g

Last-Modified: Sun, 26 Sep 2010 22:04:35 GMT

ETag: "45b6-834-49130cc1182c0"

Accept-Ranges: bytes

Content-Length: 13

Connection: close

Content-Type: text/html

**Header**

Standard FlowFile Attributes

Key: 'entryDate'   Value: 'Fri Jun 17 17:15:04 EDT 2016'

Key: 'lineageStartDate'   Value: 'Fri Jun 17 17:15:04 EDT 2016'

Key: 'fileSize'        Value: '23609'

FlowFile Attribute Map Content

Key: 'filename'    Value: '15650246997242'

Key: 'path'          Value: './'

Hello world!

**Content**

*Binary Content ***

HORTONWORKS®

# FlowFiles & Data Agnosticism

- NiFi is data agnostic!

- But, NiFi was designed understanding that users

  can care about specifics and provides tooling

  to interact with specific formats, protocols, etc.

# *Robustness principle*

Be conservative in what you do,
be liberal in what you accept from others



ISO 8601 - http://xkcd.com/1179/

# NiFi - Terminology

- FlowFile
  - Unit of data moving through the system
  - Content + Attributes (key/value pairs)

- Processor
  - Performs the work, can access FlowFiles

- Connection
  - Links between processors
  - Queues that can be dynamically prioritized

- Process Group
  - Set of processors and their connections
  - Receive data via input ports, send data via output ports

HORTONWORKS®

# 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. |

HORTONWORKS®

# Apache NiFi

## Key Features and Principles

- Guaranteed delivery
- Data buffering
  - Backpressure
  - Pressure release
- Prioritized queuing
- Flow specific QoS
  - Latency vs. throughput
  - Loss tolerance
- Data provenance

- Recovery/recording a rolling log of fine-grained history
- Visual command and control
- Flow templates
- Pluggable/multi-role security
- Designed for extension
- Clustering

HORTONWORKS®

# The need for data provenance

**For Operators**

- Traceability, lineage
- Recovery and replay

**For Compliance**

- Audit trail
- Remediation

**For Business / Mission**

- Value sources
- Value IT investment

# Data Provenance (Not just Lineage)



- View attributes and content at given points in time (before and after each processor) !!!
- Records, indexes, and makes events available for display

# Provenance



**SOURCES**

**Origin – attribution**
**Replay – recovery**

**REGIONAL**
**INFRASTRUCTURE**

**CORE**
**INFRASTRUCTURE**

**Evolution of topologies**
**Long retention**

**Types of Lineage**
- Event (runtime)
- Configuration (design time)

CONTEXT/COMMANDS

DATA/EVENTS

EDGE

CORE

- **Constrained**
- **High-latency**
- **Localized context**

- **Hybrid – cloud/on-premises**
- **Low-latency**
- **Global context**

**HORTONWORKS®**

# The need for fine-grained security and compliance

**It's not enough to say you have encrypted communications**

- Enterprise authorization services –entitlements change often

- People and systems with different roles require difference access levels

- Tagged/classified data

**HORTONWORKS**®

# Security



**SOURCES**

**Sign, encrypt, static
(data and control)**

**REGIONAL
INFRASTRUCTURE**

**CORE
INFRASTRUCTURE**

**TLS, obfuscation, dynamic entitlements
Kerberos, PKI, AD/DS, etc.**

CONTEXT/COMMANDS

DATA/EVENTS

EDGE

CORE

- **Constrained**
- **High-latency**
- **Localized context**

- **Hybrid – cloud/on-premises**
- **Low-latency**
- **Global context**

**HORTONWORKS**®

# Back Pressure

GenerateFlowFile
GenerateFlowFile 1.2.0
org.apache.nifi - nifi-standard-nar

| In | 0 (0 bytes) | 5 min |
|---|---|---|
| Read/Write | 0 bytes / 0 bytes | 5 min |
| Out | 10,000 (0 bytes) | 5 min |
| Tasks/Time | 10,000 / 00:00:00.849 | 5 min |

Name success
Queued 10,000 (0 bytes)

LogAttribute
LogAttribute 1.2.0
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 |

Back Pressure Object Threshold ❓

10000

Back Pressure Data Size Threshold ❓

1 GB

- Configure back-pressure per connection
- Based on number of FlowFiles or total size of FlowFiles
- Upstream processor no longer scheduled to run until below threshold

**HORTONWORKS**®

# Prioritization

- Configure a prioritizer per connection

- Determine what is important for your data – time based, arrival order, importance of a data set

- Funnel many connections down to a single connection to prioritize across data sets

- Develop your own prioritizer if needed



Name success
Queued 10,000 (0 bytes)

LogAttribute
LogAttribute 1.2.0
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 |

Available Prioritizers ❓

FirstInFirstOutPrioritizer

NewestFlowFileFirstPrioritizer

OldestFlowFileFirstPrioritizer

PriorityAttributePrioritizer

Selected Prioritizers ❓

**HORTONWORKS**®

# Latency vs. Throughput

- Choose between lower latency, or higher throughput on each processor

- Higher throughput allows framework to batch together all operations for the selected amount of time for improved performance

- Processor developer determines whether to support this by using @SupportsBatching annotation

# Extension / Integration Points

| NiFi Term | Description |
|---|---|
| Flow File Processor | Push/Pull behavior.  Custom UI |
| Reporting Task | Used to push data from NiFi to some external service (metrics, provenance, etc..) |
| Controller Service | Used to enable reusable components / shared services throughout the flow |
| REST API | Allows clients to connect to pull information, change behavior, etc.. |

HORTONWORKS®

# NiFi Positioning

# NiFi Positioning



Enterprise Service Bus (Fuse, Mule, etc.)

Processing Framework (Storm, Spark, etc.)

Apache NiFi / MiNiFi

ETL (Informatica, etc.)

Messaging Bus (Kafka, MQ, etc.)

HORTONWORKS®

# Apache NiFi / Processing Frameworks

## NiFi

### Simple event processing

- **Primarily feed data into processing frameworks, can process data, with a focus on simple event processing**

- **Operate on a single piece of data, or in correlation with an enrichment dataset (enrichment, parsing, splitting, and transformations)**

- **Can scale out, but scale up better to take full advantage of hardware resources, run concurrent processing tasks/threads (processing terabytes of data per day on a single node)**

⚠ Not another distributed processing framework, but to feed data into those

## Processing Frameworks (Storm, Spark, etc.)

### Complex and distributed processing

- **Complex processing from multiple streams (JOIN operations)**

- **Analyzing data across time windows (rolling window aggregation, standard deviation, etc.)**

- **Scale out to thousands of nodes if needed**

⚠ Not designed to collect data or manage data flow

HORTONWORKS®

# Apache NiFi / Messaging Bus Services

## NiFi

### Provide dataflow solution

- **Centralized management, from edge to core**

- **Great traceability, event level data provenance starting when data is born**

- **Interactive command and control – real time operational visibility**

- **Dataflow management, including prioritization, back pressure, and edge intelligence**

- **Visual representation of global dataflow**

- ⚠ Not a messaging bus, flow maintenance needed when you have frequent consumer side updates

## Messaging Bus (Kafka, JMS, etc.)

### Provide messaging bus service

- **Low latency**

- **Great data durability**

- **Decentralized management (producers & consumers)**

- **Low broker maintenance for dynamic consumer side updates**

- ⚠ Not designed to solve dataflow problems (prioritization, edge intelligence, etc.)

- ⚠ Traceability limited to in/out of topics, no lineage

- ⚠ Lack of global view of components/connectivities

HORTONWORKS®

# Apache NiFi / Integration, or ingestion, Frameworks

## NiFi

**End user facing dataflow management tool**

- Out of the box solution for dataflow management

- Interactive command and control in the core, design and deploy on the edge

- Flexible failure handling at each point of the flow

- Visual representation of global dataflow and connectivities

- Native cross data center communication

- Data provenance for traceability

- ⚠ Not a library to be embedded in other applications

## Integration framework (Spring Integration, Camel, etc), ingestion framework (Flume, etc)

**Developer facing integration tool with a focus on data ingestion**

- A set of tools to orchestrate workflow

- A fixed design and deploy pattern

- Leverage messaging bus across disconnected networks

- ⚠ Developer facing, custom coding needed to optimize

- ⚠ Pre-built failure handling, lack of flexibility

- ⚠ No holistic view of global dataflow

- ⚠ No built-in data traceability

**HORTONWORKS**®

# Apache NiFi / ETL Tools

## NiFi

### NOT schema dependent

- **Dataflow management for both structured and unstructured data, powered by separation of metadata and payload**

- **Schema is not required, but you can have schema**

- **Minimum modeling effort, just enough to manage dataflows**

- **Do the plumbing job, maximize developers' brainpower for creative work**

⚠ Not designed to do heavy lifting transformation work for DB tables (JOIN datasets, etc.). You can create custom processors to do that, but long way to go to catch up with existing ETL tools from user experience perspective (GUI for data wrangling, cleansing, etc.)

## ETL (Informatica, etc.)

### Schema dependent

- **Tailored for Databases/WH**

- **ETL operations based on schema/data modeling**

- **Highly efficient, optimized performance**

⚠ Must pre-prepare your data, time consuming to build data modeling, and maintain schemas

⚠ Not geared towards handling unstructured data, PDF, Audio, Video, etc.

⚠ Not designed to solve dataflow problems

**HORTONWORKS®**

# Apache MiNiFi

# Apache MiNiFi

**"Let me get the key parts of NiFi close to where data begins and provide bidirectional data transfer"**

● 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

HORTONWORKS®

# Realities of computing outside the comforts of the data center

- Limited
  - computing capability
  - power/network

- Restricted software library/platform availability

- No UI

- Physically inaccessible

- Not frequently updated

- Competing standards/protocols

- Scalability

- Privacy & Security

**HORTONWORKS**®

# Apache ~~NiFi~~ MiNiFi

## Key Features

- Guaranteed delivery

- Data buffering
  - Backpressure
  - Pressure release

- Prioritized queuing

- Flow specific QoS
  - Latency vs. throughput
  - Loss tolerance

- Data provenance

- Recovery/recording a rolling log of fine-grained history

- Designed for extension

- Design and Deploy

- Warm re-deploys

HORTONWORKS®

# MiNiFi: Precedent from NiFi

## A quick look at NiFi Site to Site

- Provides the semantics between two NiFi components across network boundaries
  - A custom protocol for inter-NiFi communication
  - Secure, Extensible, Load Balanced & Scalable Delivery to Cluster

- Extracted out to a client library which powers integration into popular frameworks like Apache Spark, Apache Storm, Apache Flink, and Apache Apex

- Attributes and the FlowFile format maintained

*https://nifi.apache.org/docs/nifi-docs/html/user-guide.html#site-to-site*
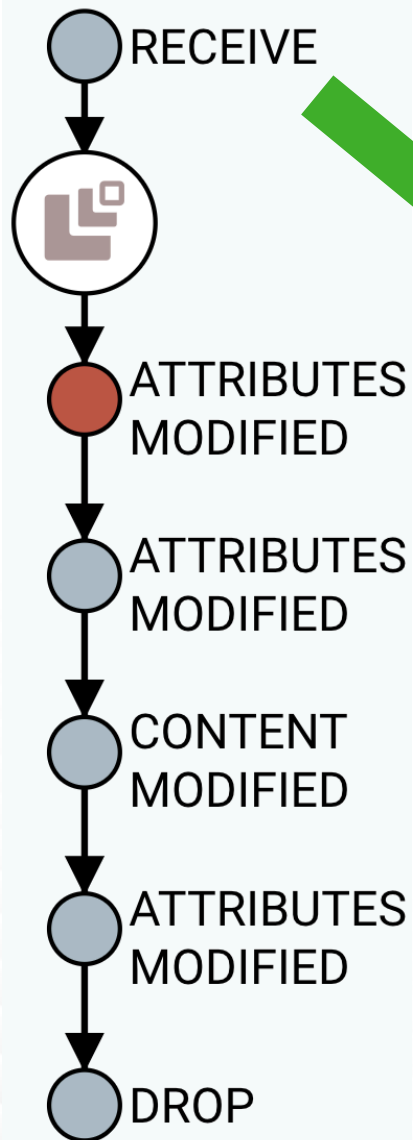
# MiNiFi: Precedent from NiFi

## A deeper dive into provenance

- Fine-grained, event level access of interactions with FlowFiles
  - CREATE, RECEIVE, FETCH, SEND, DOWNLOAD, DROP, EXPIRE, FORK, JOIN …

- Captures the associated attributes/metadata at the time of the event

- A map of a FlowFile's journey and how they relate to other FlowFiles in a system
  - MiNiFi enables us to get more and further illuminate the map of data processing

*http://nifi.apache.org/docs/nifi-docs/html/user-guide.html#data-provenance*

**HORTONWORKS®**

# MiNiFi: Precedent from NiFi

**RECEIVE event**



**Provenance Event**

| DETAILS | ATTRIBUTES | CONTENT |
|---------|-----------|---------|

Time
04/04/2017 16:59:49.642 CEST

Event Duration
< 1ms

Lineage Duration
< 1ms

Type
RECEIVE

FlowFile Uuid
e5a28106-e332-484f-9fd5-df4eb36015f0

File Size
11 bytes

Component Id
397820ce-015b-1000-b89d-7e7d65b8b671

Component Name
ListenHTTP

Component Type
ListenHTTP

**Parent FlowFiles** (0)

No parents

**Child FlowFiles** (0)

No children

# Apache MiNiFi

## Departures from NiFi in getting the right fit

- The feedback loop is longer and not guaranteed
  - Removal of Web Server and UI

- Declarative configuration
  - Lends itself well to CM processes
  - Extensible interface to support varying formats
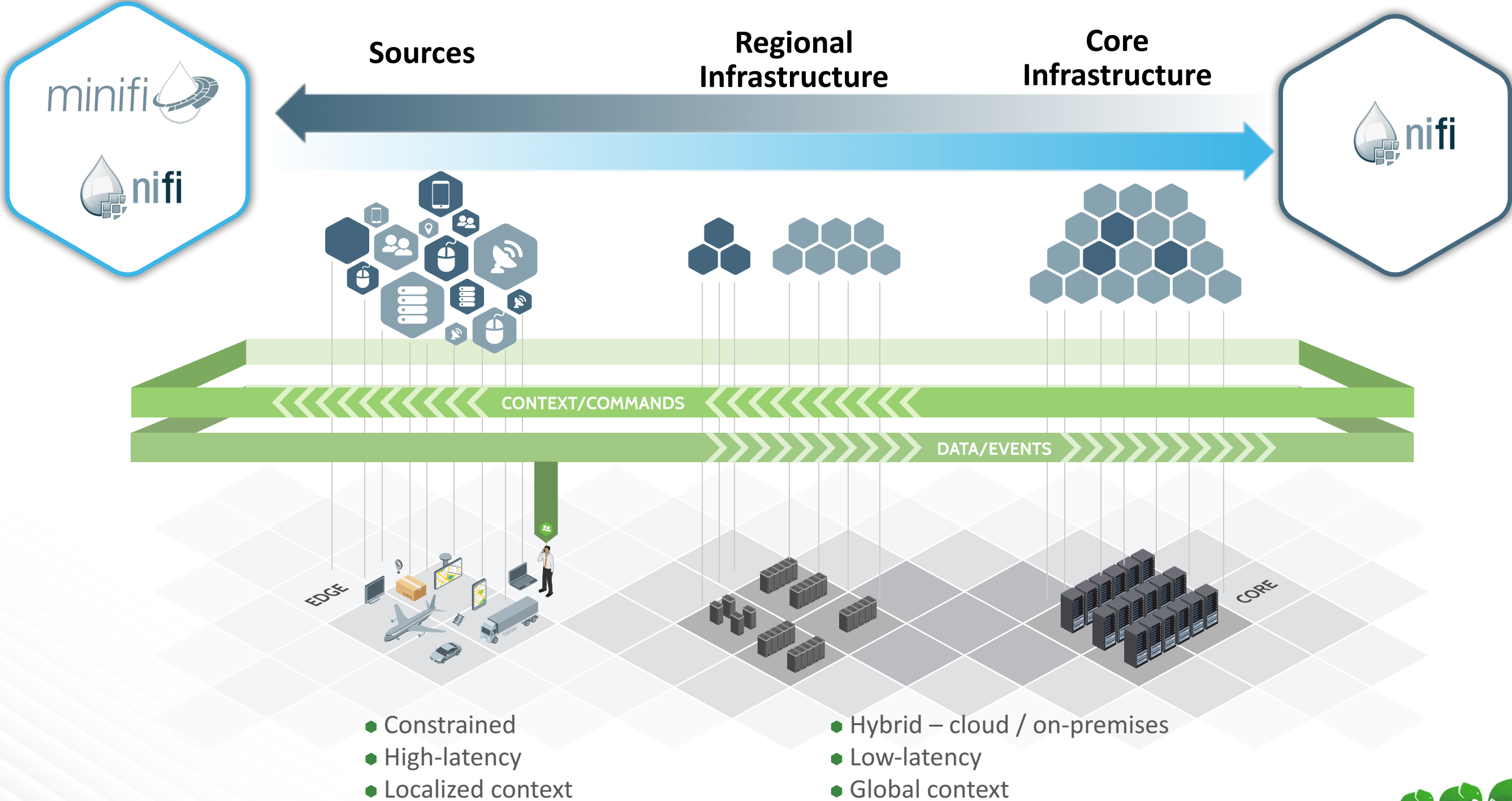    - Currently provided in YAML

- Reduced set of bundled components

**HORTONWORKS**®

# Apache MiNiFi:  Scoping

**Provide all the key principles of NiFi in varying, smaller footprints**

- ***Go small***:  Java – *Write once, run anywhere\**
  - Feature parity and reuse of core NiFi libraries

- ***Go smaller***:  C++ – *Write once\*\*, run anywhere*

- ***Go smallest***: *Write n-many times, run anywhere*
  Language libraries to support tagging, FlowFile format, Site to Site protocol, and provenance generation without a processing framework
  - Mobile:  Android & iOS
  - Language SDKs

# Harnessing Data in Motion



Sources

Regional Infrastructure

Core Infrastructure

minifi nifi

nifi

CONTEXT/COMMANDS

DATA/EVENTS

EDGE

CORE

- Constrained
- High-latency
- Localized context

- Hybrid – cloud / on-premises
- Low-latency
- Global context

HORTONWORKS®

# Demo

Hortonworks

# 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
users@nifi.apache.org

**Submit Ideas or Issues**
https://issues.apache.org/jira/browse/NIFI
https://issues.apache.org/jira/browse/MINIFI

**Follow on Twitter**
@apachenifi

HORTONWORKS®

# Thank You

**HORTONWORKS®**
POWERING THE FUTURE OF DATA™