cloudera

INTRODUCTION TO CDF

Pedro Algaba – Solution Engineer

AGENDA



Cloudera Data Flow platform

Flow & Edge management

Technical details

Stream Processing

CLOUDERA DATA FLOW

Cloudera dataflow Platform

Cloudera DataFlow Data-in-Motion Platform



EDGE DATA MANAGEMENT

Edge data collection, routing and monitoring

Apache MiNiFi

Edge Flow Manager



FLOW MANAGEMENT

Enterprise data ingestion, transformation and enrichment

Apache NiFi NiFi Registry

STREAM PROCESSING

Real-time streams processing at IoT scale

Apache Kafka Streams Messaging Manager

•••

STREAMING ANALYTICS Predictive analytics

and real-time insights

Apache Storm Streaming Analytics Manager Kafka Streams



ENTERPRISE SERVICES

Provisioning, Management and Monitoring Unified Security Edge-to-Enterprise Governance Single Sign-on

Schema Registry

HDF 3.3 Platform

Orange available in **HDP+HDF** scenario



FLOW MANAGEMENT

Data acquisition and delivery Simple transformation and data routing Simple event processing End to end provenance Edge intelligence & bi-directional communication



STREAM PROCESSING

Scalable data broker for streaming apps Scale Out Streaming Computation Engine

STREAM ANALYTICS

Pattern Matching Prescriptive & Predictive Stream Analytics Complex Event Processing Continuous Insight

🔆 kafko



ENTERPRISE SERVICES

Provisioning, Management, Monitoring, Security, Audit, Compliance, Governance, Multi-tenancy



Flow & Edge Management



Today's Data Landscape

- Teradata
- RDBMS
- Hive
- HDFS
- S3
- ADLS
- Google PubSub
- HBase
- Druid

- Memcache
- Solr
- Kafka
- Cassandra
- ElasticSearch
- AWS Redshift
- Couchbase
- JMS Queues
- MQTT

- REST
- SFTP
- MongoDB
- AeroSpike
- Redis
- Splunk
- Syslog
- InfluxDB

What tools do you use?

Sqoop

• Flume

• ETL tools

Script

Script_V2

Script_incremental_tets

SITUATION: THERE ARE 14 COMPETING STANDARDS.



Challenges with Current Frameworks

• Most frameworks on the market are focused on transformations (i.e. ETL vs EL)

- **Traditional**: Informatica, Pentaho, Talend, SSIS
- **Big Data**: MR, Tez, Pig, Hive, Spark, Storm, Flink
- EL only requires simple processing (e.g. translation)
 - i.e. it is a map-only MR job with no shuffle; a Storm topology with one bolt
- Sqoop, DistCp, and MirrorMaker are specialized to a specific route
- Kafka Connect requires Kafka as an intermediary

=> NiFi is a general purpose EtL engine

Flow Management

Enable easy ingestion, routing, management and delivery of any data anywhere (Edge, cloud, data center) to any downstream system with built in end-to-end security and provenance





Advanced tooling to industrialize flow development (Flow Development Life Cycle)



- Over 280+ Prebuilt Processors
- Easy to build your own
- Parse, Enrich & Apply Schema
- Filter, Split, Merger & Route
- Throttle & Backpressure

- Guaranteed Delivery
- Full data provenance from acquisition to delivery
- Diverse, Non-Traditional Sources
- Eco-system integration



An overview of NiFi capabilities



Data Ingest

Data Transformation



An overview of NiFi capabilities

Data Enrichment







An overview of NiFi capabilities



Routing Content Attribute Query Partition Multi-ingest Merge Prioritiz

Why should you care about NiFi?

- Accelerates development time with its UI and more than 260 OOTB processors
- **Reduces costs** by offloading expensive techs such as CFT, Talend
- Industrializes data movements with central monitoring, flow/variable registries
- Reduces risks with Security and lineage (Ranger, Atlas, NiFi provenance)
- Improves visibility with data ingest centralization and monitoring
- **Breaks silos** with large ecosystem integration
- Improves application performance with horizontal and vertical scalability
- Enables new use case such as log optimization, IoT, streaming_clotheradic. All rights reserved

TECHNICAL DETAILS

Apache NiFi High Level Capabilities

- Strong Out Of The Box
 - Flow can be modified at runtime
 - Back pressure
 - Loss tolerant vs guaranteed delivery
 - Low latency vs high throughput
 - Dynamic prioritization
- Secure
 - SSL, HTTPS, SFTP, etc.
 - End to end data provenance
- Designed for extension
 - Provider are a first class citizen
 - Build your own processors and Controller services
 - Integrate with external systems (Security, Monitoring, Governance, etc)





🍖 NiFi Flow 🗙 🍐 NiFi	× IniFi Registry × +
(←) → C'	nifi/?processGroupId=root&componentIds=567ff3f1-528f-3348 ••• 🗵 🏠 🔍 Search 👱 🛝 🗊 💩 🔾 \Xi
🗀 Hortonworks 💥 SME 🧥 Docs 🥖 Docs 🐤 Re	eports 🌍 SMM DPS 🧔 SMM HDF 🔗 CB 🧔 HDF Field 🍖 NiFi 🚏 Registry 🕀 C2 💽 Git Cloudera 🕞 CDSW 🕀 SmartSense 🧮 PDF 🤛 DWS 🔊
nifi 💿 🕹 🖙	
0 🗮 1,775 / 438.79 KB 💿 0	ⓐ 0 ⓑ 3 월 4 ▲ 0 ♀ 0 ♣ 0 ⓑ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 ♀ 0 <t< th=""></t<>
0	
 ↓ 	Dev Cluster
▶	
	Data Generation Cloud Demo
	◎ 0 ◎ 0 ▶ 0 ■ 2 ▲ 0 % 0 ◎ 0 ◎ 0 ▶ 3 ■ 0 ▲ 0 % 0
	Queued 1,681 (345.79 KB) Queued 92 (92.67 KB)
	In $0 (0 \text{ bytes}) \rightarrow 0$ 5 min In $0 (0 \text{ bytes}) \rightarrow 0$ 5 min In $0 (0 \text{ bytes}) \rightarrow 0$
	Read/Write 20.57 KB / 20.57 KB 5 min Read/Write 8.47 KB / 25.51 KB 5 min
	Out $0 \rightarrow 0$ (0 bytes)5 minOut $0 \rightarrow 0$ (0 bytes)5 min
	✓ 0 * 0 ⊙ 0 ⊕ 0 ? 0
	Big Query
	Queued 2 (327 bytes)
	In $0 (0 \text{ bytes}) \rightarrow 0$ 5 min
	Read/Write 0 bytes / 0 bytes 5 min
	Out $0 \rightarrow 0$ (0 bytes) 5 min
	✓ 0 * 0 0 0 0 ? 0

FDLC: Flow Development LifeCycle



NiFi - Extensibility

- Built from the ground up with extensions in mind
- Service-loader pattern for...
 - Processors
 - Controller Services
 - Reporting Tasks
 - Prioritizers
- Extensions packaged as NiFi Archives (NARs)
 - Deploy NiFi lib directory and restart
 - Same model as standard components



<parent>

<groupId>org.apache.nifi</groupId>
<artifactId>nifi-flatjson-bundle</artifactId>
<version>1.0-SNAPSHOT</version>
</parent>

<artifactId>nifi-flatjson-processors</artifactId>
<packaging>jar</packaging>

<dependencies>

<dependency>
 <groupId>com.github.wnameless</groupId>
 <artifactId>json-flattener</artifactId>
 <version>0.4.0</version>
</dependency>
<dependency>
 <groupId>org.apache.nifi</groupId>
 <artifactId>nifi-api</artifactId>

</dependency> <dependency>

<groupId>org.apache.nifi</groupId>
 <artifactId>nifi-utils</artifactId>
</dependency>

NiFi - Architecture



Scalable and distributed architecture

•	S/Host				
	— 0	S/Host			
Q0	₿.	•	S/Host		
	Q°	₿.	_ 0	S/Host	
5		¢°	₿.	OS/Host	ZooKeeper Server
Local	5		Q 0	G JVM Server	
	Local			C Flow Controller	
		Local	1	Processor 1 Extension N	Cluster Coordinator
			Local	FlowFile Content Provenance Repository	🜲 Primary Node
				Local Storage	ZooKeeper Client

Hub and Spoke Architecture with Site-to-Site (S2S)

- Send data from a NiFi instance/cluster to one or multiple NiFi instances/clusters
- Preferred communications protocol when NiFi on both ends
- · 2-way secure protocol, push & pull, high availability and load balancing





NiFi vs. MiNiFi Java Prcessor, Smaller Footprint ~40 MB



How Does MiNiFi Interact With NiFi?

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





Stream Processing



Stream Processing



Apache Kafka "Pub & Sub"

Massively Scalable Publish-Subscribe Message Queue



Real-time, distributed data and task processing engine



Streams Messaging Manager

eview		8* 4- B Topics / Other Events	8* -
Eldris *	27 · 10 · 4	57 60 501 Minute* 201 Minute* 501 Minute* 201 Minute*	S Led Har 2000
Advantes (27) Other Events	1234 manufactore 2349 m		and 11 Arrows
n an NewsFeed	301 1,619 - 3 Overview	A ² Ar	Lini Stal Menages ma
terrente in BidgetAtenios	282 2,555 *-	815 156727 49710 199746 82,0757 -	
rendro 104 CotalQuatomena rendro 103 CotalQuatomena Activit	257 million 2700 to 2	Other Darks Aussission Other Darks Aussission Other Darks Aussission Other Darks Aussission Other Darks	فسخطأها المسخد
ten 19 O TruckBoets	234 2321	And a second sec	
Anno D	154 memory 1,619 to 2		
Anima B a N a N A A A A A A A A A A A A A	Manager Manager Manager Manager M Manager Manager Mana Manager Manager Man Manager Manager Mana Manager Manager Manager Manager Manager Manage		

- Cure Kafka blindness and help the different • streaming personas be more productive
- End-to-end integration with Ambari, • Grafana, Ranger & Atlas
- Comprehensive REST service for open integration





Children III alla	m			1 m	-	(0) and (0)
0.000	5945.4	0.	12-	10	3	22
					1.0	
	* > 17 (m)		÷			
			9			Louis in
5.		1 8			N N	A D D
IA.	-	1	1. 18	-		Statie -
Transferrer	1000			ware a		-[
		-			Teres .	
-						

rts Dashboard is			
25.0	Salat Looks to Talka	Marries and search of the sear	
539	Product Relation Study Kee Kee	luuuuu	
		:	-

- Create streaming applications easily
- Visualize enrichment, matching, aggregations, notifications, etc

Our committer focus areas...

- Time to insights
- Security and governance
- **DevOps and SDLC Support**
- Management & Monitoring
- Stability and performance

5 Building blocks for robust and scalable real time solutions



A real-time integrated data logistics and simple event processing platform



Publish subscribe high-throughput, lowlatency broker



API for building real-time applications and microservices on top of Kafka

Kafka's Omnipresence Has Led to the Onset of "Kafka Blindness"

• What is **"Kafka Blindness**"?

 Customers who use Kafka today struggle with monitoring / "seeing"/troubleshooting what is happening in their clusters

• Who is Affected?

- Platform Operation Teams
- Developers / DevOps Teams
- Security / Governance Teams

• What are the **Symptoms**?

- Difficulty seeing who is producing and consuming data
- Difficulty understanding the flow of data from producers -> topics \rightarrow consumers
- Difficulty troubleshooting/monitoring.

Hortonworks Streams Messaging Manager (SMM)

- New Open Source project led by Hortonworks to Cure the "Kafka Blindness"
- Single Monitoring Dashboard for all your Kafka Clusters across 4 entities
 - Broker
 - Topic
 - Producer
 - Consumer
- Designed for the Enterprise
 - Support for Secure/Kerborized Kafka cluster
 - Rich Access Control Policies (ACLS)
 - Supports multiple HDP and/or HDF Kafka Clusters
- REST as a First Class Citizen
- Delivered as a DataPlane Service



•	
OVORVI	10
Uvervi	15

H

83			3		*	28	×.		Const	19	Cl
TOPICS (28) BROKERS (3)										c	2 da
Producers (83)		NAME	A V	DATA IN \$	DATA OUT 🗢	MESSAGES IN \$	CONSUMER GROUPS \$			Consumer Groups (19)	
ACTIVE (83) PASSIVE (0)	ALL	syndi	icate-transmission	55MB	294KB	0.3m	0	🧑 🚱 Q 🔳	~	ACTIVE (17) PASSIVE (2)	AL
MESS geo-critical-event-collector-i1	sages ≑ 63k									predictive-micro-service	2
route-apps	46k	syndi	icate-speed-event-json	48MB	0B	0.2m	0	🧑 🚱 Q 🔳	~	load-optimizer-micro-service	
geo-critical-event-collector-i2	32k									route-micro-service	1
minifi-eu-i1	32k	syndi	icate-speed-event-avro	29MB	0B	0.2m	0	👸 😡 Q 🔳	~	kafka-streams-analytics-geo	. 0
load-optimizer-apps	30k						-			flink-analytics-geo-event	C
geo-critical-event-collector-i3	22k									spark-streaming-analytics-g	0
geo-critical-event-collector-i6	21k	syndi	icate-oil	1MB	0B	6.1k	0	🧔 🛛 🖓 🖾	~	adjudication-micro-service	0
nifi-syndicate-speed-avro	20k									audit-micro-service	0
nifi-syndicate-geo-avro	20k	syndi	icate-geo-event-json	65MB	0B	0.2m	0	👌 😧 Q 🔳	~	supply-chain-micro-service	C
nifi-syndicate-geo-json	20k									energy-micro-service	0
nifi-syndicate-speed-json	20k				10115		-			compliance-micro-service	0
minifi-eu-i2	16k 💙	synal	icate-geo-event-avro	37MB	19MB	0.2m	3	10 V Q 🗉	~	approval-micro-service	
geo-critical-event-collector-i4	16k									nifi-truck-sensors-central	
minifi-eu-i4	16k	syndi	icate-battery	60MB	0B	0.3m	0	🧑 🥝 Q 🔳	~	nifi-truck-sensors-east	
fuel-apps	15k									nifi-truck-sensors-west	
predictive-apps	15k	syndi	icate-all-geo-critical-ev	16B	113MB	5.2m	1	10 Q Q 📼	~	ranger_entities_consumer	
minifi-truck-c1	13k		in goo on too on the	100	TISIND	0.2111	1	~~~		atlas	

Schema Registry

raw-truck_events_avro	TYPE avro	GROUP truck-sen	BRANCH 2 📽	SERIALIZER & DESER O	RIALIZER	~	
BRANCH : Dev BRANCH DESCRIPTION : dev VERSION DESCRIPTION :	<pre>1 { 2 "type": "record", 3 "namespace": "hortonwo 4 "name": "truckgeoevent 5 "fields": [6 { 7 "name": "eventTime",</pre>	rks.hdp.refaj	VERSION 2	BRANCH: Dee CHANGE LOG V2 V1 V1 V1 V1	go ITED Review P 565 ago		
adding eventime long field	raw-truck_events_avro		TYPE avro	GROUP truck-sen	branch 2 ℃	SERIALIZER & DESERIALIZER	^
	BRANCH : Dev BRANCH DESCRIPTION :	1 { 2 "ty 3 "na 4 "na 5 "f	ype": "record", amespace": "hortonwo ame": "truckgeoevent	prks.hdp.refapp.tr t",	VERSION 2	CHANGE LOG	
	dev VERSION DESCRIPTION : adding eventime long field	6 { 7 } 9 }.	'name": "eventTime", 'type": "string"	,		2d 17h 26m 26s ago Enabled	
		10 { 11 12 13 14 } 15 f	'name": "eventTimeLo 'type": "long", 'default": 0	ong",			

Schema registry is integrated with NiFi, Kafka and SAM



Header	Payload
First 14 bytes contains info to identify the schemafor the event in the Schema Registry: schemald, schemaVersion, protocol	The remaining bytes is the Avro object serialized. The Avro object will not have the full schema as a traditional Avro object contains

Configure Controller Serv	COMMENTS	
Required field		
Property		
Schema Access Strategy	Use Embedded Avro Schema 🗸	
Schema Registry	Use 'Schema Name' Property	
Schema Name	Lico 'Scheme Text' Preparty	
Schema Text	Use Schema Text Property	
Schema Access Strategy	HWX Schema Reference Attributes	
	HWX Content-Encoded Schema Reference	
	Use Embedded Avro Schema 🚱	
		04
		— ОК

THANK YOU

cloudera