



# Cloudera DataFlow - Virtual Hands-on Workshop

October, 2021

# Cloudera Team



Suresh MR  
Director - Channels &  
Alliances



Vinay Rayker  
Partner Technology  
Leader



Puneet Joshi  
Partner Solutions  
Engineer



Pannag Katti  
Partner Solutions  
Engineer



Manick Mehra  
Partner Solutions  
Engineer

---

# AGENDA

- Cloudera DataFlow - Overview (20 minutes)
- Workshop Overview & Detailed Instructions (10 minutes)
- Hands-on exercises (120 minutes)
- Q&A (throughout the workshop)

# Cloudera Data Flow

# CLOUDERA - THE ENTERPRISE DATA CLOUD COMPANY

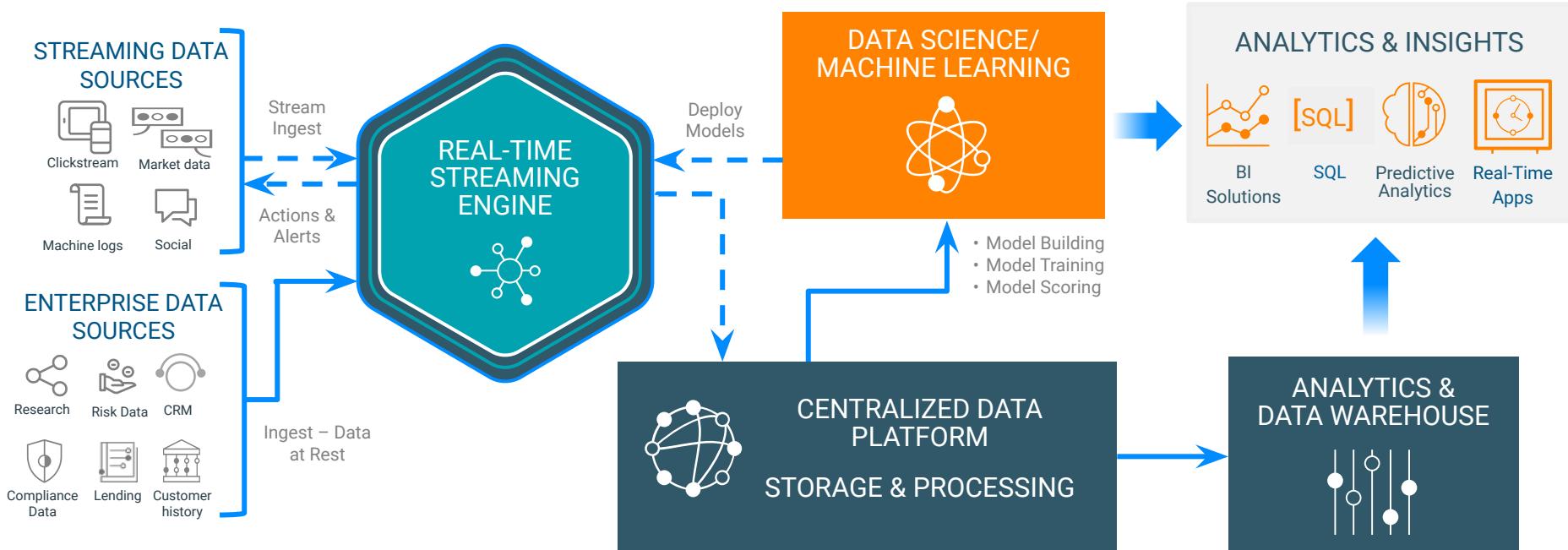
Manage and secure the data lifecycle in any cloud or datacenter



SECURITY | GOVERNANCE | LINEAGE | MANAGEMENT | AUTOMATION

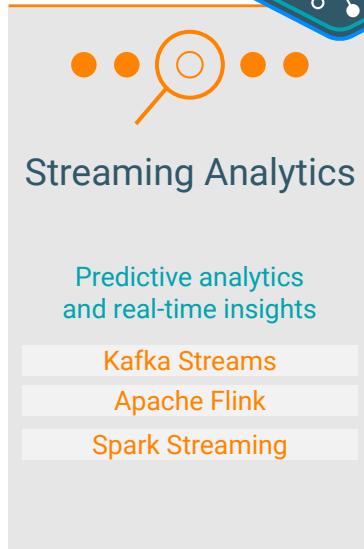
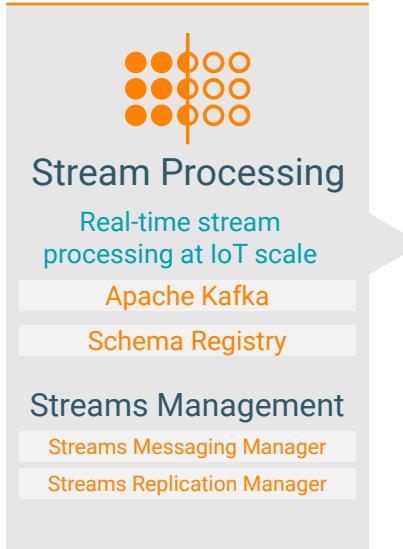
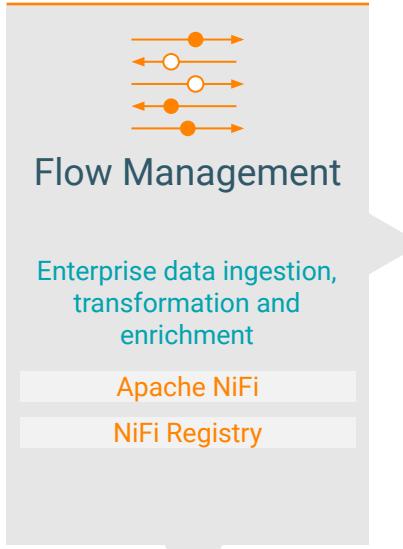
# ENABLING ANALYTICS AND INSIGHTS ANYWHERE

## Driving enterprise business value





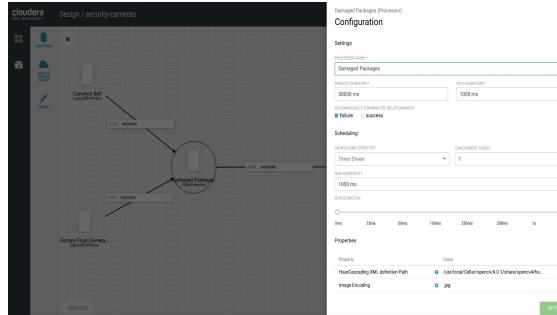
# Cloudera Data Flow



# Cloudera EDGE Management

Edge device data collection and processing with easy to use central command and control

## Edge Flow Manager

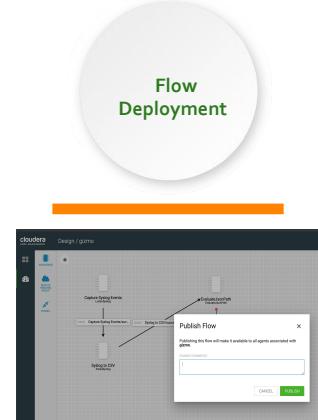


A lightweight edge agent that implements the core features of Apache NiFi, focusing on data collection and processing at the edge

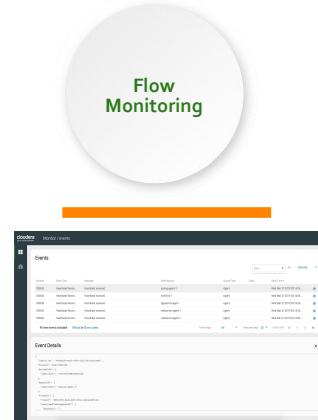
### Flow Authorship



### Flow Deployment



### Flow Monitoring

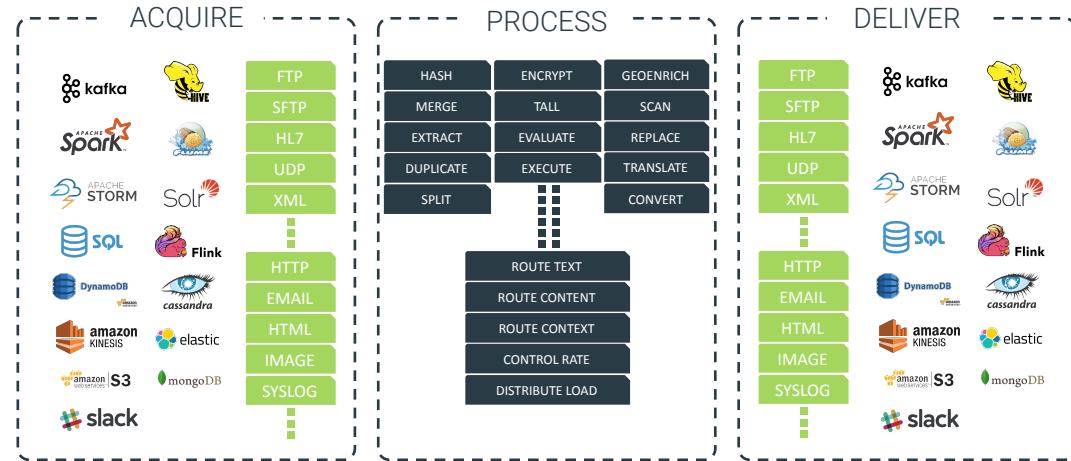
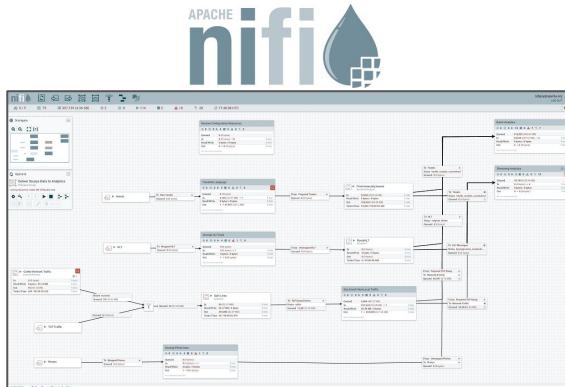


- Small footprint agent with MiNiFi
- Java and C++ agents
- Rich edge processors (edge collection & processing)
- End to end lineage and security

- Central Command and Control (C2)
- Design and deploy to thousands of agents
- Edge Applications lifecycle management
- Multitenancy with Agent classes
- Native integration with other CDF services

# Cloudera 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



- Over 300 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

# Stream processing & Analytics

Providing simpler ways of complex event processing at scale for the enterprise

 **kafka**  
+ KStreams

“Pub & Sub”

Massively Scalable Publish-Subscribe Message Queue, API for building real-time microservices with Kafka Streams

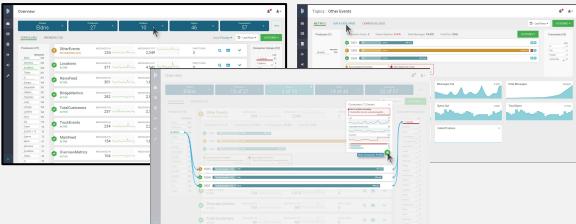
 “Process”

Reliable real-time/micro-batch data processing engines able to process millions tuples/second/node



  
**Streams Messaging Manager**

“Manage & Monitor”



- 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

  
**Streams Replication Manager (SRM)**

“Mirroring & DR”

Replicate data & configuration in Realtime, Smart clients for easy fail-over & fail-back, Support active active scenarios, Monitoring

  
**SCHEMA REGISTRY**

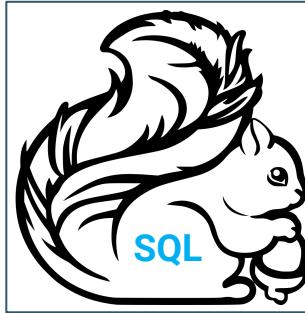


“Manage”

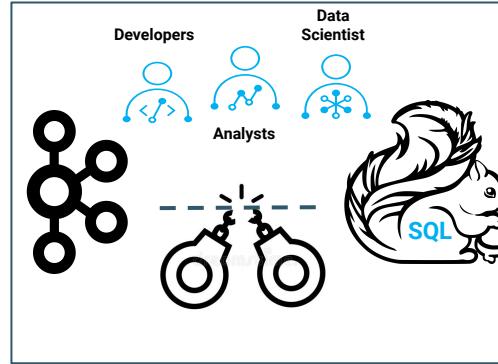
A shared schemas repository that allows applications to save, retrieve and reuse schemas and interact with each other

# CLOUDERA STREAMING ANALYTICS INCL. SQL STREAM BUILDER

Democratize access to real-time data with just SQL



Runs on  
Apache Flink



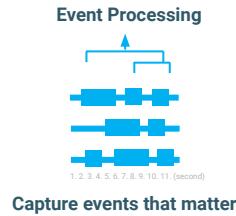
Liberates access to  
data within Kafka  
and Flink



No need for Java  
and Scala experts

# Enabling Streaming SQL (aka Continuous SQL)

Query language for processing live data streams



Streaming SQL Console

Console

Run SQL against unbounded streams of data and create persistent SQL streaming jobs

Compose Virtual Tables Functions History SQL Jobs

SQL Job Name: practical\_kepler Random Name

Sink Virtual Table: None

Advanced settings

SQL

```
1 -- detect multiple auths in a short window and
2 -- log them to topic#microservice
3 SELECT card, ts
4 MAX(amount) as theamount,
5 TUMBLE_END(eventtimestamp, interval '$1 minute') as ts
6 FROM card_tx
7 WHERE ts IS NOT NULL
8 AND lon IS NOT NULL
9 GROUP BY card, TUMBLE(eventtimestamp, interval '$1 minute')
10 HAVING COUNT(*) > 4 -- ==> fraud
```

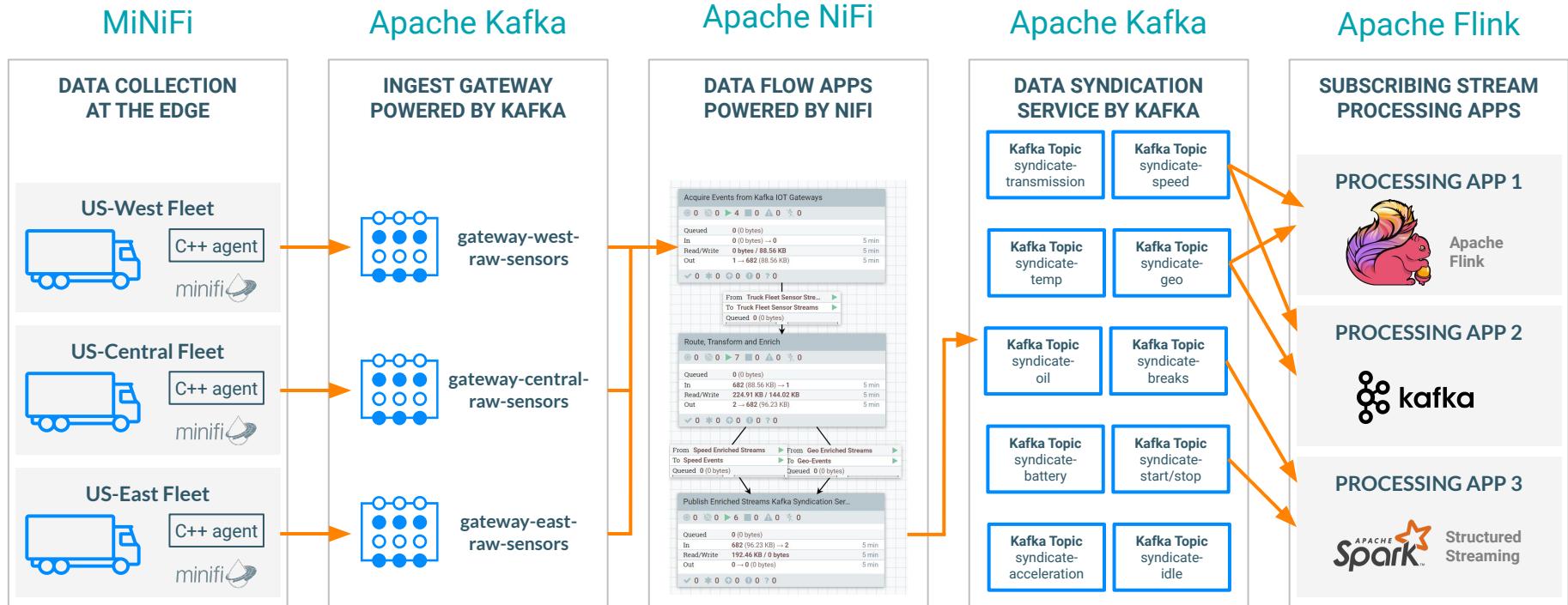
Logs Results Help

[2/23/2021, 4:11:16 PM] [INFO] StreamBuilder is ready.

This screenshot shows the Streaming SQL Console interface. It features a sidebar with 'Console', 'Data Sources', 'Teams', and 'API' sections. The main area has tabs for 'Compose', 'Virtual Tables', 'Functions', 'History', and 'SQL Jobs'. A 'SQL Job Name' field is set to 'practical\_kepler'. The 'Sink Virtual Table' dropdown is set to 'None'. Below these are 'Advanced settings' and a code editor containing a complex SQL query for detecting fraud. At the bottom, there are tabs for 'Logs', 'Results', and 'Help', and a status message: '[2/23/2021, 4:11:16 PM] [INFO] StreamBuilder is ready.'



# Data-in-Motion Reference Architecture



# Workshop Overview

# Predicting machine failures

Photos by Dominik Vanyi on Unsplash



# Sensors

Photo by Kai Dahms on Unsplash

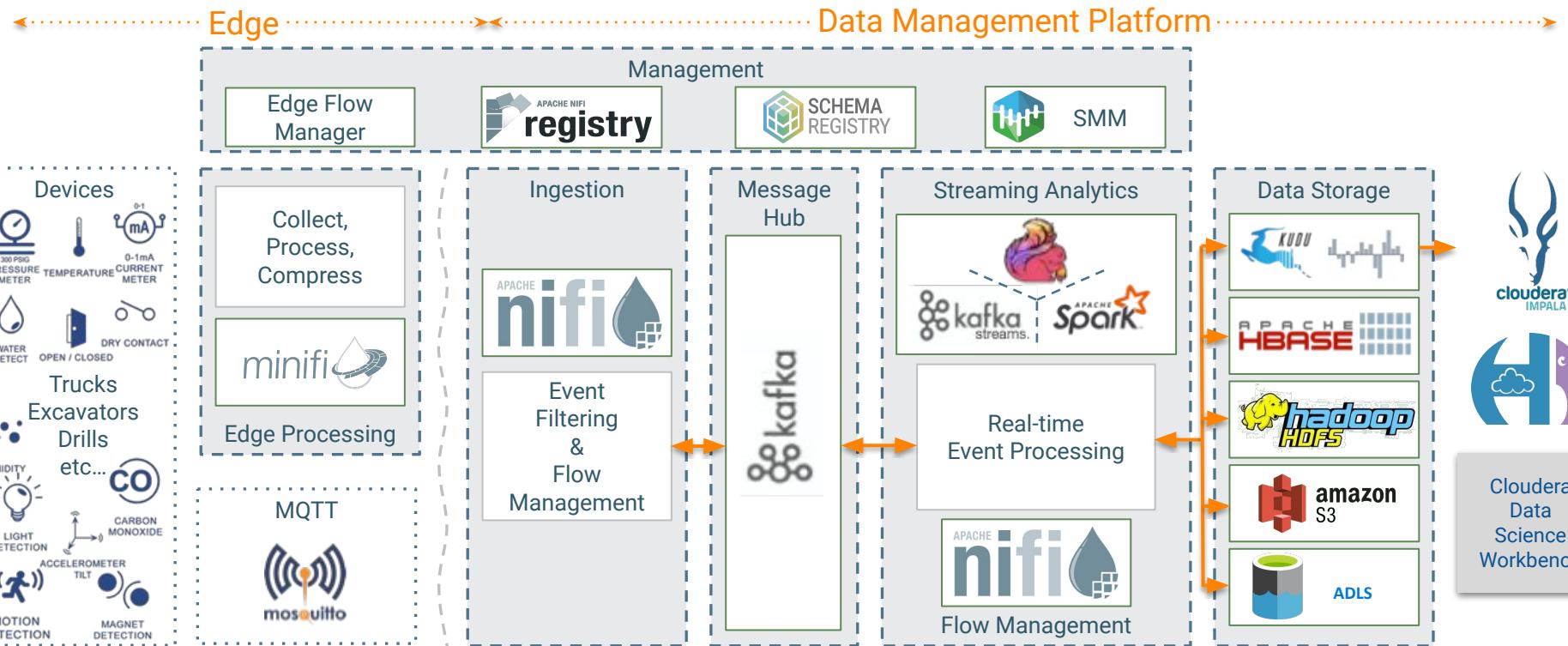


- Temperatures of different machines
- 12 per plant
- Data is published in {"JSON": "format"} to a MQTT server

# Streaming and Analytics

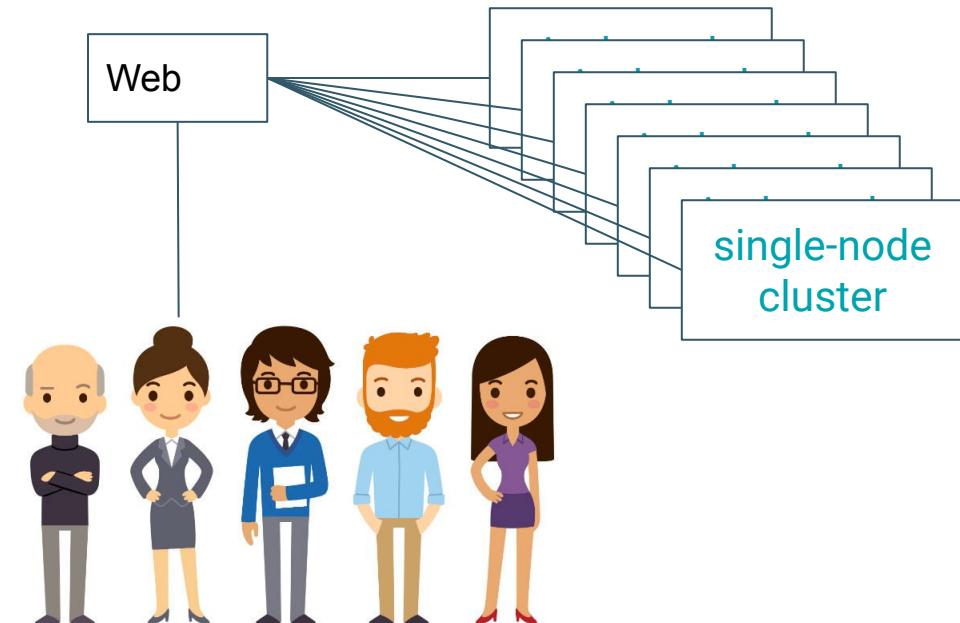
- Process the real-time data and it's collected from plants sensors
- Collect data in our central NiFi, update and forward it to our main Kafka message bus
- Read sensor readings from Kafka and persist into Kudu
- Enable fast analytics on Kudu
- Deploy a small Streaming Analytics application based on Apache Flink which read data stream from kafka topic, summarizing it and publish result into a new topic

# And that's what we are doing today



# The demo environment

Each participant receives a single-node cluster



## Hardware:

- Instance Size m5.4xlarge
- vCPU 16
- Memory (GiB) 64

## Software:

- OS CentOS Linux 7.6
- Cloudera **CDP PvC Base 7.1.6**
  - Flow Management 2.0
  - Stream Processing 3.0

# How do I get access to a cluster?

The screenshot displays the Cloudera Manager interface with several service components highlighted:

- NiFi Registry**: Shows the NiFi Registry / Administration interface with sections for Buckets, Users, and Design.
- Schema Registry**: Shows the Schema Registry interface with a list of Producers and Topics.
- Cloudera Edge Flow Manager**: Shows the Cloudera Edge Flow Manager interface with a Projects section.
- Streams Messaging Manager (SMM)**: Shows the SMM interface with a Projects section.
- Hue**: Shows the Hue interface with a Projects section.
- NiFi**: Shows the NiFi interface with a Process Group named "NFI Flow".

Each highlighted component has an orange square icon to its right.

TH<sup></sup>ON<sup></sup> Y<sup></sup>U<sup></sup>