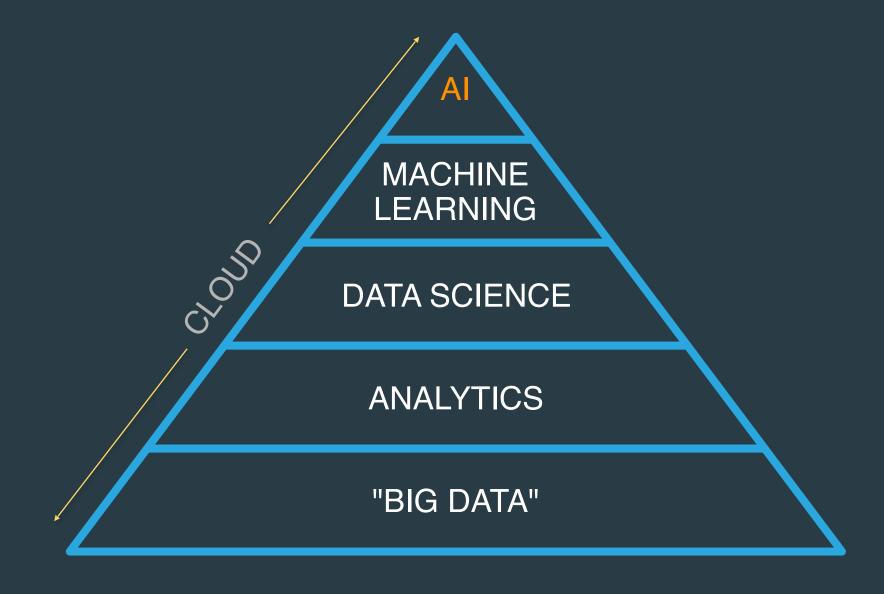
# cloudera

# Edge2AI

Johannes Muselaers I Sales Engineer <u>jmuselaers@cloudera.com</u> +46725881091



# CLOUDERA DATA SCIENCE WORKBENCH

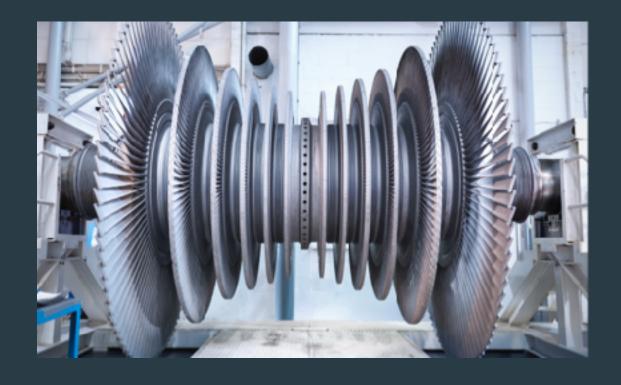
### MACHINE LEARNING AT CLOUDERA

Our philosophy

We empower our customers to run their business on data with an open platform:

- Your data
- Open algorithms
- Running anywhere

We accelerate enterprise data science.



### What is Cloudera Data Science Workbench (CDSW)

Supports your workloads in an enterprise secure way

# Accelerate data science from exploration to production using R, Python, Spark and more

#### For data scientists



Open data science, your way.

Use R, Python, or Scala with your favorite libraries and frameworks



No need to sample.

Directly access data in secure Hadoop clusters through Apache Spark and Apache Impala



Reproducible, collaborative research.

Share insights with your whole team

### For IT professionals



Bring analysis to the data.

Give your data science team the freedom to work how they want, when they want



Secure by default.

Stay compliant with out-of-the-box support for full Hadoop security



Flexible deployment.

Run on-premises or in the cloud

### Data Science & Data Scientist



# THE CHALLENGE

Balance these needs

### DATA SCIENCIST

- Access to granular data
- Flexibility
  - Preferred open source tools
- Elastic provisioning
  - Compute
  - Storage
- Reproducible research
- Path to production



### DevOps/IT

- Security
- Governance
- Standards
- Low maintenance
- Low cost
- Self-service access

WS,

### THE TYPICAL SOLUTION

### "If I can't use my favorite tools, I'll..."

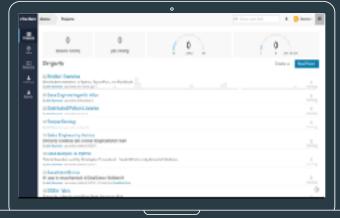
- Copy data to my laptop
- Copy data to a data science appliance
- Copy data to a cloud service

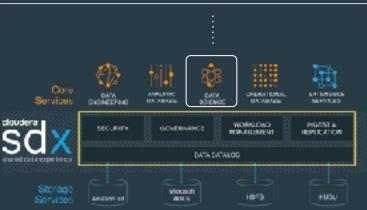
### Why this is a problem:

- Complicates security
- Breaks data governance
- Adds latency to process
- Makes collaboration more difficult
- Complicates model management and deployment
- Creates infrastructure silos

### CLOUDERA DATA SCIENCE WORKBENCH

### Accelerate Machine Learning from Research to Production





### For data scientists

- Experiment faster
   Use R, Python, or Scala with
   on-demand compute and
   secure CDH data access
- Work together
   Share reproducible research with your whole team
- Deploy with confidence
   Get to production repeatably
   and without recoding

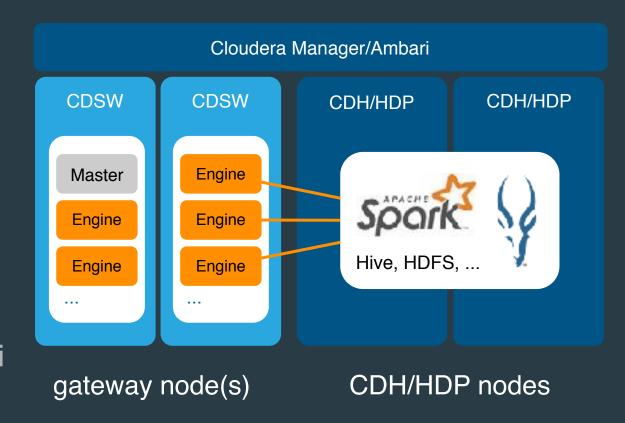
### For IT professionals

- Bring data science to the data
   Give your data science team
   more freedom while reducing
   the risk and cost of silos
- Secure by default
   Leverage common security and governance across workloads
- Run anywhere
   On-premises or in the cloud

### A MODERN DATA SCIENCE ARCHITECTURE

Containerized environments with scalable, on-demand compute

- Built with Docker and Kubernetes
  - Isolated, reproducible user environments
- Supports both big and small data
  - Local Python, R, Scala runtimes
  - Schedule & share GPU resources
  - Run Spark, Impala, and other CDH services
- Secure and governed by default
  - Easy, audited access to Kerberized clusters
  - Leverages SDX platform services
- Deployed with Cloudera Manager/Ambari



### ACCELERATED DEEP LEARNING WITH GPUS

Multi-tenant GPU support on-premises or cloud

"Our data scientists want GPUs, but we need multi-tenancy. If they go to the cloud on their own, it's expensive and we lose governance."

- Extend CDSW to deep learning
- Schedule & share GPU resources
- Train on GPUs, deploy on CPUs
- Works on-premises or cloud



single-node training



distributed training, scoring



### WHAT DATA SCIENCE TEAMS DO

### PREPARE DATA

Ingest data at scale.

Store and secure data.

Clean and transform data for analysis.

### **BUILD MODELS**

Explore data and build predictive models, offline.

Evaluate and tune models.

Develop and deliver a modeling pipeline.

### **DEPLOY MODELS**

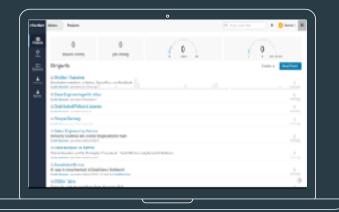
Test, verify, and approve model for deployment.

Create and maintain batch/ stream pipelines, embedded models, APIs.

Update models in production.

### CLOUDERA DATA SCIENCE WORKBENCH

Accelerate and simplify machine learning from research to production





#### **ANALYZE DATA**

 Explore data securely and share insights with the team



### TRAIN MODELS

Run, track, and compare reproducible experiments



### **DEPLOY APIs**

 Deploy and monitor models as APIs to serve predictions



#### MANAGE SHARED RESOURCES

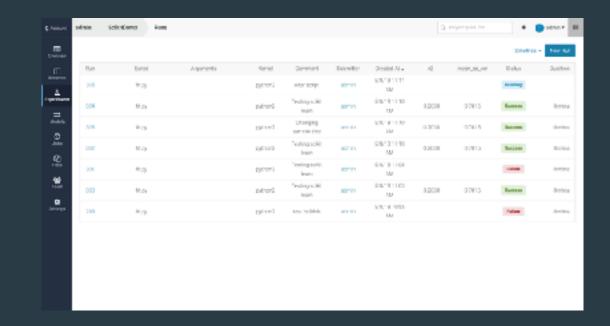
• Provide a secure, collaborative, self-service platform for your data science teams

### INTRODUCING EXPERIMENTS

Versioned model training runs for evaluation and reproducibility

### Data scientists can now...

- Create a snapshot of model code, dependencies, and configuration necessary to train the model
- Build and execute the training run in an isolated container
- Track specified model metrics, performance, and model artifacts
- Inspect, compare, or deploy prior models



### INTRODUCING MODELS

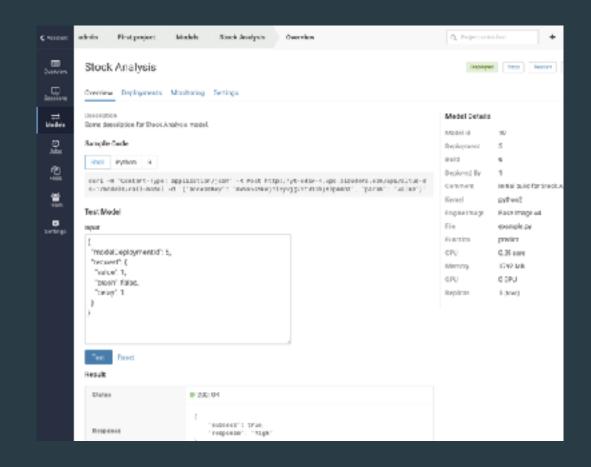
### Machine learning models as one-click microservices (REST APIs)

- 1. Choose file, e.g. score.py
- 2. Choose function, e.g. forecast

```
f = open('model.pk', 'rb')
model = pickle.load(f)
def forecast(data):
    return model.predict(data)
```

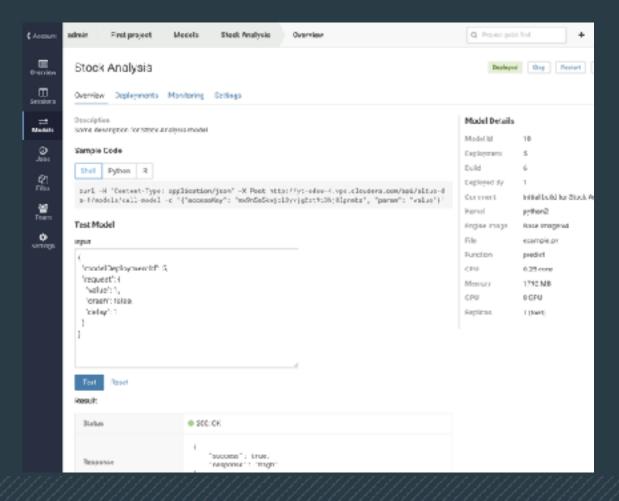
- 3. Choose resources
- 4. Deploy!

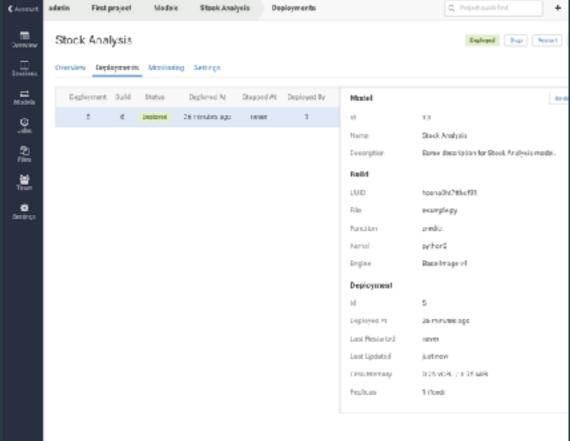
Running model containers also have access to CDH for data lookups.

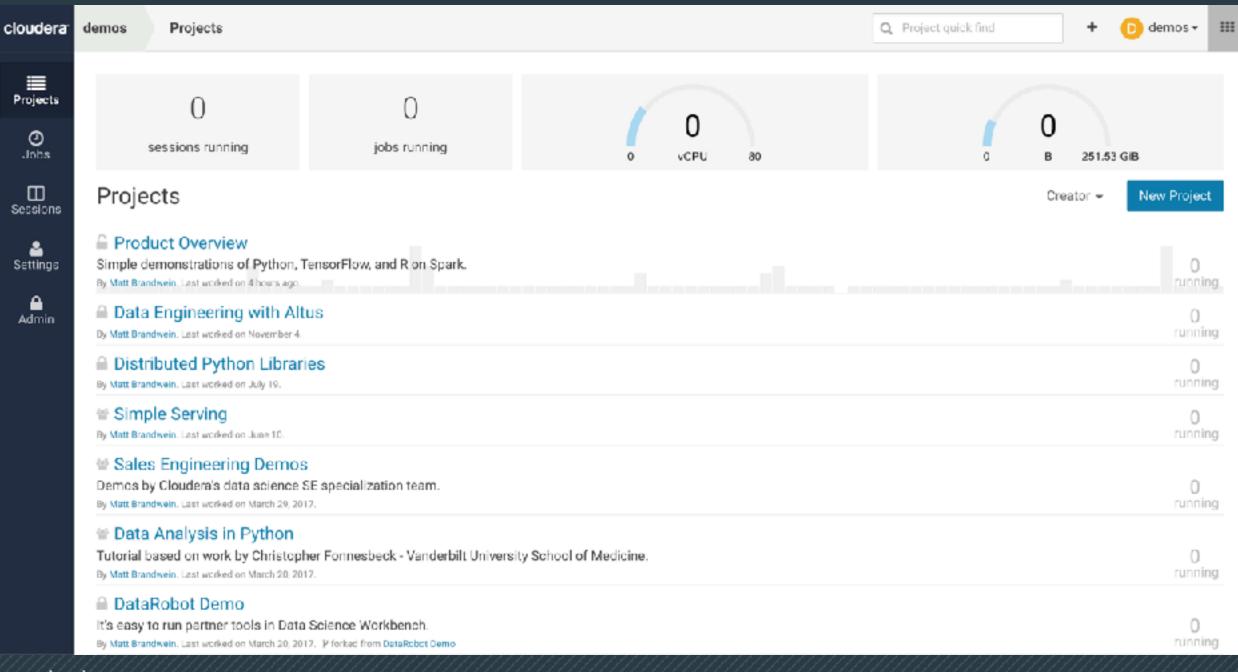


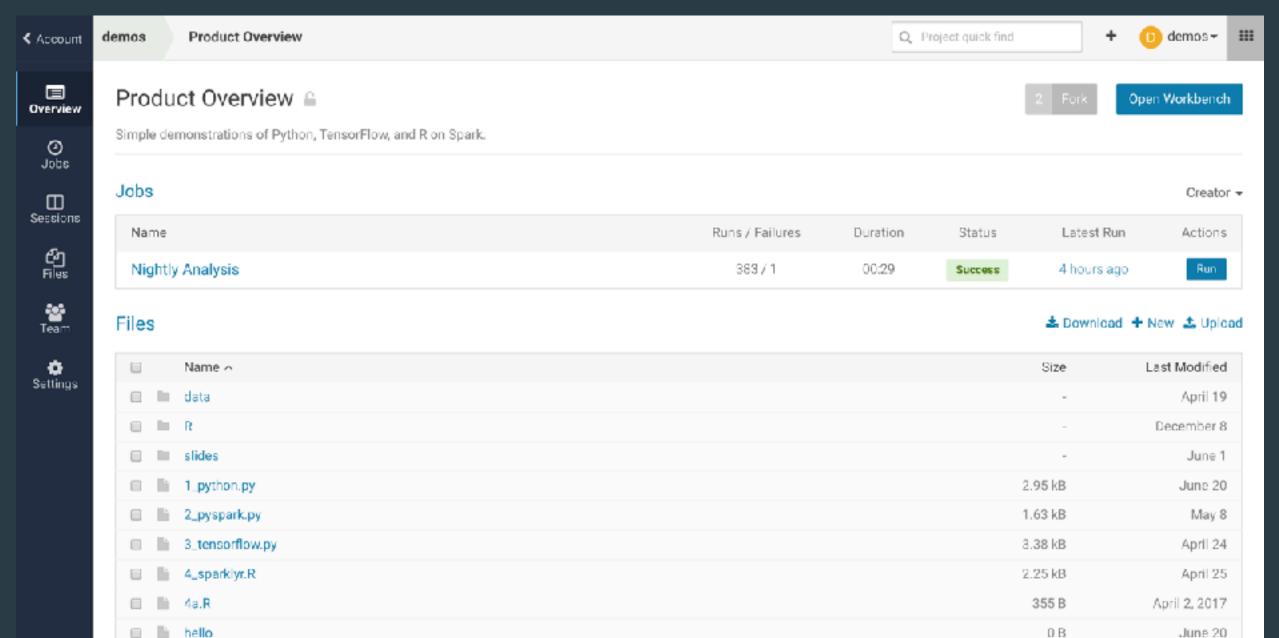
### MODEL MANAGEMENT

### View, test, monitor, and update models by team or project









1\_python.py

2\_pyspark.py

3\_tensorflow.py

4\_sparklyr.R

45.R

▼ deta.

GoogleTrendsData.csv

kmeans\_data.txt

MNIST

hello

⊩ R

README.md

⊩ slides

utils.py

```
# Goodle Stock Analytics
   # -----
   # This notebook implements a strategy that uses Google
 5 # trade the Dow Jones Industrial Average.
   import pandas as pd
   import matplotlib.pyplot as olt
   import matplotlib as mpl
10 from pandas_highcharts.display import display_charts
11 import seaborn
12 mpl.rcParams['font.family'] = 'Source Sans Pro'
   mpl.rcParams['axes.labelsize'] = '15'
14
15 # Import Data
16 # ========
17 #
18 # Load data from Google Trends.
19
   data = pd.read_csv('data/GoogleTrendsData.csv', index_
   data.head()
22
   # Show DJIA vs. debt related query volume.
   display_charts(data, chart_type="stock", title="DJIA v
   seaborn.lmplot("debt", "djia", data=data, size=7)
25
27 # Detect if search volume is increasing or decreasing
28 # any given week by forming a moving average and testil
29 # crosses the moving average of the past 3 weeks.
39 #
31 # Let's first compute the moving average.
32
33
   data['debt_mavg'] = data.debt.rolling(window=3, center)
   data.head()
34
35
35 # Since we want to see if the current value is above the
   # *preceeding* weeks, we have to shift the moving aver-
38
39 data['debt_mavq'] = data.debt_mavq.shift(1)
   data.head()
```

#### Start New Session

#### Engine Image - Configure

Base Image v1 - docker.repository.cloudera.com/cdsw/engine:1

#### Select Engine Kernel

- Python 2
- Python 3
- Scala
- R

#### Select Engine Profile

1 vCPU / 2 GiB Memory

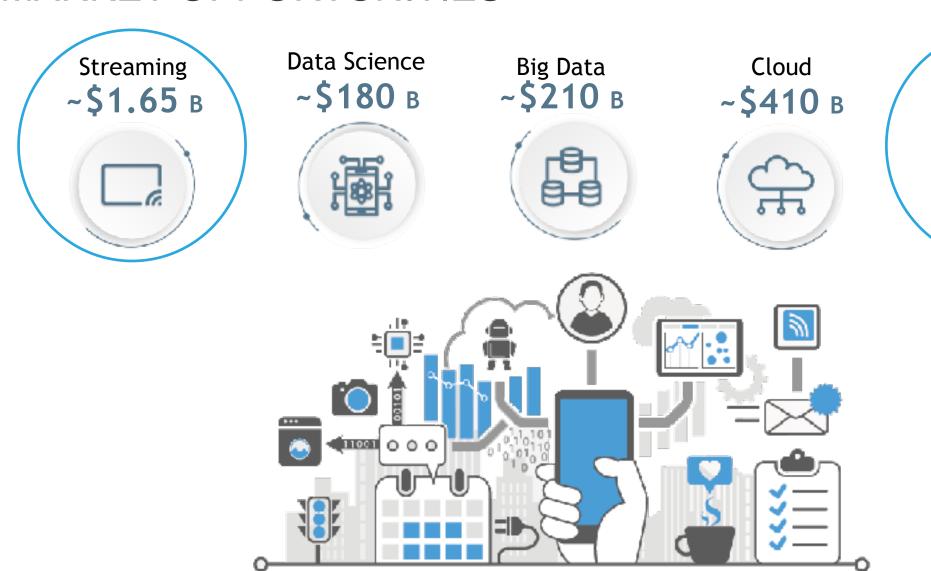
Launch Session

41

1\_bython.p... Edit View Navigate Run 2\_pyspark.py 1\_python.py # Google Stock Analytics My Python Session @ \_\_\_\_ 4\_sparklyr.R By Matt Brandwein - Python 2 Session - 1 vCPU / 2 GiB Memory -Running # This notebook implements a strategy that uses Google just now 5 # trade the Dow Jones Industrial Average. Product Overview 2 Cloudera Deta Science Workbench Terminal import pandas as pd ① tty-1gp7dhoexb4sn8fr.cdsw.edh.cloudera.com/cky3o8ggcyi0n49p/ import matplotlib.pyplot as plt GOOGLE Welcome to Cloudera Data Science Workbench 2\_pyspark.py import matplotlib as mpl 10 from pandas\_highcharts.display import display\_charts 3\_tensorflow.py This notebook Kernel: python2 import seaborn 4\_sparklyr.R 12 mpl.rcParams['font.family'] = 'Source Sans Pro' Industrial Aver mpl.rcParams['axes.labelsize'] = '16' 4a.R Project workspace: /home/cdsw 14 ▼ deta > import par > import mat
Kerberos principal: mbrandwein@CLOUDERA.LOCAL 15 # Import Data GoogleTrendsData.csv 4 ========= 17 > import matRuntimes: kmeans\_data.txt # Load data from Google Trends. R: R version 3.3.0 (2016-05-03) — "Supposedly Educational" > from panda 19 ▶ MNIST Python 2: Python 2.7.11 data = pd.read\_csv('data/GoogleTrendsData.csv', index\_d hello Python 3: Python 3.6.1 > import sea data.head() ⊪ R 22 Java: java version "1.8.0\_111" > mpl.rcPana # Show DJIA vs. debt related query volume. README.md > mpl.rcPara Git origin: http://github.mtv.cloudera.com/mbrandwein/cdsw-demo-sh display\_charts(data, chart\_type="stock", title="DJIA v seaborn.lmplot("debt", "djia", data=data, size=7) ⊪ slides 25 mport Dcdsw@1qp7dhoexb4snBfr:~\$ ls -al # Detect if search volume is increasing or decreasing total 96 # any given week by forming a moving average and testing utils.pyc drwxr-xr-x 14 cdsw cdsw 4096 Jul 14 - 2017 . 29 # crosses the moving average of the past 3 weeks. Load data fro 39 # 31 # Let's first compute the moving average. > data = pd.read\_csv('data/GoogleTrendsData.csv', index\_col='Date', parse\_dates 32 data['debt\_mavg'] = data.debt.rolling(window=3, center: > data.head() 34 data.head() djia debt 35 36 # Since we want to see if the current value is above the Date # \*preceeding\* weeks, we have to shift the moving aver: 38 2004-01-14 | 10485.18 | 0.210000 data['debt\_mavg'] = data.debt\_mavg.shift(1) data.head() 2004-01-22 | 10528.66 | 0.210000 41 # Generate Orders 43 # -----

# **CLOUDERA DATA FLOW**

### MARKET OPPORTUNITIES



IoT

# **IOT MARKET**

24.9B	By 2024 more than 24.9 Billion loT connections will be established
\$70B	An estimated \$70 billion will be spent by global manufacturers on IoT solutions in 2020
646M	An estimated 646 million healthcare devices (excluding fitness trackers and wearable devices) will be connected by 2020
78%	An estimated 78% of cars shipped globally will be built with hardware that connects to the internet by 2020
50%	50% of decision-makers in IT, services, utilities, and manufacturing have either deployed IoT, or will deploy it in the next 12-24 months

### PROBLEMS IN THE MARKET – PAIN THE CUSTOMER EXPERIENCES



Data movement



Continuous data ingestion



Streaming ETL



Streaming analytics

### COMMON USE CASES

#### **Data Movement**

Optimize resource utilization by moving data between data centers or between on-premises infrastructure and cloud infrastructure

### **Optimize Log Collection & Analysis**

Optimize log analytics solutions by using CDF as a single platform to collect and deliver multiple data sources

### **Gain key insights with Streaming Analytics**

Accelerate big data ROI by analyzing streaming data for patterns, comparing with ML models and delivering actionable intelligence

### Single view / 360° view of customer

Ingest, transform and combine customer data from multiple sources into a single data view / lake

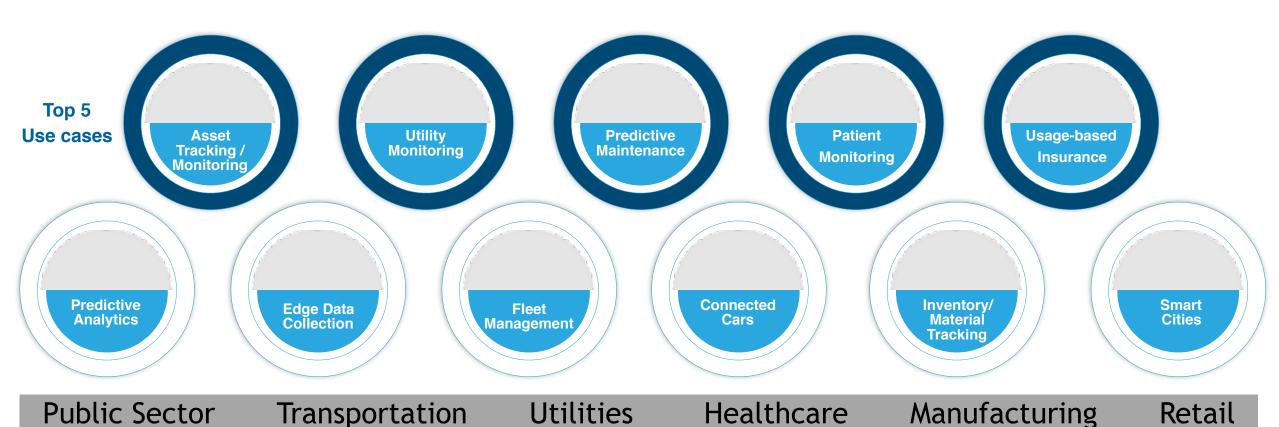
### **Stream Processing**

Combine multiple streams of data in realtime, enrich the data and route it to different end points based on rules

### **Capture and Analyze IoT Data**

Ingest sensor data from IoT devices and stream it for further processing and comprehensive analysis

### COMMON IOT USE CASES BY INDUSTRY



- IoT is a \$1.13T market opportunity in 2021.
- Americas \$329B IoT spending. Manufacturing and Transportation are top industries, accounting for 26% of total spending.
- APAC \$500B IoT spending. Manufacturing, Utilities and Transportation are top industries.
- EMEA \$264B IoT spending. Manufacturing is top industry, powered by Industry 4.0 initiatives.
- Worldwide IoT Analytics and Information Management Market = \$573M

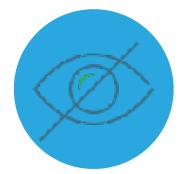
### KEY CUSTOMER CHALLENGES



**Data Ingestion:** High-volume streaming sources, multiple message formats, diverse protocols and multi-vendor devices creates data ingestion challenges



**Real-time Insights**: Analyzing continuous and rapid inflow (velocity) of streaming data at high volumes creates major challenges for gaining real-time insights



Visibility: Lack visibility of end-to-end streaming data flows, inability to troubleshoot bottlenecks, consumption patterns etc.

### **CLOUDERA DATAFLOW**

# Cloudera DataFlow Data-in-Motion Platform











#### **ENTERPRISE SERVICES**

Provisioning, Management and Monitoring Unified Security

Edge-to-Enterprise Governance

Single Sign-on

# WHAT IS CLOUDERA DATAFLOW (CDF)?

Cloudera DataFlow (CDF) is a scalable, real-time streaming data platform that collects, curates, and analyzes data so customers gain key insights for immediate actionable intelligence.



### HISTORY OF CDF

#### **Data-in-Motion:**

- Comprehensive real-time streaming data platform
- Manage data-in-motion from edge-toenterprise
- Power IoT-scale streaming architectures

#### Cloudera DataFlow Data-in-Motion Platform











**Unified Security** Edge-to-Enterprise Governance Single Sign-on

Bring this to the edge with connected platforms

**Enable next generation** Modern Data Architecture

Mid-2000's NiFi was developed and used at NSA

2015

Onyara is acquired HDF is born

2018

Strong Streaming Platform

- Support for Kafka 2.0
- SMM is introduced

2019

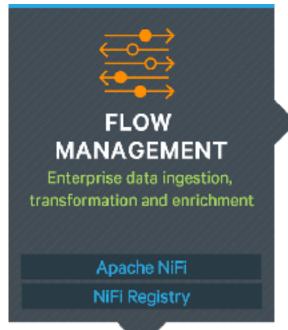
Cloudera merger Enable Edge Intelligence

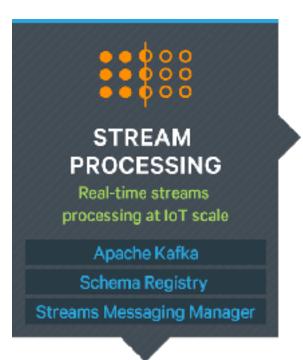
Tomorrow: **Edge-to-Al** 

# PRODUCT OVERVIEW

### **CLOUDERA DATAFLOW**











#### ENTERPRISE SERVICES

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

# **CLOUDERA DATAFLOW**



## WHAT IS CLOUDERA EDGE MANAGEMENT (CEM)?

Cloudera Edge Management (CEM) is an edge management solution made up of edge agents and an edge management hub. It manages, controls and monitors edge agents to collect data from edge devices and push intelligence back to the edge. CEM allows you to develop, deploy, run and monitor edge flow apps on thousands of edge devices.



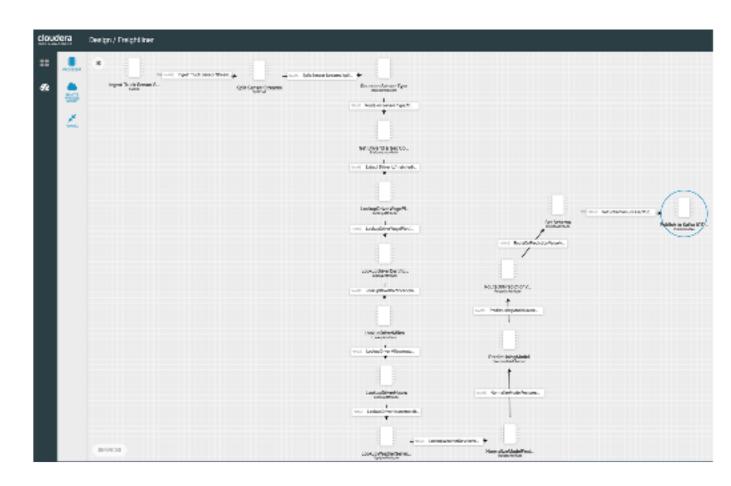
### **EDGE DATA MANAGEMENT**

- Edge data collection powered by Apache MiNiFi
- MiNiFi smaller footprint than NiFi
  - Guaranteed delivery
  - Data buffering
  - Prioritized queuing
  - Flow-specific QoS
  - Data provenance
  - Designed for extension
  - C++ / Java agents
  - TensorFlow support
- Designed for IoT



### **EDGE FLOW MANAGER**

- Edge management hub
- NiFi-like user interface to develop and deploy flow files to the edge
- Update and deploy ML model files to the edge agents
- Monitor thousands of edge agents
- Integration with NiFi Registry







## WHAT IS CLOUDERA FLOW MANAGEMENT (CFM)?

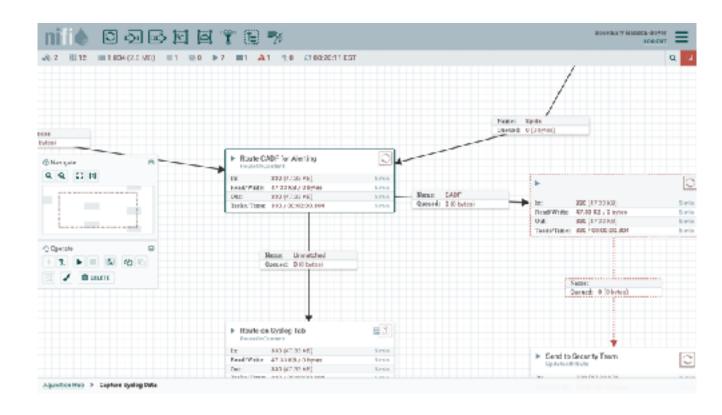
Cloudera Flow Management (CFM) is a no-code data ingestion and management solution powered by Apache NiFi. With NiFi's intuitive graphical interface and 300+ processors, CFM delivers highly scalable data movement, transformation and management capabilities to the enterprise. CFM also enables DevOps type development and deployment with its support for NiFi Registry.



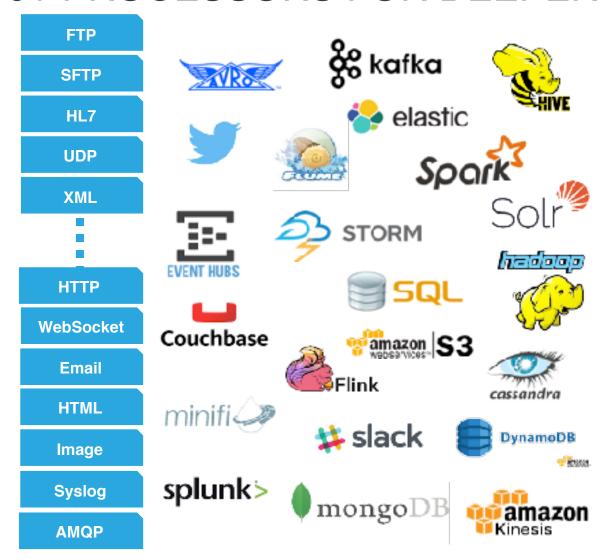
#### FLOW MANAGEMENT



- Web-based user interface
- Highly configurable
- Out-of-the-box data provenance
- Designed for extensibility
- Secure
- NiFi Registry
  - DevOps support
  - FDLC
  - Versioning
  - Deployment



## 300+ PROCESSORS FOR DEEPER ECOSYSTEM INTEGRATION



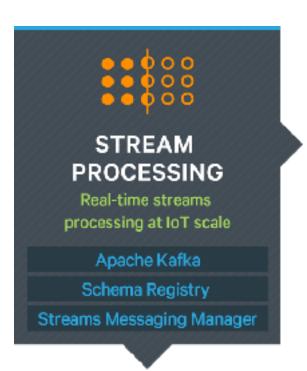
Hash	Encrypt	GeoEnrich
Merge	Tail	Scan
Extract	Evaluate	Replace
Duplicate	Execute	Translate
Split	Fetch	Convert

Route Text	Distribute Load
Route Content	Generate Table Fetch
Route Context	Jolt Transform JSON
Control Rate	Prioritized Delivery

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

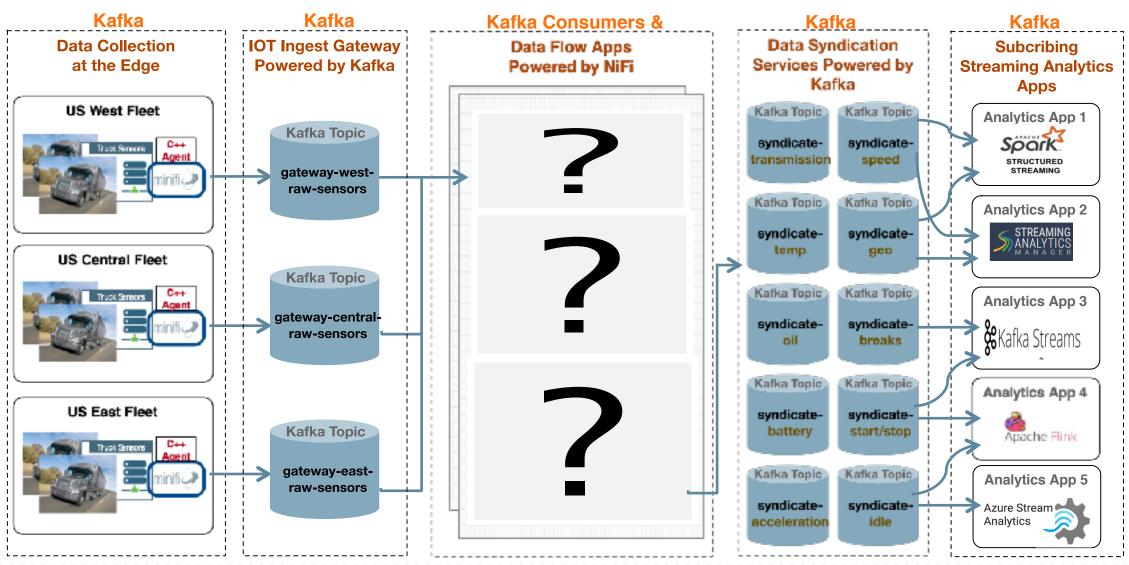






#### **Streaming Analytics Reference Architecture**

#### Kafka is Everywhere. Critical Component of Streaming

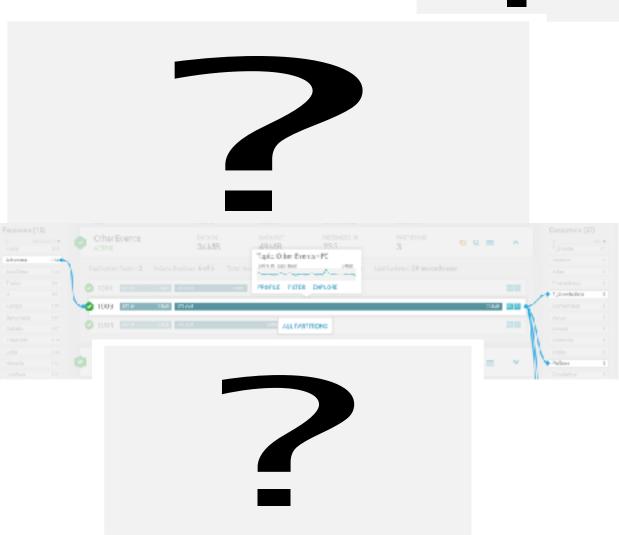


#### Cloudera Streams Messaging Manager (SMM)



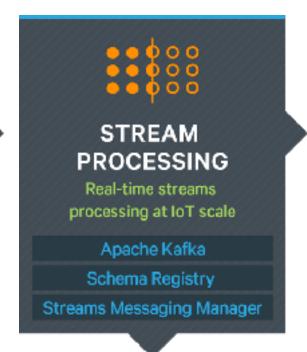
#### What is SMM?

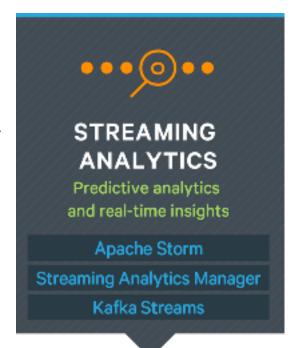
- Kafka Management and Monitoring tool
- Cure the "Kafka Blindness"
- Single Monitoring Dashboard for all your Kafka Clusters across 4 entities
  - Broker
  - Producer
  - Topic
  - Consumer
- REST as a First Class Citizen











## STREAMING ANALYTICS

- Pattern matching
- Predictive and Prescriptive Analytics
- Complex Event Processing
- Continuous & Real-time Insights



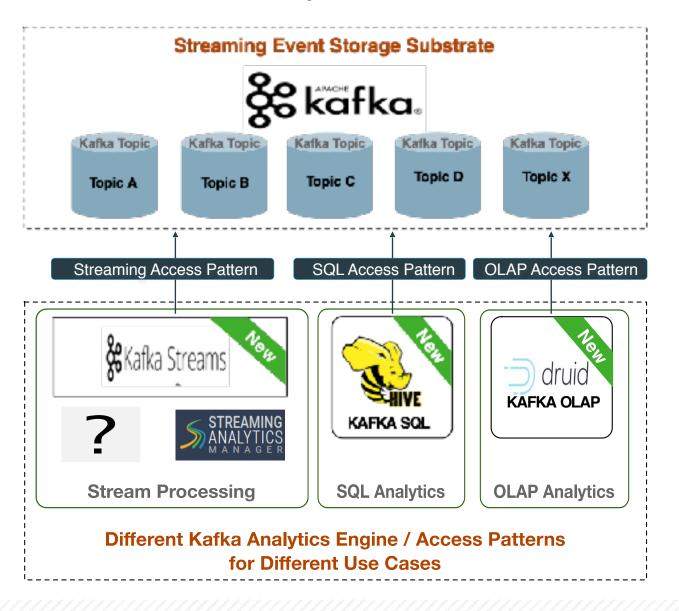






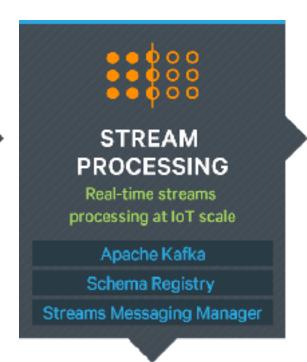


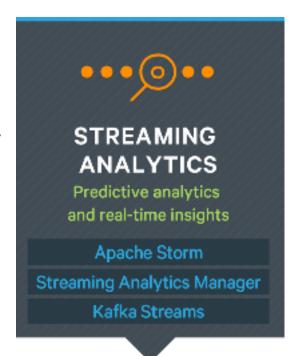
#### 3 Kafka Analytics Access Patterns













#### ENTERPRISE SERVICES

Provisioning, Management and Monitoring Unified Security

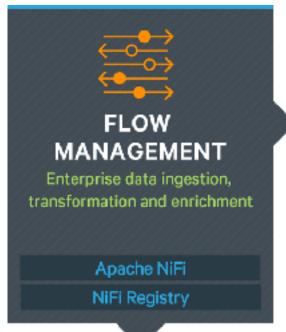
Edge-to-Enterprise Governance

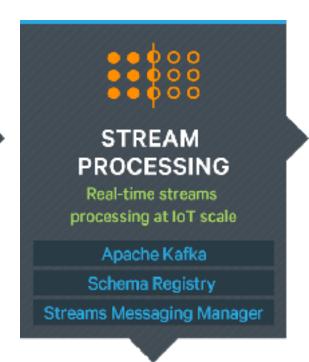
Single Sign-on

## **ENTERPRISE SERVICES**

- Provisioning
- Management
- Monitoring
- Unified Security
- Single Sign-on
- Audit
- Compliance
- Edge-to-Enterprise Governance











#### ENTERPRISE SERVICES

Provisioning, Management and Monitoring Unified Security

Edge-to-Enterprise Governance

Single Sign-on



#### **KEY DIFFERENTIATORS**

**100% open source technology** – Only vendor with this strategy; prevents vendor lock-in



**300+ pre-built processors** – Only product to offer such comprehensive connectivity from edge to enterprise



**3 Streaming analytics engines** – Only vendor to offer a choice of three streaming analytics engines to customers for all their streaming architecture needs



**Built-in data provenance** – Only product in the market to offer out-of-the-box data provenance on data-in-motion



**Comprehensive streaming platform** – Only big data vendor to offer a comprehensive streaming platform from real-time data ingestion, transformation, routing to descriptive, prescriptive and predictive analytics.



# THANK YOU

cloudera