

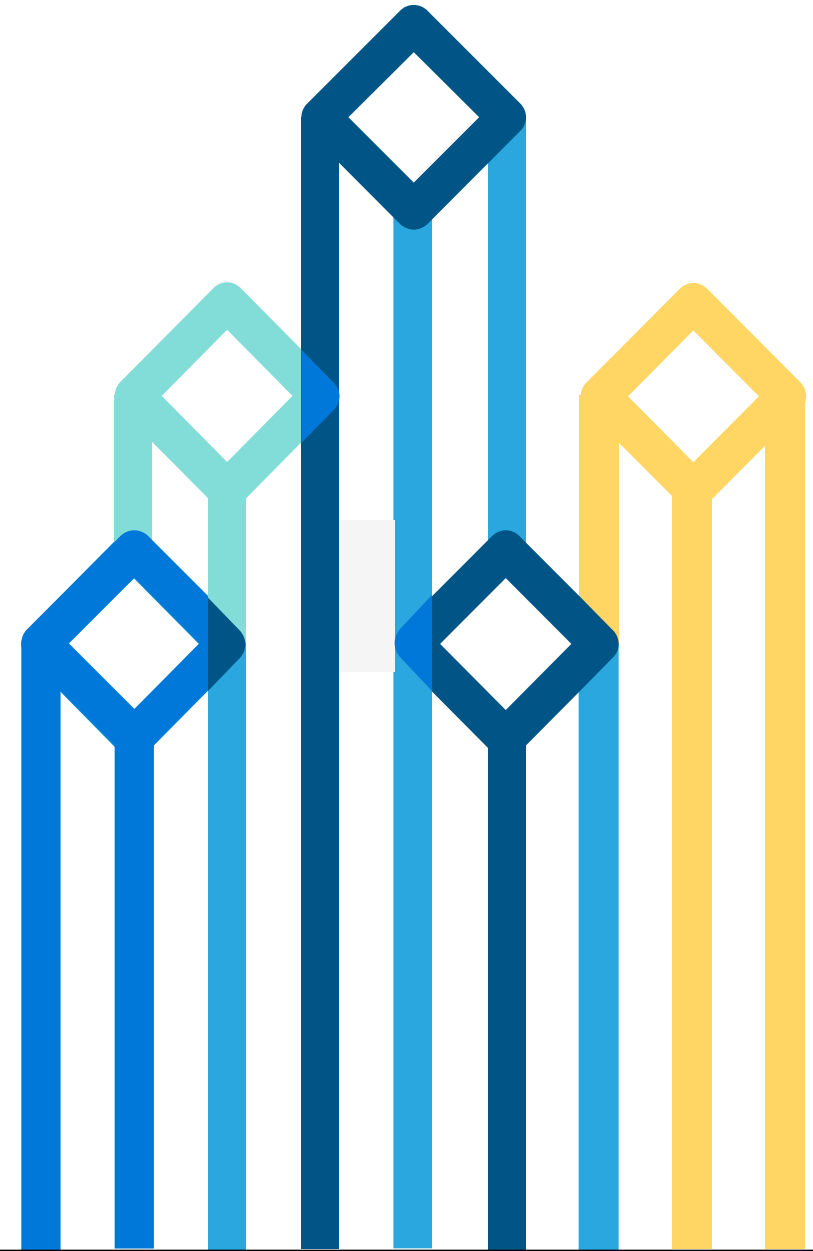


# Building the Enterprise Data Lake with Cloudera & Cisco

Prepared by :

Marilyn Tan, Country Manager Singapore

Xue Daming, Senior Systems Engineer



# Digital Transformation with Data

# DATA is Transforming the World!



## CONNECTED WORLD



### INTERNET OF THINGS INDUSTRY 4.0

- Smart Things / Devices
- New Use Cases
- New Architectures
- By 2020: 20.8b devices
- IoT: \$1.7trillion in 2020



## MORE DATA



### DATA AS 4<sup>th</sup> PRODUCTION FACTOR

- New Analytics
- Machine Learning & AI
- New data sources
- Data Virtualization
- Data Science



## SMART APPs



### THE NEW UX

- Connected Experience
- Ubiquitous computing: Everywhere & on every device (Voice, VR, AR, mobile, Wearables)



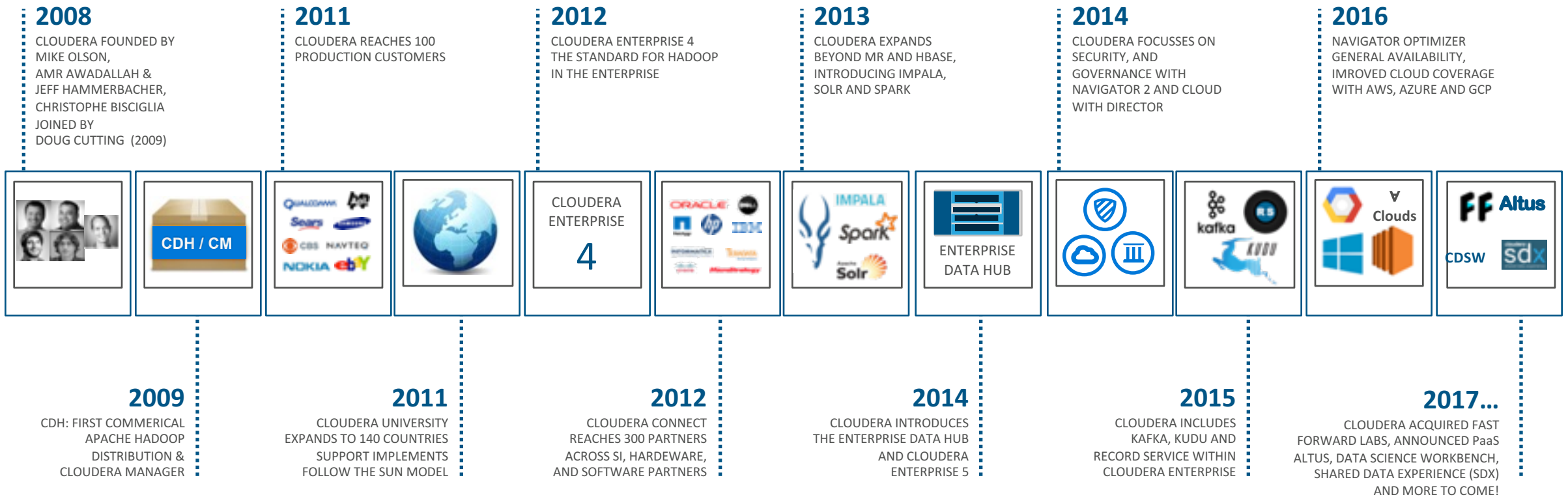
## DIGITAL DEMOCRACY & SECURITY



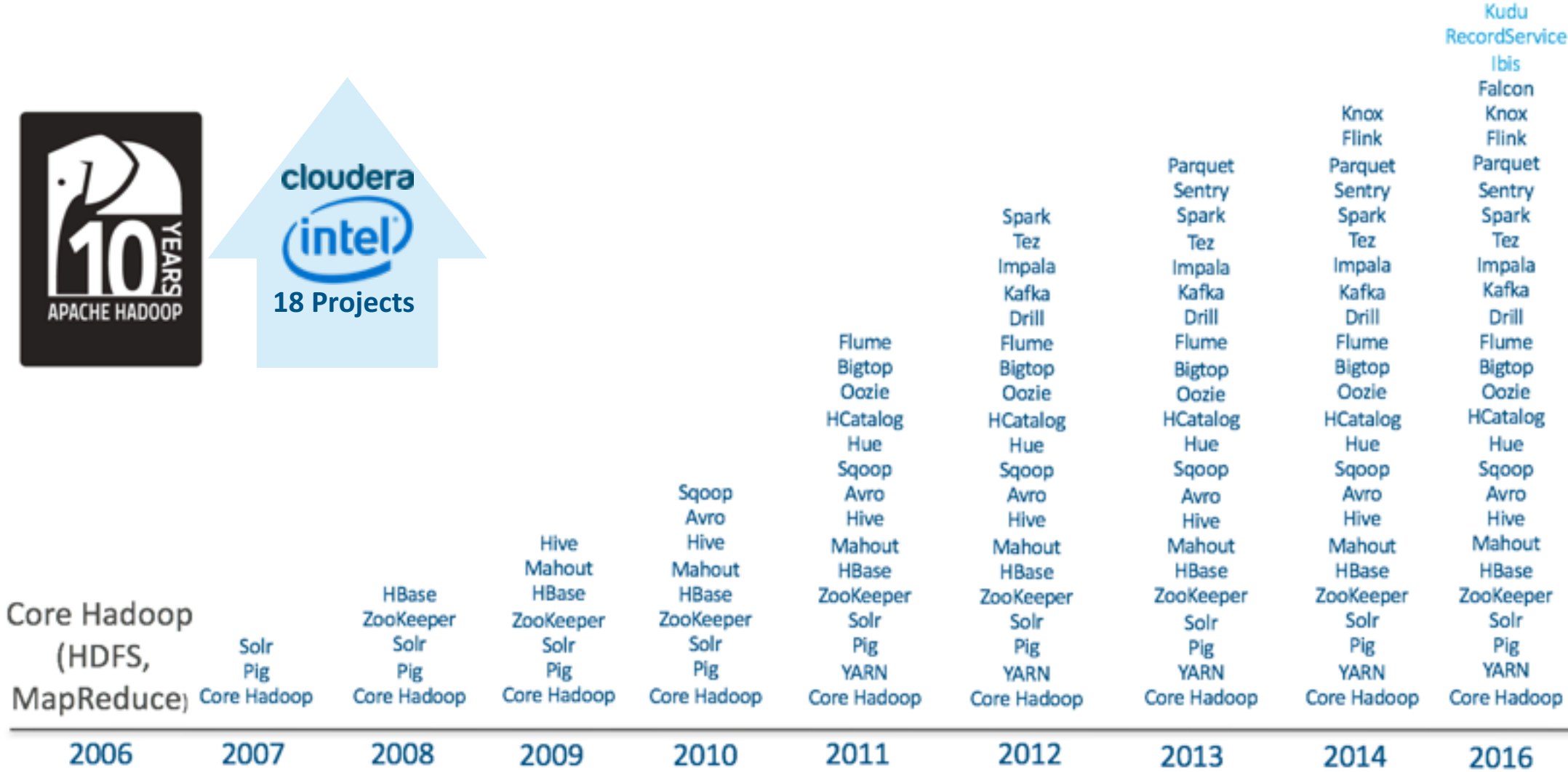
### RISE OF OPEN SOURCE

- Digital Sharing Economy: Open Data & Algorithms
- Enterprise ready Open Source (e.g. Apache)
- Digital (distributed) Trust (esp. Blockchain)

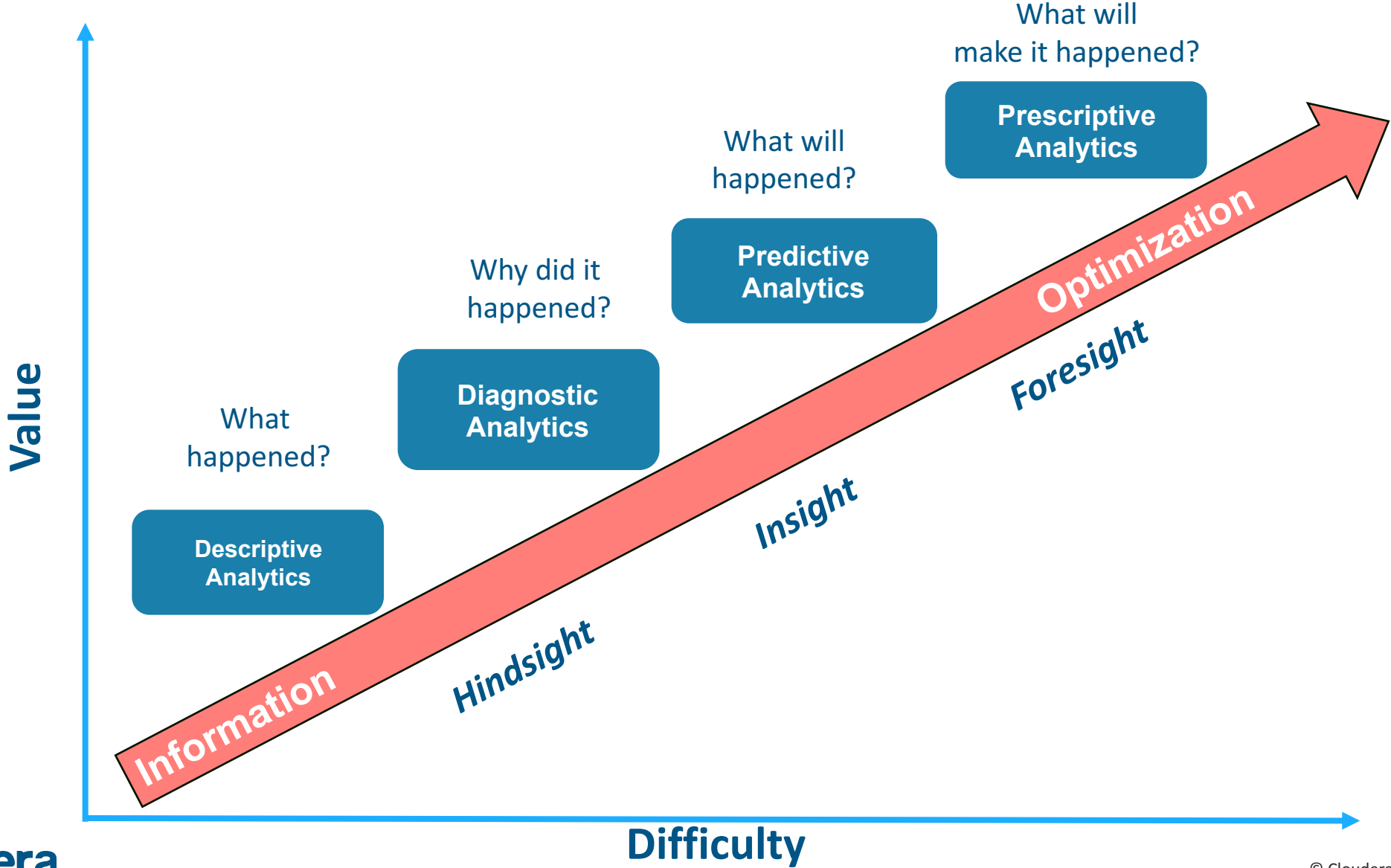
# The 9 year Cloudera journey ...



# What Happen Next: A Decade of Hadoop

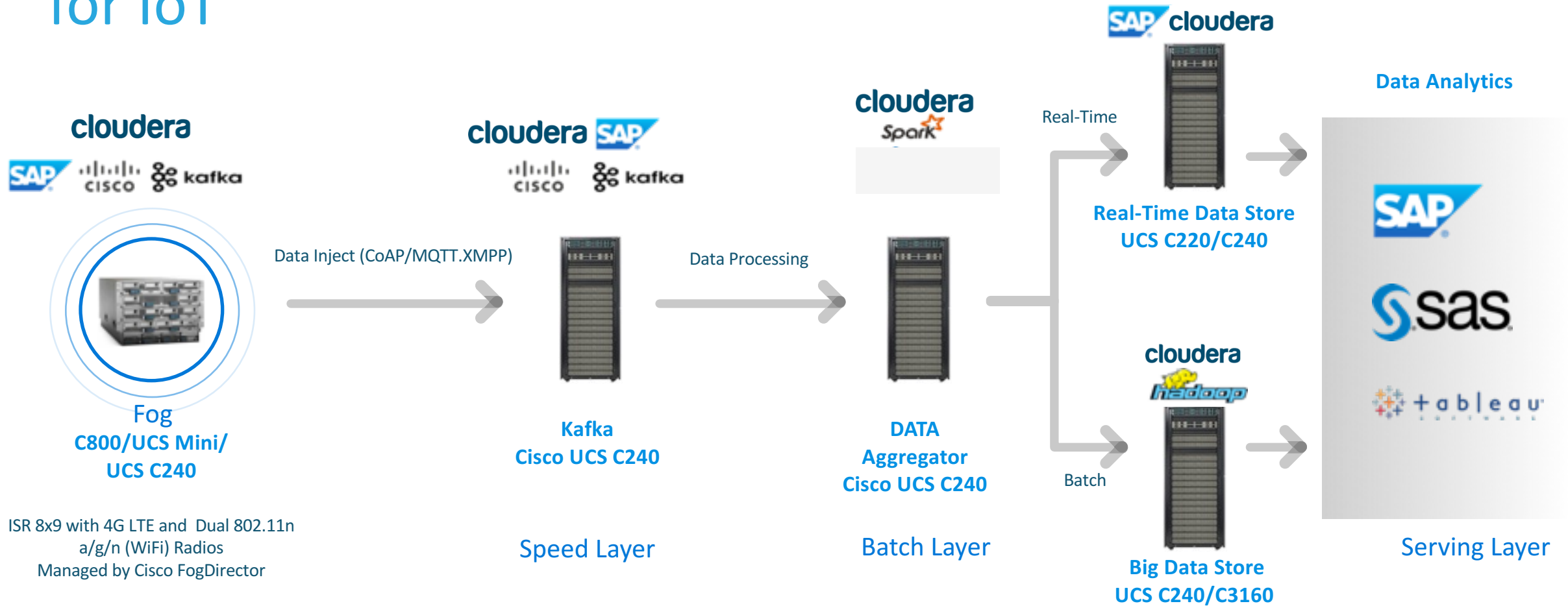


# Gartner Analytics Ascendancy Model



# Cloudera & Cisco Enterprise Data Lake Innovation

# Cisco UCS Integrated Infrastructure with Cloudera for IoT



Cisco UCS at all layers, fully validated architectures with all major players



# Fabric Centric Design

## High Performance

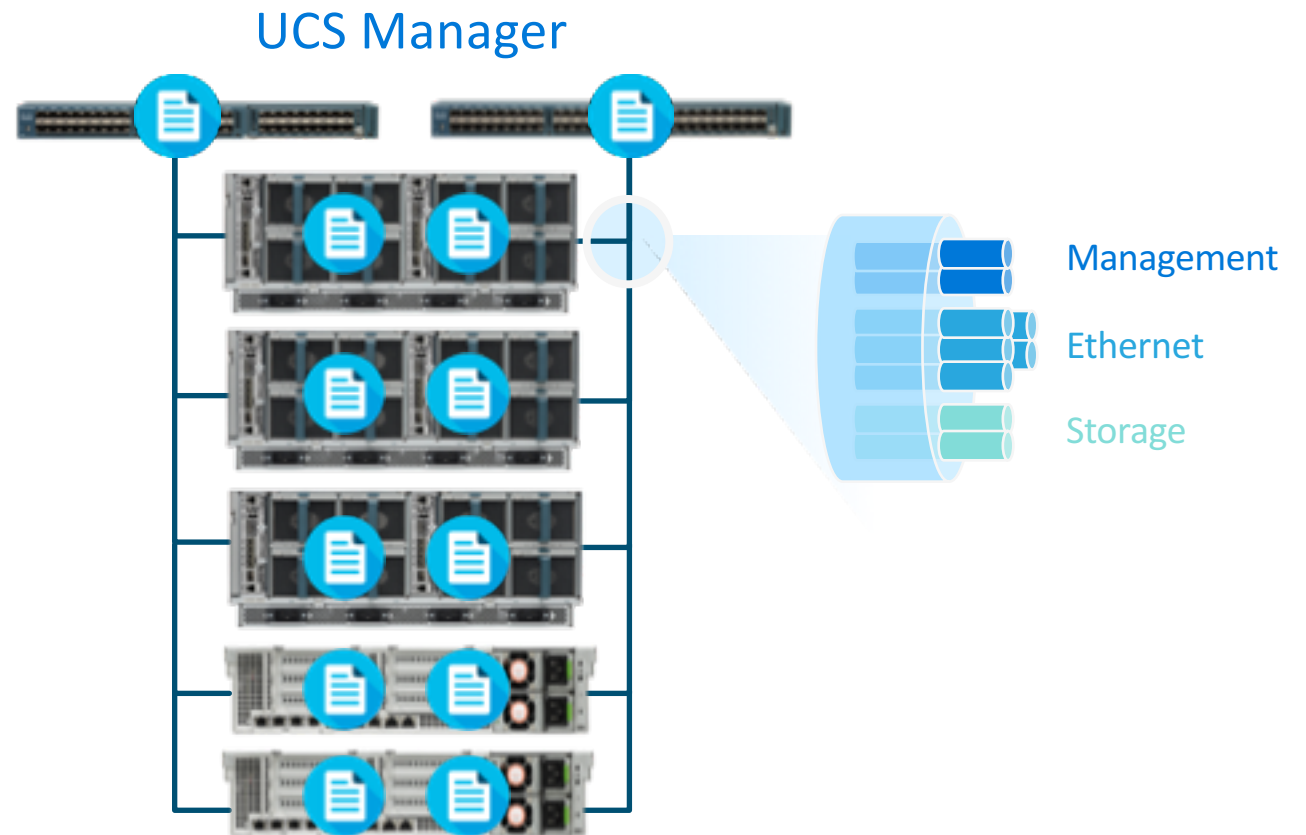
40 GB/s Ethernet; 320 GB/s per Chassis

## Unified Fabric

Single Cable for Network, Storage, and Management Traffic

## Easy to Scale

Single Point of Management: Add Cables for Bandwidth vs. Fabric Type



# Management Simplicity

## Big Data: Management Consistency

Hundreds of Servers

Thousands of management points

## Simplified Scalability

Easily Scale your infrastructure from few servers to thousands of servers with a fully Integrated Infrastructure

## Centralized Management

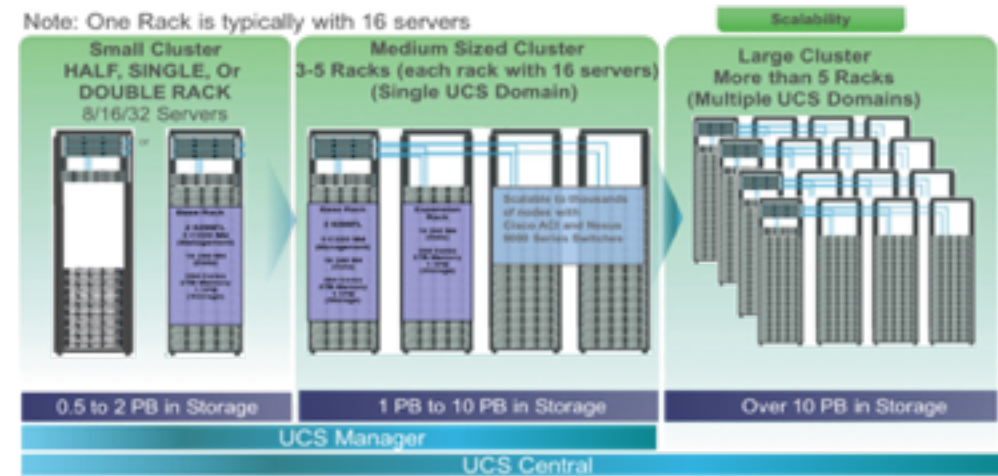
Service Profiles for Servers

- Manage all servers centrally

Application Profiles for Network

- Manage all network centrally

**cloudera**



UCS Service Profile



Cisco ACI Application Profile



# The enterprise platform for machine learning



## DRIVE CUSTOMER INSIGHTS

- Market segmentation
- Customer 360
- Next best offer
- Churn analysis & prevention

PATTERN RECOGNITION

DETECTION

500+

CUSTOMERS RUN ON

Spark cloudera

PREDICTION



CONNECT PRODUCTS & SERVICES (IoT)

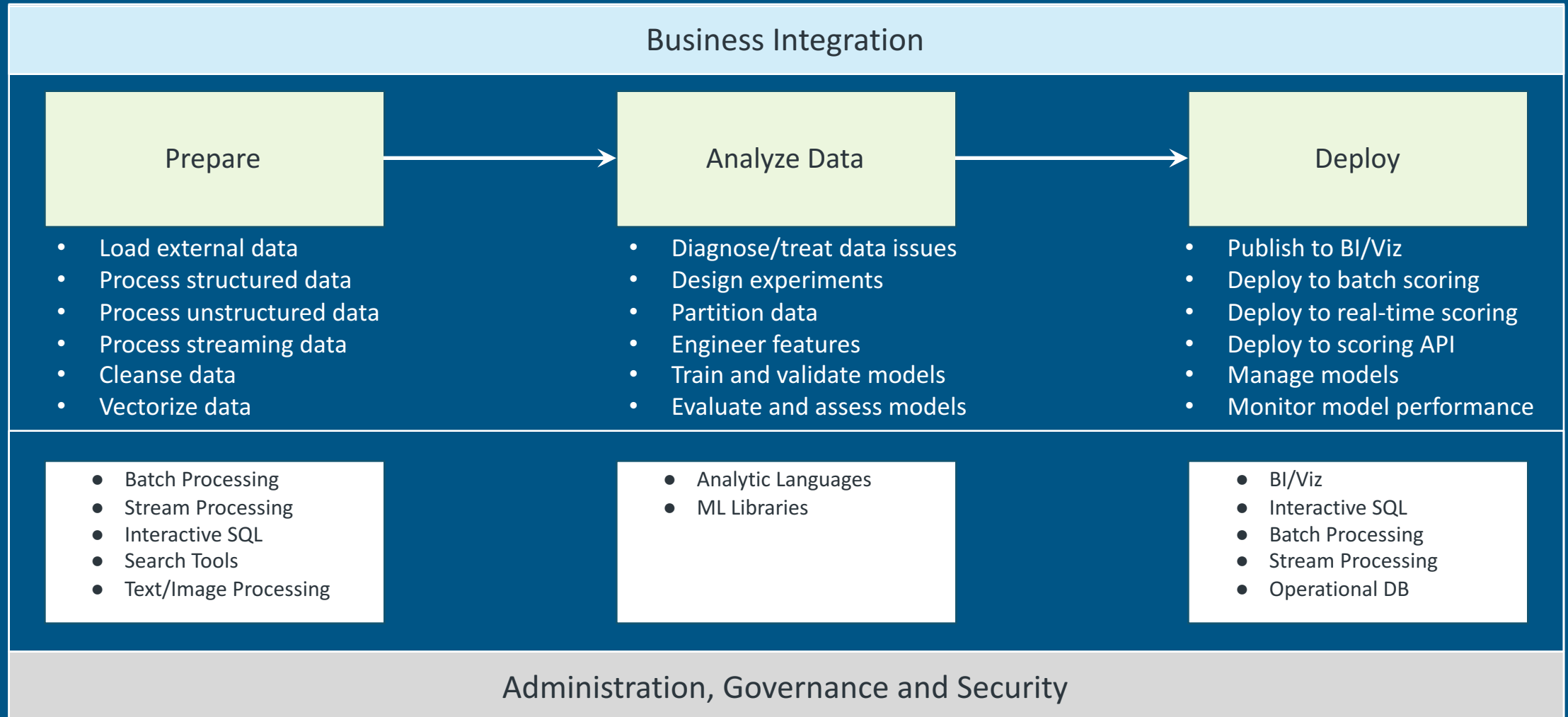
- Predictive maintenance
- Genomics & personalized medicine
- Predicting and preventing disease
- Natural language



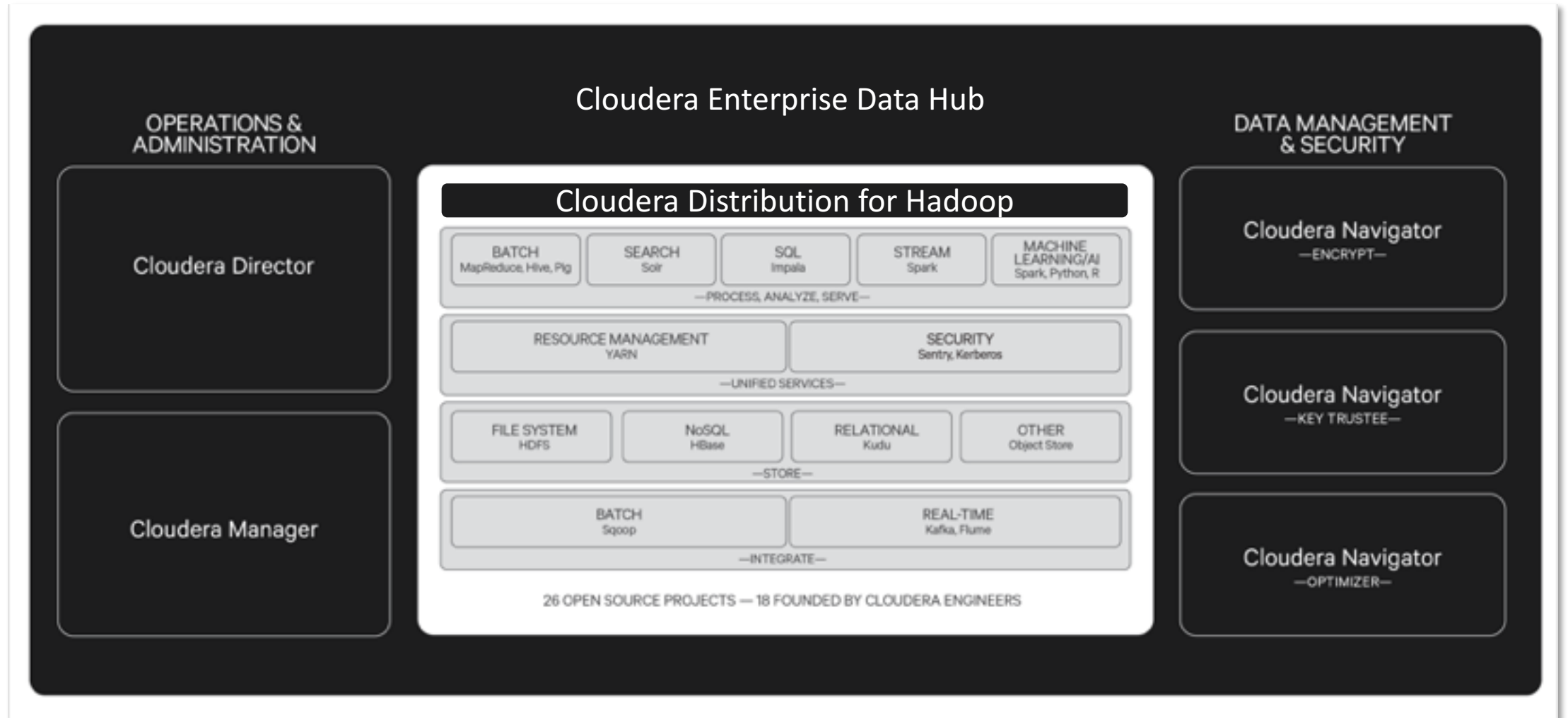
## PROTECT BUSINESS

- Cybersecurity
- Fraud
- Anti-money laundering
- Risk modeling & assessment
- SPAM detection

# Machine learning requires a complete stack.



# A complete, integrated enterprise platform



# Cloudera Data Science Workbench

The screenshot displays the Cloudera Data Science Workbench interface. At the top, there are four gauges showing system metrics: 0 sessions running, 2 jobs running, 3 vCPU usage, and 6 GB memory usage. Below these are project management options like 'New Project'. The main area is divided into a code editor on the left with Python code for data analysis and a visualization area on the right showing a table of 'dja' and 'debt' data and a corresponding line chart titled 'DJA vs. Debt Query Volume'. A 'Cluster Metadata' window is also visible in the bottom left corner.

- Supports data science end-to-end
- Full access to data
- Secure self-service provisioning
- Containerized environments
- Supports Python, R, and Scala
- Automates:
  - Workflow
  - Version control
  - Collaboration
  - Sharing

# CDSW Benefits

## Data Scientists

- Web browser, no desktop footprint
- Use R, Python, or Scala
- Install any library or framework
- Isolated project environments
- Direct access to data in secure clusters
- Share insights with team
- Reproducible, collaborative research
- Automate and monitor data pipelines
- Built-in job scheduling

## IT

- Support self-service data science
- Full platform security
- Kerberos authentication
- Run on-premises or in the cloud

# Deep learning in Cloudera with Apache Spark

## Spark Packages

- Two packages:
  - CaffeOnSpark
  - TensorFlowOnSpark
- Developed by Yahoo
- Python and Scala APIs
- All DL architectures
- Integrated pipeline
- Run on existing clusters
- Training and inference

## DL4J

- Open source DL library
- Developed by Skymind
- Built on JVMs
- Supports CPUs and GPUs
- Java, Scala, Python APIs
- Training and inference
- Imports models from:
  - TensorFlow
  - Caffe
  - Torch
  - Theano
- Runs on existing clusters

## BigDL

- Deep learning framework
- Developed by Intel
- Supports CPUs only
- Leverages Intel MKL
- Scala, Python APIs
- Imports models from:
  - TensorFlow
  - Caffe
  - Torch
- Runs on existing clusters

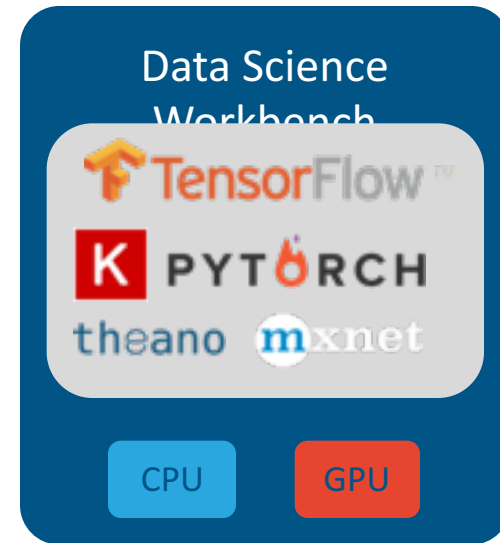


# New! Accelerated deep learning on-demand with GPUs

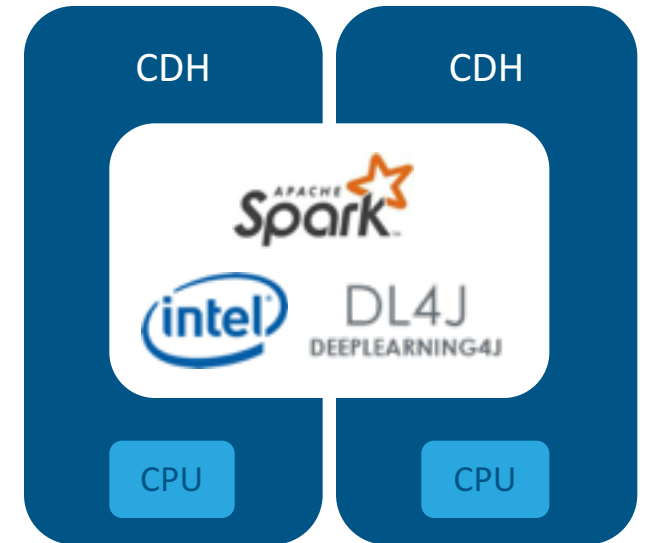
“Our data scientists want GPUs, but we can’t find a way to deliver multi-tenancy. If they go to the cloud on their own, it’s expensive and we lose governance.”

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

## Multi-tenant GPU support on-premises or cloud



single-node training



distributed training, scoring

# Enterprise Data Lake Architecture

# Canonical Ingestion & Spark Streaming Analytics with Cisco Big Data Analytics Platform

- Integrate with Apache Spark Streaming for **real-time analysis** of data
- Write back to Kafka for further processing or to send to an application layer

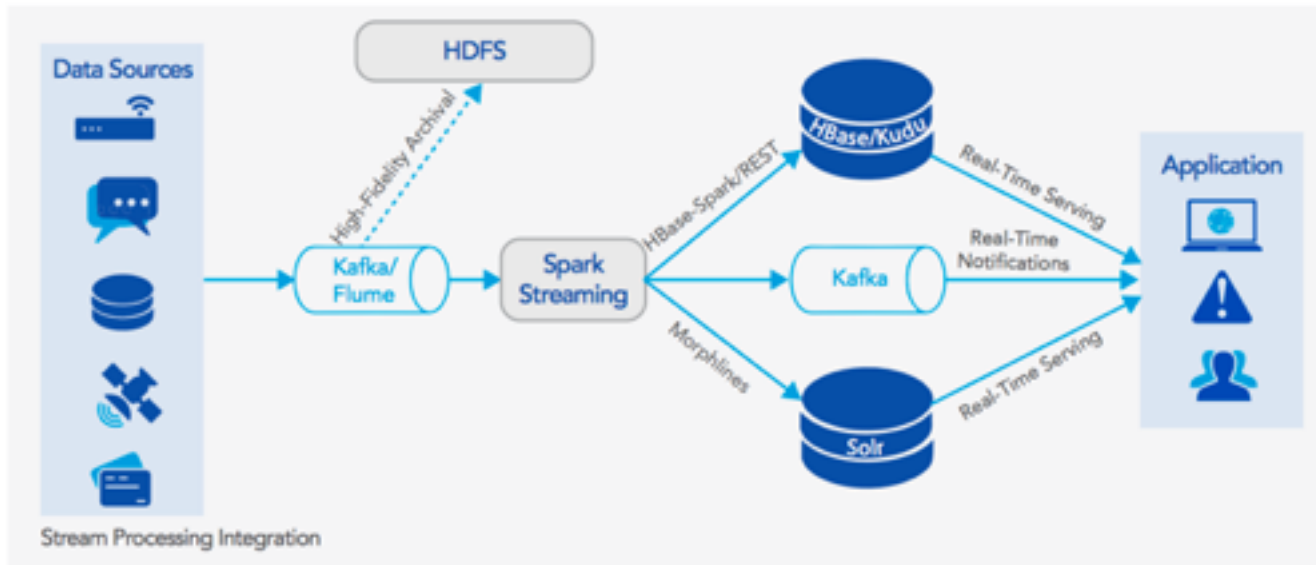
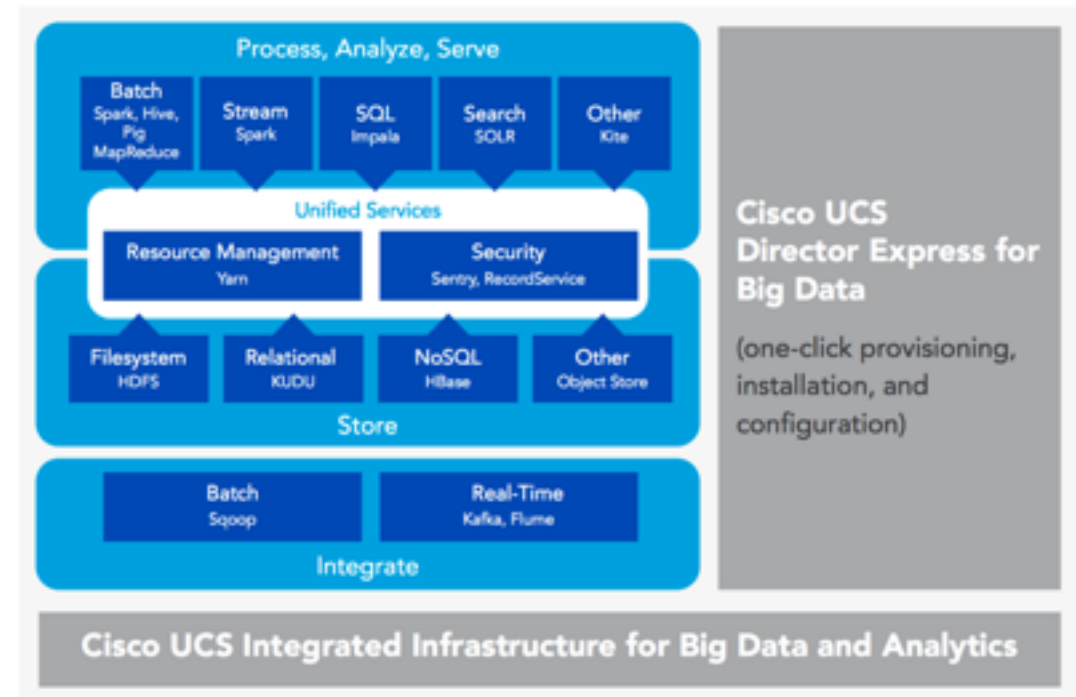
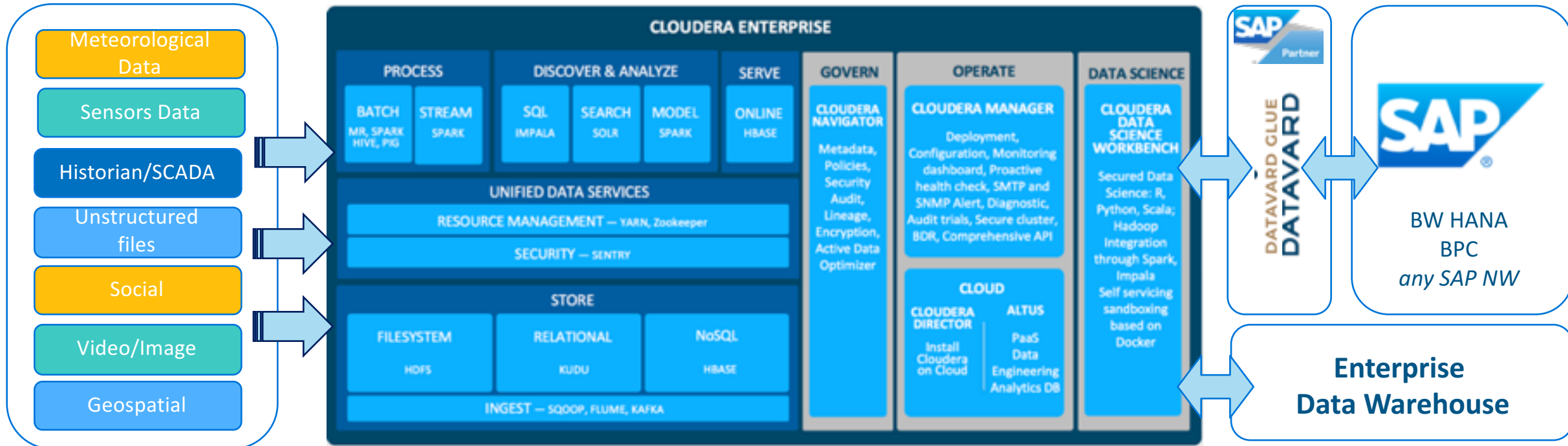


Figure 1. Typical Architecture with Cloudera Enterprise Data Hub



# Proposed Architecture for Enterprise Data Platform



# Big Data Blueprints: Cisco Validated Designs

## Designs Big Data



### Cisco® Validated Designs with Cloudera

Data Center Designs Cloud Computing  
**Design Zone for Big Data**

HOME  
SOLUTIONS  
ENTERPRISE  
PROGRAMS FOR ENTERPRISE  
DESIGN ZONE  
DESIGN ZONE FOR DATA CENTERS  
DATA CENTER DESIGNS: CLOUD COMPUTING  
**Design Zone for Big Data**

**Unlock the Value in Big Data**  
Learn the power of the Hadoop Data Platform and Cisco UCS.  
[View Design](#)

**Let Us Help**  
Locate International Contacts  
Get Technical Support  
Find a Reseller in Your Area  
Manage Your E-mail Preferences

Share

**You Might Also Like**  
Cisco Big Data Portal  
Cisco UCS Integrated Infrastructure for Big Data (plug)  
Cisco UCS Solution Accelerator Packs for Big Data (PDF - 170 KB)   
Cisco Publishes First-Ever Industry Standard Benchmark Results for Big Data Systems (PDF - 162 KB)

**Build and Deploy Applications Quickly**  
The Cisco UCS Integrated Infrastructure for Big Data, the third generation of the Cisco UCS Common Platform Architecture for Big Data, integrates industry-leading computing, network, and management capabilities into a unified fabric-based architecture. Optimized for big data workloads, it is a highly efficient, scalable, high-performance solution. With it, organizations can unlock the intelligence in their data to help create a sustainable, competitive business advantage quickly and cost-effectively. The Cisco Validated Designs for Big Data listed below consist of systems and solutions that have been designed, tested, and documented to facilitate faster, more reliable, and more predictable customer deployments and significantly lower the cost of ownership.

**Generation 3**  
Hadoop as a Service on Bare Metal with UCS Director Express for Big Data CI with ACI  
Cisco UCS Integrated Infrastructure for Big Data with SAP HANA Vora for In-memory Analytics  
Cisco UCS Integrated Infrastructure for Big Data with IBM BigInsights  
Cisco UCS Integrated Infrastructure for Big Data with MapR with Optional Multi-Tenancy Extension (PDF - 10.9 MB)



## What you get

Industry-leading partnerships

Tested and validated reference architectures to meet performance, capacity, and scale

Joint engineering lab

Extensive options for data management (Hadoop, MPP, and NoSQL) to meet your business needs

Solution bundles optimized for cost of ownership and ease of ordering

Solution designed, tested, and documented to facilitate faster, more reliable, and more predictable customer deployments.

# Our Customers' Success Stories



DATA-DRIVEN  
PRODUCTS

CASE STUDY

## TRANSPORTATION

- » PREDICTIVE MAINTENANCE
- » IMPROVED SERVICE
- » DATA DRIVEN PRODUCTS

# NAVISTAR<sup>®</sup>

## Using Predictive Maintenance to Improve Performance and Reduce Fleet Downtime

- OnCommand Connection is collecting **telematics** and **geolocation** data across the fleet
- Reduced maintenance costs to **\$.03 per mile** from **\$.12-\$.15 per mile**
- Centralizing data from **13 systems** with varying frequency and semantic definitions
- **Real-time** visibility of **250,000+** trucks in order to **improve uptime** and vehicle performance

**cloudera**





DATA-DRIVEN  
PRODUCTS

CASE STUDY

## TRAVEL & TRANSPORTATION

- » SMART BUILDINGS
- » PREDICTIVE MAINTENANCE
- » ADVANCED ANALYTICS



## Smart Buildings - Preventative Maintenance

### Using Sensors & IoT to Improve Passenger Safety and Airport Efficiency

#### Challenge:

- Improve traveler satisfaction and safety, by reducing downtime for critical operational machinery

#### Solution:

- [Cloudera on Azure](#) to capture, secure, and correlate sensor (IoT) data collected from escalators, elevators, and baggage carousels
- Provide necessary fixes to [prevent unplanned downtime](#)

**cloudera**







## 2016 Data Impact Award Winner State of Kentucky Department of Transportation

### Smart Cities

Enabling the State of Kentucky manage snow and ice events in real time

#### Challenge:

- Needed more efficient approach to inclement weather road management

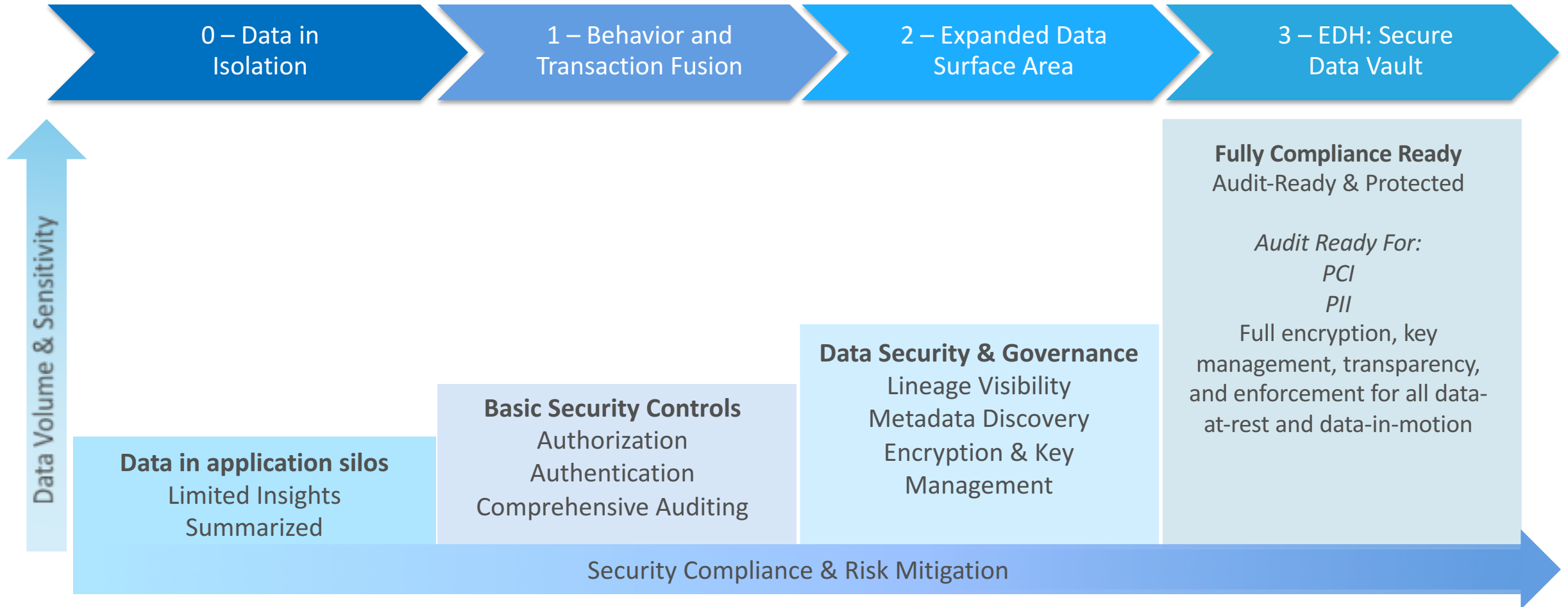
#### Solution:

- [Real-time weather response system](#) that incorporates real-time data from Waze, HERE, ESRI's GeoEvent processor, and Automatic Vehicle Locations (providing sensor data from salt trucks).
- KYTC aggregates 15-20 million records every day and process more than a [million records per second](#).

**cloudera**



# Data Protection & Governance



# Why Cloudera

## The **Platform** for Next-Generation Analytics

Cloudera Enterprise delivers the capabilities required by the largest enterprises, spanning analytics, security, governance, and management. We make Hadoop fast, easy, and secure.

## The **Experience** to Help You Succeed

**No one knows Hadoop like Cloudera.**

As the first Hadoop company, Cloudera is the world's leading contributor to and provider of enterprise Hadoop, with experience you can rely on to help you succeed.

## **Open Innovation**

Our unique hybrid open source strategy enables us to lead the enterprise expansion of the Hadoop ecosystem, driving innovative new capabilities and open standards in the community.



**cloudera**

Thank you

[marilyn@cloudera.com](mailto:marilyn@cloudera.com) | +65 9822 2338  
[daming@cloudera.com](mailto:daming@cloudera.com) | +65 9368 2316