

# cloudera

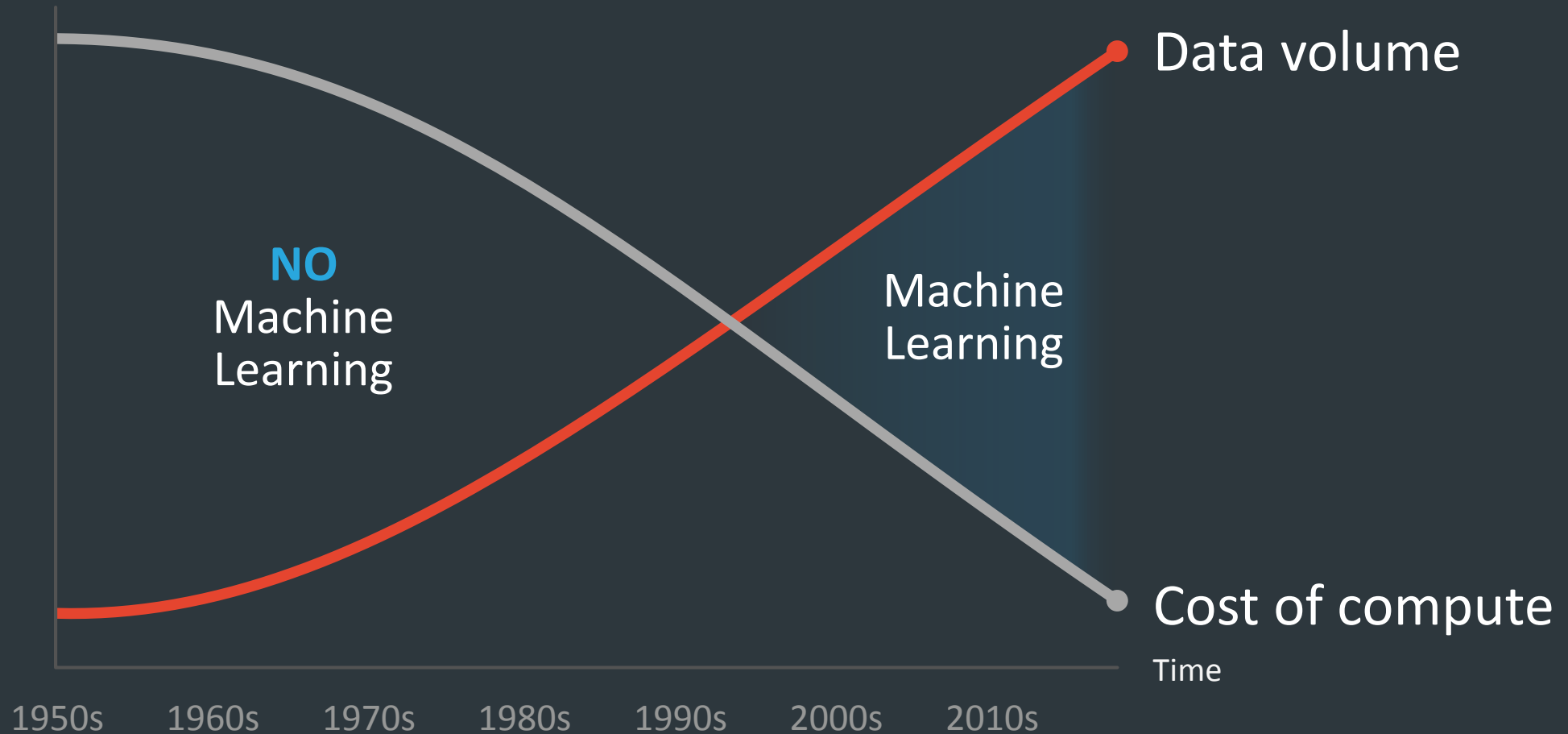
---

Cloudera Data Science and Machine Learning

Robin Harrison, Account Executive

David Kemp, Systems Engineer

# This is the age of machine learning.



# Machine learning presents a multitude of opportunities

Data



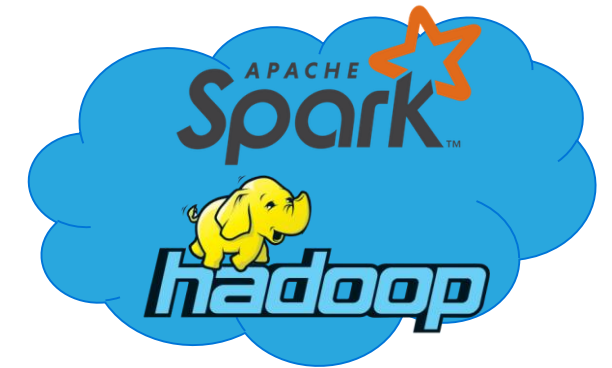
Data has never been more plentiful

Analytics



Open source data science and machine learning libraries are rapidly evolving

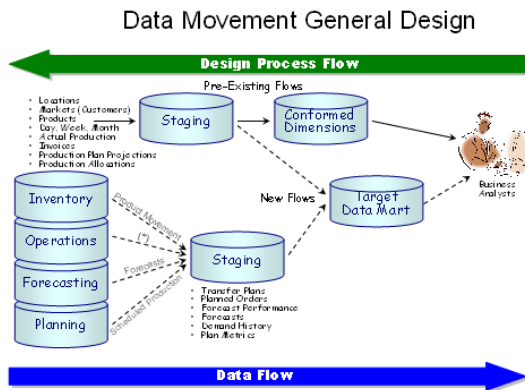
Deployment



Flexible commodity storage and compute make scalable production machine learning affordable

# But there are practical challenges

Data



Analytics



Deployment



Data volumes are increasing and it needs to move across multiple different systems

Teams have different, conflicting requests for languages & libraries

Most data science done at small scale, individually, and is difficult to replicate

Very few models reach production

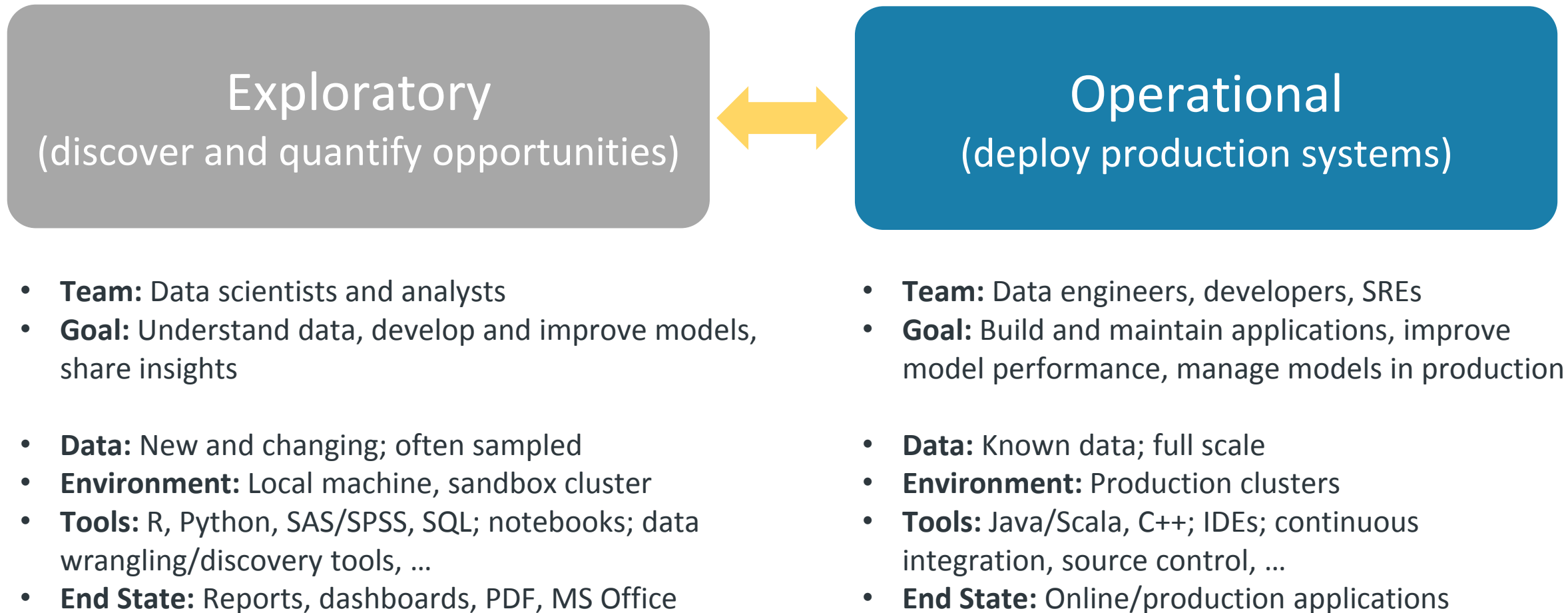
# What is Machine Learning and Data Science

- Machine Learning and Data Science: algorithms and methods that extract useful insights and patterns from data.
- These insights can drive profits, find outliers, cluster like items, predict future issues and insights, cut losses, classify different groupings and many other tasks.

# What does a Data Scientist Do?



# Types of data science



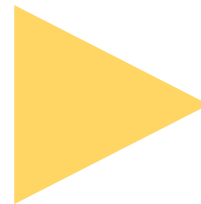
# Our goal: Open data science at enterprise scale

Help more data scientists  
use the power of Cloudera

Use a powerful, familiar  
environment with direct access  
to Cloudera data and compute



Data Scientist  
Data Engineer



Make it easy and secure to  
add new users, use cases

Offer secure self-service  
analytics and a faster path to  
production on common,  
affordable infrastructure



Enterprise Architect  
Hadoop Admin



# Balancing the needs of data scientists and IT

## Data Scientists

explore, experiment, collaborate



## IT

drive adoption, maintain compliance



# Cloudera Data Science Workbench

Self-service data science for the enterprise

Accelerates data science from development to production with:

- Secure self-service data access
- On-demand compute
- Support for Python, R, and Scala
- Project dependency isolation for multiple library versions
- Workflow automation, version control, collaboration and sharing

The screenshot displays the Cloudera Data Science Workbench interface. At the top, there are four dashboard cards showing: 0 sessions running, 2 jobs running, 3 vCPU (0/80.00), and 6 GB (0/263.86). Below these is a 'Projects' section with a 'Product Overview' card. The main area is a code editor for a Python script named '1\_python.py'. The code includes comments and uses pandas and seaborn for data analysis. To the right of the code editor is a terminal window showing the output of 'data.head()' as a table:

Date	djia	debt
2004-01-14	10485.18	0.210000
2004-01-22	10528.66	0.210000
2004-01-28	10702.51	0.210000
2004-02-04	10499.18	0.213333
2004-02-11	10579.03	0.200000

Below the terminal is a line chart titled 'DJIA vs. Debt Query Volume' showing two data series from 2004 to 2011. The chart includes zoom controls (1m, 3m, 6m, YTD, 1y, All) and date range selectors (From: Jan 12, 2004, To: Mar 2, 2011). At the bottom, there is a 'Cluster Metadata' panel showing details for a script named 'bin/transformation.py', including its schedule, engine profile (1 vCPU, 2 GB memory), and a 'Job History' chart showing duration over time.

# With Cloudera Data Science Workbench...

## Data scientists can:

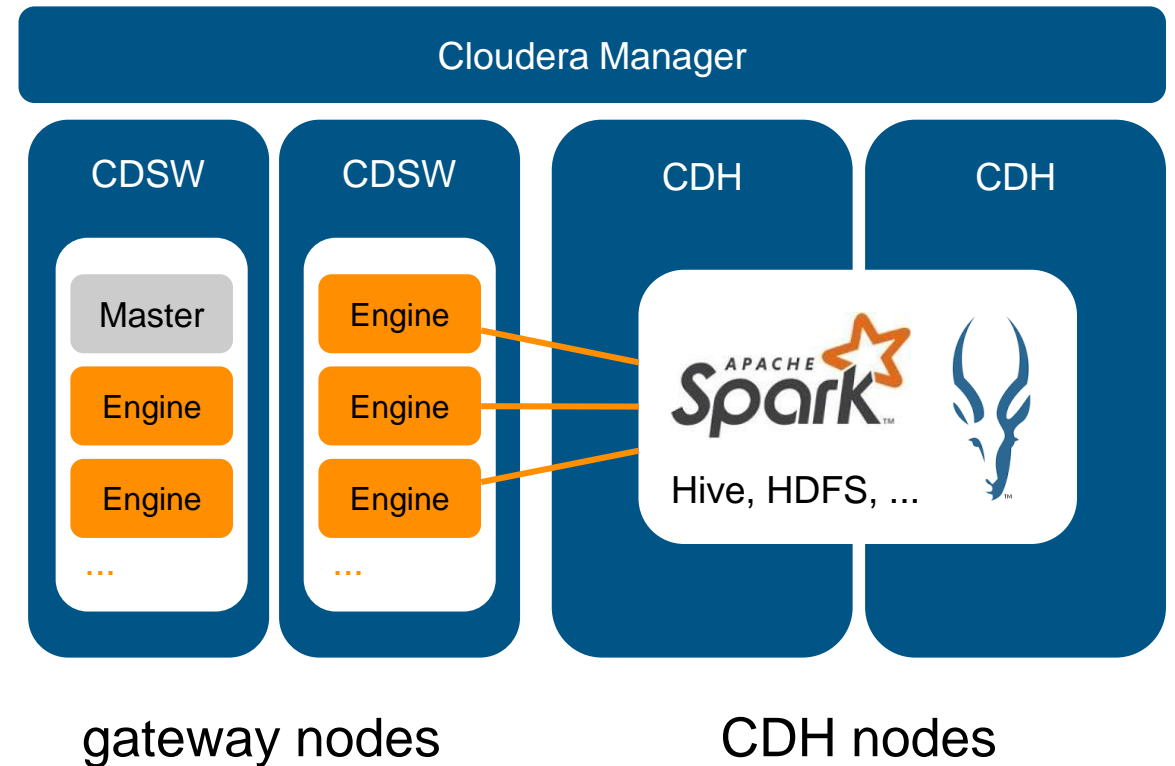
- Use R, Python, or Scala from a web browser, with no desktop footprint
- Install any library or framework within isolated project environments
- Directly access data in secure clusters with Spark and Impala
- Share insights with their team for reproducible, collaborative research
- Automate and monitor data pipelines using built-in job scheduling

## IT can:

- Give their data science team the freedom to work how they want, when they want
- Stay compliant with out-of-the-box support for full platform security, especially Kerberos
- Run on-premises or in the cloud, wherever data is managed

# A modern data science architecture

- Built on Docker and Kubernetes
- Runs on dedicated gateway nodes
- User sessions run in isolated “engine” containers which:
  - Host Kerberos-authenticated Python/R/Scala runtimes
  - Interact with Spark via YARN client mode (Driver runs in container, workers on CDH)
- Single-cluster only (for now)



CYBERSECURITY

- » THREAT DETECTION
- » DATA SECURITY
- » MACHINE LEARNING



Uncovering Zero-Day Attacks and  
Stopping Advanced Persistent Threats  
More Quickly

- Helps threat hunters obtain responses to queries magnitudes faster
- Provides access to a wider range of data that wasn't accessible before
- Increases researcher productivity by 60 percent





## MANUFACTURING

- » PREDICTIVE ANALYTICS
- » PROCESS IMPROVEMENT
- » PRODUCT INNOVATION



### Improving Flight Safety with Rapid, Data-Driven Decision Support

- Uncovers patterns in aircraft performance and parts that can help Sikorsky engineers **improve flight safety** and **optimize aircraft operations**
- **Extend useful life** of key components
- **Helps prevent unscheduled maintenance** and better prioritize repairs



## PROBLEM

Needed scalable system for real-time endpoint threat detection and response

- System couldn't handle growing number of endpoints
- No real-time processing
- Limited operational resources

## SOLUTION

Replatformed to deliver actionable security intelligence to users

- Support deployments with **>100,000 endpoints**
- Threat detection and **response in minutes** vs months
- Cloudera Predictive Support **anticipates issues before they occur** & direct connection to the experts





DATA-DRIVEN  
PRODUCTS

## AEROSPACE

- » SPACECRAFT TELEMETRY
- » REMOTE MONITORING
- » PREDICTIVE MAINTENANCE

### Aerospace – Spacecraft Telemetry

Advanced analytics on streaming data to reduce human space mission risks

#### Challenge:

- Over **2 TB/ hour** of telemetry test data streaming in from over 1200 sensors in test environment

#### Solution:

- Cloudera cluster supporting high rate of data ingest – up to **~300MB/sec**
- Advanced analytics run on the streaming data to check for issues or determine patterns and reduce risk

**cloudera**







## Connected Product Support

### Juniper Networks Monitors Thousands of In-Field Devices with Cloudera & Zoomdata

#### Challenge:

- Monitor thousands of in-field devices in real time to provide the best support experience

#### Solution:

- Cloudera Enterprise + Zoomdata provide every support client with access to both aggregate and detailed view of their devices
- Pinpoint issues at any level: network, device or application

The dashboard displays the following metrics:

- 3 Cases awaiting input
- 2k EOL parts
- 21 PHC issues
- 1 Adv Svcs credit
- 2 RMAs pending
- 5 KB recommended articles

Support Exposure:

- Contracts: 2.6K Active, 368 Expired, 80 < 30 days, 37 30-90 days
- EOL: 2.0K Not pending EOL, 2.0K EOL, 114 EOL 60-180 days

Host Name	Active	Expired	< 30 days	30-90 days
re0-ce-dmrc0i04w.delgeny.co.denver	3	0	0	0
sfc0-re0-ar04.howardcounty.red.bad	3	0	0	0
sur01.sacramento.ca.ccal.comcast.net	3	0	0	0
pe-wchclbz03w.wchicago.il.biz.comcast.net	2	0	0	0
re0-ar03.norristown.pa.panjdc	2	0	0	0
re0-ca-btvlmd1200w.baltsvilla.md.bad	2	0	0	0
re0-ce-hanob0802w.cicero.tx.houston	2	0	0	0
re0-ca-nrcsgagz04w.peachtrecr.ga.atlanta	2	0	0	0

Device Health:

- PHC: 21 Minor, 15 None
- Recommended Release: 15.8K Prior, 12.5K Current, 687 Other

Host Name	Minor	None
CIR1.Ashburn-VA	2	0
CIR1.Miami2-FL	2	0
CIR1.Seattle7-WA	2	0
MCR1.Charlotte-NC	2	0
MCR1.Cleveland-OH	2	0
MCR1.Newark-NJ	2	0
CIR1.Denver2-CO	1	0
CIR1.London2-ENG	1	0

# cloudera

---

Thank you!

Robin Harrison, [robinharrison@cloudera.com](mailto:robinharrison@cloudera.com), 703-795-4706

David Kemp, [david.kemp@cloudera.com](mailto:david.kemp@cloudera.com), 703-282-2317