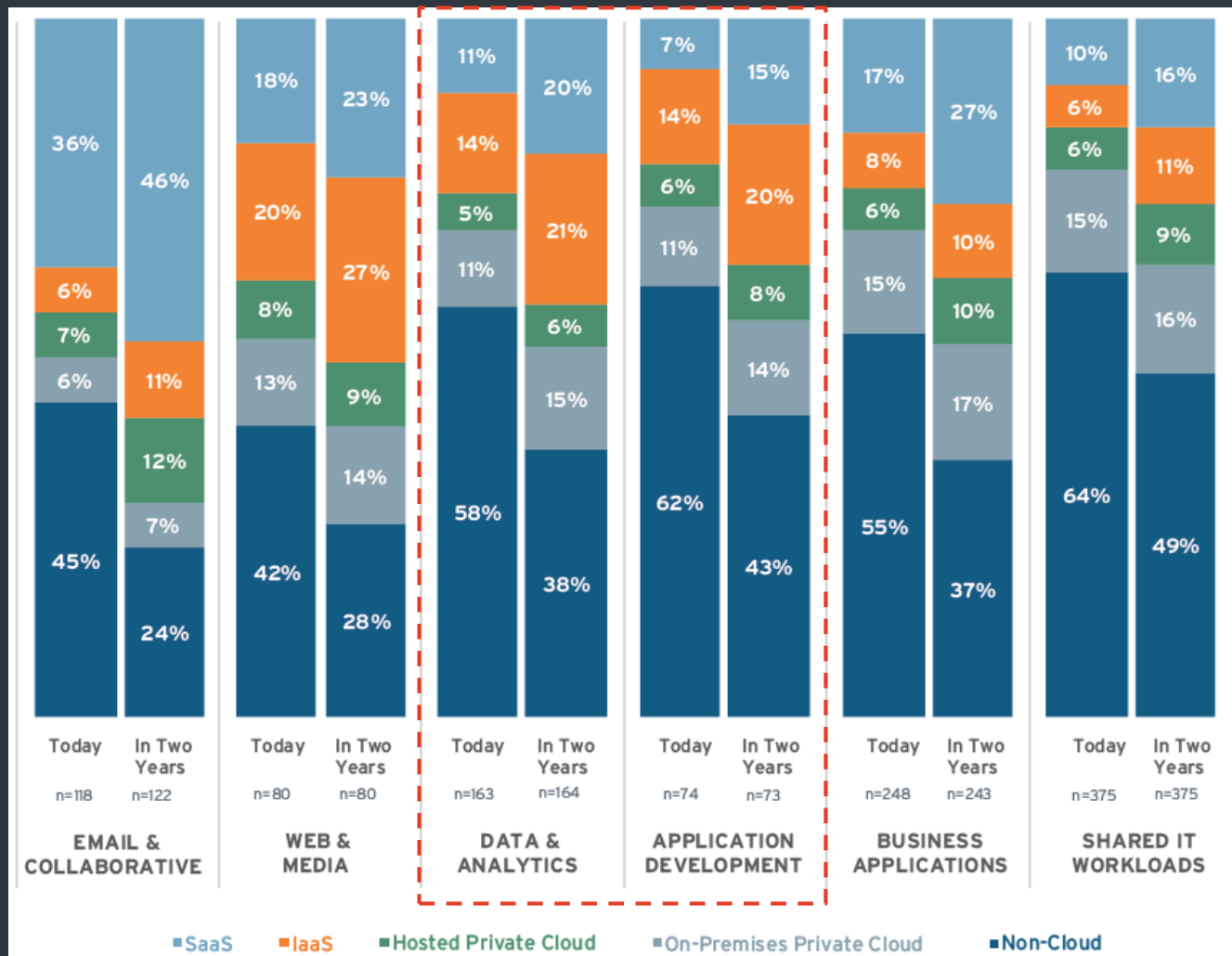# cloudera

## Cloudera Altus: Big Data in der Cloud einfach gemacht

Michael Kohs    |    Sales Engineer    |    mkohs@cloudera.com

# Shift to cloud: an analyst view
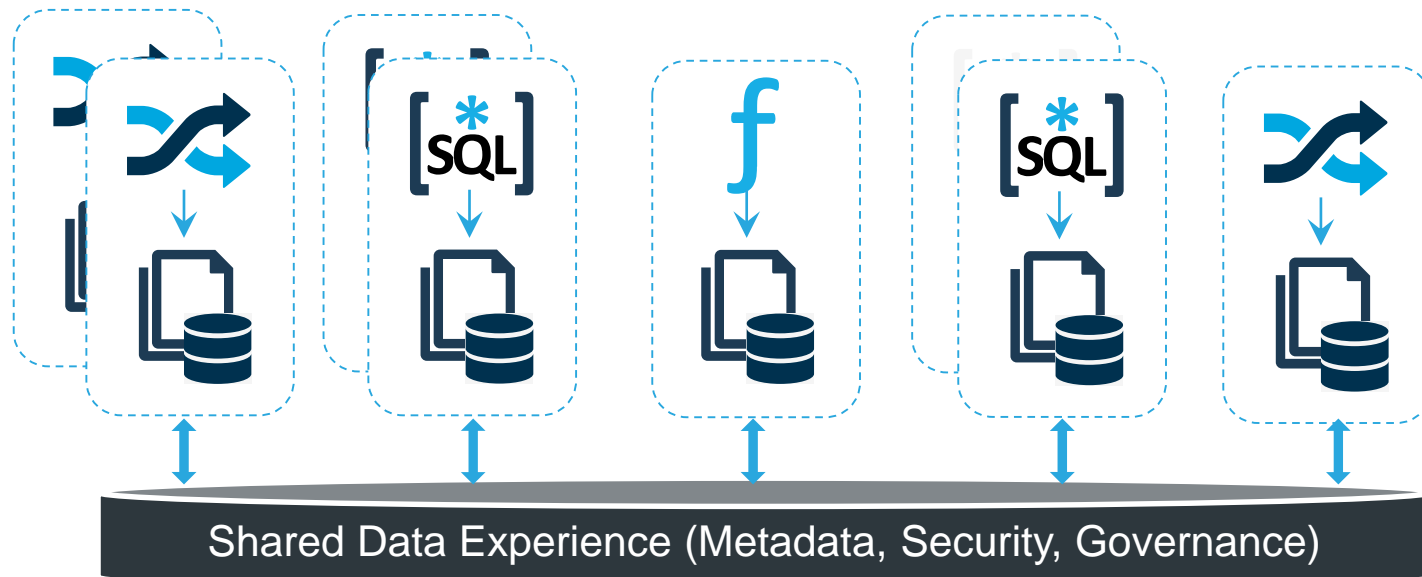


**Key Points**

- Cloud deployments will be the dominant environment in every category

- Every cloud deployment environment will see increases in every workload category

- Analytics and App Development areas expected strong gains

**Source**: 451 Research, Voice of the Enterprise: Workloads and Key Projects, Cloud Transformation, 2017.

# My organization
# is moving to the cloud,
# why should we
# consider Cloudera?

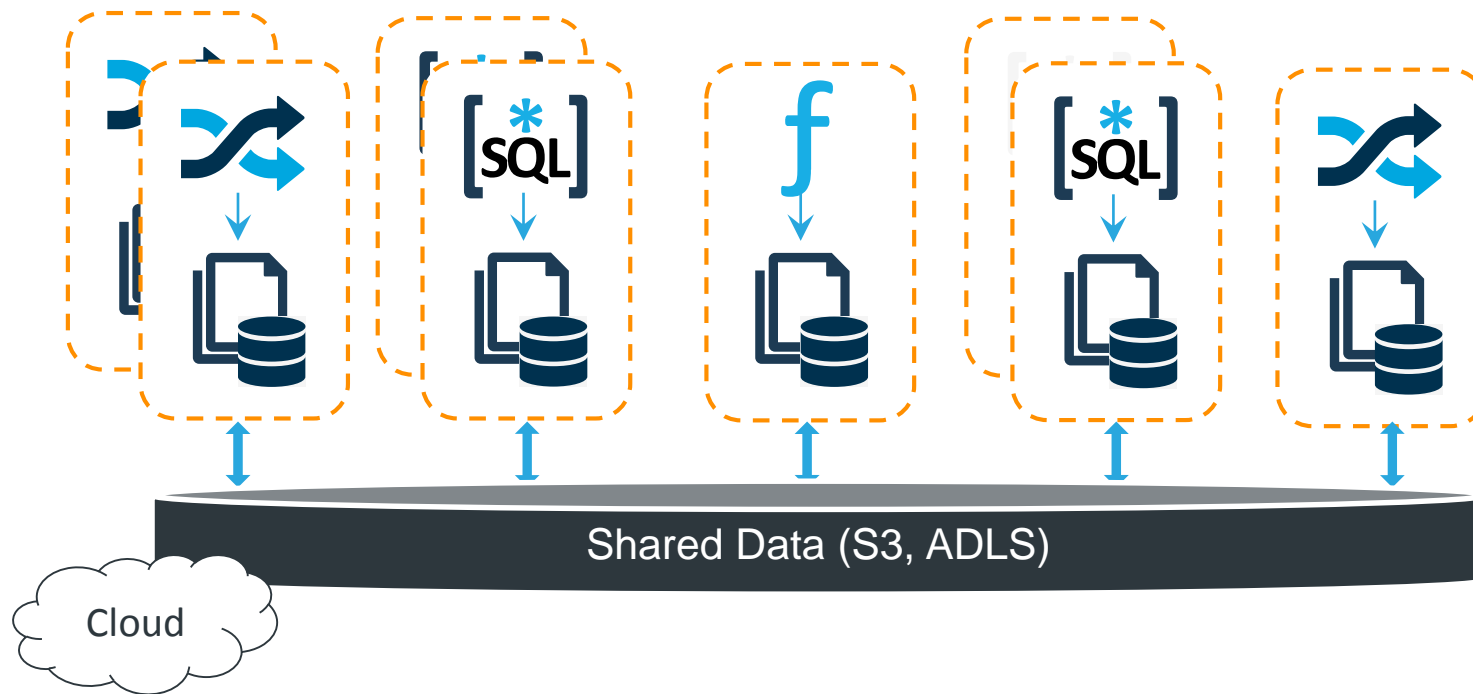# Traditional on-premises deployments perform reasonably well

One physical cluster provides a shared data experience to multiple workloads and tenants

Shared Data Experience (Metadata, Security, Governance)

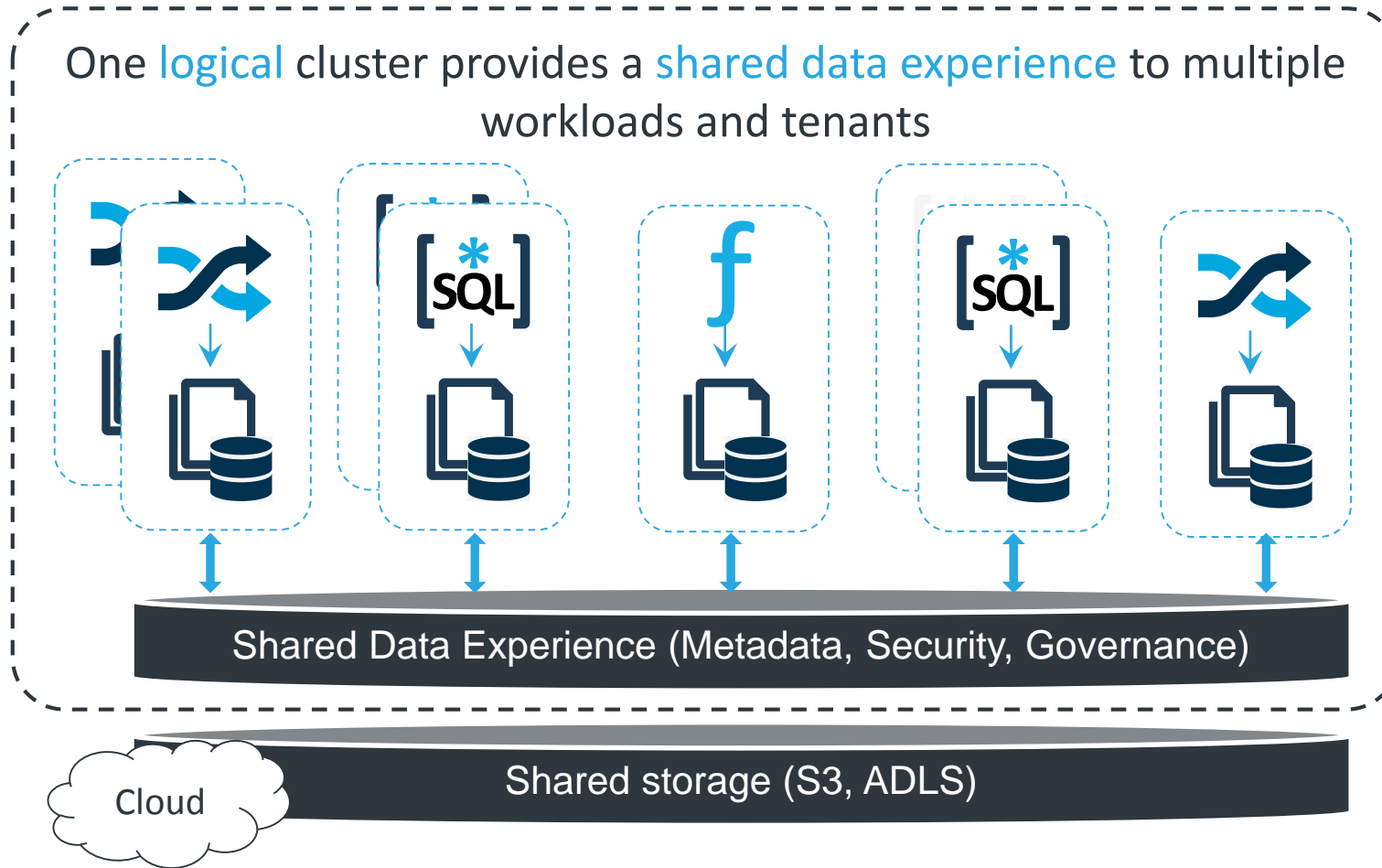| | |
|---|---|
| ● | Strong multi-function support |
| ● | Strong shared data experience |
| ● | Strong operational model |
| ◑ | Moderate cost management |
| ◑ | Moderate tenant isolation |
| ◑ | Moderate workload elasticity |
| ○ | Weak on self service |
| ○ | Weak on speed of deployment |

But not good enough for tomorrow

# Traditional cloud deployments are strong where on-premises is weak, but at the expense of creating workload silos



| | |
|---|---|
| ◐ | Moderate multi-function support |
| ○ | Weak on shared data experience |
| ○ | Weak operational model |
| ◐ | Moderate cost management |
| ● | Strong on tenant isolation |
| ● | Strong on workload elasticity |
| ● | Strong on self service |
| ● | Strong on speed of deployment |

Shared Data (S3, ADLS)

Cloud

**This is the experience of cloud house offerings**

cloudera

# Only Cloud deployments with SDX optimize for all design goals

One **logical** cluster provides a **shared data experience** to multiple workloads and tenants



Shared Data Experience (Metadata, Security, Governance)

Shared storage (S3, ADLS)

Cloud

| | |
|---|---|
| ● | Strong multi-function support |
| ● | Strong shared data experience |
| ● | Strong operational model |
| ● | Strong on cost management |
| ● | Strong on tenant isolation |
| ● | Strong on workload elasticity |
| ● | Strong on self service |
| ● | Strong on speed of deployment |

## SDX makes it possible to transfer on-premises design wins to cloud

cloudera

# Cloudera's Public Cloud Reference Architecture

**Management**

| Cloudera Manager | Cloudera Navigator | Altus | Cloudera Director |
|---|---|---|---|

Each **management service** supports all Workload and SDX services

**Compute**

Spark — Hive — Hive on Spark — Impala — Solr — …

Each workload runs in an isolated **Workload Cluster**

**cloudera sdx** shared data experience

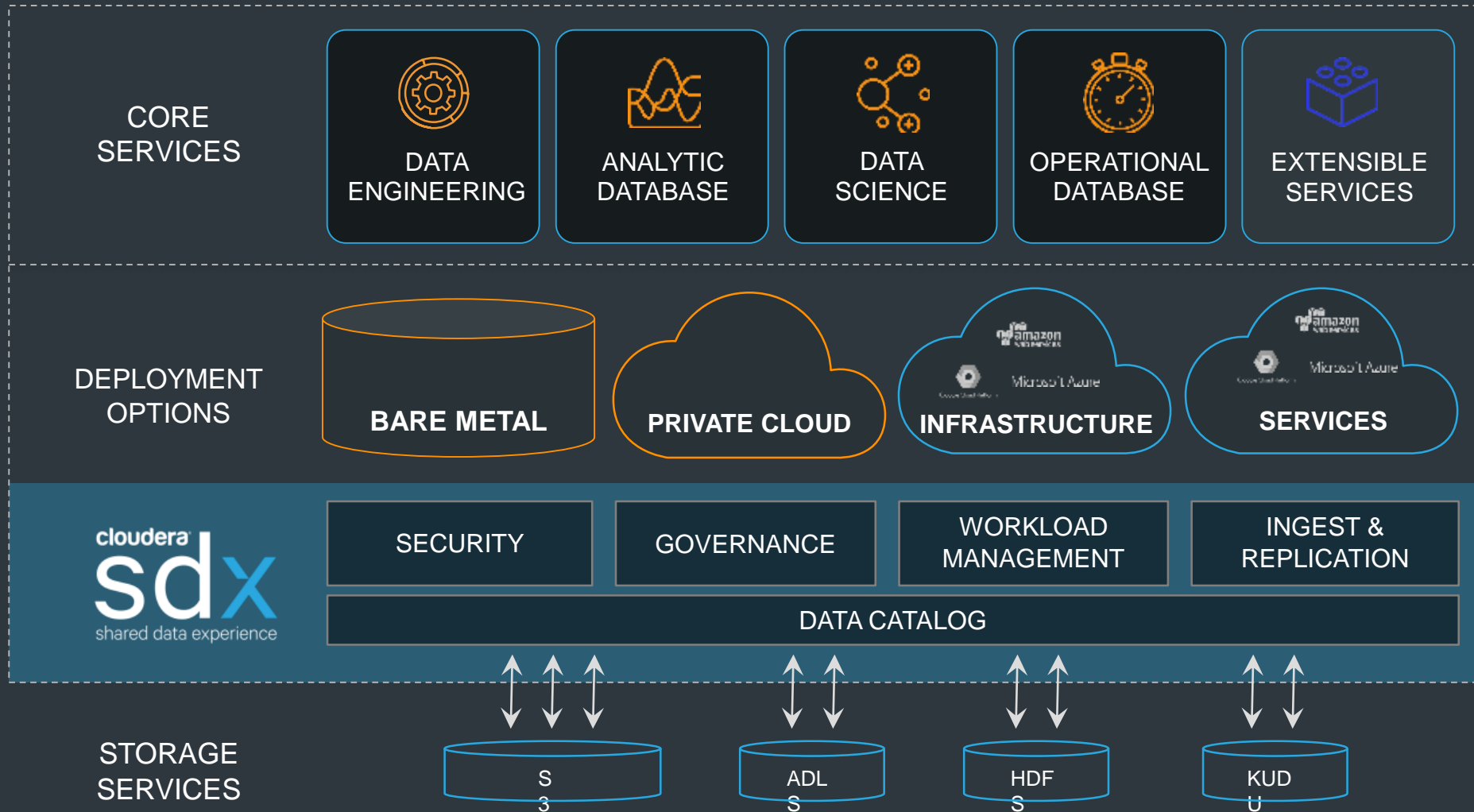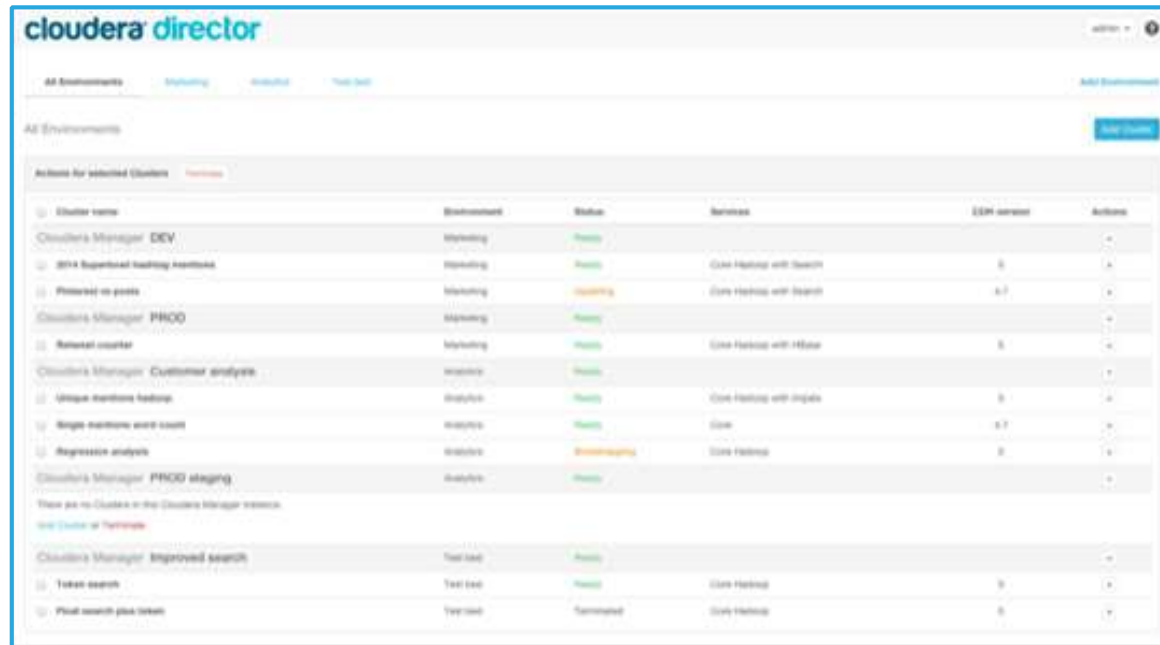| Hive Metastore | Navigator Metastore | Sentry Metastore |
|---|---|---|

SDX services run on shared RDS | MySQL

**Storage**

S3, ADLS

# Cloudera Enterprise

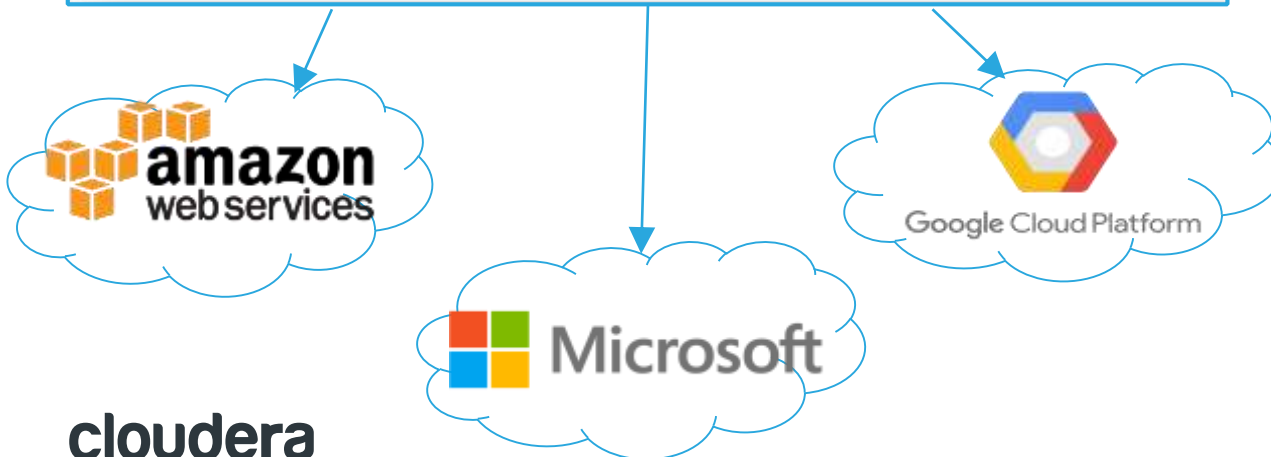The modern platform for machine learning and analytics optimized for the cloud



CORE
SERVICES

DATA
ENGINEERING

ANALYTIC
DATABASE

DATA
SCIENCE

OPERATIONAL
DATABASE

EXTENSIBLE
SERVICES

DEPLOYMENT
OPTIONS

BARE METAL

PRIVATE CLOUD

INFRASTRUCTURE

SERVICES

cloudera
sdx
shared data experience

SECURITY

GOVERNANCE

WORKLOAD
MANAGEMENT

INGEST &
REPLICATION

DATA CATALOG

STORAGE
SERVICES

S3

ADLS

HDFS

KUDU

cloudera

# Cloudera Director for cluster lifecycle management



## Easy
- Single pane of glass for all cloud infrastructure
- Create templates to run applications in a pre-optimized manner

## Flexible
- Multi-cloud: AWS, Azure, GCP
- Hourly pricing with auto billing & metering
- Spot instance/block support

## Enterprise-grade
- Integration across Cloudera Enterprise
- Management of CDH deployments at scale
- Deeply integrated with Cloudera Manager

# Cloudera Altus
## Multi-cloud foundation for building new cloud services

Altus Platform Services

DATA ENGINEERING

ANALYTIC DATABASE
(beta)

Altus PaaS Foundation

# Director vs Altus - when to use either?

| | Altus | Director |
|---|---|---|
| Automated log saving | x | |
| Automated Cluster spin up / down (no extra coding) | x | |
| Data Engineering – Hive, Spark, HoS, MR | x | x |
| Production Job Driven | x | |
| Workload Analytics | x | |
| Cluster Duration | Purely Transient | Transient OR Persistent |
| Job Development / Exploration | | x |
| 3rd Party Installations | | x |
| Full Control of CM | | x |
| Analytical / Operational - Impala, HBase, Search | | x |
| Persistent (or Transient) | | x |
| Grow / Shrink Cluster | | x |

**cloudera**

# Altus Service Architecture (AWS)



- Runs in Cloudera's secured and monitored environment
- Manages CDH clusters in customer cloud account
- Customer data does not pass to Cloudera  (Workload Analytics requires opt-in log data transfer to Cloudera)

12

# Altus Service Architecture (Azure)



- Runs in Cloudera's secured and monitored environment
- Manages CDH clusters in customer cloud account
- Customer data does not pass to Cloudera  (Workload Analytics requires opt-in log data transfer to Cloudera)

# Altus Data Engineering
## for ETL, machine learning, and data processing

- Fast, easy job submission without the cluster management

- Built-in Workload Analytics for troubleshooting and optimization

- Lower costs with transient resources and pay-per-use pricing

- Full benefits of isolation + Shared Data Experience



**cloudera**

# End-user focused with jobs as first-class objects

**Workload troubleshooting and analytics**

- Troubleshoot jobs after cluster termination through job log and configuration browsing

- Insight into causes of job failure

- Identification and root cause analysis of slow jobs

# Capture metadata spanning multiple clusters

## Persist metadata with Cloudera Navigator

- Export metadata and lineage information from Altus clusters

- Insight into full data management pipeline including transient clusters

# Three immediate use cases for Altus Data Engineering

| ETL FOR ANALYTIC DB | BATCH MACHINE LEARNING | ETL OFFLOAD |
|---|---|---|

**ETL** ➡ **Analytic DB**

**Data Science** ➡ **ML**

**ETL** ↕ **On-Prem**

Cloud-native batch preparation for Impala on IaaS or, soon, Altus Analytic DB.

Scalable compute for massively-parallel batch machine learning training, scoring, or simulation.

Offload batch processing jobs from overburdened on-premises clusters.

**cloudera**

# Altus Analytic Database

The first data warehouse cloud service to bring the warehouse to the data - delivering instant analytics to anyone

## For business analysts:

- Query with predictable performance, at any time, without risking SLAs
- Bring limitless new users and use cases with instant self-service analytic access
- Data available for broad access (SQL, BI tools, Python, R, etc)

## For IT:

- Easily and elastically provision isolated resources as and when they're needed
- Simple multi-tenant management including federated identity and consistent governance
- Eliminate data movement and copies across workloads

**cloudera**

# ETL to Cloud-Native Analytic Database

Workflow, Monitoring, and Scheduling (e.g. AutoSys, ControlM, Airflow, Talend, Informatica)

Transient Cluster(s):
Atlus Managed

**Data Engineering Workflow: ETL**

Source Data

Source Data

1+ Jobs
(e.g. Hive)

1+ Jobs
(e.g. Spark)

BI Tools
(e.g. Tableau)

SQL Editor
(e.g. Hue)

Analytic Database
(Impala)

Elastic Cluster(s):
Director Managed
/ Altus Managed

SDX: Schema (HMS), Security (Sentry), Lineage/Audit (Navigator)

Long Running Cluster:
Director Managed

Object Store (ADLS, S3)

cloudera

# The Scenario

**My Role:**  Data Analyst at DataCo – a Sports Retailer

**Business Issue:**  Experiencing lower than expected website sales. Why?

**Technical Issues:**  I have a data warehouse on premise, which contains my sales
order  data, but it is very old and slow and it is difficult to do ad hoc
queries on  it.

My clickstream data is too big to ingest into my data warehouse

**Requirements:**  Need an Analytic Database to do ad hoc queries on order data

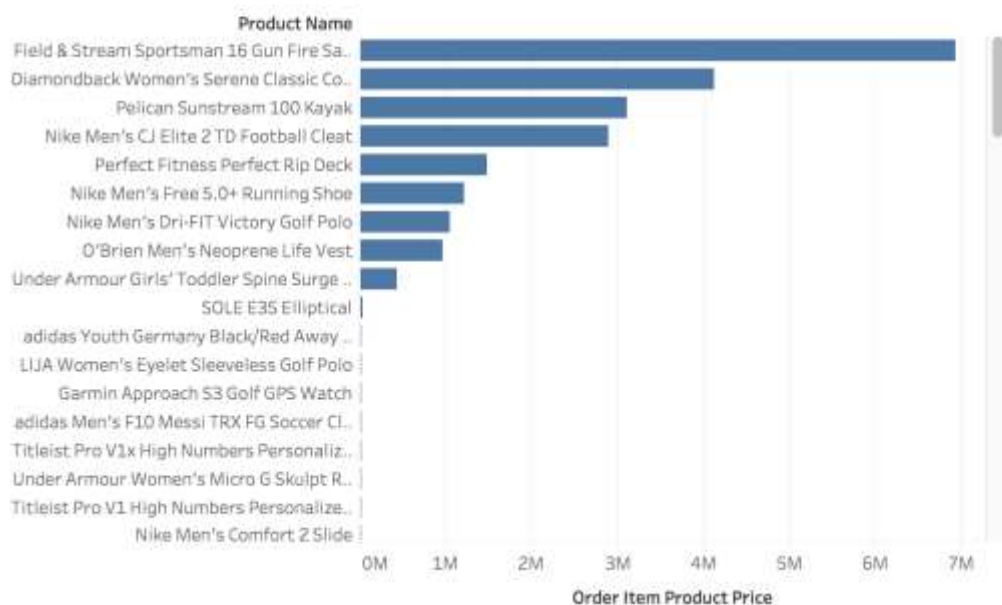Need a temporary platform to process weblogs once a day

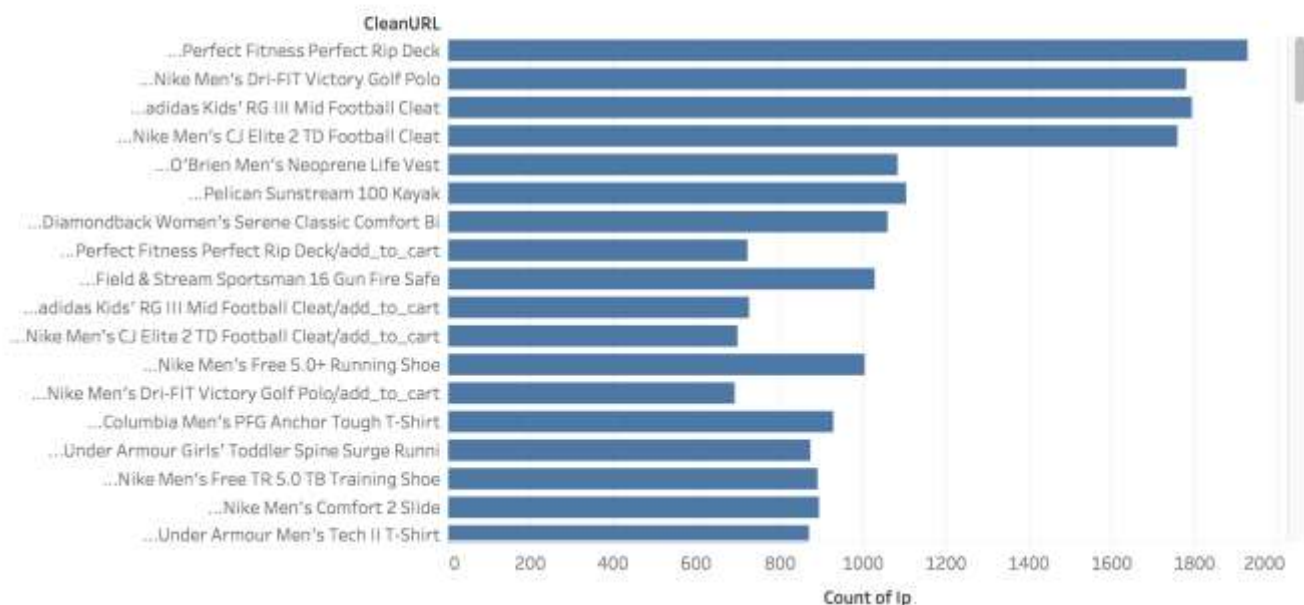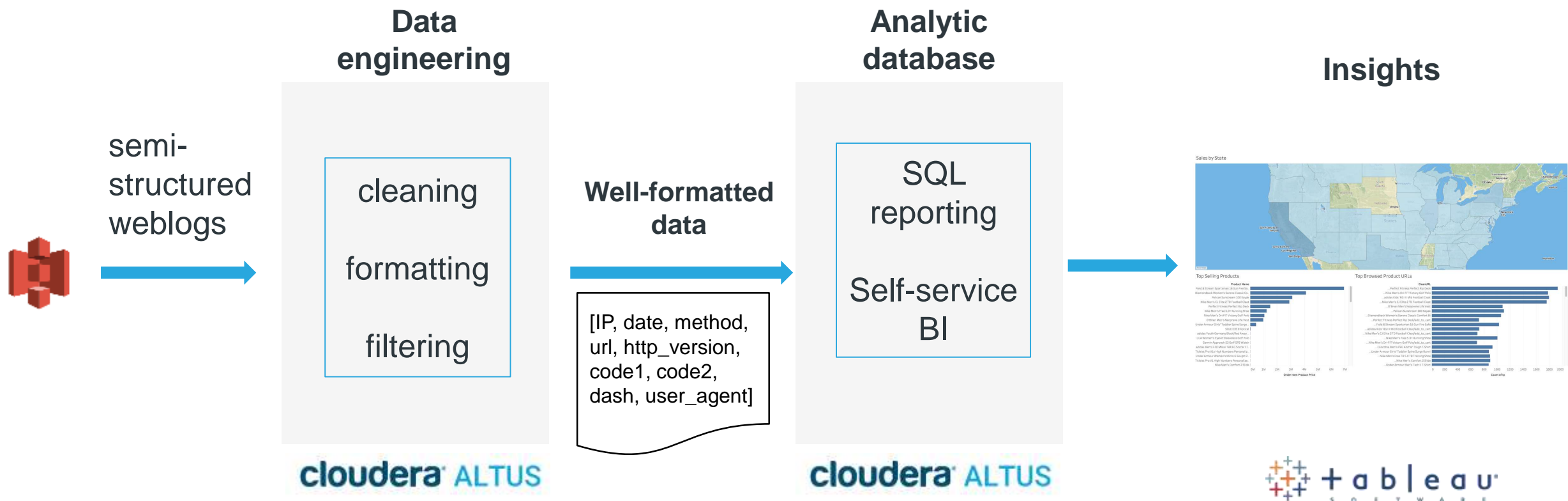Ability to join processed weblogs to order data

## Sales by State



## Top Selling Products

| Product Name | Order Item Product Price |
|---|---|
| Field & Stream Sportsman 16 Gun Fire Sa.. | |
| Diamondback Women's Serene Classic Co.. | |
| Pelican Sunstream 100 Kayak | |
| Nike Men's CJ Elite 2 TD Football Cleat | |
| Perfect Fitness Perfect Rip Deck | |
| Nike Men's Free 5.0+ Running Shoe | |
| Nike Men's Dri-FIT Victory Golf Polo | |
| O'Brien Men's Neoprene Life Vest | |
| Under Armour Girls' Toddler Spine Surge .. | |
| SOLE E35 Elliptical | |
| adidas Youth Germany Black/Red Away .. | |
| LIJA Women's Eyelet Sleeveless Golf Polo | |
| Garmin Approach S3 Golf GPS Watch | |
| adidas Men's F10 Messi TRX FG Soccer Cl.. | |
| Titleist Pro V1x High Numbers Personaliz.. | |
| Under Armour Women's Micro G Skulpt R.. | |
| Titleist Pro V1 High Numbers Personalize.. | |
| Nike Men's Comfort 2 Slide | |

*x-axis: 0M 1M 2M 3M 4M 5M 6M 7M — Order Item Product Price*

## Top Browsed Product URLs

| CleanURL | Count of Ip |
|---|---|
| ...Perfect Fitness Perfect Rip Deck | |
| ...Nike Men's Dri-FIT Victory Golf Polo | |
| ...adidas Kids' RG III Mid Football Cleat | |
| ...Nike Men's CJ Elite 2 TD Football Cleat | |
| ...O'Brien Men's Neoprene Life Vest | |
| ...Pelican Sunstream 100 Kayak | |
| ...Diamondback Women's Serene Classic Comfort Bi | |
| ...Perfect Fitness Perfect Rip Deck/add_to_cart | |
| ...Field & Stream Sportsman 16 Gun Fire Safe | |
| ...adidas Kids' RG III Mid Football Cleat/add_to_cart | |
| ...Nike Men's CJ Elite 2 TD Football Cleat/add_to_cart | |
| ...Nike Men's Free 5.0+ Running Shoe | |
| ...Nike Men's Dri-FIT Victory Golf Polo/add_to_cart | |
| ...Columbia Men's PFG Anchor Tough T-Shirt | |
| ...Under Armour Girls' Toddler Spine Surge Runni | |
| ...Nike Men's Free TR 5.0 TB Training Shoe | |
| ...Nike Men's Comfort 2 Slide | |
| ...Under Armour Men's Tech II T-Shirt | |

*x-axis: 0 200 400 600 800 1000 1200 1400 1600 1800 2000 — Count of Ip*

# Demo - Retail clickstream analysis

**Data engineering**

**Analytic database**

**Insights**

semi-structured weblogs

cleaning

formatting

filtering

**Well-formatted data**

[IP, date, method, url, http_version, code1, code2, dash, user_agent]

SQL reporting

Self-service BI

**cloudera** ALTUS

**cloudera** ALTUS

+ + + + + + + + t a b l e a u
+ + + + SOFTWARE

**cloudera**

**Cloudera Manager**

**Cloudera Director**

**Cloudera Altus**

**Cloudera Altus**

Transient Data Engineering clusters (Altus)

Self-Service Analytic DB clusters

Click stream data

Long-running Kafka cluster

Long-running stream processing cluster

Long-running Analytic DB cluster (Impala)

HDFS on premise

Shared Data Experience (Metadata, Security, Governance)

Object store

amazon web services

Microsoft Azure

cloudera

# cloudera

## Q&A

# Resources

Cloudera Altus

Cloudera Altus documentation

Cloudera Director

Cloudera Director documentation

Cloudera SDX

Try Cloudera Director on Microsoft Azure

Try Cloudera Director on AWS with AWS Quickstart

Cloudera Reference Architectures for public and private cloud deployments

cloudera