**CLOUDERA**

# ADATPLATFORM A FELHŐBEN MENNI VAGY NEM MENNI?

Gáspár Balázs | Solutions Engineer

2021 június 23.
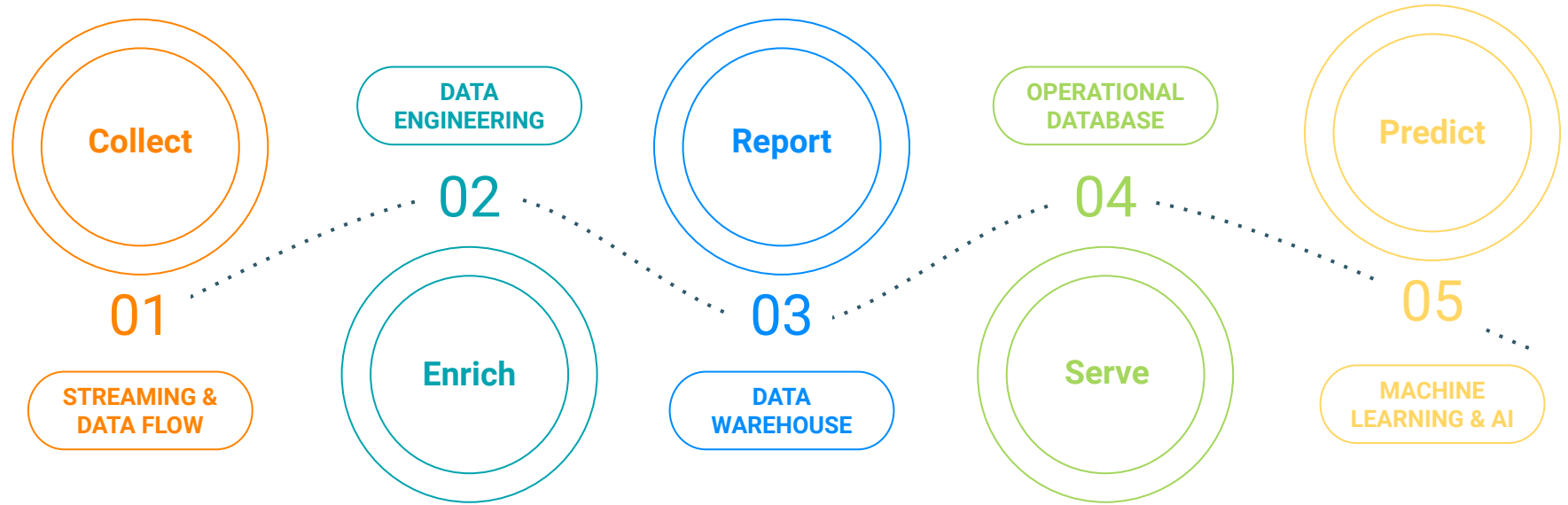
# AGENDA

**Mitől modern egy adatplatform?**
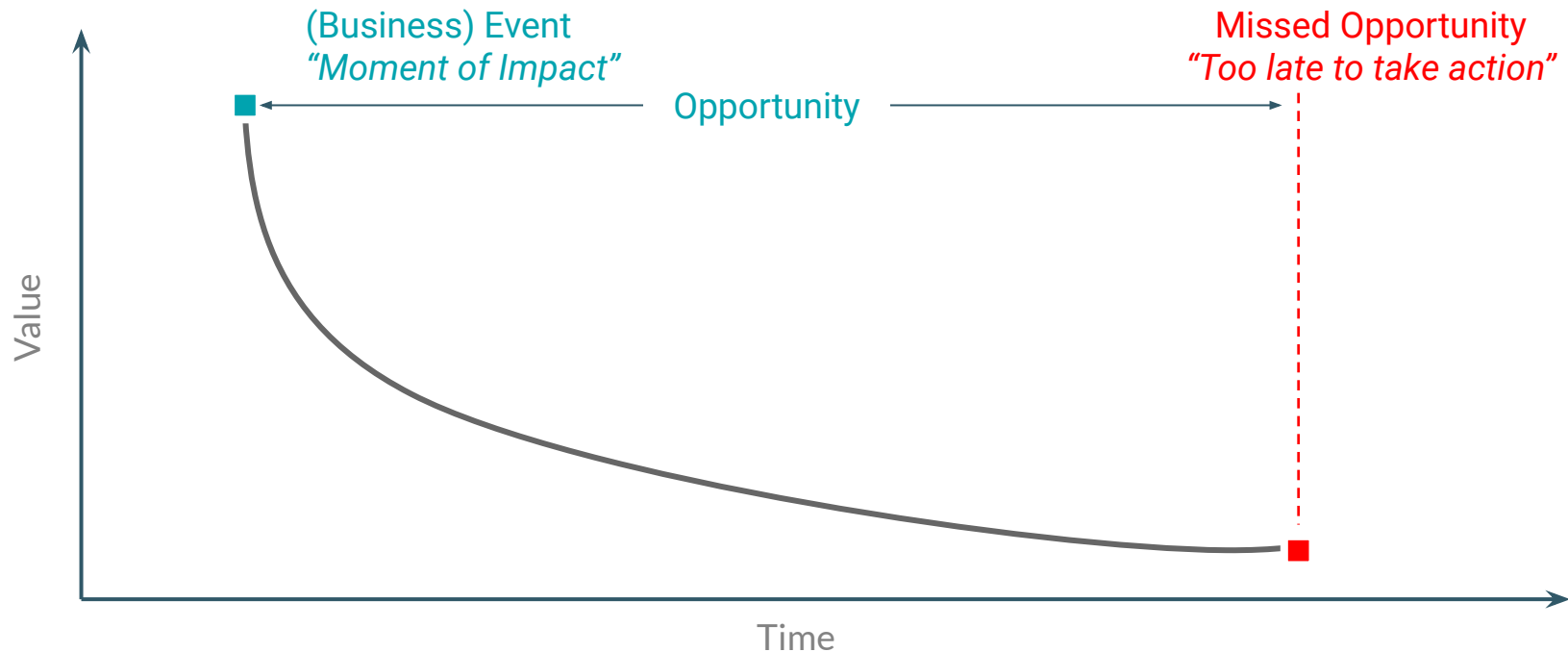
Adatmenedzsment a felhőben

Az út a felhőbe

# MANAGING THE DATA LIFECYCLE - BATCH AND STREAMING



**Collect**

01

STREAMING & DATA FLOW

**DATA ENGINEERING**

02

**Enrich**

**Report**

03

DATA WAREHOUSE

**OPERATIONAL DATABASE**

04

**Serve**

**Predict**

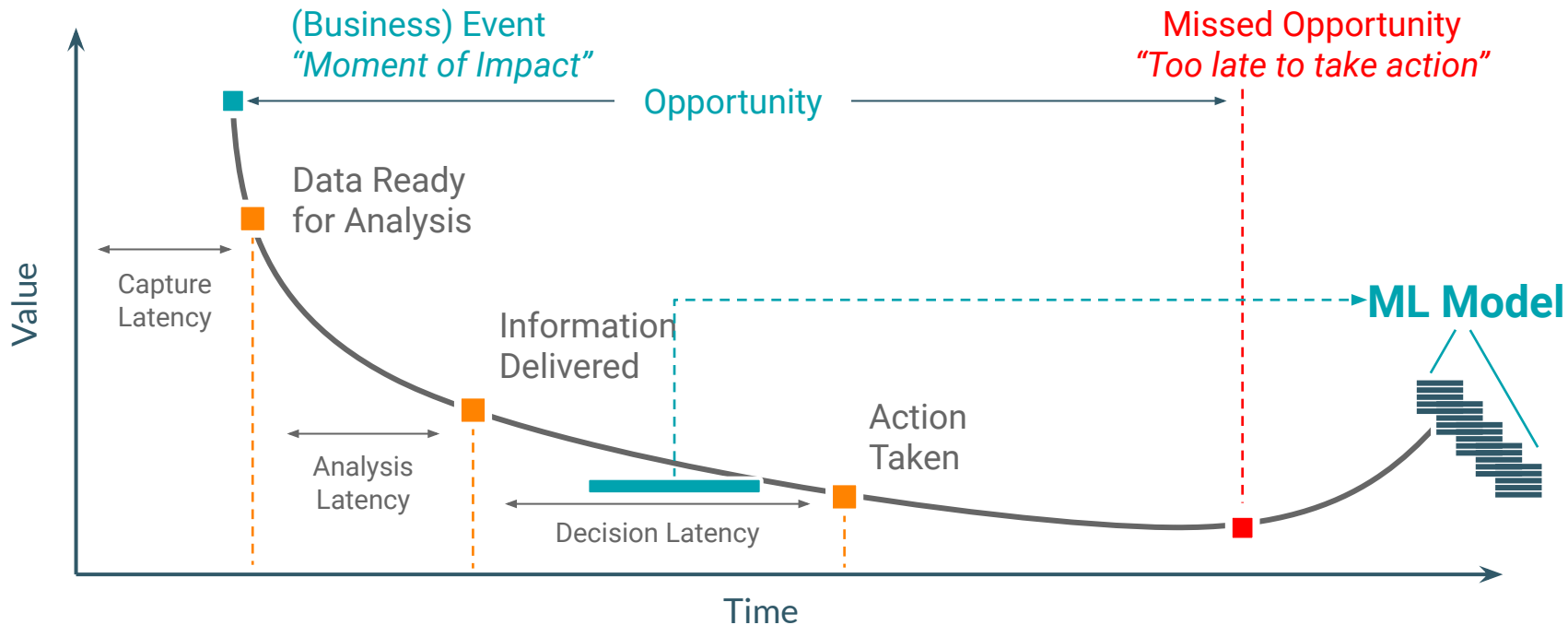05

MACHINE LEARNING & AI

**CLOUDERA SDX** SECURITY | GOVERNANCE | LINEAGE | MANAGEMENT | AUTOMATION

# Streaming and machine learning are platform problems

# CREATING VALUE FROM THE REAL-TIME FLOW OF BIG DATA

CLOUDERA

# CREATING VALUE FROM THE REAL-TIME FLOW OF BIG DATA



(Business) Event
*"Moment of Impact"*

Missed Opportunity
*"Too late to take action"*

Opportunity

Data Ready
for Analysis

Value

Capture
Latency

Information
Delivered

ML Model

Analysis
Latency

Action
Taken

Decision Latency

Time

# IMAGINE A WORLD WITH SQL ANALYTICS ON *ALL AND ANY DATA*

```
1  SELECT
2      geo_event.eventTimestamp, geo_event.driverId, geo_event.eventTime,geo_event.eventSource,
3      geo_event.truckId,geo_event.driverName,geo_event.routeId,geo_event.route,geo_event.eventType,
4      geo_event.latitude, geo_event.longitude, geo_event.correlationId, geo_event.geoAddress,
5      speed_event.speed,
6      driver.certified, driver.wage_plan,
7      timesheet.hours_logged, timesheet.miles_logged
8  FROM
9      geo_events_json as geo_event
10     join speed_events_json as speed_event
11       on (geo_event.driverId = speed_event.driverId)
12     left join CDP_Hive_Catalog.employees_hr_hive_db.driver
13         FOR  SYSTEM_TIME AS OF PROCTIME() driver
14         on driver.driverid = geo_event.driverId
15     left join `CDP_Kudu_Catalog`.`default_database`.`impala::employees_hr_kudu_impala_db.timesheet`
16         FOR  SYSTEM_TIME AS OF PROCTIME() timesheet
17         on (timesheet.driverid = geo_event.driverId  and timesheet_week = 1)
18       where
19     geo_event.eventTimestamp BETWEEN
20         speed_event.eventTimestamp - INTERVAL '1' SECOND AND
21         speed_event.eventTimestamp + INTERVAL '1' SECOND
22     AND geo_event.eventType <> 'Normal'
23     AND driver.wage_plan = 'hours'
24     AND timesheet.hours_logged > 45
```

**kafka** — Join 2 streaming event topics

**HIVE** — Enrich stream from data warehouse

**APACHE KUDU** — Enrich stream from real-time data mart

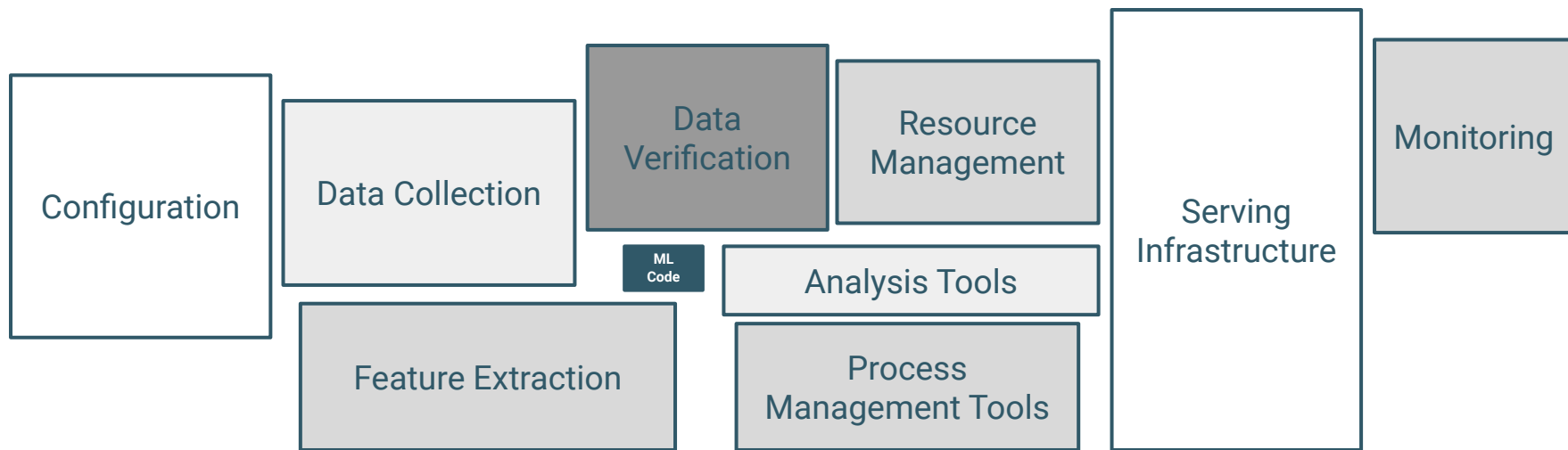Filter & transform in real-time

**APACHE KUDU** — Write streaming result to table or topic

**CLOUDERA**

# HIDDEN TECHNICAL DEBT IN MACHINE LEARNING SYSTEMS



Configuration

Data Collection

Data Verification

Resource Management

Monitoring

ML Code

Analysis Tools

Feature Extraction

Process Management Tools

Serving Infrastructure

Source: https://papers.nips.cc/paper/5656-hidden-technical-debt-in-machine-learning-systems.pdf

# 35%

## Making it to production

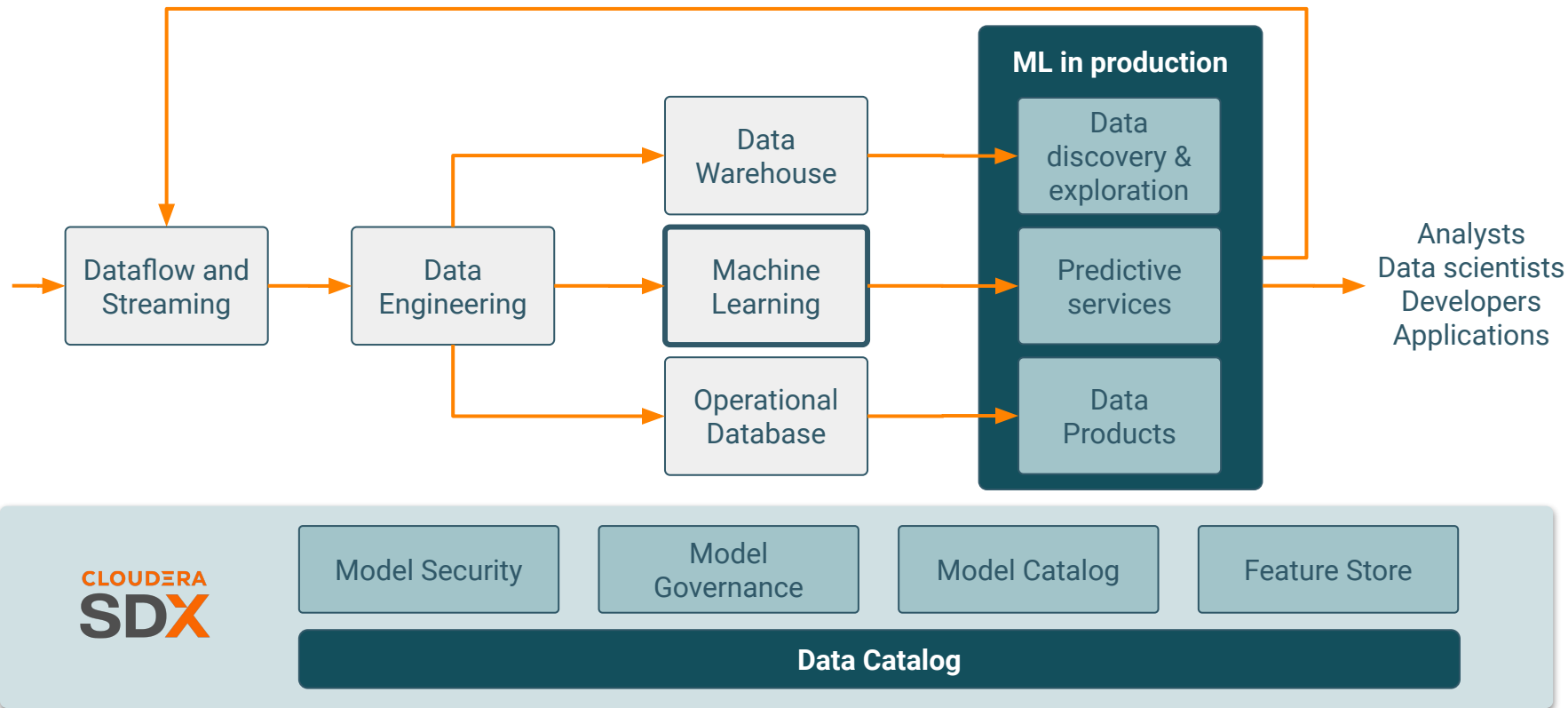Currently only 35% of organizations indicate that analytical models are fully deployed in production and are often challenged in the "last mile" of the complex and iterative ML workflow

*IDC's Advanced and Predictive Analytics survey and interviews, n = 400, 2017 – 2019*

# MACHINE LEARNING IN PRODUCTION

# MACHINE LEARNING IN PRODUCTION

# AGENDA

Mitől modern egy adatplatform?

**Adatmenedzsment a felhőben**

Az út a felhőbe

[…] on the other hand, we have the phenomenon we've outlined in this post, where the cost of cloud "takes over" at some point, locking up hundreds of billions of market cap that are now stuck in this paradox:

**You're crazy if you don't start in the cloud; you're crazy if you stay on it.**



VentureBeat

The Machine    GamesBeat    Jobs    Special Issue

Become a Member    |    Sign In

# The cost of cloud, a trillion dollar paradox

Sarah Wang, Andreessen Horowitz        Martin Casado, Andreessen Horowitz        June 4, 2021 3:05 PM

CLOUDERA

# CLOUD 101: "SEPARATION OF COMPUTE AND STORAGE"

On-premise
data lake architecture /
legacy data warehouse

Cloud-native
data lake architecture /
data services

COMPUTE

CONTEXT

STORAGE

Physical Servers

COMPUTE

Virtual Machine

STORAGE

Object Storage

# A DATA CLOUD SEPARATES STORAGE & COMPUTE *& CONTEXT*

On-prem
data lake architecture

Cloud optimized
data lake architecture

Basic cloud-native
data lake architecture

**CLOUDERA**

# CLOUD DATA LAKES REQUIRE …

## Identities

How do we easily connect to **corporate** identity?

## Schema

How do multiple personas **find** and **share** datasets across different compute engines?

## Policy

Is there a central place to protect sensitive data at a **field** and **row** level?

## Audits

Can I comply to the regulatory requirements with comprehensive **audits** and **lineage**?

Privacy, compliance & regulation

CLOUDERA

# ... AND *ACTIVE* MANAGEMENT

## Data Profiling & Stewardship

Can I detect **sensitive** data automatically, **tag** and **curate** for easy discovery and security controls?

## Replication

What happens to metadata and lineage, when I move data **across environments**?

## Workload Management

How do I track workloads for **troubleshooting** and optimization with elastic computes?

## Encryption

Is the data **protected** at rest and in-transit?

Context is required for all but the simplest use cases

- Multi-tenant
- Multi-workload
- Sensitive data
- Regulated use case
- Long Running Apps

# AGENDA

Mitől modern egy adatplatform?

Adatmenedzsment a felhőben

**Az út a felhőbe**

# CDP PUBLIC CLOUD PHILOSOPHY

Cloudera Data Platform Public Cloud enables you to:

- **leverage existing public cloud** environment including networking and data
- **migrate existing workloads** with minimal transformation effort
- **use existing user identities** to access cloud data and CDP services
- **modernize data applications** with CDP experiences

It provides a quick way to public cloud and integrates with current cloud infrastructure blueprints.

**Public cloud environment - customer VPC**

CLOUDERA

# 1) LEVERAGE AN EXISTING PUBLIC CLOUD ENVIRONMENT

Create a **centralized data lake** that provides security and governance services for all **data and services**.

The **SDX** (Shared Data Experience) layer provides unified services:

- schema / metastore
- authorization (fine-grain access to S3 as well as structured data)
- data catalog / metadata management with visual lineage
- audit service to track data access

All data and SDX metadata lives within a customer environment.



**Public cloud environment - customer VPC**

# 2) MIGRATE EXISTING WORKLOADS WITH MINIMAL EFFORT

Create fully secured and integrated **clusters-as-a-service** within minutes

This is best suited for **lift & shift** of existing on-premise Cloudera workloads or complex applications.

Data Hub uses **Cloudera Runtime 7.x** in the public cloud, the same platform codebase as CDP on-prem.

All data is **covered by SDX and is persisted on S3**, enabling seamless access from non-Cloudera apps.

**Public cloud environment - customer VPC**

**CLOUDERA**

# 3) USE EXISTING USER IDENTITIES TO ACCESS CDP SERVICES

**Public cloud environment - customer VPC**

Azure Active Directory

single sign-on

CDP provides centralized backend and management services (KDC, LDAP, DNS, certificates) for all services within a CDP **environment**.

Customers **bring existing user identities via single sign-on** (SAML).

CDP creates a unified identity for each user and automatically provides secure access to data, UIs, API endpoints and SQL interfaces.

**Environment**
- KDC
- LDAP
- DNS
- CA

Amazon EC2

S3 Storage

Amazon IAM

**SDX Data Lake**
- Ranger
- Atlas
- IDBroker
- Hive Metastore
- Cloudera Manager

Amazon EC2

Amazon RDS

S3 Storage

**KNOX**

**K N O X**

**Data Hub**

EC2    S3

Cloudera Runtime (all services)

In-built definitions

**Features:**
PaaS
Manual Scaling
Custom Templates

**Data Warehouse Experience**

EKS    S3

Hive (w LLAP)
Impala
DAS & Hue

**Features:**
PaaS
Fast Auto Scaling
JDBC connectivity

**Machine Learning Experience**

EKS    S3

ML Workbench
Spark
Python/R

**Features:**
PaaS
Fast Auto Scaling
GPU support

**Data Engineering Experience**

EKS    S3

Spark
Airflow
Safari (perf-mon)

**Features:**
PaaS
Fast Auto Scaling
Job Management

# 4) MODERNIZE DATA APPLICATIONS WITH CDP EXPERIENCES

Move beyond clusters and adopt CDP experiences, cloud services that provide **on-demand autoscaling**, **workload isolation** and a **redefined self-service user experience**:

Independent data warehouses and data marts that autoscale to meet workload demands with **CDW**.

Unified self-service data science and data engineering in a single, portable service with **CML**.

Hive and Spark jobs on an auto-scaling cluster scheduled with Apache Airflow with **CDE**.
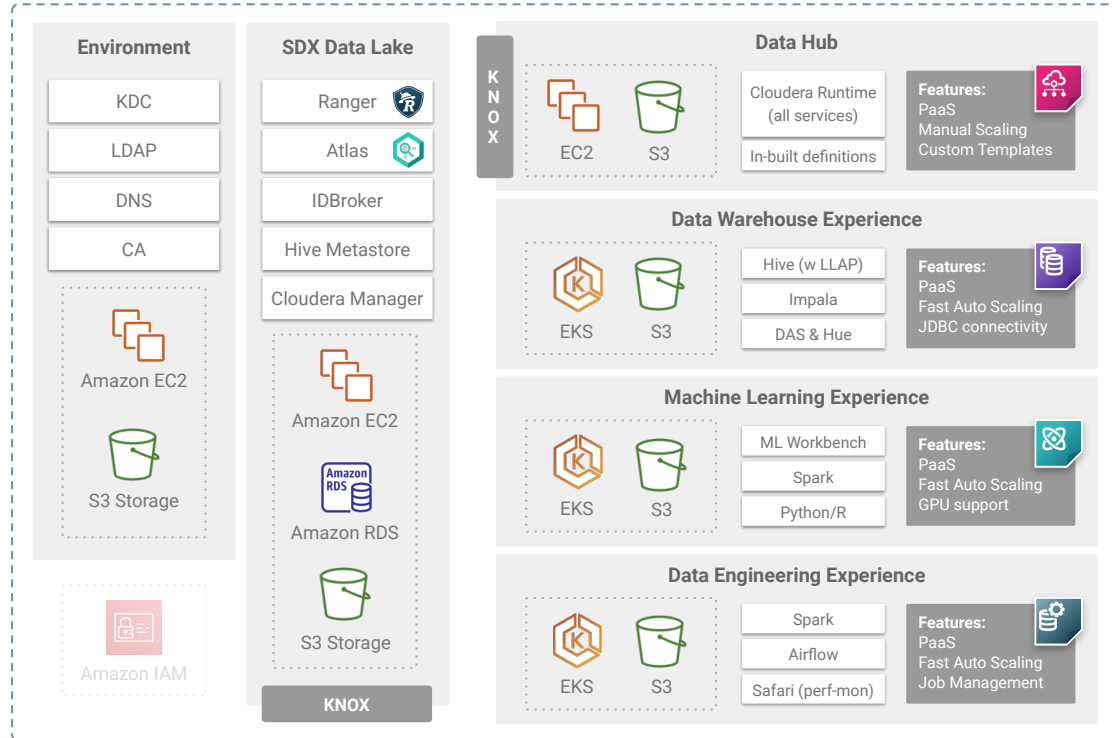
**Public cloud environment - customer VPC**



## Environment
- KDC
- LDAP
- DNS
- CA
- Amazon EC2
- S3 Storage
- Amazon IAM

## SDX Data Lake
- Ranger
- Atlas
- IDBroker
- Hive Metastore
- Cloudera Manager
- Amazon EC2
- Amazon RDS
- S3 Storage
- KNOX

## KNOX

### Data Hub
EC2  S3  Cloudera Runtime (all services)  In-built definitions
Features: PaaS Manual Scaling Custom Templates

### Data Warehouse Experience
EKS  S3  Hive (w LLAP)  Impala  DAS & Hue
Features: PaaS Fast Auto Scaling JDBC connectivity

### Machine Learning Experience
EKS  S3  ML Workbench  Spark  Python/R
Features: PaaS Fast Auto Scaling GPU support

### Data Engineering Experience
EKS  S3  Spark  Airflow  Safari (perf-mon)
Features: PaaS Fast Auto Scaling Job Management

**CLOUDERA**

# ... ALL MANAGED FROM A UNIFIED SINGLE PANE OF GLASS



**Public cloud environment - customer VPC**

User

CLI · Web UI · SDK

**Cloudera Cloud**

Data Catalog · DLM · Replication Manager · BDR · Workload Analytics · SAML Federation · Identity Management · ID Broker Mapping · Mgmt, Orchestration, & Monitoring Services · Control Plane APIs

**Cross Account Access IAM Role**

**Environment:** KDC, LDAP, DNS, CA, Amazon EC2, S3 Storage

**Amazon IAM**

**SDX Data Lake:** Ranger, Atlas, IDBroker, Hive Metastore, Cloudera Manager, Amazon EC2, Amazon RDS, S3 Storage, KNOX

**KNOX**

**Data Hub:** EC2, S3, Cloudera Runtime (all services), In-built definitions. Features: PaaS, Manual Scaling, Custom Templates

**Data Warehouse Experience:** EKS, S3, Hive (w LLAP), Impala, DAS & Hue. Features: PaaS, Fast Auto Scaling, JDBC connectivity

**Machine Learning Experience:** EKS, S3, ML Workbench, Spark, Python/R. Features: PaaS, Fast Auto Scaling, GPU support

**Data Engineering Experience:** EKS, S3, Spark, Airflow, Safari (perf-mon). Features: PaaS, Fast Auto Scaling, Job Management

**CLOUDERA**

... SUPPORTING MIGRATIONS AND A TRUE HYBRID EXPERIENCE

# ... ACROSS ANY CLOUDS

**Your AWS VPC(s)** · **Your Azure Resource Group(s)** · **Your GCP Project(s)**

**Cloudera Cloud** — Data Catalog · Replication Manager · Workload Analytics · Identity Management · Mgmt, Orchestration, & Monitoring Services · Control Plane APIs
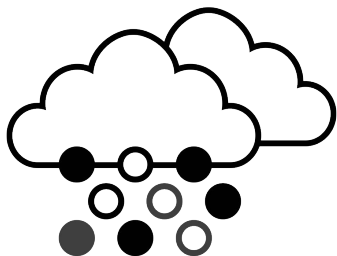
# DEMO

CLOUDERA

CLOUDERA

THE ENTERPRISE DATA CLOUD COMPANY

Any Cloud

Data Lifecycle

Secure & Governed

Open

# THANK YOU

**CLOUDERA**