Intel Accelerated Solutions
Solution Snapshot

# Modernizing the Data Pipeline with Cloudera Data Platform and Intel® Technology

## The Challenge

Businesses continue to struggle with gaining meaningful insights from their data. The explosion of data volume to the prevalence of data and analytics silos increases these challenges even more. These need to be addressed along with effectively managing and securing an organization's data to prevent costly privacy breaches.

## Cloudera Overview + Benefits

As data grows and analytics workloads get more complex, analytics workloads need to be flexible and highly responsive. Cloudera Data Platform (CDP) is a hybrid data platform designed for unmatched freedom to choose—any cloud, any analytics, any data. CDP delivers faster and easier data management and data analytics for data anywhere, with optimal performance, scalability, and security. With CDP, you get all the advantages of CDP Private Cloud and CDP Public Cloud for faster time to value and increased control.

## Cloudera Data Platform

### Use Cases

### Deploying a Private Cloud

CDP delivers an end-to-end analytics platform that supports hybrid and multi-cloud environments and eliminates costly and inefficient data silos.

### Predictive Analytics and Machine Learning

End-to-end machine learning pipeline with Cloudera and open-source frameworks.
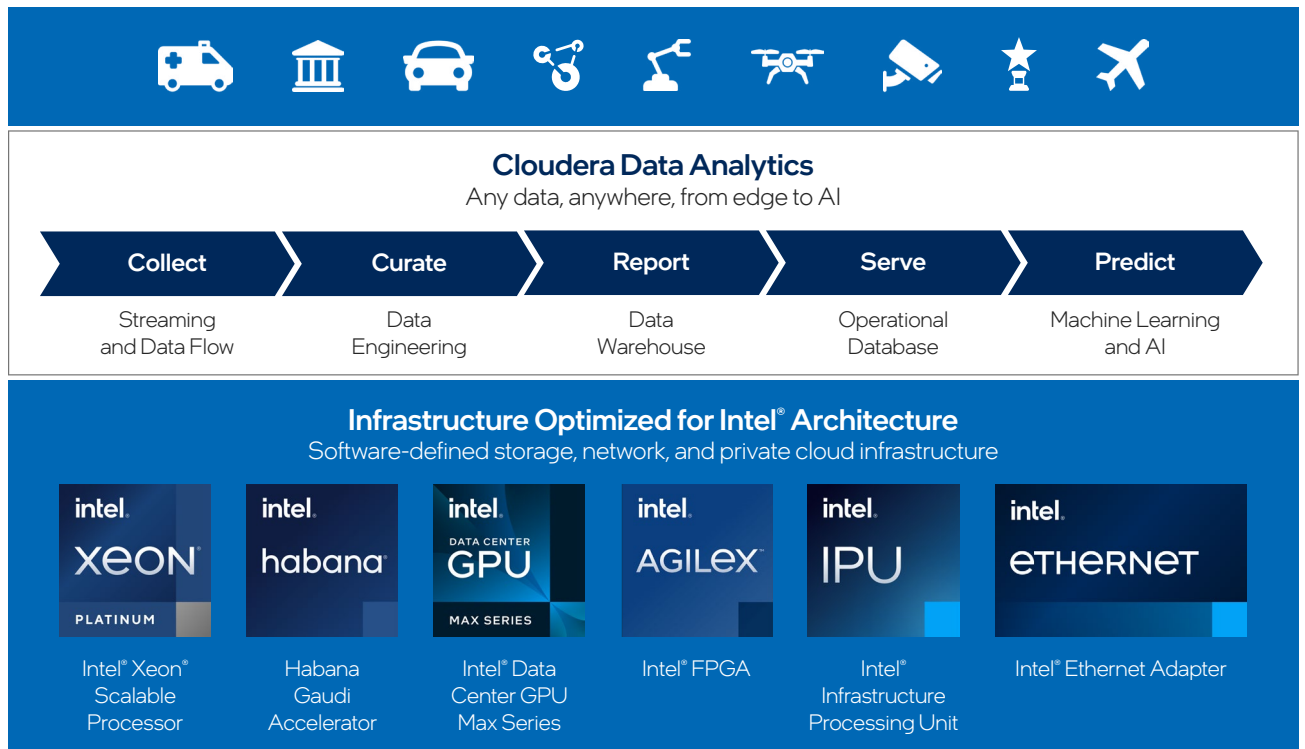
### Modern Data Pipeline

Apache Airflow-based pipeline provides orchestration for multi-step data pipelines in the cloud. The pipeline is based on a combination of Spark and Hive and generates curated datasets for downstream applications securely and efficiently.

### Data Lake/Lakehouse Cloudera

Open Data Lakehouse helps organizations run quick analytics on all data—structured and unstructured—at a massive scale.
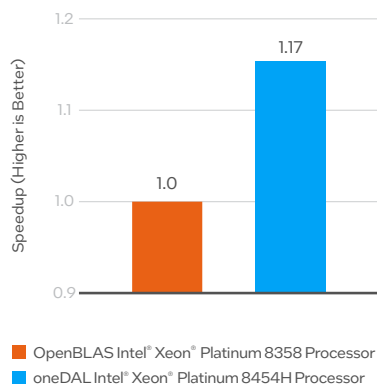
## Applications and Solutions

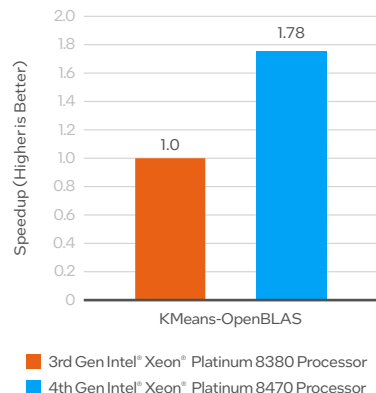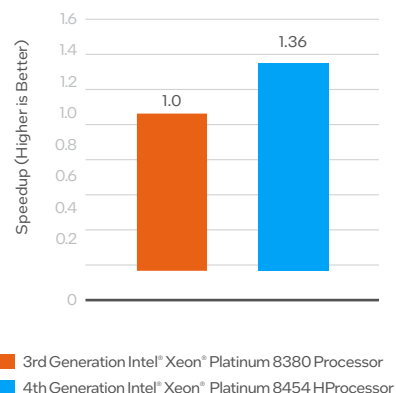Analytics-Powered Vertical and Horizontal Solutions



### Cloudera Data Analytics
Any data, anywhere, from edge to AI

| Collect | Curate | Report | Serve | Predict |
|---------|--------|--------|-------|---------|
| Streaming and Data Flow | Data Engineering | Data Warehouse | Operational Database | Machine Learning and AI |

### Infrastructure Optimized for Intel® Architecture
Software-defined storage, network, and private cloud infrastructure

| intel® XEON PLATINUM | intel® habana | intel® DATA CENTER GPU MAX SERIES | intel® AGILEX | intel® IPU | intel® ETHERNET |
|---|---|---|---|---|---|
| Intel® Xeon® Scalable Processor | Habana Gaudi Accelerator | Intel® Data Center GPU Max Series | Intel® FPGA | Intel® Infrastructure Processing Unit | Intel® Ethernet Adapter |

# Proof Points

### 3rd Gen Intel® Xeon® Platinum 8358 Processor versus 4th Gen Intel® Xeon® Platinum 8454H Processor[1]



Speedup (Higher is Better)

- 1.0
- 1.17

Legend:
- OpenBLAS Intel® Xeon® Platinum 8358 Processor
- oneDAL Intel® Xeon® Platinum 8454H Processor

### HiBench KMeans Performance Gain for 4th Gen Intel® Xeon® Scalable Processor[2]
(1 Master + 3 Worker Node)



KMeans-OpenBLAS

- 1.0
- 1.78

Legend:
- 3rd Gen Intel® Xeon® Platinum 8380 Processor
- 4th Gen Intel® Xeon® Platinum 8470 Processor

### Performance for Spark KMeans Clustering using OpenBLAS Library on Cloudera Private Cloud Base Cluster[3]



- 1.0
- 1.36

Legend:
- 3rd Generation Intel® Xeon® Platinum 8380 Processor
- 4th Generation Intel® Xeon® Platinum 8454 H Processor

**Up to 17% performance gain** on 4th Gen Intel® Xeon® Platinum 8454H processor versus 3rd Gen Intel® Xeon® Platinum 8358 processor using Spark CDP Private Cloud Base with MLlib and Intel® oneAPI Data Analytics Library[1]

**Up to 78% performance gain** on Intel® Xeon® Platinum 8470 processor when compared to Spark to Intel® Xeon® Platinum 8380 processor for Spark KMeans algorithm using Spark CDP Private Cloud Base 7.1.7 and OpenBLAS library instances[2]

**Up to 36% performance gain** for KMeans Clustering (ML) using Spark CDP Private Cloud Base 7.1.7 on 4th Gen Intel® Xeon® Scalable processor versus 3rd Gen Intel® Xeon® Scalable processor[3]

## Why Intel for Cloudera Data Platform?

### Faster Machine Learning Algorithms

Cloudera Data Platform (CDP) uses Intel® accelerators, libraries, and instruction sets to enable several predictive analytics and machine learning workloads. These include Intel® Data Streaming Accelerator, Intel® In-Memory Analytics Accelerator (IAA), Intel® QuickAssist Technology (QAT), and Intel® Advanced Vector Extensions 512. There are also innovative, open-source technologies that are a part of the overall solution.

### Reduce Complexity of Analytics Deployments

Cloudera and Intel joint solutions are workload-optimized to minimize the challenges of infrastructure evaluation and deployment. They accelerate deployment with validated hardware and software stacks. These solutions balance the needs of performance, scale, and availability to reduce deployment risk.

### Easy Data Management

Cloudera Shared Data Experience (SDX), running on Intel® technologies, combine enterprise-grade centralized security, governance, and management capabilities with shared metadata and a data catalog to eliminate costly data silos, prevent lock-in to proprietary formats, and eradicate resource contention. SDX keeps your information secure by design with an integrated set of security and governance technologies to protect data. Security is further strengthened with Intel® Software Guard Extensions (SGX), which isolates sensitive data with hardware-based memory protections and Intel® Trusted Domain Extensions (TDX), used in virtualization environments where the entire virtual machine is isolated. Cloudera users can use these features to securely migrate their sensitive CDP Private Cloud workloads to CDP Public Cloud.

### Solutions for Data-Intensive Workloads

Optimized Cloudera software applications combined with proven Intel® architecture are key to building highly performant data applications. They support all kinds of data ingested from cloud to edge and are robust enough to handle the challenges of modern workloads. The open and flexible platforms provide interoperability with existing tools, storage efficiency, and scale.

## Want More Information?

To get more information on Intel and Cloudera joint solutions, visit www.cloudera.com/partners/solutions/intel.html or www.intel.com/cloudera

**intel.**