cloudera | intel

# UNLOCK THE POWER OF
# ON-PREMISES BIG DATA ANALYTICS

**Transform complex data into clear, and actionable insights with a high-performance big data analytics solution from Cloudera and Intel**

## Is This Solution for You?

Do you...

- ✓ Need to integrate isolated data silos into a cohesive data lake?
- ✓ Have large amounts of on-premises data?
- ✓ Want to reduce total costs by replacing expensive proprietary data warehouses?
- ✓ Need a platform for easy and fast distributed data analytics and data management?
- ✓ Need to consolidate hardware sprawl?

**Learn about the Business Story →**
**Learn more about the Reference Architecture →**

# BUSINESS STORY
## Unlock the Power of On-Premises Big Data Analytics

## Authors

**Ali Bajwa**
Director, Partner Engineering, Cloudera

**Amandeep Raina**
Cloud Solution Engineer, Data Center, Intel

**Priyanka Sebastian**
Cloud Solution Engineer, Data Center, Intel

**Sandeep Togrikar**
Solutions Architect for Analytics & AI, Data Center, Intel

**Lifeng Wang**
Machine Learning Engineer, Intel Architecture Graphics & Software, Intel

**Ning Wang**
Machine Learning Engineer, Intel Architecture Graphics & Software, Intel

**cloudera**

## Executive Summary

Where are you on your data journey? Do you struggle with integrating isolated silos of data throughout your data center? Is decision making hampered by slow and cumbersome data analytics? You're not alone. According to Gartner, 91 percent of organizations struggle to reach data maturity.[1] But don't despair—collaboration between Intel and Cloudera has created a big data analytics platform specifically designed for large-scale on-premises workloads.

Cloudera Enterprise is characterized by a shared data experience, is powered by open source software, and offers multi-function analytics on a unified platform that eliminates silos and speeds the discovery of data-driven insights. This big data analytics platform helps streamline data management and workload orchestration. The Intel® architecture underlying Cloudera Enterprise provides the power that big data analytics demand. Intel and Cloudera have worked together to improve compute performance, storage efficiency, artificial intelligence (AI) acceleration, and more. The result is an on-premises data platform that is built to meet today's big data analytics needs and that can scale to meet the needs of the future. Tests show that a modern version of Cloudera Enterprise running on the latest Intel® hardware can **improve big data analytics performance by up to 79 percent.**[2] Extending system memory with Intel® Optane™ persistent memory provides **additional performance improvements.**[3]

This document provides a business-level overview of Cloudera Enterprise,[4] describes a reference architecture for deployment, and highlights the platform's performance and scalability.

## Cloudera Enterprise Solution Benefits

- Open source platform provides flexibility and interoperability with existing tools.
- Comprehensive set of management tools simplifies cluster configuration and scaling.
- The combination of Cloudera Enterprise and the latest Intel® technology improves big data analytics performance:
  - Boost performance by up to 79 percent, compared to an older version of Cloudera Enterprise running on a previous-generation Intel® Xeon® processor.[5]
  - Use Intel® Optane™ persistent memory to improve Spark SQL performance/$ by 1.26x to 3x, depending on node count and workload.[6]
- Nearly linear scaling of the platform means that as your data grows, Cloudera Enterprise can keep up.

## About Cloudera

Cloudera, Inc. is a U.S.-based software company that provides a software platform for data engineering, data warehousing, machine learning, and analytics that runs in the cloud or on-premises. Founded in 2008, the company is committed to accelerating enterprise-class data management innovation to make what is impossible today, possible tomorrow. Cloudera solutions empower enterprises to transform complex data into clear and actionable insights.

Powered by the relentless innovation of the open source community, Cloudera has offices around the globe and is headquartered in Palo Alto, California.

## Business Challenge: Transforming Data Silos into a Comprehensive Data Lake

Data is everywhere. It's constantly being generated by machines, customers, and applications. It piles up in various data warehouses. And hidden in that data are insights about your business—information about network security, customer preferences, supply chain dynamics, and more. But when data exists in silos, and data analytics runs in those same silos, it is nearly impossible to establish a coherent approach to data mining, data privacy, and intellectual property protection. What's more, most machine-learning and analytics tools require proprietary storage and algorithms—leading to vendor lock-in and lack of agility. And as data continues to grow and big data analytics workloads increase accordingly, your infrastructure must be able to scale and maintain the required cluster performance without driving up costs. Cloudera Enterprise is an open source, scalable data platform that is optimized to run on high-performance Intel® hardware. It can help you quickly and cost-efficiently extract value from your data (see the Solution Value section).

## Use Cases Abound for an On-Premises Data Analytics Platform

Industries, from manufacturing to healthcare to retail, from transportation to hospitality to life science, are under pressure to turn their data into an asset. Cloudera Enterprise can be used across several broad use cases, including data lakes; extract, load, and transform (ELT) applications; and offloading analytics from expensive proprietary databases. And while the trend is to take data to the cloud (and Cloudera does support cloud deployments), in some cases it makes more sense to keep data on-premises. For example, the data may be sensitive (such as intellectual property). Or the analytics use case may require extremely low latency (such as real-time fraud detection).

The following examples barely scratch the surface of how you can put Cloudera Enterprise to work in your data center, but they illustrate the flexibility of the solution:

- **Manufacturing.** Send Internet of Things (IoT) data—much of it semi-structured—into a data lake, then run analytics on that data for predictive maintenance, real-time production line changes, and more.
- **Sales and marketing.** Use your data to optimize marketing campaigns and create advanced recommendation engines.
- **Healthcare and life sciences.** Use a data lake to store unstructured data (phone call recordings, webinar transcripts, imaging data, etc.) and make that data available to academic researchers.
- **Financial services.** Use ELT to detect fraud, predict customer churn, or perform risk management.
- **Supply chain management.** Use data science to identify cost savings opportunities, speed the costing cycle, perform historical pricing analysis, and conduct "what if" analysis for procurement and planning.
- **Transportation.** Gather data from sensors and run real-time analytics that provide input to autonomous cars, or analyze images for surveillance and safety systems.
- **Cybersecurity.** Pour all of your networking data into a data lake and run analytics to detect vulnerabilities and threats.

## Solution Value: High Performance, Ultimate Flexibility, and Excellent Scalability

Cloudera Enterprise delivers an integrated suite of analytic engines ranging from stream and batch data processing to data warehousing, operational database, and machine learning (see Figure 1). By using Cloudera, enterprises gain end-to-end big data capabilities—ingest, store, process, and analyze for insights. Cloudera SDX applies consistent security and governance, enabling users to share and discover data for use across many demanding big data analytics workloads.
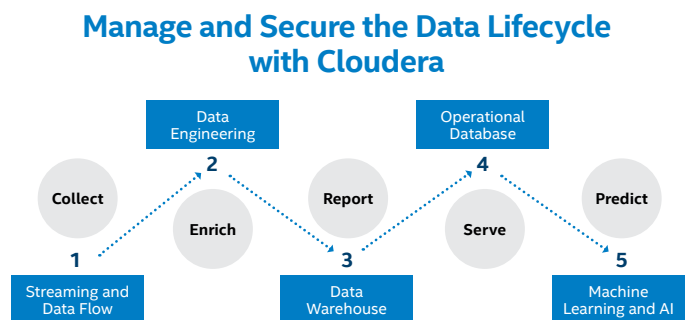
### Manage and Secure the Data Lifecycle with Cloudera



**Figure 1.** Cloudera Enterprise simplifies data management, offers a wide variety of analytic engines, and provides a shared data experience (SDX).

Some specific business benefits from Cloudera include:

- **Ultimate flexibility.** It's easy to integrate Cloudera Enterprise with existing infrastructure and tools, and the platform offers robust security, governance, data protection, and data management features. Cloudera Enterprise is open source, which allows for flexibility in choosing the technologies you want to use, without vendor lock-in. What's more, the open source community drives the evolution of data management and advanced analytics. Intel and Cloudera's strong relationships with a broad portfolio of data center solution providers can help streamline the process for building solutions.

- **Excellent scalability.** The Cloudera Enterprise reference architecture from Intel is the foundation for a highly scalable big data analytics platform that can store any amount or type of data in its original form—and keep it for as long as it's needed. Intel testing shows that Cloudera Enterprise performance on 2nd Generation Intel® Xeon® Scalable processors scales nearly linearly.[7] Also, Cloudera Enterprise includes software that simplifies scaling by automating node configuration.

Intel and Cloudera have participated in joint engineering, resulting in optimizations of Cloudera Enterprise pertaining to faster compute performance, increased storage efficiency, enhanced security, superb AI support, and great query performance. These optimizations let Cloudera take advantage of Intel® Optane™ persistent memory, Intel Xeon Scalable processors, Intel® FPGAs, and Intel® QuickAssist Technology, using Intel® Advanced Vector Extensions 512, Intel® Math Kernel Library, Intel® AES-NI, and Intel® Intelligent Storage Acceleration Library. With Cloudera Enterprise running on 2nd Gen Intel Xeon Scalable processors, **big data analytics job performance can improve by as much as 79 percent**, compared to an older version of Cloudera Enterprise running on a previous-generation Intel Xeon processor.[8] What's more, you can increase **Spark SQL performance by up to by extending system memory with Intel Optane persistent memory.**[9]

Additional benefits of Cloudera Enterprise include security-by-design—with encryption, access controls, and governance and lineage—as well as a single pane of glass for cluster administration, automation, management, and security.

## Solution Architecture: Scalable and Versatile Big Data Analytics Platform

Cloudera Enterprise is comprehensive; it can handle all the stages of data analytics, from data ingestion through storage and processing, to delivery of insights that fuel intelligent action. Cloudera Enterprise also offers the following:

- **Enhanced security features,** including perimeter security, authentication, granular authorization, and data protection

- **Governance features,** such as enterprise-grade data auditing, data lineage, and data discovery

- **Manageability features,** including native high-availability; fault-tolerance and self-healing storage; automated backup and disaster recovery; and advanced system and data management
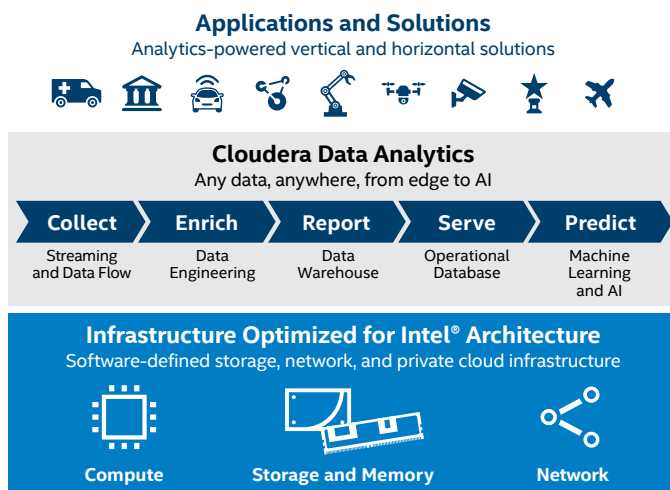


**Figure 2.** This high-level diagram shows the basic Cloudera Enterprise solution architecture.

Cloudera Enterprise (see Figure 2) is built on and optimized for Intel® compute, storage, and networking architecture. The reference architecture for Cloudera Enterprise incorporates the following Intel technologies:

- **2nd Generation Intel Xeon Scalable processors.** These processors are optimized for big data analytics workloads like Hadoop. They incorporate architecture improvements and enhancements for compute-intensive and data-intensive workloads, making them well suited for ingesting and analyzing massive quantities of data.

- **Intel® 3D NAND SSDs.** These PCIe/NVMe-based SSDs deliver scalable, cost-effective performance and low latency. The SSDs also offer outstanding quality, reliability, advanced manageability, and serviceability to minimize service disruptions.

- **Intel® Ethernet network connection.** Intel Ethernet network controllers, adapters, and accessories enable agility in the data center to deliver services efficiently and cost-effectively. Compatible with the Open Compute Platform, these high-performance connectors support high throughput, reliability, and compatibility.

- **Intel® Optane™ technology.** The first breakthrough in memory and storage in 25 years, Intel Optane persistent memory and Intel Optane SSDs are a unique innovation that bridges critical gaps in the storage and memory hierarchy, delivering persistent memory, large memory pools, fast caching, and fast storage. For example, Intel Optane persistent memory can be used as the Hadoop Distributed File System (HDFS) read cache, replacing expensive DRAM cache and providing a larger capacity.

➔ **Ready to learn more?** Turn the page for the detailed Reference Architecture discussion.

# REFERENCE ARCHITECTURE
## Unlock the Power of On-Premises Big Data Analytics

## Summary of Findings

- Cloudera Enterprise performance improves significantly—by as much as 79 percent—when running on the latest Intel® processors and Intel® Solid State Drives, compared to an older version of CDH running a previous-generation Intel® Xeon® processor.[10]

- Performance scales linearly as dataset size and the number of worker nodes increase.

- By adding Intel® Optane™ persistent memory to the platform, we observed a 1.26x to 3x improvement in Spark SQL performance/$ (depending on node count and workload), compared to an all-DRAM configuration.[11]

## Table of Contents

## Overview of Cloudera Enterprise

Cloudera Enterprise provides a scalable, versatile, and integrated platform that simplifies managing the growing volume and variety of data in your enterprise. Using Cloudera products and solutions can help you deploy and manage Apache Hadoop and related projects, manipulate and analyze your data, and enhance data security.[12]

At the center of Cloudera Enterprise is the Cloudera Distribution of Hadoop (CDH). Hadoop is an open source Apache project built for distributed big data analytics. Besides CDH, Cloudera Enterprise includes many other popular open source software components with enterprise capabilities for advanced system and data management. The platform also offers dedicated support and community advocacy from its team of Hadoop developers and experts. See Table 1 for a full list of Cloudera Enterprise software components.

**Table 1.** Cloudera Enterprise Software Components

| PRODUCT/SERVICE | DESCRIPTION |
|---|---|
| Cloudera Distribution of Hadoop (CDH) 6.2 | CDH is Cloudera's fully open source platform distribution—including Apache Hadoop—that's built specifically to meet enterprise demands. CDH delivers everything you need for enterprise use right out of the box. By integrating Hadoop with more than a dozen other critical open source projects, Cloudera has created a functionally advanced system that helps you perform end-to-end big data workflows. |
| Cloudera Manager | Cloudera Manager helps you more easily deploy, manage, monitor, and troubleshoot issues with a Hadoop cluster, to enhance platform scalability. |
| Cloudera Support | This feature provides technical support for Hadoop, which can help increase uptime and diagnose problems quickly. (Note: Not available in the out-of-the-box free version of CDH.) |
| Apache Hive | Hive is a data warehouse software project built on top of Hadoop for providing data query and analysis. Hive uses an SQL-like interface to query data stored in databases and file systems that integrate with Hadoop. |
| YARN | YARN is the resource management and job scheduling component in the Hadoop framework. It allocates system resources to applications running in a Hadoop cluster and schedules tasks to be executed on different cluster nodes. |
| Apache ZooKeeper | ZooKeeper is a centralized service for maintaining configuration information, naming, providing distributed synchronization, and providing group services. |
| Apache Spark | Spark is a unified analytics engine for large-scale data processing. It delivers fast, in-memory analytics and real-time stream processing. |
| HBase | HBase is a column-oriented non-relational database management system that runs on top of HDFS. |

## Reference Architecture Overview

Overall, this reference solution is an effective big data extension to an enterprise data warehouse (EDW) analytics platform that offers the following:

- Scalability and flexibility

- Excellent performance and total cost of ownership (TCO)

- Ability to satisfy business requirements for service-level agreements (SLAs), multiple users, and future growth

The solution can be applied to both green-field and refresh deployment scenarios, and is designed to be used to its fullest extent before scaling out by adding more nodes.

### Node Definitions

We have defined a highly available reference architecture that includes the following nodes:

- **Three management nodes.** One management node is for Cloudera Manager; the other two are master nodes (NameNode and Resource Manager). Management nodes support services that are needed for the cluster operation. Even if one management node fails, the other two can pick up the extra work.

- **10 worker nodes.** This reference architecture uses an Intel SSD on every worker node. Worker nodes handle the bulk of the Hadoop processing. The number of worker nodes necessary depends on data set size. Depending your scalability needs, you can replicate this configuration to 20 or 30 worker nodes, or even cut it in half to five worker nodes. Some of our testing was performed with four worker nodes, but 10 worker nodes represent a nominal workload configuration.

## Solution Architecture

Figure 3 illustrates the entire solution architecture, while Table 2 provides the bill of materials; see Appendix A for software requirements. The Base configuration uses all DRAM; the Plus configuration adds 1.5 TB of Intel® Optane™ persistent memory for additional performance enhancement.
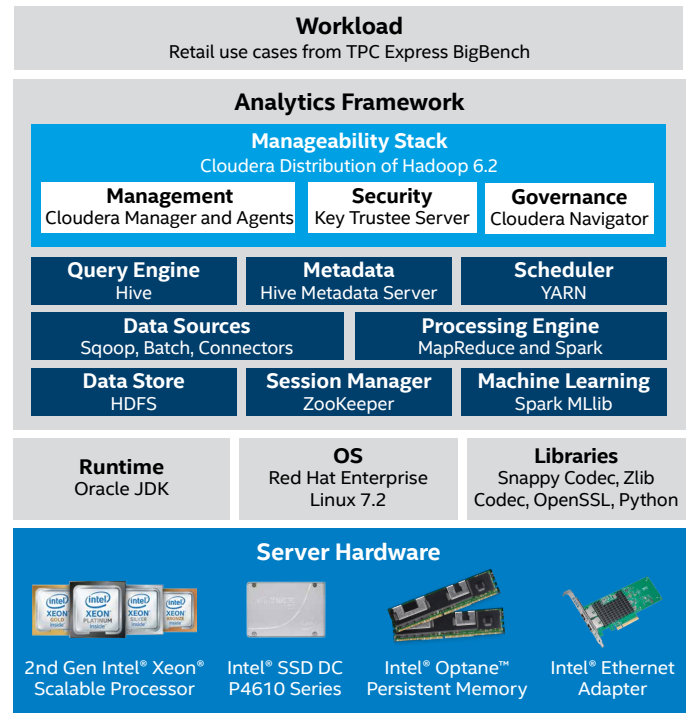


**Figure 3.** Reference architecture for Cloudera Enterprise, running on Intel® hardware.

**Table 2.** Bill of Materials

| COMPONENT (PER NODE) | DESCRIPTION | COMMENTS | REQUIRED/ RECOMMENDED |
|---|---|---|---|
| PROCESSOR | | | |
| **Management Nodes** | Intel® Xeon® Gold 6248 processor @ 2.5/3.9 GHz | – 1x Master Node/Active NameNode<br>– 1x Utility Node/Standby NameNode<br>– 1x Utility Node/Management Node | REQUIRED |
| **Worker Nodes** | Intel Xeon Gold 6248 processor @ 2.5/3.9 GHz | At least 4 worker nodes, can scale as necessary | REQUIRED |
| MEMORY | | | |
| **Base Configuration** | 12x 32 GB DDR4 @ 2666 MHz | Total 384 GB per node | REQUIRED |
| **Plus Configuration** | – 12x 32 GB DDR4 @ 2666 MHz *and*<br>– 12x 128 GB Intel® Optane™ persistent memory | Total 1.884 TB per node | RECOMMENDED |
| NETWORK | | | |
| **Ethernet Adapter** | 1x Intel® Ethernet Adapter X722 (10 GbE) | | REQUIRED |
| STORAGE | | | |
| **OS per Management Node** | Intel® SSD DC S4500 Series 960 GB | | REQUIRED |
| **HDFS per Management Node** | 2x 4 TB HDD 3.5-inch SAS3 12 Gb/s 7200RPM Seagate Enterprise V.5 ST4000NM0095 | Total 8 TB per management node | REQUIRED |
| **HDFS per Worker Node** | 8x 4 TB HDD 3.5-inch SAS3 12 Gb/s 7200RPM Seagate Enterprise V.5 ST4000NM0095 | Total 32 TB per worker node | REQUIRED |
| **One per Worker Node** | Intel® SSD DC P4610 Series 1.6 TB, 2.5-inch PCIe 3.1 x4, 3D2, TLC | For YARN tmp files and Spark shuffles | RECOMMENDED |

The appropriate number of worker nodes in the cluster depends on the dataset size. For smaller datasets, four worker nodes may be sufficient; for larger datasets, you may need 10 worker nodes. For distributed processing, you could scale out to 20 or 30 worker nodes, if necessary. According to our observations, a HDD plus one 3D NAND SSD per worker node is more cost effective compared to using a dense, all-SSD-based solution.

## Big Data Analytics Process

❶ **Ingest.** This reference architecture supports high-volume data ingestion of structured and unstructured data from various sources such as transactional relational database management systems, operations data, web logs, click streams, and other external sources.

❷ **Prepare.** Once ingested, the data is cleaned and formatted with metadata and schema.

❸ **Analyze.** Next, the data is loaded into the analytical data warehouse with shared local storage applied with pre-defined use-case logic for sorting and combining functions on distributed computing nodes.

❹ **Act.** Finally, results in the form of compressed datasets are made available for business needs for consumption to run reporting, machine-learning models, and business intelligence. In addition, the solution is flexible enough to support ad-hoc analysis of data when there are no pre-defined use cases.

## Test Results: Intel® Technology Significantly Improves Big Data Analysis Performance

Our tests used the BigBench benchmark, which is derived from the TPCx-BB Express Benchmark BB (TPCx-BB)—an industry-standard benchmark licensed under TPC that measures the performance of Hadoop-based big data systems. BigBench measures the performance of both hardware and software components by executing 30 frequently performed analytical queries in the context of retailers with physical and online store presence (see Appendix B). The queries are expressed in SQL for structured data and in machine-learning algorithms for semi-structured and unstructured data. The SQL queries can use Hive or Spark, while the machine-learning algorithms use machine-learning libraries, user-defined functions, and procedural programs.

## Newer Software and Hardware Improves CDH Throughput

When we compared the performance of CDH version 6.2 running on the Intel Xeon Gold 6248 processor to CDH 5.16.1 running on the Intel Xeon processor E5-2680 v3, normalized performance improved by up to 65 percent for two streams and by up to 79 percent for four streams (see Figure 4).[13] Improved throughput on more streams is due to the higher number of cores being able to support the parallel queries/multiple streams.

**BigBench Performance**[13]
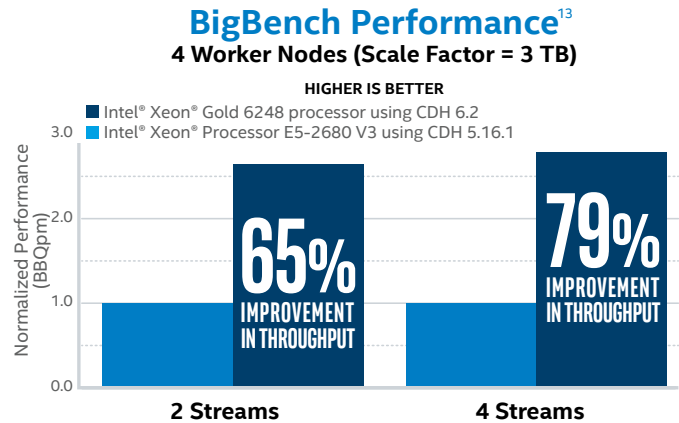**4 Worker Nodes (Scale Factor = 3 TB)**

**Figure 4.** Performance improvements on a four-worker-node cluster range from 65 to 79 percent when running a newer version of CDH on the latest Intel® Xeon® Scalable processor and Intel® 3D NAND solid state drives.

As we increased the number of worker nodes from four to 10, to accommodate an increased BigBench scale factor (from 3 TB to 10 TB), we found that the performance scaled almost linearly, as shown in Figure 5.[14]
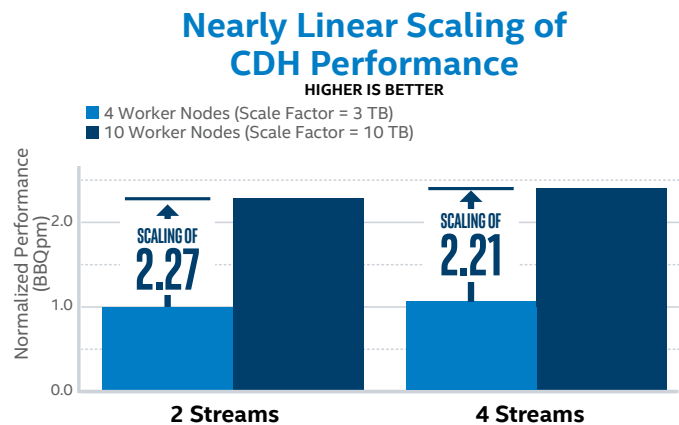
**Nearly Linear Scaling of CDH Performance**

**Figure 5.** CDH performance scales almost linearly as we add worker nodes.

**Table 3.** Performance and Performance/$ Benefits from Adding Intel® Optane™ Persistent Memory to the Spark Cluster (normalized results)

|  | DRAM<br>10 Nodes, 4 TB Total Memory | INTEL® OPTANE™ PERSISTENT MEMORY<br>4 Nodes, 8 TB Total Memory | 10 Nodes, 20 TB Total Memory | 4 Nodes, 8 TB Total Memory |
|---|---|---|---|---|
| Queries | Baseline | 91 Standard Spark Queries | 91 Standard Spark Queries | 9 I/O-intensive Spark Queries |
| Performance | 1.0 | 0.75 | 2.1 | 1.87 |
| Cost | 1.0 | 0.6 | 1.5 | 0.6 |
| Performance/$ | 1.0 | 1.26 | 1.41 | 3.0 |
| Dataset size | 10 TB | 10 TB | 10 TB+ | 10 TB |

## Intel Optane Persistent Memory Improves System Performance

Adding Intel Optane persistent memory to your Cloudera clusters can improve system performance, both for Spark SQL and for HBase.

### Spark SQL with Optimized Analytics Package for Spark Platform (OAP) Performance

Many big data workloads are memory-bound; in an ideal world, you could just add more DRAM to avoid memory bottlenecks. However, DRAM is an expensive resource and adding more may not be financially feasible. The good news is that Intel Optane persistent memory can be an affordable method of expanding system memory. If you are running real-time queries, expanding memory can improve query response time. If you are running batch jobs, expanded memory can help ensure the job completes within a certain window of time.

For real-time workloads, we compared Spark SQL performance on an all-DRAM configuration to a configuration that included 1.5 TB of Intel Optane persistent memory. We observed benefits from the addition of Intel Optane persistent memory in three scenarios[15] (see Table 3).

- **Consolidation.** You can shrink the cluster size by 60 percent (from 10 nodes to 4 nodes) and still get better performance/$.

- **Standard Spark queries.** You can keep the 10-node cluster and double Spark performance and increase Spark performance/$ by up to 1.41x.

- **I/O-intensive Spark queries.** With a 10-node cluster, performance/$ increases by almost 3x.

### HBase Performance

Starting with CDH 6.2, Cloudera now includes the ability to use Intel Optane persistent memory as an alternate destination for the second tier of the bucket cache. This deployment configuration enables you to have approximately three times more cache for a similar cost (as compared to off-heap cache on DRAM). It does incur some additional latency compared to the traditional off-heap configuration, but our testing indicates that by allowing more (if not all) of the data's working set to fit in the cache, this system configuration results in a net performance improvement when the data is ultimately stored on HDFS (using hard disk drives). For a full discussion of using Intel Optane persistent memory with HBase, read the Cloudera blog and the Intel Builders article.

Find the solution that is right for your organization. Contact your Intel representative or visit **intel.com/AI.**

### Learn More

You may find the following references helpful:

**Intel**
- 2nd Generation Intel® Xeon® Scalable processors
- Intel® Optane™ persistent memory
- Intel® Optane™ solid state drives
- Intel® Solid State Drives Data Center Family
- Intel® Ethernet Network Connection

**Cloudera**
- Cloudera Enterprise
- Cloudera Distribution of Hadoop

## Appendix A: Additional Configuration Details

### Required Software

Table A1 provides the software required for the reference architecture.

**Table A1.** Required Software

| COMPONENT | VERSION | DESCRIPTION |
|---|---|---|
| **Cloudera Distribution of Hadoop** | 6.2 | |
| **HDFS** | 3.0.0 | File System |
| **Hive** | 2.1.1 | Query Engine |
| **Spark** | 2.4.0 | Execution Engine |
| **YARN** | 3.0.0 | Resource Allocator/Scheduler |
| **ZooKeeper** | 3.4.5 | Cluster Management |

### Cloudera Enterprise Role Distribution

Table A2 provides information about which Cloudera Enterprise roles are present on which nodes. For more information, read the Cloudera article.

**Table A2.** Cloudera Enterprise Role Distribution

| CLOUDERA ROLE | MASTER NODE/ACTIVE NAMENODE | UTILITY NODE/STANDBY NAMENODE | UTILITY NODE/ MANAGEMENT NODE | WORKER NODES |
|---|---|---|---|---|
| **HDFS NameNode** | Y | | | |
| **HDFS Secondary NameNode** | | Y | | |
| **HDFS DataNode** | | | | Y |
| **YARN Resource Manager** | Y | | | |
| **YARN Job History Server** | Y | | | |
| **YARN Node Manager** | | | | Y |
| **HiveServer2** | Y | | | |
| **Hive Metastore Server** | | Y | | |
| **ZooKeeper** | Y | Y | Y | |
| **Hive Gateway** | Y | Y | Y | Y |
| **Spark Gateway** | Y | Y | Y | Y |
| **Spark History Server** | Y | | | |
| **Cloudera Manager** | | | Y | |
| **Cloudera Manager Management Service** | | | Y | |

## Appendix B: BigBench Use Case Descriptions

Refer to the following resources for descriptions of the parameters listed in Table B1:
- https://hadoop.apache.org/docs/r3.0.0/hadoop-mapreduce-client/hadoop-mapreduce-client-core/mapred-default.xml
- https://cwiki.apache.org/confluence/display/Hive/Configuration+Properties

**Table B1.** Workload Use Cases and Optimizations

| | USE CASE | METHOD | PRIMARY DATA TYPE | TUNINGS |
|---|---|---|---|---|
| 1 | Find top 100 products that are sold together frequently in given stores. | UDF or UDTF | Structured | mapreduce.input.fileinputformat.split.maxsize = 134217728 |
| 2 | Find the top 30 products that are mostly viewed together with a given product in online store. | MapReduce | Semi-structured | hive.exec.reducers.bytes.per.reducer = 512000000 |
| 3 | For a given product, get a top-30 list sorted by number of views in descending order of the last five products that are mostly viewed before the product was purchased online. | MapReduce | Semi-structured | hive.exec.reducers.max = 1000000000<br>hive.exec.reducers.bytes.per.reducer = 512000000 |
| 4 | Shopping cart abandonment analysis: For users who added products in their shopping carts but did not check out in the online store during their session, find the average number of pages they visited during their sessions. | MapReduce | Semi-structured | mapreduce.input.fileinputformat.split.maxsize = 536870912<br>hive.exec.reducers.bytes.per.reducer = 128000000<br>hive.exec.reducers.max = 1000000000<br>hive.optimize.correlation = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |
| 5 | Build a model using logistic regression for a visitor to an online store, based on existing users' online activities (interest in items of different categories) and demographics. | Machine Learning | Semi-structured | hive.auto.convert.join.noconditionaltask.size = 5000000000<br>hive.exec.reducers.bytes.per.reducer = 2048000000<br>hive.exec.reducers.max = 1000000000 |
| 6 | Identify customers shifting their purchase habit from physical store to Web sales. | Pure Query Language only | Structured | mapreduce.input.fileinputformat.split.maxsize = 134217728<br>hive.exec.parallel = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000<br>hive.optimize.index.filter = TRUE |
| 7 | List top-10 states in descending order with at least 10 customers, who during a given month bought products with the price at least 20 percent higher than the average price of products in the same category. | Pure Query Language only | Structured | mapreduce.input.fileinputformat.split.maxsize = 536870912<br>hive.exec.reducers.bytes.per.reducer = 536870912<br>hive.auto.convert.join.noconditionaltask.size = 5000000000 |
| 8 | For online sales, compare the total sales amount for which customers checked online reviews before making the purchase and that of sales for which customers did not read reviews. Consider only online sales for a specific category in a given year. | MapReduce | Semi-structured | mapreduce.input.fileinputformat.split.maxsize = 134217728<br>hive.exec.reducers.bytes.per.reducer = 256000000 |
| 9 | Aggregate total amount of sold items over different combinations of customers based on selected groups of marital status, education status, sales price, and different combinations of state and sales profit. | Pure Query Language only | Structured | mapreduce.input.fileinputformat.split.maxsize = 134217728<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |
| 10 | For all products, extract sentences from reviews that have positive or negative sentiment. | UDF, UDTF, NLP | Unstructured | hive.auto.convert.join.noconditionaltask.size = 10000000000 |
| 11 | For a given product, measure the correlation of sentiments, including the number of reviews and average review ratings, on product monthly revenues within a given time frame. | Pure Query Language only | Semi-structured | hive.exec.parallel = TRUE |
| 12 | Find all customers who viewed items of a given category on the Web in a given month and year that was followed by an in-store purchase of an item from the same category in the three consecutive months. | Pure Query Language only | Semi-structured | mapreduce.input.fileinputformat.split.maxsize = 536870912<br>hive.exec.parallel = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |
| 13 | Display customers with both store and Web sales in consecutive years for which the increase in Web sales exceeds the increase in-store sales for a specified year. | Pure Query Language only | Structured | mapreduce.input.fileinputformat.split.maxsize = 134217728<br>hive.exec.reducers.bytes.per.reducer = 128000000<br>hive.auto.convert.sortmerge.join = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |
| 14 | Calculate the ratio between the number of items sold over the internet in the morning (7 to 8 AM) to the number of items sold in the evening (7 to 8 PM) for customers with a specified number of dependents. | Pure Query Language only | Structured | No optimizations |
| 15 | Find the categories with flat or declining sales for in-store purchases during a given year for a given store. | Pure Query Language only | Structured | mapreduce.input.fileinputformat.split.maxsize = 536870912<br>hive.exec.reducers.bytes.per.reducer = 128000000<br>hive.auto.convert.sortmerge.join = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 5000000000 |
| 16 | Compute the impact of an item price change on store sales by computing the total sales for items in a 30-day period before and after the price change. | Pure Query Language only | Structured | hive.auto.convert.join.noconditionaltask.size = 5000000000 |
| 17 | Find the ratio of items sold with and without promotions in a given month and year. Only items in certain categories sold to customers living in a specific time zone are considered. | Pure Query Language only | Structured | mapreduce.input.fileinputformat.split.maxsize = 134217728<br>hive.exec.reducers.bytes.per.reducer = 128000000<br>hive.auto.convert.sortmerge.join = TRUE<br>hive.exec.parallel = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |

| | USE CASE | METHOD | PRIMARY DATA TYPE | TUNINGS |
|---|---|---|---|---|
| 18 | Identify the stores with flat or declining sales in three consecutive months, and check if there are any negative reviews regarding these stores available online. Analyze the online reviews for these items to determine if there are any major negative reviews. | UDF, UDTF, NLP | Unstructured | mapreduce.input.fileinputformat.split.maxsize = 33554432<br>hive.exec.reducers.bytes.per.reducer = 512000000<br>hive.auto.convert.sortmerge.join = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 5000000000 |
| 19 | Retrieve the items with the highest number of returns, where the number of returns was approximately equivalent across all store and web channels. | UDF, UDTF, NLP | Unstructured | mapreduce.input.fileinputformat.split.maxsize = 33554432<br>hive.exec.reducers.bytes.per.reducer = 16777216<br>hive.exec.parallel = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |
| 20 | Perform customer segmentation for return analysis. Segment customers according to varying criteria such as return frequency, returns to order ratio, and return amount ratio. | Machine Learning | Structured | mapreduce.input.fileinputformat.split.maxsize = 134217728<br>mapreduce.task.io.sort.factor = 100<br>mapreduce.task.io.sort.mb = 512<br>mapreduce.map.sort.spill.percent = 0.99 |
| 21 | Get all items that were sold in stores in a given month and year, and which were returned in the next six months and repurchased by the returning customer afterwards through the Web sales channel in the following three years. | Pure Query Language only | Structured | hive.auto.convert.join.noconditionaltask.size = 5000000000 |
| 22 | Compute the percentage change in inventory between the 30-day period before the price change and the 30-day period after the change. | Pure Query Language only | Structured | mapreduce.input.fileinputformat.split.maxsize = 33554432<br>hive.exec.reducers.bytes.per.reducer = 33554432<br>hive.exec.parallel = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |
| 23 | Calculate the coefficient of variation and mean of every item and warehouse of the given and the consecutive months. | Pure Query Language only | Structured | mapreduce.input.fileinputformat.split.maxsize = 33554432<br>hive.exec.reducers.bytes.per.reducer = 5368709120<br>hive.exec.reducers.max = 1000000000<br>hive.exec.parallel = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |
| 24 | For a given product, measure the effect of competitors' prices on the product's in-store and online sales. | Pure Query Language only | Structured | hive.exec.parallel = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 1000000000 |
| 25 | Perform customer segmentation analysis. Segment customers according to key shopping dimensions such as how recent the last visit was, frequency of visits, and monetary amount. | Machine Learning | Structured | hive.exec.reducers.bytes.per.reducer = 512000000<br>hive.auto.convert.join.noconditionaltask.size = 1000000000<br>mapreduce.input.fileinputformat.split.maxsize = 33554432<br>hive.exec.reducers.bytes.per.reducer = 32000000<br>hive.exec.mode.local.auto = TRUE<br>hive.exec.mode.local.auto.inputbytes.max = 1500000000<br>mapreduce.job.ubertask.enable = TRUE |
| 26 | Cluster customers into book buddies or groups based on their in-store book purchasing histories. | Machine Learning | Structured | mapreduce.input.fileinputformat.split.maxsize = 134217728<br>hive.exec.reducers.bytes.per.reducer = 256000000<br>hive.auto.convert.join.noconditionaltask.size = 10000000000 |
| 27 | Extract competitor product names and model names (if any) from online product reviews for a given product. | UDF, UDTF, NLP | Unstructured | mapreduce.input.fileinputformat.split.maxsize = 33554432<br>hive.exec.reducers.bytes.per.reducer = 32000000<br>hive.exec.mode.local.auto = TRUE<br>hive.exec.mode.local.auto.inputbytes.max = 1500000000<br>mapreduce.job.ubertask.enable = TRUE |
| 28 | Build text classifier for online review sentiment classification (Positive, Negative, Neutral). | Machine Learning | Unstructured | bigbench.hive.optimize.sampling.orderby = FALSE<br>mapreduce.input.fileinputformat.split.maxsize = 16777216<br>hive.vectorized.execution.enabled = FALSE |
| 29 | Perform category affinity analysis for products purchased together online. | UDF or UDTF | Structured | hive.exec.reducers.bytes.per.reducer = 256000000<br>hive.exec.parallel = TRUE<br>hive.auto.convert.join.noconditionaltask.size = 10000000000 |
| 30 | Perform category affinity analysis for products viewed together online. | UDF, UDTF, MapReduce | Semi-structured | mapreduce.input.fileinputformat.split.maxsize = 134217728<br>hive.exec.parallel = TRUE<br>hive.exec.reducers.bytes.per.reducer = 256000000<br>hive.exec.reducers.max = 1000000000<br>hive.auto.convert.join.noconditionaltask.size = 10000000000 |

**NLP** – natural language processing; **UDF** – user-defined function; **UDTF** – user-defined table function

## Appendix C: BigBench Tuning Parameters

Cloudera Enterprise includes many default parameter settings. Table C1 shows only the tuning parameters we changed from the default setting to achieve the best performance from the reference architecture cluster. Table C2 provides the recommended values for workload-specific parameters when running the use cases listed in Appendix B.

**Table C1.** BigBench Tuning Parameters

| PROPERTY | VALUE |
| --- | --- |
| **YARN** | |
| yarn.scheduler.maximum-allocation-mb | 360 GB |
| yarn.scheduler.maximum-allocation-vcores | 80 |
| yarn.nodemanager.resource.memory-mb | 360 GB |
| yarn.nodemanager.resource.cpu-vcores | 80 |
| Client Java Heap Size in Bytes | 4 GB |
| **Hadoop Distributed File System (HDFS)** | |
| dfs.namenode.handler.count | 80 |
| dfs.namenode.service.handler.count | 80 |
| dfs.datanode.handler.count | 30 |
| dfs.permissions | FALSE |
| dfs.datanode.socket.write.timeout | 630000 |
| dfs.socket.timeout | 630000 |
| dfs.datanode.max.transfer.threads | 65536 |
| **MapReduce** | |
| mapreduce.output.fileoutputformat.compress | true |
| mapreduce.output.fileoutputformat.compress.codec | org.apache.hadoop.io.compress.SnappyCodec |
| mapreduce.map.output.compress.codec | org.apache.hadoop.io.compress.SnappyCodec |
| mapreduce.map.memory.mb | 4.5 GB |
| mapreduce.reduce.memory.mb | 4.5 GB |
| mapreduce.map.java.opts.max.heap | 4 GB |
| mapreduce.reduce.java.opts.max.heap | 4 GB |
| mapreduce.reduce.shuffle.parallelcopies | 80 |
| mapreduce.job.reduce.slowstart.completedmaps | 0.9 |
| **Hive and Spark** | |
| hive.execution.engine | spark |
| spark.master | yarn-client |
| hive.spark.job.monitor.timeout | 7200 seconds |
| spark.network.timeout | 9000 seconds |
| spark.eventLog.enabled | TRUE |
| spark.eventLog.dir | hdfs://localhost:8020/user/spark/applicationHistory |
| spark.serializer | org.apache.spark.serializer.KryoSerializer |
| hive.merge.sparkfiles | FALSE |
| spark.kryo.referenceTracking | FALSE |
| spark.io.compression.codec | lzf |
| spark.storage.memoryFraction | 0.01 |
| spark.executor.extraJavaOptions | -XX:+UseParallelOldGC -XX:ParallelGCThreads=5 -XX:NewRatio=1 -XX:SurvivorRatio=1 -XX:+UseCompressedOops |
| spark.yarn.maxAppAttempts | 1 |
| spark.driver.memory | 20 GB |
| spark.executor.cores | 5 |
| spark.executor.memory | 18475 MB |
| spark.yarn.executor.memoryOverhead | 4096 MB |

1 Source: gartner.com/en/newsroom/press-releases/2018-02-05-gartner-survey-shows-organizations-are-slow-to-advance-in-data-and-analytics

2 Testing by Intel. No product or component can be absolutely secure.

**Baseline Configuration:** Testing by Intel as of June 6, 2019. Cluster: 1 NameNode and 4 data nodes. 2x Intel® Xeon® processor E5-2680 v3 @ 2.5 GHz (12 cores/24 threads, 30 MB SmartCache); 24 x 16 GB DDR4 @ 2133 MHz (operating @ 1600 MHz); total memory = 384 GB; boot drive = 1x Intel® SSD DC S3500 120 GB (3 nodes) and 1x Intel® SSD DC S3500 240 GB (1 node); worker node storage = 8 x 3 TB 7200 RPM HDD (total storage = 24 TB); management node storage = 2 x 3 TB 7200 RPM HDD; microcode = 0x43; Intel® Hyper-threading Technology ON; Intel® Turbo Boost Technology ON; NIC = Intel® Ethernet Adapter X540-AT2 (10 GbE); BIOS = SE5C610.86B.01.01.0028.121720182203; OS = Red Hat Enterprise Linux (RHEL) 7.6; kernel = 3.10.0-957.12.2.el7.x86_64; CDH v5.16.1; Java version jdk1.7.0_67-cloudera; workload = BigBench (based on TPCx-BBv1.3.1) with scale factor of 3 TB; Hadoop v2.6.0-cdh5.16.1; Hive v1.1.0-cdh5.16.1; Spark v1.6.0-cdh5.16.1.
Results: 2 parallel streams: 375 BigBench BBQpm ; 4 parallel streams: 395 BigBench BBQpm

**Device-under-Test Configuration:** Testing by Intel as of September 22, 2019. Cluster: 3 management nodes and 4 worker nodes. 2x Intel® Xeon® Gold 6248 processor @ 2.5 GHz (20 cores/40 threads, 25 MB SmartCache); 12 x 32 GB DDR4 @ 2666 MHz; total memory = 384 GB; boot drive = 960 GB SSD 2.5in SATA 3.0 6Gb/s Intel Youngsville SE SSDSC2KB019T701 DC S4500 Series; worker node storage = 8 x 4 TB 7200 RPM HDD plus 1x 1.6 TB Intel® SSD DC P4610 for YARN tmp files/Spark shuffle (total storage = 32 TB); microcode = 0x5000021; Intel® Hyper-threading Technology ON; Intel® Turbo Boost Technology ON; NIC = 1x Intel® Ethernet Adapter X722 (10 GbE); BIOS = SE5C620.86B.02.01.0008.031920191559; OS = RHEL 7.6; kernel = 3.10.0-957.27.2.el7.x86_64; CDH v6.2.0; Java version jdk1.8.0_181-cloudera; workload = BigBench (based on TPCx-BBv1.3.1) with scale factor of 3 TB; Hadoop v3.0.0-cdh6.2.0; Hive v2.1.1-cdh6.2.0; Spark v2.4.0-cdh6.2.0.
Results: 2 parallel streams: 618 BigBench BBQpm; 4 parallel streams: 675 BigBench BBQpm (higher is better)

3 Testing by Intel. No product or component can be absolutely secure.

**Baseline Configuration (all-DRAM):** Testing by Intel as of January 3, 2020. 10 worker nodes, 2x 2nd Generation Intel® Xeon® Gold 6248 processor (20 cores); Intel® Hyper-threading Technology ON; Intel® Turbo Boost Technology ON; total memory 384 GB (12 x 32 GB @ 2666 MHz); BIOS = SE5C620.86B.02.01.0008.031920191559 (ucode:0x500002c); OS = CentOS 7.6, 5.3.11-1.el7.elrepo.x86_64
Results: 91 standard queries: 38.716; 9 I/O-intensive queries: 22.838

**DUT Configuration (with Intel® Optane™ persistent memory):** Testing by Intel as of January 3, 2020. 4 worker nodes and 10 worker nodes, 2x 2nd Generation Intel® Xeon® Gold 6248 processor (20 cores), Intel® Hyper-threading Technology ON; Intel® Turbo Boost Technology ON; DRAM = 384 GB (12 x 32 GB @ 2666 MHz); Intel Optane persistent memory 1.5 TB (12 x 128 GB @ 2666 MHz); BIOS = SE5C620.8 6B.02.01.0008.031920191559 (ucode:0x500002c), OS = CentOS 7.6, 5.3.11-1.el7.elrepo.x86_64
Results: 4-node, 91 standard queries: 4-node: 51.460; 10-node: 18.701. 4-node, 9 I/O-intensive queries: 4-node: 12.226; 10-node: 8.107

**Software Configuration:** OpenJDK 1.8.0_222; Apache Hadoop 2.7.5; Apache Spark 2.3.2; Optimized Analytics Package for Spark Platform (OAP) Commit a4e64ca. **Workload:** 91 queries from a decision-support benchmark (derived from TPC-DS) and 9 I/O-intensive queries; data scale = 10 TB; metric = GEOMEAN. **Spark Configuration (2 executors per node):** executor number = 20; executor core = 40 vCores; executor memory = 150 GB. **Performance/$** figures in Table 3 were derived using an internal approved tool and were approved by PDT and performance forums.

4 While this reference architecture is created for Cloudera Enterprise and CDH 6.2, we anticipate that the solution architecture recommendations will be applicable to and compatible with future Cloudera product releases.

5 See endnote 2.

6 See endnote 3.

7 Same configurations as listed in endnote 2, except the number of worker nodes for the DUT configuration increased from 4 to 10.
Results: 2 parallel streams @ 3 TB, 618 BigBench BBQpm; @ 10 TB, 1401 BigBench BBQpm. 4 parallel streams @ 3 TB, 675 BigBench BBQpm; @ 10 TB, 1481 BigBench BBQpm

8 See endnote 2.

9 See endnote 3.

10 See endnote 2.

11 See endnote 3.

12 For more information about Cloudera Enterprise, visit docs.cloudera.com/documentation/enterprise/5-10-x/PDF/cloudera-introduction.pdf.

13 See endnote 2.

14 See endnote 7.

15 See endnote 3.

**cloudera** | (intel®)