

Unlock the Power of Private Cloud Big Data Analytics

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

Is this solution for you? Do you...

- ✓ Need more performance from your Cloudera distribution?
- ✓ 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?
- ✓ Need to increase data center efficiency and flexibility?

Learn about the [Business Story](#) →

Learn more about the [Reference Architecture](#) →

Intel Authors

Esther Baldwin

Cloud & Enterprise Solutions Lead

VijayKumar Bandari

Cloud Solutions Architect Sales

Christopher Bodine

Cloud Solution Engineer

Gregory Fields

Cloud Solution Engineer

Jude Lee

Cloudera Alliance Manager

Amandeep Raina

Cloud Solution Engineer

Priyanka Sebastian

Cloud Solution Engineer

Merritte Stidston

Cloud Solutions Architect Sales

Iyer Venkatasan

Enterprise Solutions Marketing

Cloudera Author

Ali Bajwa

Partner Engineering Director

Mo Amao

Partner Solution Engineer

Business Story

Unlock the Power of Private Cloud Big Data Analytics

Solution Benefits

- Open-source platforms provide flexibility and interoperability with existing tools.
- Disaggregation of compute and storage leads to more flexibility and efficiency.
- A comprehensive set of management tools simplifies cluster configuration and scaling.
- Upgrading from a 1st Generation Intel® Xeon® Scalable processor to Cloudera Data Platform (CDP) Private Cloud Base running on 3rd Generation Intel Xeon Scalable processors provides up to a 2x throughput improvement.¹

Executive Summary

It is well understood that enterprises can't extract business value from all the large volumes of data they generate. The difficulty lies in integrating isolated silos of data throughout the enterprise and managing that data efficiently. One common hurdle is competing business units that control access to the data but refuse to share it. Also, internal policies that were intended to protect the company and its customers can hamper data analysis and create a significant burden on extracting business value and improving decision making. The velocity of the data flood adds additional stress on business units, IT and executive management, who must work through the complexities created in a data-driven world. Every enterprise has data management problems. The difficulty is identifying and then correcting them promptly. Simply having a better product or service isn't sufficient in today's world. Getting your product or service to the right customer, whether new or existing, is paramount. Spotting trends before your biggest rival could mean the difference between success and failure.

If these challenges sound familiar, you're not alone. According to Gartner, 91 percent of organizations struggle to reach data maturity.² But don't despair—a collaboration between Intel and Cloudera has created a big data analytics platform specifically designed for large-scale on-premises workloads.

Cloudera Data Platform (CDP) Private Cloud powers on-premises, data-driven decision making by easily, quickly and safely connecting and securing the business's entire data lifecycle. This big data analytics platform helps business leaders modernize their data center by streamlining data management and workload orchestration.

The separation of compute and storage leads to improved flexibility and efficiency. Enterprises can use CDP Private Cloud to migrate to a container-based environment and take advantage of the agility and scalability of containers.

Time to business value is the main issue once data volume, velocity and access issues are resolved. Solving a problem one minute too late is lost business value. Said another way, "time is money." With the Intel® architecture underlying CDP Private Cloud, we provide the power that big data analytics demand. Intel and Cloudera collaborated to improve compute performance, storage efficiency, artificial intelligence (AI) acceleration and more. The result is a private cloud data platform built to meet current big data analytics needs that can scale to meet your business needs today and into the future. Our tests show that a modern version of CDP Private Cloud running on the latest Intel hardware can improve the data analysis performance of several aspects of the CDP Private Cloud system. For example, **running a data warehouse workload on 3rd Generation Intel® Xeon® Scalable processors can improve throughput by up to 2x.**³

This document provides a business-level overview of CDP Private Cloud, describes a reference solution for deployment and highlights the platform's performance and scalability. And if you are already using a legacy distribution of a Cloudera product, this document also describes best practices for migrating to the latest distribution.

Business Challenge: Extracting Business Insights from a Jumble of Data

Data is everywhere in your enterprise. It is constantly being generated by machines, customers and applications. It piles up in various data warehouses. The significant problem is that hidden in that data are insights about your business—information about network security, customer preferences, supply chain dynamics and more. When data exists in silos and data analytics runs in those same silos, it is nearly impossible to unlock the business value within this volume of data. And as data continues to grow and big data analytics workloads increase accordingly, your infrastructure must be able to scale appropriately. However, storage requirements often outpace compute needs, so the ability to scale these resources independently is crucial to data center efficiency.

CDP Private Cloud is an open-source, scalable data platform optimized to run on high-performance Intel hardware. It also supports the disaggregation of compute and storage. Using CDP Private Cloud, you can quickly and cost-efficiently extract value from your data without leaving the data center (see the [Solution Value](#) section).

Use Cases Abound for a Private Cloud Data Analytics Platform

Many industries, from manufacturing to healthcare to retail, and from transportation to hospitality to life sciences, are under pressure to turn their data into business value. CDP Private Cloud 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 (Cloudera supports cloud deployments), keeping data on-premises in a private cloud sometimes makes more sense. For example, 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 illustrate how CDP Private Cloud can help solve your business's data processing issues. These use cases demonstrate the flexibility of the CDP Private Cloud 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. Intel manufacturing uses large datasets containing billions of data points per day per factory. For high-volume manufacturing companies like Intel, advanced analytics using the latest Intel® processors and CDP Private Cloud can help reduce data processing time and can lead to improved yield, fast time to market, accelerated insights, increased productivity, excursion control and more.
- **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. IQVIA, a global provider of advanced analytics, technology solutions and contract research services to the life sciences industry, has used CDP Private Cloud to accelerate query responses from days to seconds.⁴
- **Sales and marketing:** Use data to optimize marketing campaigns and create advanced recommendation engines.
- **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 a “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 networking data into a data lake and run analytics to detect vulnerabilities and threats.

Solution Value: High Performance, Ultimate Flexibility and Excellent Scalability

Enterprises need answers—and they need them now! At its core, CDP Private Cloud is a next-generation, cloud-native architecture for on-premises deployments that drives real-time and batch analytics processing. It speeds time to value by separating compute and storage, integrating a suite of security and governance tools, and packaging all pieces together with a simpler and intuitive management console.

“ Boost query latency and throughput by up to 2x by running Cloudera Private Cloud Base on 3rd Gen Intel® Xeon® Scalable processors.”¹

CDP Private Cloud delivers a suite of analytic engines ranging from stream and batch data processing to data warehousing, operational database and machine learning (see Figure 1). Enterprises can use CDP Private Cloud to gain end-to-end big data capabilities—ingest, store, process and analyze for insights. The CDP Private Cloud shared data experience (SDX) applies consistent security and governance, enabling users to share and discover data for use across many demanding big data analytics workloads.

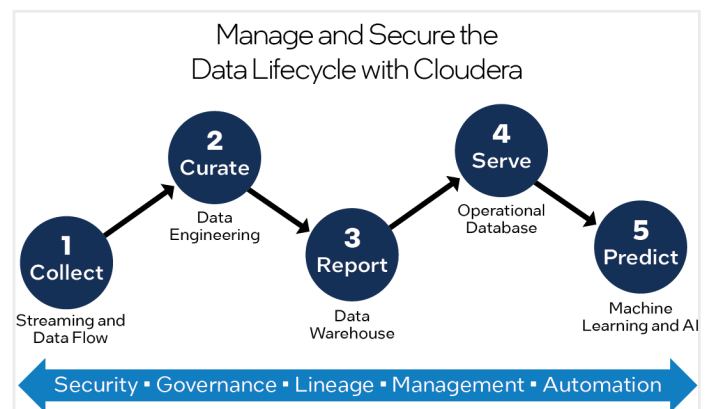


Figure 1. CDP Private Cloud simplifies data management, offers a wide variety of analytic engines and provides a shared data experience (SDX).

Some specific business benefits from CDP Private Cloud include:

- **An end-to-end data analytics processing platform**, in contrast to stand-alone solutions that require stitching together for the data flow.
- **Ultimate agility and flexibility:** It’s easy to integrate CDP Private Cloud with existing infrastructure and tools, and the platform offers robust security, governance, data protection and data management features. CDP Private Cloud is built using open-source technologies, 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 CDP Private Cloud reference solution 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. Also, CDP Private Cloud includes software that simplifies scaling by automating node configuration.
- **High efficiency:** Support for independent scaling of compute and storage resources lets data centers invest in precisely the resources they need for their specific workloads. Avoiding overprovisioning of either compute or storage can lead to cost savings, a simpler infrastructure and lower maintenance efforts.
- **Enhanced security and governance:** CDP Private Cloud has a complete suite of security and governance capabilities. These services regulate what end users can do through the analytic experiences, but operate independently of these experiences. This means that the security and governance tools can be independently configured, managed and upgraded, and these changes will automatically reflect in the analytic experiences.

Intel and Cloudera’s participation in joint engineering results in optimizations of CDP Private Cloud pertaining to faster compute performance, increased storage efficiency, enhanced security, excellent AI support and great query performance. These optimizations let CDP Private Cloud take advantage of Intel® Optane™ persistent memory (PMem), Intel Xeon Scalable processors, Intel® FPGAs, Intel® QuickAssist Technology, Intel® Advanced Vector Extensions 512, Intel® Math Kernel Library, Intel® AES-NI and Intel® Intelligent Storage Acceleration Library. **Intel testing shows that running CDP Private Cloud Base on 3rd Gen Intel Xeon Scalable processors can boost throughput by up to 2x.**⁵ Additional benefits of CDP Private Cloud include manageability features, including native high-availability; fault-tolerance and self-healing storage; automated backup and disaster recovery; advanced system and data management; and a single pane of glass for cluster administration, automation, management and security.

Solution Architecture: Scalable and Agile Big Data Analytics Platform

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

- **3rd Gen 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® 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 PMem and Intel® Optane™ SSDs are unique innovations that bridge critical gaps in the storage and memory hierarchy, delivering persistent memory, large memory pools, fast caching and fast storage.
- **Solidigm SSDs:** Solidigm storage solutions are optimized to unlock data’s virtually unlimited potential. Solidigm draws on decades of technical innovation to offer a broad portfolio of SSDs. Choose from a wide selection of performance-optimized and value-optimized NAND SSDs that offer high density, performance and reliability. The testing documented in this reference architecture was conducted using the Solidigm D7-P5600 SSD as the fast storage drive for the CDP Private Cloud Base general-purpose worker nodes. An alternative would be to use the Solidigm D7-P5620 SSD—both drives are currently available.

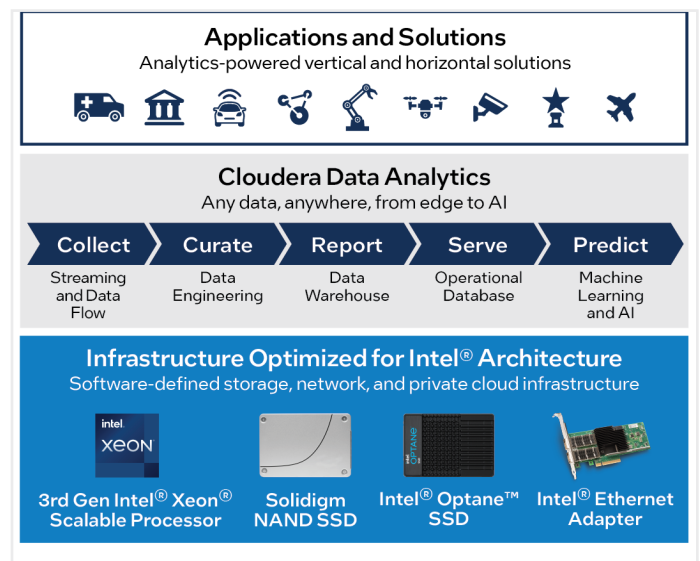


Figure 2. High-level diagram of the basic CDP Private Cloud solution architecture.

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

Reference Architecture

Unlock the Power of Private Cloud Big Data Analytics

Get better performance by upgrading to a newer version of Cloudera Data Platform with 3rd Gen Intel® Xeon® Scalable processors.⁶

Big Data Analytics Process

1. 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, weblogs, click streams and other external sources.

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

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

4. Act. Finally, results in the form of compressed datasets are made available for business 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 predefined use cases.

Table of Contents

- Overview of CDP Private Cloud 5
- Reference Architecture Considerations..... 6
- BigBench Results 9
- Migration Path from Legacy Cloudera Distributions to CDP Private Cloud... 10
- Additional Infrastructure Considerations 10
- Learn More 11
- Appendices A-C 12-14

Overview of CDP Private Cloud

CDP Private Cloud provides a scalable, versatile and integrated platform that helps modernize the data center. This platform simplifies managing the growing volume and variety of data in your enterprise, unleashing the business value of that data. By disaggregating compute and storage and supporting a container-based environment, CDP Private Cloud helps enhance business agility and flexibility. 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. The platform also includes secure user access and data governance features.⁷

CDP Private Cloud consists of CDP Private Cloud Base and CDP Private Cloud Data Services (see Figure 3).

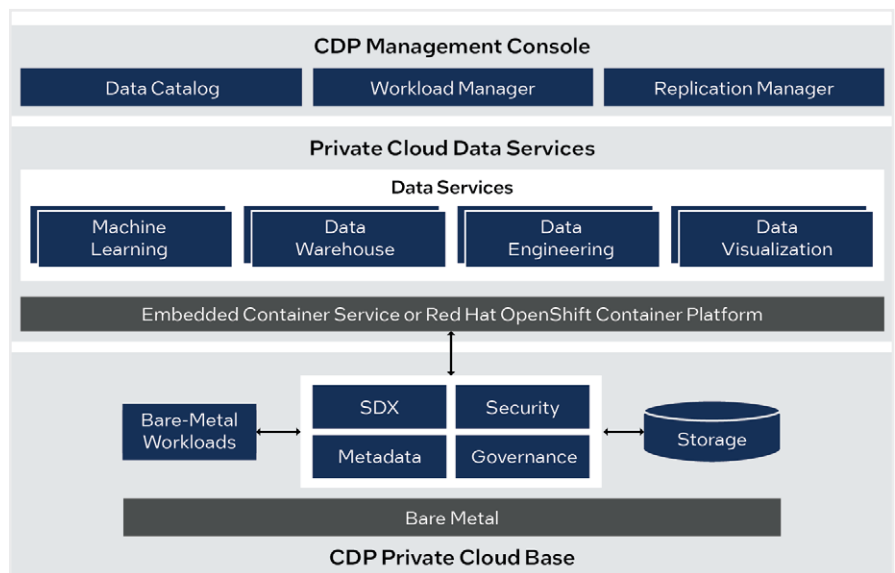


Figure 3. Overview of the Cloudera Data Platform.

CDP Private Cloud Base

CDP Private Cloud Base is the foundation of CDP Private Cloud. It supports a variety of hybrid solutions where compute tasks are separated from data storage and where data can be accessed from remote clusters, including workloads created using CDP Private Cloud Data Services. This hybrid approach provides a foundation for containerized applications by managing storage, table schema, authentication, authorization and governance.

CDP Private Cloud Base is composed of a variety of components such as Apache HDFS, Apache Hive 3, Apache Spark, Apache HBase and Apache Impala, along with many other components for specialized workloads. You can select any combination of these services to create clusters that address your business requirements and workloads. Several preconfigured service packages are also available for common workloads. It also includes SDX, storage management, and metadata and governance features. It replaces—and is equivalent to—CDP Data Center.

CDP Private Cloud Data Services

CDP Private Cloud Data Services has access to all the Base features (including SDX) and brings many of the benefits of the public cloud to the data center. Through the use of containers deployed on Kubernetes, CDP Private Cloud Data Services brings both agility and predictable performance to analytic applications. You can either use a dedicated Red Hat OpenShift cluster or deploy an Embedded Container Service (ECS) for the containers. You can use CDP Private Cloud Data Services to rapidly provision and deploy services such as Cloudera Data Warehousing, Cloudera Machine Learning and Cloudera Data Engineering, and easily scale them up or down as required.

Reference Architecture Considerations

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

- Scalability, agility and flexibility
- Excellent performance and lower total cost of ownership
- 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.

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

- **Edge nodes:** One edge node is for Cloudera Manager; the other is for Hadoop Clients. Edge nodes support services that are needed for cluster operation. If one edge node fails, the other one can pick up the extra work.
- **Master nodes:** Up to four master nodes are available to manage important services in the Hadoop cluster: NameNode, Resource Manager, Secondary NameNode and the HBase Master.

- **Worker nodes:** Worker nodes handle the bulk of the Hadoop processing. The number of worker nodes necessary depends on dataset size. Depending on your scalability needs, you can replicate this configuration to 20 or 30 worker nodes or decrease it to four or five worker nodes. Some testing was performed with four worker nodes, but 10 worker nodes represent a nominal workload configuration.

The following sections provide information on recommended configurations for the various types of nodes associated with CDP Private Cloud. Considerations for networking are also included. Table 1 summarizes the software services that each node provides.

Table 1. CDP Private Cloud Base and CDP Private Cloud Data Services Nodes and Roles

Physical Node	Software Function
Private Cloud Base	
Edge Nodes	<ul style="list-style-type: none"> ▪ Hadoop Clients ▪ Cloudera Manager
Master Nodes 1–4	<ul style="list-style-type: none"> ▪ NameNode ▪ Resource Manager ▪ ZooKeeper
Worker Nodes	<ul style="list-style-type: none"> ▪ DataNode ▪ NodeManager ▪ YARN workloads
Private Cloud Data Services	
OpenShift Controller Nodes (minimum of 3)	<ul style="list-style-type: none"> ▪ OpenShift services (More details on OpenShift)
OpenShift Cluster Admin Node	
Bootstrap Node	
Worker Nodes (minimum of 10)	<ul style="list-style-type: none"> ▪ Kubernetes operators ▪ Workload pods ▪ Kubernetes services (More details on Kubernetes)

CDP Private Cloud Base Master Nodes

At least three master nodes are required. More nodes with the same configuration can be used for edge nodes. This configuration is sized for approximately 1 PB of cluster storage or 250 worker nodes. Table 2 lists the recommended hardware for the CDP Private Cloud Base master nodes.

Table 2. CDP Private Cloud Base Master Node Configuration

Recommended Component	
Processor	2x Intel® Xeon® Gold 6326 processor (16 cores, 2.9 GHz, 24 MB cache)
Memory	256 GB (16x 16 GB 3200 MT/s)
Network	Intel® Ethernet Network Adapter E810-XXVDA2 for OCP 3.0 (dual-port 25 GbE)
Storage (Data)	2x 3.84 TB Solidigm D7-P5520 SSD NVMe
Storage (OS)	2x 480 GB Solidigm D3-S4520 SSD

CDP Private Cloud Base Worker Nodes

A minimum of three worker nodes is required, but we recommend at least five to support high availability and downtime for maintenance. Table 3 lists the recommended hardware for general-purpose CDP Private Cloud Base worker nodes. This processor and memory configuration is ideal for nodes that are also running other services, such as Cloudera Data Warehouse, Cloudera Machine Learning and Cloudera Data Hub.

Table 3. CDP Private Cloud Base General-Purpose Worker Node Configuration

Recommended Component	
Processor	2x Intel® Xeon® Gold 6348 processor (28 cores, 2.6GHz, 40 MB cache)
Memory	512 GB (16 x 32 GB 3200 MT/s)
Network	Intel® Ethernet Network Adapter E810-XXVDA2 for OCP 3.0 (dual-port 25 GbE)
Storage (Data)	12x 4 TB 7.2K RPM SATA 6 Gbps
Fast Storage	2x 3.2 TB Solidigm D7-P5600 SSD NVMe
Storage (OS)	2x 480 GB Solidigm D3-S4520 SSD

In contrast, if most workloads will be running on the CDP Private Cloud Data Services cluster and the CDP Private Cloud Base cluster will be providing only HDFS storage, then the storage-only worker node configuration shown in Table 4 may be more appropriate.

Table 4. CDP Private Cloud Base Data or Storage-Only Worker Node Configuration

Recommended Component	
Processor	2x Intel® Xeon® Gold 6326 processor (16 cores, 2.9 GHz, 24 MB cache)
Memory	256 GB (16 x 16 GB, 3200 MT/s)
Network	Intel® Ethernet Network Adapter E810-XXVDA2 for OCP 3.0 (dual-port 25 GbE)
Storage (Data)	16x 4 TB 7.2K RPM SATA 6 Gbps
Storage (OS)	2x 480 GB Solidigm D3-S4520 SSD

Note: We recommend Intel® Xeon® Gold processors for best performance when using erasure coding with HDFS. The HDFS erasure coding feature uses the Intel® Storage Acceleration Library, which uses the Intel® AES-NI, SSE, Intel® AVX, Intel® AVX2, and Intel® AVX-512 instruction sets that Intel Xeon Gold processors support.

A slightly different configuration is recommended for CDP Private Cloud Base worker nodes that are running memory-intensive workloads (see Table 5). This configuration may be suitable for workloads that benefit from keeping large datasets in memory. The Intel Optane PMem that is used in the large-memory configuration enables large memory capacity.

Table 5. CDP Private Cloud Base Worker Node Configuration for Memory-Intensive Workloads

Recommended Component	
Processor	Two-socket Intel® Xeon® Gold 6348 processor (28 cores, 2.6 GHz, 40 MB cache)
Memory (DRAM)	512 GB (16x 32 GB, 3200 MT/s, Dual Rank)
Persistent Memory	12x 128 GB Intel® Optane™ PMem
Network	Intel® Ethernet Network Adapter E810-XXVDA2 for OCP 3.0 (dual-port 25 GbE)
Storage (Data)	1x 7.68 TB Solidigm D7-P5520 SSD (2.5in, U.2)

CDP Private Cloud Data Services Master Nodes

The CDP Private Cloud Data Services requires a dedicated OpenShift cluster or Embedded Container Service (ECS). The OpenShift or ECS cluster consists of a number of master nodes (for managing OpenShift/ECS) and a number of worker nodes (for running CDP applications). Table 6 specifies the hardware requirements for the CDP Private Cloud Data Services master nodes.

Table 6. CDP Private Cloud Data Services Master Node Configuration

Recommended Component	
Processor	2x Intel® Xeon® Gold 6326 processor (16 cores, 2.9 GHz, 24 MB cache)
Memory	256 GB (16x 16 GB 2933 MT/s)
Network	Intel® Ethernet Network Adapter E810-XXVDA2 for OCP 3.0 (dual-port 25 GbE)
Storage (OS and Data)	1x 480 GB Solidigm D3-S4520 SSD

The following master nodes are required: three OpenShift master nodes; one bootstrap node (it can be converted into an OpenShift worker node after initial deployment); and one Cluster System Admin Host node. This master node configuration is adequate for OpenShift container clusters up to 250 nodes and rarely needs to be customized.

CDP Private Cloud Data Services Worker Nodes

A minimum of four worker nodes are required. For a new deployment, we recommend 10–20 worker nodes. OpenShift supports heterogeneous node configurations, and it is possible to create specialized node configurations. We provide two recommended configurations for CDP Private Data Services worker nodes: typical CDP cloud workloads and memory-intensive CDP workloads.

Typical CDP Cloud Workloads

The recommendations in Table 7 provide a good balance of compute and memory for typical CDP cloud workloads. The configuration assumes that the Cloudera Data Warehouse and Cloudera Machine Learning experiences are running in the cloud while HDFS provides the primary big data storage on the CDP Private Cloud Base cluster. Therefore, local storage is limited to just enough to support temporary files, disk caches and storage for other applications. Contact your Intel account representative for assistance in sizing and customizing any of these configurations.

Table 7. CDP Private Cloud Data Services Worker Node Configuration for Typical Cloud Workloads

Recommended Component	
Processor	2x Intel® Xeon® Gold 6348 processor (24 cores, 3.0 GHz, 40 MB cache)
Memory	512 GB (16x 32 GB 3200 MT/s, Dual Rank)
Network	Intel® Ethernet Network Adapter E810 CQDA2
Storage (OS and Data)	1x 7.68 TB Solidigm D7-P5520 SSD

Network Considerations

At a high level, the network requirements for CDP Private Cloud are similar to those of a Cloudera Manager Virtual Private Cloud (CM VPC) deployment. The model of deployment in both cases is similar—a base cluster that houses HDFS storage and remote compute-only clusters that can read from and write to the base HDFS cluster. However, the network bandwidth requirements for CDP Private Cloud are less stringent than those of the CM VPC because data caching technology is introduced at the compute layer, which is not available in CM VPCs.

This implies that the initial load of data from the remote storage would require significant bandwidth between the compute and storage clusters, subject to the quantity of data being ingested. However, subsequent network bandwidth requirements should be significantly lower.

As shown in Figure 4, a minimum of 10 Gbps (recommended 25 Gbps) guaranteed bandwidth is required between each OpenShift worker node and each CDP Private Cloud Base Data Node. To test the worst-case bandwidth scenario, the network should be fully stressed with all OpenShift nodes trying to read or write simultaneously from the CDP Private Cloud Base nodes.

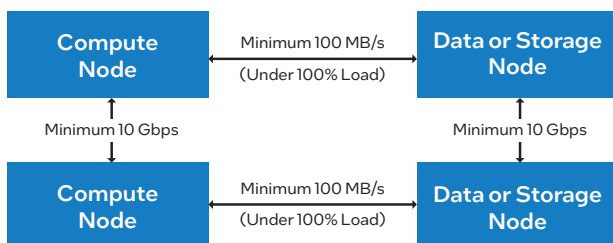


Figure 4. Network bandwidth recommendations.

The recommended network architecture is Spine-Leaf (Figure 5), with no more than a 4:1 oversubscription between the spine and leaf switches. For more information, visit [Cloudera’s networking documentation](#) for CDP Private Cloud.

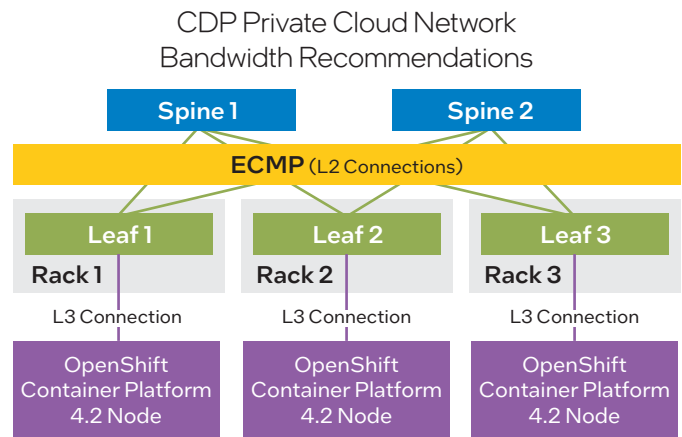


Figure 5. CDP Private Cloud network configuration recommendations.

Figure 6 illustrates the networks that are used for the CDP Private Cloud Base and CDP Private Cloud Data Services clusters, including the interconnect. Connectivity between the clusters and existing network infrastructure can be adapted to specific installations.

A common scenario is when the cluster data network is exposed to an existing network. In this scenario, the edge network is either unused or is used for application access or ingest processing.

CDP Private Cloud Cluster Networks

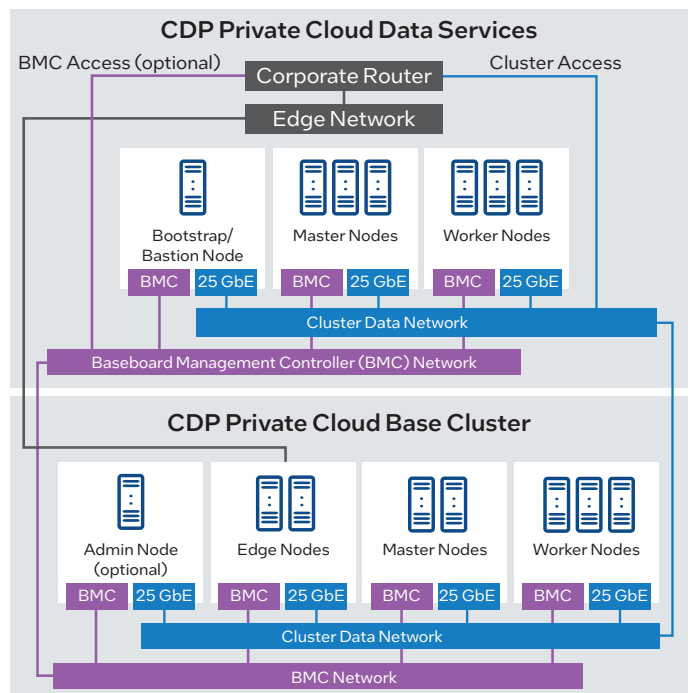


Figure 6. CDP Private Cloud cluster networks.

Network Functions

The CDP Private Cloud Cluster network functions include:

- CDP Private Cloud Base Cluster Data Network:**
 The data network carries the bulk of the traffic within the cluster. This network is aggregated within each point of delivery (POD), and pods are aggregated into the cluster switch. The CDP Private Cloud Base services are available on this network. **Note:** The CDP Private Cloud Base services do not support multihoming and are only accessible on the cluster's data network.
- CDP Private Cloud Data Services Cluster Data Network:**
 The data network carries the bulk of the traffic within the cluster. This network is aggregated within each POD, and PODs are aggregated into the cluster switch. The CDP Private Cloud Data Services are available on this network.
- Baseboard Management Controller (BMC) Network:**
 The BMC network connects the BMC ports and the out-of-band management ports of the switches. It is used for hardware provisioning and management. This network is aggregated into a management switch in each rack. This network provides access to the BMC functionality on the servers. It also provides access to the management ports of the cluster switches.
- Edge Network:** The edge network provides connectivity from one or more edge nodes to an existing on-premises network, either directly or by the POD or cluster aggregation switches. SSH access to one or more edge nodes is available on this network, and other application services may be configured and available.

Cluster Interconnect Sizing

For cluster interconnect sizing information, see [Cloudera's networking documentation](#).

BigBench Results: Newer Hardware Increases Query Performance by 2x⁸

The BigBench benchmark is derived from the TPCx-BB Express Benchmark 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.

When we compared the performance (see Figure 7) of CDP Private Cloud Base version 7.1.7 running on the Intel Xeon Gold 6348 processor to the same software running on a first-generation Intel Xeon Scalable processor, normalized big data batch analytics performance for the BigBench retail use cases improved by up to 2x for two streams.⁹ Table 8 provides the bill of materials for the tested cluster; see the appendices for additional configuration details.

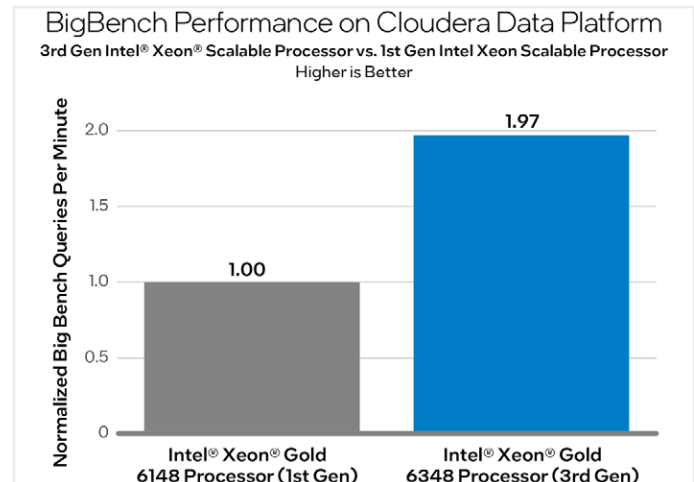


Figure 7. Double your throughput by upgrading to CDP Private Cloud Base on the latest Intel Xeon Scalable processor.

Table 8. Bill of Materials for Tested Cluster

Recommended Component	
Processor (per node)	
Management	Intel Xeon Gold 6326 processor @ 2.9/3.5 GHz (1x Master Node/Active NameNode)
Worker	Intel Xeon Gold 6348 processor @ 2.6/3.5 GHz (at least 4 worker nodes; can scale as necessary)
Memory	16x 32 GB DDR4 @ 3200 MHz (total 512 GB per node)
Network	Intel Ethernet Network Adapter E810-XXVDA2 (dual-port 25 GbE)
Storage	
Management Nodes (per node)	
— OS	1x 960 GB Solidigm D3-S4610 SSD
— HDFS	2x 1.92 TB Solidigm D7-P5500 SSD
Worker Nodes (per node)	
— OS ^a	1x 960 GB Solidigm D3-S4510 SSD
— HDFS	6x 1.92 TB Solidigm DC-P4510 SSD
— Fast Storage	1x 3.2 TB Solidigm D7-P5600 SSD NVMe

^aFor OS, the D3-P4510 SSD was used during testing. Going forward, Solidigm recommends the D3-P4520. For HDFS, the 6x 1.92 TB Solidigm DC-P4510 SSD was used during testing. Going forward, Solidigm recommends D7-P5520 SSD. All Solidigm SSDs referenced were previously known as Intel SSDs.

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 ten (10) worker nodes. You could scale out to 20 or 30 worker nodes for distributed processing, if necessary. An HDD plus one NAND SSD per worker node is expected to be more cost-effective than a dense, all-SSD-based solution.

Migration Path from Legacy Cloudera Distributions to CDP Private Cloud

Cloudera no longer supports Hortonworks Data Platform (HDP) 2.6.x and CDH 5.x. Therefore, customers using these products need to map a path for upgrading to a supported version.

Customers using Cloudera Enterprise 5.13+ and HDP 2.6.5, or Cloudera Enterprise 6 and HDP 3, should first migrate their workloads to the latest version of CDP Private Cloud Base 7.1.x (an SDX environment). Existing applications, data and hardware can be brought into the new environment; there is also an opportunity for hardware refreshes that can improve performance as noted in the [Solution Architecture](#) section. Once the migration to CDP Private Cloud Base is complete, then organizations can expand to new experiences by upgrading to CDP Private Cloud Data Services (see Figure 8). Refer to the [Cloudera Upgrade Companion](#) for more details.

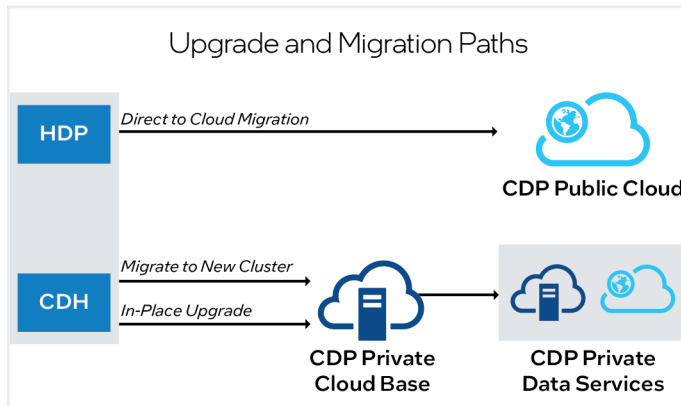


Figure 8. Hortonworks Data Platform or Cloudera Enterprise Data Hub users can follow one of the several upgrade or migration paths to CDP.

Before upgrading a cluster to CDP Private Cloud Base, the existing infrastructure should be evaluated to determine whether it meets the needs of a full CDP Private Cloud Base environment. The primary requirement for in-place upgrades is that the cluster infrastructure meets the same requirements as those specified for CDH or HDP in the [Cloudera Reference Architecture Documentation](#).

Beyond those recommendations, our general guidance is:

- After CDP Private Cloud is deployed, most of the CDP Private Cloud Base compute workloads will begin migrating to the new infrastructure. The CDP Private Cloud Base infrastructure becomes a storage-centric cluster. This configuration means that smaller memory and lighter processor configurations may be adequate for the CDP Private Cloud Base cluster.
- CDP Private Cloud relies heavily on the network infrastructure for data access. We recommend a minimum of two bonded 10 GbE connections, with 25 GbE or faster preferred. Clusters with 1 GbE network should probably not be upgraded. Although faster network cards can be installed, consider the infrastructure and labor costs.
- If the cluster is due for an infrastructure refresh, we suggest deploying a new CDP Private Cloud environment and then migrating data and workloads to the new environment.

Additional Infrastructure Considerations

This section provides information about infrastructure topics that may be useful as organizations deploy CDP Private Cloud.

Repurposing Infrastructure

Repurposing infrastructure is a common topic in upgrade conversations; this document provides some suggestions based on experience.

The recommended CDP Private Cloud Base worker nodes are storage-heavy configurations, while the CDP Private Cloud Data Services worker nodes are memory- and compute-heavy. Older CDH, HDP or CDP-DC worker nodes generally cannot be reused for cloud workloads in the CDP Private Cloud Data Services cluster. However, those nodes are potentially useful in another environment for storage-heavy applications. In some instances, it may be possible to add them to the OpenShift cluster for non-CDP Private Cloud Data Services workloads.

After a cluster refresh, customers often repurpose older CDH or HDP nodes for development, test or disaster recovery clusters. In many instances, those nodes are added to existing clusters that are not performance-critical.

Heterogeneous Nodes

The core OpenShift platform supports heterogeneous node types, including compute-, memory- and accelerator-optimized configurations. The necessary support for seamless use of heterogeneous nodes is not currently available in either OpenShift 4.2 or CDP Private Cloud Data Services. They are not recommended for initial deployments. Software support in this area is evolving rapidly, so if these configurations are of interest, contact your Intel or Cloudera representative for the latest status.

Using an Existing OpenShift Cluster

Depending on the OpenShift node configurations, this reuse may be possible. Contact your Intel or Cloudera representative for the latest status.

Scaling with Apache Ozone Object Storage

Apache Ozone is a scalable, redundant and distributed object store that is optimized for big data workloads. It addresses the scale limitation of HDFS with respect to small files and the total number of file system objects. Ozone can scale to billions of objects of varying sizes and can function effectively in containerized environments such as Kubernetes and YARN. Additionally, S3 Gateway support makes it scalable to various cloud-native architectures.

We recommend using at least one Intel® Optane™ SSD P5800X Series (800 GB, 2.5in PCIe x4) per node, and expect that adding another of these next-generation SSDs for the Ozone cache will further improve performance. Ozone is typically available as a service in the CDP Private Cloud Base distribution. (7.1+) For customers interested in Ozone, visit [Cloudera's Ozone web page](#).

Data Warehousing Using Apache Hive and Tez

Apache Hive data warehouse software enables reading, writing and managing large datasets in distributed storage. Using the Hive query language (HiveQL), which is very similar to SQL, queries are converted into a series of jobs that execute on a Hadoop cluster through Apache Tez. Users can run batch processing workloads with Hive while also analyzing the same data for interactive SQL or machine-learning workloads using tools like Apache Impala or Apache Spark—all within a single platform.

Apache Tez replaces MapReduce as the default Hive execution engine. MapReduce is no longer supported, and Tez stability is proven. With expressions of directed acyclic graphs (DAGs) and data transfer primitives, execution of Hive queries under Tez improves performance. SQL queries submitted to Hive are executed as follows:

- Hive compiles the query.
- Tez executes the query.
- YARN allocates resources for applications across the cluster and enables authorization for Hive jobs in YARN queues.
- Hive updates the data in HDFS or the Hive data warehouse, depending on the table type.
- Hive returns the query results over a JDBC connection.

Sizing CDP Private Cloud Base Clusters

Cluster nodes are broadly described as master nodes, utility nodes, gateway nodes or worker nodes.

- Master nodes run Hadoop controller processes such as the HDFS NameNode and YARN Resource Manager.
- Utility nodes run other cluster processes that are not controller processes, such as Cloudera Manager and the HMS.
- Gateway nodes are client access points for launching jobs in the cluster. The number of gateway nodes required varies depending on the type and size of the workloads.
- Worker nodes primarily run DataNodes and other distributed processes such as Impala.

Cloudera provides [guidance on sizing clusters](#), recommending role allocations for different cluster sizes ranging from three to 10 worker nodes up to 500 to 1,000 worker nodes.

Note: Cloudera recommends that you always enable high availability when using Runtime in a production environment.

Learn More

You may find the following references helpful:

Intel

- [3rd Generation Intel® Xeon® Scalable processors](#)
- [Intel® Optane™ persistent memory](#)
- [Intel® Optane™ solid-state drives](#)
- [Solidigm SSDs](#) (formerly Intel® Solid State Drives Data Center Family)
- [Intel® Ethernet Network Adapter E810-XXVDA2 for OCP 3.0 \(dual-port 25 GbE\)](#)

Cloudera

- [Cloudera Data Platform \(CDP\) Private Cloud](#)

Find the solution that is right for your organization. Visit intel.com/AI or contact an Intel representative.

Appendix A: Additional Configuration Details for BigBench Testing

Tables A1-A3 provide configuration details for the BigBench cluster testing. Table A4 shows which CDP Private Cloud roles are present on which nodes. For more information, read the [Cloudera article on this topic](#).

Table A1. Software Configuration

	Intel® Xeon® Gold 6148 Processor	Intel® Xeon® Gold 6348 Processor
OS	CentOS Linux release 7.9.2009 (Core)	CentOS Linux release 7.9.2009 (Core)
Kernel	3.10.0-1160.45.1.el7.x86_64	3.10.0-1160.45.1.el7.x86_64
Java	jdk1.7.0_67-cloudera	jdk1.8.0_232-cloudera
Workload	BigBench (based on TPCx-BB v1.6)	BigBench (based on TPCx-BB v1.6)
Cloudera Distribution	CDP Private Cloud Base 7.1.7	CDP Private Cloud Base 7.1.7

Table A2. Benchmark Configuration

Benchmark Configuration	
Total Nodes	<ul style="list-style-type: none"> 5x (1x Master node, 4 Worker nodes) for Intel® Xeon® Gold 6148 – 3 TB Scale Factor 5x (1x Master node, 4 Worker nodes) for Intel® Xeon® Gold 6348 – 3 TB Scale Factor
Workload	BigBench
Storage	HDFS, RF = 3
SQL Engine	Hive on Tez
Scale Factor	3 TB
Streams	2 parallel streams

Table A3. Spark Configuration

Spark Configuration Details	
Spark Driver memory	39 GB
Spark Executor memory	39 GB
spark.executor.memoryOverhead	10240
Spark executor cores	12
spark.dynamicAllocation.enabled	true
spark.dynamicAllocation.maxExecutors	1

Table A4. Cloudera Data Platform Private Cloud Role Distribution

Cloudera Role	Controller Node 1	Utility Node 1	Utility Node 2	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
Hive on Tez (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			
Tez Gateway	Y	Y	Y	Y
Cloudera Manager			Y	
Cloudera Manager Management Service			Y	
Apache Ranger	Y			
Apache Atlas	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	Primary Data Type
1 Find top 100 products that are sold together frequently in given stores.	Structured
2 Find the top 30 products that are mostly viewed together with a given product in the online store.	Semi-structured
3 For a given product, get a top-30 list sorted by the number of views in descending order of the last five products that were mostly viewed before the product was purchased online.	Semi-structured
4 Shopping cart abandonment analysis: For users who added products to 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.	Semi-structured
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.	Semi-structured
6 Identify customers shifting their purchase habits from a physical store to web sales.	Structured
7 List top-10 states in descending order with at least 10 customers, who during a given month, bought products with a price at least 20 percent higher than the average price of products in the same category.	Structured
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.	Semi-structured
9 The aggregate total number 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.	Structured
10 For all products, extract sentences from reviews that have a positive or negative sentiment.	Unstructured
11 For a given product, measure the correlation of sentiments, including the number of reviews and average review ratings, on monthly product revenues within a given time frame.	Semi-structured
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.	Semi-structured
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.	Structured
14 Calculate the ratio between the number of items sold on the Web 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.	Structured
15 Find the categories with flat or declining sales for in-store purchases during a given year for a given store.	Structured
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.	Structured
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.	Structured
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.	Unstructured
19 Retrieve the items with the highest number of returns, where the number of returns was approximately equivalent across all store and Web channels.	Unstructured
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.	Structured
21 Get all items sold in stores in a given month and year, which were returned in the next six months and repurchased by the returning customer afterward through the Web sales channel in the following three years.	Structured
22 Compute the percentage change in inventory between the 30-day period before the price change and the 30-day period after the change.	Structured
23 Calculate the coefficient of variation and mean of every item and warehouse of two consecutive months. Find items that had a coefficient of variation in the first months of 1.5 or larger.	Structured
24 For a given product, measure the effect of competitors' prices on the product's in-store and online sales.	Structured
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.	Structured
26 Cluster customers into book buddies or groups based on their in-store book purchasing histories.	Structured
27 Extract competitor product names and model names (if any) from online product reviews for a given product.	Unstructured
28 Build text classifier for online review sentiment classification (Positive, Negative, Neutral).	Unstructured
29 Perform category affinity analysis for products purchased together online.	Structured
30 Perform category affinity analysis for products viewed together online.	Semi-structured

Appendix C: Benchmark Tuning Parameters

This appendix provides the detailed tuning parameters used in the benchmark tests described earlier in this document.

BigBench Tuning Parameters

CDP Private Cloud Base includes many default parameter settings. Table C1 shows only the tuning parameters that we changed from the default setting to achieve the best performance from the testing cluster. Table C2 lists the optimizations for each query from [Appendix B](#).

Table C1. BigBench Tuning Parameters

Property	Value 1st Gen Intel® Xeon® Scalable Processor	Value 3rd Gen Intel® Xeon® Scalable Processor
YARN		
yarn.nodemanager.resource.memory-mb	370 GB	500 GB
yarn.scheduler.maximum-allocation-mb	370 GB	500 GB
yarn.nodemanager.resource.cpu-vcores	80	112
yarn.scheduler.maximum-allocation-vcores	80	112
hive_client_java_heapsize	4 GB	4 GB
node_manager_java_heapsize	4 GB	4 GB
Hadoop Distributed File System (HDFS)		
hdfs.block.size	256 MB	256 MB
dfs.socket.timeout	63000	63000
dfs.datanode.socket.write.timeout	63000	63000
dfs.datanode.handler.count	30	30
dfs.namenode.handler.count	300	300
dfs.namenode.service.handler.count	300	300
MapReduce		
mapreduce.map.memory.mb	4352 MB	4352 MB
mapreduce.reduce.memory.mb	4352 MB	4352 MB
mapreduce_map_java_opts_max_heap	4 GB	4 GB
mapreduce_reduce_java_opts_max_heap	4 GB	4 GB
io.sort.mb	512	
Tez		
tez.am.resource.memory.mb	3 GB	3 GB
tez.task.resource.memory.mb	4 GB	4 GB
tez.runtime.unordered.output.buffer.size-mb	409 MB	409 MB
hive.tez.container.size	4096	4096
tez.runtime.io.sort.mb	1 GB	1 GB

Table C2. Query Optimizations for 1st Gen Intel® Xeon® Scalable Processors

Query	1st Gen Intel® Xeon® Scalable Processor
1	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=536870912; set hive.exec.reducers.bytes.per.reducer=8000000;
2	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=13217728; set hive.exec.reducers.bytes.per.reducer=1024000000;
3	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=138435456; --set hive.exec.reducers.bytes.per.reducer=16096000000; set hive.exec.reducers.bytes.per.reducer=896000000;
4	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=268435456; set hive.exec.reducers.bytes.per.reducer=2048000000;
5	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=268435456; set hive.exec.reducers.bytes.per.reducer=4096000000;
6	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=268435456; set hive.exec.reducers.bytes.per.reducer=512000000;
7	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=536870912; set hive.exec.reducers.bytes.per.reducer=128000;
8	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=268435456; set hive.exec.reducers.bytes.per.reducer=256000000;
9	set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=67108864; set hive.exec.reducers.bytes.per.reducer=67108864;
10	None
11	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=3221225472; set hive.exec.reducers.bytes.per.reducer=2000000;
12	set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set tez.grouping.max-size=9934592; set tez.grouping.min-size=435456; set hive.exec.reducers.bytes.per.reducer=200000;
13	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=536870912; set hive.exec.reducers.bytes.per.reducer=256000000;
14	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=2147483648; set hive.exec.reducers.bytes.per.reducer=16000;
15	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=2147483648; set hive.exec.reducers.bytes.per.reducer=2000000;
16	set tez.grouping.max-size=8589934592; set tez.grouping.min-size=268435456; set hive.exec.reducers.bytes.per.reducer=2048000000;

Query	1st Gen Intel® Xeon® Scalable Processor
17	<pre>set tez.grouping.max-size=8589934592; set tez.grouping.min-size=1073741824; set hive.exec.reducers.bytes.per.reducer=1024000;</pre>
18	<pre>set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=16777216; set hive.exec.reducers.bytes.per.reducer=256000000;</pre>
19	<pre>set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set hive.tez.container.size=14240; set hive.tez.java.opts=-Xmx10192m; set tez.runtime.io.sort.mb=2560; set tez.runtime.unordered.output.buffer.size-mb=1024; set hive.auto.convert.join.noconditionaltask.size=3579139413; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=32108864; set hive.exec.reducers.bytes.per.reducer=32000;</pre>
20	<pre>set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=536870912; set hive.exec.reducers.bytes.per.reducer=1024000000;</pre>
21	<pre>set tez.grouping.max-size=8589934592; set tez.grouping.min-size=1073741825; set hive.exec.reducers.bytes.per.reducer=128000000;</pre>
22	<pre>set tez.grouping.max-size=8589934592; set tez.grouping.min-size=67108864; set hive.exec.reducers.bytes.per.reducer=256000000;</pre>
23	<pre>set tez.grouping.max-size=8589934592; set tez.grouping.min-size=33554432; set hive.exec.reducers.bytes.per.reducer=128000000;</pre>
24	<pre>set tez.grouping.max-size=8589934592; set tez.grouping.min-size=2147483648; set hive.exec.reducers.bytes.per.reducer=16000;</pre>
25	<pre>set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=536870912; set hive.exec.reducers.bytes.per.reducer=512000000;</pre>
26	<pre>set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=536870912; set hive.exec.reducers.bytes.per.reducer=64000000;</pre>
27	<pre>set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=536870912;</pre>

Query	1st Gen Intel® Xeon® Scalable Processor
28	<pre>set bigbench.hive.optimize.sampling.orderby=\${hiveconf:bigbench.hive.optimize.sampling.orderby}; set bigbench.hive.optimize.sampling.orderby.number=\${hiveconf:bigbench.hive.optimize.sampling.orderby.number}; set bigbench.hive.optimize.sampling.orderby.percent=\${hiveconf:bigbench.hive.optimize.sampling.orderby.percent}; set bigbench.hive.optimize.sampling.orderby=false; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=536870912; set hive.exec.reducers.bytes.per.reducer=67108864;</pre>
29	<pre>set hive.tez.container.size=6144; set tez.grouping.max-size=8589934592; set tez.grouping.min-size=276870912; set hive.exec.reducers.bytes.per.reducer=256000000;</pre>
30	<pre>set tez.grouping.max-size=8589934592; set tez.grouping.min-size=268435456; set hive.exec.reducers.bytes.per.reducer=512000000;</pre>

Solution Provided By:

CLOUDERA

SOLIDIGM

intel

¹ Tested by Intel June 21, 2022.

Workload: BigBench v1.6.0 @ 3TB – 2 streams (BBqPM), warm run method, minimum of 2 iterations

Software: Cloudera Data Platform v7.1.7 – CDP Private Cloud Base, Hive v 3.1.3000.7.1.7.0-551, Tez v0.9.1, Hadoop v 3.1.1.7.1.7.0-551, Jdk v enjdk version 1.8.0_312, Spark v 2.4.7.1.7.0-551

1st Gen Intel® Xeon® Scalable processor configuration (baseline): 4 worker nodes and 1 master node.

2x Intel Xeon Gold 6148 (20 cores, 2.4 GHz), 384 GB (12x 32 GB @ 2666 MT/s DDR4), Intel® Hyper-Threading Technology = ON, Intel® Turbo Boost Technology = ON, OS = CentOS Linux release 7.9.2009 (Core), kernel = 3.10.0-1160.el7.x86_64, BIOS = SE5C620.86B.02.01.0015.032120220358, microcode = 0x2006c0a, 7x 1.5 TB SSD DC S3610 (SATA), Intel® Ethernet Network Adapter X722 for 10GBASE-T (10 GbE), 80 vCores per node.

3rd Gen Intel Xeon Scalable processor configuration (upgraded): 4 worker nodes and 1 master node.

Common settings across worker and master nodes: Intel Hyper-Threading Technology = ON, Intel Turbo Boost Technology = ON, OS = CentOS Linux release 7.9.2009 (Core), kernel = 3.10.0-1160.45.1.el7.x86_64, BIOS = WLYDCRB1.SYS.0020.P21.2012150710, microcode = 0xd0002a0, Intel Ethernet Adapter E810-XXVDA2 (25 GbE).

Worker nodes: 2x Intel Xeon Gold 6326 processor (16 cores, 2.90 GHz), 256 GB (16x 16 GB @ 3200 MT/s DDR4), 2x 1.6 TB P4510 for NameNode, 64 vCores per node.

Master node: 2x Intel Xeon Gold 6348 processor (28 cores, 2.60 GHz), 512 GB (16x 32 GB @ 3200 MT/s DDR4), 6x 1.6 TB P4510 for DataNode, 1x 3.2 TB P5600 for YARN cache, 112 vCores per node.

² Source: gartner.com/en/newsroom/press-releases/2018-02-05-gartner-survey-shows-organizations-are-slow-to-advance-in-data-and-analytics

³ See endnote 1.

⁴ Source: cloudera.com/about/customers/iqvia.html

⁵ See endnote 1.

⁶ See endnote 1.

⁷ For more information about CDP Private Cloud, visit docs.cloudera.com/cdp-private-cloud/latest/index.html.

⁸ See endnote 1.

⁹ See endnote 1.

See [Appendix A](#) and [Appendix C](#) for additional configuration information.

All Solidigm SSDs referenced were previously known as Intel SSDs.

Performance varies by use, configuration and other factors. Learn more at intel.com/PerformanceIndex. Your costs and results may vary. Intel technologies may require enabled hardware, software or service activation. Performance results are based on testing as of June 21, 2022 and may not reflect all publicly available updates. See configuration disclosure for details. No product or component can be absolutely secure. All references to NAND and SATA SSDs in this document are references to SSDs formerly manufactured by Intel. Notice Revision #20110804. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest Intel product specifications and roadmaps. © Intel Corporation. All rights reserved. Intel, the Intel logo and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others. 1022/IVEN/KC/PDF 341705-004US