

# RED HAT CEPH STORAGE ON SERVERS WITH INTEL® PROCESSORS AND SSDs

Configuring scalable workload-optimized Ceph clusters



Ceph has been developed to deliver object, file, and block storage in one self-managing, self-healing platform with no single point of failure.

Red Hat® Ceph Storage offers multi-petabyte software-defined storage for the enterprise, across a range of industry-standard hardware.

Intel® Xeon® Processor-based servers equipped with Intel® SSD DC Series are ideally suited for Red Hat Ceph Storage clusters.

With proper configuration, Red Hat Ceph Storage clusters can be designed for IOPS-optimized, throughput-optimized, or cost/capacity-optimized workloads.



facebook.com/redhatinc  
@redhatnews  
linkedin.com/company/red-hat

## EXECUTIVE SUMMARY

Ceph users frequently request simple, optimized cluster configurations for different workload types. Common requests are for throughput-optimized and capacity-optimized workloads, but IOPS-intensive workloads on Ceph are also emerging. Based on extensive testing by Red Hat and Intel with a variety of hardware providers, this document provides general performance, capacity, and sizing guidance for servers based on Intel® Xeon® processors, optionally equipped with Intel® Solid State Drive Data Center (Intel® SSD DC) Series.

## TABLE OF CONTENTS

<b>1 INTRODUCTION</b> .....	<b>2</b>
<b>2 CEPH ARCHITECTURE OVERVIEW</b> .....	<b>3</b>
<b>3 CLUSTER CONFIGURATION GUIDANCE</b> .....	<b>4</b>
3.1 Qualifying the need for scale-out storage .....	4
3.2 Identifying target workload I/O profiles .....	5
3.3 Choosing a storage access method .....	6
3.4 Identifying capacity needs .....	7
3.5 Determining fault domain risk tolerance .....	8
3.6 Selecting a data protection method .....	9
<b>4 INTEL HARDWARE CONFIGURATION GUIDELINES</b> .....	<b>10</b>
4.1 Monitor nodes .....	10
4.2 OSD hosts .....	11
4.3 Configuration guidance for Intel processor-based servers .....	14
<b>5 CONCLUSION</b> .....	<b>15</b>

## INTRODUCTION

Storage infrastructure is undergoing tremendous change, particularly as organizations deploy storage to support big data and private clouds. Traditional scale-up arrays are limited in scalability, and their complexity can compromise cost-effectiveness. In contrast, software-defined storage infrastructure based on clustered storage servers has emerged as a way to deploy cost-effective and manageable storage at scale, with Ceph among the leading solutions. In fact, cloud storage companies are already using Ceph at near exabyte scale, with expected continual growth. For example, Yahoo estimates that their Ceph-based cloud object store will grow 20-25% annually.<sup>1</sup>

Red Hat Ceph Storage on Intel Xeon processor- and Intel SSD DC Series-based servers significantly lowers the cost of storing enterprise data and helps organizations manage exponential data growth. The software is a robust, petabyte-scale storage platform for those deploying public or private clouds. As a modern storage system for cloud deployments, Red Hat Ceph Storage offers mature interfaces for enterprise block and object storage, making it well suited for active archive, rich media, and cloud infrastructure workloads like OpenStack®.<sup>2</sup> Delivered in a unified self-healing and self-managing platform with no single point of failure, Red Hat Ceph Storage handles data management so businesses can focus on improving application availability, with properties that include:

- Scaling to exabytes
- No single point of failure in the cluster
- Lower capital expenses (CapEx) by running on industry-standard server hardware
- Lower operational expenses (OpEx) by self-managing and self-healing

Many organizations are trying to understand how to configure Intel Xeon processor-based servers for optimized Ceph clusters that meet their unique needs. Red Hat Ceph Storage is able to run on myriad diverse hardware configurations, but designing a successful Ceph cluster requires careful analysis of issues related to application, capacity, and workload. The ability to address dramatically different kinds of I/O workloads within a single Ceph cluster makes understanding these issues paramount to a successful deployment.

After extensive performance and server scalability evaluation and testing with many vendors, Red Hat has developed a proven methodology that helps ask and answer key questions that lead to properly sized and configured scale-out storage clusters based on Red Hat Ceph Storage. Described in greater detail in this guide, the methodology includes:

- Qualifying the need for scale-out storage
- Identifying target workload I/O profiles
- Choosing a storage access method
- Identifying capacity
- Determining fault domain risk tolerance
- Selecting a data protection method

---

<sup>1</sup> [yahooeng.tumblr.com/post/116391291701/yahoo-cloud-object-store-object-storage-at](http://yahooeng.tumblr.com/post/116391291701/yahoo-cloud-object-store-object-storage-at)

<sup>2</sup> [superuser.openstack.org/articles/openstack-users-share-how-their-deployments-stack-up](http://superuser.openstack.org/articles/openstack-users-share-how-their-deployments-stack-up)

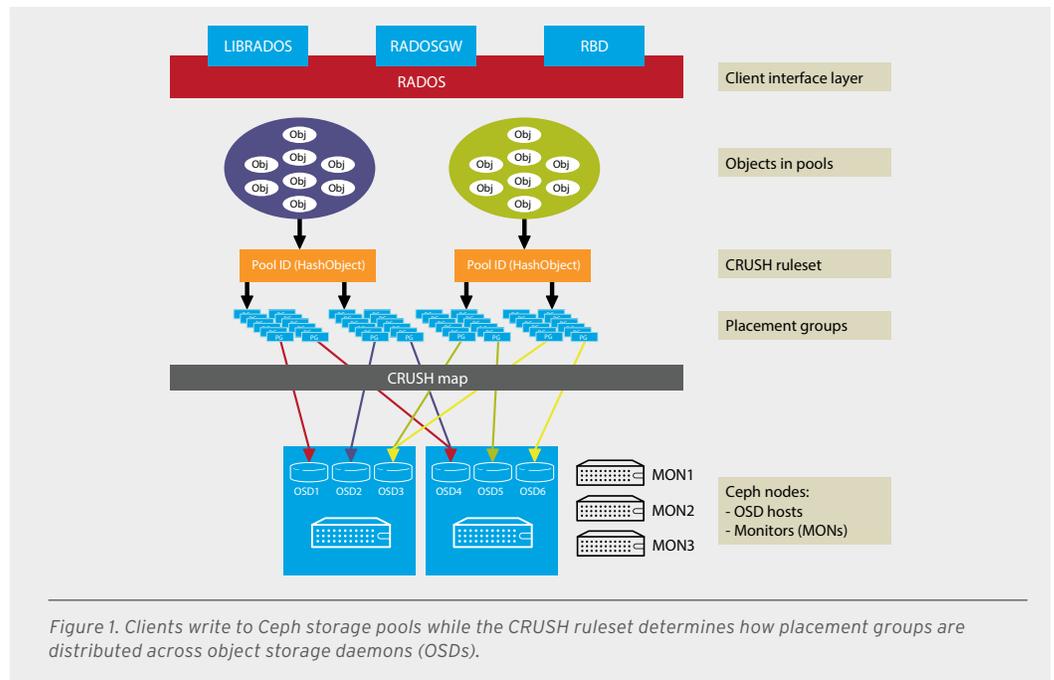
## CEPH ARCHITECTURE OVERVIEW

A Ceph storage cluster is built from large numbers of Ceph nodes for scalability, fault-tolerance, and performance. Each node is based on industry-stand hardware and uses intelligent Ceph daemons that communicate with each other to:

- Store and retrieve data
- Replicate data
- Monitor and report on cluster health
- Redistribute data dynamically (remap and backfill)
- Ensure data integrity (scrubbing)
- Detect and recover from faults and failures

To the Ceph client interface that reads and writes data, a Ceph storage cluster appears as a simple pool where data is stored. However, the storage cluster performs many complex operations in a manner that is completely transparent to the client interface. Ceph clients and Ceph object storage daemons (Ceph OSD daemons, or OSDs) both use the CRUSH (controlled replication under scalable hashing) algorithm for storage and retrieval of objects.

For a Ceph client, the storage cluster is very simple. When a Ceph client reads or writes data (referred to as an I/O context), it connects to a logical storage pool in the Ceph cluster. Figure 1 illustrates the overall Ceph architecture, with concepts that are described in the sections that follow.



- **Pools.** A Ceph storage cluster stores data objects in logical dynamic partitions called pools. Pools can be created for particular data types, such as for block devices, object gateways, or simply to separate user groups. The Ceph pool configuration dictates the number of object replicas and the number of placement groups (PGs) in the pool. Ceph storage pools can be either replicated or erasure coded for added reliability and fault tolerance, as appropriate for the application and cost model. Additionally, pools can “take root” at any position in the CRUSH hierarchy. This property allows placement on groups of servers with differing performance characteristics—allowing optimization for diverse workloads.
- **Placement groups (PGs).** Ceph maps objects to placement groups (PGs). PGs are shards or fragments of a logical object pool that are composed of a group of Ceph OSD daemons that are in a peering relationship. Placement groups provide a way to create replication or erasure coding groups of coarser granularity than on a per object basis. A larger number of PGs (for example, 200 per OSD or more) leads to better balancing.
- **CRUSH ruleset.** The CRUSH algorithm provides controlled, scalable, and declustered placement of replicated or erasure-coded data within Ceph and determines how to store and retrieve data by computing data storage locations. CRUSH empowers Ceph clients to communicate with OSDs directly, rather than through a centralized server or broker. By determining a method of storing and retrieving data by algorithm, Ceph avoids a single point of failure, as well as a performance bottleneck and a physical limit to scalability.
- **Ceph monitors (MONs).** Before Ceph clients can read or write data, they must contact a Ceph MON to obtain the current cluster map. A Ceph storage cluster can operate with a single monitor, but this introduces a single point of failure. For added reliability and fault tolerance, Ceph supports an odd number of monitors in a quorum (typically three or five for small to mid-sized clusters). Consensus among various monitor instances ensures consistent knowledge about the state of the cluster.
- **Ceph OSD daemons.** In a Ceph cluster, Ceph OSD daemons store data and handle data replication, recovery, backfilling, and rebalancing. They also provide some cluster state information to Ceph monitors by checking other Ceph OSD daemons with a heartbeat mechanism. A Ceph storage cluster configured to keep three replicas of every object requires a minimum of three Ceph OSD daemons, two of which must be operational to successfully process write requests. Ceph OSD daemons roughly correspond to a file system on a physical hard disk drive.

## CLUSTER CONFIGURATION GUIDANCE

Despite the flexibility of Ceph, no one cluster or pool configuration fits all applications or situations. Instead, successful configuration of a Ceph cluster requires answering key questions about how the cluster will be used and the applications it will serve.

## QUALIFYING THE NEED FOR SOFTWARE-DESIGNED STORAGE

Not every storage situation calls for scale-out storage. When requirements include several of the following needs, scale-out storage is likely the best solution.

- **Dynamic storage provisioning.** By dynamically provisioning capacity from a pool of storage, organizations are typically building a private storage cloud, emulating services such as Amazon Simple Storage Service (S3) for object storage or Amazon Elastic Block Store (EBS).

- **Standard storage servers.** Scale-out storage employs storage clusters built from industry-standard Intel Xeon processor-based servers rather than proprietary storage appliances, allowing incremental growth of storage capacity and/or performance without forklift appliance upgrades.
- **Unified name spaces.** Scale-out storage allows pooling storage across tens, hundreds, or even thousands of storage servers in one or more unified namespaces.
- **High data availability.** Scale-out storage provides high-availability of data across what would otherwise be “server storage islands” within the storage cluster.
- **Independent scalability of performance and capacity.** Unlike typical scale-up NAS and SAN devices that frequently run out of performance before running out of capacity, scale-out storage allows organizations to add storage performance or capacity incrementally by independently adding more storage servers or disks as required.

### IDENTIFYING TARGET WORKLOAD I/O PROFILES

Accommodating the target workload I/O profile is perhaps the most crucial design consideration. As a first approximation, organizations need to understand if they are simply deploying low-cost archive storage or if their storage needs to meet specific performance requirements. If the lowest cost per terabyte is the overriding need, a Ceph cluster architecture can be designed at dramatically lower costs. For example, Ceph object archives with erasure-coded pools and without dedicated SSD write journals can be dramatically lower in cost than Ceph block devices on 3x-replicated pools with dedicated flash write journals.

For performance-oriented Ceph clusters, however, IOPS, throughput, and latency requirements must be clearly defined. Historically, Ceph has performed very well with high-throughput workloads, and has been widely deployed for these use cases. Applications are frequently characterized by large-block, asynchronous, sequential I/O (such as digital media performance nodes). In contrast, high IOPS workloads are frequently characterized by small-block synchronous random I/O (for example, 4KB random I/O). The use of Ceph for high IOPS open source database workloads is emerging with MySQL, MariaDB, PostgreSQL and other offerings. Moreover, when Ceph is deployed as Cinder block storage for OpenStack virtual machines (VMs), it typically serves a mix of IOPS- and throughput-intensive I/O patterns.

Additionally, understanding the workload read/write mix can affect architecture design decisions. For example, erasure-coded pools can perform better than replicated pools for sequential writes, but perform worse than replicated pools for sequential reads. As a result, a write-mostly object archive workload (such as video surveillance archival) may perform similarly between erasure-coded pools and replicated pools, although erasure-coded pools may be significantly less expensive.

To simplify configuration and testing options and optimally structure optimized cluster configurations, Red Hat categorizes workload profiles as:

- IOPS-optimized
- Throughput-optimized
- Cost/capacity-optimized

Table 1 provides the criteria used to identify optimal Red Hat Ceph Storage cluster configurations, including their properties and example uses. These categories are provided as general guidelines for hardware purchasing and configuration decisions, and can be adjusted to satisfy unique workload blends. As the workload mix varies from organization to organization, actual hardware configurations chosen will vary.

As previously mentioned, a single Ceph cluster can be configured to have multiple pools that serve different workloads. For example, OSDs on IOPS-optimized servers can be configured into a pool serving MySQL workloads, while OSDs on throughput-optimized servers can be configured into a pool serving digital media performance workloads.

**TABLE 1. CEPH CLUSTER OPTIMIZATION CRITERIA.**

OPTIMIZATION CRITERIA	PROPERTIES	EXAMPLE USES
<b>IOPS-OPTIMIZED</b>	<ul style="list-style-type: none"> <li>• Lowest cost per IOPS</li> <li>• Highest IOPS</li> <li>• Meets minimum fault domain recommendation (single server is less than or equal to 10% of the cluster)</li> </ul>	<ul style="list-style-type: none"> <li>• Typically block storage</li> <li>• 3x replication on hard disk drive (HDD) or 2x replication on Intel® SSD DC Series</li> <li>• MySQL on OpenStack clouds</li> </ul>
<b>THROUGHPUT-OPTIMIZED</b>	<ul style="list-style-type: none"> <li>• Lowest cost per given unit of throughput</li> <li>• Highest throughput</li> <li>• Highest throughput per BTU</li> <li>• Highest throughput per watt</li> <li>• Meets minimum fault domain recommendation (single server is less than or equal to 10% of the cluster)</li> </ul>	<ul style="list-style-type: none"> <li>• Block or object storage</li> <li>• 3x replication</li> <li>• Active performance storage for video, audio, and images</li> <li>• Streaming media</li> </ul>
<b>CAPACITY-OPTIMIZED</b>	<ul style="list-style-type: none"> <li>• Lowest cost per TB</li> <li>• Lowest BTU per TB</li> <li>• Lowest watt per TB</li> <li>• Meets minimum fault domain recommendation (single server is less than or equal to 15% of the cluster)</li> </ul>	<ul style="list-style-type: none"> <li>• Typically object storage</li> <li>• Erasure coding common for maximizing usable capacity</li> <li>• Object archive</li> <li>• Video, audio, and image object archive repositories</li> </ul>

### CHOOSING A STORAGE ACCESS METHOD

Choosing a storage access method is another important design consideration. All data in Ceph is stored in pools—regardless of type. The data itself is stored in the form of objects via the Reliable Autonomic Distributed Object Store (RADOS) layer (Figure 2) to:

- Avoid a single point of failure.
- Provides data consistency and reliability.
- Enable data replication and migration.
- Offer automatic fault detection and recovery.

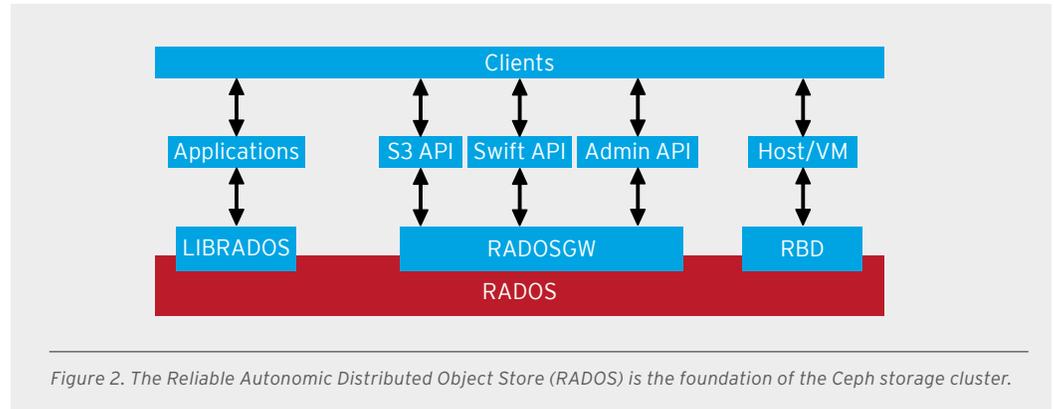


Figure 2. The Reliable Autonomic Distributed Object Store (RADOS) is the foundation of the Ceph storage cluster.

Writing and reading data in a Ceph storage cluster is accomplished using the Ceph client architecture. Ceph clients differ from competitive offerings in how they present data storage interfaces. A wide range of access methods are supported, including:

- **RADOSGW.** RADOSGW is a bucket-based object storage gateway service with S3 compatible and OpenStack Swift compatible RESTful interfaces.
- **LIBRADOS.** LIBRADOS provides direct access to RADOS with libraries for most programming languages, including C, C++, Java™, Python, Ruby, and PHP.
- **RADOS Block Device (RBD).** RBD offers a Ceph block storage device that mounts like a physical storage drive for use by both physical and virtual systems (with a Linux® kernel driver, KVM/QEMU storage backend, or user-space libraries).

Storage access method and data protection method (discussed later in this document) are interrelated. For example, Ceph block storage is currently only supported on replicated pools, while Ceph object storage is supported on both erasure-coded and replicated pools. Due to a significant difference in media costs, replicated architectures are categorically more expensive than erasure-coded architectures. Note that while CephFS distributed file storage is not yet supported on Red Hat Ceph Storage<sup>3</sup>, file systems are routinely created on top of Ceph block devices.

### IDENTIFYING CAPACITY NEEDS

Identifying storage capacity may seem trivial, but it can have a significant effect on the chosen target server architecture. In particular, storage cluster’s predicted capacity needs must be considered together with fault domain risk tolerance and other capabilities. For example, minimum server fault domain recommendations for a small half-petabyte cluster will prevent the use of ultra-dense storage servers in the architecture. Doing so avoids unacceptable fault domain risk on a small number of very large nodes. Table 2 lists broad server sizing trends, with typical types of servers categorized by both workload optimization and overall cluster size.

<sup>3</sup> As of this document’s publication.

TABLE 2. BROAD SERVER SIZING TRENDS.

OPTIMIZATION CRITERIA	OPENSTACK STARTER (64TB)	SMALL (250TB)	MEDIUM (1PB)	LARGE (2PB)
IOPS-OPTIMIZED	<ul style="list-style-type: none"> <li>• Servers with 2-4x PCIe/NVMe slots, or</li> <li>• Servers with 8-12x 2.5-inch SSD bays (SAS/SATA)</li> </ul>		<ul style="list-style-type: none"> <li>• Not typical</li> </ul>	<ul style="list-style-type: none"> <li>• Not typical</li> </ul>
THROUGHPUT-OPTIMIZED	<ul style="list-style-type: none"> <li>• Servers with 12-16x 3.5-inch drive bays</li> </ul>		<ul style="list-style-type: none"> <li>• Servers with 24-36x 3.5-inch drive bays</li> </ul>	<ul style="list-style-type: none"> <li>• Servers with 24-36x 3.5-inch drive bays</li> </ul>
CAPACITY-OPTIMIZED				<ul style="list-style-type: none"> <li>• Servers with 60-72x 3.5-inch drive bays</li> </ul>

### DETERMINING FAULT DOMAIN RISK TOLERANCE

It may be tempting to deploy the largest servers possible in the interest of economics. However, production environments need to provide reliability and availability for the applications they serve, and this necessity extends to the scale-out storage upon which they depend. The fault domain that a single OSD server represents is key to cluster design. As a result, dense servers should be reserved for multi-petabyte clusters where the capacity of an individual server accounts for less than 10-15% of the total cluster capacity. This recommendation may be relaxed for less critical pilot projects. Primary factors for weighing fault domain risks include:

- **Reserving capacity for self-healing.** When a storage node fails, Ceph self-healing begins after a configured time period. For successful self-healing, the unused storage capacity of the surviving cluster nodes must be greater than the used capacity of the failed server. For example, in a 10-node cluster, each node should reserve 10% unused capacity for self-healing of a failed node (in addition to reserving 10% for statistical deviation due to using algorithmic placement). As a result, each node in a cluster should operate at less than 80% of total capacity.
- **Accommodating impact on performance.** During self-healing, a percentage of cluster throughput capacity will be diverted to reconstituting object copies from the failed node on the surviving nodes. The percentage of cluster performance degradation is a function of the number of nodes in the cluster and how Ceph is configured. More nodes in the cluster results in less impact per node.

Ceph will automatically recover by re-replicating the data from the failed node using secondary copies on other nodes in the cluster. As a result, a node failure has several effects:

- Total cluster capacity is reduced by some fraction.
- Total cluster throughput is reduced by some fraction.
- The cluster enters an I/O-heavy recovery process, temporarily diverting an additional fraction of the available throughput.

The time required for the recovery process is directly proportional to how much data was on the failed node and how much throughput the rest of the cluster can sustain. A general formula for calculating recovery time in a Ceph cluster given one disk per OSD is:

$$\text{Recovery time seconds} = (\text{disk capacity in gigabits} / \text{network speed}) / (\text{nodes} - 1)$$

For example, if a 2TB OSD node fails in a 10-node cluster with a 10 GbE (Gigabit Ethernet) back-end, the cluster will take approximately three minutes to recover with 100% of the network bandwidth and no CPU overhead. In practice, using 20% of the available 10 GbE network, the cluster will take approximately 15 minutes to recover, and that time will double with a 4TB drive.

Red Hat recommends the following minimum cluster sizes:

- **Supported minimum cluster size:** Three storage (OSD) servers, suitable for use cases with higher risk tolerance for performance degradation during node failure recovery
- **Recommended minimum cluster size (IOPS- and throughput-optimized cluster):** 10 storage (OSD) servers
- **Recommended minimum cluster size (cost/capacity-optimized cluster):** Seven storage (OSD) servers

There are also other considerations related to fault domain risk and performance that must be considered. Ceph replicates objects across multiple nodes in a storage cluster to provide data redundancy and higher data availability. When designing a cluster, it is important to ask:

- Should the replicated node be in the same rack or multiple racks to avoid a single rack failure?
- Should Ceph OSD traffic stay within the rack or span across racks in a dedicated or shared network?
- Are the application servers in the rack or datacenter proximate to the storage nodes?
- How many concurrent node failures can be tolerated?

Automatic and intelligent placement of object replicas across server, rack, row, and datacenter fault domains can be governed by CRUSH ruleset configuration parameters.

## SELECTING A DATA PROTECTION METHOD

As a design decision, choosing the data protection method can affect the solution's total cost of ownership (TCO) more than any other factor. The chosen data protection method strongly affects the amount of raw storage capacity that must be purchased to yield the desired amount of usable storage capacity.

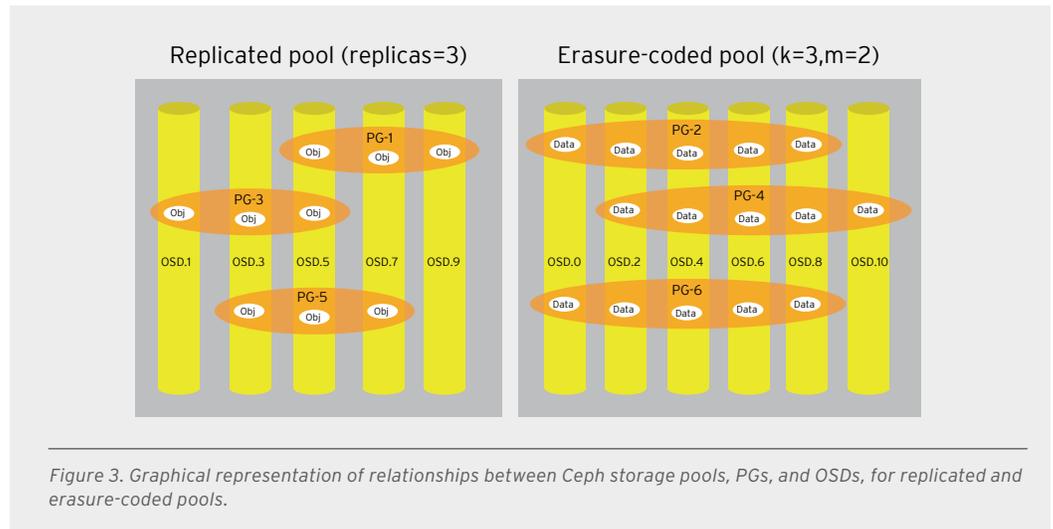
Applications have diverse needs for performance and availability. As a result, Ceph provides data protection at the storage pool level.

- **Replicated storage pools.** Replication makes full copies of stored objects, and is ideal for quick recovery. In a replicated storage pool, Ceph configuration defaults to a replication factor of three, involving a primary OSD and two secondary OSDs. If two of the three OSDs in a placement group become unavailable, data may be read, but write operations will be suspended until at least two OSDs are operational.
- **Erasured-coded storage pools.** Erasure coding provides a single copy of data plus parity and is useful for archive storage and cost-effective durability and availability. With erasure coding, storage pool objects are divided into chunks using the  $n=k+m$  notation, where  $k$  is the number data chunks that are created,  $m$  is the number of coding chunks that will be created to provide data protection, and  $n$  is the total number of chunks placed by CRUSH after the erasure coding process.

For more information on Ceph architecture, see the Ceph documentation at [docs.ceph.com/docs/master/architecture/](https://docs.ceph.com/docs/master/architecture/).

Ceph block storage is typically configured with 3x-replicated pools and is currently not supported directly on erasure-coded pools. Ceph object storage is supported on both replicated and erasure-coded pools. Depending on the performance needs and read/write mix of an object storage workload, an erasure-coded pool can provide an extremely cost-effective solution that meets performance requirements.

Figure 3 illustrates the relationships among Ceph storage pools, PGs, and OSDs for both replicated and erasure-coded pools.



## INTEL HARDWARE CONFIGURATION GUIDELINES

The sections that follow provide broad guidelines for the selection of monitor nodes and OSD hosts. Actual configuration of OSD servers can vary based on application and workload optimization.

### MONITOR NODES

The Ceph monitor is a datastore for the health of the entire cluster, and contains the cluster log. A minimum of three monitors are recommended for a cluster quorum. Monitor nodes typically have fairly modest CPU and memory requirements. A single rack unit (1U) server with a low-cost Intel processor (such as an Intel Xeon Processor E5-2603), 16GB of RAM, and GbE networking should suffice in most cases. Since logs are stored on local disk(s) on the monitor node, it is important to make sure that sufficient disk space is provisioned. In addition, the monitor store should be placed on an Intel SSD DC Series, because the leveldb store can become I/O bound.

For example, when monitoring 100 OSDs in a healthy Ceph cluster, each monitor will collect data for all of the OSDs until all the monitors are synchronized with the same Ceph cluster information. The size of the datastore can vary, but 200MB up to 1-2GB are to be expected, depending not only on the size of the cluster, but the state change (churn) that the cluster undergoes. Logs for recovering clusters can grow quickly, reaching dozens and even hundreds of gigabytes. Abnormal monitor datastore growth should be investigated by an operator, as there is usually an underlying condition that should be remedied. Log rotation is a good practice that can guarantee that available disk space is not

blindly consumed, especially if verbose debugging output is set on the monitors, since they will generate a large amount of logging information<sup>4</sup>. In most situations, monitors should be run on distinct nodes or on VMs that reside on physically separate machines to prevent a single point of failure.

## OSD HOSTS

Ceph OSD hosts are configured differently depending on both workload optimization and the data devices installed: HDDs, SSDs, or NVMe SSDs.

### CPU specifications

CPU recommendations for OSD hosts differ depending on the media that is employed.

- For HDD-based OSDs, one core-GHz is recommended for each OSD. For example, 16 HDD-based OSDs can be supported with an Intel® Xeon® Processor E5-2620 v4: **8 cores \* 2.10 GHz = 16.8 core-GHz**.
- For OSDs based on Intel® SSD Data Center P3700 Series NVMe drives, 20 core-GHz are recommended for each drive: **5 core-GHz per OSD \* 4 OSDs per SSD = 20 core-GHz per SSD**. For example, an Intel Xeon Processor E5-2630 v4 (**10 cores \* 2.20 GHz = 22 core-GHz**) is recommended to drive the OSDs for a single Intel SSD DC P3700 Series drive.

### Memory specifications

Red Hat typically recommends a baseline of 16GB of RAM, with an additional 2GB of RAM per OSD. When sizing memory requirements, it is important to consider:

- The number of OSDs per node.
- The number of memory banks available.
- The number of memory channels per bank.
- The cost of DRAM (DIMMs).

### Data devices

OSDs and OSD data drives are independently configurable and Ceph OSD performance is naturally dependent on the throughput and latency of the underlying media. The actual number of OSDs configured per drive depends on the type of media configured on the OSD host. For magnetic storage media, one OSD should be configured per HDD. On the other hand, an IOPS-optimized OSD host with a smaller number of high-speed SSDs might be configured with 2-4 OSDs per SSD to exploit the available I/O bandwidth.

Ceph uses a journal to allow it to create atomic updates, which are required to ensure data consistency. To complete these updates, an OSD writes the data payload and metadata to a write journal before writing to the OSD's data partition. A write is acknowledged to the client after all OSD peers in the PG—in replicated or erasure-coded pools—have successfully written their assigned replica or shard to their write journal. A beneficial side-effect of write journaling may be some potential coalescing of small writes. For performance-optimized clusters, journals are typically located on a partition of a faster media type than the OSD media. For example, a throughput-optimized OSD server typically has HDD-based OSDs, and a dedicated write journal based on an SSD.

---

<sup>4</sup> Refer to Ceph documentation on monitor log settings for additional details.

The suggested ratio between HDD-OSDs and flash write journals are:

- 4-5 HDD-OSDs for each Intel® SSD Data Center S3510 or S3610 400GB journal drive.
- 12-18 HDD-OSDs for each Intel SSD DC P3700 800GB NVMe journal drive.

To help decrease cost for cost/capacity-optimized clusters, journals can be co-located with OSDs on the same HDD via partitions.

### Storage Media

Performance and economics for Ceph clusters both depend heavily on an effective choice of storage media. For throughput- and cost/capacity-optimized clusters, magnetic media currently accounts for the bulk of the deployed storage capacity, but the economics of SSDs is changing rapidly. As previously discussed, SSDs are typically used for Ceph write journaling for throughput-optimized configurations. For cost/capacity-optimized configurations, write journaling is often co-resident on the HDDs.

- **Magnetic media.** Enterprise-, or cloud-class HDDs should be used for Ceph clusters. Desktop-class disk drives are not well suited for Ceph deployments as they lack sufficient rotational vibration (RV) compensation for high-density, high-duty-cycle applications and use cases. When dozens or hundreds of rotating HDDs are installed in close proximity, RV quickly becomes a challenge. Failures, errors, and even overall cluster performance can be adversely affected by the rotation of neighboring disks interfering with the rapidly spinning platters in high density storage enclosures. Enterprise-class HDDs contain higher quality bearings and RV compensation circuitry to mitigate these issues in multi-spindle applications and use cases—especially in densities above 4-6 HDDs in a single enclosure. Both SAS and SATA interface types are acceptable.
- **Solid state media (currently flash).** Ceph is strongly consistent storage, so every write to the Ceph cluster must be written to Ceph journals before the write is acknowledged to the client. The data remain in the journal until all replicas or shards are acknowledged as fully written. Only then will the next write happen. SSD journals let the OSDs write faster, reducing the time before a write acknowledgment is sent to the client. In some cases, several small writes can be coalesced during a single journal flush, which can also improve performance. SSDs can also be used for OSD data as well, as recommended in IOPS-optimized Ceph configurations.

SSD choice is an essential factor, and a key area of focus for Red Hat and Intel. Important criteria to consider when selecting solid state media for Ceph include:

- **Classification.** Only enterprise-class SSDs should be deployed with Ceph and consumer-class SSDs should not be used. Intel recommends Intel SSD Data Center Series drives.
- **Endurance.** Write endurance is important as Ceph write journals are heavily used and could exceed recommended program/erase cycles of an SSD rated for lower endurance. Current popular choices include devices rated at greater than 10 device writes per day (DWPD) for 5 years, translating to a lifetime total of 28PB written (PBW).
- **Power fail protection.** Supercapacitors for power fail protection are vital. In the event of a power failure, supercapacitors must be properly sized to allow the drive to persist all in-flight writes to non-volatile NAND storage. Intel SSD Data Center Series drives offer greater than 2 million Power Loss Imminent (PLI) cycles.<sup>5</sup>

---

<sup>5</sup> [www.intel.com/content/www/us/en/solid-state-drives/benefits-of-ssd.html](http://www.intel.com/content/www/us/en/solid-state-drives/benefits-of-ssd.html)

- **Performance.** For Ceph write journaling, the write throughput rating of the journal device should exceed the aggregate write throughput rating of all underlying OSD devices that are served by that journal device.
- **Reliability.** Intel SSD Data Center Series drives offer self-testing and trusted protection from data loss. With  $10^{17}$  uncorrectable bit-error rate (UBER), these drives are proven to be 100-fold more reliable than consumer SSDs, helping to prevent silent data corruption (SDC).<sup>6</sup> These devices also supercede Joint Electronic Device Engineering Council (JEDEC) annual failure rate (AFR).<sup>7</sup>

### I/O controllers

Servers with JBOD (just a bunch of disks) host bus adapters (HBAs) are generally appropriate for OSD hosts, depending on workload and application expectations. For large block sequential I/O workload patterns, HDDs typically perform better when configured in JBOD mode than as RAID volumes. However, for small block random I/O workload patterns, HDD-based OSDs configured on single-drive RAID 0 volumes provide higher IOPS than when configured in JBOD mode.

Many modern systems that house more than eight drives have SAS expander chips on the drive hot swap backplane. Similar to network switches, SAS expanders often allow connection of many SAS devices to a controller with a limited number of SAS lanes. Ceph node configurations with SAS expanders are well suited for large capacity-optimized clusters. However, when selecting hardware with SAS expanders, consider how the following will affect performance:

- Adding extra latency.
- Oversubscribed SAS lanes.
- Spanning Tree Protocol (STP) overhead of tunneling SATA over SAS

Due to backplane oversubscription or poor design, some servers used for Ceph deployments may encounter sub-par performance in systems that use SAS expanders. The type of controller, expander, and even brand of drive and firmware all play a part in determining performance.

### Network interfaces

Providing sufficient network bandwidth is essential for an effective and performant Ceph cluster. Fortunately, network technology is improving rapidly. Standard Ethernet-based interfaces are now available with ever-increasing bandwidth. In servers employed as OSD hosts, network capacity should generally relate to storage capacity. For smaller OSD hosts with 12-16 drive bays, 10 GbE is typically sufficient. For larger OSD hosts with 24-72 drive bays, 25 or 40 GbE may be preferred to provide the required bandwidth and throughput.

Physical deployment characteristics must also be taken into account. If the nodes are spread across multiple racks in the datacenter, the network design should ensure high bisectional bandwidth and minimal network diameter. Ideally, each OSD server should have two network interfaces for data traffic: one connected to the client systems and another for the private network connecting the OSD servers.

---

<sup>6</sup> [intel.com/content/www/us/en/solid-state-drives/benefits-of-ssd.html](http://intel.com/content/www/us/en/solid-state-drives/benefits-of-ssd.html)

<sup>7</sup> [Intel Data Center SSDs deliver 10<sup>17</sup> UBER, see jedec.org/standards-documents/focus/flash/solid-state-drives.](http://jedec.org/standards-documents/focus/flash/solid-state-drives)

**CONFIGURATION GUIDANCE FOR INTEL PROCESSOR-BASED SERVERS**

Table 3 provides general guidance for configuring Intel-based OSD hosts for Red Hat Ceph Storage.

**TABLE 3. CONFIGURING INTEL-BASED SERVERS FOR RED HAT CEPH STORAGE.**

OPTIMIZATION CRITERIA	OPENSTACK STARTER (100TB)	SMALL (250TB)	MEDIUM (1PB)	LARGE (2PB)
<b>IOPS-OPTIMIZED</b>	<ul style="list-style-type: none"> <li>• Ceph RBD (block) pools</li> <li>• OSDs on 1-4 Intel® SSD Data Center P3700 Series per server with write journals co-located on different partitions</li> <li>• 1x Intel SSD DC P3700 per server: single-socket Intel® Xeon® Processor E5-2630v4 (10 cores)</li> <li>• 2x Intel SSD DC P3700 per server: dual-socket Intel Xeon Processor E5-2630v4 (20 cores)</li> <li>• 4x Intel SSD DC P3700 per server: dual-socket Intel Xeon Processor E5-2695v4 (36 cores)</li> <li>• Data protection: Replication (2x on SSD-based OSDs) with regular backups to the object storage pool</li> <li>• 2-4 OSDs per SSD or NVMe drive</li> </ul>		<ul style="list-style-type: none"> <li>• Not typical</li> </ul>	<ul style="list-style-type: none"> <li>• Not typical</li> </ul>
<b>THROUGHPUT-OPTIMIZED</b>	<ul style="list-style-type: none"> <li>• Ceph RBD (block) or Ceph RGW (object) pools</li> <li>• OSDs on HDDs:               <ul style="list-style-type: none"> <li>• Good: Write journals on Intel® SSD Data Center S3710 400GB drives, with a ratio of 4-5 HDDs to each SSD</li> <li>• Better: Write journals on Intel SSD DC P3700 800GB NVMe drives, with a ratio of 12-18 HDDs to each SSD</li> </ul> </li> <li>• One CPU core-GHz per OSD. For example:               <ul style="list-style-type: none"> <li>• 12 OSD/HDDs/server: Single-socket Intel Xeon Processor E5-2620v4 (8 cores*2.1 GHz)</li> <li>• 36 OSD/HDDs/server: Dual-socket Intel Xeon Processor E5-2630v4 (20 cores*2.2 GHz)</li> <li>• 60 OSD/HDDs/server: Dual-socket Intel Xeon E5-2683v4 (32 cores*2.1 GHz)</li> </ul> </li> <li>• Data protection: Replication (read-intensive or mixed read/write) or erasure-coded (write-intensive)</li> <li>• High-bandwidth networking: Greater than 10 GbE for servers with more than 12-16 drives</li> </ul>			

OPTIMIZATION CRITERIA	OPENSTACK STARTER (100TB)	SMALL (250TB)	MEDIUM (1PB)	LARGE (2PB)
CAPACITY-OPTIMIZED	<ul style="list-style-type: none"> <li>Not typical</li> </ul>	<ul style="list-style-type: none"> <li>Ceph RGW (object) pools</li> <li>OSDs on HDDs with write journals co-located on HDDs in separate partition.</li> <li>One CPU core-GHz per OSD. See throughput-optimized section above for examples.</li> <li>Data protection: Erasure-coded</li> </ul>		

\* All SSDs should be enterprise-class, meeting the requirements noted above.

## CONCLUSION

Intel Xeon processors power a wide range of industry-standard server platforms with diverse capabilities. Selecting the right hardware for target Ceph workloads can be a challenge, and this is especially true for software-defined storage solutions that run on industry-standard hardware. Because every environment differs, the general guidelines for sizing CPU, memory, and storage media per node in this document should be mapped to a preferred vendor's product portfolio for determining appropriate server hardware. Additionally, the guidelines and best practices highlighted in this document are not a substitute for running baseline benchmarks before going into production.

Red Hat and Intel have conducted extensive testing with a number of vendors who supply hardware optimized for Ceph workloads. For specific information on selecting servers for running Red Hat Ceph Storage, refer to the tested configurations documented in the "Red Hat Ceph Storage Hardware Selection Guide." Detailed information including Red Hat Ceph Storage test results can be found in performance and sizing guides for popular hardware vendors.

## TECHNOLOGY DETAIL Red Hat Ceph Storage on Servers with Intel Processors and SSDs

Intel disclaimer: Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. No computer system can be absolutely secure. Check with your system manufacturer or retailer or learn more at [www.intel.com](http://www.intel.com).

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more information go to [www.intel.com/performance](http://www.intel.com/performance).



[facebook.com/redhatinc](https://facebook.com/redhatinc)  
[@redhatnews](https://twitter.com/redhatnews)  
[linkedin.com/company/red-hat](https://linkedin.com/company/red-hat)

[redhat.com](http://redhat.com)  
#INC0406282-0616

## ABOUT RED HAT

Red Hat is the world's leading provider of open source solutions, using a community-powered approach to provide reliable and high-performing cloud, virtualization, storage, Linux, and middleware technologies. Red Hat also offers award-winning support, training, and consulting services. Red Hat is an S&P company with more than 80 offices spanning the globe, empowering its customers' businesses.

**NORTH AMERICA**  
1 888 REDHAT1

**EUROPE, MIDDLE EAST,  
AND AFRICA**  
00800 7334 2835  
[europe@redhat.com](mailto:europe@redhat.com)

**ASIA PACIFIC**  
+65 6490 4200  
[apac@redhat.com](mailto:apac@redhat.com)

**LATIN AMERICA**  
+54 11 4329 7300  
[info-latam@redhat.com](mailto:info-latam@redhat.com)

Copyright © 2016 Red Hat, Inc. Red Hat, Red Hat Enterprise Linux, the Shadowman logo, and JBoss are trademarks of Red Hat, Inc., registered in the U.S. and other countries. The OpenStack® Word Mark and OpenStack Logo are either registered trademarks / service marks or trademarks / service marks of the OpenStack Foundation, in the United States and other countries and are used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation or the OpenStack community. Linux® is the registered trademark of Linus Torvalds in the U.S. and other countries. Intel, the Intel Logo, and Xeon are trademarks of Intel Corporation in the U.S. and/or other countries.