



# ceph File System

---

MINA RADY, [MINARADY@STUDENT.LUT.FI](mailto:MINARADY@STUDENT.LUT.FI)

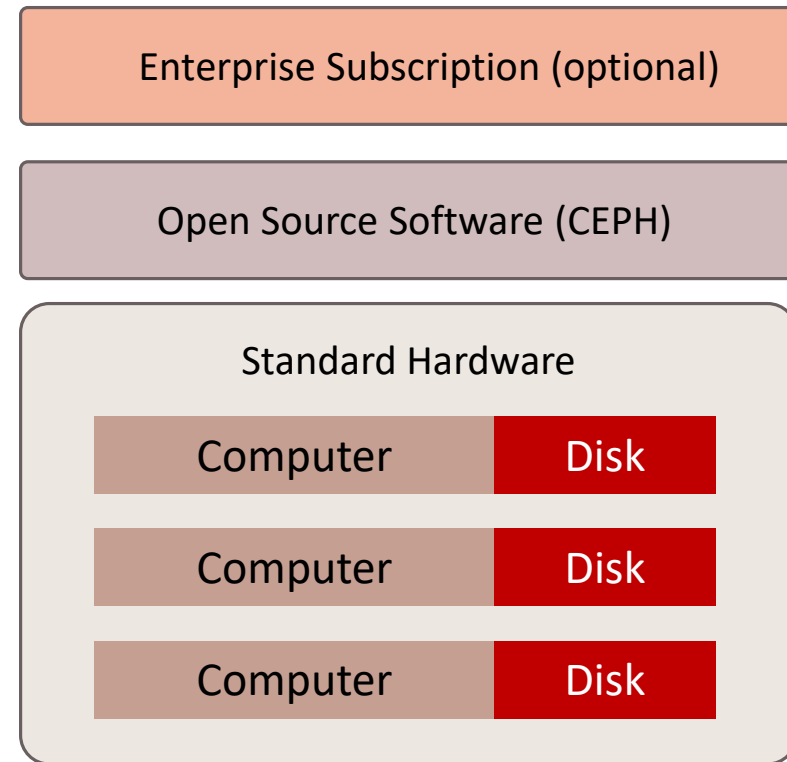
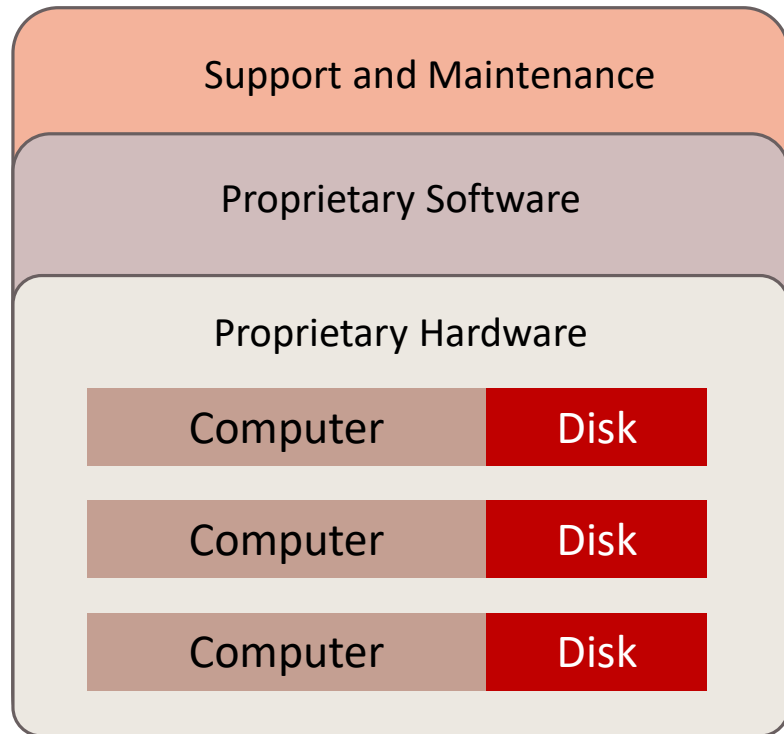
COMPUTING CLUSTERS, COMPUTING GRIDS, COMPUTING CLOUDS COURSE

PROFESSOR ANDREY SHEVEL, ITMO UNIVERSITY

5 JUNE 2017, ST. PETERSBURG, RUSSIA

# Proprietary vs Open Cluster Storage

---



# What is CEPH?

---

- Software Defined Storage Cluster solution that is:
  - Self-healing.
  - Self-managed.
  - “No” bottlenecks.
- Three Interfaces:
  - Object Access
  - Block Access
  - Distributed File System Access.

# Ceph Features

---

- CephFS Offers a traditional file system with interface with POSIX compatibility.
- Importance of POSIX (Portable Operating System Interface):
  - It implements system calls usually used across multiple operating systems (Unix, mainly).
  - Allows safe transition of applications deploying POSIX calls such as: `read`, `write`, `seekdir`.
- Provides stronger data safety for mission-critical applications.
- RADOS: Ceph core is Reliabile Autonomic Distributed Object Store
- No single point of failure
- Software Based

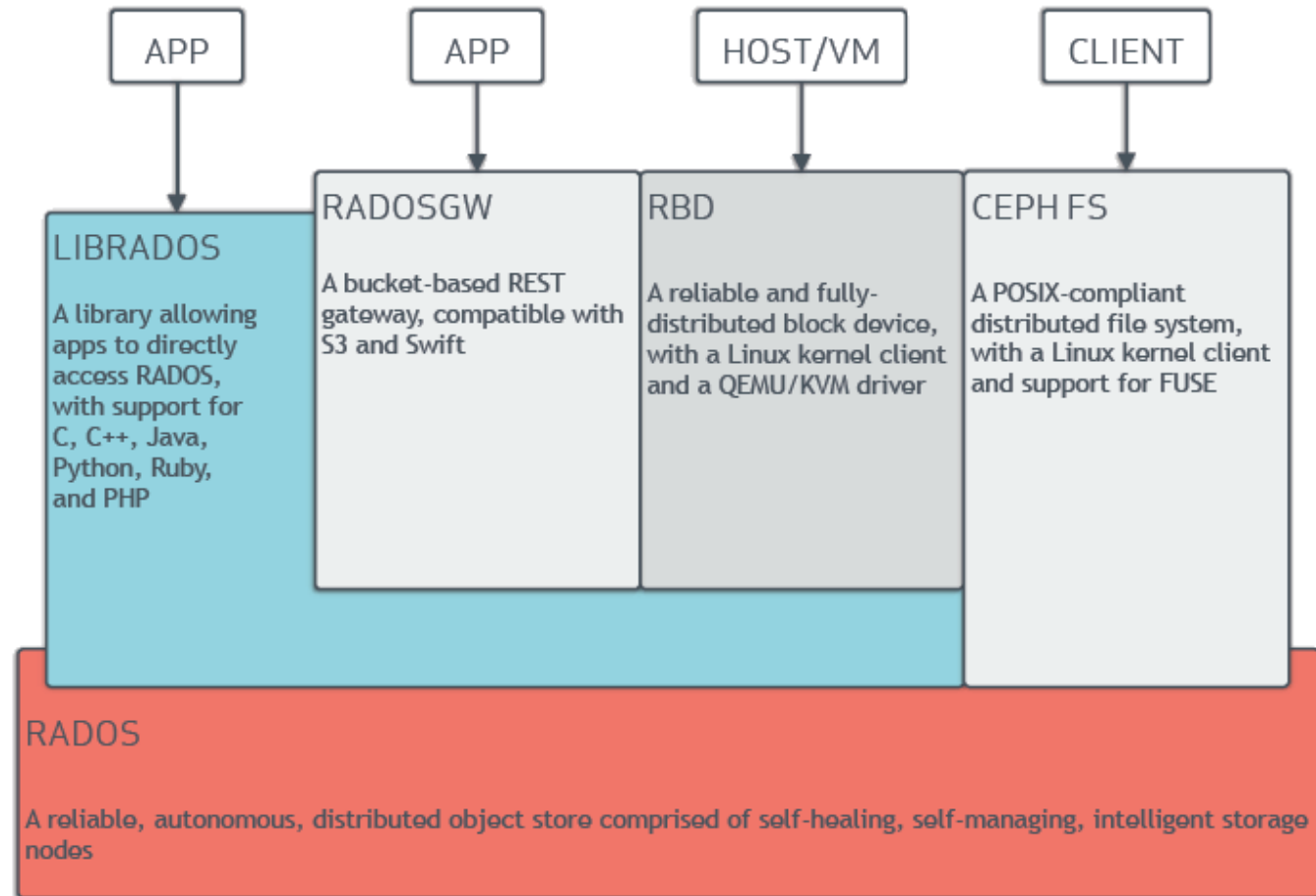
# What is Ceph?: Main Concepts

---

- Any Ceph deployment consists of:
  - OSD:
    - Object Storage Daemon to run on top of each storage node.
    - Each OSD points out to a path in the storage disk and makes it a storage location for Ceph.
    - The path can even point out to a RAID.
    - OSDs responsible for peering, replication, and recovery tasks.
  - Monitor Nodes:
    - First point of contact to clients, after that, client talks to OSDs.
    - Odd number of nodes are needed for consensus and distributed decision making.
  - Librados:
    - Library for accessing RADOS.
    - Uses native sockets to talk to RADOS or to receive commands from C, C++ (no HTTP overload).
    - Offers REST interface.
    - Offers Support for C, C++, Java, PHP, and Python.

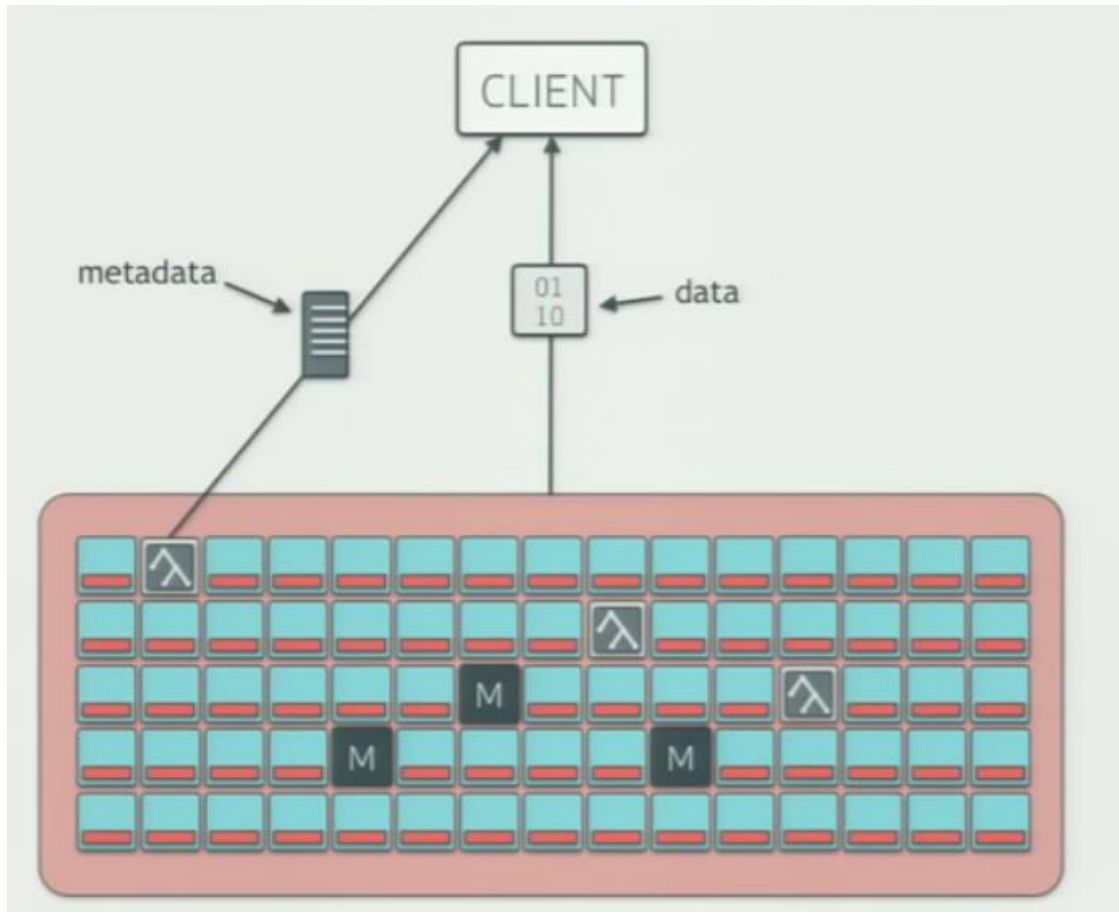
# Ceph Technical Outline

---



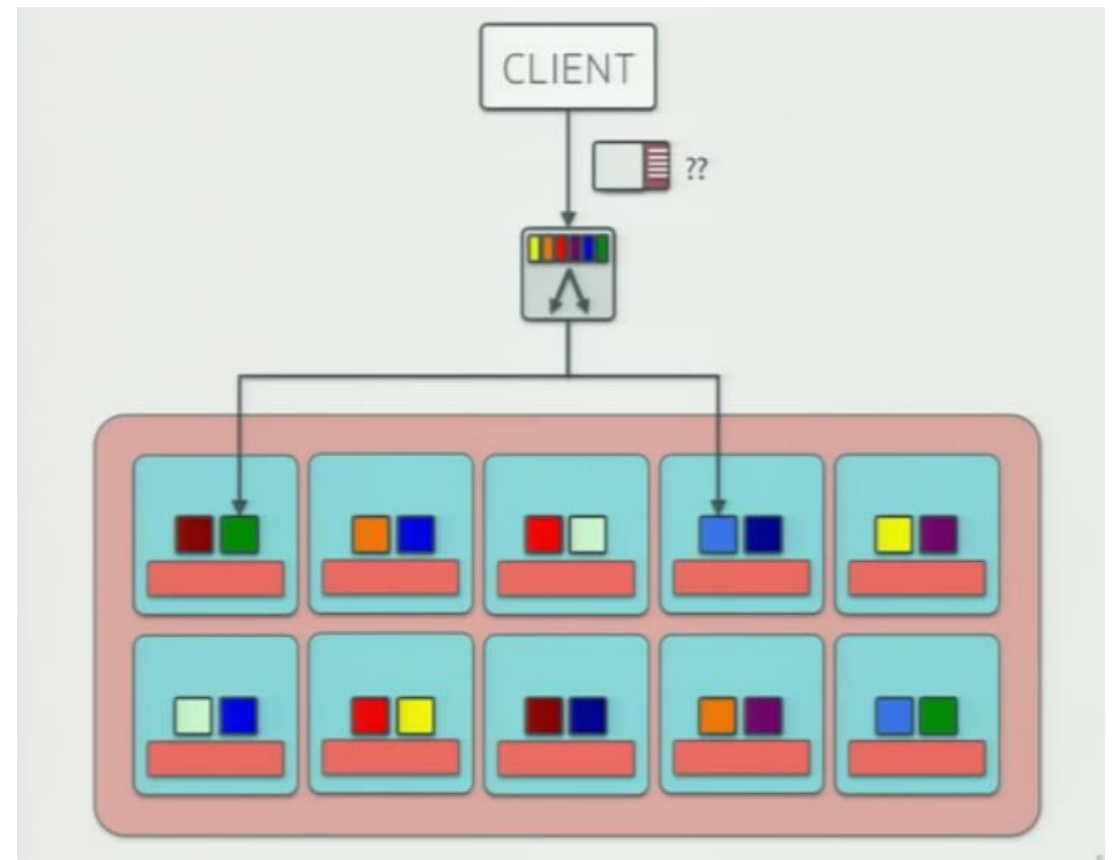
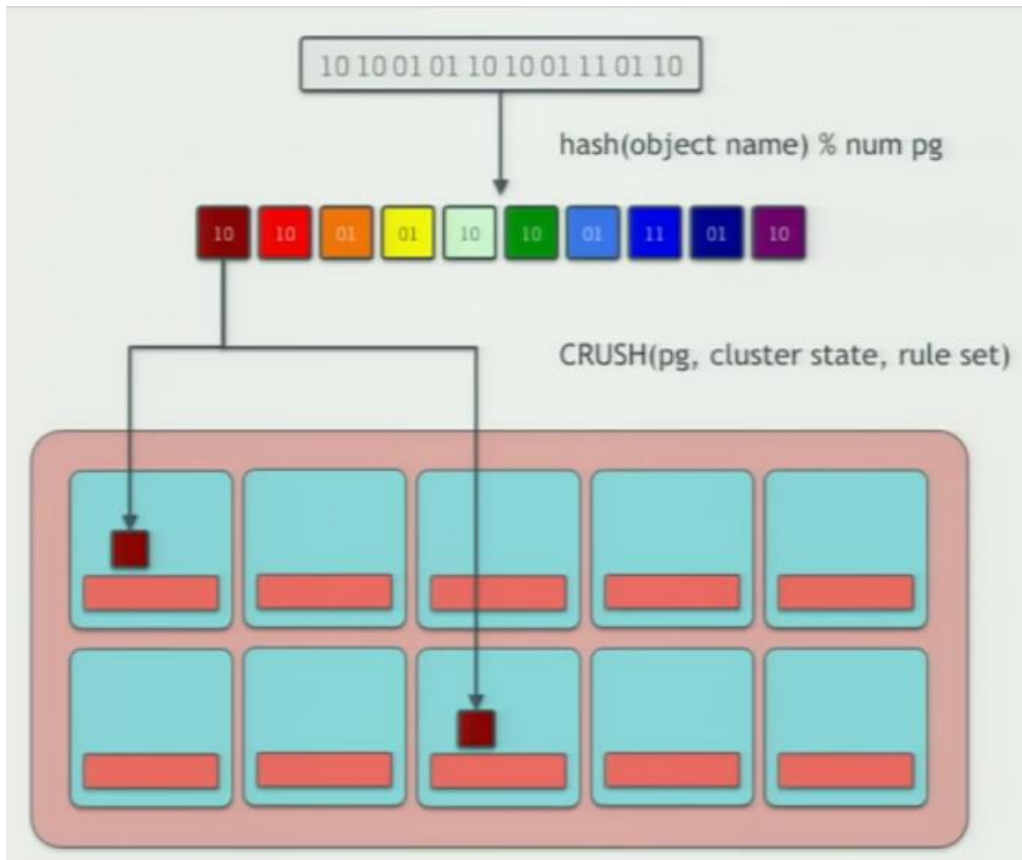
# Ceph Technical Outline

---



# Ceph Uniqueness: CRUSH Algorithm

CRUSH: Controlled Replication Under Scalable Hashing





# CephFS Main Concepts

---

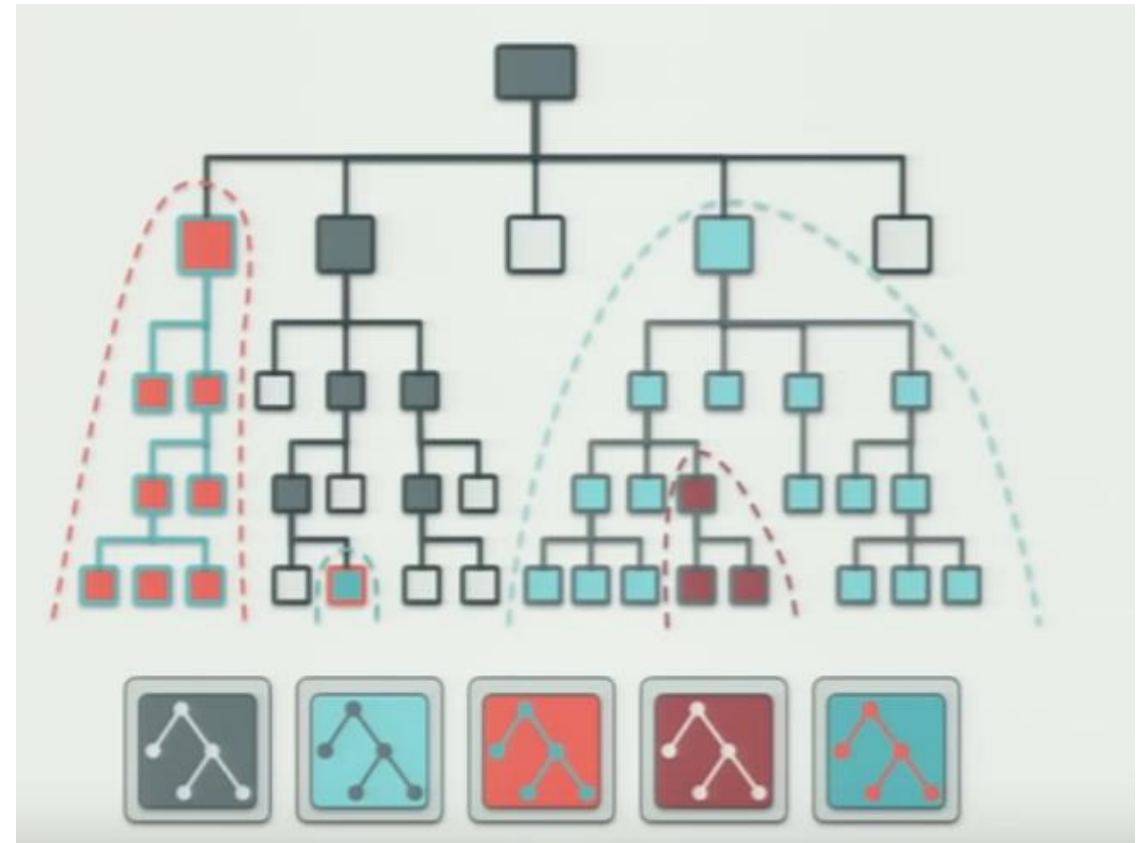
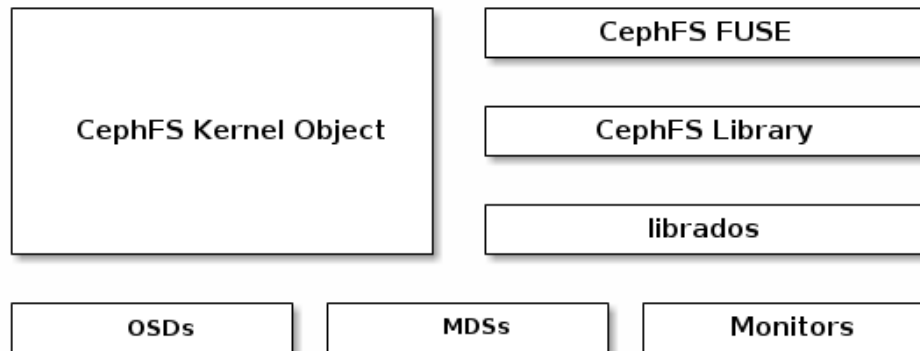
Pools: Logical Containers for Storage Objects

Placement Groups: You don't talk to them directly but you care how many there are, they are hash buckets. It is computationally intensive to map objects to OSDs, so you map them to a placement groups and then you get which OSD a pool is on

One PG spans several OSDs, one OSD serves many PGs

# CephFS Architecture

- Using CephFs require at least one Metadata Server.
- Metadata Servers perform replication using Dynamic Subtree Partitioning. It could even reach one server per file in case of peak demand.



# Deployment Considerations:

---

- Pricing: free source!
- Professional Support: By Kernel Software Inc, Wisconsin, USA: \$3700 storage architecture. \$1700 over existing Open Stack Deployment.
- Intended for large enterprise operations. Not recommended for physically small storage deployments
- It automatically balances storage spaces. Therefore not recommended if special operations are required in the storage clusters.
- Preferable separate public and cluster networks
- Cluster network should be 2x public network bandwidth

# Design Decisions

---

- Hard Disk Parameters:
  - Higher RPM, more energy cost but higher IOPS and throughput
  - Higher Density contribute positively to performance but negatively to maintainability
- Redundancy:
  - Replication: Increased read performances
  - Increases cluster utilization for writes
- Erasure Coding:
  - Customized data and parity bit distribution
  - Better Space and Network Efficiency
  - Higher CPU Overhead
- Scalability
  - Adding more nodes, Ceph will make sure data are evenly distributed across nodes until you hit network topology limitations.

# Further Questions

---

- Smart Grid Metrics in load storage balancing
- Assessment of reliability of Ceph in Grid paradigm for large enterprise to alleviate the need for dedicated storage space.
- Optimization model for data distribution depending on access patterns.

# Demo

---

- Deployment on virtual machine with one main storage disk of 8 GB (dynamic) and 3 storage disks of 1 TB each (dynamic).
- Installed 3 OSDs, 1 Monitor and 1 Metadata Server
- Networking Installation Caveat, not mentioned in installation guidelines: only real internal IP addresses need to be used to refer to hosts.
- Possible improvements?:
  - Adding an option of local host access to avoid load on network interface when possible.
  - Dashboard for analytics, visualization, and integration. JSON API for control?

```
ceph-deploy@ceph-single-node:~/my-cluster$ ceph -s
cluster 71f0420a-ffff-4bf8-ba5c-8ea68b5a48a9
health HEALTH_OK
monmap e1: 1 mons at {ceph-single-node=10.0.2.15:6789/0}
election epoch 4, quorum 0 ceph-single-node
osdmap e15: 3 osds: 3 up, 3 in
flags sortbitwise,require_jewel_osds
pgmap v33: 64 pgs, 1 pools, 0 bytes data, 0 objects
100 MB used, 3055 GB / 3055 GB avail
64 active+clean
```

# Demo

- Health check: `ceph -s`
- Stats: `ceph osd stat`
- OSD List: `ceph osd tree -f json`
- Service Orientation, JSON API: `Ceph osd dump -f json`

```
ceph-single-node [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
ceph-deploy@ceph-single-node:/mnt/mycephfs$ ceph osd stat
osdmap e47: 3 osds: 3 up, 3 in
      flags sortbitwise,require_jewel_osds
ceph-deploy@ceph-single-node:/mnt/mycephfs$ ceph osd tree
ID WEIGHT  TYPE NAME                UP/DOWN REWEIGHT PRIMARY-AFFINITY
-1 2.98380  root default
-2 2.98380  host ceph-single-node
 0 0.99460  osd.0                up 1.00000 1.00000
 1 0.99460  osd.1                up 1.00000 1.00000
 2 0.99460  osd.2                up 1.00000 1.00000
ceph-deploy@ceph-single-node:/mnt/mycephfs$ ceph -s
cluster 71f0420a-ffffb-4bf8-ba5c-8ea68b5a48a9
health HEALTH_OK
monmap e1: 1 mons at {ceph-single-node=10.0.2.15:6789/0}
election epoch 7, quorum 0 ceph-single-node
fsmap e22: 1/1/1 up {0=ceph-single-node=up:active}
osdmap e47: 3 osds: 3 up, 3 in
      flags sortbitwise,require_jewel_osds
pgmap v555: 368 pgs, 9 pools, 4356 bytes data, 191 objects
      121 MB used, 3055 GB / 3055 GB avail
      368 active+clean
ceph-deploy@ceph-single-node:/mnt/mycephfs$ _
```