# Ceph Storage on SSD for Container
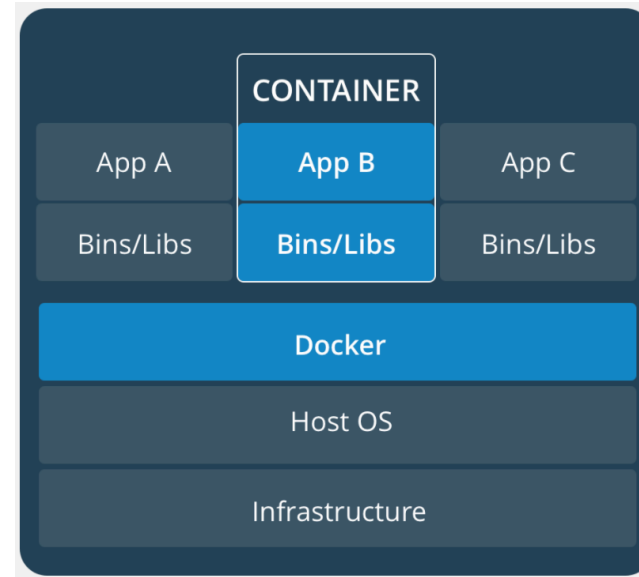
**JANGSEON RYU**
**NAVER**

**NAVER**

# Agenda

- **What is Container?**

- **Persistent Storage**

- **What is Ceph Storage?**

- **Write Process in Ceph**
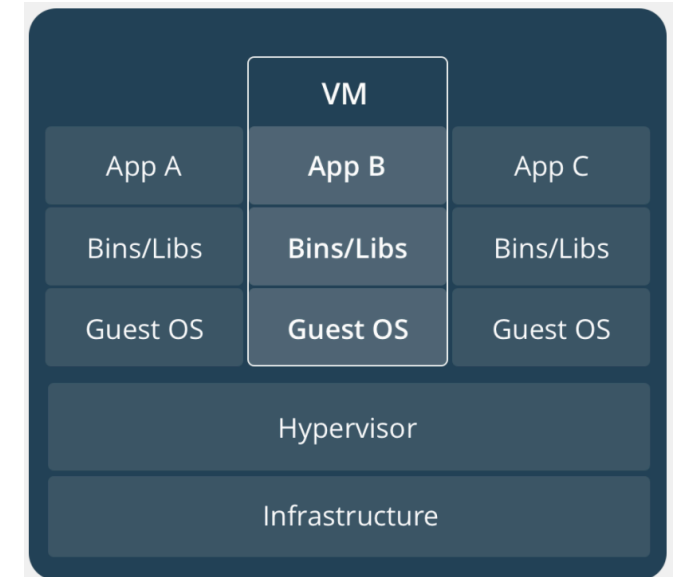
- **Performance Test**

**NAVER**

# What is Container ?

# What is Container?

- **VM vs Container**
- OS Level Virtualization
- Performance
- Run Everywhere – Portable
- Scalable – lightweight
- Environment Consistency
- Easy to manage Image
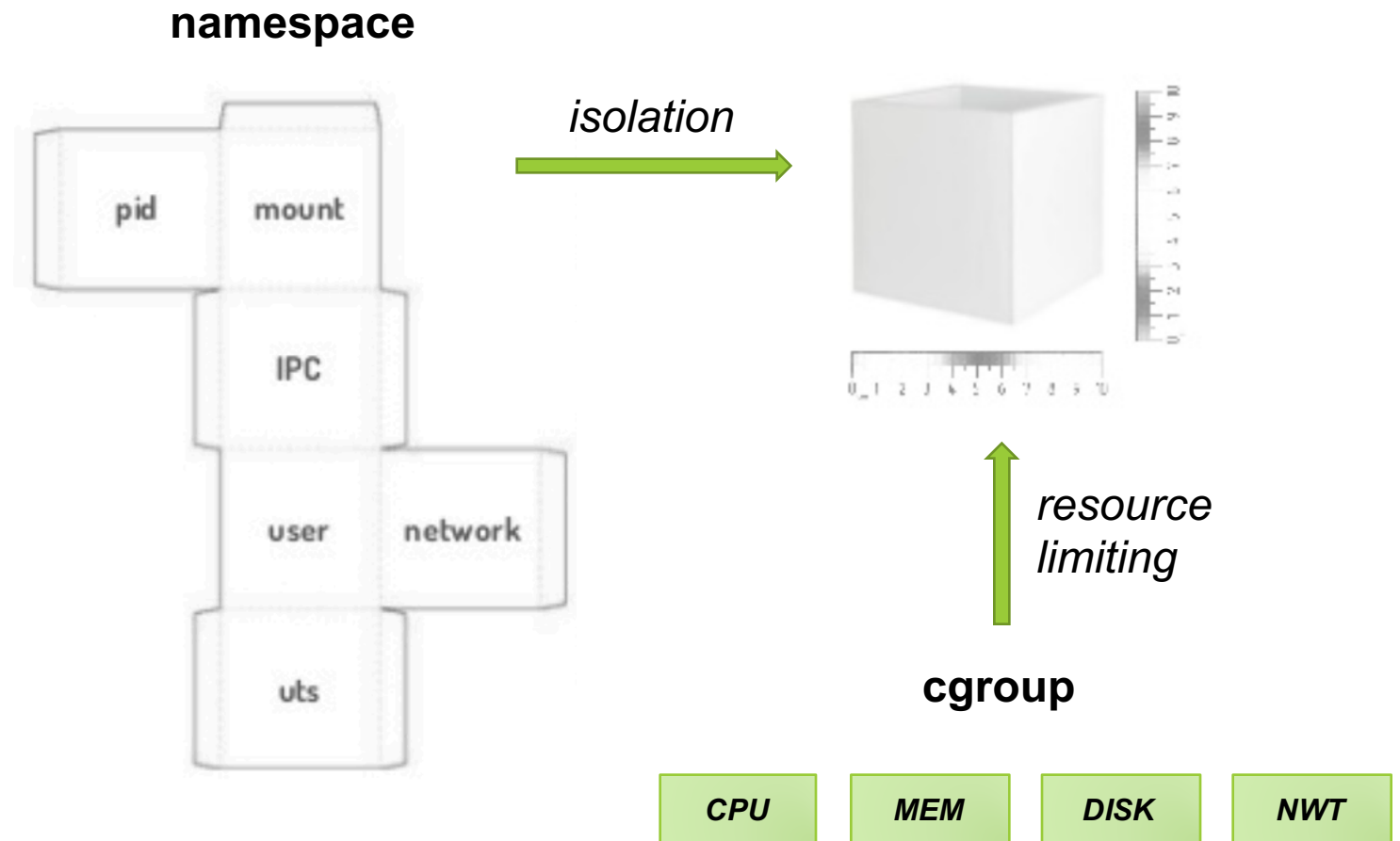- Easy to deploy



Docker Container



Virtual Machine

*Reference : https://www.docker.com/what-container*

**NAVER**

# What is Container?

- VM vs Container
- **OS Level Virtualization**
- Performance
- Run Everywhere – Portable
- Scalable – lightweight
- Environment Consistency
- Easy to manage Image
- Easy to deploy

**namespace**

*isolation*

*resource limiting*

**cgroup**

| CPU | MEM | DISK | NWT |

pid  mount

IPC

user  network

uts

*Reference : https://www.slideshare.net/PhilEstes/docker-london-container-security*
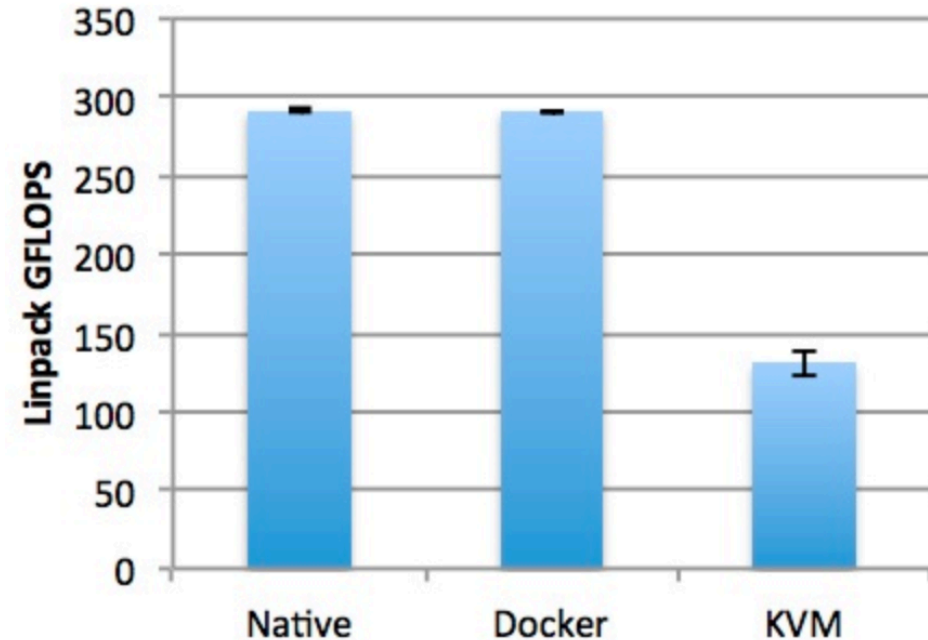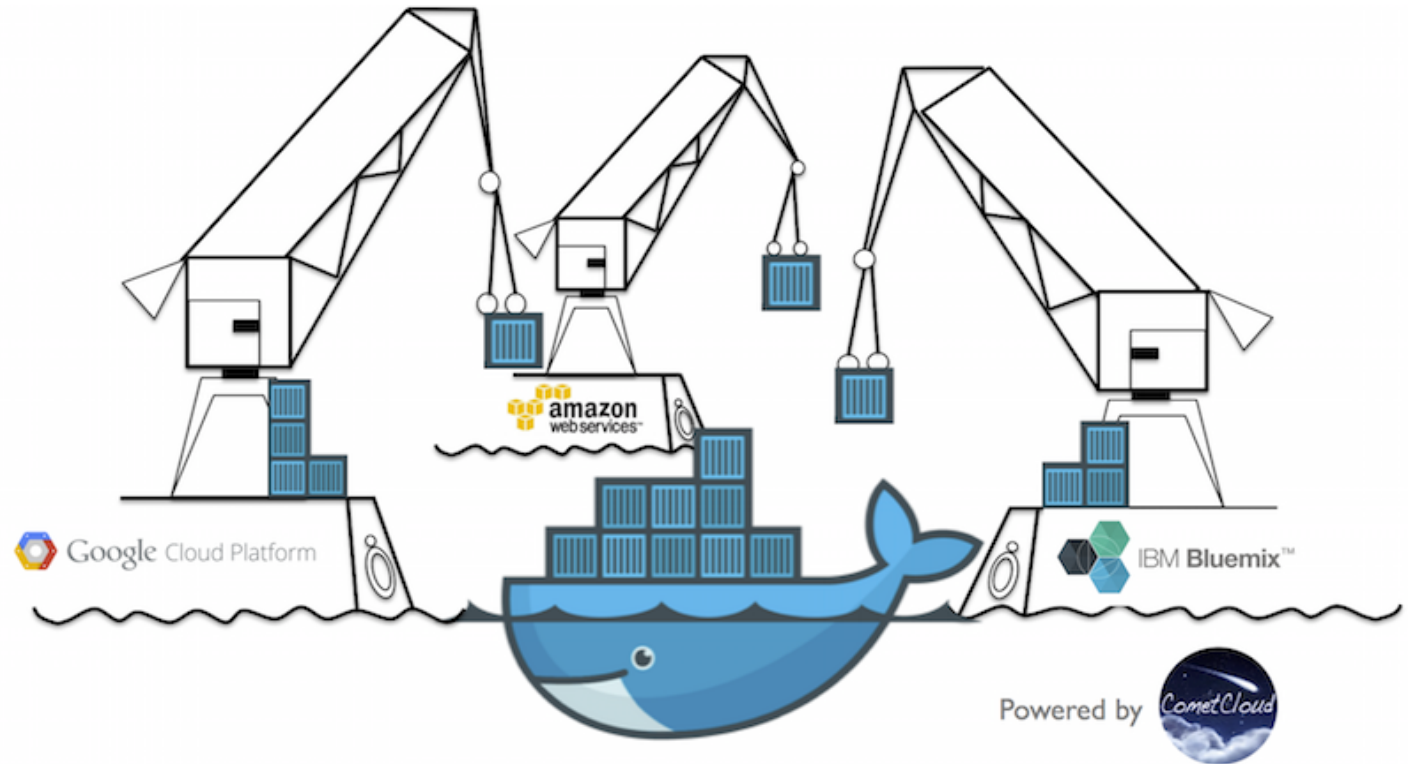
# What is Container?

- VM vs Container
- OS Level Virtualization
- **Performance**
- Run Everywhere – Portable
- Scalable – lightweight
- Environment Consistency
- Easy to manage Image
- Easy to deploy



*Reference : https://www.theregister.co.uk/2014/08/18/docker_kicks_kvms_butt_in_ibm_tests/*

NAVER

# What is Container?

- VM vs Container
- OS Level Virtualization
- Performance
- **Run Everywhere – Portable**
- Scalable – lightweight
- Environment Consistency
- Easy to manage Image
- Easy to deploy



**PM / VM / Cloud**

**NAVER**

# What is Container?

- VM vs Container
- OS Level Virtualization
- Performance
- Run Everywhere – Portable
- **Scalable – lightweight**
- Environment Consistency
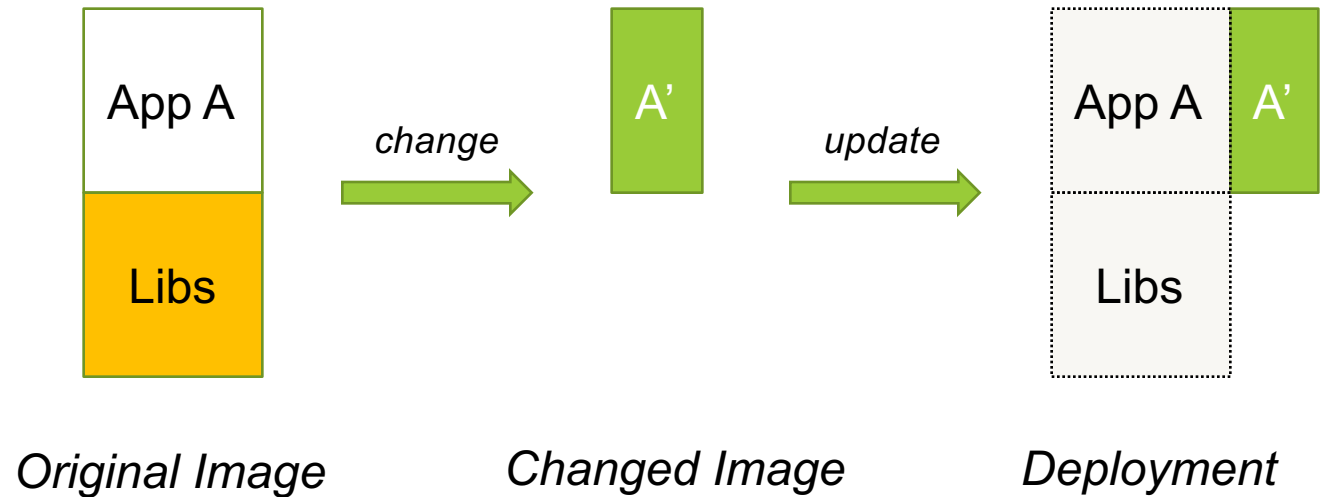- Easy to manage Image
- Easy to deploy

App A

Libs

*change*

A'

*update*

App A   A'

Libs

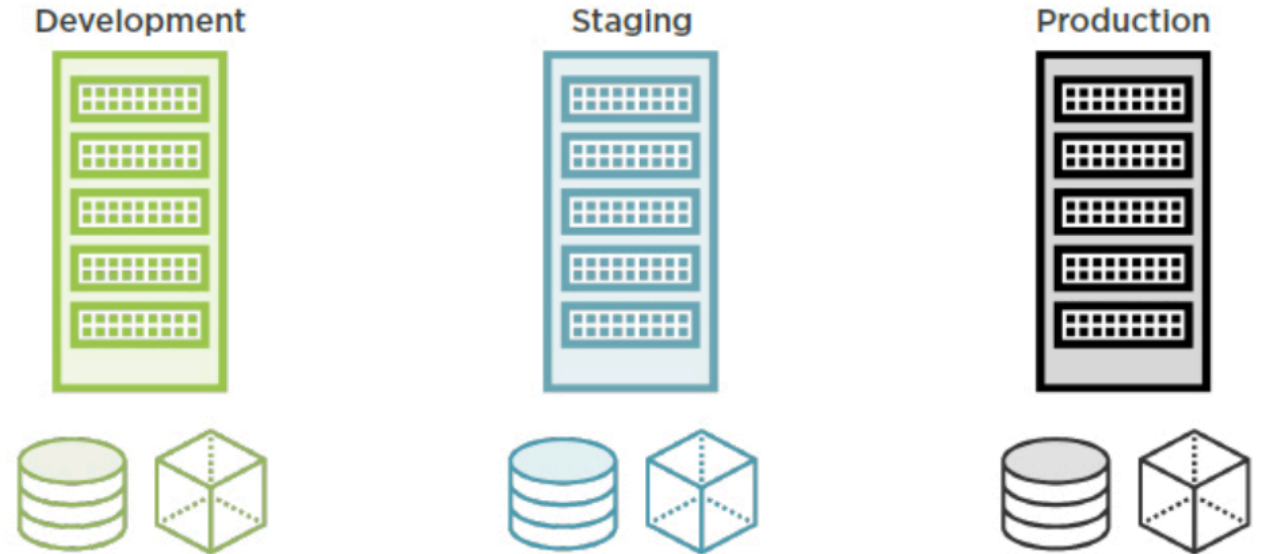*Original Image*          *Changed Image*          *Deployment*

**NAVER**

# What is Container?

- VM vs Container
- OS Level Virtualization
- Performance
- Run Everywhere – Portable
- Scalable – lightweight
- **Environment Consistency**
- Easy to manage Image
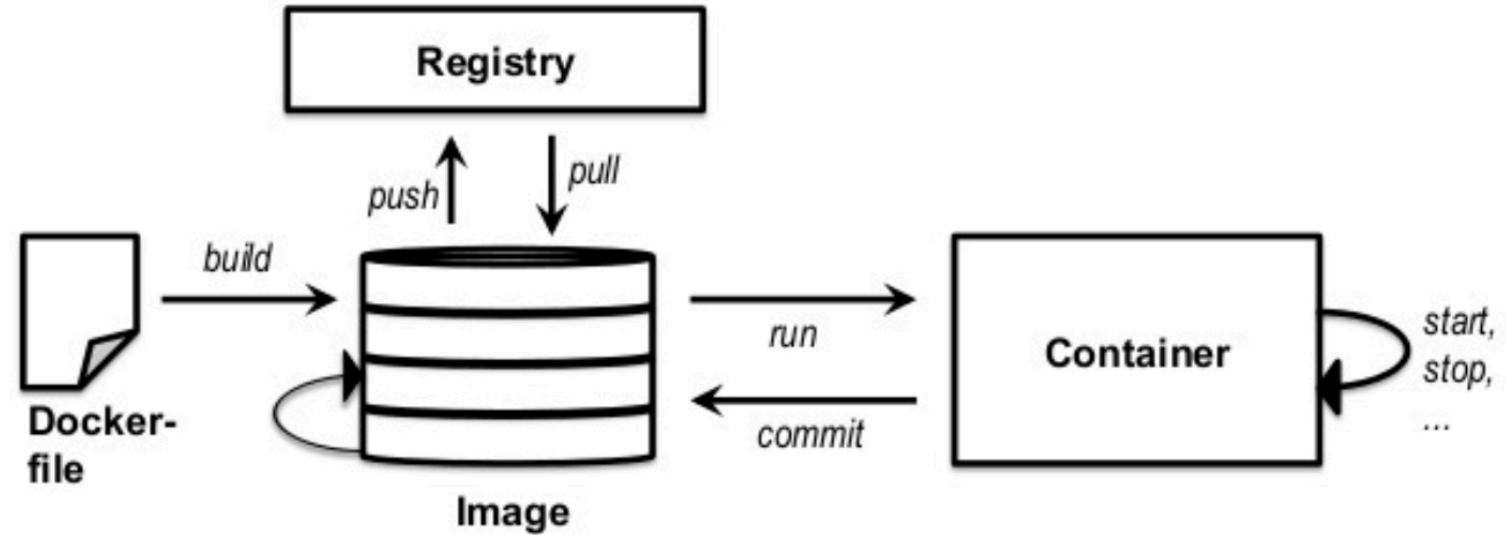- Easy to deploy

**NAVER**

# What is Container?

- VM vs Container
- OS Level Virtualization
- Performance
- Run Everywhere – Portable
- Scalable – lightweight
- Environment Consistency
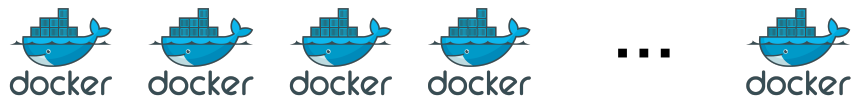- **Easy to manage Image**
- **Easy to deploy**

NAVER

# Persistent Storage

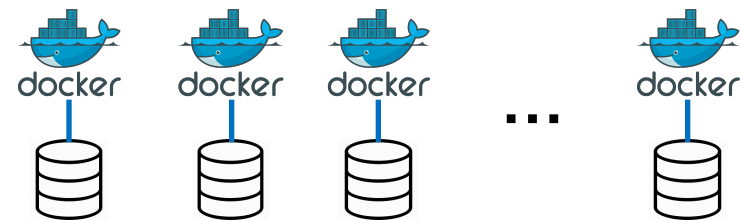# Stateless vs Stateful Container

## Stateless

- Nothing to write on disk
- Web (Front-end)
- Easy to scale in/out
- Container is ephemeral
- If delete, will be lost data



*Easy to scale out*

## Stateful

- Needs storage to write
- Database
- Logs
- CI config / repo data
- Secret Keys



*Hard to scale out*

# Ephemeral vs Persistent

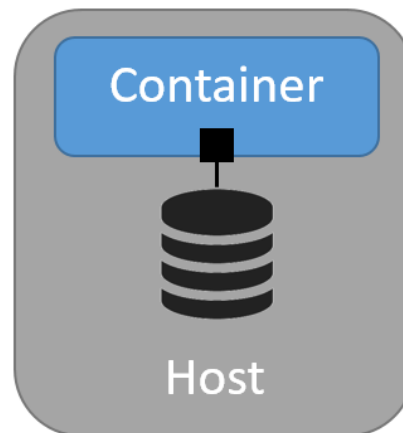| **Ephemeral Storage** | | **Persistent Storage** |
|---|---|---|
| • Data Lost<br>• Local Storage<br>• Stateless Container | **VS** | • Data Save<br>• Network Storage<br>• Stateful Container |



Data in the container
Lost when the container terminates

Data in a Host Volume
Lost when the host terminates

Networked Volume / File System
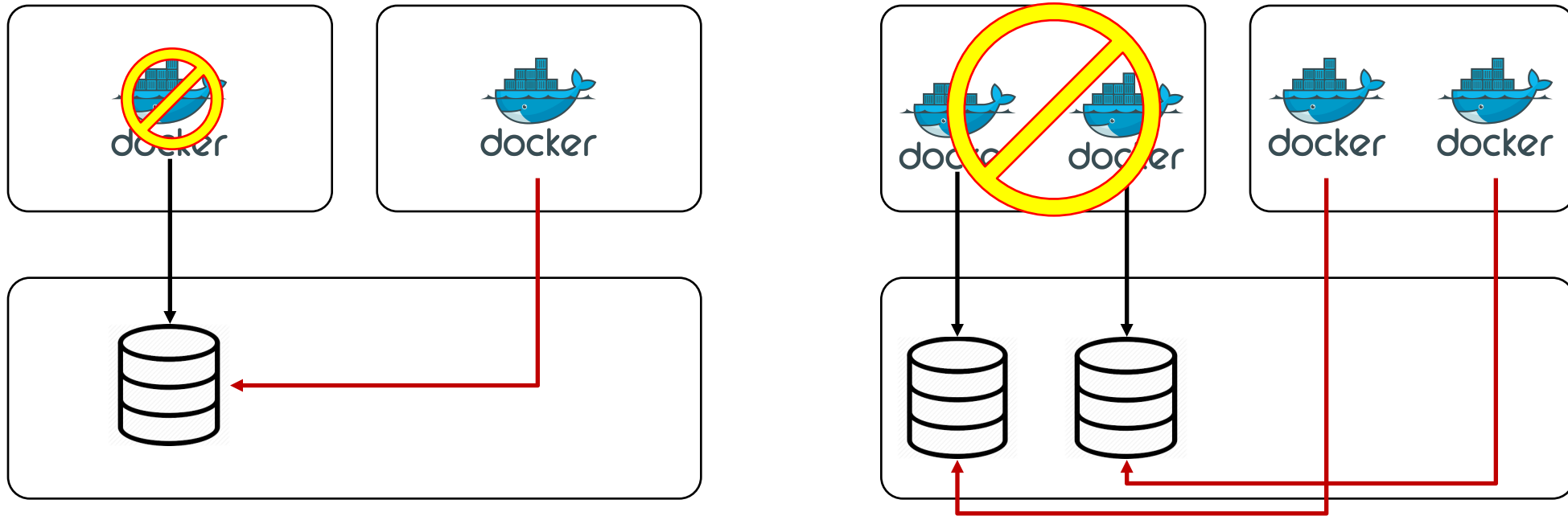Independent of host and container

NAVER

# Needs for Persistent Storage

Our mission is to provide **"Persistence Storage Service"**
while maintaining **"Agility" & "Automation"** of Docker Container

# Container in NAVER

# What is Ceph Storage?

## What is Ceph Storage?

- **Open Source Distributed Storage Solution**
- Massive Scalable / Efficient Scale-Out
- Unified Storage
- Runs on commodity hardware
- Integrations into Linux Kernel
     / QEMU/KVM Driver / OpenStack
- Self-managing / Self-healing
- Peer-to-Peer Storage Nodes
- RESTful API
- No metadata bottleneck (no lookup)
- CRUSH algorithm determines data placement
- Replicated / Erasure Coding
- Architecture

# Vender
# NO Lock-in

**NAVER**

# What is Ceph Storage?

- Open Source Distributed Storage Solution
- **Massive Scalable / Efficient Scale-Out**
- Unified Storage
- Runs on commodity hardware
- Integrations into Linux Kernel
      / QEMU/KVM Driver / OpenStack
- Self-managing / Self-healing
- Peer-to-Peer Storage Nodes
- RESTful API
- No metadata bottleneck (no lookup)
- CRUSH algorithm determines data placement
- Replicated / Erasure Coding
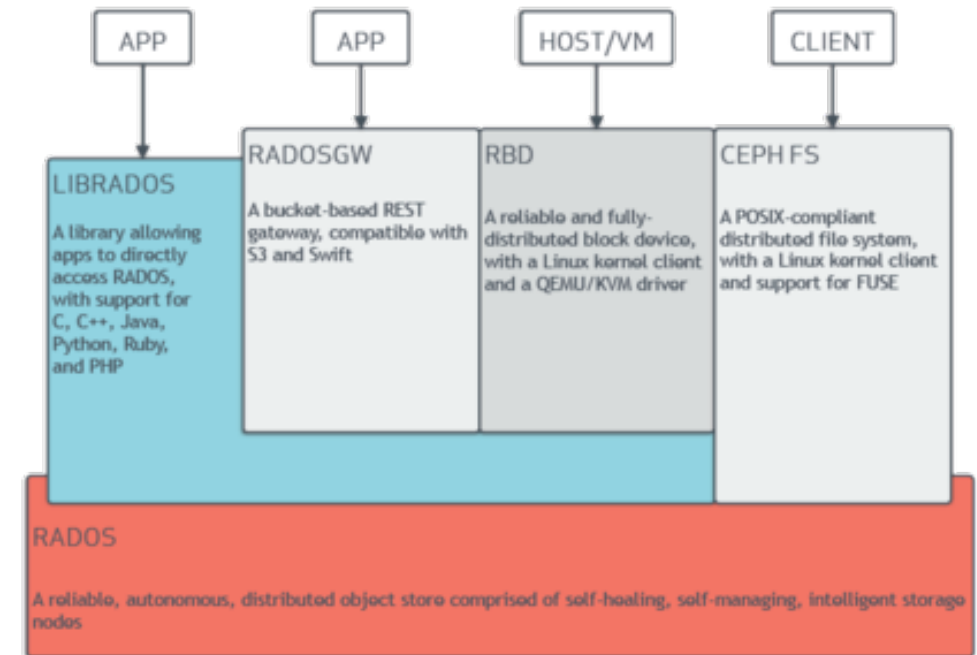- Architecture

# Up-to 16 Exa

**NAVER**

# What is Ceph Storage?

- Open Source Distributed Storage Solution
- Massive Scalable / Efficient Scale-Out
- **Unified Storage**
- Runs on commodity hardware
- Integrations into Linux Kernel
    / QEMU/KVM Driver / OpenStack
- Self-managing / Self-healing
- Peer-to-Peer Storage Nodes
- RESTful API
- No metadata bottleneck (no lookup)
- CRUSH algorithm determines data placement
- Replicated / Erasure Coding
- Architecture



**APP** → **LIBRADOS**
A library allowing apps to directly access RADOS, with support for C, C++, Java, Python, Ruby, and PHP

**APP** → **RADOSGW**
A bucket-based REST gateway, compatible with S3 and Swift

**HOST/VM** → **RBD**
A reliable and fully-distributed block device, with a Linux kernel client and a QEMU/KVM driver

**CLIENT** → **CEPH FS**
A POSIX-compliant distributed file system, with a Linux kernel client and support for FUSE

**RADOS**
A reliable, autonomous, distributed object store comprised of self-healing, self-managing, intelligent storage nodes

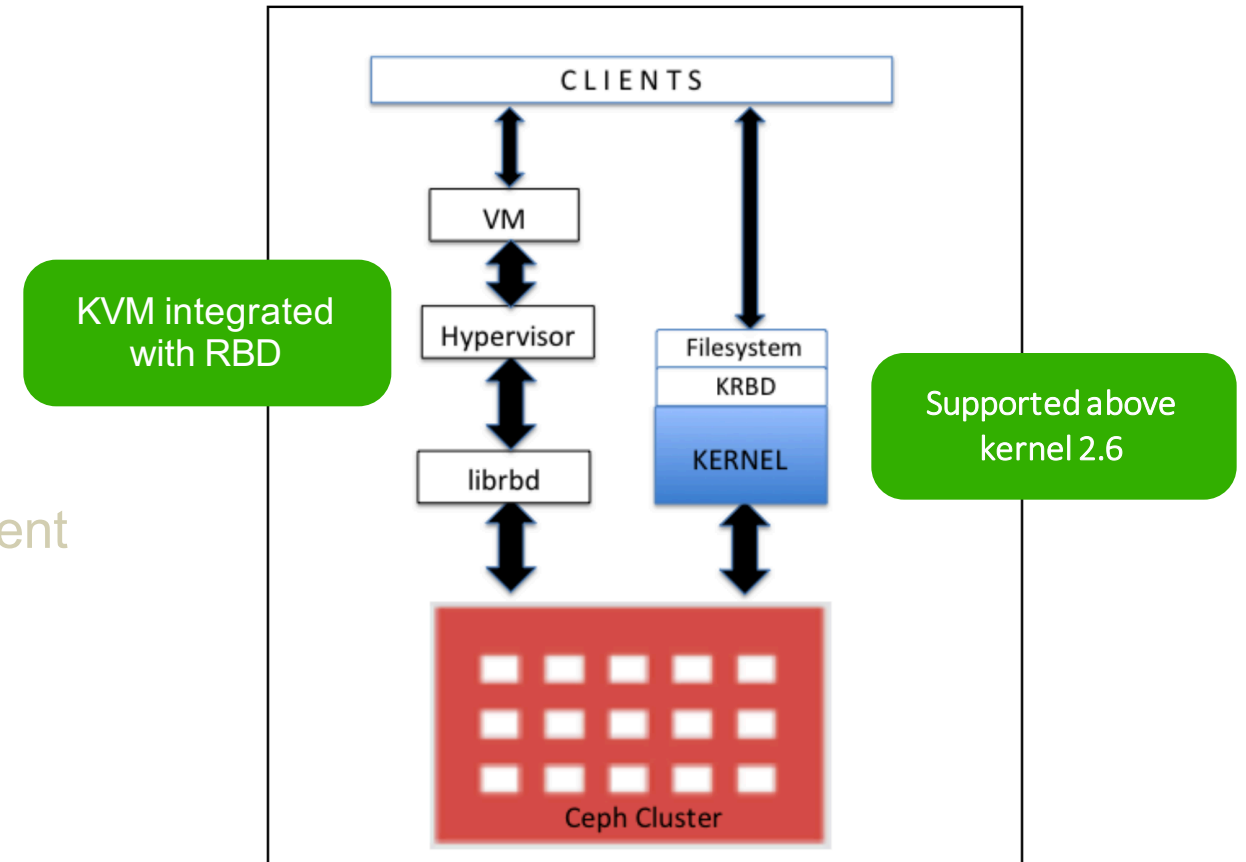Reference : https://en.wikipedia.org/wiki/Ceph_(software)

NAVER

# What is Ceph Storage?

- Open Source Distributed Storage Solution
- Massive Scalable / Efficient Scale-Out
- Unified Storage
- **Runs on commodity hardware**
- Integrations into Linux Kernel
        / QEMU/KVM Driver / OpenStack
- Self-managing / Self-healing
- Peer-to-Peer Storage Nodes
- RESTful API
- No metadata bottleneck (no lookup)
- CRUSH algorithm determines data placement
- Replicated / Erasure Coding
- Architecture

# Low TCO

**NAVER**

# What is Ceph Storage?

- Open Source Distributed Storage Solution
- Massive Scalable / Efficient Scale-Out
- Unified Storage
- Runs on commodity hardware
- **Integrations into Linux Kernel
  / QEMU/KVM Driver / OpenStack**
- Self-managing / Self-healing
- Peer-to-Peer Storage Nodes
- RESTful API
- No metadata bottleneck (no lookup)
- CRUSH algorithm determines data placement
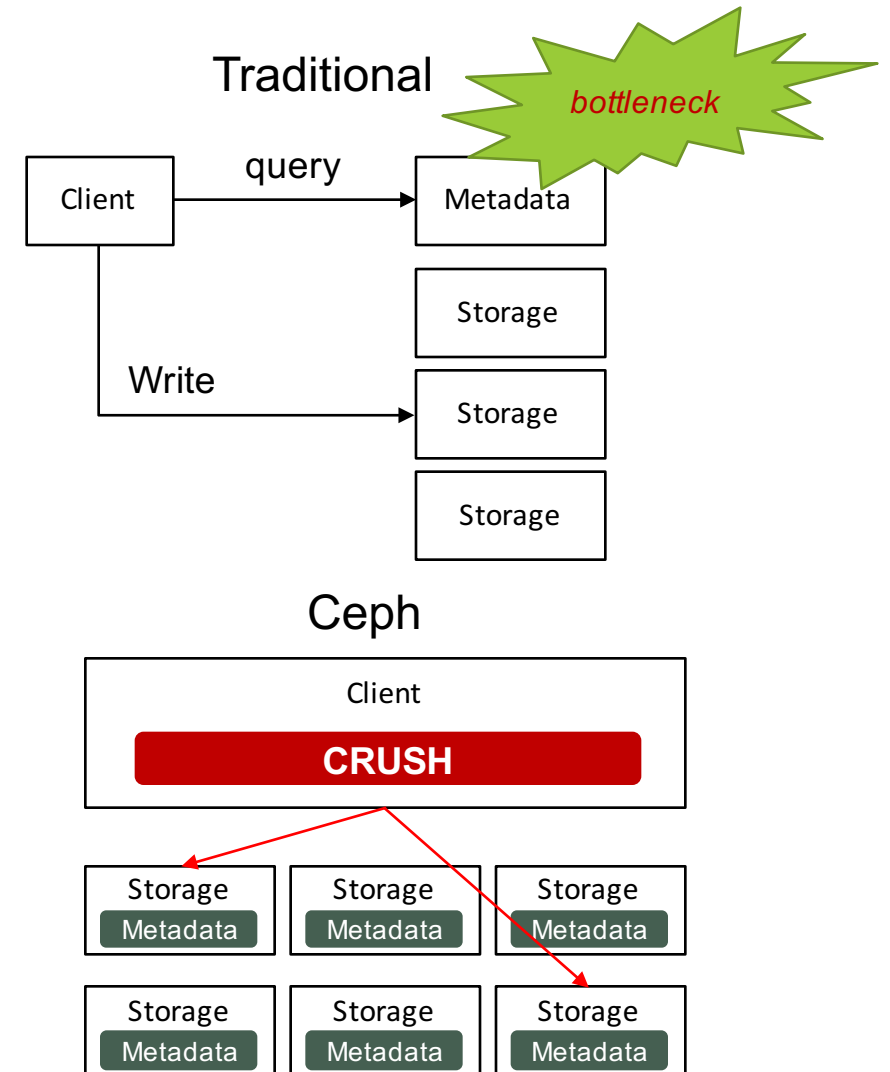- Replicated / Erasure Coding
- Architecture

# What is Ceph Storage?

- Open Source Distributed Storage Solution
- Massive Scalable / Efficient Scale-Out
- Unified Storage
- Runs on commodity hardware
- Integrations into Linux Kernel
    / QEMU/KVM Driver / OpenStack
- **Self-managing / Self-healing**
- **Peer-to-Peer Storage Nodes**
- **RESTful API**
- No metadata bottleneck (no lookup)
- CRUSH algorithm determines data placement
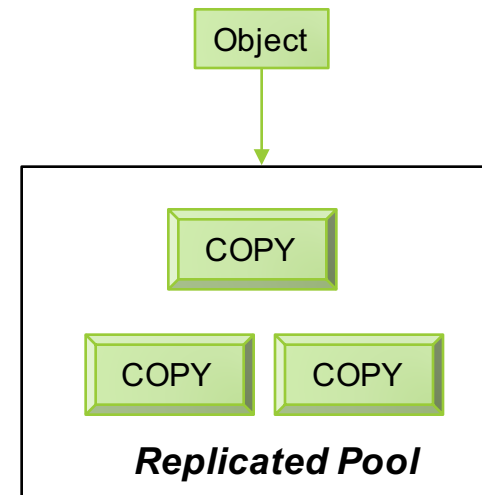- Replicated / Erasure Coding
- Architecture

# What is Ceph Storage?

- Open Source Distributed Storage Solution
- Massive Scalable / Efficient Scale-Out
- Unified Storage
- Runs on commodity hardware
- Integrations into Linux Kernel
      / QEMU/KVM Driver / OpenStack
- Self-managing / Self-healing
- Peer-to-Peer Storage Nodes
- RESTful API
- **No metadata bottleneck (no lookup)**
- **CRUSH  algorithm determines data placement**
- Replicated / Erasure Coding
- Architecture

Traditional

bottleneck

Client — query → Metadata

Storage

Write → Storage

Storage

Ceph

Client

**CRUSH**

| Storage | Storage | Storage |
| Metadata | Metadata | Metadata |

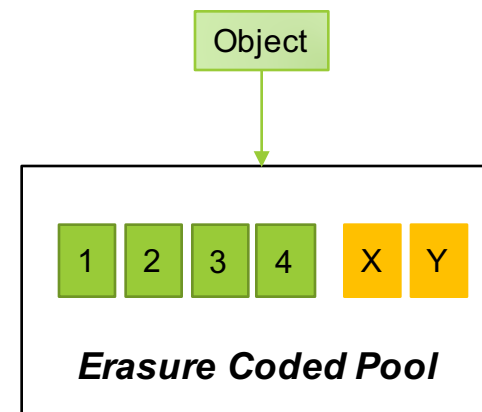| Storage | Storage | Storage |
| Metadata | Metadata | Metadata |

NAVER

# What is Ceph Storage?

- Open Source Distributed Storage Solution
- Massive Scalable / Efficient Scale-Out
- Unified Storage
- Runs on commodity hardware
- Integrations into Linux Kernel
      / QEMU/KVM Driver / OpenStack
- Self-managing / Self-healing
- Peer-to-Peer Storage Nodes
- RESTful API
- No metadata bottleneck (no lookup)
- CRUSH algorithm determines data placement
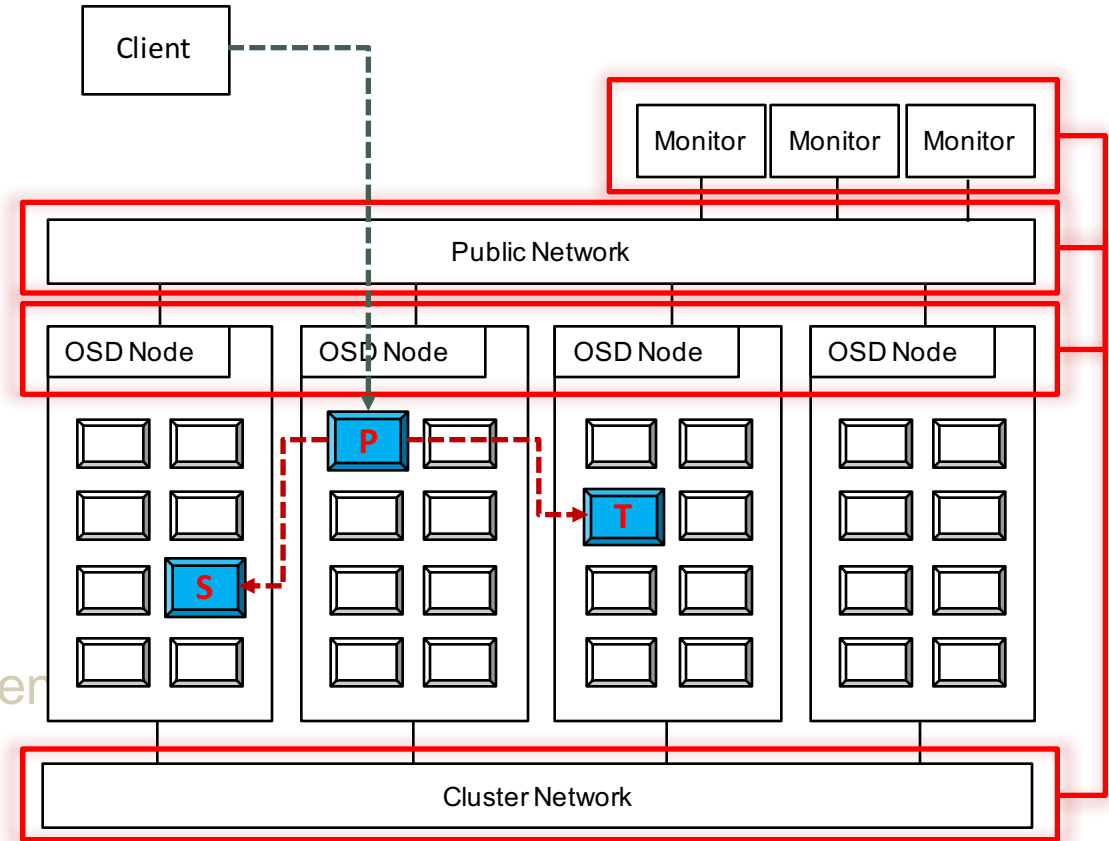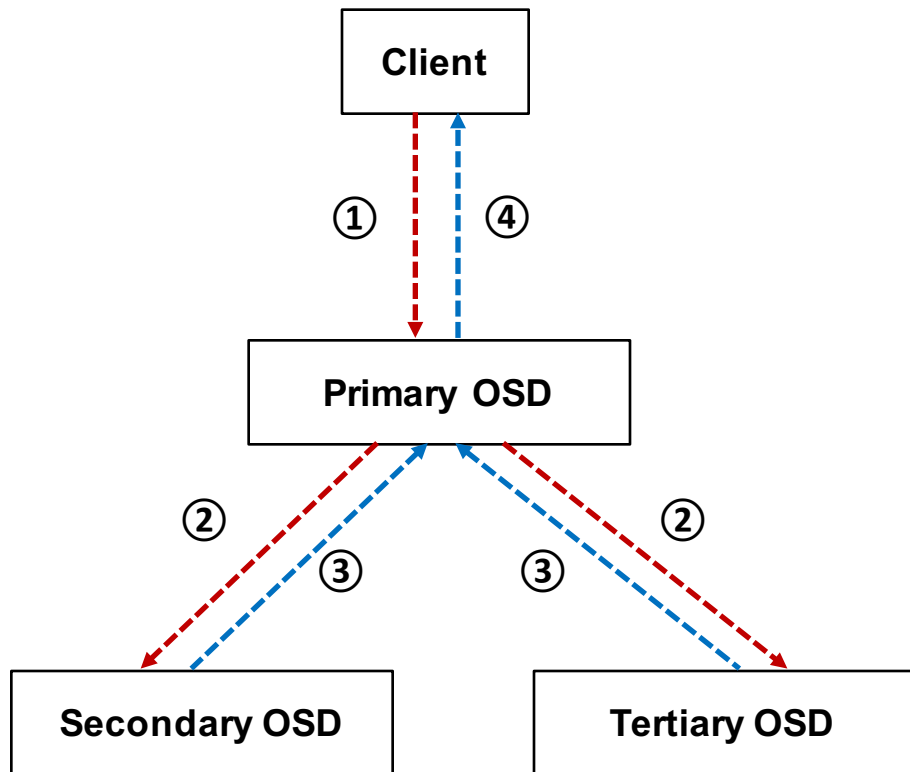- **Replicated / Erasure Coding**
- Architecture

Object

COPY

COPY    COPY

**Replicated Pool**

- Very high durability
- 200 % overhead
- Quick recovery

Object

| 1 | 2 | 3 | 4 | X | Y |

**Erasure Coded Pool**

- Cost-Effective
- 50 % overhead
- Expensive Recovery

NAVER

# What is Ceph Storage?

- Open Source Distributed Storage Solution
- Massive Scalable / Efficient Scale-Out
- Unified Storage
- Runs on commodity hardware
- Integrations into Linux Kernel
        / QEMU/KVM Driver / OpenStack
- Self-managing / Self-healing
- Peer-to-Peer Storage Nodes
- RESTful API
- No metadata bottleneck (no lookup)
- CRUSH algorithm determines data placement
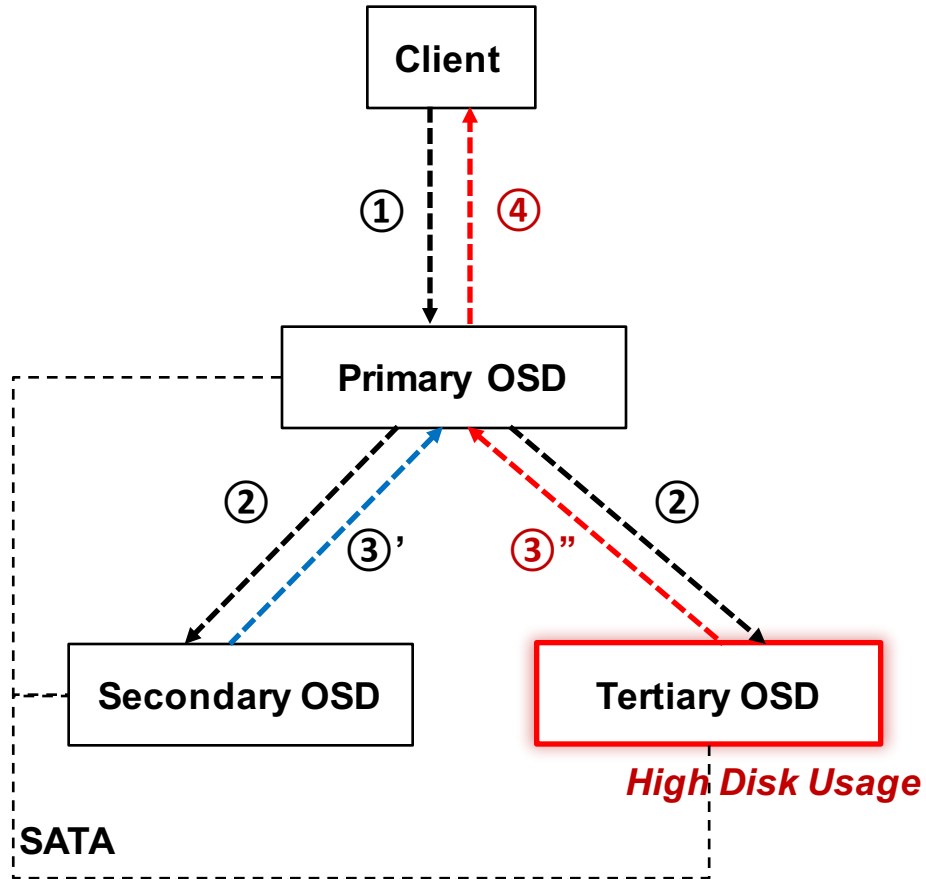- Replicated / Erasure Coding
- **Architecture**

Client

Monitor  Monitor  Monitor

Public Network

OSD Node  OSD Node  OSD Node  OSD Node

P

T

S

Cluster Network

# Write Process in Ceph

# Write Flow



- **Strong Consistency**
- **CAP Theorem : CP system**
  **(Consistency / network Partition)**

① **Client writes data to primary osd.**

② **Primary OSD sends data to replica OSDs, write data to local disk.**

③ **Replica OSDs write data to local disk, signal completion to primary.**

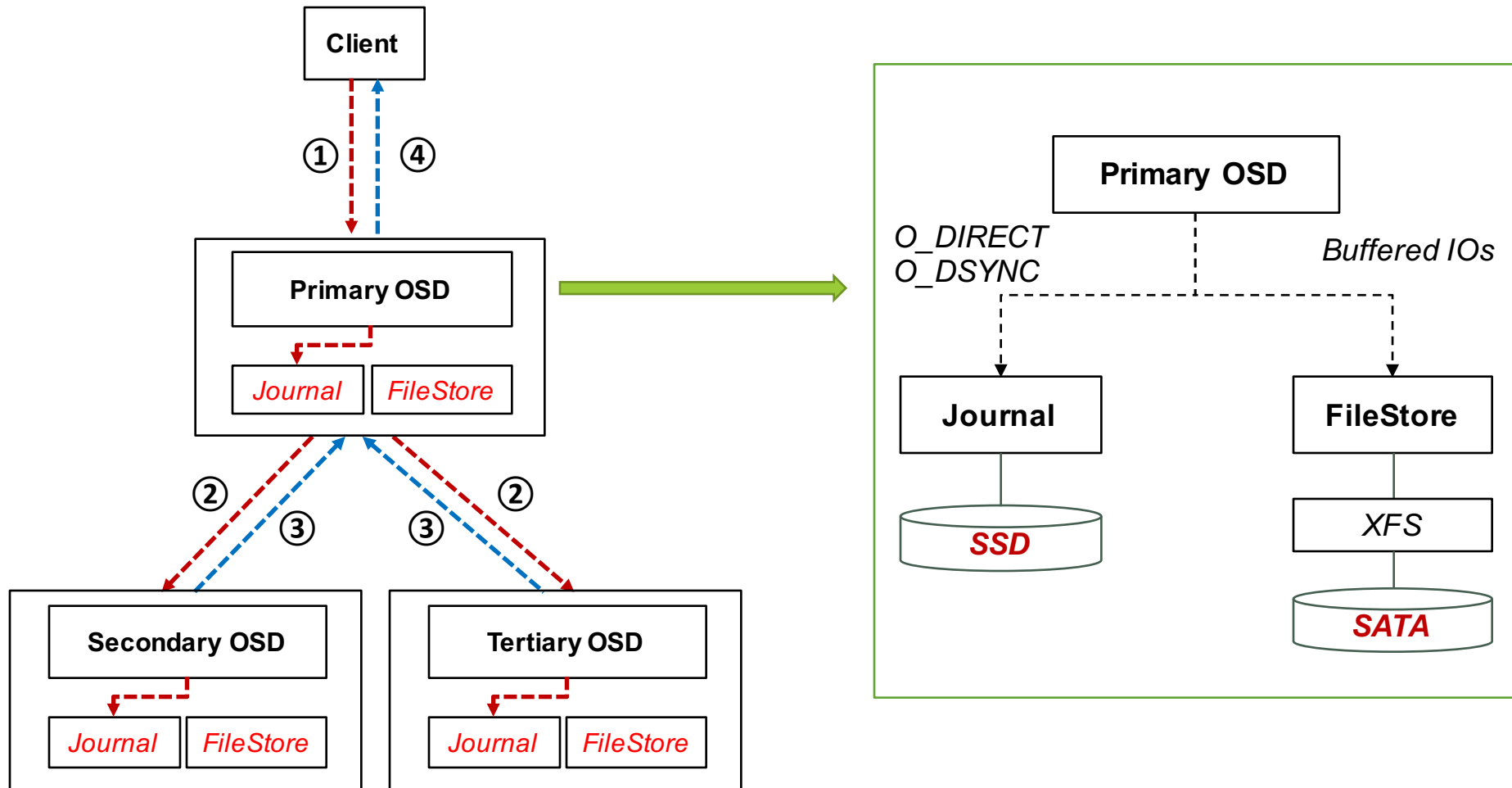④ **Primary OSD signals completion to client.**

# Slow Requests



- ■ **All of Data Disks are SATA.**
- ■ **7.2k rpm SATA : ~ 75 ~ 100 IOPS**

① **Client writes data to primary osd.**

② **Primary OSD sends data to replica OSDs, write data to local disk.**

③' **Write data to local disk, Send ack to primary.**

③" **Slow write data to local disk, Send ack to primary.**

④ **Slow send ack to client.**

# Journal on SSD

# Performance Test

# Test Environment



**Clients**

**Network Switch (10G)**

*Public Networks*          *Cluster Networks*

**Ceph**

**KRBD**
- /dev/rbd0
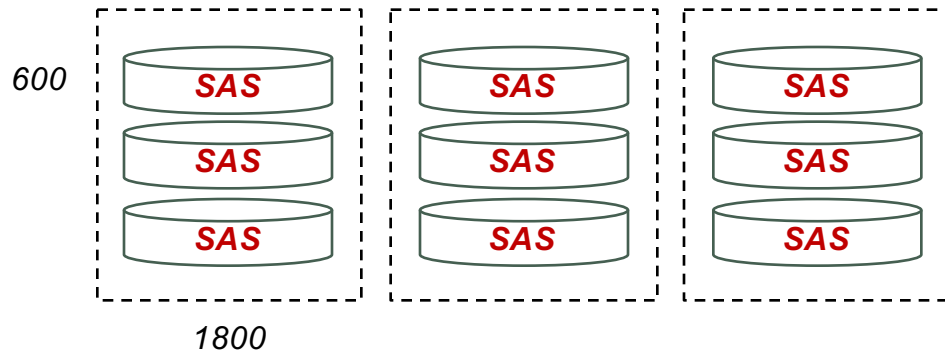- mkfs / mount
- fio

**Ceph**
- FileStore
- Luminous (12.2.0)

**Server**
- Intel® Xeon CPU L5640 2.27 GHz
- 16 GB Memory
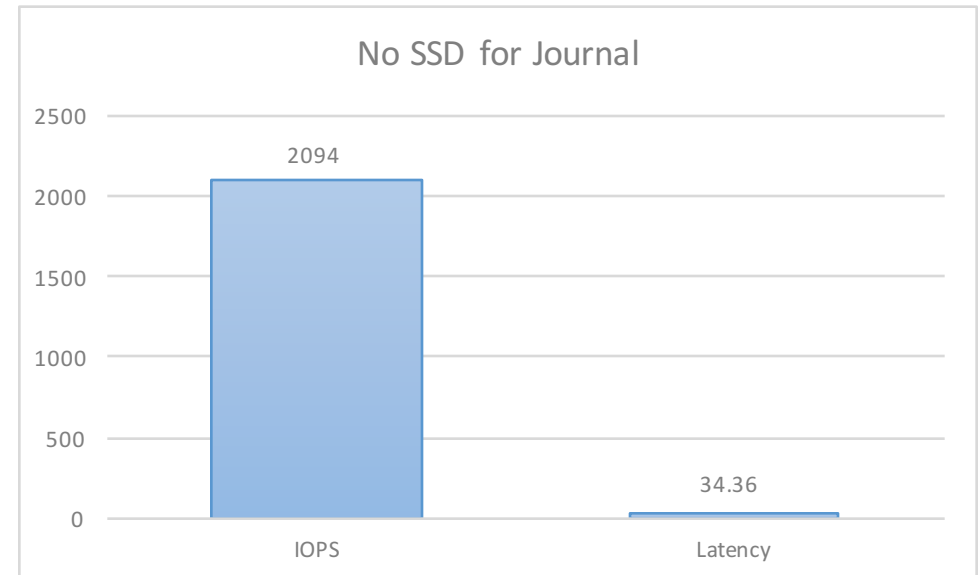- 480GB SAS 10K x 5
- 480GB SAS SSD x 1
- Centos 7

NAVER

# Case #1. Only Use SAS DISK

## Expect



*600*

*1800*

SAS : 10,000 rpm / 600 IOPS
Per Node : SAS * 3 = 600 * 3 = 1,800 IOPS
Total : Node * 3 = 1,800 * 3 = 5,400 IOPS
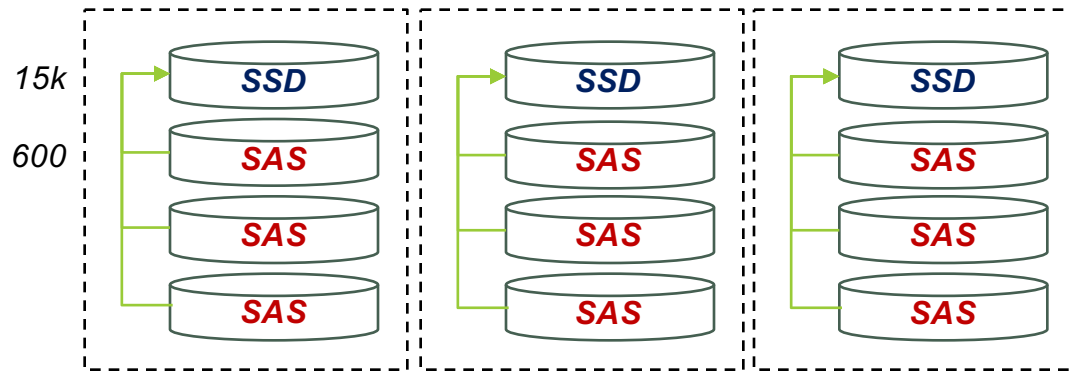Replicas 3 = 5,400 / 3 = **1,800 IOPS**
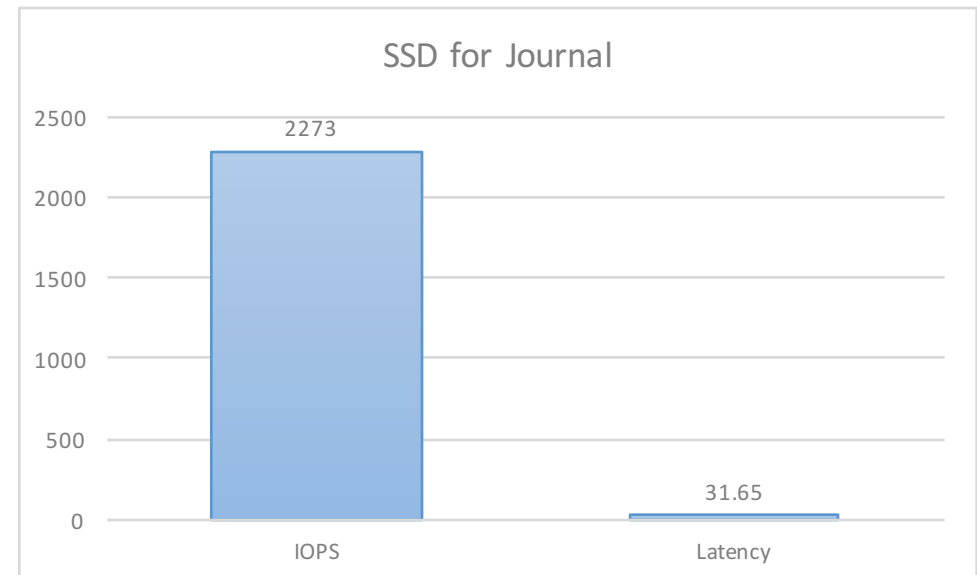
## Result : 4K Rand Write



No SSD for Journal

| | IOPS | Latency |

2094

34.36

**2094 IOPS/s**

NAVER

# Case #2. Use SSD for Journal

**Expect**



15k
600

SSD : 15,000 IOPS
Total : Node * 3 = 15,000* 3 = 45,000 IOPS
Replicas 3 = 45,000 / 3 = **15,000 IOPS**
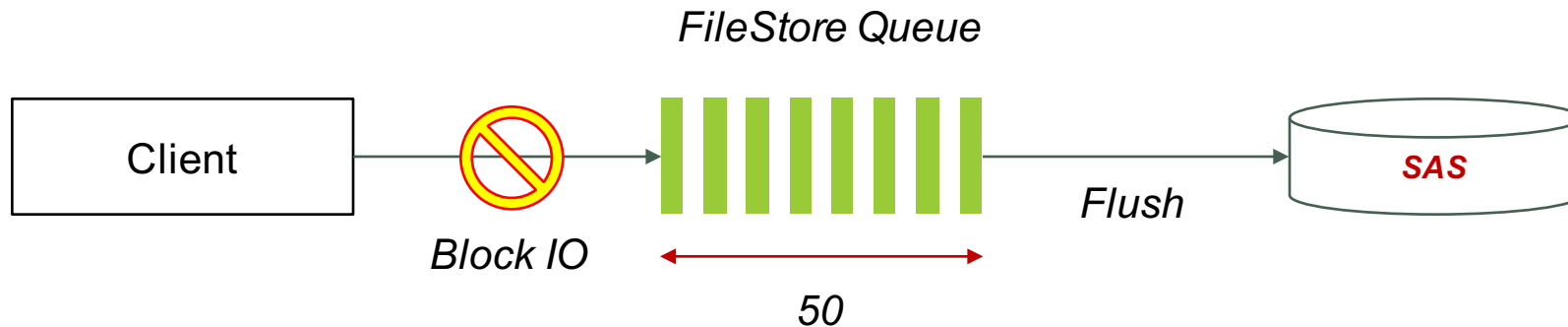
**Result : 4K Rand Write**



**2094 → 2273**

# Analysis

# ceph --admin-daemon /var/run/ceph/ceph-osd.2.asok  perf dump | grep -A 15 throttle-filestore_ops
"throttle-filestore_ops": {
"val": 50,          ← limitation
"max": 50,
…

# ceph --admin-daemon /var/run/ceph/ceph-osd.2.asok  config show | grep queue_max
…
"filestore_queue_max_bytes": "104857600",
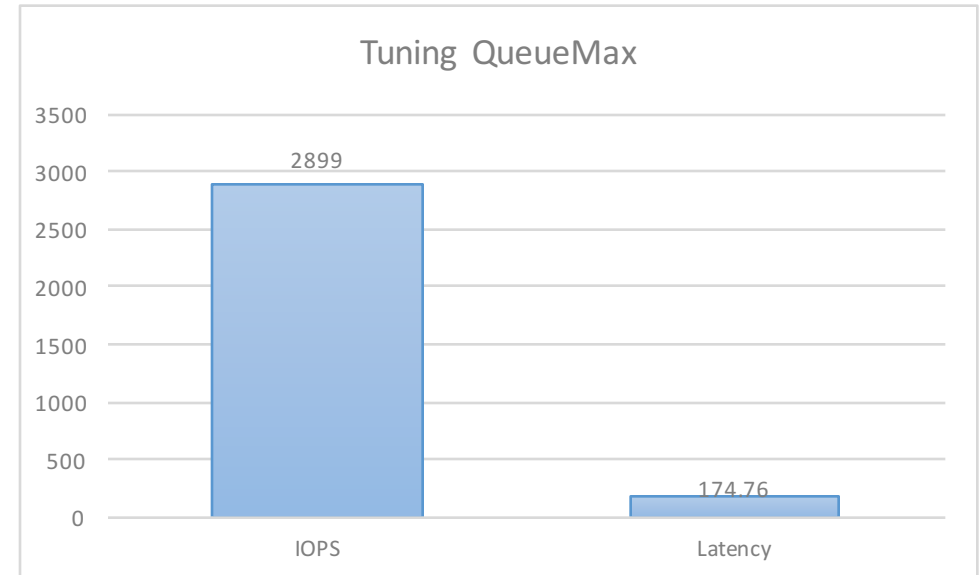"filestore_queue_max_ops": "50",     ← default value
…

FileStore Queue

# Case #3. Tuning FileStore

## Tuning

*filestore_queue_max_ops = 500000*
*filestore_queue_max_bytes = 42949672960*
*journal_max_write_bytes = 42949672960*
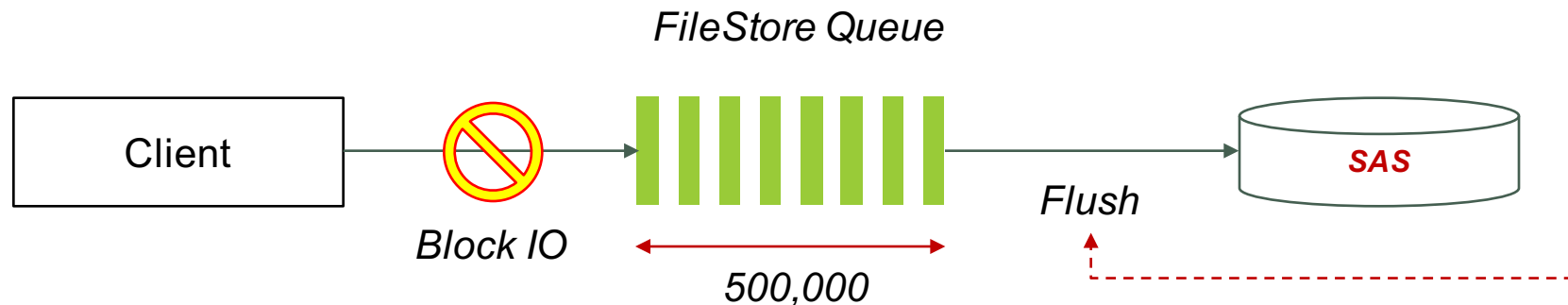*journal_max_write_entries = 5000000*

## Result



**2094 → 2899 (38%)**

# Analysis

# ceph --admin-daemon /var/run/ceph/ceph-osd.2.asok perf dump | grep -A6 WBThrottle
"WBThrottle": {
"bytes_dirtied": 21049344,
"bytes_wb": 197390336,
"ios_dirtied": 5017,     ← limitation
"ios_wb": 40920,
"inodes_dirtied": 1195,
"inodes_wb": 20602

# ceph --admin-daemon /var/run/ceph/ceph-osd.2.asok config show | grep wbthrottle_xfs
"filestore_wbthrottle_xfs_bytes_hard_limit": "419430400",
"filestore_wbthrottle_xfs_bytes_start_flusher": "41943040",
"filestore_wbthrottle_xfs_inodes_hard_limit": "5000",
"filestore_wbthrottle_xfs_inodes_start_flusher": "500",
"filestore_wbthrottle_xfs_ios_hard_limit": "5000",
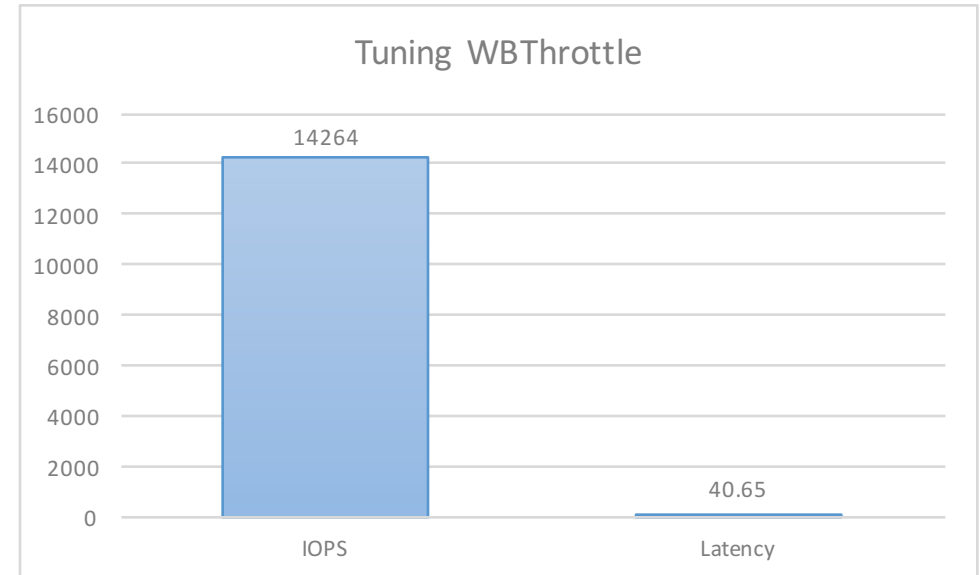"filestore_wbthrottle_xfs_ios_start_flusher": "500",

*FileStore Queue*

Client

*Block IO*

*500,000*

*Flush*

SAS

# Case #4. Tuning WBThrottle

**Tuning**

**Result**

"filestore_wbthrottle_enable": "false",
or
"filestore_wbthrottle_xfs_bytes_hard_limit": "4194304000",
"filestore_wbthrottle_xfs_bytes_start_flusher": "419430400",
"filestore_wbthrottle_xfs_inodes_hard_limit": "500000",
"filestore_wbthrottle_xfs_inodes_start_flusher": "5000",
"filestore_wbthrottle_xfs_ios_hard_limit": "500000",
"filestore_wbthrottle_xfs_ios_start_flusher": "5000",



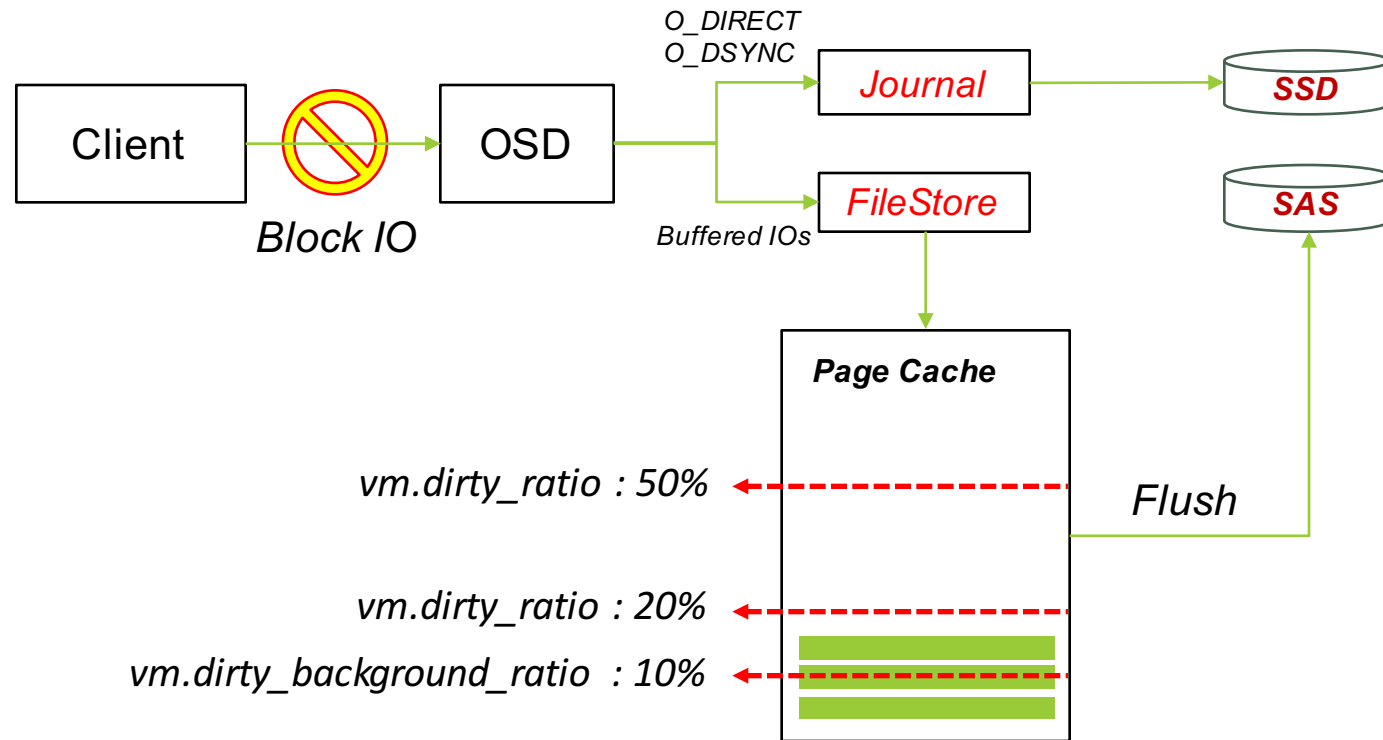**2,094 → 14,264 (x7)**

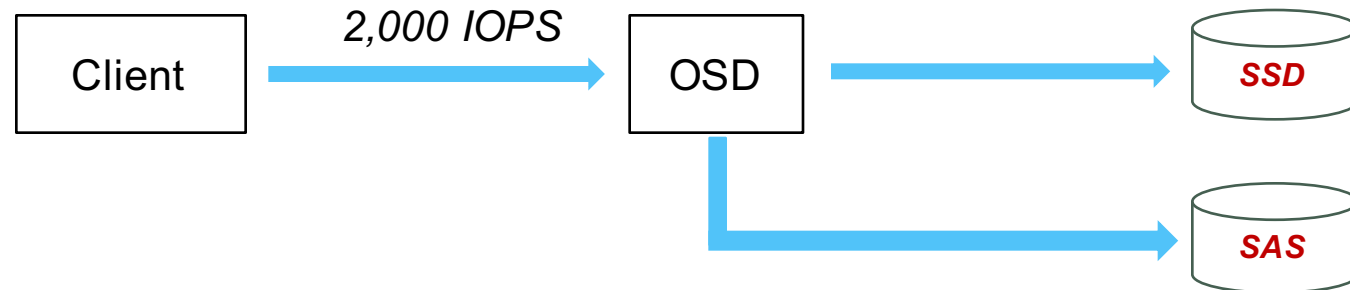# Analysis

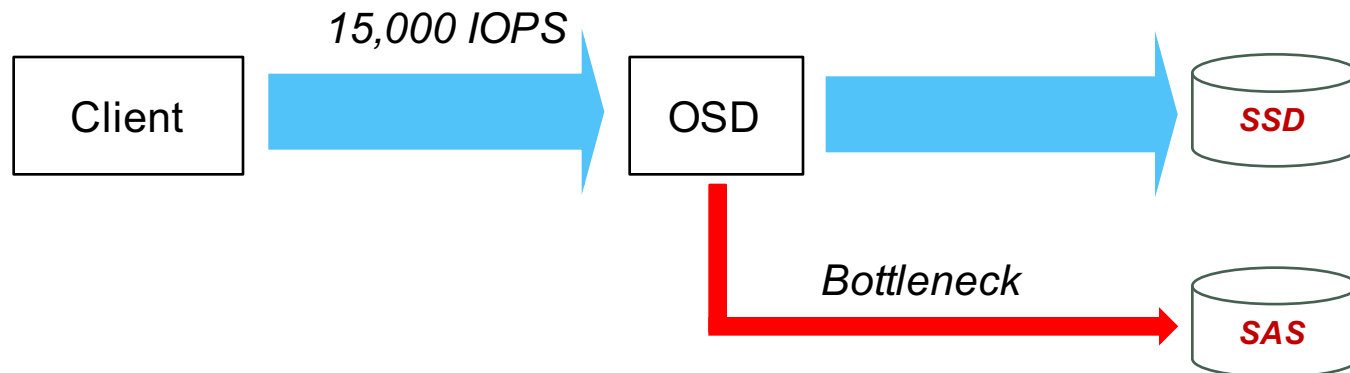*After 35 secs, Performance (IPOS) drops below 100 ~ 200 IOPS...*

# Problems of Throttle

# Dynamic Throttle

**Dynamic Throttle**

Client — 500 IOPS → OSD → SSD, SAS

- *Burst (~ 60 secs)*
- *Throttle From 80%*

*filestore_queue_max_ops*
*filestore_queue_max_bytes*

*filestore_expected_throughput_ops*
*filestore_expected_throughput_bytes*

*filestore_queue_low_threshhold*
*filestore_queue_high_threshhold*
*filestore_queue_high_delay_multiple*
*filestore_queue_max_delay_multiple*

```
r = current_op / max_ops
high_delay_per_count = high_multiple / expected_throughput_ops
max_delay_per_count = max_multiple / expected_throughput_ops
s0 = high_delay_per_count / (high_threshhold - low_threshhold)
s1 = (max_delay_per_count - high_delay_per_count) / (1 -
high_threshhold)
 if r < low_threshhold:
    delay = 0
elif r < high_threshhold:
    delay = (r - low_threshhold) * s0
else:
    delay = high_delay_per_count + ((r - high_threshhold) * s1)
```

Reference : http://blog.wjin.org/posts/ceph-dynamic-throttle.html

NAVER

# Conculsion

- ■ High performance improvements with SSD : **2,094 → 14,264 (x7)**

- ■ Must need Throttling for stable storage operation : Dynamic Throttle

- ■ Need to tune OS(page cache, io scheduler), Ceph config

**NAVER**

# QnA