# A New Partitioning Scheme for CephFS

**LINE**

**오용석**

**Ceph Days Korea 2023**

# Outline of Contents

- Introduction and background

- CephFS subtree partitioning

- Combining static and dynamic partitioning schemes
  - Workload based static partitioner with bal_rank_mask
  - A new CephFS MDS balancer with ceph.dir.bal.mask
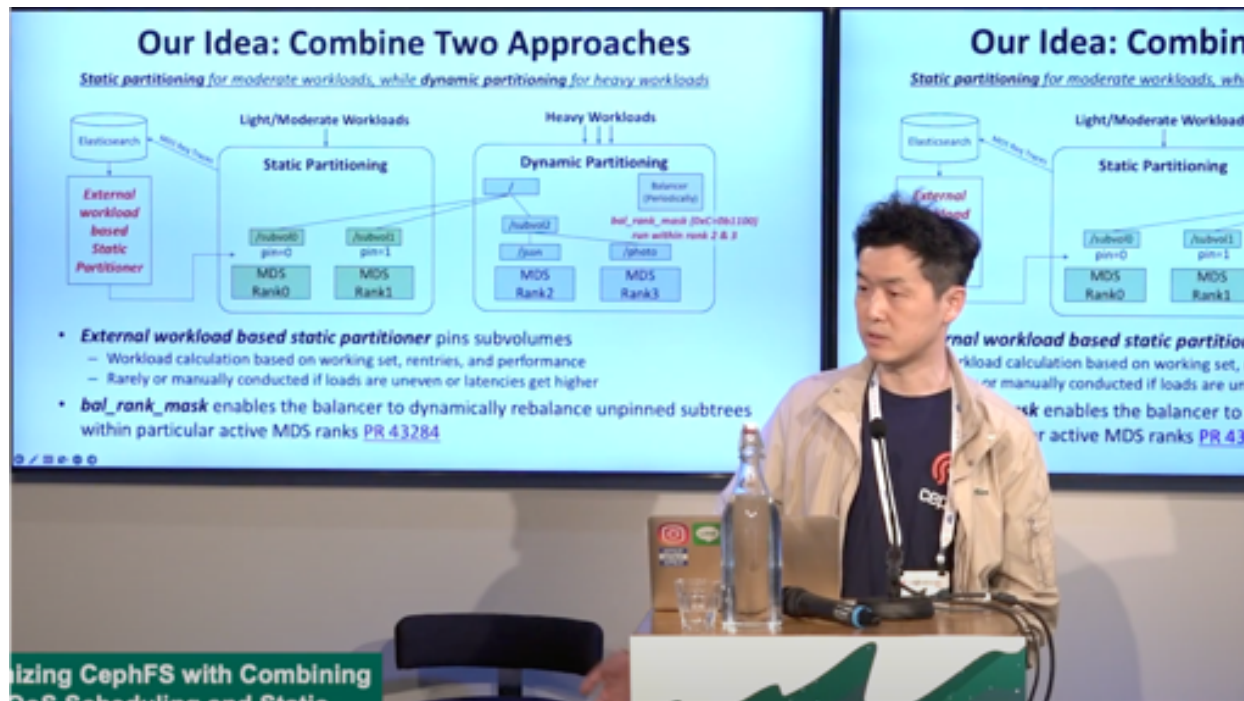
- Future works and conclusions

# Who Am I?

- Yongseok Oh (오용석)
- 2015 Ph.D. in Computer Science
  - Flash based storage systems
  - Research papers in FAST, MSST, Systor (over 400 citations)
- 2015 SKT - SW-defined Storage Lab.
  - Enterprise NVMe SSD FW
  - SPDK based file system
- 2017 SK hynix - FW Group
  - Enterprise NVMe SSD FW
  - ZNS, global wear leveling, SR-IOV, NVMe-MI
- 2020 ~ LINE Plus - Cloud Storage
  - CephFS based shared file service
  - NVMeoF based storage service
- Hobby: camping, boxing, jogging
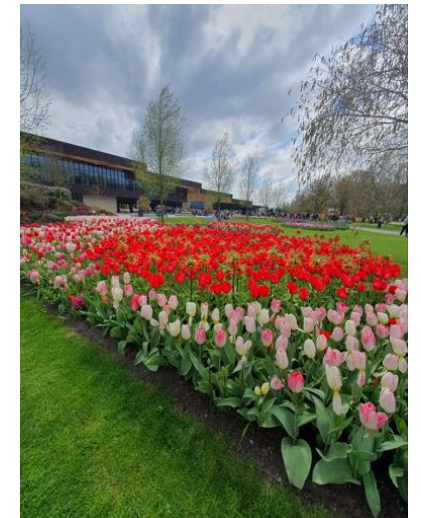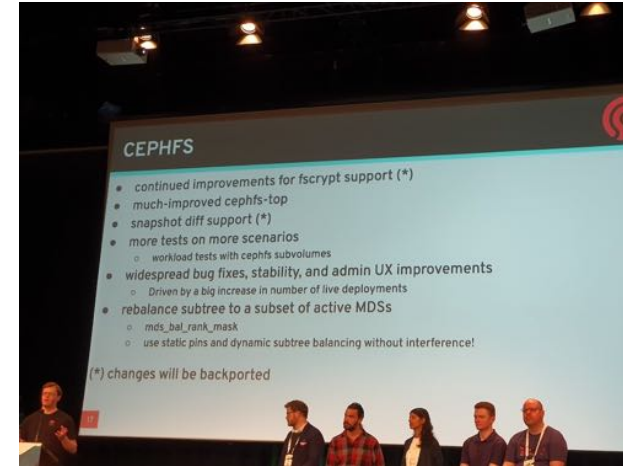
# Cephalocon2023

- **Optimizing CephFS with Combining MDS QoS Scheduling and Static-Dynamic Subtree Partitioning**
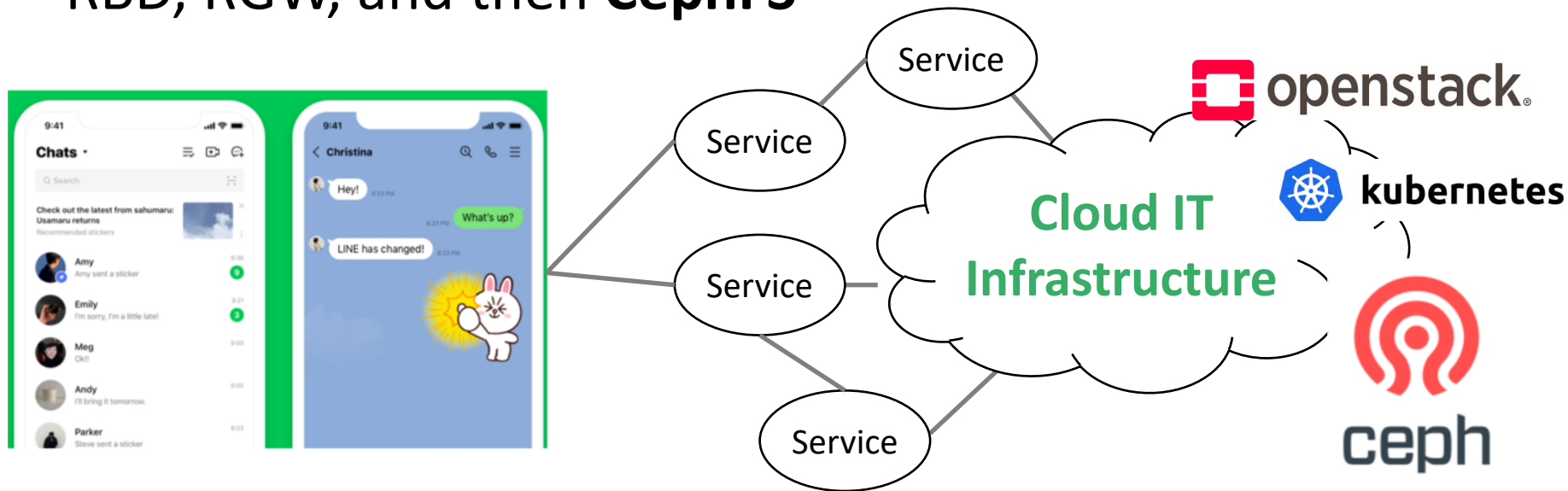


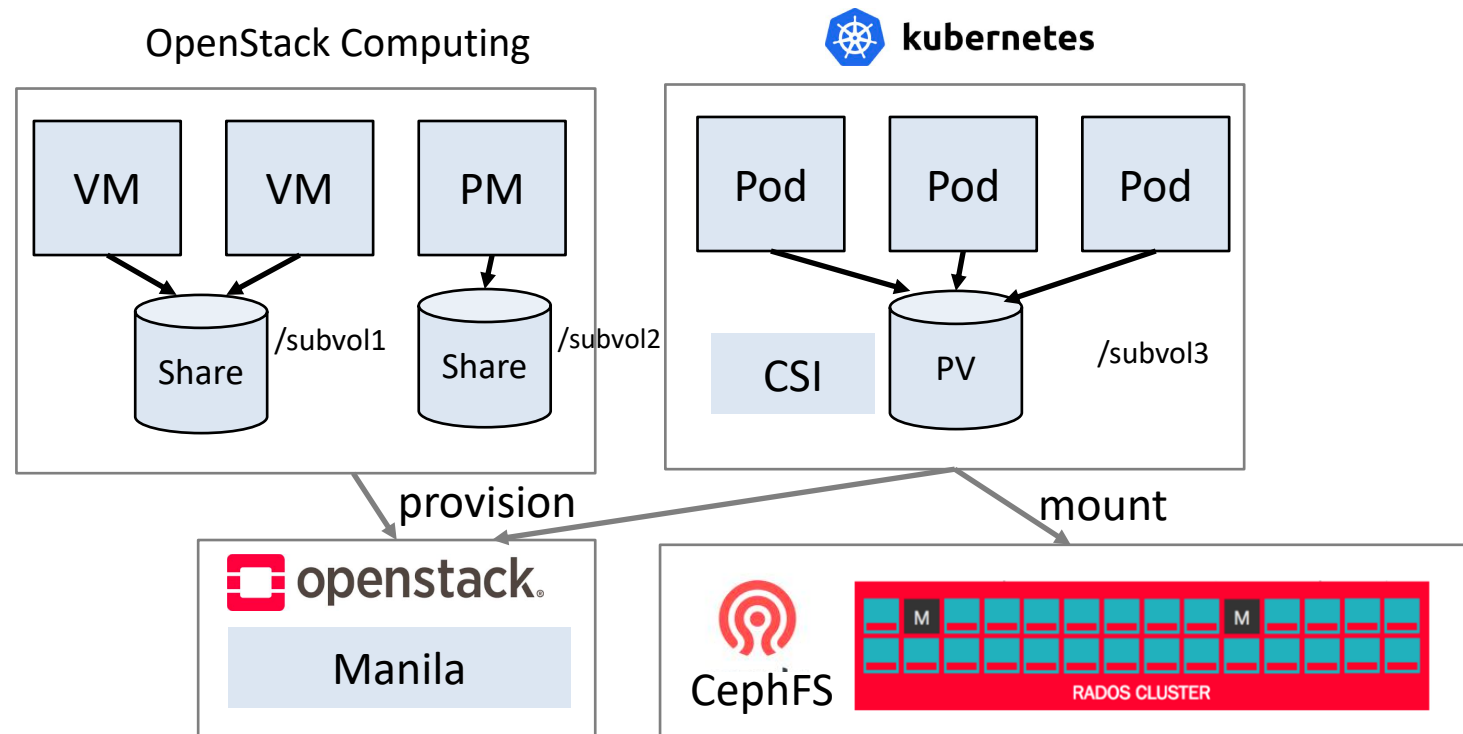https://www.youtube.com/watch?v=pDURll6Y-Ug

# Cephalocon2023

# LINE

- LINE, a messaging and communication service with over **150 million** daily active users

- OpenStack and Kubernetes as a cloud Infrastructure

- Ceph as a software defined storage
  - RBD, RGW, and then **CephFS**

# Use Cases of CephFS in LINE

- CephFS as a shared file system for clouds
  - OpenStack Manila to manage subvolumes of CephFS
- File system can be shared to VMs/Pods (major benefit)
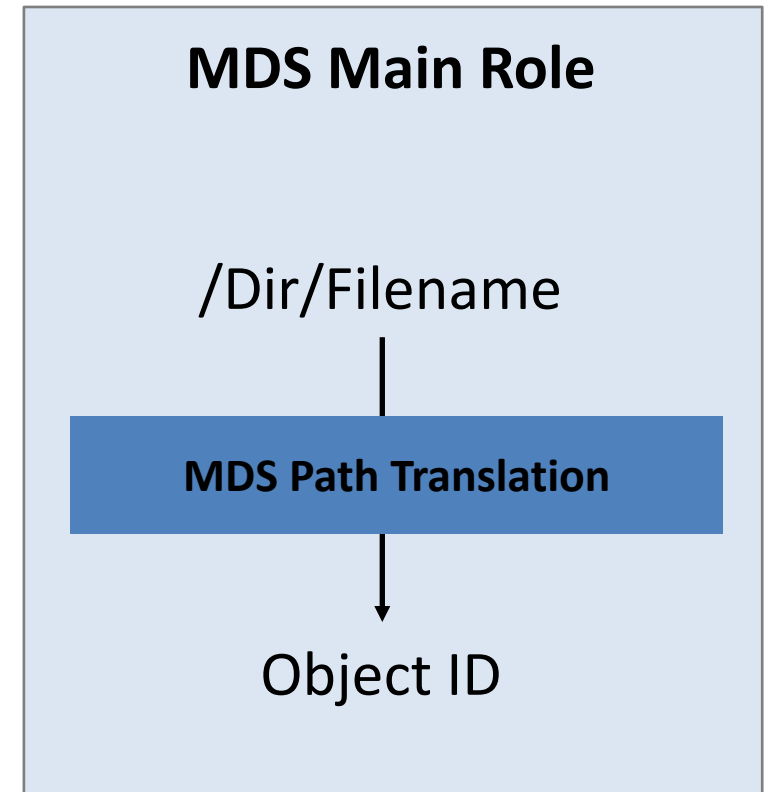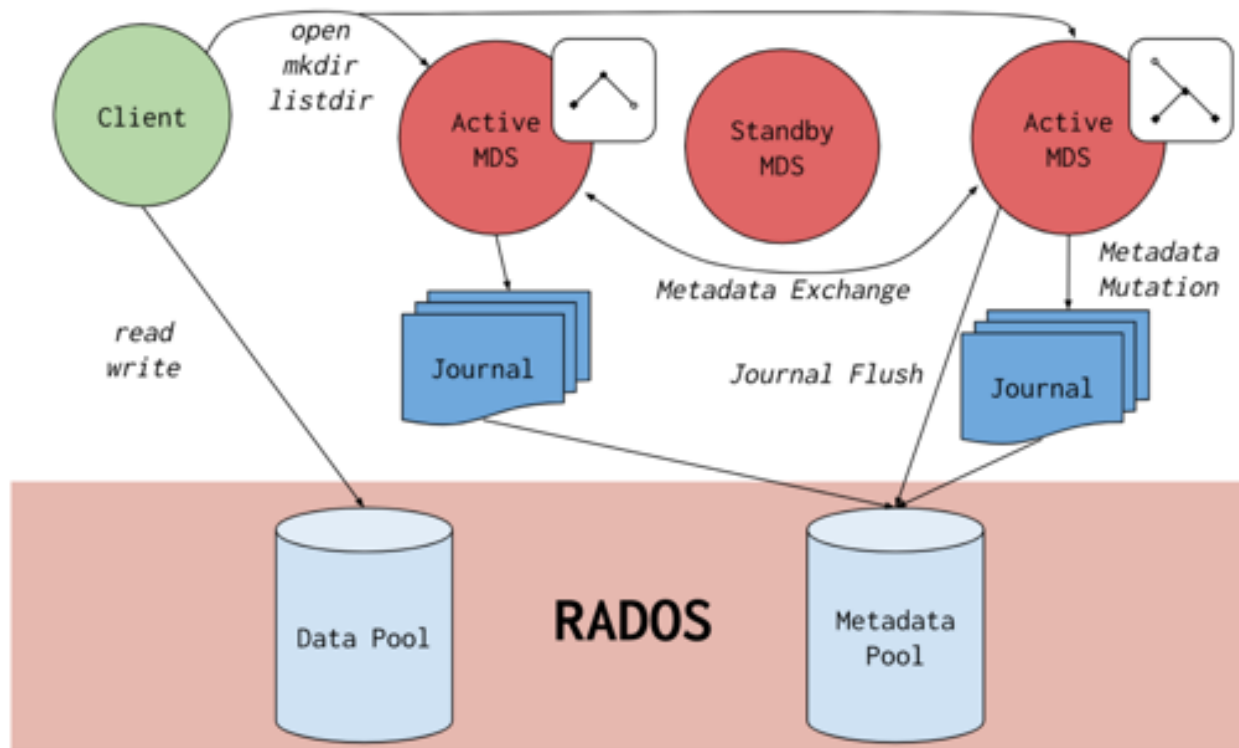  - NAS, ML training, speech data, backup

# Key Features of CephFS

- POSIX compliant distributed file system
- Multi active MDSs (e.g., scalability)
- Dynamic/static sub directory partitioning
- Standby/standby-replay MDSs (e.g., rapid failover)
- Journaling (e.g., guaranteeing metadata consistency)
- Kernel/FUSE/libcephfs client support
- Subvolume/snapshot/quota management
- QoS (not support)

# CephFS Architecture

# CephFS Clusters in LINE

- Major two clusters of our environment

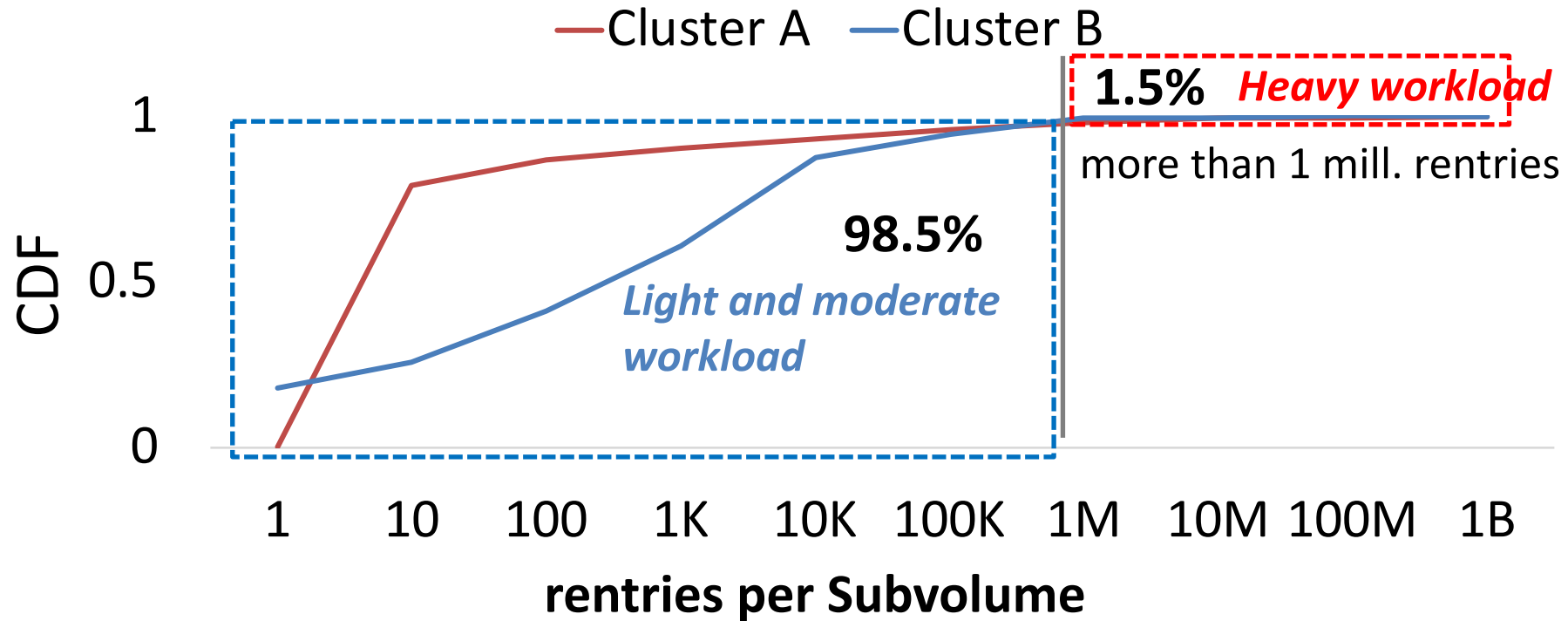|  | Cluster A | Cluster B |
|---|---|---|
| **Workload** | General Purpose | ML, Data Processing |
| **Active MDSs (on PMs + VMs)** | 30 | 36 |
| **OSDs (SSD based)** | PB scale | PB scale |
| **Sessions** | > 5,000 | > 5,000 |
| **Subvolumes** | > 3,000 | > 3,000 |

# On Deploying CephFS in Production

- Sep 2020: CephFS service with two active MDSs + standby

- Jan 2021: technical issues and outages with growing users
  - Noisy neighbor, MDS slow request by balancer, increased recovery time

- Feb 2021: simple static partitioning has been applied

- Sep 2021: proposed MDS QoS scheduler (e.g., PR 38506)

- Dec 2021: enhanced MDS recovery time (e.g., PR 44246 , PR 41358)
  - Becoming up:active took 2 hours because of heartbeat timed out
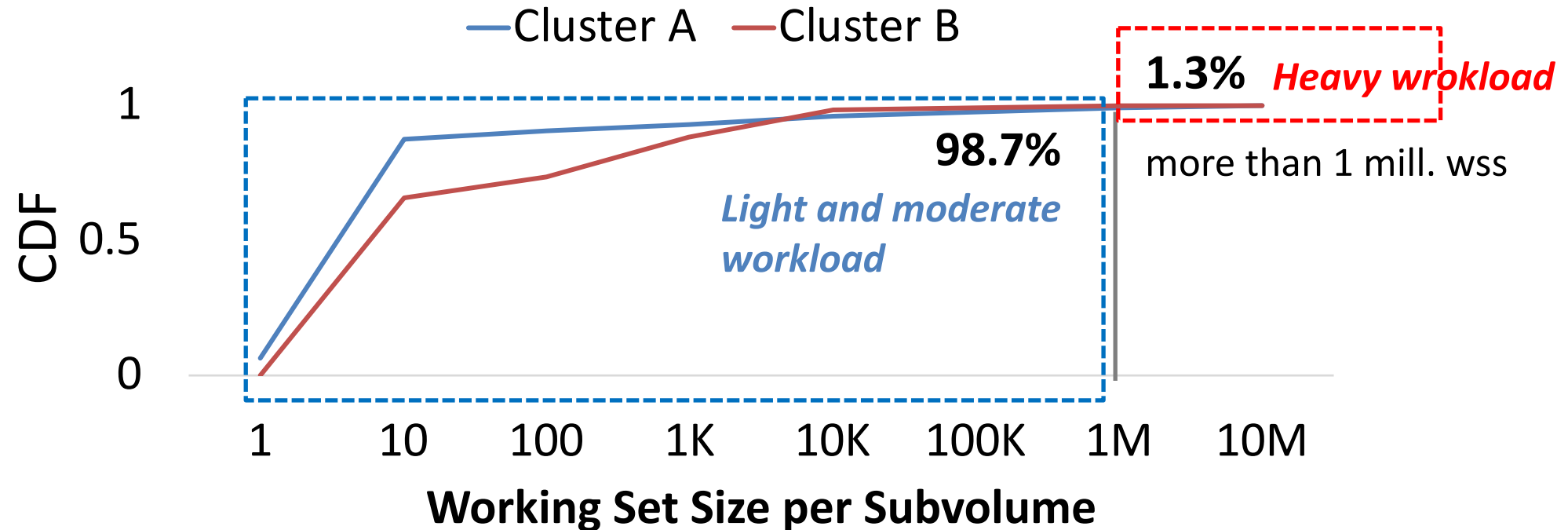
- Aug 2022: the number of MDSs has increased to up to 36

Now, we are working on *subtree partitioning!*

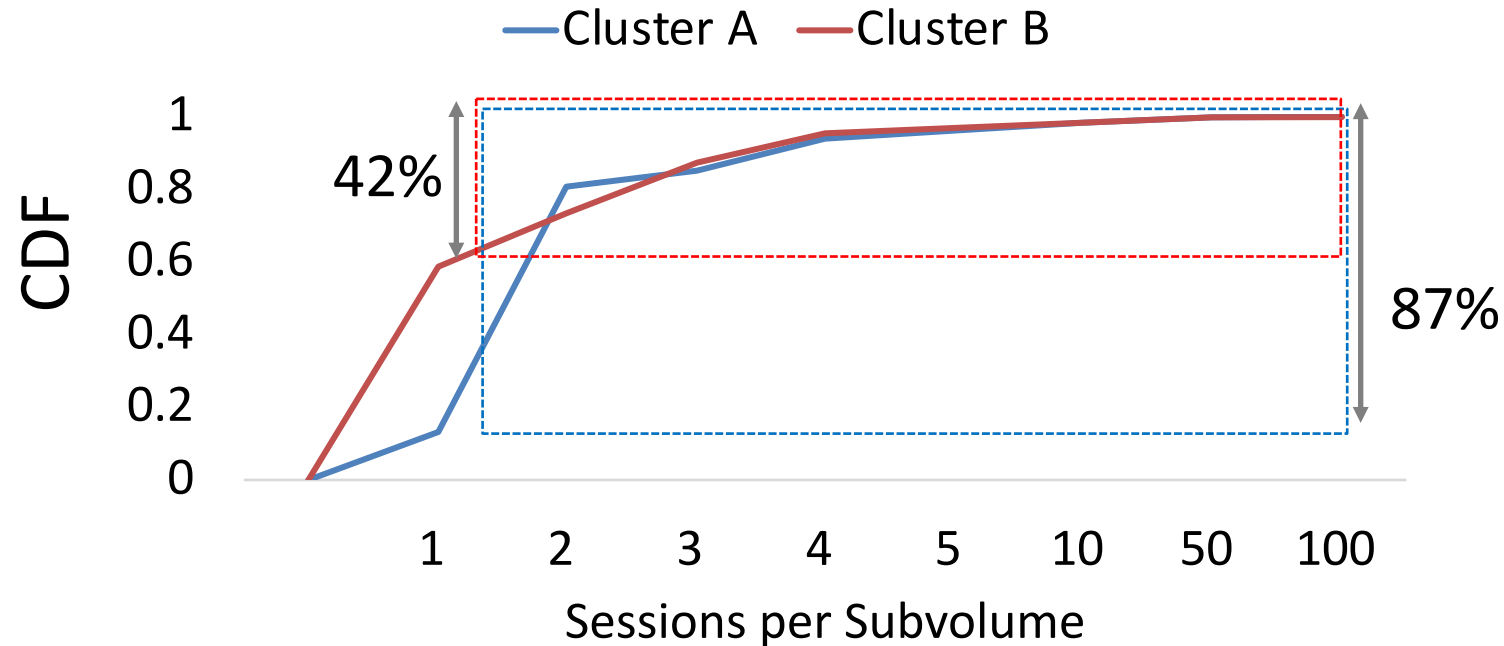# MDS Workloads: rentries Distribution



- **rentries** (e.g., files + dirs) per subvolume (measured via ceph.dir.rentries)
  - 98.5% of subvolumes contain rentries less than or equal to 1 mil.
  - Remaining of subvolumes have more than 1 mil rentries
- **Heavy workload** may affect the cache hit ratio, leading to performance degradation

# MDS Workloads: Working Set Distribution



- **Working set size** (WSS) per subvolume (collected via uniq() of Elasticsearch)
  - 98.7% subvolumes have WSS of less than 1 mil.
- Heavy workload requires the large cache memory

# MDS Workloads: Sessions per Subvolume



- 42% and 87% subvolumes are shared with sessions that have at least two sessions
  - Many sessions have the potential to generate heavy workloads
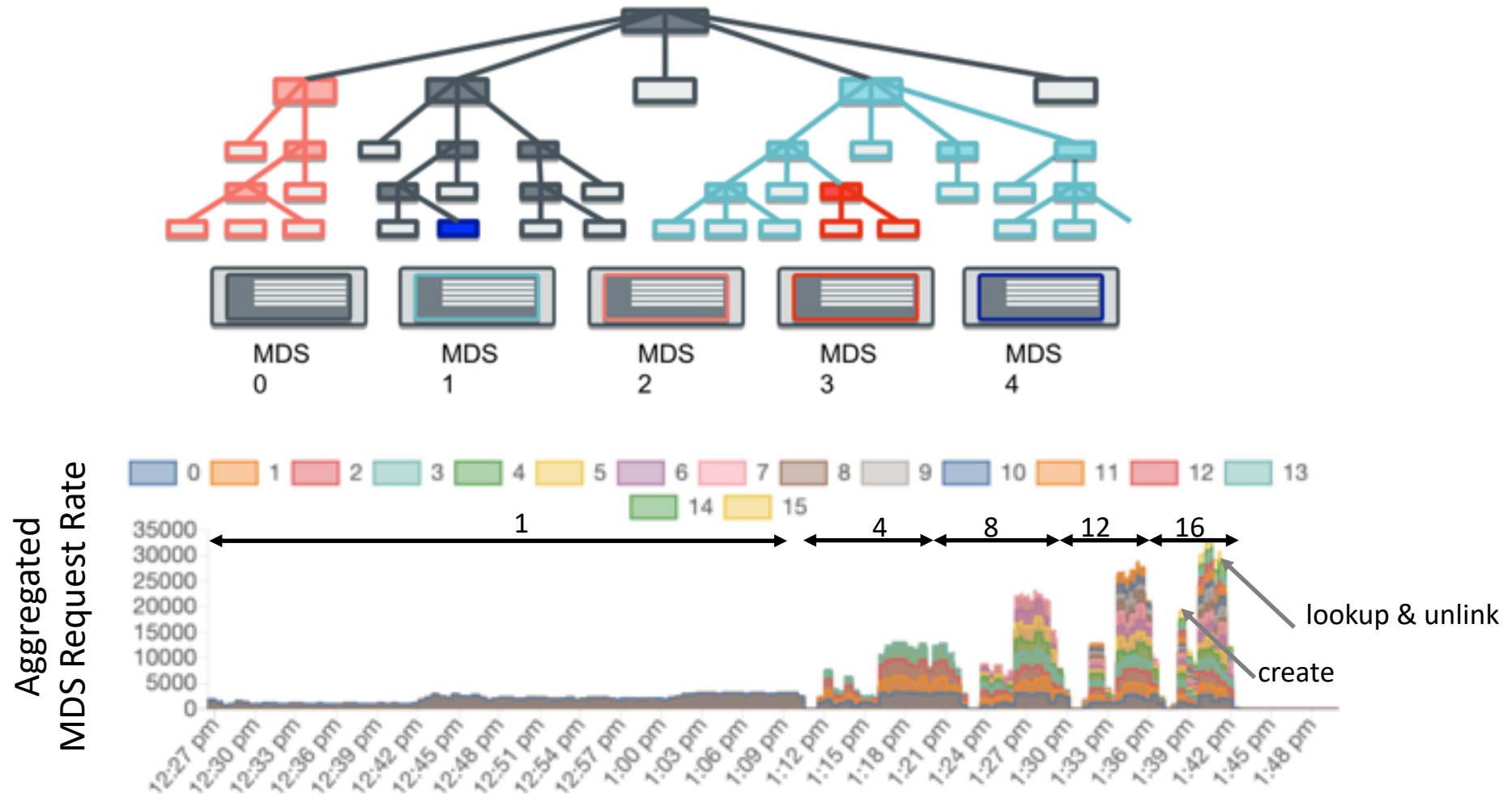
# Outline of Contents

- Introduction and background

- **CephFS subtree partitioning**

- Combining static and dynamic partitioning schemes
  - Workload based static partitioner with bal_rank_mask
  - A new CephFS MDS balancer with ceph.dir.bal.mask

- Future works and conclusions
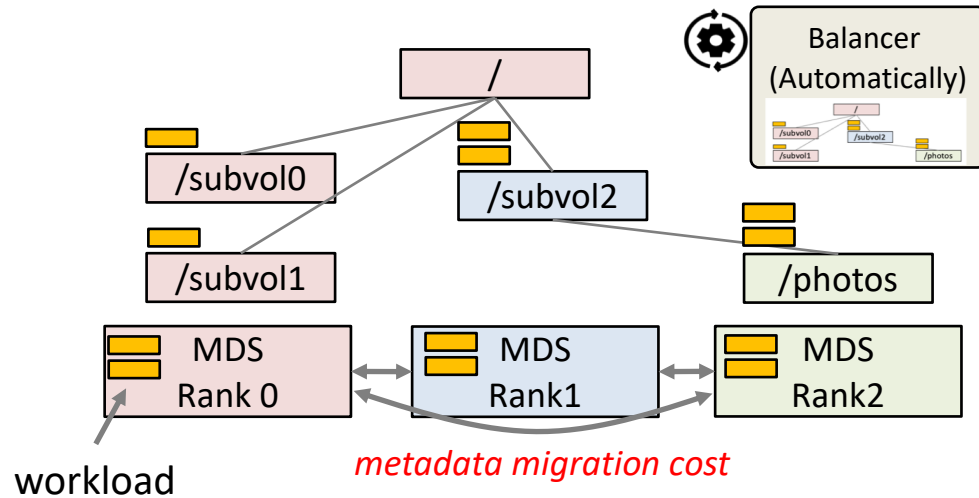
# MDS Scalability with Subtree Partitioning

Can MDS performance scale with multiple MDSs?

# Dynamic vs Static Partitioning

## Dynamic Subtree Partitioning
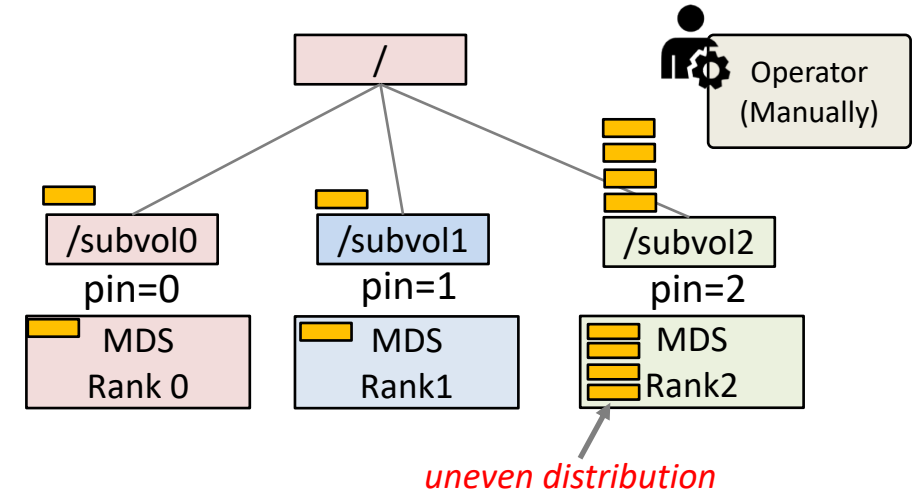Periodically redistribute subtrees



Pros
- MDS horizontal scalability

Cons
- Metadata migration cost

## Static Subtree Partitioning
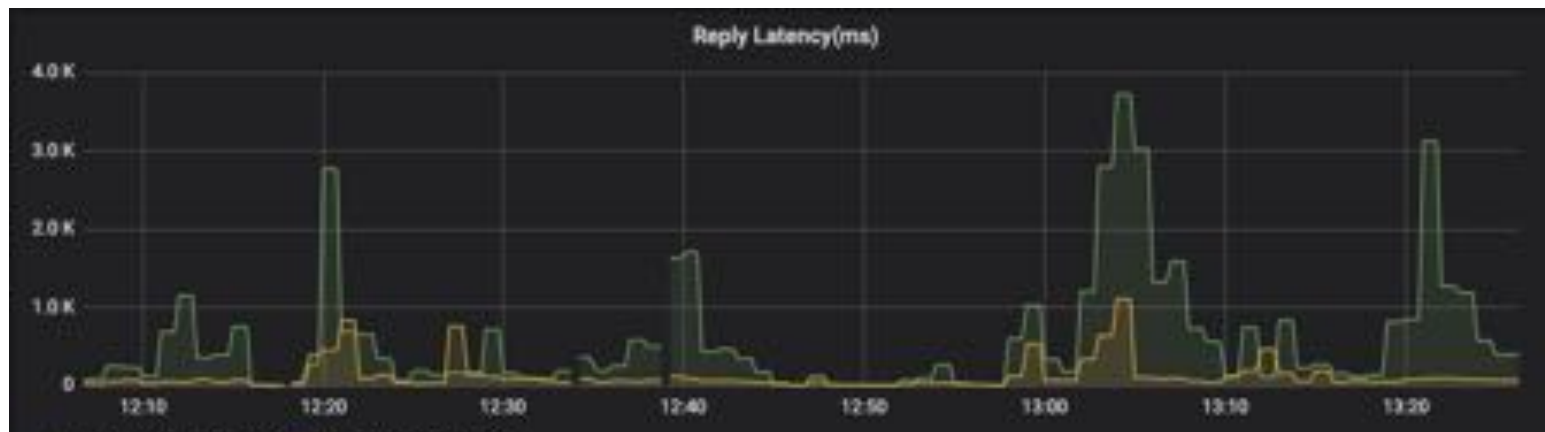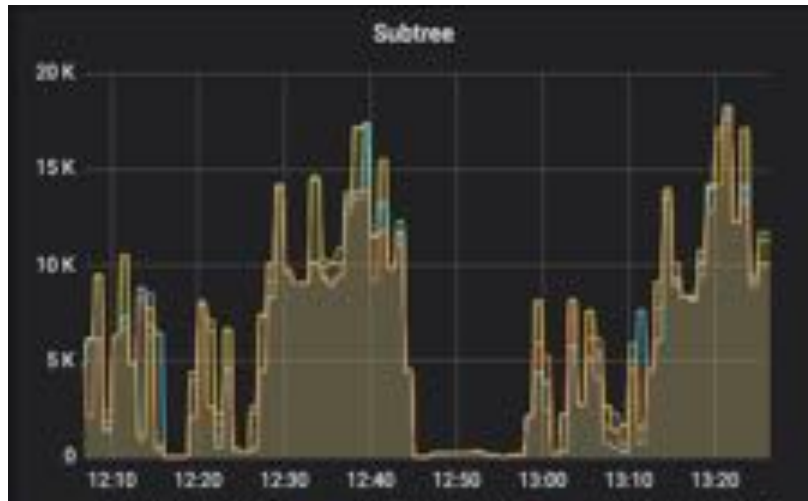Simply pin subtrees to their own ranks



Pros
- Negligible migration cost
- Provisioned performance

Cons
- Uneven workload distribution
- Additional operator efforts

# Negative Impact of Dynamic Partitioning

- Increased subtree changes and inode migrations

# Migrating Subtrees Incurs MDS Slow Requests

```
"MDS_SLOW_REQUEST": {
        "severity": "HEALTH_WARN",
        "summary": {
          "message": "2 MDSs report slow requests"
        },
        "detail": [
          {
            "message": "mds(mds.1): 1253 slow requests are blocked > 30 secs"
          },
          {
            "message": "mds(mds.0): 1 slow requests are blocked > 30 secs"
          }
        ]
}
```

We have observed slow requests are resolved since employing static partitioning!

# Outline of Contents

- Introduction and background

- CephFS subtree partitioning

- Combining static and dynamic partitioning schemes
  - Workload based static partitioner with bal_rank_mask
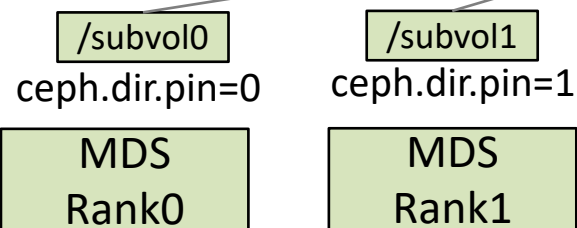  - A new CephFS MDS balancer with ceph.dir.bal.mask

- Future works and conclusions

# Our Idea: Combine Two Approaches

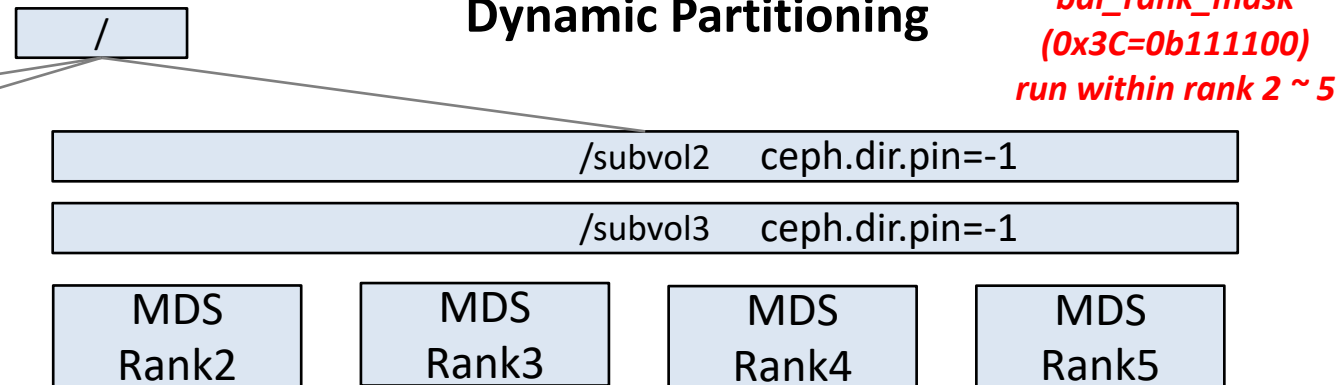*Static partitioning* for moderate workloads, while *dynamic partitioning* for heavy workloads



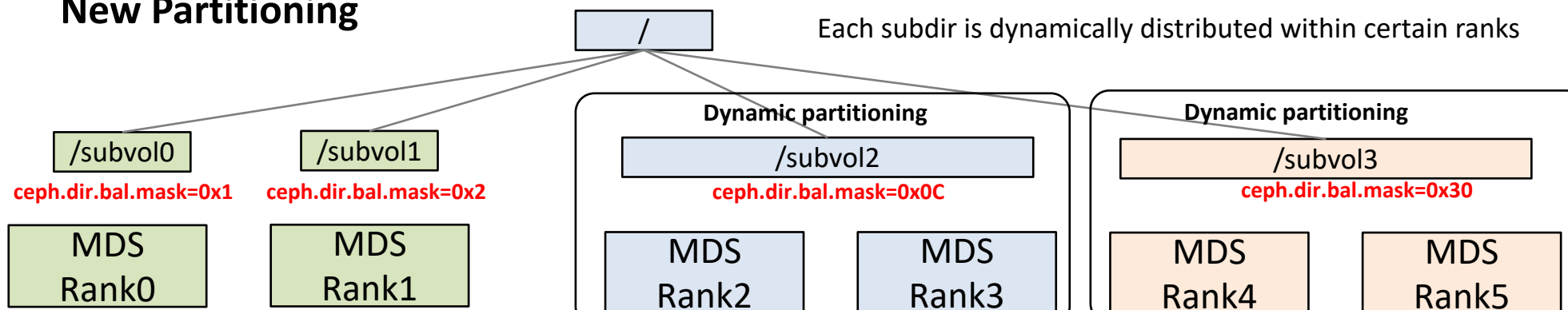**Light/Moderate Workloads**

**Heavy Workloads**

Current Version (presented at Cephalocon'23)

### Static Partitioning

/

/subvol0 — ceph.dir.pin=0

/subvol1 — ceph.dir.pin=1

MDS Rank0

MDS Rank1

### Dynamic Partitioning

**bal_rank_mask (0x3C=0b111100) run within rank 2 ~ 5**

/subvol2 — ceph.dir.pin=-1

/subvol3 — ceph.dir.pin=-1

MDS Rank2

MDS Rank3

MDS Rank4

MDS Rank5

New Version (in-progress)

### New Partitioning

/

Each subdir is dynamically distributed within certain ranks

/subvol0 — ceph.dir.bal.mask=0x1

/subvol1 — ceph.dir.bal.mask=0x2

MDS Rank0

MDS Rank1

**Dynamic partitioning**

/subvol2 — ceph.dir.bal.mask=0x0C

MDS Rank2

MDS Rank3

**Dynamic partitioning**

/subvol3 — ceph.dir.bal.mask=0x30

MDS Rank4

MDS Rank5

20

# Workload Based Static Partitioner with bal_rank_mask



- ***Workload based static partitioner*** pins subvolumes
  - Workload calculation based on working set, rentries, and performance
  - Rarely or manually conducted if loads are uneven or latencies get higher
- ***bal_rank_mask*** enables the balancer to dynamically rebalance unpinned subtrees within particular active MDS ranks PR 43284

# How to use bal_rank_mask

Change from 0b**111000** to 0b**000111**

```
$ ceph fs set cephfs bal_rank_mask 0x38
setting the metadata balancer rank mask to 0x38
$ sleep 180
$ ceph fs set cephfs bal_rank_mask 0x7
setting the metadata balancer rank mask to 0x7
```

```
$ ceph fs status
cephfs - 8 clients
======
RANK STATE MDS    ACTIVITY     DNS   INOS  DIRS  CAPS
 0   active mds001 Reqs:  0 /s  5259  5095   84    14
 1   active mds017 Reqs:  0 /s  463k  463k 31.1k  21
 2   active mds018 Reqs:  0 /s  295k  295k 18.3k 38.3k
 3   active mds016 Reqs:  2 /s 53.9k 53.9k 3573  53.6k
 4   active mds015 Reqs: 405 /s 65.4k 65.5k 4018  65.0k
 5   active mds010 Reqs: 398 /s 81.9k 81.9k 4906  81.8k
```

```
$ ceph fs status
cephfs - 8 clients
======
RANK STATE MDS     ACTIVITY     DNS  INOS  DIRS  CAPS
 0   active mds001 Reqs: 364 /s 75.8k 75.4k 5943 70.9k
 1   active mds017 Reqs: 259 /s  588k  587k 39.7k 123k
 2   active mds018 Reqs: 366 /s  395k  394k 26.5k 137k
 3   active mds016 Reqs:  5 /s   678   683  637   36
 4   active mds015 Reqs:  8 /s   504   509  500   83
 5   active mds010 Reqs:  4 /s   970   975  831   84
```

# Evaluation Environment

- Ceph Pacific 16.2.10 (integrated with our PRs) installed in Rocky8
  - bal_rank_mask - PR 43284 (merged)
  - MDS QoS Scheduler - PR 38506 (ready to review)

- VDBench tool generates workloads
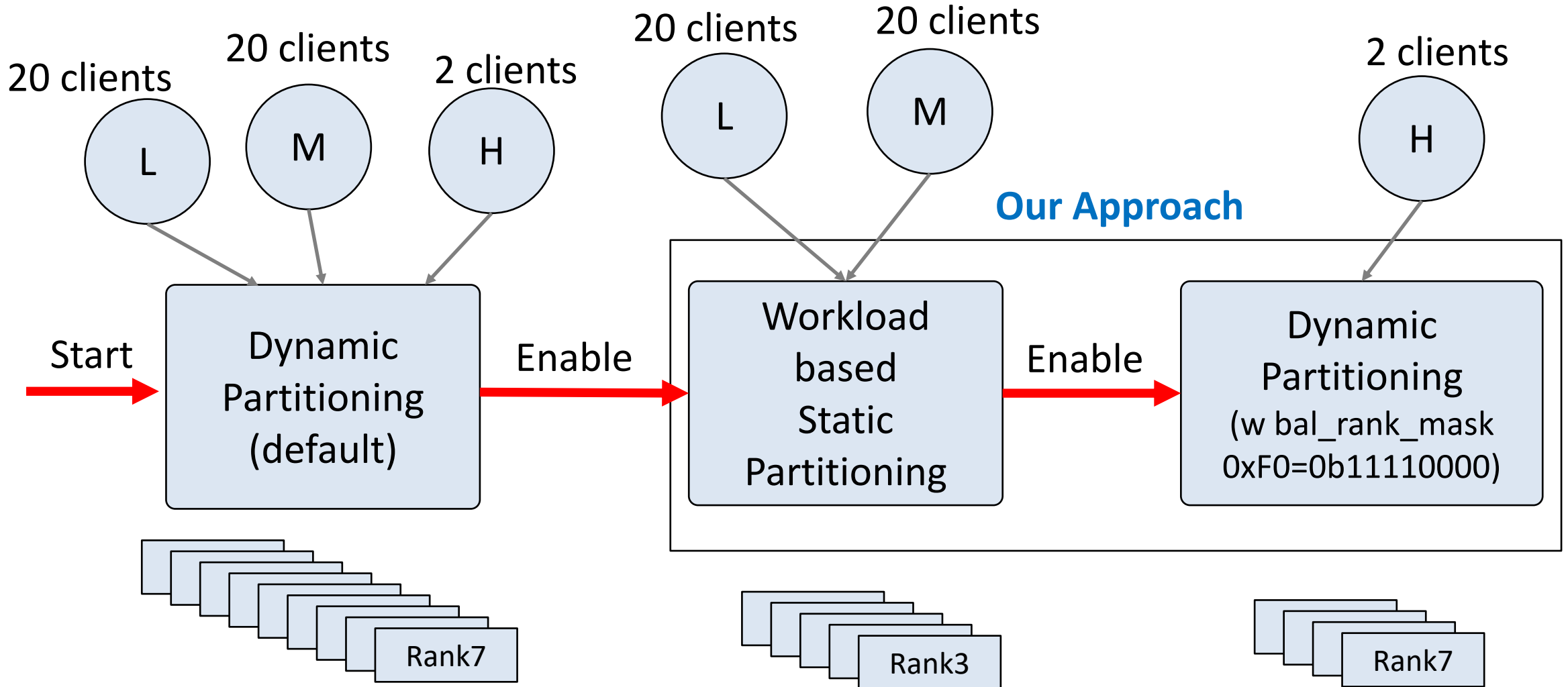  - Each client has its own subvolume (e.g., /volumes/_nogroup/$subvol)

### Workload Spec.

| Type | Files | File Size | Clients | Threads per Client |
|---|---|---|---|---|
| **Light** | 50K | 4K | 20 | 1 |
| **Moderate** | 500K | 4K | 20 | 1 |
| **Heavy** | 5,000K | 4K | 2 | 16 |

### HW Spec.

| Type | Server Spec | Count |
|---|---|---|
| MDS | 4 * vCPU 32GB RAM | active 8 (mds_cache_memory_limt: 12GB) |
| OSD | 40 * pCPU 128GB RAM | 6 Servers * 6 SATA SSDs = 36 |
| Mon | 2 * vCPU 4GB RAM | 3 |
| Mgr | 4 * vCPU 8GB RAM | 3 |
| Client | 4 * vCPU 8GB RAM | 42 |
| Network | 10Gbps | |

# Evaluation Sequences



20 clients — L

20 clients — M

2 clients — H

20 clients — L

20 clients — M

2 clients — H

**Our Approach**

Start → Dynamic Partitioning (default) → Enable → Workload based Static Partitioning → Enable → Dynamic Partitioning (w bal_rank_mask 0xF0=0b11110000)
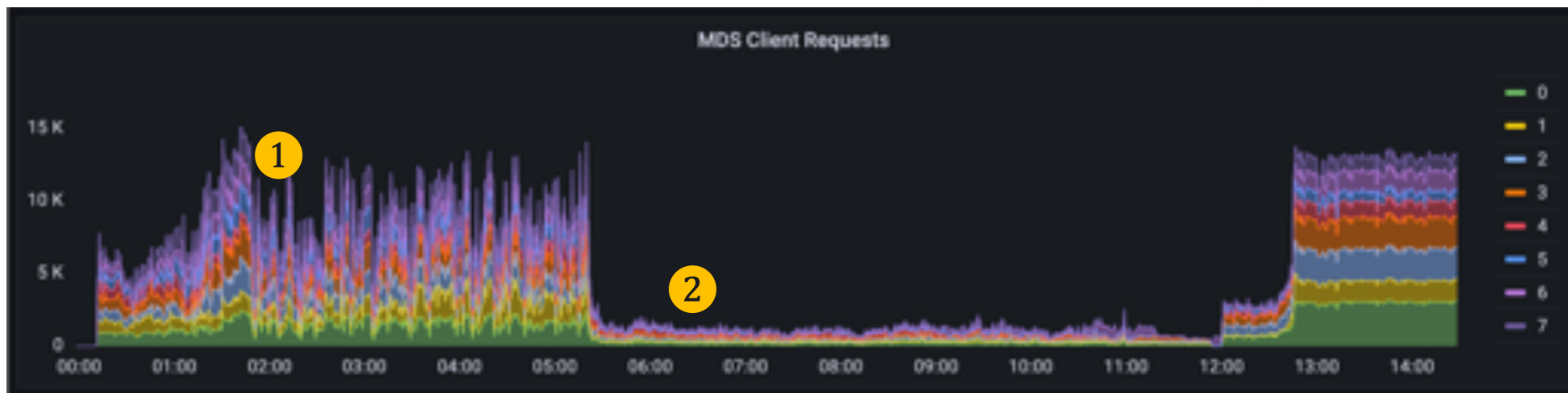
Rank7

Rank3

Rank7

24

# Evaluation Results (1/3)



**Balancer default**

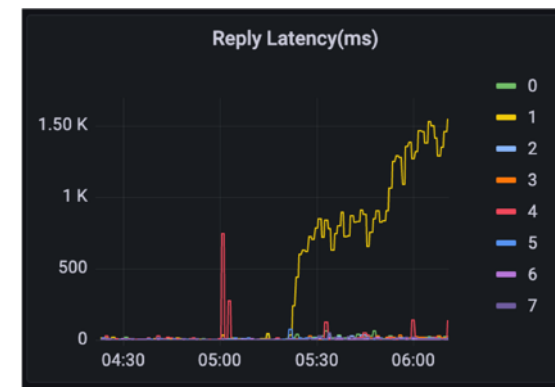Our Static Partitioning

Dynamic Partitioning bal_rank_mask 0xf0

① Performances of MDSs are fluctuated due to exporting/importing inodes

② inodes are unevenly distributed rank1 keeps inodes more than 4mil.
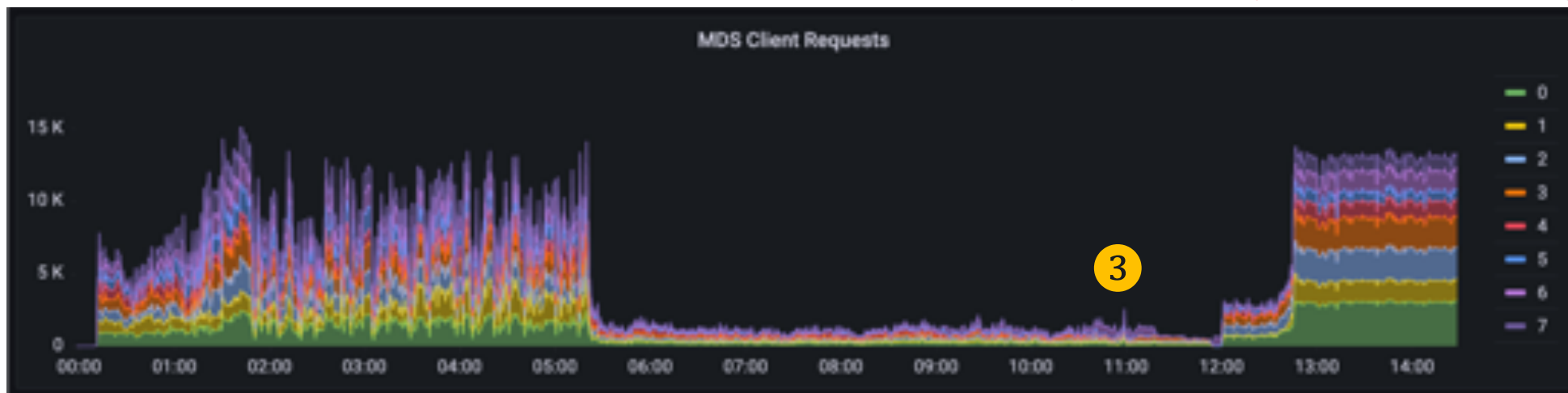
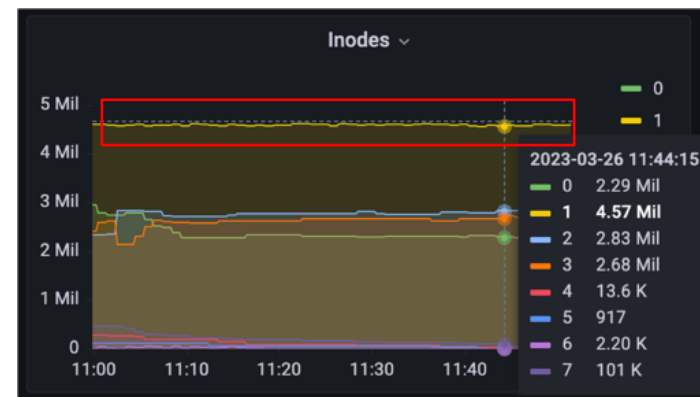Reply latencies highly increased

# Evaluation Results (2/3)



subvolumes of light/moderate workloads are moved from all ranks to rank0~3

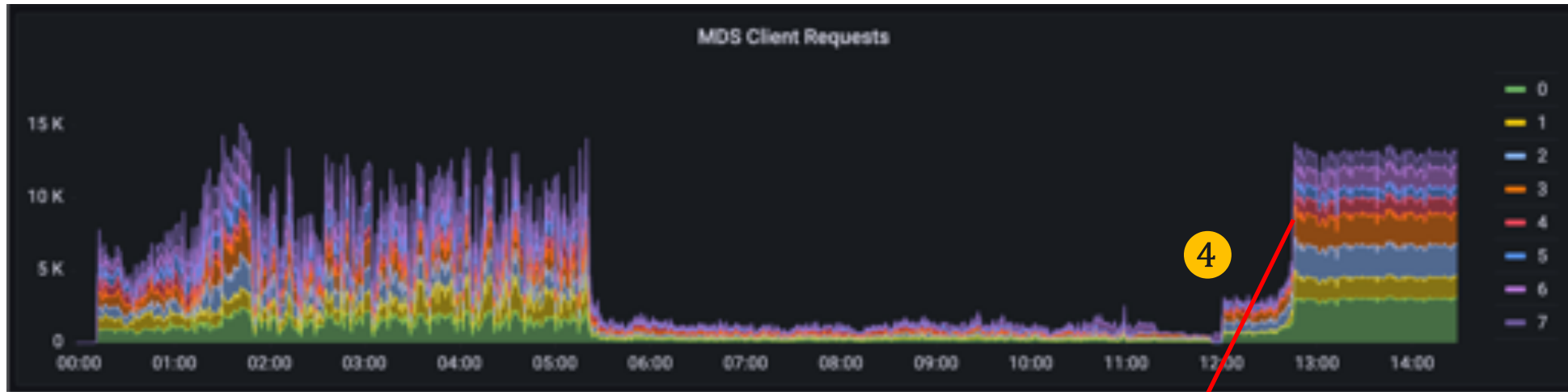Performance is still not recovered as rank 1 has a lot of inodes
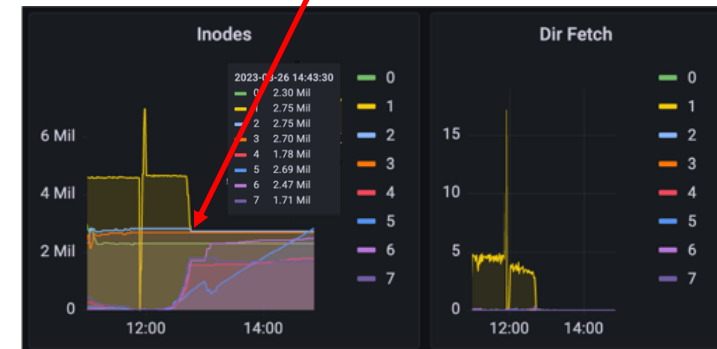
# Evaluation Results (3/3)

Balancer default

Our
Static
Partitioning

**Dynamic Partitioning**
***bal_rank_mask* 0xf0**



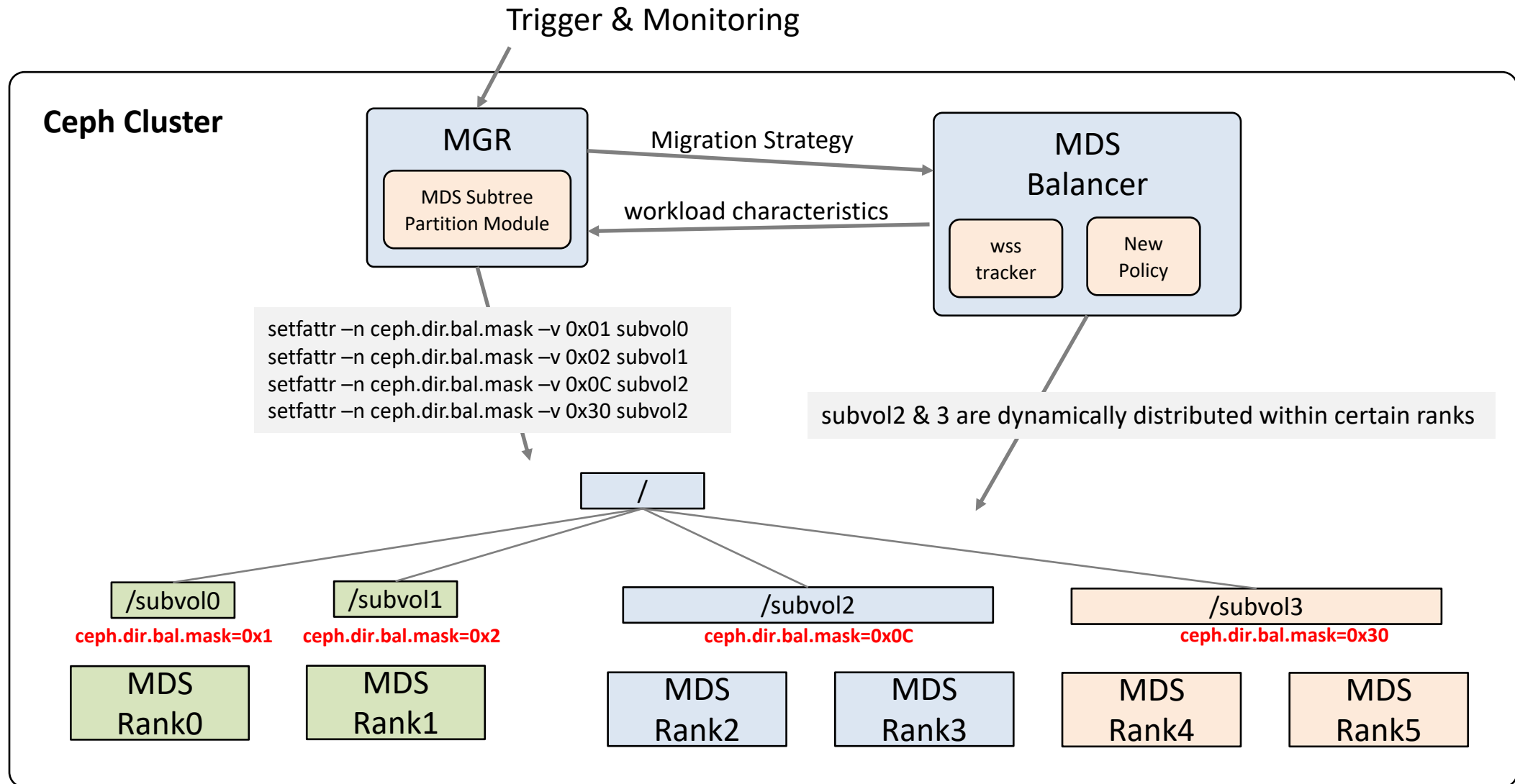④ subvolumes of heavy workloads are migrated to rank 4~7

inodes are balanced and dir fetch count decreases
and performance increases

# Outline of Contents

- Introduction and background

- CephFS subtree partitioning

- Combining static and dynamic partitioning schemes
  - Workload based static partitioner with bal_rank_mask
  - **A new CephFS MDS balancer with ceph.dir.bal.mask**

- Future works and conclusions

# A New MDS Balancer with ceph.dir.bal.mask

# Implementation

- rank mask option per subdir as a virtual extended attribute
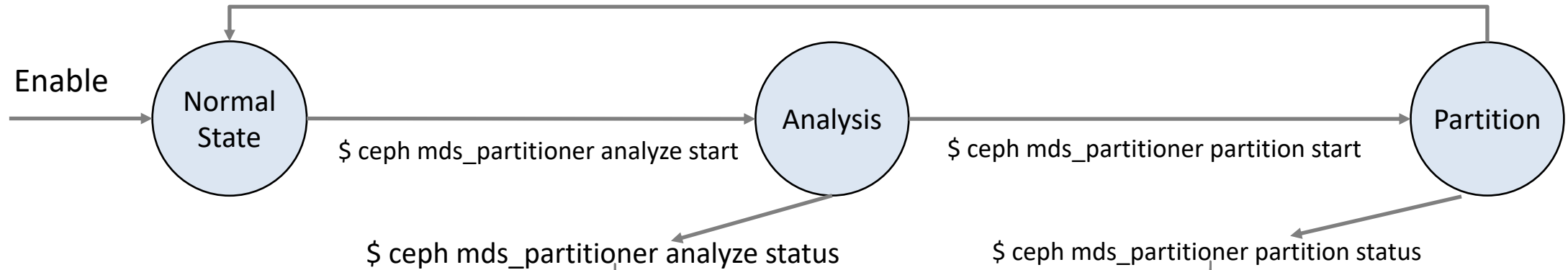  - A target subdir is dynamically within certain MDS ranks (e.g., rank0 and 1)

    setfattr –n ceph.dir.**bal.mask** –v 0x3 /cephfs/home/yongseok

- MDS Subtree Partition Module in MGR

    ceph mgr module enable mds_partitioner
    ceph mds_partitioner analyze start   # analyze client workloads ontained from MDSs
    ceph mds_partitioner analyze status # report analysis results and recommend optimal the number of MDSs
    ceph mds_partitioner partition start # start partitioning
    ceph mds_partitioner partition status # report partitioning status

- MDS Balancer modifications
  - Working set size tracker
  - Migrate subdirs based on ceph.dir.bal.mask values of subdirs
  - Minimize MDS slow requests

# Example of Operation Flow



Enable → Normal State

$ ceph mds_partitioner analyze start

Analysis

$ ceph mds_partitioner partition start

Partition

$ ceph mds_partitioner analyze status

$ ceph mds_partitioner partition status

| Name | wss | reqs | rentries | workload | Current Ranks | New Ranks | Migration Progress | Migration Status |
|------|-----|------|----------|----------|---------------|-----------|--------------------|------------------|
| Subvol1 | 1,000,000 | 1,203,030 | 50,000,000 | 339 | 0 | 0,1 | 20% | In-progress |
| Subvol2 | 700,000 | 500,000 | 1,000,000 | 137 | 1 | 2 | 0% | Ready |
| subvol3 | 100,000 | 5,000 | 5,000,000 | 21 | 2 | 2 | 100% | Done |
| subvol4 | 3,000 | 20,000 | 70,000 | 3 | 2 | 2 | 100% | Done |
| Total | 1,803,000 | 1,728,030 | 56,070,000 | 500 | | | | |

wss: working set size
reqs: requests
rentries: files + dirs

**workload** = (working_set_size / total_working_set *2
+ requests / total_requests * 2
+ rentries / total_rentries) * 100

# Conclusions

- We employ CephFS as a shared file service in LINE's cloud
- We compared static and dynamic partitioning schemes
  - Dynamic partitioning incurs metadata migration cost
  - Static partitioning present uneven workload distribution
- We presented how to combine both static and dynamic partitioning scheme effectively
- We will contribute our work on a new partitioning to the community

# Thank you!

Any Questions?

# References

- Ceph: A Scalable, High-Performance Distributed File System, Sage Weil, OSDI'06
- Overview and Status of the Ceph File System, Patrick Donnelly, 2018, https://indico.cern.ch/event/644915/
- ceph-linode for CephFS testing, Patrick Donnelly, https://github.com/batrick/ceph-linode
- CephFS with OpenStack Manila based on BlueStore and Erasure Code, 2018, https://cutt.ly/dc7Qnn7
- Revisiting CephFS MDS and mClock QoS Scheduler, Yongseok Oh, 2021, https://cutt.ly/A4s8AsC
- Ephemeral Pinning: A Dynamic Metadata Management Strategy for CephFS, Sidharth Anupkrishnan, 2020 https://www.youtube.com/watch?v=zimAEm_8efA
- Optimizing CephFS with Combining MDS QoS Scheduling and Static-Dynamic Subtree Partitioning, Yongseok Oh, https://ceph2023.sched.com/event/1JKas/optimizing-cephfs-with-combining-mds-qos-scheduling-and-static-dynamic-subtree-partitioning-yongseok-oh-jinmyeong-lee-line

# VDBench Parameter Configs

**client 01~02**

```
fsd=fsd1,anchor=/mnt,depth=1,width=1,files=50000,size=4k
fwd=fwd1,fsd=fsd1,operation=getattr,xfersize=4k,fileio=sequential,fileselect=random,threads=1
fwd=fwd2,fsd=fsd1,operation=read,xfersize=4k,fileio=sequential,fileselect=random,threads=1
rd=rd1,fwd=(fwd1,fwd2),fwdrate=max,elapsed=86400,format=yes,interval=10
```

**client 03~20**

```
fsd=fsd1,anchor=/mnt,depth=1,width=1,files=50000,size=4k
fwd=fwd1,fsd=fsd1,operation=getattr,xfersize=4k,fileio=sequential,fileselect=random,threads=1
fwd=fwd2,fsd=fsd1,operation=create,xfersize=4k,fileio=sequential,fileselect=random,threads=1
fwd=fwd3,fsd=fsd1,operation=delete,xfersize=4k,fileio=sequential,fileselect=random,threads=1
fwd=fwd4,fsd=fsd1,operation=write,xfersize=4k,fileio=sequential,fileselect=random,threads=1
rd=rd1,fwd=(fwd1,fwd2,fwd3,fwd4),fwdrate=max,elapsed=86400,format=yes,interval=1
```
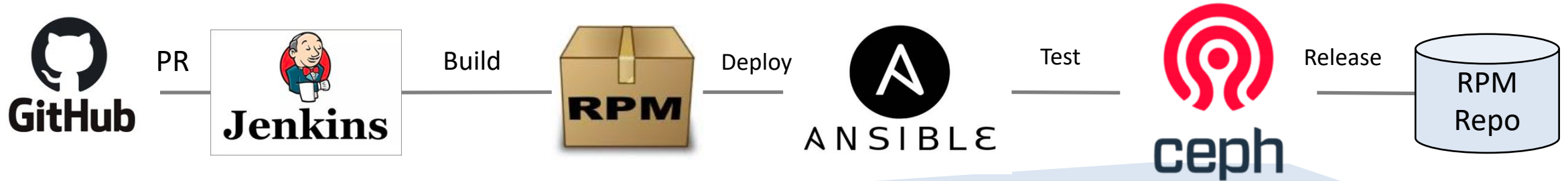
**client 21~22**

```
fsd=fsd1,anchor=/mnt,depth=1,width=10,files=50000,size=4k
fwd=fwd1,fsd=fsd1,operation=getattr,xfersize=4k,fileio=sequential,fileselect=random,threads=1
fwd=fwd2,fsd=fsd1,operation=read,xfersize=4k,fileio=sequential,fileselect=random,threads=1
rd=rd1,fwd=(fwd1,fwd2),fwdrate=max,elapsed=86400,format=yes,interval=10
```

**client 23~40**

```
fsd=fsd1,anchor=/mnt,depth=1,width=10,files=50000,size=4k
fwd=fwd1,fsd=fsd1,operation=getattr,xfersize=4k,fileio=sequential,fileselect=random,threads=1
fwd=fwd2,fsd=fsd1,operation=create,xfersize=4k,fileio=sequential,fileselect=random,threads=1
fwd=fwd3,fsd=fsd1,operation=delete,xfersize=4k,fileio=sequential,fileselect=random,threads=1
fwd=fwd4,fsd=fsd1,operation=write,xfersize=4k,fileio=sequential,fileselect=random,threads=1
rd=rd1,fwd=(fwd1,fwd2,fwd3,fwd4),fwdrate=max,elapsed=86400,format=yes,interval=1
```

**client 41**

```
fsd=fsd1,anchor=/mnt,depth=1,width=100,files=50000,size=4k
fwd=fwd1,fsd=fsd1,operation=getattr,xfersize=4k,fileio=sequential,fileselect=random,threads=16
fwd=fwd2,fsd=fsd1,operation=read,xfersize=4k,fileio=sequential,fileselect=random,threads=16
rd=rd1,fwd=(fwd1,fwd2),fwdrate=max,elapsed=86400,format=yes,interval=10
```

**client 42**

```
fsd=fsd1,anchor=/mnt,depth=1,width=100,files=50000,size=4k
fwd=fwd1,fsd=fsd1,operation=getattr,xfersize=4k,fileio=sequential,fileselect=random,threads=16
fwd=fwd2,fsd=fsd1,operation=create,xfersize=4k,fileio=sequential,fileselect=random,threads=16
fwd=fwd3,fsd=fsd1,operation=delete,xfersize=4k,fileio=sequential,fileselect=random,threads=16
fwd=fwd4,fsd=fsd1,operation=write,xfersize=4k,fileio=sequential,fileselect=random,threads=16
rd=rd1,fwd=(fwd1,fwd2,fwd3,fwd4),fwdrate=max,elapsed=86400,format=yes,interval=10
```
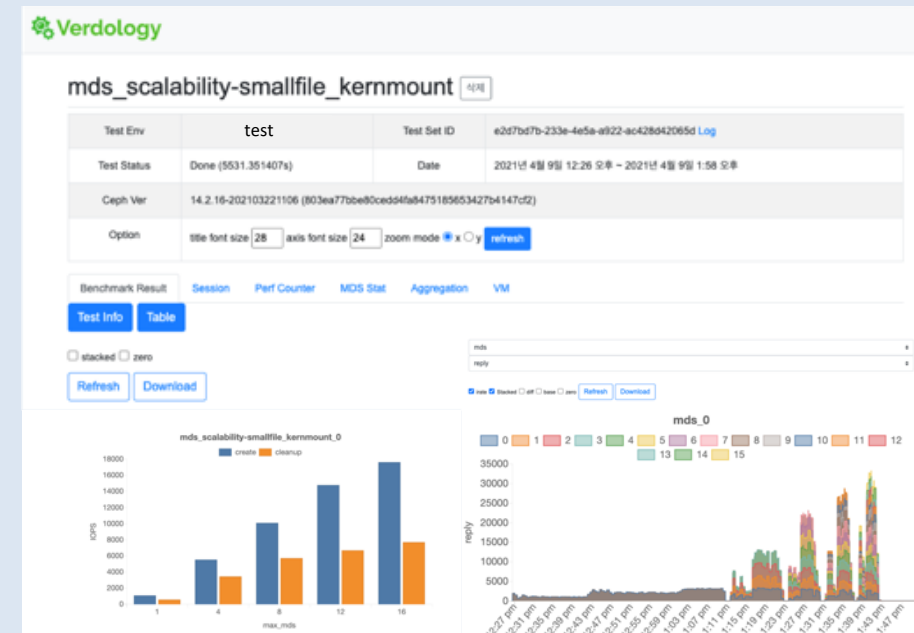
# Ceph CI/CD System in LINE



PR — Jenkins — Build — RPM — Deploy — ANSIBLE — Test — ceph — Release — RPM Repo

**Ceph Test Frameworks**

[teuthology]  [verdology]

# CephFS Test Flow of Verdology



OSD/FS Creation

VM Provision

Pkg & Tool Installation

Subvolume Allocation

Subvolume Mount

## Benchmark Tools

Smallfile

FIO

VDBench

Kernel Compile

Purge Queue Flush

Subvolume Umount

Subvolume Release

VM Release

OSD/FS Removal

Result JSON Post