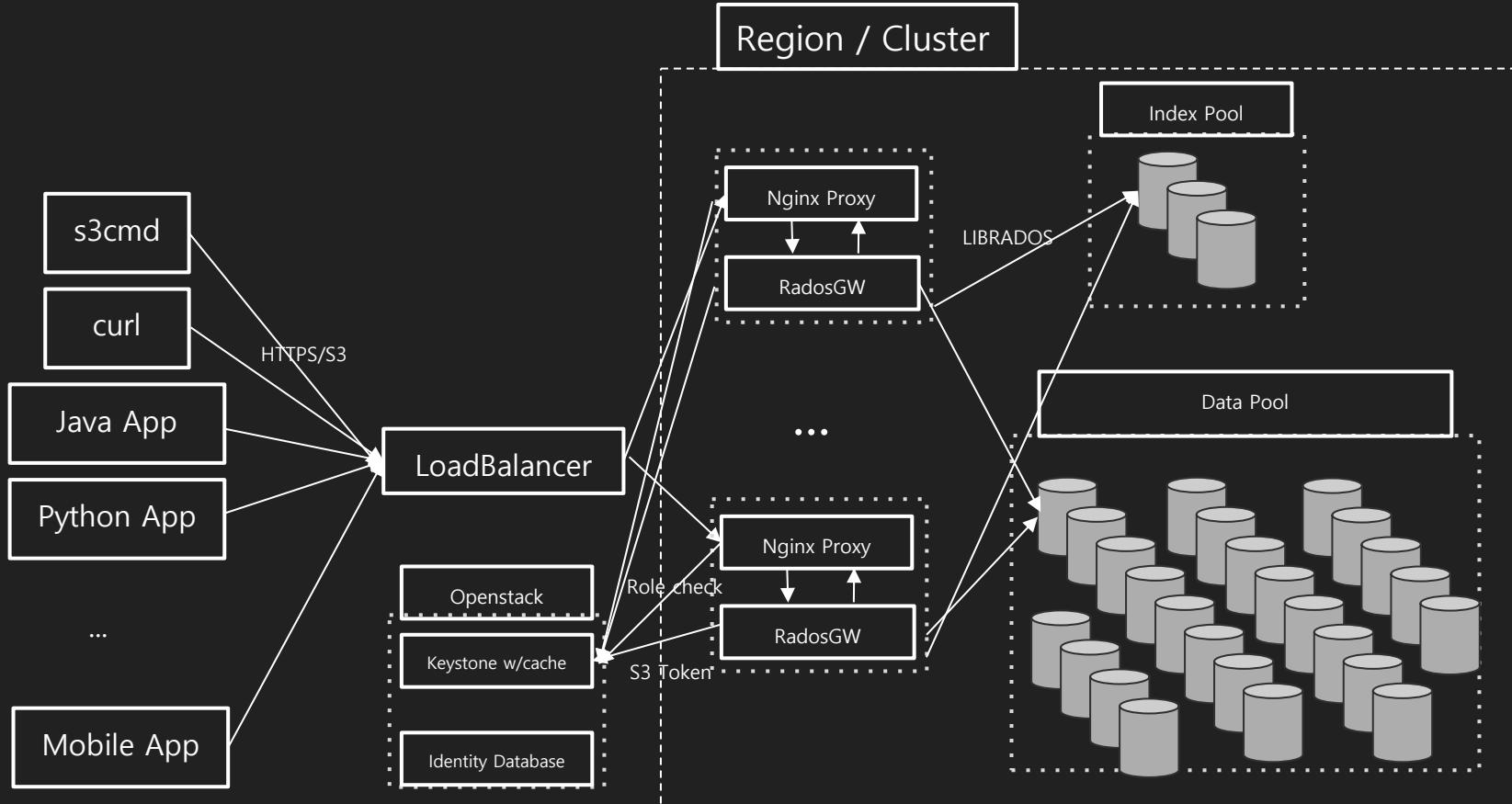


# Online migration process of 6 billion objects from HDD-based OSD to SSD-based OSD

Open Infra Community Day Korea 2020  
Dongsoo Song, LINE+

# Object Storage in LINE



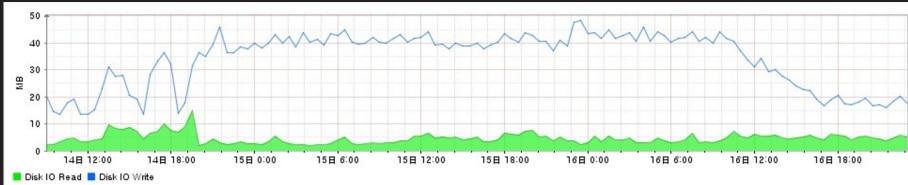
# Before Migration

Category	#	Description
<b>Number of Objects</b>	<b>50 Billions – 125 TB</b>	<ul style="list-style-type: none"><li>average object size = 25KB</li><li>objects per OSD = 104 millions</li><li>objects per PG = 610 thousands</li><li>PGs per OSD = 170</li><li>Total objects = 50 Billions * 3 replicas = 150 Billions (375 TB)</li></ul>
<b>Storage Type</b>	<b>File Store</b>	<ul style="list-style-type: none"><li>Journal on NVMe</li><li> RocksDB on xfs</li><li>Data on xfs , 100M+ inode / xfs</li></ul>
<b>Storage Media</b>	<b>NVMe SSD HDD</b>	<ul style="list-style-type: none"><li>Index Pool = NVMe</li><li>Data Pool = 7.2RPM SAS HDD</li></ul>
<b>Traffic</b>	<b>Up - 20 Millions Down - 40 Millions</b>	<ul style="list-style-type: none"><li>Contents - mostly web contents as CDN Origin</li><li>CDN traffic - 1~2 Billions / Day</li></ul>
<b>Maintenance</b>	<b>patch daemon &amp; restart daemon</b>	<ul style="list-style-type: none"><li>20~30 mins to restart a single OSD servers</li><li>12 servers * 12 OSDs =&gt; 6+ hours to patch daemon</li></ul>

# Motivation

- 1 HDD failure
  - 8 hours to recovery by self healing  
(ex, if 1 disk failed, data in failed disk are rebalanced to the rest of 143 disks) → IOs are spread into 143 disks.
  - at least 4 days to backfill data when failed disk is replaced  
(ex. data need to be recovered to 1 disk from 143 disks) ← IOs are concentrated into 1 disk.

	<b>recovery time (30% full) - 50 Millions</b>	<b>recovery time (60% full) - 100 Millions</b>
Self Healing	4 hours	8 hours
Replacement	52 hours	104 hours



← backfilling started from 13:15 14/05 ~ until 17:20 16/05  
= ~ 52 Hours

# Initial Plan

Category	Plan	Description
<b>Servers to be replaced</b>	<p>HDD 8TB * 12ea * 12 2U servers ==&gt; SSD 3.8TB * 24ea * 12 2U servers</p>	<p>raw 1PB, usable 340TB 5 billions objects → total 15 billions objects need to be migrated</p>
<b>Estimated migration window</b>	3 months	2 months - backfilling 1 month - decommission HDD servers
<b>Steps</b>	<ol style="list-style-type: none"><li>1. insert 3 SSD servers and backfill objects</li><li>2. repeat step#1 3 more times</li><li>3. decommission 3 HDD servers</li><li>4. repeat step#3 3 more times</li></ol>	expected least service impaction during backfilling

## In real world (1)

1 pg backfilling = 600,000 objects are read from HDD + 600,000 objects are removed from HDD + GET/PUT to from HDD

→ cause disk utilization up to 100%

If `osd_max_backfills` = 1 → estimated 1 year to complete



## In real world (2)

`osd_max_backfills = 5` → estimated 3 months to complete, but ...

```
cluster:
  id: 2d07cbf3-123b-42ec-a73d-2a7a7a54f021
  health: HEALTH_WARN
    noout,noscrub,nodeep-scrub flag(s) set
    16 large omap objects
    2852437587/15024729363 objects misplaced (18.985%)

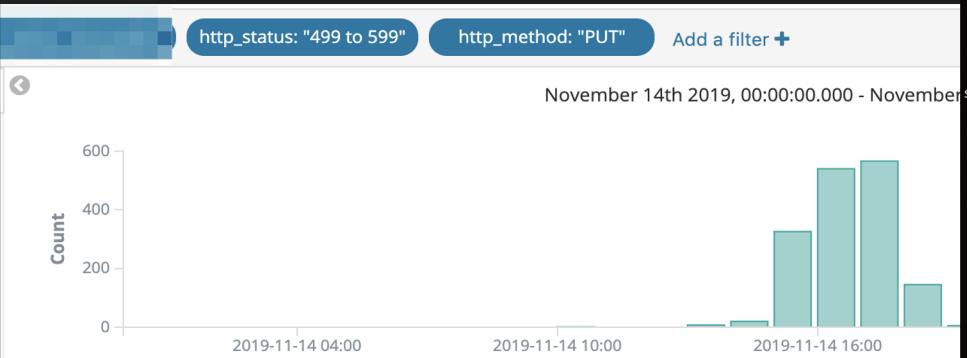
services:
  mon: [REDACTED]
  mgr: [REDACTED]
  osd: 720 osds: 720 up, 720 in; 3999 remapped pgs
    flags noout,noscrub,nodeep-scrub
  rgw: 12 daemons active

data:
  pools: 13 pools, 13632 pgs
  objects: 5.01G objects, 159TiB
  usage: 552TiB used, 1.53PiB / 2.07PiB avail
  pgs: 2852437587/15024729363 objects misplaced (18.985%)
    9633 active+clean
    3732 active+remapped+backfill_wait
    267 active+remapped+backfilling
```

## In real world (3)

OSD failed to respond health check inspite of alive → flapping osd

⇒ **ceph osd set nodown**



cluster:  
id: 2d07cbf3-123b-42ec-a73d-2a7a7a54f021  
health: HEALTH\_WARN  
nodown, oout,noscrub,nodeep-scrub flag(s) set  
16 large omap objects  
2716297836/15041207850 objects misplaced (18.05%)  
Degraded data redundancy: 36208250/15041207850 objects degraded (0.241%), 117 pgs degraded, 63 pgs under  
sized  
75 slow requests are blocked > 32 sec. Implicated osds 62,80,89,91,94,98,102,124,128,142,145,158,162  
Ticket 즉 서비스에서 문제가 있는지 확인  
services: help\_verda  
mon:  
mgr:  
osd: 720 osds: 720 up, 720 in; 3916 remapped pgs  
flags nodown,noout,noscrub,nodeep-scrub  
rgw: 12 daemons active  
data:  
pools: 13 pools, 13632 pgs  
objects: 5.01G objects, 160TiB  
usage: 556TiB used, 1.53PiB / 2.07PiB avail  
pgs: 36208250/15041207850 objects degraded (0.241%)  
2716297836/15041207850 objects misplaced (18.05%)  
[가 발생한 시간대]  
GIT L  
2019  
9679 active+clean  
3761 active+remapped+backfill\_wait  
65 active+remapped+backfilling  
62 active+undersized+degraded+remapped+backfill\_wait  
Pyth  
2019  
30 active+recovery\_wait+degraded  
23 active+recovery\_wait+degraded+remapped  
6 active+recovery\_wait  
4 active+recovery\_wait+remapped  
S3 M  
2019  
1 active+recovery\_wait+forced\_recovery+degraded  
1 active+undersized+degraded+remapped+backfilling  
NoPD - SeungHeun, Noh 2:15 PM  
[SHORT NOTICE] I'm visiting TH office  
18th fl.  
# 1C  
verda Shm... 84 replies Last reply today at shm  
oc3660af47534c2bd387d857c4... overlay  
40f79605281e0ab21082ee47/m... overlay  
Ticket 즉 서비스에서 문제가 있는지 확인  
overl... 9416fe51bb9c3e3e83cc025dc9/m... shm  
562e996002a104dd890a0c/F883/m... overlay  
55bd321a9fca2399ad5d4300e/m... tmpfs  
es/kubernetes.io-secret/longer... overlay  
c8257595e3127fb0131d4c9b5/m... tmpfs  
es/kubernetes.io-secret/longer... overlay  
c8257595e3127fb0131d4c9b5/m... tmpfs  
Failed, "target": "https://ticket... shm  
98e056df841750453ed4e5c06a/m... overlay  
e87854b390d4973810d82ca/m... tmpfs  
es/kubernetes.io-secret/cattl... overlay  
e87854b390d4973810d82ca/m... tmpfs  
es/kubernetes.io-secret/cattl... overlay  
9227f661784f314f82b50b0803/m... shm  
8e17460113e1b24cd94dFec3f1e/m... overlay  
out4ba063f7Seed2b41581/m... overlay

# In real world (4)

## Super Outage for 2 hours

- **Problem 1**

All operations should be checked whether they have proper authorizations.

Authorization requires its owner's ACL which is stored in a user meta object in the form of key/value pairs.

Once a user's ACL is accessed, it is stored in a cache to serve following requests.

But in case of an anonymous user, his ACL is not cached, so every requests from anonymous users involve reading ACLs from user meta objects in the default.rgw.users.uid pool every time.

Unfortunately the default.rgw.users.uid pool was on the HDD based pool

- **Cause**

Disk util 100% → filestore\_op\_thread\_suicide\_timeout (180sec) expired → osd kill himself → restart osd

→ rocksdb journal check on startup (it took around 9 mins)

→ gw objecter\_inflight (either 24576 or 100MB contents) threshold was over

→ rgw OP IO threads were throttled and waited until the OP Q was run out.



The screenshot shows a terminal window with several lines of command-line output and a file browser window in the background.

```
~]# ceph osd map default.rgw.users.uid anonymous
osdmap e42025 pool 'default.rgw.users.uid' (10) object 'anonymous' -> pg 10.eac94372 (1
15Z], p8Z)
~]# ceph osd find 82
{
    "osd": 82,
    "ip": "■■■■■ 6809/382106",
    "osd_fsid": "7f7aa037-ce94-4ca1-b026-0e90d68a2d4f",
    "crush_location": {
        "host": "■■■■■",
        "rack": "pb01pb02-data",
        "root": "data"
    }
}
```

The file browser window in the background lists various system services and files, including launchservicesd, LINE, dsAccessService, bluetoothd, NetDrive, Finder Integration, fontd, and suggestd.

# In real world (5)

## Super Outage for 2 hours

- **Problem 2**

LB periodically sends an operation as an anonymous user(GET / HTTP/1.0), which means every request from LB involves reading ACL from objects.

In this situation, one of OSDs which has anonymous user meta objects was stalled to respond.

Because of this, every LB operations were also blocked.

this cause health check failure of LB and LB would not send traffic to real servers which are radosgw.

this cause whole cluster is not operational for some time even though radosgw and osd servers are not busy.

- **Cause**

osd **nodown** flag was set because heavy disk IO may result in osd health check fail (if 2 or more peers report to ceph-mon, then ceph-mon mark the osd down.)

← it is recommended to turn on nodown flag under backfilling state with heavy loaded DISK IO and this prevent osd flapping but... all the osds were waiting for the troubled osd to answer the health check (the osd was bootstrapped, have no response of health check message but it took 10 minutes due to roscksdb was on HDD and it was pressed by disk IO)

# In real world (6)

## Patch & Tune

- anonymous user → pass authorization procedure
- move bucket meta pools to SSD based pool
  - `ceph osd pool set default.rgw.users.uid crush_rule SSD`
- tune osd thread\_suicide\_timeout <https://access.redhat.com/solutions/2971581>

```
osd_op_thread_timeout = 90 #default is 15
```

```
osd_op_thread_suicide_timeout = 2000 #default is 150
```

```
osd_backfill_scan_max = 1 # default 512
```

```
osd_backfill_scan_min = 1 # default 64
```

```
osd_peering_wq_threads = 4 # default 2
```

- If recovery thread is hitting the timeout, increase recovery op threads timeout

```
osd_recovery_thread_timeout = 120 #default is 30
```

```
osd_recovery_thread_suicide_timeout = 2000 #default 300
```

- If filestore op threads are hitting the timeout

```
filestore_op_thread_timeout = 180 #default is 60
```

```
filestore_op_thread_suicide_timeout = 2000 #default is 180
```

# Outage Prevention

- slow down backfilling speed & control backfilling
  - set osd\_max\_backfills = 1
  - if slow\_requests are found, then stop backfilling (set osd\_max\_backfills = 0)
- make HDD based osd secondary or tertiary if possible
  - ceph osd primary-affinity osd.\$id 0
- Prevent osd flapping & prepare disk failure
  - ceph osd set nodown
  - if osd commit suicide due to hight disk utilization or osd die due to HW failure, then 'ceph osd unset nodown'
- gently update crush reweight of new osd by 0.01
  - <https://github.com/cernceph/ceph-scripts/blob/master/tools/ceph-gentle-reweight>
- patch anonymous user not to check authorization
  - <https://github.com/ceph/ceph/pull/34287>

# Final results - 10 months later

	<b>Before</b>	<b>After</b>
<b>Disks and Servers</b>	HDD 8TB * 12ea * 12 servers	SSD 3.8TB * 24ea * 30 servers
<b>SIZE</b>	raw 1PB	raw 2.5PB
<b>OSDs</b>	144	720
<b>Objects</b>	5 billions	6 billions
<b>Placement groups</b>	8192	16384
<b>Traffic</b>	60 millions / day	140 millions / day
<b>CDN hit</b>	2 billions / day	3.5 billions / day
<b>Rebalance</b>	8 hours	2 hours
<b>Recovery</b>	4.5 days	8 hours
<b>Restart osd</b>	20 minutes	2 minutes

# Lesson and Learn - capacity expansion, migration, ...

**create new crush root for new osds**

```
ceph osd crush rule create-replicated <RULE_NAME> <ROOT> <FAILURE_DOMAIN> <DEVICE_CLASS>
```



**move new osds to crush root or deploy osds with new crush root**

```
ceph osd crush move <FAILURE_DOMAIN> root=<ROOT>
```



**set crush reweight 0 of new osds**

```
ceph osd crush reweight-subtree <ROOT> 0.0
```



**move new osds to service crush root**

```
ceph osd crush move <BUCKET> <FAILURE_DOMAIN>
```



**update crush weight of new osds by 0.01**

<https://github.com/cernceph/ceph-scripts/blob/master/tools/ceph-gentle-reweight>

THANK YOU

Q & A

dongsoo.song@linecorp.com