FUJITSU

shaping tomorrow with you

Igor Podoski / Zofia Domaradzka

# Ghost Cluster and Throwing Fireballs

# Agenda

- Ghost Cluster

- Throwing Fireballs

- Questions??

# ETERNUS CD10000

■ A Fujitsu software-defined storage system based on Ceph and RHEL7

■ Appliance fully integrated with and automated on Fujitsu Primergy Servers

■ Provides custom tools for installation, configuration, monitoring, diagnostics etc.

# ETERNUS CD10000 - monitoring

- Custom monitoring system using CD10000 snmp agents

- Monitoring of PGs, OSDs, monitors and overall cluster state

- Active polling and traps

- Responsiveness for a large cluster must be tested (e.g. cluster with 224 nodes)

# Testing monitoring system responsiveness

Challenges:

Testing for different ceph configurations and cluster sizes, e.g.:
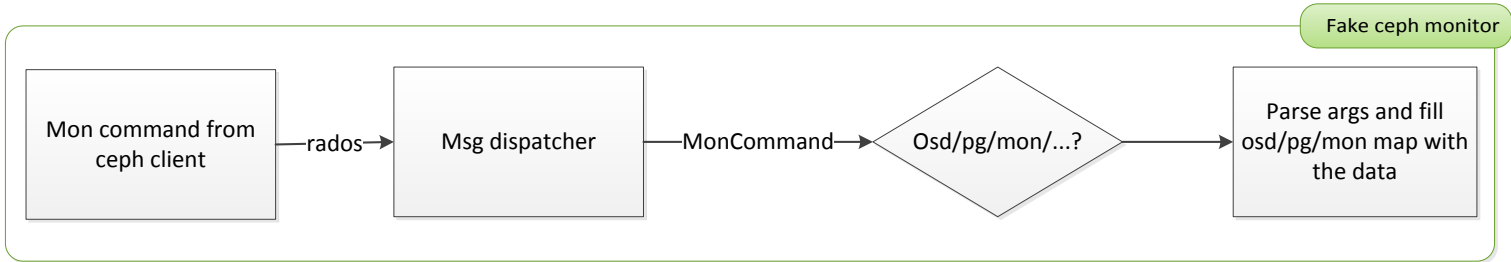
- Large number of PGs

- Cluster at full/near full state

- Change of state of specific MOs at given moment or several states at once

# Solution
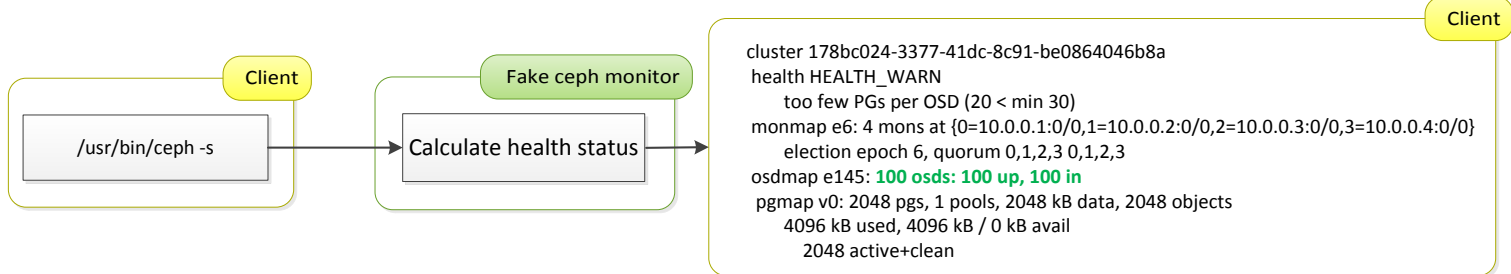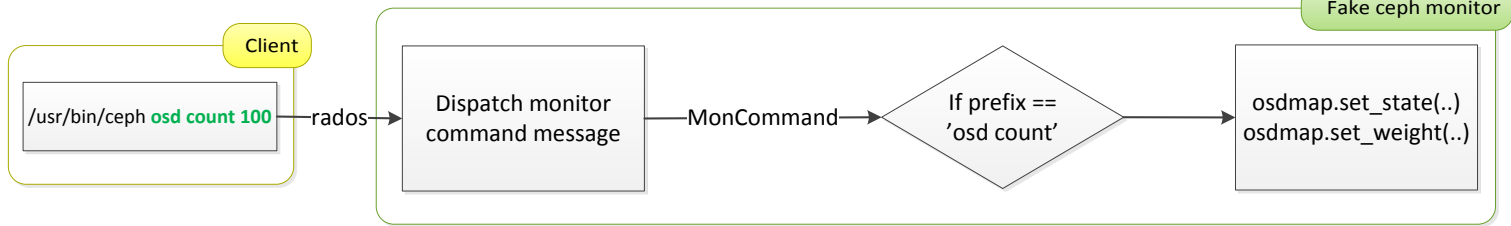
FUJITSU

# GHOST CLUSTER

# How does it work?

**FUJITSU**

- Work in progress stand alone process simulating ceph monitor

- Based on the fragments of ceph monitor code

- Uses RADOS protocol, so it is compatible with all Ceph clients (/usr/bin/ceph etc...)

- Uses MON / OSD / PG maps for storing fake objects in the process memory

# How does it work?



**Example**: Add 100 fake OSDs and get Ghost Cluster status

# Current functionality

- **MONmap manipulations:**
    - add/remove monitors
    - change quorum status
- **OSDmap manipulations:**
    - change number of OSDs
    - change state of each OSD (full, nearfull, etc.)
- **PGmap manipulations:**
    - change numer of PGs
    - change state of each PG (active+clean, degraded, etc.)
- **Overall cluster status manipulations:**
    - change cluster flags (full, etc.)

# Benefits

■ Time and resource saver for:

  ■ reconfiguring and filling up real Ceph cluster

  ■ no need to use physical cluster, single process can be run on any virtual environment

■ Flexible configuration and resposiveness:

  ■ change of parameters on the fly

  ■ immediate cluser state response

■ Easier automation of test scenarios:

  ■ predefined configuration profiles can be used

  ■ state transision can be also emulated e.g. long PGs recovery time

# Demo

- Fake Ceph monitor was started on localhost
- Currently PG / MON / OSD maps are empty:

```
[root@localhost build]# ceph -s
    cluster 1bb821e7-4550-4f1b-baec-259e2809261a
     health HEALTH_ERR
            no osds
     monmap e0: 0 mons at {}
            election epoch 0, quorum
     osdmap e1: 0 osds: 0 up, 0 in
      pgmap v0: 0 pgs, 0 pools, 0 bytes data, 0 objects
            0 kB used, 0 kB / 0 kB avail
```

# Demo

**FUJITSU**

■ Let's have a look at allowed options:

```
[root@localhost build]# ceph –h
...
mon add <name> <IPaddr[:port]>                          add new monitor with <name> ip <ip:[port]>
mon quorum <quorum> [<quorum>...]                       set quorum <0 1 2>
mon rm <name>                                           remove monitor <name>
mon skew <int[0-]>                                      set skew <seconds>
osd <int[0-]> <state>                                   set <state> on osd <num>
osd count <int[0-]>                                     set <num> osds
osd set <flag>                                          set <flag> on osdmap
osd unset <flag>                                        unset <flag> on osdmap
osd usage <int> <int[0-]> <int[0-]> <int[0-]>           add usage to <num> osd with <kb> <kb_used> <kb_avail>
pg count <int[0-]> <int[0-]> <int[0-]> <int[0-]>        set <pool> <size> <obj> <num> pgs
status                                                  show cluster status
```

# Demo

- **Set predefined profile:**

```
#!/bin/bash
ceph mon add 0 10.0.0.1
ceph mon add 1 10.0.0.2
ceph mon add 2 10.0.0.3
ceph mon quorum 0 1 2


ceph osd count 20


ceph pg count 0 1024 1 2048


ceph osd usage 0 0 1024 1024
ceph osd usage 1 0 1024 1024
ceph osd usage 2 0 1024 1024
ceph osd usage 3 0 1024 1024
```

- **Get ghost cluster status:**

```
[root@localhost build]# ceph -s
    cluster 1bb821e7-4550-4f1b-baec-259e2809261a
     health HEALTH_OK
     monmap e4: 3 mons at {0=10.0.0.1:0/0,1=10.0.0.2:0/0,2=10.0.0.3:0/0}
            election epoch 4, quorum 0,1,2 0,1,2
     osdmap e21: 20 osds: 20 up, 20 in
      pgmap v0: 2048 pgs, 1 pools, 2048 kB data, 2048 objects
            4096 kB used, 4096 kB / 0 kB avail
                2048 active+clean
```

# Demo

■ Reduce quorum to mon.0 and mon.1

```
[root@localhost build]# ceph mon quorum 0 1


[root@localhost build]# ceph -s
    cluster 1bb821e7-4550-4f1b-baec-259e2809261a
     health HEALTH_WARN
            1 mons down, quorum 0,1 0,1
     monmap e5: 3 mons at {0=10.0.0.1:0/0,1=10.0.0.2:0/0,2=10.0.0.3:0/0}
            election epoch 5, quorum 0,1 0,1
     osdmap e21: 20 osds: 20 up, 20 in
      pgmap v0: 2048 pgs, 1 pools, 2048 kB data, 2048 objects
            4096 kB used, 4096 kB / 0 kB avail
                2048 active+clean
```

# Demo

## ■ Emulate OSD.0 down

```
[root@localhost build]# ceph osd 0 down

[root@localhost build]# ceph -s
    cluster 1bb821e7-4550-4f1b-baec-259e2809261a
     health HEALTH_WARN
            1 mons down, quorum 0,1 0,1
            1/20 in osds are down
    monmap e5: 3 mons at {0=10.0.0.1:0/0,1=10.0.0.2:0/0,2=10.0.0.3:0/0}
            election epoch 5, quorum 0,1 0,1
    osdmap e23: 20 osds: 19 up, 20 in
     pgmap v0: 2048 pgs, 1 pools, 2048 kB data, 2048 objects
            4096 kB used, 4096 kB / 0 kB avail
                2048 active+clean
```

# Demo

**FUJITSU**

## ■ Add monitor clock skew

```
[root@localhost build]# ceph mon skew 2

[root@localhost build]# ceph -s
cluster 1bb821e7-4550-4f1b-baec-259e2809261a
    health HEALTH_WARN
            clock skew detected on mon.0
            1 mons down, quorum 0,1 0,1
            1/20 in osds are down
            Monitor clock skew detected
    monmap e5: 3 mons at {0=10.0.0.1:0/0,1=10.0.0.2:0/0,2=10.0.0.3:0/0}
            election epoch 5, quorum 0,1 0,1
    osdmap e23: 20 osds: 19 up, 20 in
     pgmap v0: 2048 pgs, 1 pools, 2048 kB data, 2048 objects
            4096 kB used, 4096 kB / 0 kB avail
                2048 active+clean
```

# Demo

**FUJITSU**

## ■ Emulate nearfull OSD.1

```
[root@localhost build]# ceph osd 1 nearfull

[root@localhost build]# ceph -s
cluster 1bb821e7-4550-4f1b-baec-259e2809261a
    health HEALTH_WARN
            clock skew detected on mon.0
            1 near full osd(s)
            1 mons down, quorum 0,1 0,1
            1/20 in osds are down
            Monitor clock skew detected
    monmap e5: 3 mons at {0=10.0.0.1:0/0,1=10.0.0.2:0/0,2=10.0.0.3:0/0}
            election epoch 5, quorum 0,1 0,1
    osdmap e24: 20 osds: 19 up, 20 in
     pgmap v0: 2048 pgs, 1 pools, 2048 kB data, 2048 objects
            4096 kB used, 4096 kB / 0 kB avail
                2048 active+clean
```

# Demo

**FUJITSU**

## ■ Set cluster full flag

```
[root@localhost build]# ceph osd set full

[root@localhost build]# ceph -s
cluster 1bb821e7-4550-4f1b-baec-259e2809261a
     health HEALTH_WARN
             clock skew detected on mon.0
             1 near full osd(s)
             1 mons down, quorum 0,1 0,1
             1/20 in osds are down
             Monitor clock skew detected
     monmap e5: 3 mons at {0=10.0.0.1:0/0,1=10.0.0.2:0/0,2=10.0.0.3:0/0}
             election epoch 5, quorum 0,1 0,1
     osdmap e24: 20 osds: 19 up, 20 in
      pgmap v0: 2048 pgs, 1 pools, 2048 kB data, 2048 objects
             4096 kB used, 4096 kB / 0 kB avail
                 2048 active+clean
```

# Demo

## ■ Emulate OSD.2 full

```
[root@localhost build]# ceph osd 2 full
[root@localhost build]# ceph -s
    cluster 1bb821e7-4550-4f1b-baec-259e2809261a
     health HEALTH_ERR
            clock skew detected on mon.0
            1 full osd(s)
            1 near full osd(s)
            1 mons down, quorum 0,1 0,1
            1/20 in osds are down
            full flag(s) set
            Monitor clock skew detected
     monmap e9: 3 mons at {0=10.0.0.1:0/0,1=10.0.0.2:0/0,2=10.0.0.3:0/0}
            election epoch 9, quorum 0,1 0,1
     osdmap e113: 20 osds: 19 up, 20 in
            flags full
      pgmap v0: 2048 pgs, 1 pools, 2048 kB data, 2048 objects
            4096 kB used, 4096 kB / 0 kB avail
                2048 active+clean
```

# Throwing fireballs

# What it is?

Throwing fireballs into Ceph means:

- Break stuff e.g.
  - Add 10% packet drop to public interface for node with mon0
  - Add 100ms network delay to cluster interface on different node
  - Pin all ceph-osd processes from node with mon1 to one logical CPU core
  - Move all ms_dispatch threads from all ceph-osds on node without monitors to one logical CPU core
  - Misconfigure OSD parameters in resobanble way
  - Filling up OSD partitions with non PG stuff (e.g. using dd)

- See Ceph reaction:
  - When / where / how it breaks

- Create a cure for newly created dissease:
  - Analyse ceph logs and potential core dumps
  - Deduce probablity of newly created conditions and prepare a solution

# Tools for throwing fireballs

- Ceph configuration poisoning:
  - Injecting args at runtime
  - Permanent changes in ceph.conf
- System tools:
  - tuned, tc, /proc files, iptables, changing XFS properties, etc.
- Dedicated tools:
  - Newly created CPM (Ceph Process Manager)
  - Dedicated scripts and code snippets

# CPM - Ceph Process Manager

- Uses python and salt to interact with Ceph cluster

- Manages Ceph processes at higher level:
  - Doesn't matter on which node ceph-* are running
  - Keeps configuration in a flat JSON file
  - Uses regular expressions to match process and thread names

- Can tune several things (for processes and individual threads)
  - Set any logical CPUs on which can run
  - Change nice priority of processes and threads
  - Change scheduling and real-time priority
  - Change I/O scheduling policy and priority

- Uses python and custom salt module
  - It will be released soon

# CPM demo

CPM can be started with GUI or in batch mode only.

# CPM demo

Processes can be filtred by reglar expressions.

```
──────────────────────────── Ceph Process Manager ────────────────────────────
{p,P} -- list all Ceph processes              {m,M} -- toggle MONs
enter -- show process details                 {d,D} -- toggle MDS
{q,Q} -- quit                                 {o,O} -- toggle OSDs

──────────────────────────────── Filter processes ────────────────────────────
Filter: ceph-osd-2.              < Create configuration for filter            >

MON: ON                    OSD: ON                    MDS: ON

   24584 ceph-osd-28
   24681 ceph-osd-29
   70447 ceph-osd-20
   70477 ceph-osd-21
   70481 ceph-osd-22
   70498 ceph-osd-23
   70578 ceph-osd-26
   70580 ceph-osd-25
   70581 ceph-osd-24
   70583 ceph-osd-27




< Apply currently saved settings!                                            >
```

# CPM demo

In process view serveral options can be chosen.

Settings will be saved in JSON format.

# CPM demo

- JSON config example:
  - Pin every osd process on the whole cluster to logical cpu core 0 and 1
  - Change will be made on all nodes in the cluster

```json
{
    "ceph-osd-*": {
        "scheduling": {
            "policy": "OTHER",
            "priority": 0
        },
        "ionice": {
            "policy": "REAL_TIME",
            "priority": 4
        },
        "enable": {
            "io_sched": false,
            "sched": false,
            "cpu": true
        },
        "taskset": [0,1],
        "thread_name": "",
        "nice": 0
    }
}
```

# CPM demo

- JSON config example:
  - Move ms_dispatch thread for every ceph-osd process to logical cpu cores:  3,4,6,18
  - Change will be made on all nodes in the cluster

- To apply JSON profile:

```
> python cpm.py --apply profile.json
```

```
{
    "ceph-osd-*": {
        "scheduling": {
            "policy": "OTHER",
            "priority": 0
        },
        "ionice": {
            "policy": "REAL_TIME",
            "priority": 4
        },
        "enable": {
            "io_sched": false,
            "sched": false,
            "cpu": true
        },
        "taskset": [3,4,6,18],
        "thread_name": "ms_dispatch",
        "nice": 0
    }
}
```

# Throwing fireballs in the wild

- How to present this technique?

- Is there a way to:
  - Make it more real than just flat files and terminal commands?
  - Move it to different level of abstraction?
  - Make it more fun?

- Blender comes for the rescue!
  - Game-like interface for throwing fireballs
  - Realtime logs and Ceph status on HUD display
  - True interaction with physical servers
  - Interaction through librados and salt

# Let's play!

**FUJITSU**

- **Controls:**
  - Mouse look
  - W S A D keyboard for movement
- **Graphics:**
  - 3D models of ETERNUS CD10000 appliance

# Heads-Up Display

- **Left**
  - OSD count, UP vs IN
  - PG count
  - PG states
- **Center**
  - Cluster usage in GB
- **Right**
  - Health status
  - Health summary

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

11.86 GB used, 46670.96 / 46682.82 GB avail

HEALTH_OK

ETERNUS CD10000

ETERNUS CD10000

ETERNUS CD10000

# Monitor log wall

Realtime update from
monitor log callback
(python)

# What is inisde?

Starting from top:
- public network switch

- node4

- node3

- node2

- node1

- cluster network switch

- *management node

- *admin network switch

* Management node and
admin network is an
additional part of
ETERNUS CD10000
appliance.



OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

11.85 GB used, 46670.97 / 46682.82 GB avail

HEALTH_OK

D10000

# Every object has its own menu

Let's start rados bench from node1.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

11.85 GB used, 46670.97 / 46682.82 GB avail          HEALTH_OK

[ node1 ]
bench: 4M 10 sec wr
bench: 4M 60 sec wr
bench: 4k 10 sec wr
bench: 4k 60 sec wr
bench: 4k 60 sec wr no cleanup
bench: 60 sec rand
bench: 60 sec seq
bench: 4k 60 sec 16t wr no cleanup
bench: 60 sec 16t rand
bench: 60 sec 16t seq
bench: cleanup last
bench: cleanup all

# Logical CPU usage 0-31

Orange blocks are scaling from 0 to 100% just as logical core usage on server after starting rados bench test.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

11.92 GB used, 46670.90 / 46682.82 GB avail
write: 8.47 MB/s
2168 w/s
0 r/s

HEALTH_OK

# Rados bench wall

Current rados bench results are displayed on another wall of 3D server room.

# Let's look on the back side



Network cabling:

Blue – public

Yellow  – cluster

Red – admin

# First fireball – public network delay

Add 100ms delay for public interface of node2.



OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

15.10 GB used, 46667.72 / 46682.82 GB avail

HEALTH_OK

[ node2_public ]
tc: 10% data loss
tc: 100% data loss
○ tc: 100ms delay
tc: 200ms delay
tc: clear
tc: list

# First fireball – public network delay

Add another delay, this time 200ms for public interface of node3.



OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

15.11 GB used, 46667.72 / 46682.82 GB avail

HEALTH_OK

[ node3_public ]
tc: 10% data loss
tc: 100% data loss
tc: 100ms delay
○ tc: 200ms delay
tc: clear
tc: list

# First fireball – public network delay

Have a global look on what we've set on the public network switch.



OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

15.11 GB used, 46667.72 / 46682.82 GB avail

HEALTH_OK

[ switch_public ]
tc: 10% data loss
tc: 25% data loss
tc: 50% data loss
tc: 100% data loss
tc: 100ms delay
tc: 200ms delay
tc: 500ms delay
tc: clear
tc: list

# First fireball – public network delay

Public network delays:

node3: 200ms

node2: 100ms



OSD: 56, UP: 56, IN: 56        15.11 GB used, 46667.72 / 46682.82 GB avail        HEALTH_OK
PGS: 1866
active+clean: 1866

node1:qdisc mq 0: root
node3:qdisc netem 8047: root refcnt 65 limit 1000 delay 200.0ms
node2:qdisc netem 8054: root refcnt 65 limit 1000 delay 100.0ms
node4:qdisc mq 0: root

# First fireball – public network delay

Since we slighty broke the public network, let's run a read test.



```
OSD: 56, UP: 56, IN: 56          15.11 GB used, 46667.72 / 46682.82 GB avail    HEALTH_OK
PGS: 1866
active+clean: 1866

                              [ node1 ]
                    bench: 4M 10 sec wr
                    bench: 4M 60 sec wr
                    bench: 4k 10 sec wr
                    bench: 4k 60 sec wr
                    bench: 4k 60 sec wr no cleanup
                    bench: 60 sec rand
                    bench: 60 sec seq
                    bench: 4k 60 sec 16t wr no cleanup
                  o bench: 60 sec 16t rand
                    bench: 60 sec 16t seq
                    bench: cleanup last
                    bench: cleanup all
```

# First fireball – public network delay



Look closely on the latency, which sometimes is very low, but sometimes reaches above 100ms and 200ms.

These are the values we have set as delays of node2 and node3.

# Second fireball – kill one of cluster network interfaces

Start rados bench first to see what is going to happen after we kill one network card.



OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

15.56 GB used, 46667.26 / 46682.82 GB avail
write: 8.64 MB/s
2210 w/s
0 r/s

HEALTH_OK

mon.0 [INF] pgma
mon.0 [INF] pgma
mon.0 [INF] pgma
mon.0 [INF] pgma

[ node1 ]
bench: 4M 10 sec wr
bench: 4M 60 sec wr
bench: 4k 10 sec wr
bench: 4k 60 sec wr
bench: 4k 60 sec wr no cleanup
bench: 60 sec rand
bench: 60 sec seq
bench: 4k 60 sec 16t wr no cleanup
bench: 60 sec 16t rand
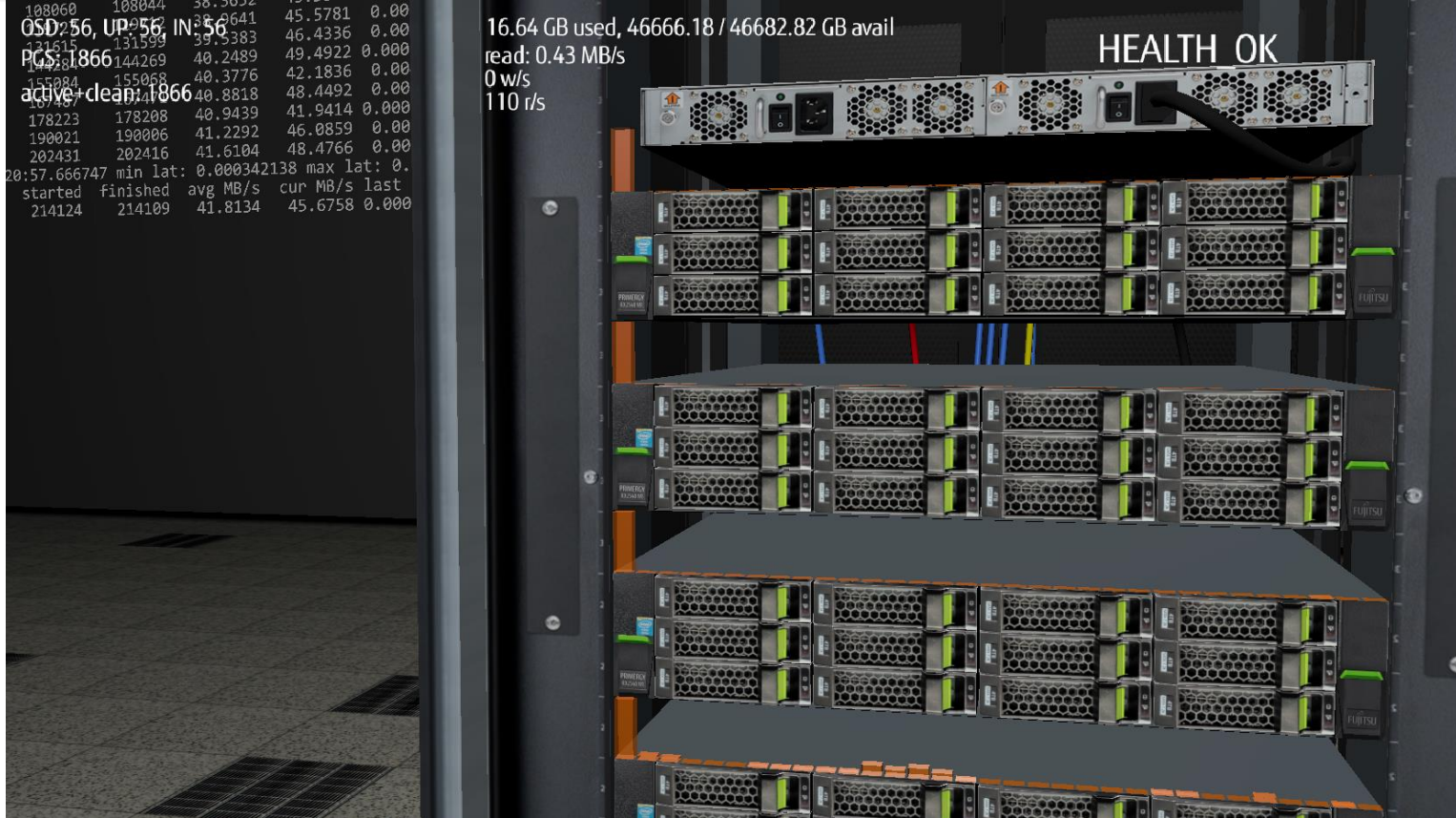bench: 60 sec 16t seq
bench: cleanup last
bench: cleanup all

# Second fireball – kill one of cluster network interfaces



Adding 100% data loss to a network interface could simulate e.g. NIC overheat.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

15.90 GB used, 46666.93 / 46682.82 GB avail
write: 8.48 MB/s
2170 w/s
0 r/s

HEALTH_OK

[ node2_cluster ]
tc: 10% data loss
tc: 100% data loss
tc: 100ms delay
tc: clear
tc: list

# Second fireball – kill one of cluster network interfaces

Writes are now blocked, because there is no communication via cluster network interface of node2.

# Second fireball – kill one of cluster network interfaces

Ok, it is enough, clear 100% packet drop on this interface and let's check if writes to the cluster will start working again.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

16.41 GB used, 46666.41 / 46682.82 GB avail

HEALTH_OK

[ node2_cluster ]
tc: 10% data loss
tc: 100% data loss
tc: 100ms delay
tc: clear
tc: list

# Second fireball – kill one of cluster network interfaces

At the end of rados bench log, you can see that writes were unblocked.

Cluster went into the HEALTH_WARN state, and reported slow requests.

Probably some threads were constantly waiting for the cluster network response to finish write operation.

# Third fireball – CPU time



Have a look on tuned profiles.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

16.64 GB used, 46666.18 / 46682.82 GB avail

HEALTH_OK

[ pmgmt ]
tuned: latency-performance
tuned: throughput-performance
tuned: network-latency
tuned: network-throughput
tuned: balanced
tuned: powersave
tuned: show active profiles
cpm: set all osd on one cpu
cpm: set all osd on all cpus
cpm: set all osd sched ilde
cpm: split thread groups
cpm: show thread split

# Third fireball – CPU time

The settings are telling us that the cluster is profiled to achieve max performance.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

16.64 GB used, 46666.18 / 46682.82 GB avail

HEALTH_OK

node1:Current active profile: throughput-performance
node3:Current active profile: throughput-performance
node2:Current active profile: throughput-performance
node4:Current active profile: throughput-performance

# Third fireball – CPU time

Using CPM (our newly developed tool), pin all ceph-osd processes to first logical CPU core on every node in the cluster.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

16.64 GB used, 46666.18 / 46682.82 GB avail

HEALTH_OK

[ pmgmt ]
tuned: latency-performance
tuned: throughput-performance
tuned: network-latency
tuned: network-throughput
tuned: balanced
tuned: powersave
tuned: show active profiles
cpm: set all osd on one cpu
cpm: set all osd on all cpus
cpm: set all osd sched ilde
cpm: split thread groups
cpm: show thread split

# Third fireball – CPU time



Start rados bench test with16 threads to see cpu usage.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

Average latency(s): 0.0152345
0.426569
27.3506
0.0012833

16.64 GB used, 46666.18 / 46682.82 GB avail

HEALTH_OK

[ node1 ]
bench: 4M 10 sec wr
bench: 4M 60 sec wr
bench: 4k 10 sec wr
bench: 4k 60 sec wr
bench: 4k 60 sec wr no cleanup
bench: 60 sec rand
bench: 60 sec seq
bench: 4k 60 sec 16t wr no cleanup
bench: 60 sec 16t rand
bench: 60 sec 16t seq
bench: cleanup last
bench: cleanup all

# Third fireball – CPU time



OSDs are allowed only to use first logical CPU core, bandwith dropped twice.

# Third fireball – CPU time

Allow each OSD to use every logical core.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

16.64 GB used, 46666.18 / 46682.82 GB avail
read: 0.67 MB/s
0 w/s
171 r/s

HEALTH_OK

[ pmgmt ]
tuned: latency-performance
tuned: throughput-performance
tuned: network-latency
tuned: network-throughput
tuned: balanced
tuned: powersave
tuned: show active profiles
cpm: set all osd on one cpu
cpm: set all osd on all cpus
cpm: set all osd sched ilde
cpm: split thread groups
cpm: show thread split

# Third fireball – CPU time

Check rados bench wall – bandwith jumped back to normal level.

OSDs are now using every core (0-31).

# Third fireball – CPU time



OSD has about 28 uniqe thread names.

Let's pin each OSD thread name to each of logical cpu cores (0-27).

Thanks to this we could easily check which thread groups need more CPU time than the others.

Unique thread names for OSD:

admin_socket, ceph-osd, filestore_sync, fn_anonymous, fn_appl_fstore, fn_jrn_objstore, fn_odsk_fstore, journal_write, journal_wrt_fin, log, ms_accepter, ms_dispatch, ms_local, ms_pipe_read, ms_pipe_write, ms_reaper, osd_srv_agent, osd_srv_heartbt, safe_timer, service, sginal_handler, tp_fstore_op, tp_osd, tp_osd_cmd, tp_osd_disk, tp_osd_recov, tp_osd_tp, wb_throttle

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

16.64 GB used, 46666.18 / 46682.82 GB avail
read: 0.87 MB/s
0 w/s
222 r/s

HEALTH_OK

[ pmgmt ]
tuned: latency-performance
tuned: throughput-performance
tuned: network-latency
tuned: network-throughput
tuned: balanced
tuned: powersave
tuned: show active profiles
cpm: set all osd on one cpu
cpm: set all osd on all cpus
cpm: set all osd sched ilde
cpm: split thread groups
cpm: show thread split

55

# Third fireball – CPU time



The ones with 100% usage are:
- ms_accepter
- tp_osd_tp

# Multiple fireballs at once

Add 25% data loss to every NIC connected to cluster switch.

# Multiple fireballs at once



Start rados bench write to see cluster reaction.

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

16.65 GB used, 46666.17 / 46682.82 GB avail
read: 1.02 MB/s
0 w/s
261 r/s

HEALTH_OK

[ node1 ]
bench: 4M 10 sec wr
bench: 4M 60 sec wr
bench: 4k 10 sec wr
bench: 4k 60 sec wr
bench: 4k 60 sec wr no cleanup
bench: 60 sec rand
bench: 60 sec seq
bench: 4k 60 sec 16t wr no cleanup
bench: 60 sec 16t rand
bench: 60 sec 16t seq
bench: cleanup last
bench: cleanup all



58

# Multiple fireballs at once

As we suspected, writes are starting to be unstable.

# Multiple fireballs at once

Stop OSD.0 gracefully, this should create small rebalance.

# Multiple fireballs at once

Writes were blocked again.

# Multiple fireballs at once



Allow OSDs only to use single logical CPU core.

OSD: 56, UP: 55, IN: 56
PGS: 1866

active+clean: 1737
active+undersized+degraded: 118
active+clean+scrubbing: 6
activating+undersized+degraded: 3
peering: 1
stale+active+clean: 1

17.30 GB used, 46665.52 / 46682.82 GB avail
write: 0.00 MB/s
0 w/s
0 r/s

[ pmgmt ]
tuned: latency-performance
tuned: throughput-performance
tuned: network-latency
tuned: network-throughput
tuned: balanced
tuned: powersave
tuned: show active profiles
cpm: set all osd on one cpu
cpm: set all osd on all cpus
cpm: set all osd sched ilde
cpm: split thread groups
cpm: show thread split

HEALTH_WARN

121 pgs degraded
1 pgs peering
1 pgs stale
121 pgs undersized
13 requests are blocked > 32 sec
recovery 14296/682920 objects degraded (2.093%
1/56 in osds are down

# Multiple fireballs at once

Writes are still blocked, and since OSDs are working only on one cpu, recovery process now runs slower.

# Multiple fireballs at once

Let's check if cluster will recover after cluster network switch will be healed.

# Multiple fireballs at once

Writes are still blocked, but there is some bigger movement on cluster.

# Multiple fireballs at once

Cluster starts healing.

# Multiple fireballs at once

This time Ceph won!

OSD: 56, UP: 56, IN: 56
PGS: 1866
active+clean: 1866

17.34 GB used, 46665.48 / 46682.82 GB avail

HEALTH_OK

# **Thank you for your attention** ☺

Igor.Podoski@ts.fujitsu.com

Zofia.Domaradzka@ts.fujitsu.com