



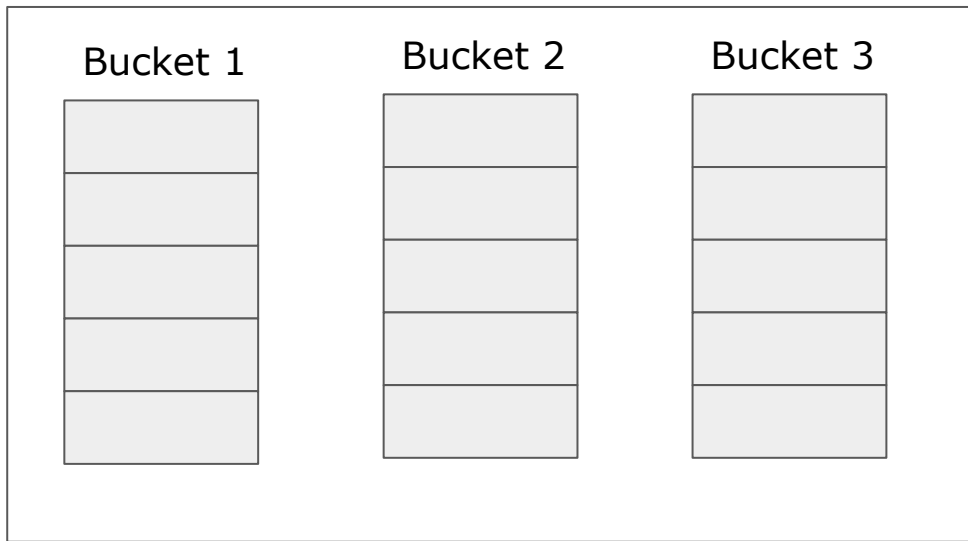
Интенсив СЕРН

5ти-дневный интенсив
День №3

Программа занятия

- Различные “крутилки” в serf
- Связка с Openstack и варианты деплоя
- Ускорение производительности (разные параметры messenger, stripe, malloc)
- Мониторинг состояния
- Вывод ноды из кластера. Действия при падении ноды
- Ознакомление с авторизацией в CEPH
- Начальный дебаг проблем

Bucket type = rack



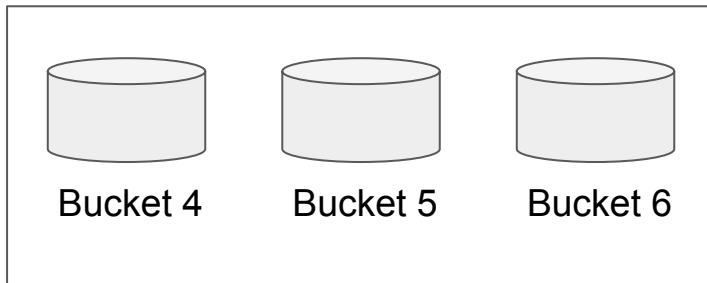
pool pool_name
replication_size=3



FD=rack

pool pool_name2
replication_size=3

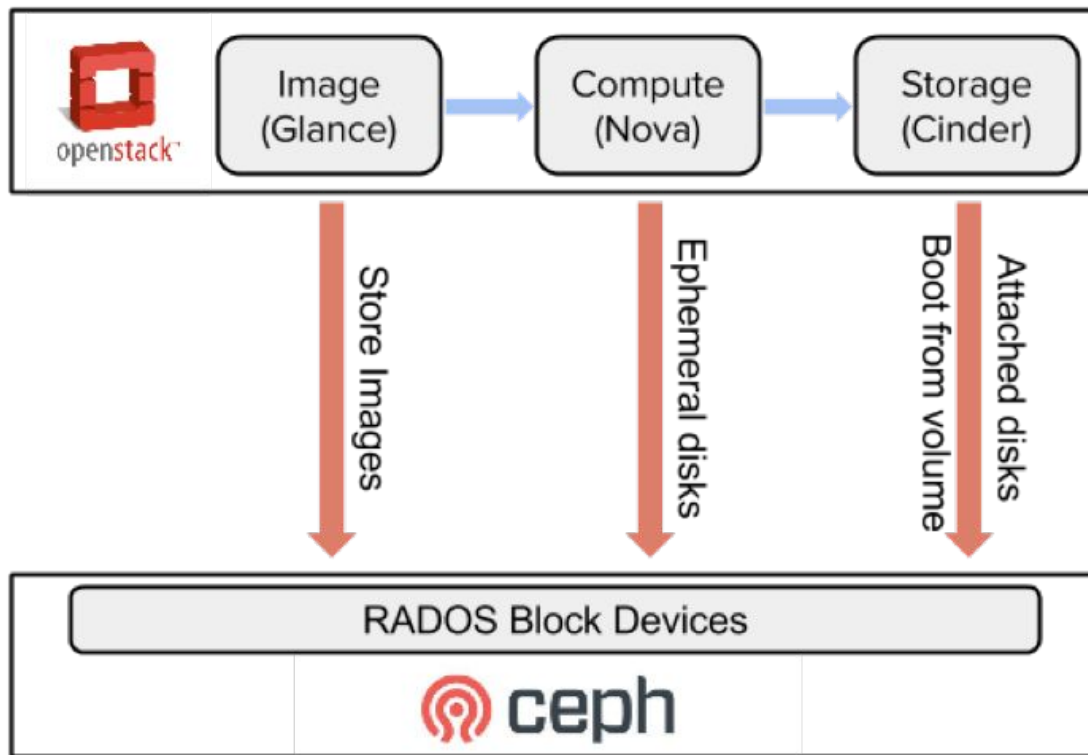
Bucket type = osd



FD=osd



OpenStack + Ceph



OpenStack + Ceph

- 2 разных кластера не удобно
- 1 зона OpenStack - 1 Pool CEPH
- Ограничить IOPS. base image 200 IOPS block device 500 IOPS
- Снэпшоты != бекапы

Алгоритмы распределения

Uniform – все веса строго одинаковы. Подходит, когда кластер состоит из совершенно одинаковых машин и дисков

List – перемещаемые данные с некоторой вероятностью попадают в новое или старое хранилище. Expanding cluster

Tree – бинарные деревья, оптимизация скорости помещения объектов в хранилище

Straw – комбинация стратегий List и Tree для реализации принципа «разделяй и властвуй». Обеспечивает быстрое размещение, но иногда создает проблемы для реорганизации

"Крутки" scrub

osd_scrub_chunk_max = 1

osd_scrub_priority = 1

osd_scrub_begin_hour = 0

osd_scrub_end_hour = 4

osd_scrub_min_interval = 86400

osd_scrub_chunk_min = 1

osd_scrub_max_interval = 2419200.0

osd_scrub_sleep = 0.1

osd_deep_scrub_interval = 2419200.0

Крутилки кешей

`rbd_cache_max_dirty = 33554432`

`rbd_cache_size = 67108864`

`rbd_cache_max_dirty_age = 5`

`rbd_cache = True`

RBD object map feature

layering

striping

exclusive-lock

object-map

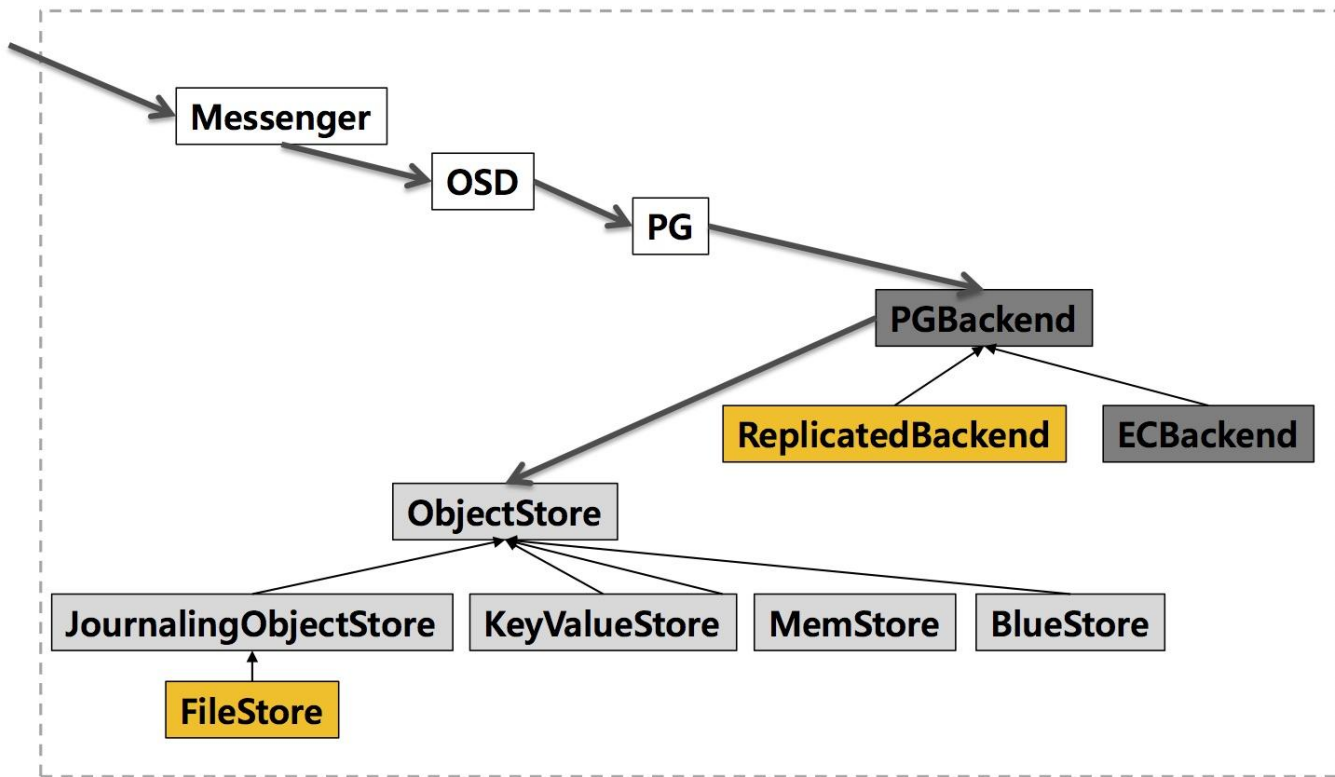
Включение через `rbid_default_features = num`

Биты для установки rbd feature

В параметрах `rbd_default_features` указывается сумма бит.

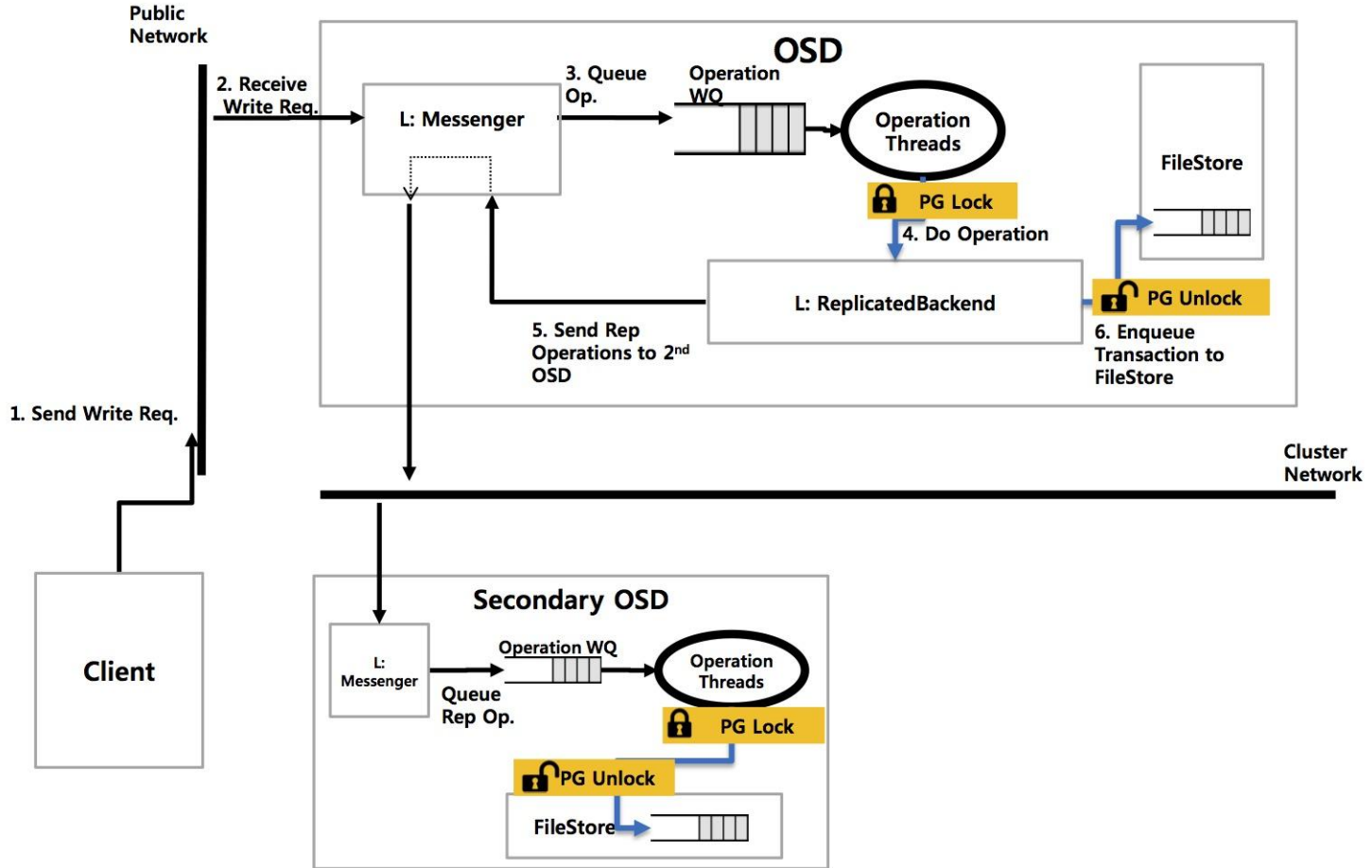
属性	BIT码
layering	1
striping	2
exclusive-lock	4
object-map	8
fast-diff	16
deep-flatten	32

Ceph IO Flow in OSD

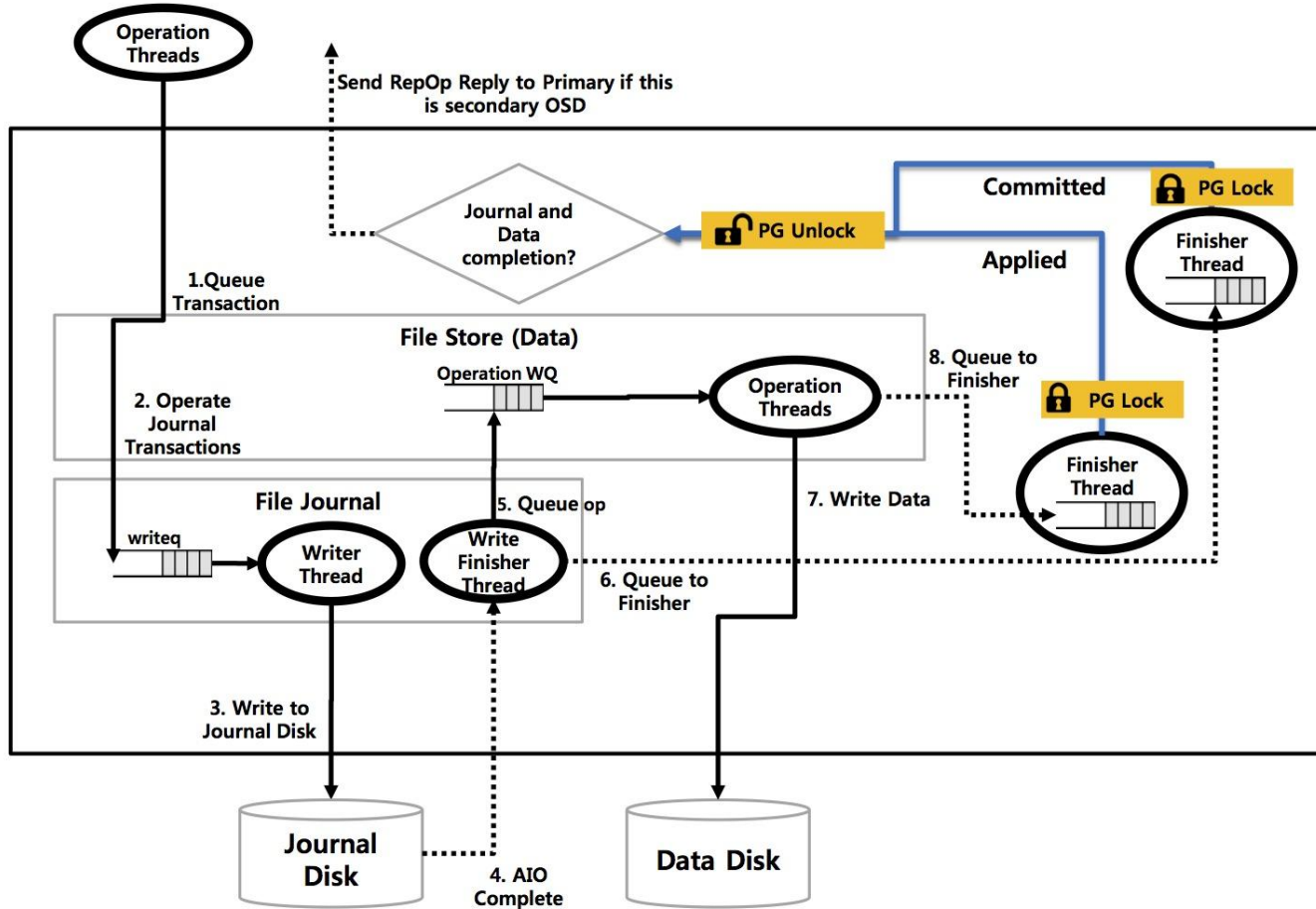


1. Journal: LIBAIO (O_DIRECT && O_DSYNC) → Committed
2. Data: Buffered IO and syncfs() later → Applied

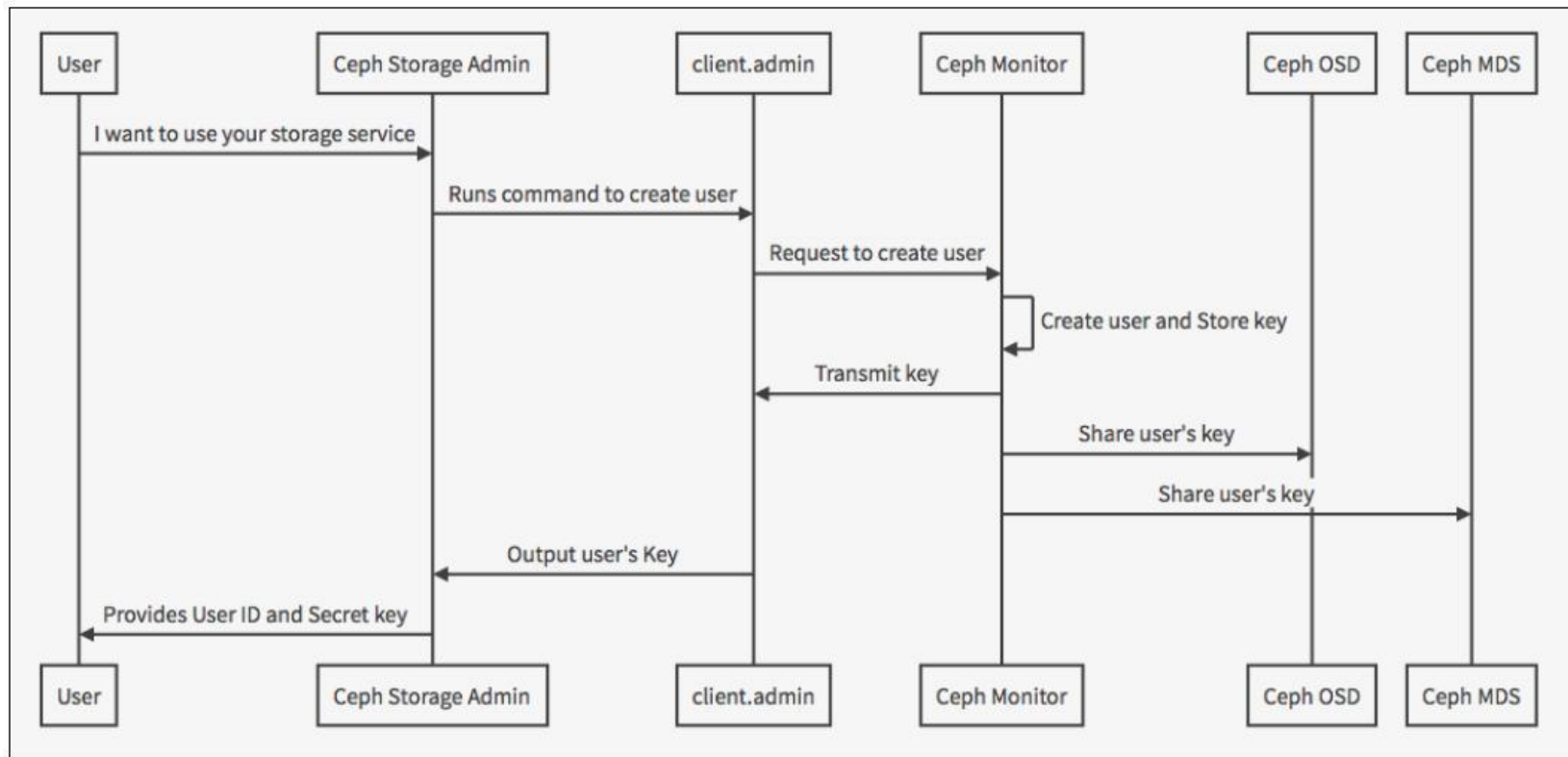
Ceph Write IO Flow: Receiving Request



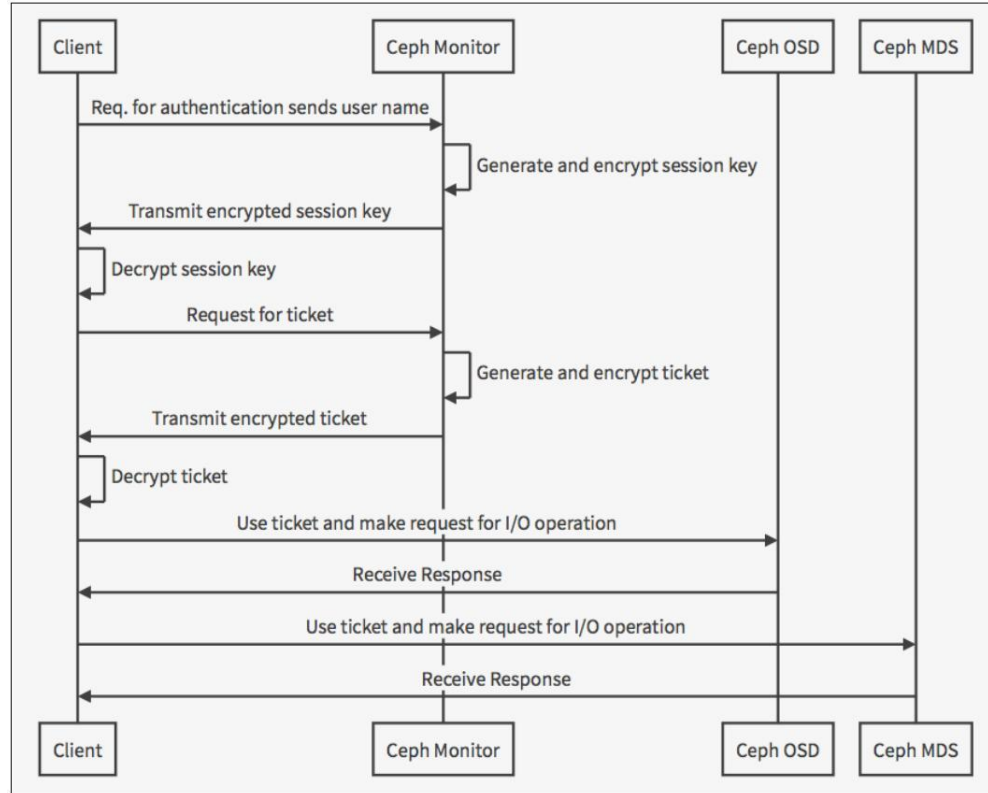
Ceph Write IO Flow: in File Store



Процесс создания пользователя и его ключа безопасности



Процесс аутентификации



Авторизация в CEPH

- Общий вид `{daemon-type} 'allow {capability}' [{daemon-type} 'allow {capability}]`
- Monitor caps: Включает параметры `r`, `w`, `x`, а также `allow profiles {cap}`
- OSD caps: Включает параметры `r`, `w`, `x`, `class-read`, `class-write` и `profile osd`.
- MDS caps: Требуется только `allow`

Что делать при падении

- Смотрим загрузку сети
- Смотрим потребление памяти
- Мониторим иопсы у клиентов
- Смотрим в логи ceph
- Регулируем
 - `osd_recovery_op_priority`
 - `osd_recovery_threads`
 - `osd_client_op_priority`

Вывод ноды из кластера.

- Выставить weight на OSD ноды в 0
- Выполнить следующие команды:
 - `ceph osd out osd.$i`
 - `systemctl stop ceph-osd.target`
 - `ceph osd crush remove osd.$i`
 - `umount /var/lib/ceph/osd/`
 - `ceph auth del osd.$i`
 - `ceph osd rm osd.$i`

Мониторинг CEPH

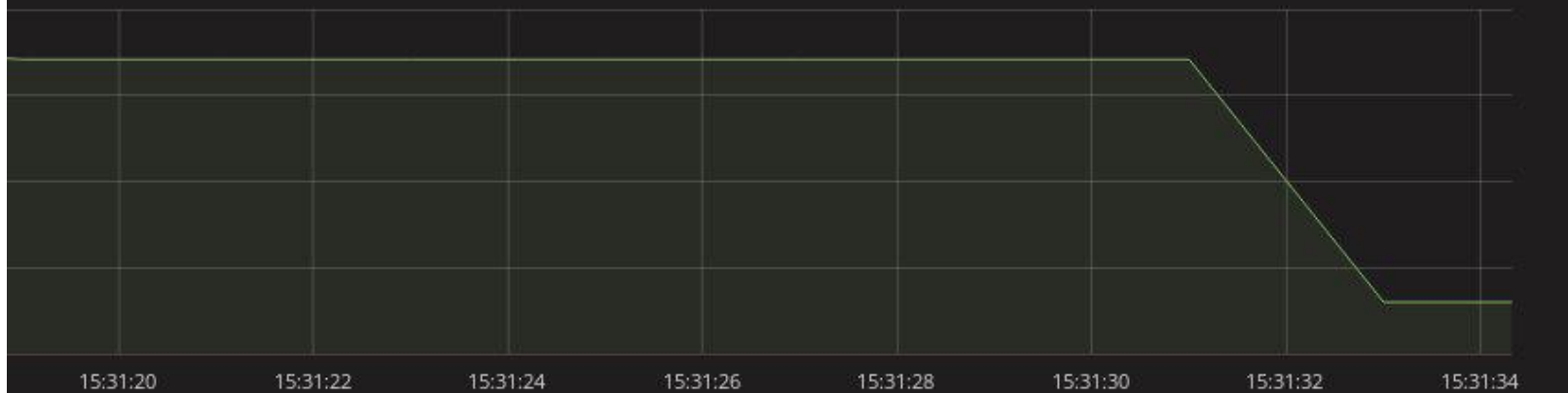
- Статус кластера OK Health Warn Health err
- warning на clock skew
- fs apply latency
- fs commit latency
- CEPH Cluster IOPS
- CEPH cluster IO (Mb\sec)
- Мониторить в логах наличие слова "slow request"

Мониторинг сервера с CERN

- Загрузка CPU
- Загрузка RAM
- Загрузка сетевых каналов.
- Заполненность дисков. warning на 85% Critical 92%

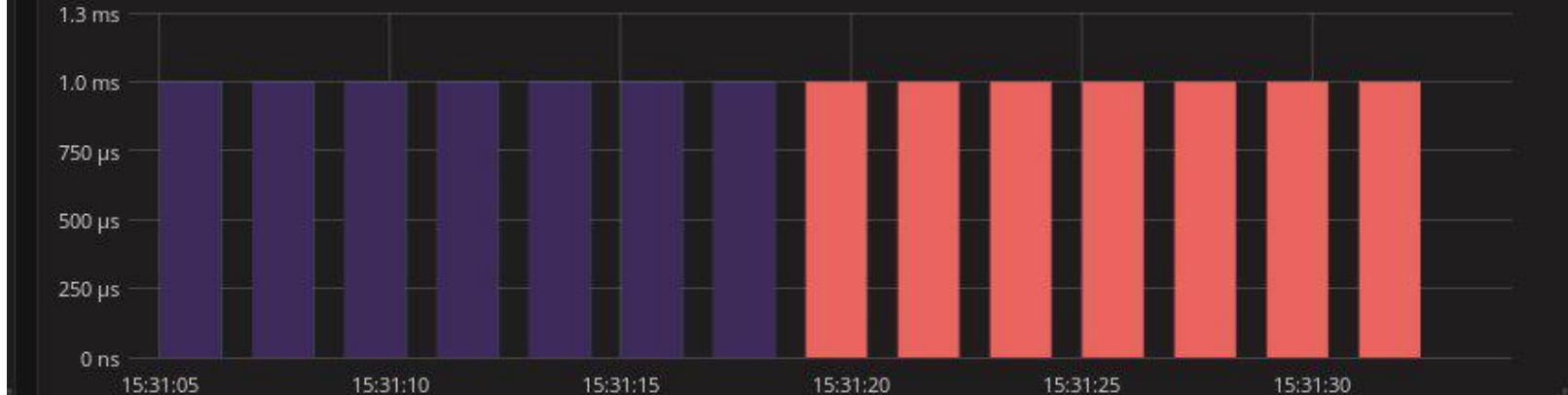
write_bytes SGK4

Last 30 seconds



ceph_osd_perf_commit_latency_seconds

Last 30 seconds





CEPH-SGK4

Today R...

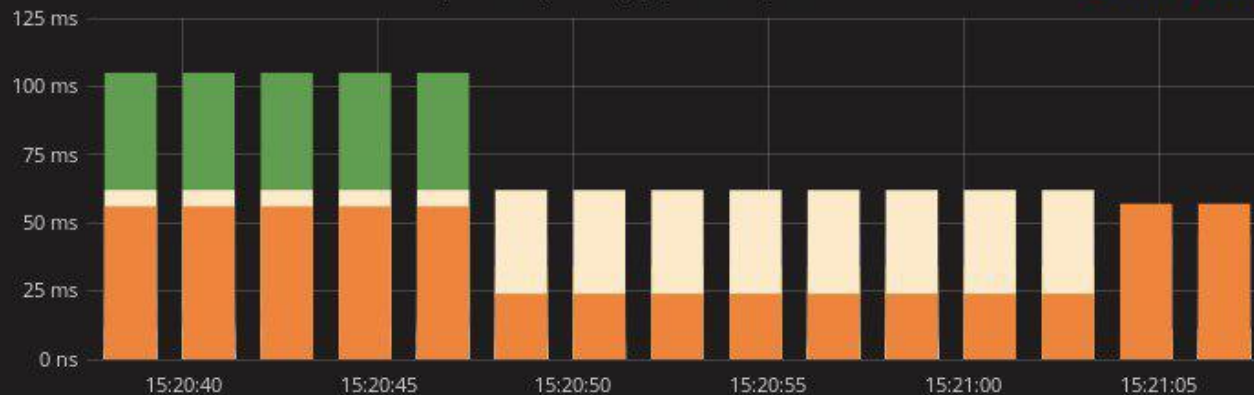
SGK4-write-iops

Last 30 seconds



ceph_osd_perf_apply_latency_seconds

Last 30 seconds



Начальный дебаг проблем. (FAQ)

Логи:

- `/var/log/ceph/ceph.log`
- `/var/log/ceph/ceph-mon.NAME.log`
- `/var/log/ceph/ceph-osd.ID.log`

Начальный дебаг проблем. (FAQ)

B) unable to mount OSD ...

O) Проверяем права на блочное устройство и существование точки монтирования

B) unable to link journal...

O) Проверяем права на симлинк journal в папке OSD. проверяем права на блочное устройство

B) CEPH тормозит.

O) Смотрим мониторинг, slow request, правильную работу журналов, iostat на нодах.

B) У меня тут демон OSD падает с жутким трейсом!

O) Вспоминаем, что меняли перед этими событиями, смотрим не было ли OOM, гуглим, идем в чат ceph_ru

Вопросы?

