# Benchmarking Ceph's BlueStore

## Niklaus Hofer
## Open Cloud Day 2018
## 2018-05-30

# stoney cloud

# Where we are coming from

- Small cloud company
  - Mostly PaaS

- FOSS-Cloud.org based cloud
  - libvirt based
  - No longer actively maintained
  - GlusterFS storage backend

# State of the (OpenStack) cloud

- OpenStack based

- Live: next month

  - A lot of expensive hardware is currently idly

  - Benchmarks!
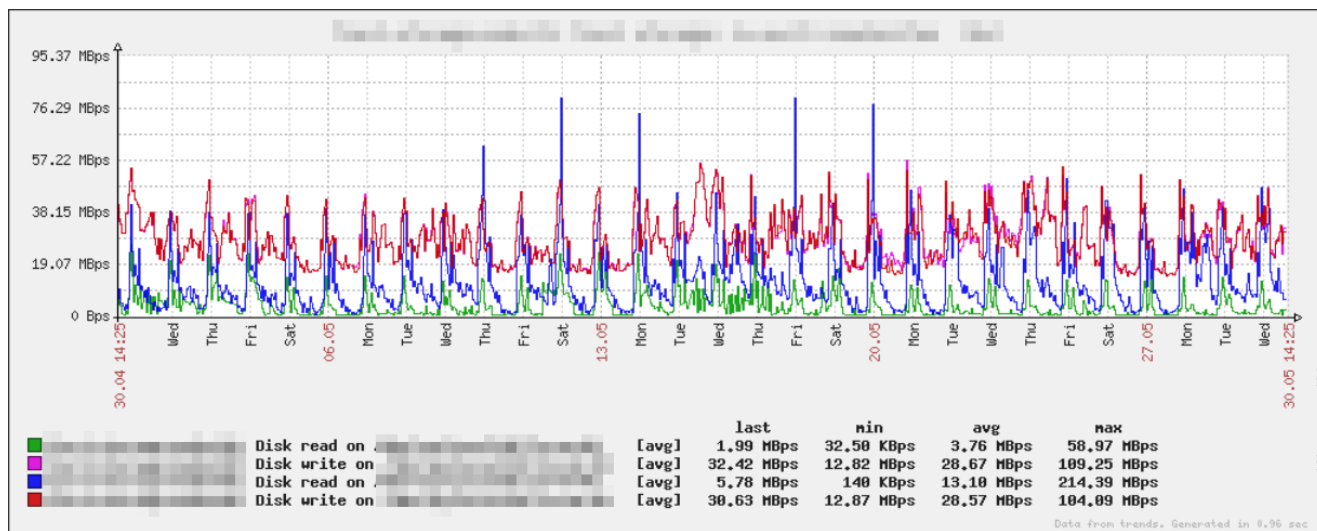
# Early benchmarks

- Direct SSD vs Ceph vs GlusterFS

- Later: GlusterFS is NOT an option

- Variations within Ceph
  - BlueStore vs Filestore

# Ceph hardware

- Dedicated nodes

- 3 nodes

- SATA SSDs
  - Micron m5100 max

- NVME SSDs Test

Benchmarking Ceph's BlueStore

# Ceph setup

- Ceph Luminous

- BlueStore

- Replica 3

# Ceph

# Ceph journal

- Limitations of POSIX filesystem

- All data gets written twice

- Optimization: Journal on separate disk

```
#         Start           End   Size  Type            Name
1      10487808    3750748814   1.8T  Ceph OSD        ceph data
2          2048      10487807     5G  Ceph Journal    ceph journal



/dev/sdd :
 /dev/sdd1 ceph data, active, cluster ceph, osd.0, journal /dev/sdd2
 /dev/sdd2 ceph journal, for /dev/sdd1
```
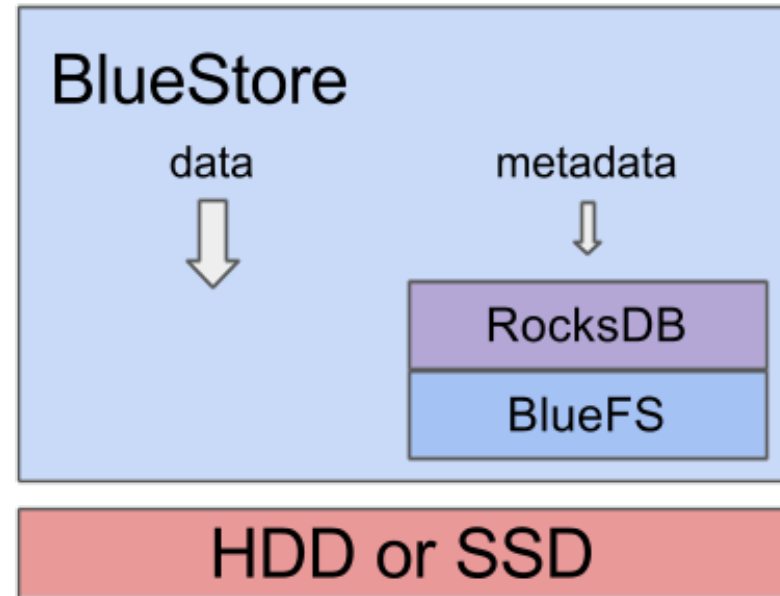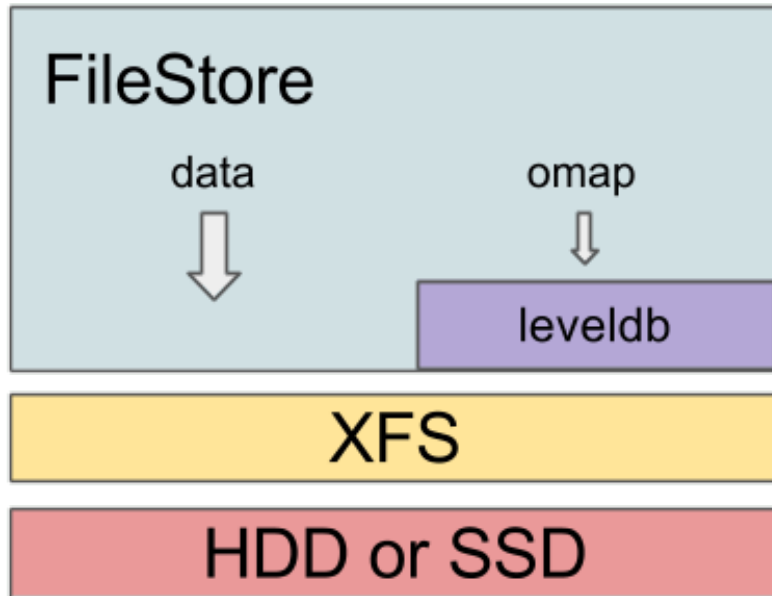
# Ceph journal optimization

# Ceph BlueStore

# Healing performance

```
--dsk/sdc-- --dsk/sdd-- --dsk/sde--          --dsk/sdc-- --dsk/sdd-- --dsk/sde--
 read  writ: read  writ: read  writ          read  writ: read  writ: read  writ
3905k 4397k:3469k 3666k:3380k 5890k             0  132M: 128M    0 :   0     0
    0  193M:    0     0 : 111M    0             0  128M: 132M    0 :   0     0
    0  264M:    0     0 : 123M    0             0  128M: 128M    0 :   0     0
    0  241M:    0     0 : 114M    0             0  112M: 108M    0 :   0     0
    0  255M:    0     0 : 114M    0             0  125M: 128M    0 :   0     0
    0  251M:    0     0 : 121M    0             0  128M: 128M    0 :   0     0
    0  241M:    0     0 : 116M    0             0  120M: 117M    0 :   0     0
    0  247M:    0     0 : 120M    0             0  111M: 112M    0 :   0     0
    0  245M:    0     0 : 122M    0             0  125M: 128M    0 :   0     0
    0  244M:    0     0 : 122M    0             0  128M: 127M 652k:   0     0
    0  246M:    0     0 : 110M    0             0  130M: 129M    0 :   0     0
    0  236M:    0     0 : 118M    0             0  132M: 132M    0 :   0     0
    0  253M:    0     0 : 124M    0             0  131M: 132M    0 :   0     0
    0  251M:    0     0 : 126M    0             0  130M: 128M    0 :   0     0
    0  251M:    0     0 : 122M    0             0  131M: 132M    0 :   0     0
    0  247M:    0     0 : 114M    0             0  129M: 128M    0 :   0     0
```

# Benchmarking

# Objectives

- Single VM

- VM located on separate storage

- Phoronix Test Suite

  – pts/disk

# Benchmarking results

# AIO-Stress



AIO-Stress v0.21
Random Write

MB/s, More Is Better

PHORONIX-TEST-SUITE.COM

| | |
|---|---|
| local filesystem<br>SE +/- 73.02 | 1478.38 |
| CEPH Jewel 3 OSDs replica 1<br>SE +/- 25.85 | 1721.82 |
| Direct SSD io=native cache=none<br>SE +/- 55.02 | 1802.66 |
| CEPH Jewel 1 OSD w/ external Journal<br>SE +/- 109.84 | 1822.87 |
| CEPH Jewel 1 OSD<br>SE +/- 13.54 | 1340.55 |
| CEPH Jewel 3 OSDs replica 3<br>SE +/- 24.90 | 1818.64 |
| CEPH luminous bluestore 3 OSDs replica 3<br>SE +/- 25.18 | 1690.54 |
| CEPH luminous bluestore 3 OSDs replica 1<br>SE +/- 96.24 | 1754.61 |
| CEPH luminous bluestore 1 OSD<br>SE +/- 68.62 | 1773.67 |

400   800   1200   1600   2000

Phoronix Test Suite 7.8.0

1. (CC) gcc options: -pthread -laio

# SQLite



**SQLite v3.22**
Timed SQLite Insertions

◄ Seconds, Less Is Better    PHORONIX-TEST-SUITE.COM

| | |
|---|---|
| **local filesystem** SE +/- 0.06 | 20.61 |
| **CEPH Jewel 3 OSDs replica 1** SE +/- 0.77 | 52.75 |
| **Direct SSD io=native cache=none** SE +/- 0.28 | 17.29 |
| **CEPH Jewel 1 OSD w/ external Journal** SE +/- 0.34 | 45.10 |
| **CEPH Jewel 1 OSD** SE +/- 0.10 | 46.21 |
| **CEPH Jewel 3 OSDs replica 3** SE +/- 0.38 | 98.30 |
| **CEPH luminous bluestore 3 OSDs replica 3** SE +/- 0.93 | 109.48 |
| **CEPH luminous bluestore 3 OSDs replica 1** SE +/- 1.07 | 69.95 |
| **CEPH luminous bluestore 1 OSD** SE +/- 0.78 | 65.14 |

Phoronix Test Suite 7.8.0

1. (CC) gcc options: -O2 -ldl -lpthread

# FS-Mark



**FS-Mark v3.3**
**1000 Files, 1MB Size**

ptsh.

▶ Files/s, More Is Better                    PHORONIX-TEST-SUITE.COM

| Configuration | Value | SE |
|---|---|---|
| local filesystem | 152.13 | SE +/- 4.84 |
| CEPH Jewel 3 OSDs replica 1 | 87.98 | SE +/- 1.36 |
| Direct SSD io=native cache=none | 159.03 | SE +/- 1.29 |
| CEPH Jewel 1 OSD w/ external Journal | 95.50 | SE +/- 1.35 |
| CEPH Jewel 1 OSD | 83.53 | SE +/- 0.80 |
| CEPH Jewel 3 OSDs replica 3 | 61.93 | SE +/- 0.29 |
| CEPH luminous bluestore 3 OSDs replica 3 | 66.07 | SE +/- 0.64 |
| CEPH luminous bluestore 3 OSDs replica 1 | 82.60 | SE +/- 1.25 |
| CEPH luminous bluestore 1 OSD | 83.60 | SE +/- 0.46 |

Phoronix Test Suite 7.8.0

1. (CC) gcc options: -static

# Dbench



Dbench v4.0
48 Clients

MB/s, More Is Better

PHORONIX-TEST-SUITE.COM

| | |
|---|---|
| local filesystem<br>SE +/- 98.18 | 812.43 |
| CEPH Jewel 3 OSDs replica 1<br>SE +/- 7.36 | 968.60 |
| Direct SSD io=native cache=none<br>SE +/- 2.82 | 1220.01 |
| CEPH Jewel 1 OSD w/ external Journal<br>SE +/- 2.21 | 1055.65 |
| CEPH Jewel 1 OSD<br>SE +/- 8.61 | 938.32 |
| CEPH Jewel 3 OSDs replica 3<br>SE +/- 1.19 | 712.22 |
| CEPH luminous bluestore 3 OSDs replica 3<br>SE +/- 2.02 | 679.26 |
| CEPH luminous bluestore 3 OSDs replica 1<br>SE +/- 3.97 | 768.75 |
| CEPH luminous bluestore 1 OSD<br>SE +/- 12.96 | 842.77 |

300    600    900    1200    1500

Phoronix Test Suite 7.8.0

1. (CC) gcc options: -lpopt -O2

# Threaded I/O tester



**Threaded I/O Tester v20170503**
64MB Random Write - 32 Threads

MB/s, More Is Better

PHORONIX-TEST-SUITE.COM

| | |
|---|---|
| **local filesystem** SE +/- 27.51 | 958.96 |
| **CEPH Jewel 3 OSDs replica 1** SE +/- 5.79 | 337.00 |
| **Direct SSD io=native cache=none** SE +/- 10.53 | 555.54 |
| **CEPH Jewel 1 OSD w/ external Journal** SE +/- 1.00 | 300.23 |
| **CEPH Jewel 1 OSD** SE +/- 5.80 | 299.60 |
| **CEPH Jewel 3 OSDs replica 3** SE +/- 2.37 | 214.01 |
| **CEPH luminous bluestore 3 OSDs replica 3** SE +/- 9.35 | 151.00 |
| **CEPH luminous bluestore 3 OSDs replica 1** SE +/- 3.91 | 255.32 |
| **CEPH luminous bluestore 1 OSD** SE +/- 3.19 | 229.61 |

Phoronix Test Suite 7.8.0

1. (CC) gcc options: -O2

stoney cloud
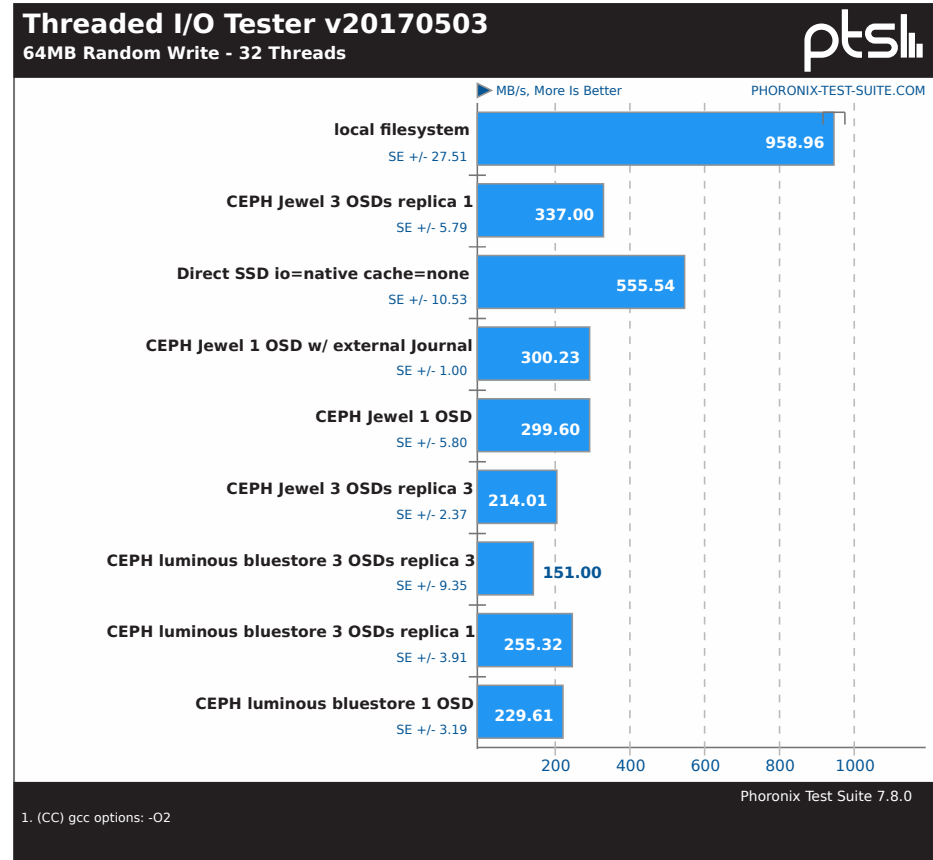
# Analysis and conclusion

# What went wrong?

- Largely guesswork for now

- Spectre patches?

- Regression?

- Configuration problem?

# Future work

- Optimization

- NVME vs SATA

    - Especially for database setups

- More precise IOPS measurements

    - fio, highly workload dependant

# Questions?

**stepping stone GmbH**
Wasserwerkgasse 7
CH-3011 Bern

Telefon: +41 31 332 53 63
www.stepping-stone.ch
info@stepping-stone.ch