

# Understanding Ceph

## One Performance Counter at a Time

Marcel Lauhoff | Staff Software Engineer

December, 2024

## Scope

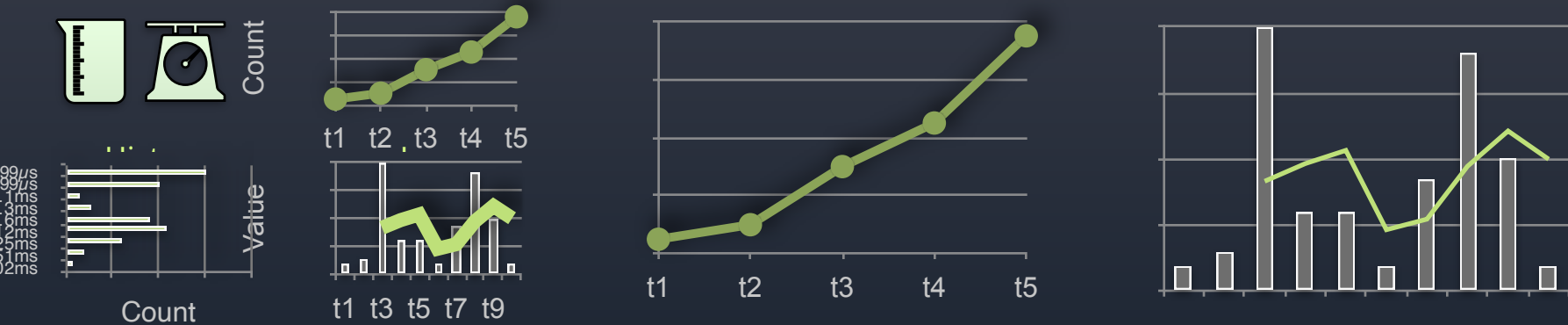
**What can we learn by looking at Perf Counters?**

**Let's develop an intuition!**

We won't go into analysis methods like the USE method

**We look at *raw* performance counters.**

No Prometheus, Grafana, etc.



## What is a Perf Counter anyway?

A **metric** is a **measurement** of a service captured at runtime  
(OpenTelemetry)

Values **captured** at **strategic** places  
(me)

Examples:

librados sends a write operation → increment osdop\_write counter

OSD processed an operation → update op\_latency

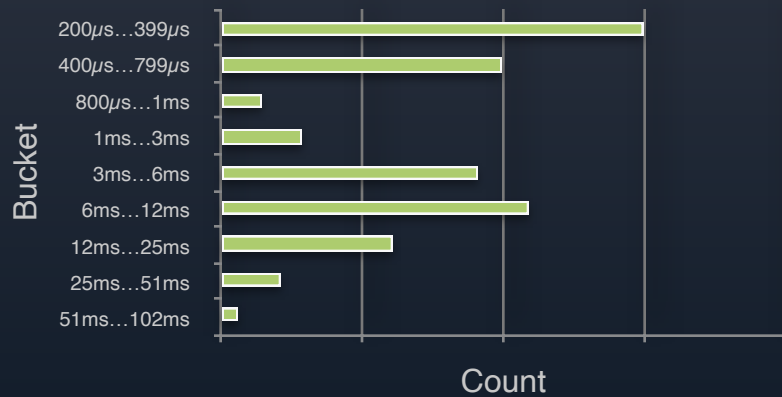
Gauge



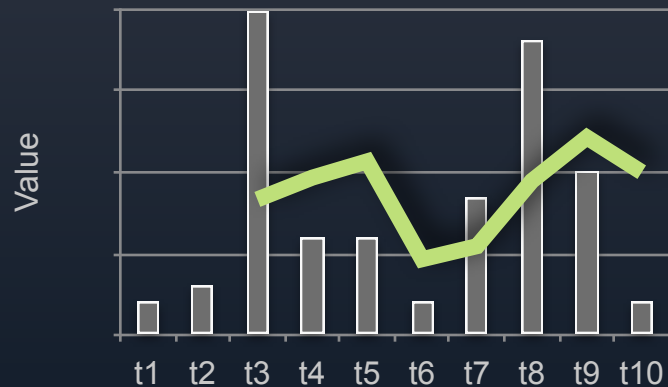
Counter



Histogram



Long Running Average



## How to read Perf Counters?

```
# ceph daemon $out/radosgw.8000.asok perf dump
```

```
# ceph tell osd.0 perf dump
```

```
# ceph tell osd.0 perf schema
```

```
# ceph tell osd.0 counter dump
```

```
# ceph tell osd.0 counter schema
```

# 17k JSON

## PERF COUNTER BROWSER

Found 880 perf counters in 41 groups

Group	Name	Description	Type	Tags
async messenger	msggr_connection_ready_timeouts	Number of not yet ready connections declared as dead	integer	mds, mgr, mon, client, osd
async messenger	msggr_connection_idle_timeouts	Number of connections closed due to idleness	integer	mds, mgr, mon, client, osd
async messenger	msggr_rcv_messages	Network received messages	integer	mds, mgr, mon, client, osd
async messenger	msggr_send_messages	Network sent messages	integer	mds, mgr, mon, client, osd
async messenger	msggr_rcv_bytes	Network received bytes	integer	mds, mgr, mon, client, osd
async messenger	msggr_send_bytes	Network sent bytes	integer	mds, mgr, mon, client, osd
async messenger	msggr_created_connections	Created connection number	integer	mds, mgr, mon, client, osd
async messenger	msggr_active_connections	Active connection number	integer	mds, mgr, mon, client, osd
async messenger	msggr_running_total_time	The <small>Calendar</small> thread running	real	mds, mgr, mon,

MON=3 OSD=3 MGR=1  
 RGW=1 MDS=2 NFS=0  
 vstart.sh

880 counters

in

41 groups



CL'so

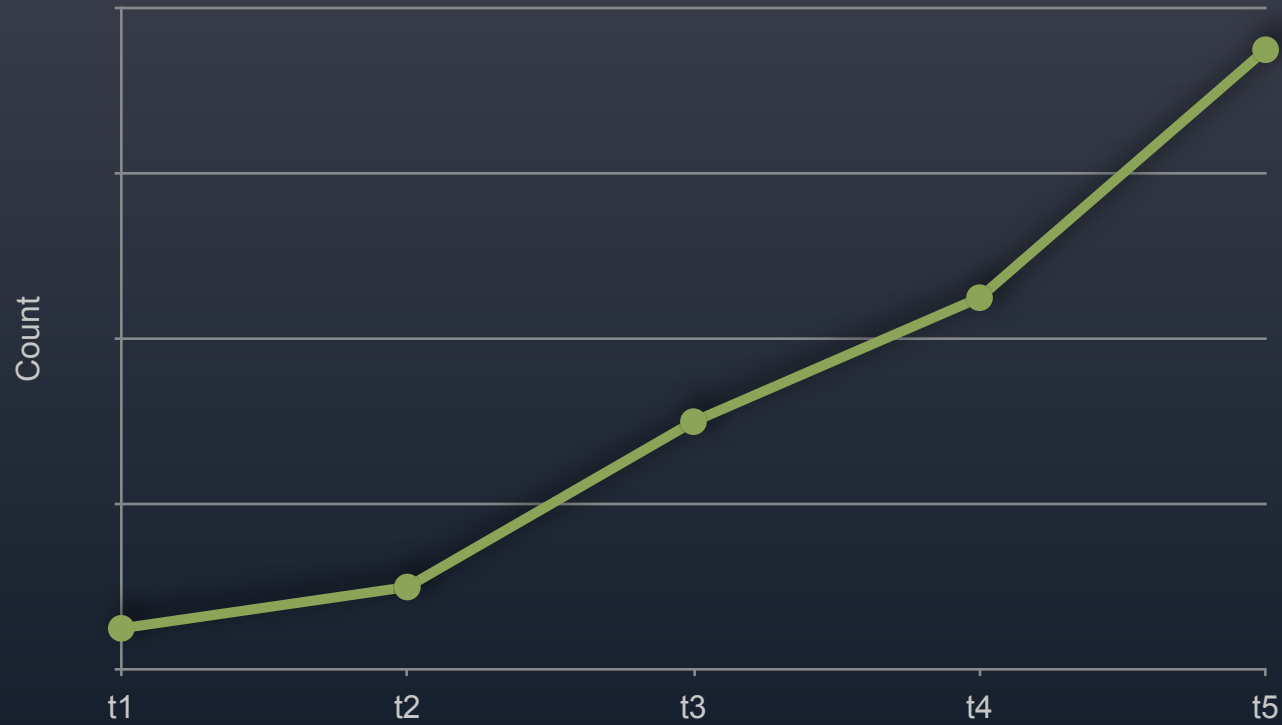
We focus on..

objecter (librados)

AsyncMessenger

osd

# Counter



## RGW: Single 4MB S3 PUT

2024-11-29T12:42:45.821+0100 7f0be538f6c0 1 beast: 0x7f0c53e61200: ::1 - testid [29/Nov/  
2024:12:42:45.758 +0100] "PUT /testbucket/13359 HTTP/1.1" 200 **4194304** - -  
**latency=0.063001677s**

```
"rgw_op": {  
  "put_obj_ops": 1,  
  "put_obj_bytes": 4194304,  
  ...  
}
```

```
"rgw": {  
  "req": 1,  
  ...  
}
```

CL'so

RGW: Single 4MB S3 PUT

objecter.{osdop|omap}

Table 1

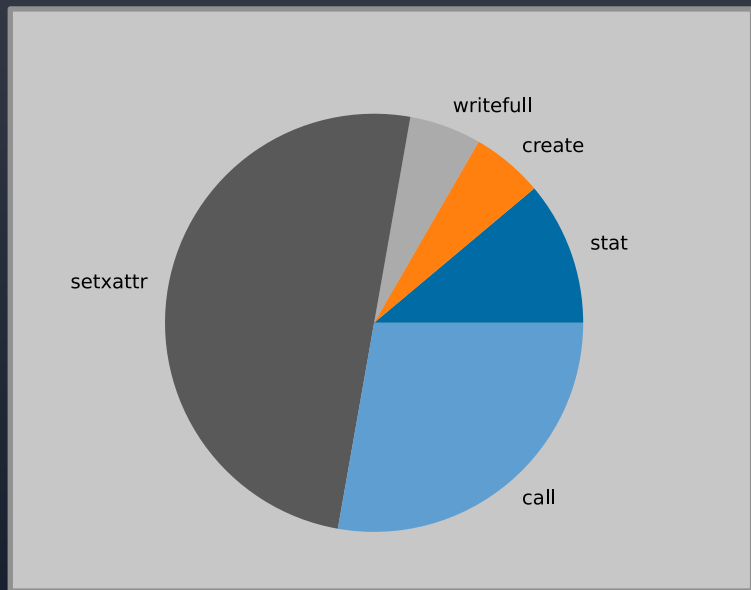
"osdop_stat"	2
"osdop_create"	1
"osdop_read"	0
"osdop_write"	0
"osdop_writefull"	1
"osdop_writesame"	0
"osdop_append"	0
"osdop_zero"	0
"osdop_truncate"	0
"osdop_delete"	0
"osdop_mapext"	0
"osdop_sparse_read"	0
"osdop_clonerange"	0
"osdop_getxattr"	0
"osdop_setxattr"	9
"osdop_cmpxattr"	0
"osdop_rmxattr"	0
"osdop_resetxattrs"	0
"osdop_call"	5
"osdop_watch"	0
"osdop_notify"	0
"osdop_src_cmpxattr"	0
"osdop_pgls"	0
"osdop_pgls_filter"	0
"osdop_other"	0
"omap_wr"	0
"omap_rd"	0
"omap_del"	0

RGW: Single 4MB S3 PUT

objecter.{osdop|omap}\_\*

stat	2
create	1
writefull	1
setxattr	9
call	5

$$\sum \dots = 18$$



RGW: Single 4MB S3 PUT

AsyncMessenger::Worker-\*

```
{  
  "msgr_recv_messages": 3,  
  "msgr_send_messages": 3,  
  "msgr_recv_bytes": 1481,  
  "msgr_send_bytes": 4197637,  
}
```

RGW: Single 4MB S3 PUT

$\Sigma$  (osdop\_\*, omap\_\*) = 18

$\Sigma$  msggr\_send\_messages = 3

objecter.op = 3

RGW: Single 4MB S3 PUT

$\sum \text{msg\_send\_bytes} - 4\text{MB}$

$\sim 4\text{k}$



## RGW: Single 4MB S3 PUT

```
osd.0
{
  "op": 6,
  "op_in_bytes": 0,
  "op_out_bytes": 224,
  "subop": 3,
  "subop_in_bytes": 4198251
}
```

```
osd.1
{
  "op": 4,
  "op_in_bytes": 0,
  "op_out_bytes": 0,
  "subop": 1,
  "subop_in_bytes": 4196187
}
```

```
osd.2
{
  "op": 4,
  "op_in_bytes": 4194304,
  "op_out_bytes": 0,
  "subop": 2,
  "subop_in_bytes": 2064
}
```

What does this teach us?

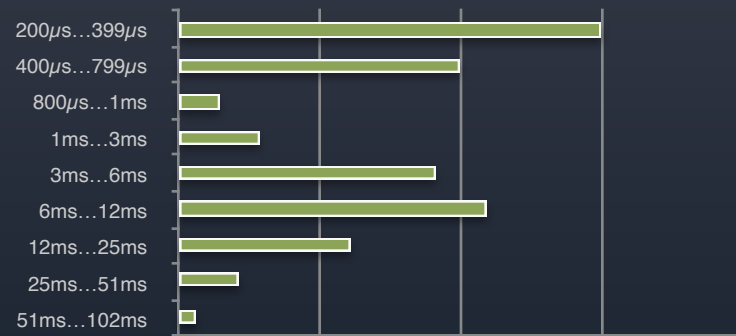
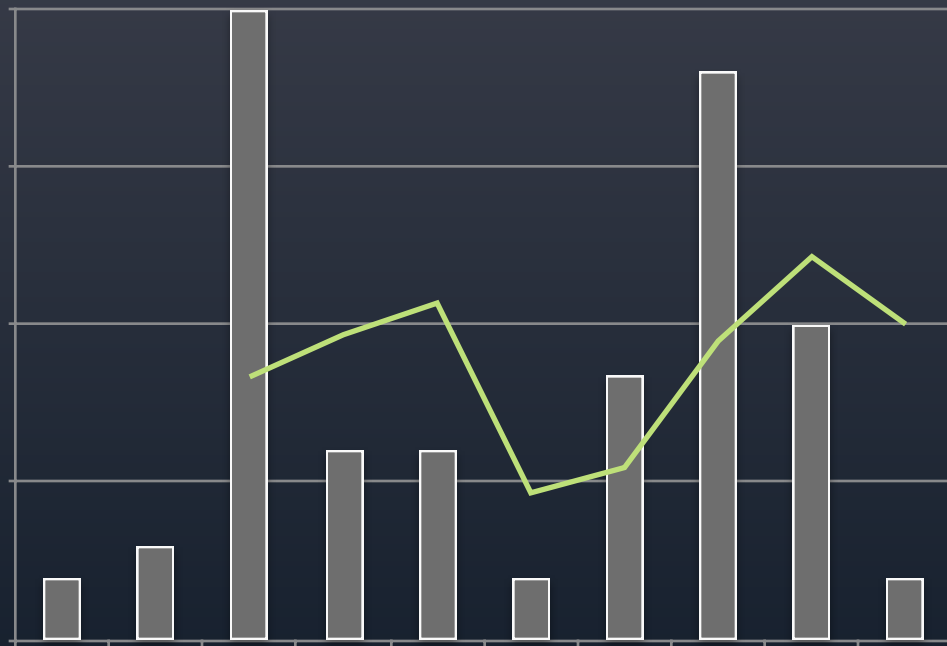
Op Mix

Op = [op, ..., op]

3 messages, 3 ops, 4k overhead

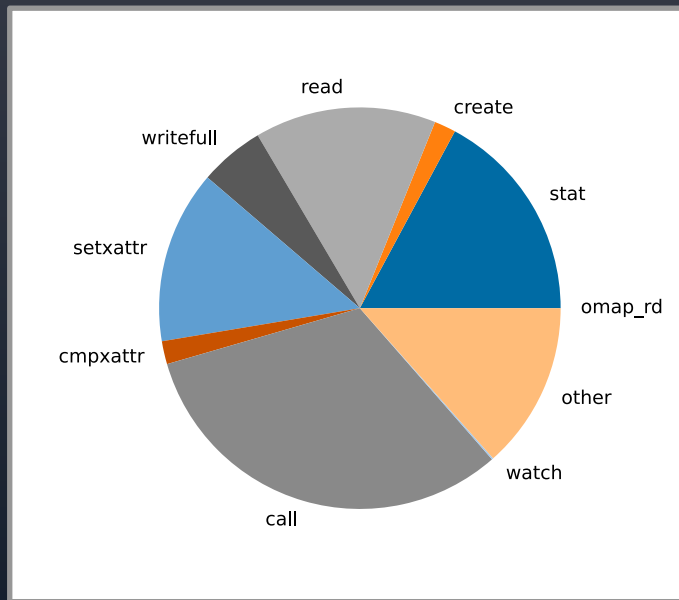
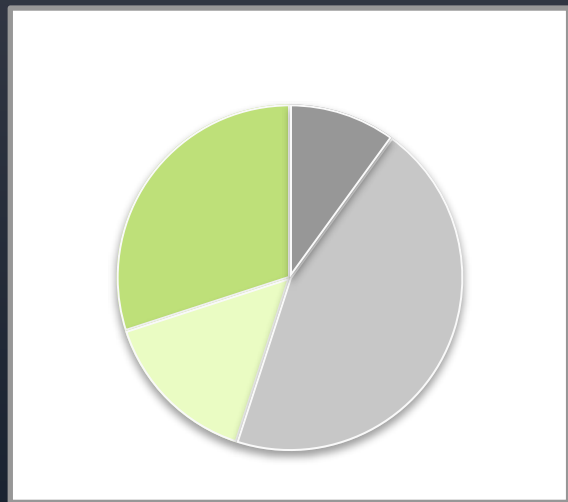
-> 18 queries and modifications to and from persistent data

# Latency



# S3 Mixed Workload Benchmark

10% DELETE, 45% GET, 15% PUT, 30% STAT



stat	20947
create	2125
read	17836
writefull	6375
setxattr	17000
cmpxattr	2265
call	39004
watch	150
other	16417
omap <sub>rd</sub>	2

## S3 Mixed Workload Benchmark

Client / RGW:

objecter.op\_latency 43ms

OSDs:

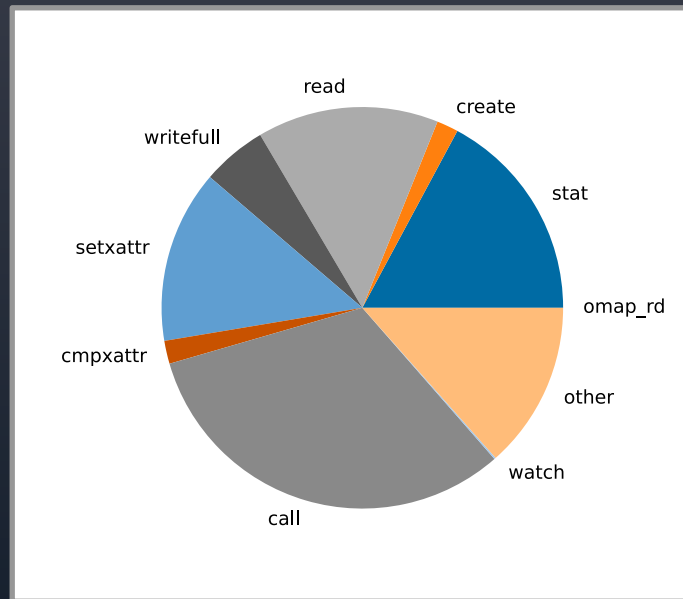
osd.op\_latency

0: 41ms, 1: 44ms, 2: 39ms

## S3 Mixed Workload Benchmark

~40 ms average latency

Is this actually meaningful?



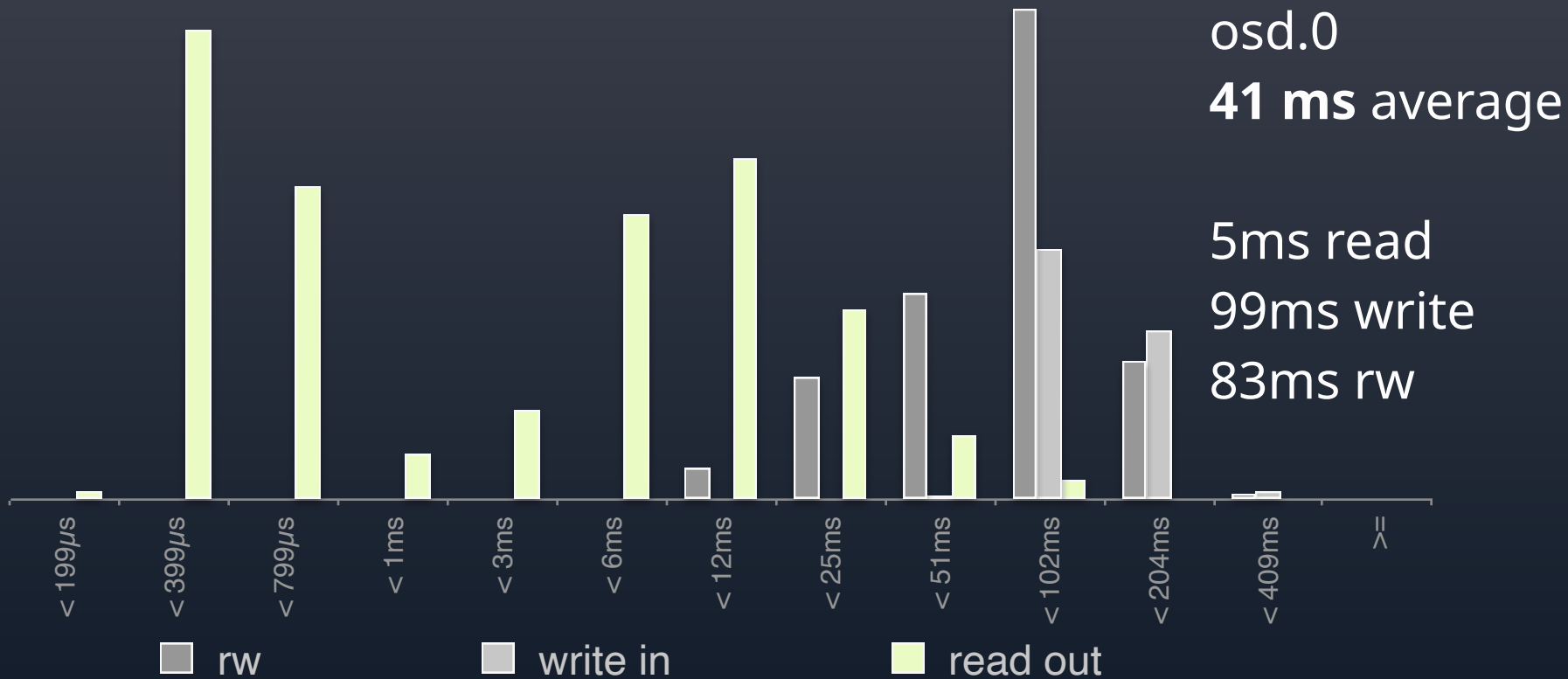
## S3 Mixed Workload Benchmark

osd.op\_latency family:

osd.op\_w\_latency osd.op\_rw\_latency osd.op\_r\_latency

target	op_latency [ms]	op_r_latency [ms]	op_w_latency [ms]	op_rw_latency [ms]
osd.0	41	5	99	83
osd.1	44	5	98	82
osd.2	39	5	100	82

# Histograms to the rescue





# Histograms to the rescue

**op\_latency captures every operation**

op_r_latency_out_bytes_histogram	Histogram of operation latency (including queue time) + data read
op_w_latency_in_bytes_histogram	Histogram of operation latency (including queue time) + data written
op_rw_latency_in_bytes_histogram	Histogram of rw operation latency (including queue time) + data written
op_rw_latency_out_bytes_histogram	Histogram of rw operation latency (including queue time) + data read

# Ceph 2D Histograms: latency x bytes

	0...511	512...1023	1K...2K	2K...4K	4K...8K	8K...16K	16K...32K	32K...64K	64K...128K	128K...256K	256K...512K	512K...1024K	1M...2M	2M...4M	4M...8M	8M...16M
<0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0ns...99µs	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
100µs...199µs	0	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0
200µs...399µs	31	1421	0	0	0	0	0	0	0	0	0	0	0	333	32	0
400µs...799µs	45	413	9	1	0	0	0	0	0	0	0	0	0	185	567	0
800µs...1ms	17	39	26	3	0	0	0	0	0	0	0	0	0	59	104	0
1ms...3ms	3	21	5	0	0	0	0	0	0	0	0	0	0	382	308	0
3ms...6ms	1	33	0	0	0	0	0	0	0	0	0	0	0	530	1168	0
6ms...12ms	2	38	0	0	0	0	0	0	0	0	0	0	0	313	1020	0
12ms...25ms	3	26	0	0	0	0	0	0	0	0	0	0	0	135	466	0
25ms...51ms	0	21	0	0	0	0	0	0	0	0	0	0	0	35	120	0
51ms...102ms	0	8	0	0	0	0	0	0	0	0	0	0	0	10	26	0
102ms...204ms	0	3	0	0	0	0	0	0	0	0	0	0	0	2	7	0

osd.op\_r\_latency\_out\_bytes\_histogram

# Ceph 2D Histograms: latency x bytes

	0...511	512...1023	1K...2K	2K...4K	4K...8K	8K...16K	16K...32K	32K...64K	64K...128K	128K...256K	256K...512K	512K...1024K	1M...2M	2M...4M	4M...8M	8M...16M
<0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0ns...99µs	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
100µs...199µs	0	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0
200µs...399µs	31	1421	0	0	0	0	0	0	0	0	0	0	0	333	32	0
400µs...799µs	45	413	9	1	0	0	0	0	0	0	0	0	0	185	567	0
800µs...1ms	17	39	26	3	0	0	0	0	0	0	0	0	0	59	104	0
1ms...3ms	3	21	5	0	0	0	0	0	0	0	0	0	0	382	308	0
3ms...6ms	1	33	0	0	0	0	0	0	0	0	0	0	0	530	1168	0
6ms...12ms	2	38	0	0	0	0	0	0	0	0	0	0	0	313	1020	0
12ms...25ms	3	26	0	0	0	0	0	0	0	0	0	0	0	135	466	0
25ms...51ms	0	21	0	0	0	0	0	0	0	0	0	0	0	35	120	0
51ms...102ms	0	8	0	0	0	0	0	0	0	0	0	0	0	10	26	0
102ms...204ms	0	3	0	0	0	0	0	0	0	0	0	0	0	2	7	0

osd.op\_r\_latency\_out\_bytes\_histogram

# Ceph 2D Histograms: latency x bytes

	0...511	512...1023	1K...2K	2K...4K	4K...8K	8K...16K	16K...32K	32K...64K	64K...128K	128K...256K	256K...512K	512K...1024K	1M...2M	2M...4M	4M...8M	8M...16M
<0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0ns...99µs	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
100µs...199µs	0	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0
200µs...399µs	31	1421	0	0	0	0	0	0	0	0	0	0	0	333	32	0
400µs...799µs	45	413	9	1	0	0	0	0	0	0	0	0	0	185	567	0
800µs...1ms	17	39	26	3	0	0	0	0	0	0	0	0	0	59	104	0
1ms...3ms	3	21	5	0	0	0	0	0	0	0	0	0	0	382	308	0
3ms...6ms	1	33	0	0	0	0	0	0	0	0	0	0	0	530	1168	0
6ms...12ms	2	38	0	0	0	0	0	0	0	0	0	0	0	313	1020	0
12ms...25ms	3	26	0	0	0	0	0	0	0	0	0	0	0	135	466	0
25ms...51ms	0	21	0	0	0	0	0	0	0	0	0	0	0	35	120	0
51ms...102ms	0	8	0	0	0	0	0	0	0	0	0	0	0	10	26	0
102ms...204ms	0	3	0	0	0	0	0	0	0	0	0	0	0	2	7	0

osd.op\_r\_latency\_out\_bytes\_histogram

# Recap

## Counters - What is a S3 PUT?



## Latency - Mixed S3 Workload



## What's next?

Tracing (e.g bpftrace)

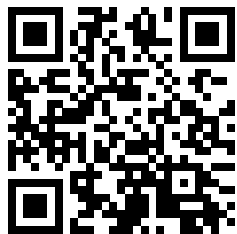
OSD Perf Queries

Collection at the Mgr: collection, Prometheus, etc

Bluestore metrics


Combining Ceph with system node exporter metrics

Messenger stats: cntop <https://github.com/irq0/cntop>



[https://github.com/irq0/talk\\_ceph\\_perf\\_counters](https://github.com/irq0/talk_ceph_perf_counters)

CLY'SO

Powered by  ceph

# Thank you!

Marcel Lauhoff <[marcel.lauhoff@clyso.com](mailto:marcel.lauhoff@clyso.com)>