

Identifying latency outliers in workload testing

Sayali Shirode
Systems Performance Engineer



Problem statement

When running and collecting workload traces of RocksDB, we sometimes see large latency spikes. To debug this issue, we used FIO, which is relatively “simple tool” and we still see the latency spikes.

Is this the SSD misbehaving, or is it the system?

Real world application – RocksDB YCSB

- RocksDB is a storage focused NoSQL database maintained by Meta.
- After executing multiple runs of RocksDB YCSB read-heavy workload, started observing read latency of 113ms, expected was <5ms

Are these high latencies from the SSD or the system?

Why latency is critical for SSDs?

- Lower latency means faster access to data
- Impact the performance of applications and user experience.
- The factors affecting the latency include various hardware components, network stack, workload characteristics, storage architecture, software stack.

Input



Output

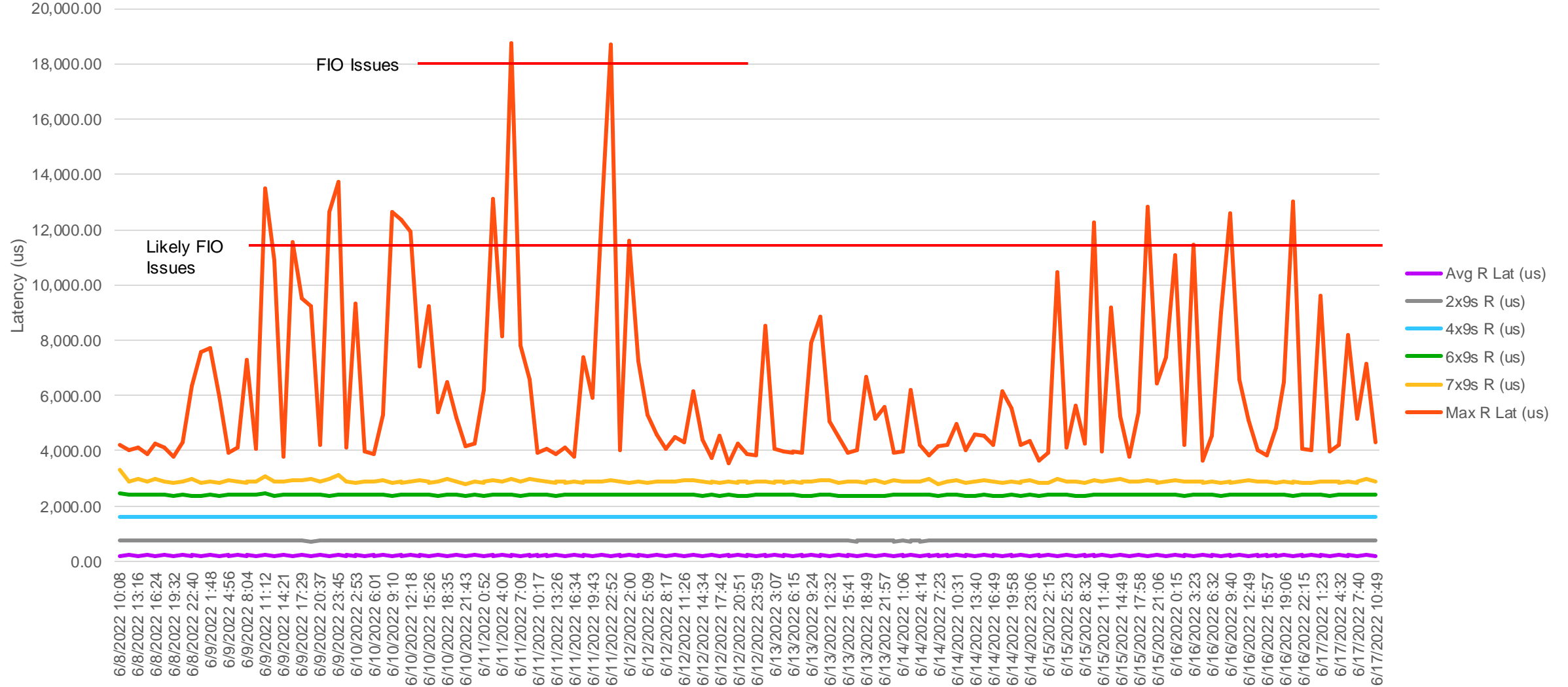


The time in between...
Latency

Synthetic workload - FIO

- Max latency of 18ms
- FIO job is stuck, because the core was busy with some other operation and not able to reach the completion command
- Takeaway – Latency excursion is not a drive issue, but likely the behavior of layer above it
 - Even “simple” tools like FIO can experience system effects that result in the reporting of high latencies
 - Noise on the system can affect max latency reporting in FIO.
 - QoS latencies out to 7x9's look consistent

70% 4k Random Read - Long Running - QoS over Time



How was the latency measured?

- BPF trace scripts
- Bpfttrace – tracing framework for Linux that allows to trace and profile various aspects of the system at runtime
- Biosnoop – details for each disk I/O with latency and look for time-ordered patterns
- OCP Latency monitor

```

1  #!/usr/bin/env bpfttrace
2  /*
3   * biosnoop.bt    Block I/O tracing tool, showing per I/O latency.
4   *                For Linux, uses bpfttrace, eBPF.
5   *
6   * TODO: switch to block tracepoints. Add offset and size columns.
7   *
8   * This is a bpfttrace version of the bcc tool of the same name.
9   *
10  * 15-Nov-2017  Brendan Gregg  Created this.
11  */
12
13 #ifndef BPFTRACE_HAVE_BTF
14 #include <linux/bkdev.h>
15 #include <linux/bkdev.h>
16 #endif
17
18 BEGIN
19 {
20     printf("%-12s %-7s %-16s %-6s %7s\n", "TIME(ms)", "DISK", "COMM", "PID", "LAT(ms)");
21 }
22
23 kprobe:blk_account_io_start,
24 kprobe:__blk_account_io_start
25 {
26     @start[arg0] = nsecs;
27     @iopid[arg0] = pid;
28     @iocomm[arg0] = comm;
29     @disk[arg0] = ((struct request *)arg0)->q->disk->disk_name;
30 }
31
32 kprobe:blk_account_io_done,
33 kprobe:__blk_account_io_done
34 /@start[arg0] != 0 && @iopid[arg0] != 0 && @iocomm[arg0] != ""/
35 {
36     $now = nsecs;
37     printf("%-12s %-7s %-16s %-6s %7s\n",
38           "elapsed / 1e6, @disk[arg0], @iocomm[arg0], @iopid[arg0],
39           ($now - @start[arg0]) / 1e6);
40
41     delete(@start[arg0]);
42     delete(@iopid[arg0]);
43     delete(@iocomm[arg0]);
44     delete(@disk[arg0]);
45 }
46
47 END
48 {
49     clear(@start);
50     clear(@iopid);
51     clear(@iocomm);
52     clear(@disk);
53 }

```

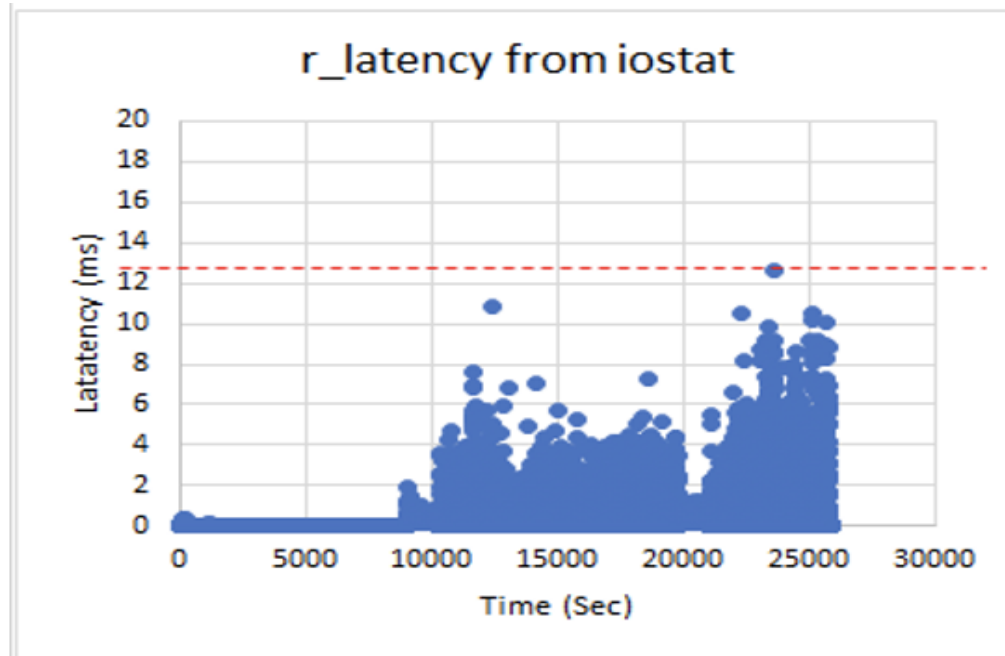
```

# biosnoop
TIME(s)      COMM          PID    DISK    T  SECTOR    BYTES    LAT (ms)
0.000004001  supervise    1950   xvda1   W  13092560  4096     0.74
0.000178002  supervise    1950   xvda1   W  13092432  4096     0.61
0.001469001  supervise    1956   xvda1   W  13092440  4096     1.24
0.001588002  supervise    1956   xvda1   W  13115128  4096     1.09
1.022346001  supervise    1950   xvda1   W  13115272  4096     0.98
1.022568002  supervise    1950   xvda1   W  13188496  4096     0.93
1.023534000  supervise    1956   xvda1   W  13188520  4096     0.79
1.023585003  supervise    1956   xvda1   W  13189512  4096     0.60
2.003920000  xfsaild/md0  456    xvdc    W  62901512  8192     0.23
2.003931001  xfsaild/md0  456    xvdb    W  62901513  512      0.25
2.004034001  xfsaild/md0  456    xvdb    W  62901520  8192     0.35
2.004042000  xfsaild/md0  456    xvdb    W  63542016  4096     0.36
2.004204001  kworker/0:3  26040  xvdb    W  41950344  65536    0.34
2.044352002  supervise    1950   xvda1   W  13192672  4096     0.65
[...]

```

Experiment

- iostat - system input/output statistics for devices and partitions
- Bus analyzer – Captures and analyzes the data transmitted on PCIe bus
 - Individual transaction information – like transaction type, address, data payload
 - Shows the signal timings



Transactions View

Book...	Down	Up	Age(us)	Pen...	Addr	LBA	Su...	Expected E...	Tag	Requester	Completer	mm:ss.ms_us_ns_ps	Delta Time	Port	Errors/
	Read	Good	86.2330	3		277F723		2000 5300	0000	Micron:51C3		00:18.056_601_056_0			Expert Down(1,1,1)
Bookmark	mm:ss.ms_us_ns_ps	Delta Time	Summary										Errors/Warnings		
	00:18.056_479_308_4		SubQBel; SQTail = 0x00000092;												
	00:18.056_479_449_3		Ack												
	00:18.056_479_582_2	0.1329	MRd(32); SubQ; Len = 0x40;												
	00:18.056_479_717_9	0.1358	Ack												
	00:18.056_490_839_1	11.1212	SubQ; Read; NSID = 0x00000001; LBA = 0x0277F723; NbBlocks = 1												
	00:18.056_490_986_1	0.1470	Ack												
	00:18.056_549_405_0	58.4190	Data; Len = 0x0200;												
	00:18.056_549_474_4	0.0694	Data(0x200); Len = 0x0200;												
	00:18.056_549_543_7	0.0694	Data(0x400); Len = 0x0200;												
	00:18.056_549_605_6	0.0620	Ack												
	00:18.056_549_619_2	0.0136	Data(0x600); Len = 0x0200;												
	00:18.056_549_688_5	0.0694	Data(0x800); Len = 0x0200;												
	00:18.056_549_688_6	0.0001	Ack												
	00:18.056_549_744_0	0.0555	Ack												
	00:18.056_549_757_9	0.0139	Data(0xA00); Len = 0x0200;												
	00:18.056_549_823_0	0.0652	Ack												
	00:18.056_549_827_2	0.0043	Data(0xC00); Len = 0x0200;												

Conclusion

- Latency between submission and completion on PCIe link is not seen to be 100ms
- Latency monitor does not show high latencies either
- We can conclude that these high latencies are on the system stack

References

- [BPF Performance Tools \(Book\) \(brendangregg.com\)](https://brendangregg.com)
- [GitHub - brendangregg/bpf-perf-tools-book: Official repository for the BPF Performance Tools book](https://github.com/brendangregg/bpf-perf-tools-book)
- <https://www.opencompute.org/>