

September 23-26, 2019 Santa Clara, CA

Average Storage Latency is Lacking. Let's Fix It!

Richard Elling Viking Enterprise Solutions a Division of Sanmina Corporation



Storage system latency does not have a normal distribution

- HDDs seek & rotate
- Flash SSDs garbage collection
- Fabrics
- Ordering & Re-ordering
- Multiple, serial queues
- Error recovery and retry

SD @

KING

Flash SSD Latency Anomaly







Ordering Problem





2019 Storage Developer Conference. © Viking Enterprise Solutions a Division of Sanmina Corporation

4

SD[®]

VIKING[™] Enterprise Solutions





Ubuntu



OSX

\$ iostat -0 disk0 disk2 cpu load average sps tps msp s sps tp: msps us sy id 1m 5m 15m 3092 177 0.0 1042 10 0.0 7 6 87 1.46 1.89 1.93

Note: ignore Linux svctm, it is broken

Linux /proc/diskstats



cat /proc/diskstats

•••

•••

8 0 sda 94376 18486 1919722 71620 70919 167617 5696296 1017520 0 71476 1102048

Field	Description	Field	Description			
I	major number	8	# of writes completed			
2	minor number	9	# of writes merged			
3	device name	10	# of sectors written			
4	# of completed reads	П	# of milliseconds spent writing			
5	# of reads merged	12	# of I/Os in progress			
6	# of sectors read	13	# of milliseconds spent doing I/Os			
7	# of milliseconds spent reading	14	weighted # of milliseconds spent doing I/Os			

Popular Tools: prometheus



2019 Storage Developer Conference. © Viking Enterprise Solutions a Division of Sanmina Corporation

SD[®]

VIKING[™] Enterprise Solutions

prometheus Example





Sampling: dtrace

- Dynamic Tracing (dtrace)
- Ported from Solaris to FreeBSD and OSX
- Halfway ported to OEL, see BPF for another path forward



Brendan Gregg • Jim M

VIKING

race

SD®

dtrace Latency Example



```
# cat iolatency.d
#!/usr/sbin/dtrace -s
io:::start
{
              start[arg0] = timestamp;
}
io:::done
/start[arg0]/
{
              @time["disk I/O latency (ns)"] = quantize(timestamp - start[arg0]);
              start[arg0] = 0;
```

dtrace Latency Example



SD©

# ./iolatency.d dtrace: script './iolatency.d' ma ^C	atched 2 probes						
disk I/O latency (ns)							
value Distribu	tion count						
32768	0						
65536	94						
131072 @@	361						
262144 @@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@							
524288 @@@@@@@	@ 1617						
1048576 @@	366						
2097152 @	178						
4194304 @	197						
8388608 @	296						
16777216 @@	313						
33554432 @	273						
67108864	34						
134217728	0						
268435456	2						
536870912	0						

Sampling: BPF

- No complete port of dtrace to Linux, so...
- Berkeley Packet Filter (BPF) meets the Linux kernel
 - Also known as extended BPF (eBPF)
- Similar to dtrace, systemtap, and friends
- Vibrant community with many good scripts to learn from and adapt

www.iovisor.org/technology/ebpf





Performance Tools

BPF

Linux System

VIKING Enterprise Solutions SD®

2 Good Ideas: eBPF_exporter

VIKING[™] Enterprise Solutions

- Descibe eBPF as YAML
- ebpf_exporter deamon
 - reads YAML
 - runs BPF
 - results are available to prometheus
- Use to sample latency distributions on almost anything... userland or kernel

SD @





- When the going gets tough, samplers are adversely affected
 - Can miss samples
 - Can commit suicide (protect the system)
- Require elevated roles -- may conflict with production security policies
- Can impact overall performance

Good idea: use samplers for debug and prototyping then add useful latency distributions in the code

Exemplar: chrome

VIKING[™] Enterprise Solutions

SD @

Try it! chrome://histograms

Histogram: Event.Latency.OS.MOUSE WHEEL recorded 18505 samples, mean = 1176.9 (flags = 0x41) 0 ... $258 O (1 = 0.0\%) \{0.0\%\}$ $340 - 0(52 = 0.3\%) \{0.0\%\}$ 448 -----O (2524 = 13.6%) {0.3%} 590 ------O (4909 = 26.5%) {13.9%} 777 -----O (5178 = 28.0%) {40.5%} 1023 -----O (2234 = 12.1%) {68.4%} 1347 -----O (763 = 4.1%) {80.5%} 1774 -----O (684 = 3.7%) {84.6%} 2336 -----O (931 = 5.0%) {88.3%} 3077 -----O (987 = 5.3%) {93.4%} 4053 --- 0 (194 = 1.0%) {98.7%} 5338 O (21 = 0.1%) {99.7%} 7031 O (14 = 0.1%) {99.9%} 9260 O (7 = 0.0%) {99.9%} $12196 O (2 = 0.0\%) \{100.0\%\}$ 16063 O (1 = 0.0%) {100.0%} 21156 O (0 = 0.0%) {100.0%} 27864 O (3 = 0.0%) {100.0%} 36699 ...



chrome is cool, but javascript uses float64

- Number.MAX_SAFE_INTEGER (2⁵³ -1)
- Number.MIN_SAFE_INTEGER (-(2⁵³ 1))
- Kernels tend to only do integer math
 - Integer math is more CPU efficient
 - Many performance counters are uint64
 - elapsed nanoseconds * # of I/Os = big number

SD @

IG

Example: ZFSonLinux

zpool iostat -w

testpool latency	tota read	al_wa d writ	it d e re	lisk_wa ad wri	it sy te re	yncq ead	_wait write 	asy rea	/ncq_ d wri	_wait ite scrub	trim
1ns	0	0	0 0) 0	0	0	0	0	0		
2us	0	0	0 0) 37	31	0	36	2	0		
4us	0	0	0 0) 33	25	0	117	39	0		
8us	0	0	0 0) 7	6	0	11	16	0		
16us	0	0	0	0 1	0	0	6	2	0		
32us	0	0	0	0 1	0	0	7	5	0		
65us	0	0	0	0 0	0	0	20	4	0		
131us	6	7	6	31 2	2 0	0	30	2	0		
262us	23	73	25	381	2	0	0	47	8	0	
524us	36	158	44	465	0	0	0	88	21	0	
1ms	40	194	56	187	0	0	0	143	37	0	
2ms	48	282	65	57	0	0	0 2	42	33	0	
4ms	56	320	99	9	0	0	0 2	57	50	0	
8ms	94	96	41	7	0 0) () 58	3 4	0 (C	
16ms	30	8	6	1 (0 (0	4	5	0		
33ms	9	0	5	0 0	0	0	0	6	0		

SD®

VIKING[™] Enterprise Solutions

Choosing Buckets



SD©

Simplest and fastest is log2

uint64_t elapsed = NOW() - start; bucket = 64 - __builtin_clzll(elapsed); // 64 - count leading zeroes

- ok for storage or many orders of magnitude
- sometimes more precision is desired
- Linear
 - ok when range is 1 or 2 orders of magnitude
 - division is 10-40x slower than addition/subtraction
 - range checking can cause branches (slow)
- Log-linear
 - mix of log10 with 10 buckets per decade
 - slower than linear, but good resolution
- Natural log
 - floating point not suitable for kernel





- Simplest implementation is a global array for latency buckets
 - used by OpenZFS
 - can lead to cache line stalls measure instructions per cycle (IPC)
 - can be space efficient
 - for consistency, use atomics but IPC can suffer
- Better implementation is one array per processor
 - eliminates cache line stalls = better IPC
 - no need for atomics
 - requires reader to sum across processors





- Getting time can be expensive
- uint64_t counters are convenient in kernel, but some tools, like prometheus, only support float64 values
- In-flight requests can't be easily measured
 - requests that do not complete have no end time
- Most iostat implementations are for HDDs
 - lack resolution for SSDs

One More Thing....

VIKING Enterprise Solutions

If you can do only one more thing...

- ... at least add a latency threshold counter
- Pick a (tunable) unexpected large latency threshold value
- Increment a counter when latency exceeds the threshold
- Make the counter accessible to instrumentation tools (node_exporter, telegraf, et.al.)

SD (9





- Average latencies are readily available, but are lacking in information
- Latency histograms can reveal latency problems
- Samplers can measure elapsed time between function calls or function entry/return
- But samplers are not suitable for running continuously
- Better solution is to collect latency histograms in kernel modules or applications
- We can do this! Let's do it!



Richard.Elling@VikingEntperise.com