

OPEN SOURCE DATABASES AND NON-VOLATILE MEMORY

Frank Ober

)@fxober @IntelStorage

May 29, 2019

NVM Solutions Group



Primer on Intel[®] Optane[™] technology and Intel[®] NVM

Benchmarks, Tools and Process

What can you do with it?



NVM Solutions Group



WHAT IS INTEL DOING WITH NVM?

INTEL® OPTANE[™] TECHNOLOGY AND INTEL® 3D NAND SSD Changing the Hierarchy





Latency in Hardware and Human Terms

- Memory Operations 64B cache line like getting an apple from the refrigerator
- **Storage Operations** 512B or 4096B typical storage block, which is like getting a case of apples from the cellar
- Most storage solutions are unfortunately in the milliseconds still today. Applications are typically characterized in milliseconds.

A blink of an eye is about 100 milliseconds

Sources – Other – these numbers will and can vary by generation of hardware, this Is purely an exercise in understanding relative performance at the synthetic Hardware level.

HW Operation	Latency
L1 CPU Cache	1 ns (nanoseconds)
L2 CPU Cache	3 ns
L3 CPU Cache	10-20 ns
DRAM	70-100 ns
Intel® Optane™ DC PMM	< 1 us
Intel® Optane™ SSD	< 10 us (micro secs)
SSD SLC NAND TLC/QLC is slower	25 us
7200RPM HDD	8 - 10 ms
Network Latency New York to Europe	40 milliseconds



INTRODUCTION TO PERSISTENT MEMORY

- What is Persistent Memory?
 - Byte Addressable
 - Cache Coherent
 - Load /Store Access
 - No page caching
 - Memory-like Performance

- Why Does it Matter Now?
 - Adds a new tier between DRAM and Block Storage (SSD/HDD)
 - Large Capacity, High Endurance, Consistent low latency
 - Ability to do in-place persistence
 - No Paging, No Context Switching, No Interrupts, No Kernel Code
 - Ability to do DMA & RDMA

- What's the Impact on the Applications?
 - Need Ways to Enable Access to New High Performance Tier
 - May need to rearchitect to unlock the new features and performance

http://pmem.io/



INTEL® OPTANE™ DC PERSISTENT MEMORY – INTENDED VALUES & BENEFITS

Plenty and affordable memory

High performance storage (latency, bandwidth, QoS, endurance)

Application managed memory



More and extended VMs Targeting > 1.2X VM at cost parity₁

Capacity for In-Memory Database at near-DRAM performance

Super-fast storage Targeting > 3X NVMe* performance

Targeting larger memory pools **Up to 3TB (not including DRAM)**

PROVIDING LOWER AND CONSISTENT LATENCY WITH MORE CAPACITY PER DOLLAR

*other names and brands my be claimed as the property of others. † Not all VM's are same and pending on usage and application, the results may differ.



INTEL® OPTANE™ DC PERSISTENT MEMORY - PRODUCT OVERVIEW

(Intel[®] Optane[™] SSD-based Memory Module for the Data Center)

Flexible, Usage Specific Partitions





NOW FOR INTEL® OPTANE[™] SSD

CACHING: INTEL® OPTANE™ SSD DC P4800X. THE IDEAL CACHING SOLUTION.



Low latency + high endurance greater SDS system efficiency

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit <u>www.intel.com/benchmarks.</u>

Source – Intel-tested: Average read latency measured at queue depth 1 during 4k random write workload. Measured using FIO 3.1. Common Configuration – Intel 2U Server System, OS CentOS 7.5, kernel 4.17.6-1.el7.x86_64, CPU 2 x Intel® Xeon® 6154 Gold @ 3.0GHz (18 cores), RAM 256GB DDR4 @ 2666MHz. Configuration – Intel® Optane™ SSD DC P4800X 375GB and Intel® SSD DC P4600 1.6TB. Latency – Average read latency measured at QD1 during 4K Random Write operations using FIO 3.1. Intel Microcode: 0x2000043; System BIOS: 00.01.0013; ME Firmware: 04.00.04.294; BMC Firmware: 1.43.91f76955; FRUSDR: 1.43. SSDs tested were commercially available at time of test. Performance results are based on testing as of July 24, 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No product can be absolutely secure.

2. Source – Intel: Endurance ratings available at https://www.intel.com/content/www/us/en/products/docs/memory-storage/solid-state-drives/data-center-ssds/optane-ssd-dc-p4800x-p4801x-brief.html

3. Source – Intel: General proportions shown for illustrative purposes.



USAGE MODELS FOR OPTANE IN DATABASES

Create a Caching drive over SATA or QLC SSDs (Intel[®] Cache Acceleration Software)

Tier in Intel[®] Optane[™] SSDs for specific cases

Write Logs or High Ingest drive

Temp Data drive





BREAK TIME: WHAT ABOUT FSYNC ON INTEL® OPTANE™ SSD?

•• Fsync Performance on Storage Devices

🔒 By Yves Trudeau 🞥 Hardware and Storage, InnoDB, Insight for DBAs, MySQL

Feb

2018 🔷 🗞 fsync, fsync performance, InnoDB, MySQL, rotating drives, SSD, Storage 🧔 8 Comments

While preparing a post on the design of ZFS based servers for use with MySQL, I stumbled on the topic of fsync call performance. The fsync call is very expensive, but it is essential to databases as it allows for *durability* (the "D" of the ACID acronym).

Let's first review the type of disk IO operations executed by InnoDB in MySQL. I'll assume the default InnoDB variable values.

The first and most obvious type of IO are pages reads and writes from the tablespaces. The pages are most often read one at a time, as 16KB random read operations. Writes to the tablespaces are also typically 16KB random operations, but they are done in batches. After every batch, fsync is called on the tablespace file handle.



Source: Percona – used with permission of Percona, *Other names and brands may be claimed as the property of others.





SYSBENCH AND OTHER TESTING

IDENTIFY STORAGE-BOUND WORKLOADS TO OPTIMIZE YOUR PLATFORM

Where average data center workloads spends their time:

CPU 13	BN21 —	→ « CPU IS NI	JI RO2 Å	·>
Арр	Sys	IOwait	Idle	DISK OR Network

Running the application
Running the OS/VMM

• Waiting for Disk

- Limited application parallelism
- Imbalances
- Waiting for network
- Indirectly waiting for network (hidden IOwaits)

INTEL VTUNE PLATFORM PROFILER 2019 (FREE WITHOUT SUPPORT)



TPCC ON NVME* DEVICES.... INTEL® SSD DC P4510 (NVME NAND)

1 [23.0%] 19 [20.8%] 37 [22.7%] 55 [23.3%
2 [22.7%] 20 [20.3%] 38 [20.8%] 56 [22.7%
3 [22.7%] 21 [22.4%] 39 [21.1%] 57 [22.0%
4 [21.9%] 22 [23.8%] 40 [21.7%] 58 [21.1%
5 [20.1%] 23 [20.8%] 41 [20.7%] 59 [19.9%
6 [18.7%] 24 [20.5%] 42 [21.9%] 60 [22.7%
7 [20.5%] 25 [19.3%] 43 [22.0%] 61 [19.3%]
8 [22.9%] 26 [23.2%] 44 [21.5%] 62 [21.7%]
9 [20.1%] 27 [22.2%] 45 [21.5%] 63 [23.8%]
10 [21.7%] 28 [20.0%] 46 [20.1%] 64 [21.2%
11 [22.9%] 29 [21.7%] 47 [18.1%] 65 [22.5%
12 [23.0%] 30 [19.9%] 48 [21.7%] 66 [22.7%
13 [22.1%] 31 [20.5%] 49 [22.9%] 67 [23.8%]
14 [20.3%] 32 [20.5%] 50 [19.9%] 68 [20.7%
15 [22.5%] 33 [22.5%] 51 [21.2%] 69 [21.2%
16 [22.4%] 34 [21.1%] 52 [21.5%] 70 [19.2%
17 [23.8%] 35 [23.3%] 53 [22.0%] 71 [22.1%]
18 [24.2%] 36 [19.9%] 54 [20.5%] 72 [20.0%
Mem[<pre> 34.5G/188G] Tasks: 141, 583 thr; 30</pre>	6 running	
Swp[0K/0K] Load average: 28.23 7.4	47 2.94	
		Uptime: 2 days, 05:10:3	39	

PID USER	R PRI	. NI	L VIRI	RES	SHR	S CPU%	ME MP8	IIME+	Command
3482 myso	ql 20) (9 40.5G	32.5G	18760	S 1407	17.3	13:15.25	/usr/sbin/mysqlddaemonizepid-file=/var/run/mysqld/mysqld.pid
3557 root	t 20) (6	9 4306M	61276	<u>6</u> 624	S 185.	Θ.Θ	1:47.93	sysbench tpccthreads=64time=300tables=10percentile=99db-driver=mysqlmysql-db=

avg-cpu:	%user 27.15	%nice 0.00	e %system) 5.85	%iowait 20.04	%steal 0.00	%idle 46.96								
Device:		rrqm/s	wrqm/s	r/s	w/s	rMB/s	wMB/s	avgrq-sz	avgqu-sz	await r	await w	await	svctm	%util
nvmelnl		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
nvme0n1		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
nvme2n1		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
nvme3n1		0.00	14468.00	27727.00	29082.00	433.23	704.1	0 41.0	0 41.63	0.40	0.24	0.56	6 0.0	2 101.10

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visitwww.intel.com/benchmarks.

¹ System configuration: Intel Sourced: Tested: October 3, 2018. Server Intel[®] Server System , 2x Intel[®] Xeon[®] Scalable 6154, 384 GB DDR4 DRAM, database drives- 3x Intel[®] SSD DC P4800X Series (375 GB) and 1x Intel[®] SSD DC P4510 Series, CentOS 7.5 (kernel 4.18.8 (from elrepo)), BIOS: SE5C620.86B.00.01.0014.070920180847 system product type: S2600WFT

Percona MySQL Server 5.7.23, Sysbench 1.0.15 configured TPCC / sysbench per github project https://github.com/Percona-Lab/sysbench-tpcc 100GB database size. 30% Database Memory provided to MySQL (30GB). Performance results are based on testing as of Oct 3, 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No component or product can be absolutely secure. *Other names and brands names may be claimed as the property of others



TPCC ON INTEL® OPTANE™ DC SSDs (NVME* INTEL® OPTANE™ MEMORY MEDIA)

1 [2 [80.19] 19 [] 20 [80. 	.4%] 37 .0%] 38				78.9%	55 [56 [79	9.9%] 0.0%]
3 [79.39	21 [79.	7%] 39				77.2%	57 []				78	8.5%
4			78.09				77.	2% 40				77.9%	58 []					8.0%
5 1111			80.4	6] 23 [] 24 [80.					78.0%	59 []					8./%]
				l 24 [l 25 [20	/s] 42				70 08	61 1				7	7 28
			77 59	25 [] 26 [78	43 42 28 44				77 4%	62 []					7.3%] 8.1%]
g iiii			80.19	1 27			78	5%] 45				77.5%	63 []				7	8.4%
10 [78.99	28 [1111180	1%] 46				78.5%	64 [İ	iiiiiii			1111 78	8.5%
11 [79.69	29 [79.	6% 47				78.3%]	65 []	iiiiii			78	8.3%]
12 [78.39	6] 30 [79.	.3%] 48				77.9%	66 []				11179	9.3%]
13 [78.49	31 [79.	5%] 49				77.6%	67 []				78	8.0%
14 [79.5	32 [80.	.3%] 50				78.8%	68 []				7	7.6%
15 []]]			80.49	6 33 L			79.	5% 51				78.2%	69 LI				78	8.9%
			79.9	34 [] 35 [79.	5% 52				77.6%	70 []					7.9%
10			79.93	6] 30 [] 26 [80.	1%] D3				70.0%	72 []					8.4% 7.6%
Mom								286] Te	eke + 120	592 thre	72 rupp	70.0%]	/2 []					/.0%
							11104.40/10	1001	iaka, 100,	JO2 CHI,	72 TUITIT	Ling						
Swn									ad average	e: 29.25	16.79 13	.09						
Swp [/0K] Lo Up	ad average time: 2 d	e: 29.25 avs. 05:2	16.79 13. 1:56	.09						
Swp [/0K] Lo Up	ad average time: 2 d	e: 29.25 ays, 05:2	16.79 13 1:56	.09						
Swp[PID USEF	R_ PRI M	NI VIRT	RES	SHR S C	PU%s <mark>MEM%s</mark>	TIME+	өк, Command	/0K] Lo Up	ad average time: 2 da	e: 29.25 ays, 05:2	16.79 13 1:56	.09						
Swp[PID USEF 4077 mysc	R PRI M	VI VIRT 0 40.5G	RES 32.5G	SHR S CI 18760 S 50	PU% <mark>MEM%</mark> 066 17.3	TIME+ 26:39.81	OK/ Command /usr/sbin/n	/OK] Lo Up nysqld	ad average otime: 2 da	e: 29.25 ays, 05:2	16.79 13 1:56 le=/var/1	.09 run/mysq	ld/mysql	d.pid				
Swp[PID USEF 4077 mysc 4140 root	R PRI M ql 20 t 20	VI VIRT 0 40.5G 0 4304M	RES 32.5G 1 60876	SHR S CI 18760 S 50 6944 S 65	PU% <mark>MEM%</mark> 966 17.3 3 59. 0.0	TIME+ 26:39.81 3:31.89	OK/ Command /usr/sbin/n sysbench tp	VOK] Lo Up nysqld occthr	oad averago otime: 2 da daemonize eads=64 -	e: 29.25 ays, 05:2 pid-fi -time=300	16.79 13 1:56 le=/var/1 tables	.09 run/mysq s=10pe	<mark>ld/mysql</mark> ercentil	<mark>d.pid</mark> e=99 -	-db-dr:	iver=mys	qlmy	ysql-d
Swp[PID USEF 4077 mysc 4140 root	R PRI M ql 20 t 20 avg-cpu	NI VIRT 0 40.5G 0 4304M : %user	RES 32.56 1 60876 %nice	<mark>SHR S</mark> C 18760 S 50 6944 S 63 e %svstem	PU% <mark>MEM%</mark> 066 17.3 59. 0.0 %iowait	TIME+ 26:39.81 3:31.89 %steal	Command /usr/sbin/n sysbench tp %idle	YOK] Lo Up nysqld pccthr	oad averago otime: 2 da daemonize eads=64 -	e: 29.25 ays, 05:2 pid-fi -time=300	16.79 13 1:56 <mark>le=/var/1</mark> tables	.09 run/mysq s=10pe	ld/mysql ercentil	<mark>d.pid</mark> e=99 -	-db-dr	iver=mys	qlmy	ysql-d
Swp[PID USEF 4077 myso 4140 root	R PRI M ql 20 t 20 avg-cpu:	VI VIRT 0 40.5G 0 4304M : %user 64.29	RES 32.5G 1 60876 %nice 0.00	<mark>SHR S</mark> CF 18760 S 50 6944 S 65 e %system 0 11.54	PU% <mark>MEM%</mark> 966 17.3 3 59. 0.0 %iowait 1.13	TIME+ 26:39.81 3:31.89 %steal 0.00	Command /usr/sbin/n sysbench tp %idle 23.04	YOK] Lo Up nysqld Doccthr	ad averago time: 2 da daemonize reads=64 -	e: 29.25 ays, 05:2 pid-fi -time=300	16.79 13 1:56 le=/var/1 tables	.09 <mark>run/mysq</mark> s=10pe	ld/mysql ercentil	<mark>d.pid</mark> e=99 -	-db-dr:	iver=mys	qlmy	ysql-d
Swp[PID USEF 4077 myso 4140 root	R PRI M ql 20 t 20 avg-cpu:	VI VIRT 0 40.5G 0 4304M %user 64.29	RES 32.5G 1 60876 %nice 0.00	<mark>SHR S CI 18760 S 50</mark> 6944 S 65 e %system 0 11.54	9 <mark>0% MEM%</mark> 966 17.3 3 59. 0.0 %iowait 1.13	TIME+ 26:39.81 3:31.89 %steal 0.00	Command /usr/sbin/n sysbench tp %idle 23.04	YOK] Lo Up nysqld Doccthr	oad averago otime: 2 da daemonize eads=64 -	e: 29.25 ays, 05:2 pid-fi -time=300	16.79 13 1:56 le=/var/f tables	.09 <mark>run/mysq</mark> s=10pe	<mark>ld/mysql</mark> ercentil	<mark>d.pid</mark> e=99	-db-dr:	iver=mys	qlmy	ysql-d
Swp[PID USEF 4077 mysc 4140 root	R PRI M qL 20 t 20 avg-cpu: Device:	VI VIRT 0 40.5G 0 4304M %user 64.29	RES 32.5G 1 60876 %nice 0.00 rrqm/s	SHR S CF 18760 S 50 6944 S 65 e %system 0 11.54 wrqm/s	PU% MEM% 966 17.3 59. 0.0 %iowait 1.13 r/s	TIME+ 26:39.81 3:31.89 %steal 0.00 w/s	Command /usr/sbin/n sysbench tp %idle 23.04 rMB/s	wMB/s a	oad averago otime: 2 da daemonize eads=64 -	e: 29.25 ays, 05:2 pid-fi -time=300 /gqu-sz	16.79 13 1:56 le=/var/n tables await r_	.09 <mark>run/mysq</mark> s=10pe await w_	<mark>ld/mysql</mark> ercentil await	<mark>d.pid</mark> e=99 svctm	-db-dr: %util	iver=mys	qlmy	ysql-d
Swp[PID USEF 4077 mysc 4140 root	R PRI M ql 20 t 20 avg-cpu: Device: nvmelnl	VIRT 0 40.5G 0 4304M %user 64.29	RES 32.5G 1 60876 %nice 0.00 rrqm/s 0.00	SHR S CF 18760 S 50 6944 S 65 e %system 0 11.54 wrqm/s 9300.00	PU% MEM% 966 17.3 3 59. 0.0 %iowait 1.13 r/s 14917.00	TIME+ 26:39.81 3:31.89 %steal 0.00 w/s 19332.00	Command /usr/sbin/m sysbench tp %idle 23.04 rMB/s 233.08	wMB/s a 470.99	vgrq-sz av 42.10	e: 29.25 ays, 05:2 pid-fi -time=300 /gqu-sz 1.64	16.79 13 1:56 tables await r_ 0.00	.09 <mark>run/mysq</mark> s=10pe await w_ 0.00	<mark>ld/mysql</mark> ercentil await 0.00	<mark>d.pid</mark> e=99 svctm 0.02	-db-dr: %util 2 54.0	iver=mys	qlmy	ysql-d
Swp[PID USEF 4077 mysc 4140 root	R PRI M ql 20 t 20 avg-cpu: Device: nvmeln1 nvme0n1	VIRT 0 40.5G 0 4304M %user 64.29	RES 32.5G 1 60876 %nice 0.00 rrqm/s 0.00 0.00	SHR S C 18760 S 50 6944 S 65 e %system 0 11.54 wrqm/s 9300.00 9344.00	PU% MEM% 966 17.3 59. 0.0 %iowait 1.13 r/s 14917.00 14801.00	TIME+ 26:39.81 3:31.89 %steal 0.00 w/s 19332.00 19229.00	Command /usr/sbin/m sysbench tp %idle 23.04 rMB/s 233.08 231.27	WMB/s a 470.10	vgrq-sz av 42.10 42.21	e: 29.25 ays, 05:2 pid-fi -time=300 /gqu-sz 1.64 1.58	16.79 13 1:56 tables await r_ 0.00 0.00	.09 <mark>run/mysq</mark> s=10pe await w_ 0.00 0.00	<mark>ld/mysql</mark> ercentil await 0.00 0.00	<mark>d.pid</mark> e=99 svctm 0.02 0.02	-db-dr: %util 2 54.0 2 56.1	iver=mys 00 10	qlmy	ysql-d
Swp[PID USEF 4077 mysc 4140 root	R PRI M ql 20 t 20 avg-cpu Device: nvmeln1 nvme0n1 nvme2n1	VI VIRT 0 40.5G 0 4304M %user 64.29	RES 32.5G 1 60876 %nice 0.00 rrqm/s 0.00 0.00 0.00	SHR S CF 18760 S 50 6944 S 65 e %system 9 11.54 wrqm/s 9300.00 9344.00 9204.00	PU% MEM% 966 17.3 59. 0.0 %iowait 1.13 r/s 14917.00 14801.00 14714.00	TIME+ 26:39.81 3:31.89 %steal 0.00 w/s 19332.00 19229.00 19224.00	Command /usr/sbin/m sysbench tp %idle 23.04 rMB/s 233.08 231.27 229.91	<pre>v0K] Lo Up mysqld occthr wMB/s ar 470.99 470.10 470.96</pre>	vgrq-sz av 42.10 42.29	e: 29.25 ays, 05:2 pid-fi -time=300 /gqu-sz 1.64 1.58 1.60	16.79 13 1:56 tables await r_ 0.00 0.00 0.00	.09 run/mysq s=10pe await w_ 0.00 0.00 0.00	await 0.00 0.00 0.00	<mark>d.pid</mark> e=99 - svctm 0.02 0.02	-db-dr: %util 2 54.0 2 56.1 2 55.5	iver=mys 00 10	qlmy	ysql-d
Swp[PID USEF 4077 mysc 4140 root	R PRI M ql 20 avg-cpu Device: nvmelnl nvme0nl nvme2nl nvme3nl	VI VIRT 0 40.5G 0 4304M : %user 64.29	RES 32.5G 1 60876 %nice 0.00 rrqm/s 0.00 0.00 0.00 0.00 0.00	SHR S CF 18760 S 50 6944 S 65 e %system 9 11.54 wrqm/s 9300.00 9344.00 9204.00 0.00	PU% MEM% 966 17.3 59. 0.0 %iowait 1.13 r/s 14917.00 14801.00 14714.00 0.00	TIME+ 26:39.81 3:31.89 %steal 0.00 w/s 19332.00 19229.00 19224.00 0.00	Command /usr/sbin/m sysbench tp %idle 23.04 rMB/s 233.08 231.27 229.91 0.00	<pre>wMB/s a 470.99 470.96 0.00</pre>	vgrq-sz av 42.10 42.21 0.00	e: 29.25 ays, 05:2 pid-fi -time=300 /gqu-sz 1.64 1.58 1.60 0.00	16.79 13 1:56 le=/var// tables await r_ 0.00 0.00 0.00 0.00	.09 run/mysq s=10pe .00 0.00 0.00 0.00	await 0.00 0.00 0.00 0.00	d.pid e=99 - svctm 0.02 0.02 0.02 0.02	-db-dr: %util 2 54.0 2 56.1 2 55.5 0.00	iver=mys 00 10 50	qlmy	ysql-d
Swp[PID USEF 4077 mysc 4140 root	R PRI M ql 20 t 20 avg-cpu Device: nvmelnl nvme0nl nvme2nl nvme3nl sda	VI VIRT 0 40.5G 0 4304M : %user 64.29	RES 32.56 60876 %nice 0.00 rrqm/s 0.00 0.00 0.00 0.00 0.00 0.00	SHR S CF 18760 S 50 6944 S 63 e %system 0 11.54 wrqm/s 9300.00 9344.00 9204.00 0.00 0.00	PU% MEM% 966 17.3 59. 0.0 %iowait 1.13 r/s 14917.00 14801.00 14714.00 0.00 0.00	TIME+ 26:39.81 3:31.89 %steal 0.00 w/s 19332.00 19229.00 19224.00 0.00 0.00	Command /usr/sbin/m sysbench tp %idle 23.04 rMB/s 233.08 231.27 229.91 0.00 0.00 0.00	WMB/s a 470.99 470.10 0.00 0.00	vgrq-sz av 42.10 42.21 0.00 0.00	e: 29.25 ays, 05:2 pid-fi -time=300 /gqu-sz 1.64 1.58 1.60 0.00 0.00	16.79 13 1:56 await r_ 0.00 0.00 0.00 0.00 0.00	.09 run/mysq s=10pe .00 0.00 0.00 0.00 0.00	await 0.00 0.00 0.00 0.00 0.00	d.pid e=99 - svctm 0.02 0.02 0.02 0.00 0.00	-db-dr: %util 2 54.0 2 56.1 2 55.5 0.00 0.00	iver=mys 00 10 50	qlmy	ysql-d
Swp[PID USEF 4077 myso 4140 root	R PRI 1 ql 20 t 20 avg-cpus Device: nvmeln1 nvme0n1 nvme2n1 nvme3n1 sda sdb	VIRT 0 40.5G 0 4304M %user 64.29	RES 32.5G 1 60876 %nice 0.00 rrqm/s 0.00 0.00 0.00 0.00 0.00 0.00 0.00	SHR S CF 18760 S 50 6944 S 63 e %system 0 11.54 wrqm/s 9300.00 9344.00 9204.00 0.00 0.00 0.00	DU% MEM% 366 17.3 59. 0.0 %iowait 1.13 r/s 14917.00 14801.00 14714.00 0.00 0.00 0.00	TIME+ 26:39.81 3:31.89 %steal 0.00 w/s 19332.00 19229.00 19224.00 0.00 0.00 0.00	Command /usr/sbin/m sysbench tr %idle 23.04 rMB/s 233.08 231.27 229.91 0.00 0.00 0.00	WMB/s a 470.99 470.10 470.90 0.00 0.00	vgrq-sz av 42.10 42.21 0.00 0.00 0.00	e: 29.25 ays, 05:2 pid-fi -time=300 /gqu-sz 1.64 1.58 1.60 0.00 0.00 0.00 0.00	16.79 13 1:56 await r_ 0.00 0.00 0.00 0.00 0.00 0.00 0.00	.09 run/mysq s=10pe 0.00 0.00 0.00 0.00 0.00 0.00	await 0.00 0.00 0.00 0.00 0.00 0.00	d.pid e=99 - svctm 0.02 0.02 0.02 0.00 0.00 0.00	-db-dr: %util 2 54.0 2 56.1 2 55.5 0.00 0.00 0.00	iver=mys 00 10 50	qlmy	ysql-d

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that products. For more complete information visit<u>www.intel.com/benchmarks</u>.

¹ System configuration: Intel Sourced: Tested: October 3, 2018. Server Intel[®] Server System , 2x Intel[®] Xeon[®] Scalable 6154, 384 GB DDR4 DRAM, database drives- 3x Intel[®] SSD DC P4800X Series (375 GB) and 1x Intel[®] SSD DC P4510 Series, CentOS 7.5 (kernel 4.18.8 (from elrepo)), BIOS: SE5C620.86B.00.01.0014.070920180847 system product type: S2600WFT

Percona MySQL Server 5.7.23, Sysbench 1.0.15 configured TPCC / sysbench per github project https://github.com/Percona-Lab/sysbench-tpcc 100GB database size. 30% Database Memory provided to MySQL (30GB). Performance results are based on testing as of Oct 3, 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No component or product can be absolutely secure. *Other names and brands may be claimed as the property of others.

17

o=tpc

PERCONA* MySQL* 5.7 TPCC TEST (64 THREADS, 10 TABLES, SCALE 100)



Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit<u>www.intel.com/benchmarks</u>.

¹ System configuration: Intel Sourced: Tested: October 3, 2018. Server Intel® Server System , 2x Intel® Xeon® Scalable 6154, 384 GB DDR4 DRAM, database drives- 3x Intel® SSD DC P4800X Series (375 GB) and 1x Intel® SSD DC P4510 Series, CentOS 7.5 (kernel 4.18.8 (from elrepo)), BIOS: SE5C620.86B.00.01.0014.070920180847 system product type: S2600WFT

Percona MySQL Server 5.7.23, Sysbench 1.0.15 configured TPCC / sysbench per github project https://github.com/Percona-Lab/sysbench-tpcc 100GB database size. 30% Database Memory provided to MySQL (30GB). Performance results are based on testing as of Oct 3, 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No component or product can be absolutely secure.
*Other names and brands names may be claimed as the property of others

ORACLE* MySQL* 8.0 OLTP RW TEST (32 THREADS, 8 TABLES X 50M ROWS) INTEL® CACHE ACCELERATION SOFTWARE (INTEL® CAS) FOR LINUX 3.6



TESTED WITH PARETO 80/20 RANDOMIZATION

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information wisit<u>www.intel.com/benchmarks</u>.

¹ System configuration: Intel Sourced: Tested: October 20, 2018. Server Intel[®] Server System, 2x Intel[®] Scalable 6154, 384 GB DDR4 DRAM, database drives- 2x Intel[®] SSD DC P4800X Series (375 GB) and 1x Intel[®] SSD DC P4510 Series, 1x Intel[®] SSD DC S4510 Series, CentOS 7.5 (kernel 4.18 (elrepo)), BIOS: SE5C620.86B.00.01.0014.070920180847 system product type: S2600WFT, Intel version of Intel CAS used version 3.6 set in write back mode over a single S4510 SSD.

MySQL Server 8.0.13, Sysbench 1.0,15 configured for 70/30 Read/Write OLTP transaction split using a 100GB database. 30% Database Memory provided to MySQL (30GB). Performance results are based on testing as of [Oct 20, 2018] and may not reflect all publicly available security updates. See configuration disclosure for details. No component or product can be absolutely secure. *Other names and brands names may be claimed as the property of others



TIDB – LOOKING AT AN OPEN HTAP DB



TiKV Architecture



Source: TiDB – used with permission of TiDB. *Other names and brands may be claimed as the property of others.







Source – TiDB – used by permission from TiDB - *Other names and brands may be claimed as the property of others.



Performance of TiKV (Txn Gets and Puts)



threads

Source – TiDB – used by permission from TiDB - *Other names and brands may be claimed as the property of others.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that products. For more complete information visit<u>www.intel.com/benchmarks</u>.

¹ System configuration: PingCap Sourced: Tested: November 20, 2018. 5 Server Cluster - Server Intel® Server System , 2x Intel® Xeon® Scalable 6142 (64 vcpus per server), BIOS: SE5C620.86B.00.01.0014.070920180847 system product type: S2600WFT 384 GB DDR4 DRAM, database drives- 2x Intel® SSD DC P4800X Series (375 GB) CentOS 7.5 (kernel 4.18.5 (elrepo)), TiDB Server v3.0.0-beta-204-gc0d4632fc, 3 PD Servers in the cluster, 3 DB engines across 5 Nodes (15 DBs)



ACCESS TO INTEL® OPTANE[™] SSDs AND OTHER TOOLS

WHAT BENCHMARKS MY TEAM USES...

sysbench -

– <u>https://github.com/akopytov/sysbench</u>

HammerDB - https://www.hammerdb.com/download.html

Redis we focus mostly on using memtier from https://github.com/RedisLabs/memtier_benchmark

What about YCSB ?? Sure, also capable ...

Use your own top 10 queries in SQL... and data schema that matches your goals

TOOLS AND GUIDES...

How to efficiently test Intel[®] Optane[™] SSDs (Intel[®] Optane[™] SSD Optimization Guide)

<u>https://itpeernetwork.intel.com/tuning-performance-intel-optane-ssds-linux-operating-systems/</u>

How to start with SPDK and VHOST (Virtualization with low latency SSDs)

– <u>http://spdk.io/</u>

More info CAS-Linux and Intel[®] Memory Drive Technology

<u>http://intel.com/cas</u> and the Open Source version, <u>https://open-cas.github.io/</u>

Persistent Memory Programming - http://pmem.io/

Performance Analysis Tool - <u>https://github.com/intel-hadoop/PAT</u>



Thanks for listening!

frank.ober@intel.com @fxober http://intel.com/optane

😂 Wednesday 4:15 PM - 5:05 PM

@ Hill Country A

An Open-Source, Cloud Native Database (CNDB)

David Cohen (Intel), Steve Shaw - Intel, Yekesa Kosuru - State Street, Jiten Vaidya - PlanetScale

INTEL® VTUNE AMPLIFIER - PLATFORM PROFILER 2019

Finds

em0

0000000 _______eth1

nvme0n1

- Configuration issues
- Poorly tuned software

Target Users

Infrastructure Architects

Memory Controller 1

Socket 0

Software Architects & QA

Server Topology Overview

Memory Controller 0

Socket

Memory Controller



- Extended capture (minutes to hours)
- Low overhead coarse grain metrics
- Sampling OS & hardware performance counters
- RESTful API for easy analysis by scripts





Memory Controller 0

