

Evaluating NVMe drives for accelerating HBase

Nicolas Poggi and David Grier

BSC Data Centric Computing – Rackspace collaboration



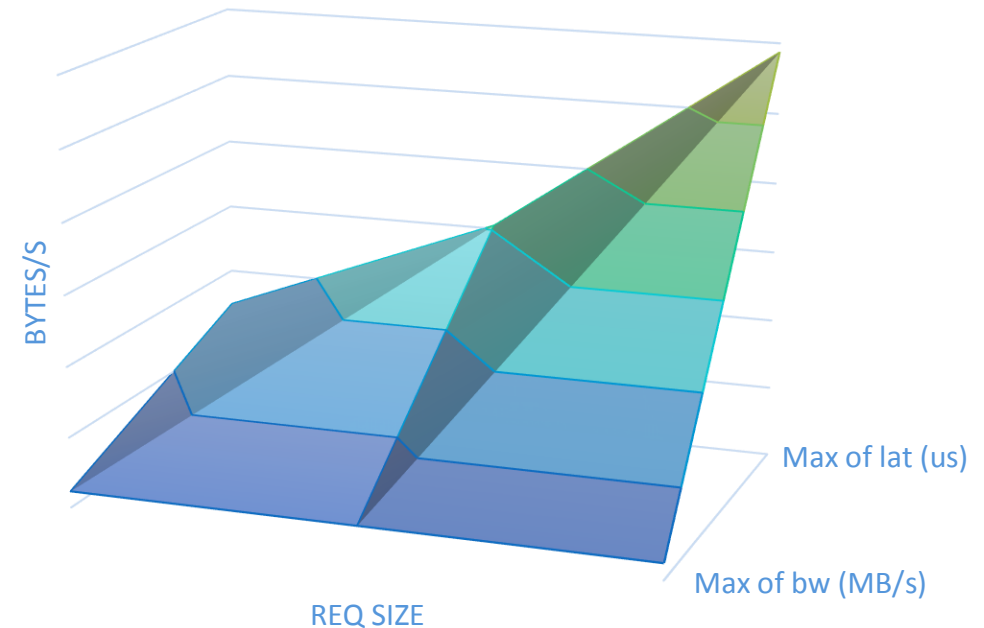
**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación

FOSDEM Jan 2017

aloja 

Outline

1. Intro on BSC and ALOJA
2. Cluster specs and disk benchmarks
- 3. HBase use case with NVE**
 1. Read-only workload
 - Different strategies
 2. Mixed workload
4. Summary



Barcelona Supercomputing Center (BSC)



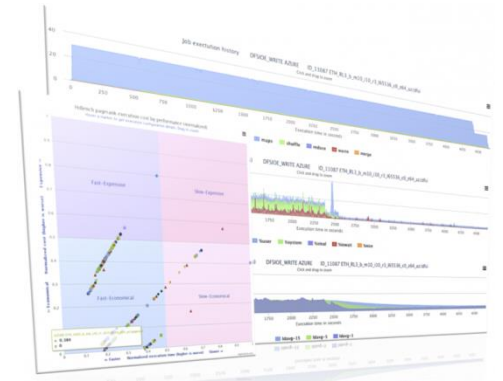
- Spanish national supercomputing center 22 years history in:
 - Computer Architecture, networking and distributed systems research
 - Based at BarcelonaTech University (UPC)
 - Large ongoing life science computational projects
- Prominent body of research activity around **Hadoop**
 - 2008-2013: SLA Adaptive Scheduler, Accelerators, Locality Awareness, Performance Management. **7+ publications**
 - 2013-Present: **Cost-efficient** upcoming Big Data architectures (**ALOJA**) **6+ publications**



ALOJA: towards cost-effective Big Data



- Research project for automating characterization and optimization of **Big Data** deployments
- Open source *Benchmarking-to-Insights* platform and tools



<http://aloja.bsc.es>

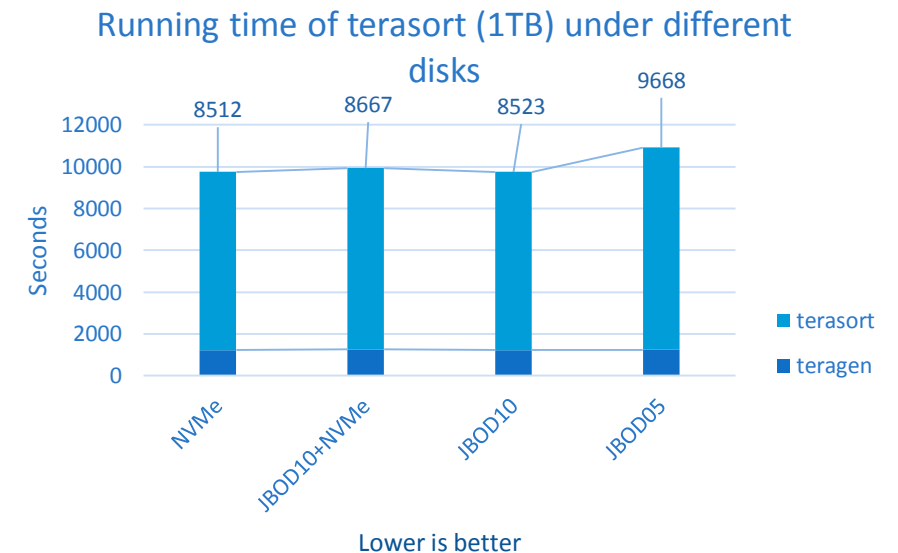


- Largest Big Data public repository (70,000+ jobs)
 - Community collaboration with industry and academia

Motivation and objectives

- Explore use cases where NVMe devices can speedup Big Data apps
 - **Poor initial results...**
 - HBase (this study) based on Intel report
- Measure the possibilities of NVMe devices
 - System level benchmarks (FIO, IO Meter)
 - WiP towards tiered-storage for Big data status
 - **Extend ALOJA into low-level I/O**
- Challenge
 - benchmark and stress high-end Big Data clusters
 - In reasonable amount of time (and cost)

First tests:



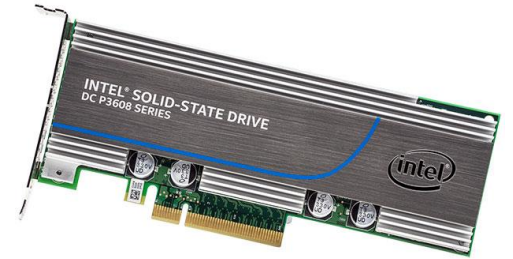
Marginal improvement!!!

Cluster and drive specs

All nodes (x5)	
Operating System	CentOS 7.2
Memory	128GB
CPU	Single Octo-core (16 threads)
Disk Config OS	2x600GB SAS RAID1
Network	10Gb/10Gb redundant
Master node (x1, extra nodes for HA not used in these tests)	
Disk Config Master storage	4x600GB SAS RAID10 XSF partition
Data Nodes (x4)	
NVMe (cache)	1.6TB Intel DC P3608 NVMe SSD (PCIe)
Disk config HDFS data storage	10x 0.6TB NL-SAS/SATA JBOD PCIe3 x8, 12Gb/s SAS RAID SEAGATE ST3600057SS 15K (XFS partition)

1. Intel DC P3608 (current 2015)

- PCI-Express 3.0 x8 lanes
- Capacity: 1.6TB (two drives)
- Seq R/W BW:
 - 5000/2000 MB/s, (128k req)
- Random 4k R/W IOPS:
 - 850k/150k
- Random 8k R/W IOPS:
 - 500k/60k, 8 workers
- Price **\$10,780**
 - (online search 02/2017)



2. LSI Nytro WarpDrive 4-400 (old gen 2012)

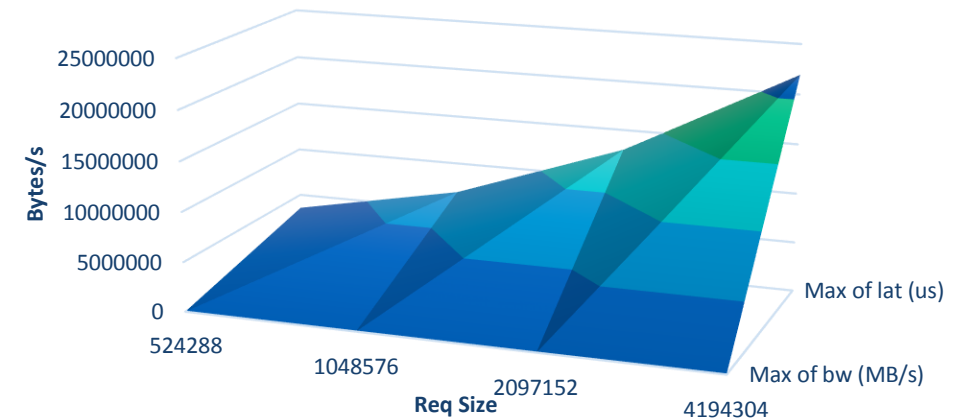
- PCI-Express 2.0 x8 lanes
- Capacity: 1.6TB (two drives)
- Seq. R/W BW:
 - 2000/1000 MB/s (256k req)
- R/W IOPS:
 - 185k/120k (8k req)
- Price **\$4,096**
 - (online search 02/2017)
 - \$12,195 MSRP 2012



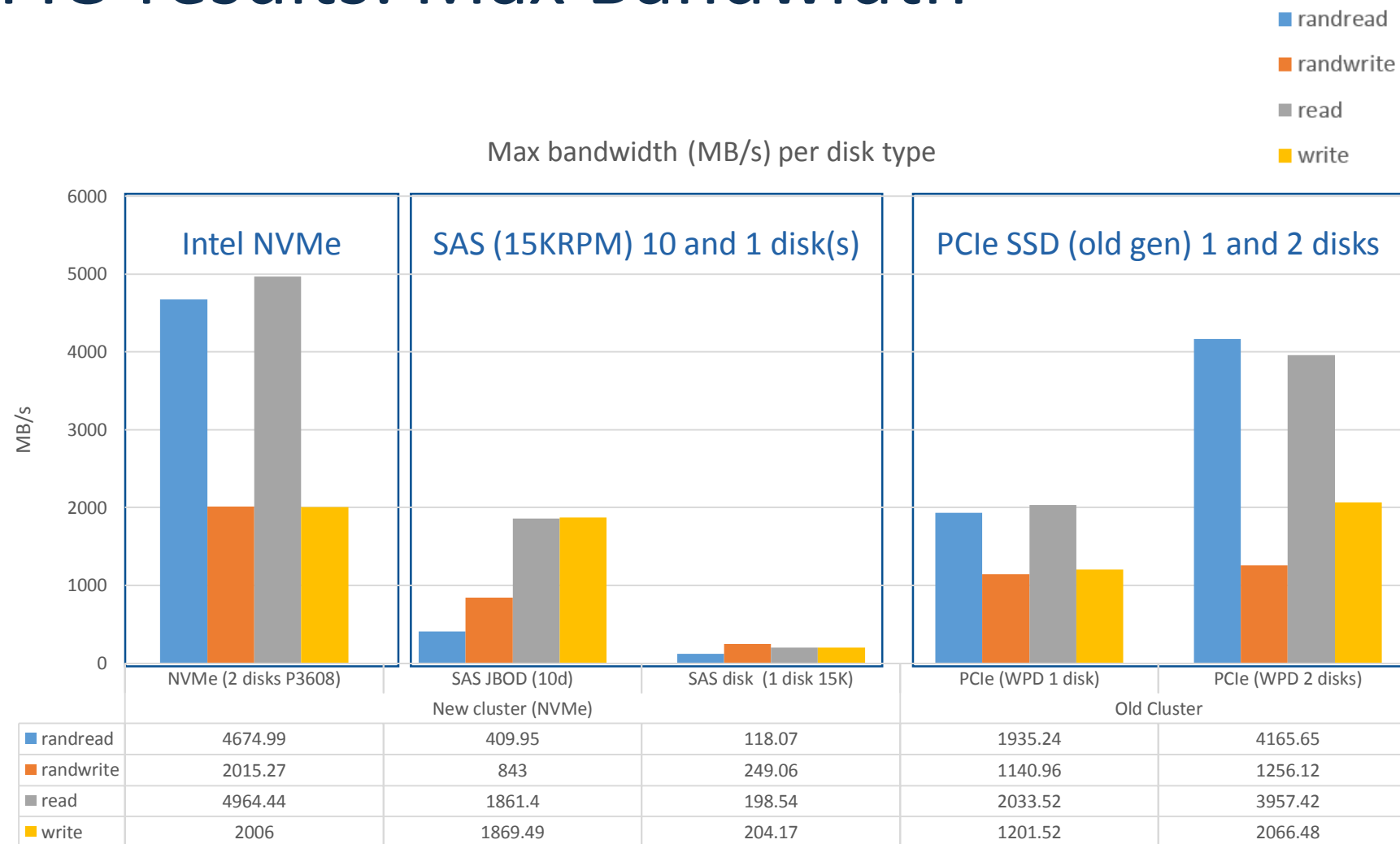
FIO Benchmarks

Objectives:

- Assert vendor specs (BW, IOPS, Latency)
- Seq R/W, Random R/W
- Verify driver/firmware and OS
- Set performance expectations



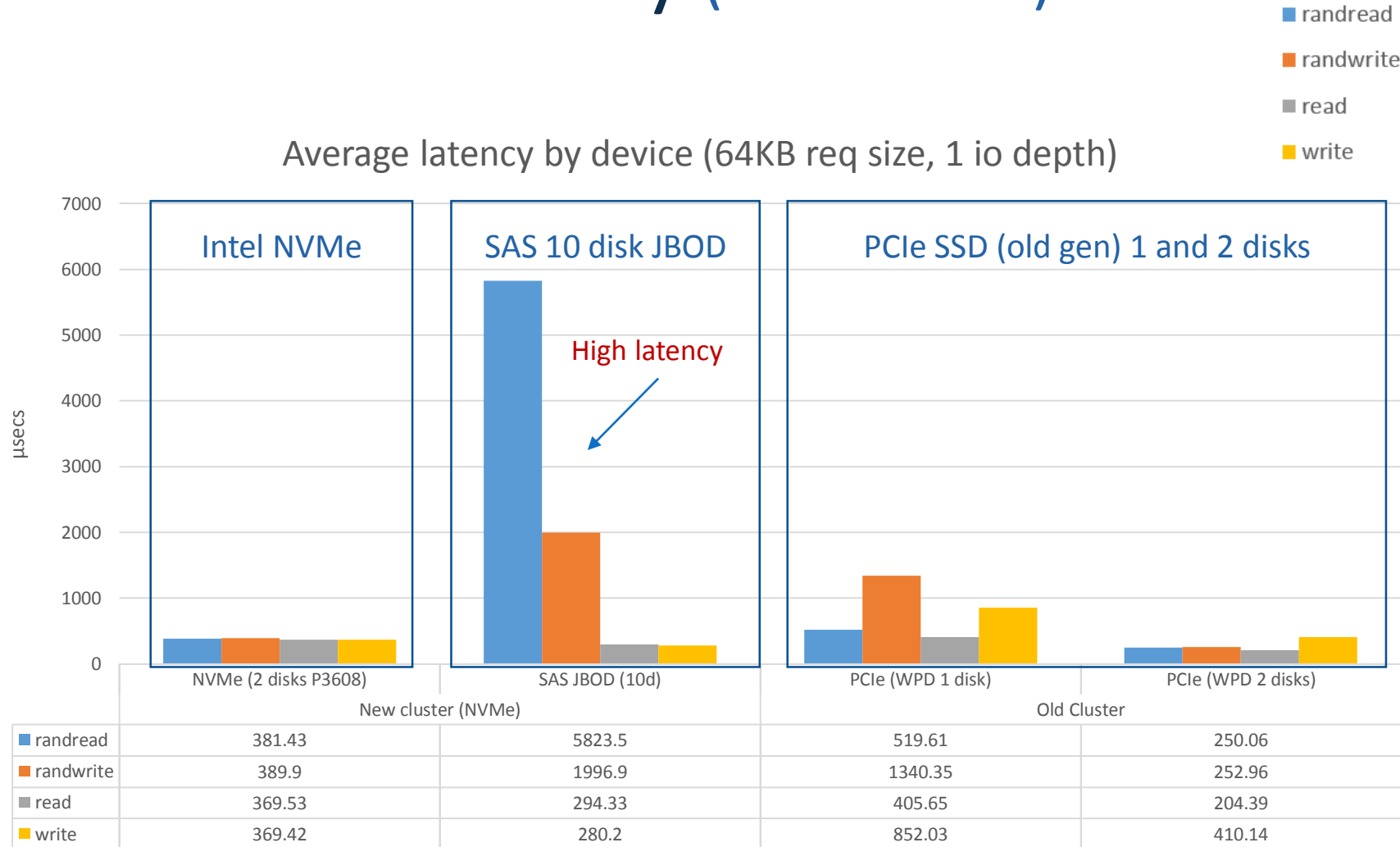
FIO results: Max Bandwidth



Results:

- Random R/W similar in both PCIe SSDs
 - But not for the SAS JBOD
- SAS JBOD achieves high both seq R/W
 - 10 disks 2GB/s
- Achieved both PCIe vendor numbers
 - Also on IOPS
- Combined WPD disks only improve in W performance

FIO results: Latency (smoke test)



Results:

- JBOD has highest latency for random R/W (as expected)
 - But very low for seq
- Combined WPD disks lower the latency
 - Lower than P3608 disks.

Notes:

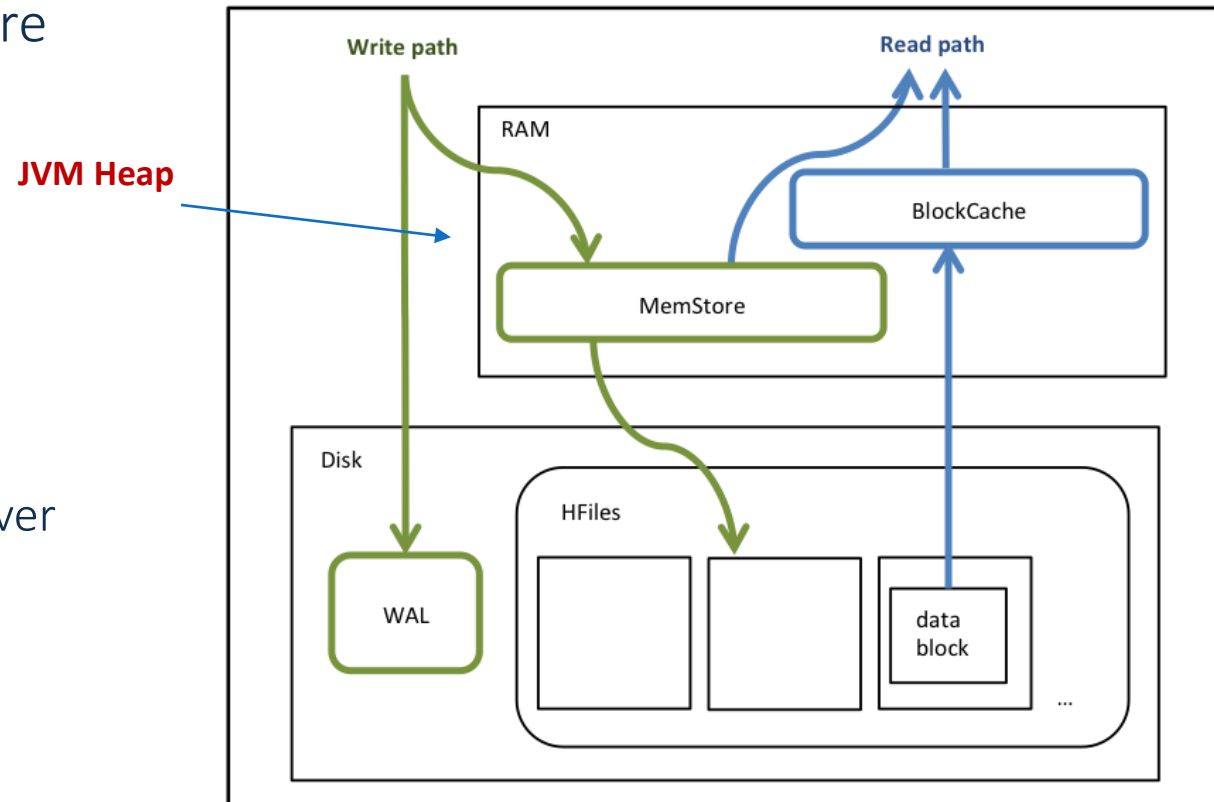
- Need to more thorough comparison and at different settings.



HBase in a nutshell

- Highly scalable Big data key-value store
 - On top of Hadoop (HDFS)
 - Based on Google's Bigtable
- Real-time and **random access**
 - Indexed
 - **Low-latency**
 - Block cache and Bloom Filters
- Linear, modular scalability
 - Automatic sharding of tables and failover
- Strictly consistent reads and writes.
 - failover support
- Production ready and battle tested
 - Building block of other projects

HBase R/W architecture

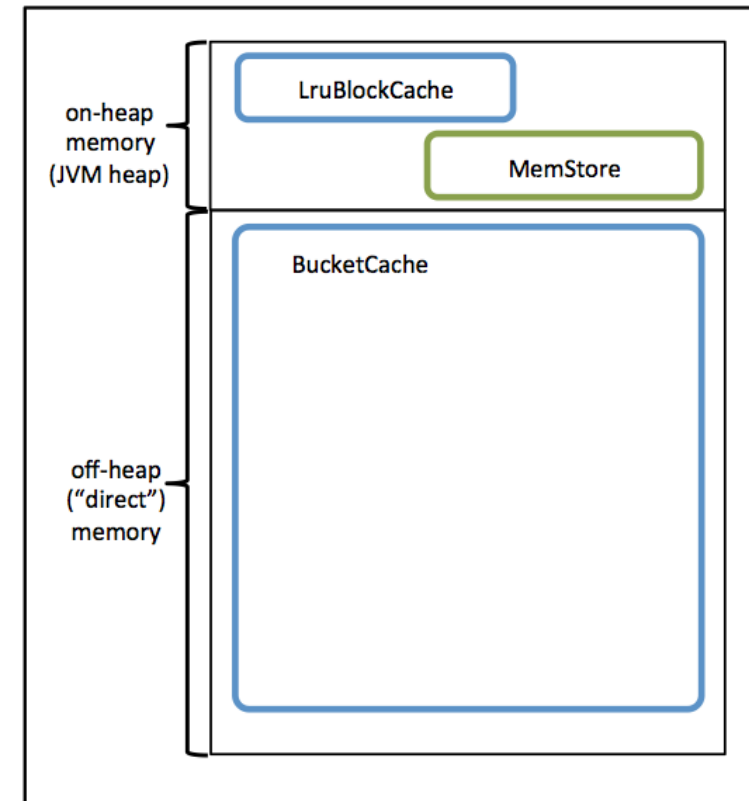


Source: HDP doc

L2 Bucket Cache (BC) in HBase

- Adds a second “block” storage for HFiles
- Use case: **L2 cache** and replaces OS buffer cache
 - Does copy-on-read
- Fixed sized, reserved on startup
- 3 different modes:
 - **Heap**
 - Marginal improvement
 - Divides mem with the block cache
 - **Offheap** (in RAM)
 - Uses Java NIO’s Direct ByteBuffer
 - **File**
 - Any local file / device
 - **Bypasses HDFS**
 - **Saves RAM**

Region server (worker) memory with BC



L2-BucketCache experiments summary

Tested configurations for HBase v1.24

1. HBase default (**baseline**)
2. HBase w/ Bucket cache **Offheap**
 1. Size: **32GB** /work node
3. HBase w/ Bucket cache in **RAM disk**
 1. Size: **32GB** /work node
4. HBase w/ Bucket cache in **NVMe disk**
 1. Size: **250GB** / worker node

- All using same Hadoop and HDFS configuration
 - On JBOD (10 SAS disks, /grid/{0,9})
 - 1 Replica, short-circuit reads

Experiments

1. Read-only (workload C)
 1. RAM at 128GB / node
 2. RAM at 32GB / node
 3. Clearing buffer cache
2. Full YCSB (workloads A-F)
 4. RAM at 128GB / node
 5. RAM at 32GB / node

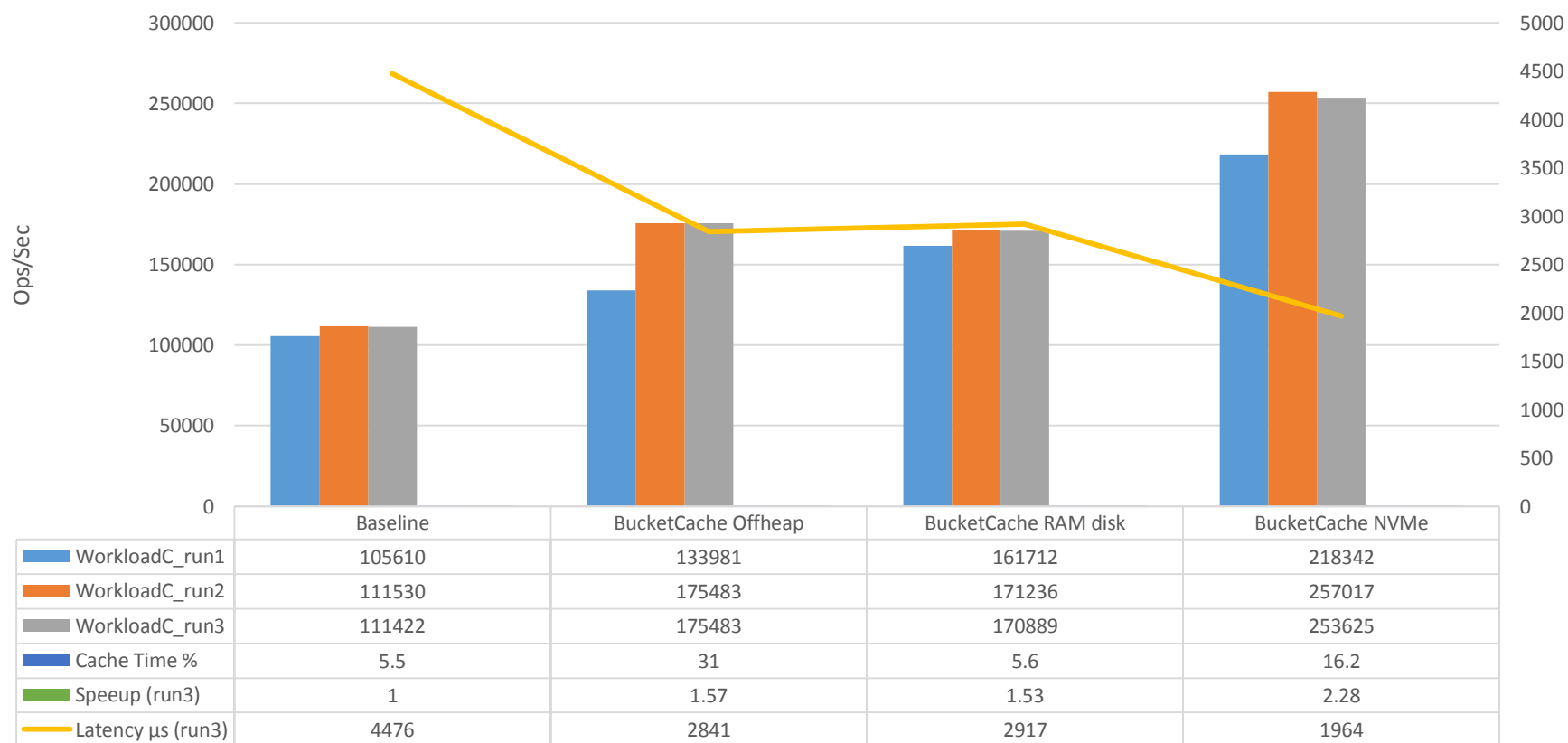
- Payload:
 - YCSB 250M records.
 - ~2TB HDFS raw

Experiment 1: Read-only (gets)

YCSB Workload C: 25M records, 500 threads, ~2TB HDFS

E 1.1: Throughput of the 4 configurations (128GB RAM)

Throughput of 3 consecutive iterations of Workload C (128GB)

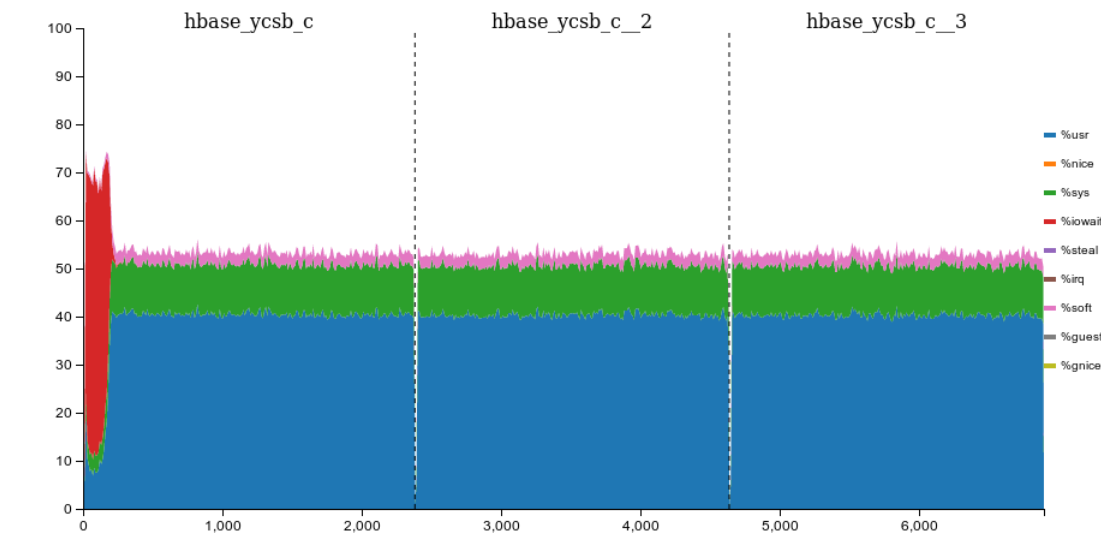


Results:

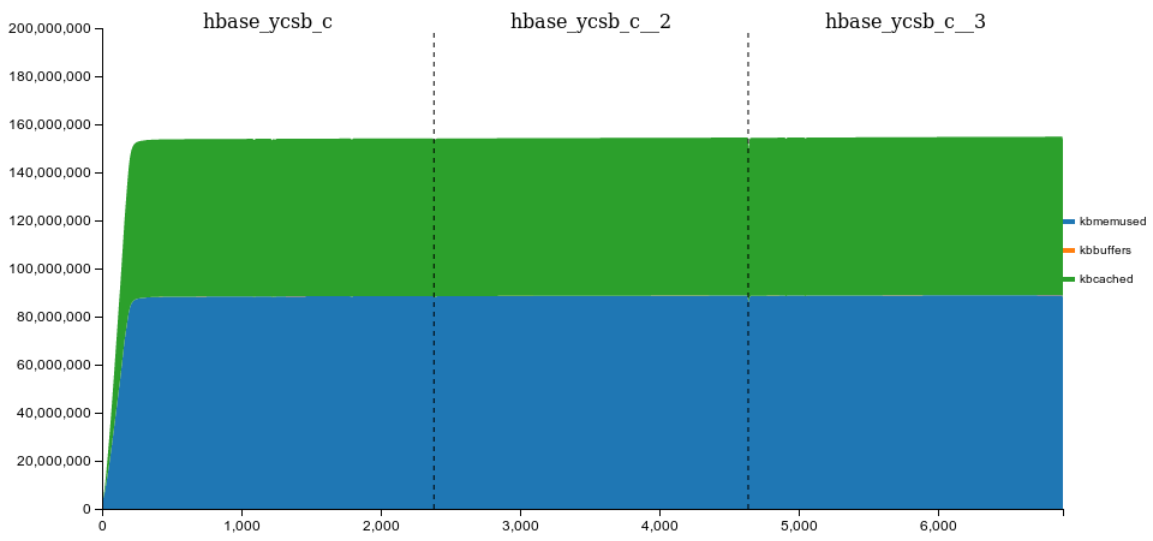
- Ops/Sec improve with BC
 - Offheap 1.6x
 - RAMd 1.5x
 - NVMe 2.3x
- AVG latency as well
 - (req time)
- First run slower on 4 cases (writes OS cache and BC)
 - Baseline and RAMd only 6% faster after
 - NVMe 16%
- 3rd run not faster, cache already loaded
- Tested onheap config:
 - > 8GB failed
 - 8GB slower than baseline

E 1.1 Cluster resource consumption: Baseline

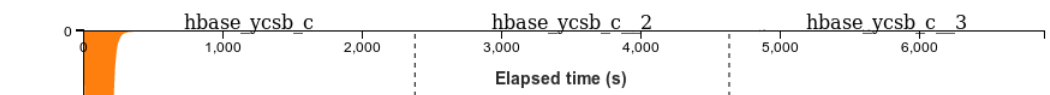
CPU % (AVG)



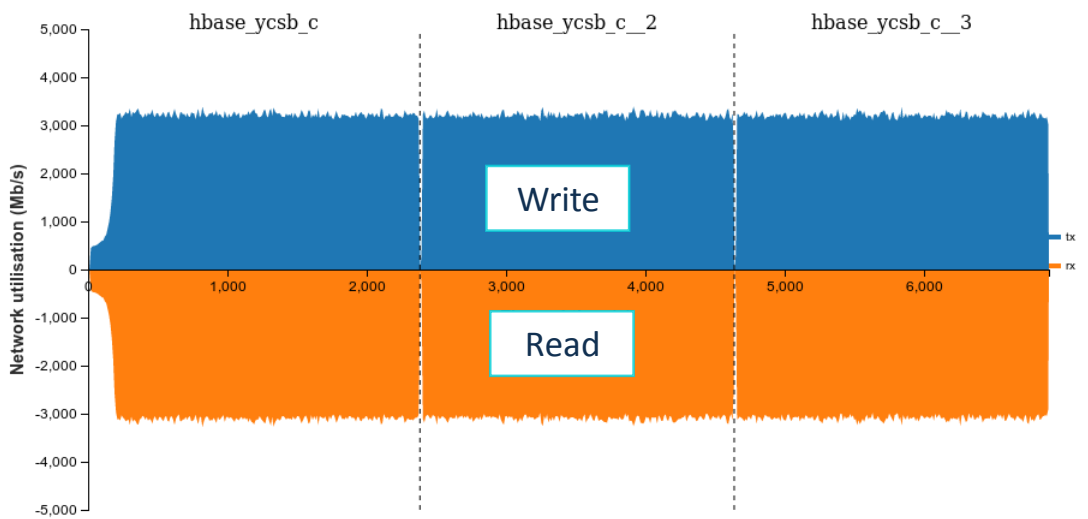
Mem Usage KB (AVG)



Disk R/W MB/s (SUM)



NET R/W Mb/s (SUM)

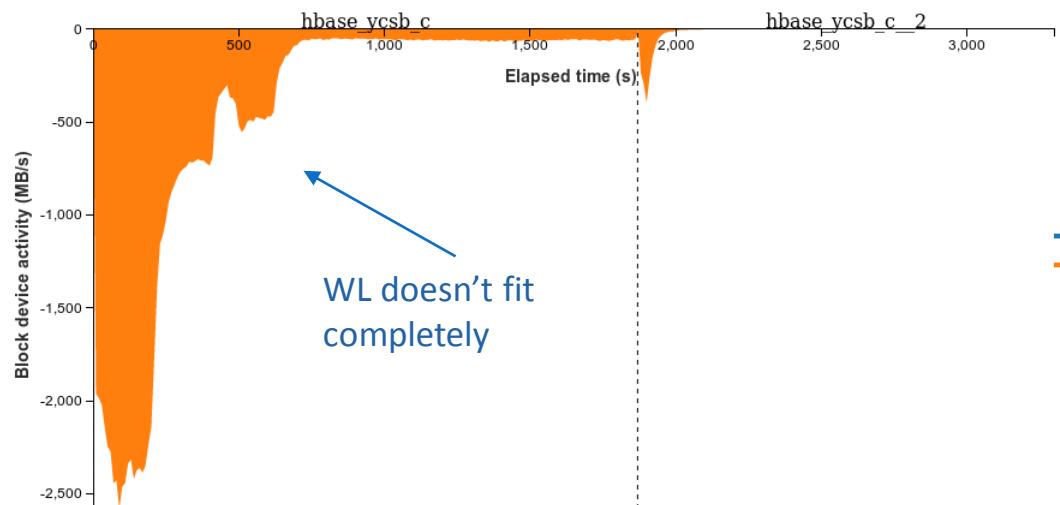


Notes:

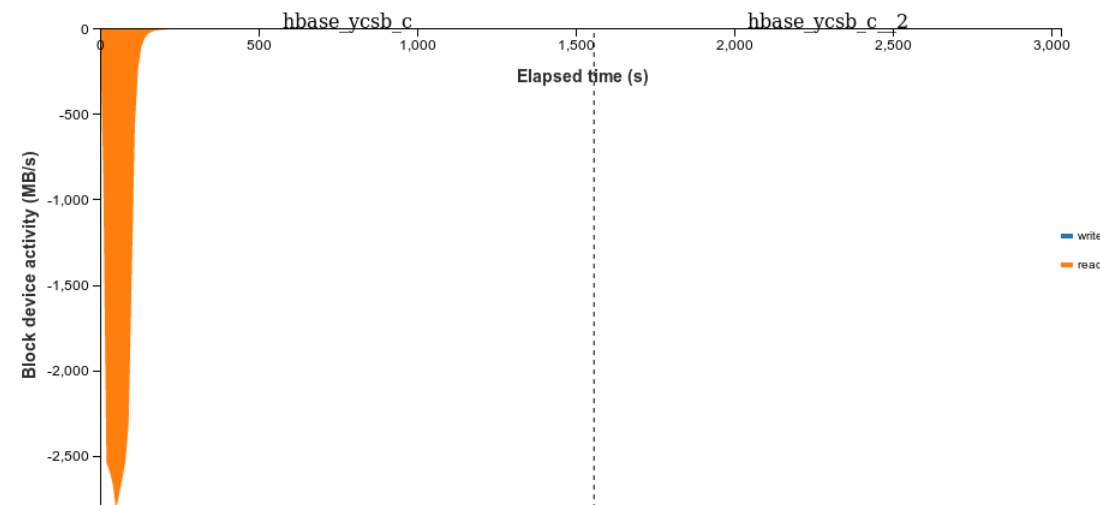
- Java heap and OS buffer cache holds 100% of WL
- Data is read from disks (HDFS) only in first part of the run,
- then throughput stabilizes (see NET and CPU)
- Free resources,
- Bottleneck in application and OS path (not shown)

E 1.1 Cluster resource consumption: Bucket Cache strategies

Offheap (32GB) Disk R/W



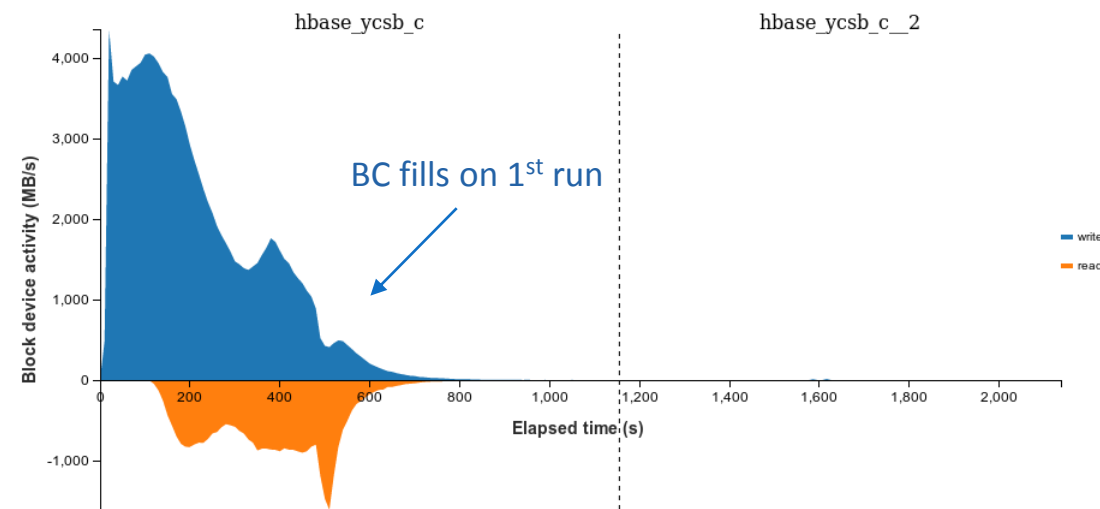
RAM disk (tmpfs 32GB) Disk R/W



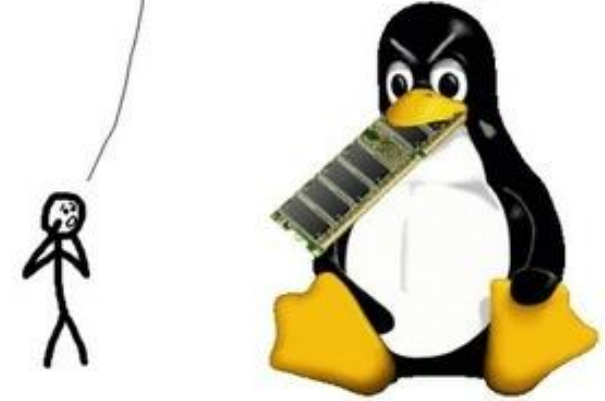
Notes:

- 3 BC strategies faster than baseline
- **BC LRU more effective than OS buffer**
- Offheap slightly more efficient than RAMd (same size)
 - But seems to take longer to fill (different per node)
 - And more capacity for same payload (plus java heap)
- NVM can hold the complete WL in the BC
- Read and Writes to MEM not captured by charts

NVMe (250GB) Disk R/W



Linux ate my ram!



E 1.2-3

Challenge:

Limit OS buffer cache effect on experiments

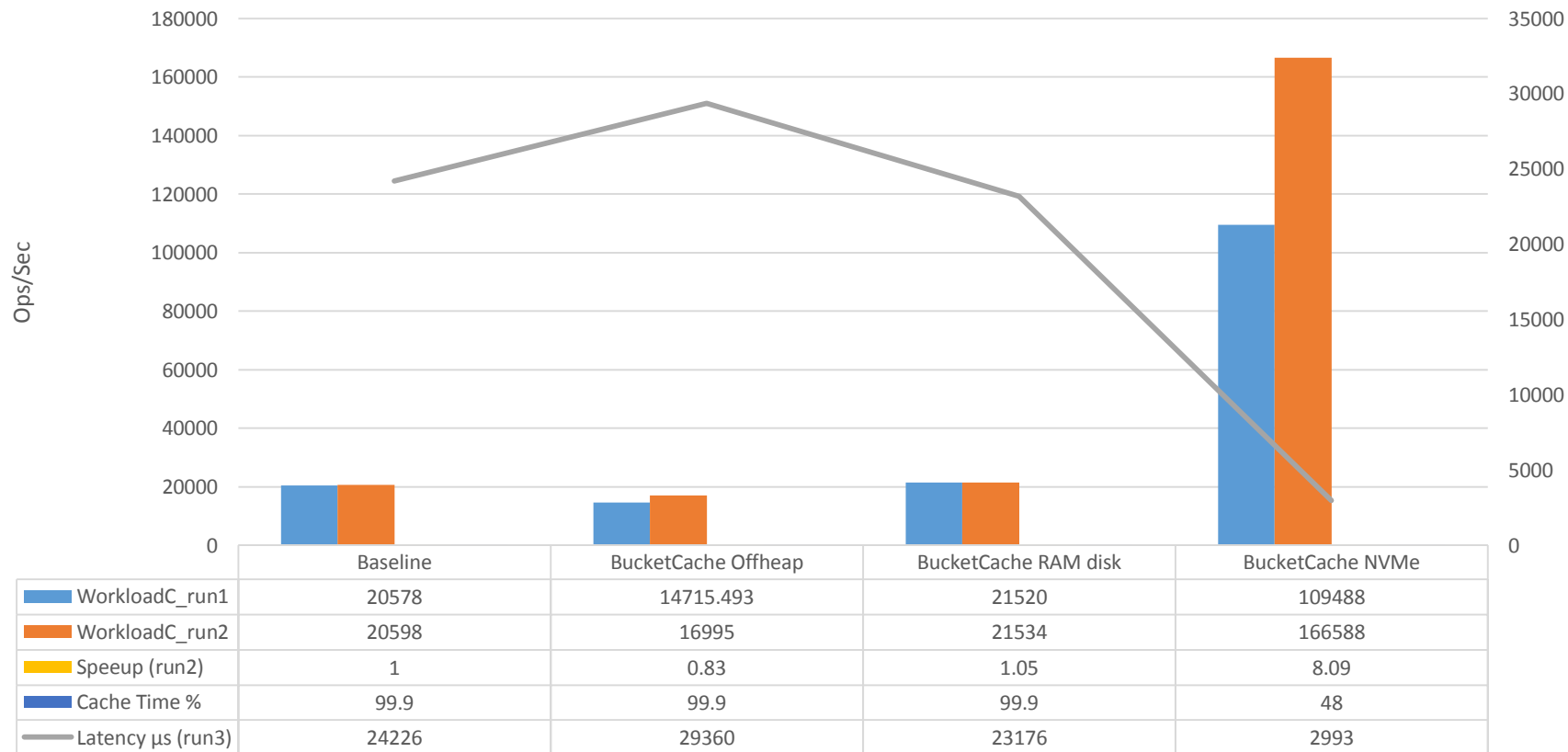
1st approach, larger payload. Cons: high execution time

2nd limit available RAM (using *stress* tool)

3rd clear buffer cache periodically (drop caches)

E1.2: Throughput of the 4 configurations (32GB RAM)

Throughput of 2 consecutive iterations of Workload C (32GB)

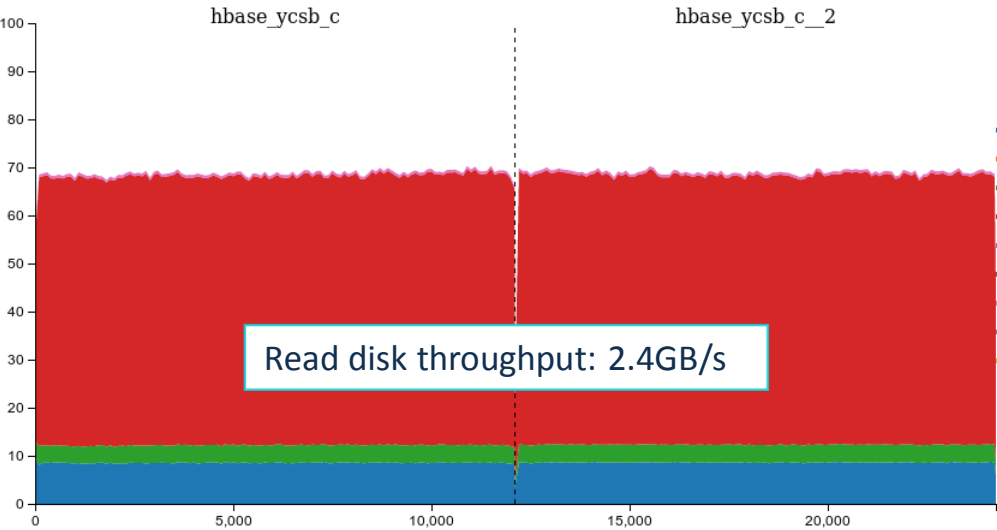


Results:

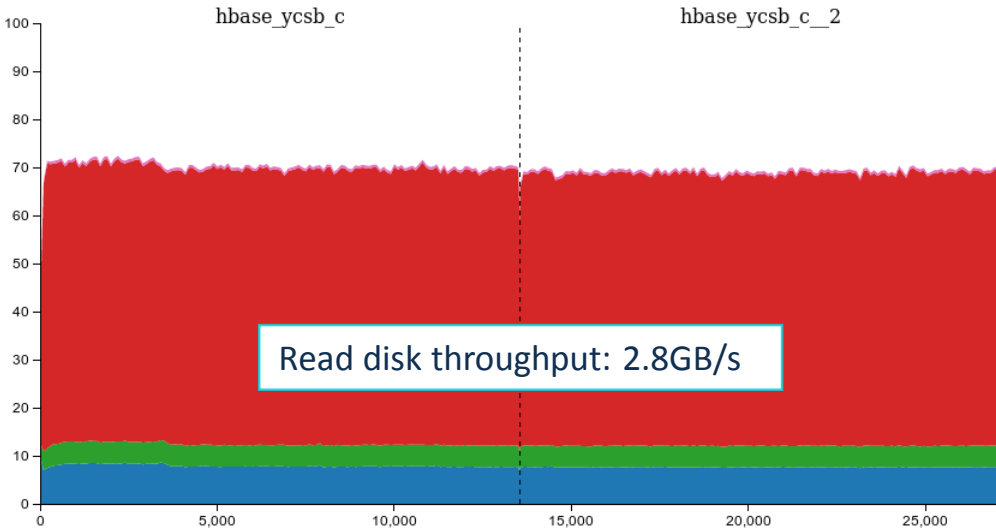
- Ops/Sec improve only with NVMe up to 8X
 - RAMd performs close to baseline
- First run same on baseline and RAMd
 - tmpfs “blocks” as RAM is needed
- At lower capacity, external BC shows more improvement

E 1.1 Cluster CPU% AVG: Bucket Cache (32GB RAM)

Baseline

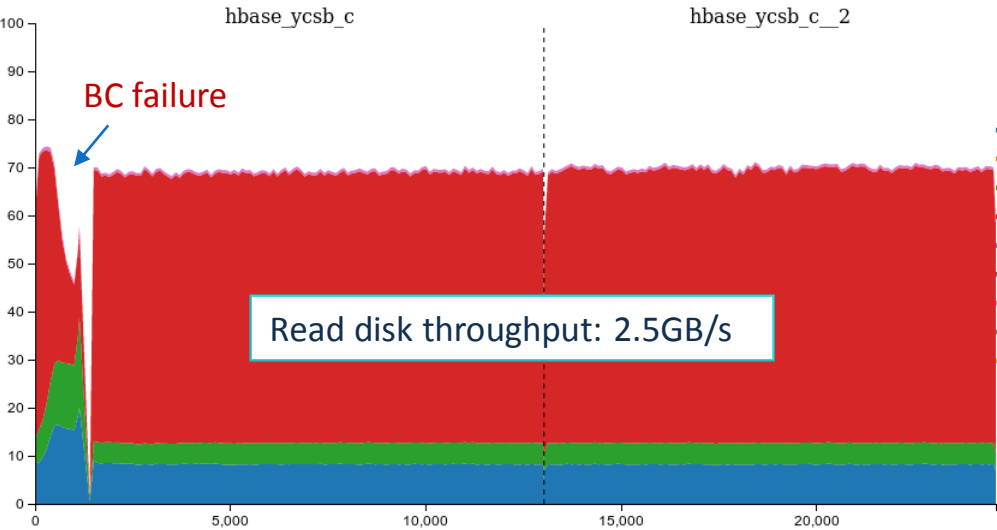


Offheap (4GB)

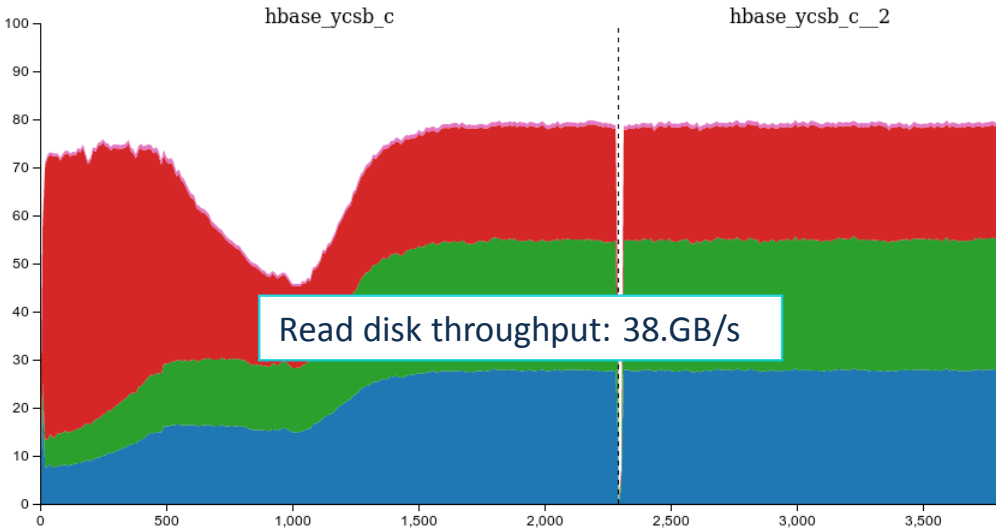


Slowest

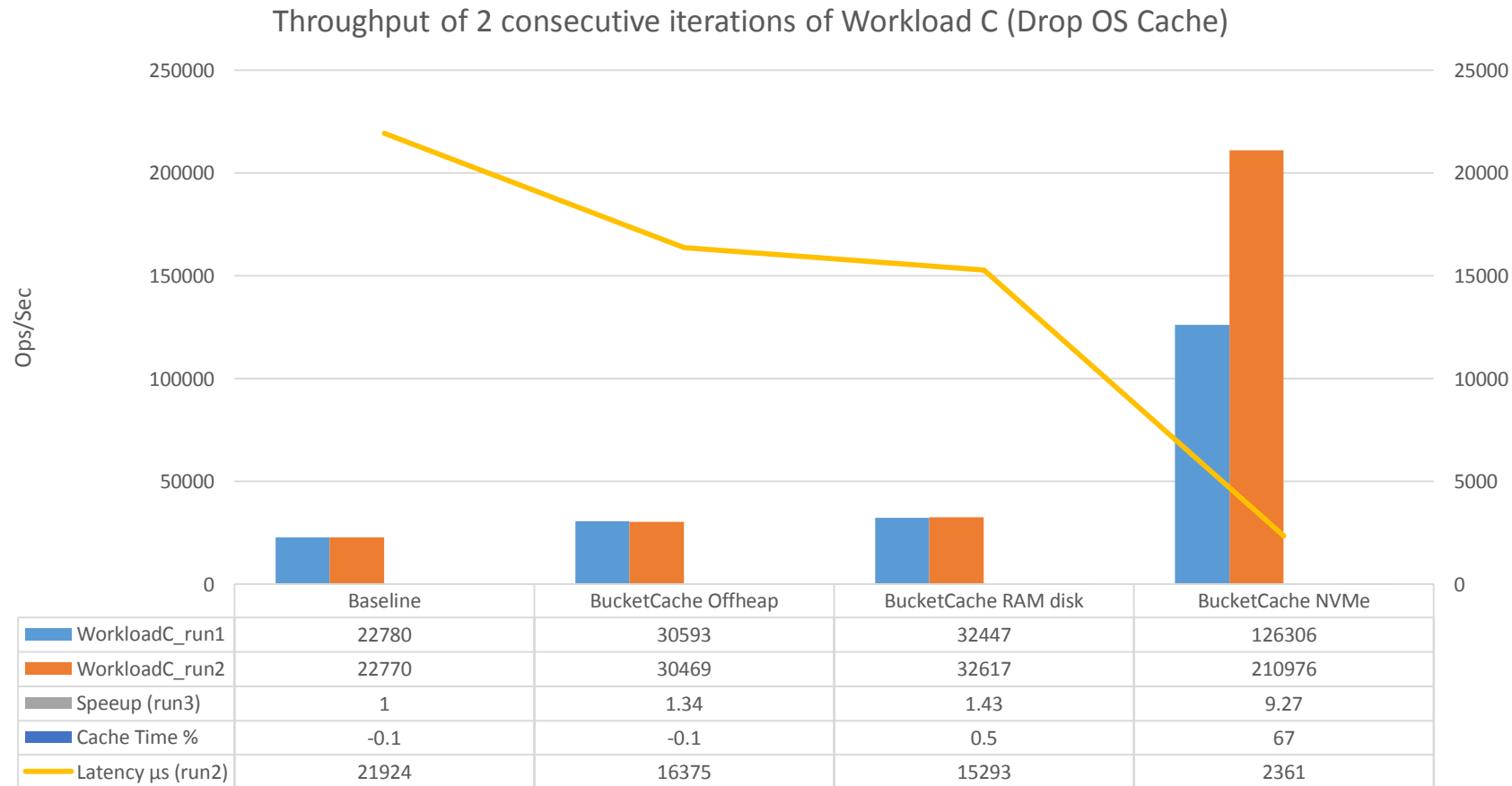
RAM disk (/dev/shm 8GB)



NVMe (250GB)



E1.3: Throughput of the 4 configurations (Drop OS buffer cache)



Results:

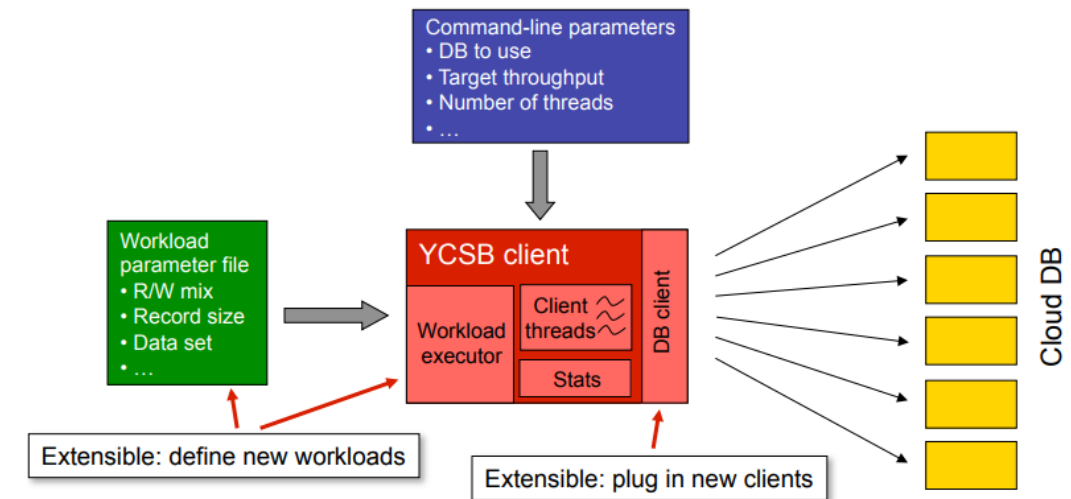
- Ops/Sec improve only with NVMe up to 9X
- RAMd performs 1.43X better this time
- First run same on baseline and RAMd
 - But RAMd worked fine
- Having a larger sized BC improves performance over RAMd

Experiment 2: All workloads (-E)

YCSB workloads A-D, F: 25M records, 1000 threads

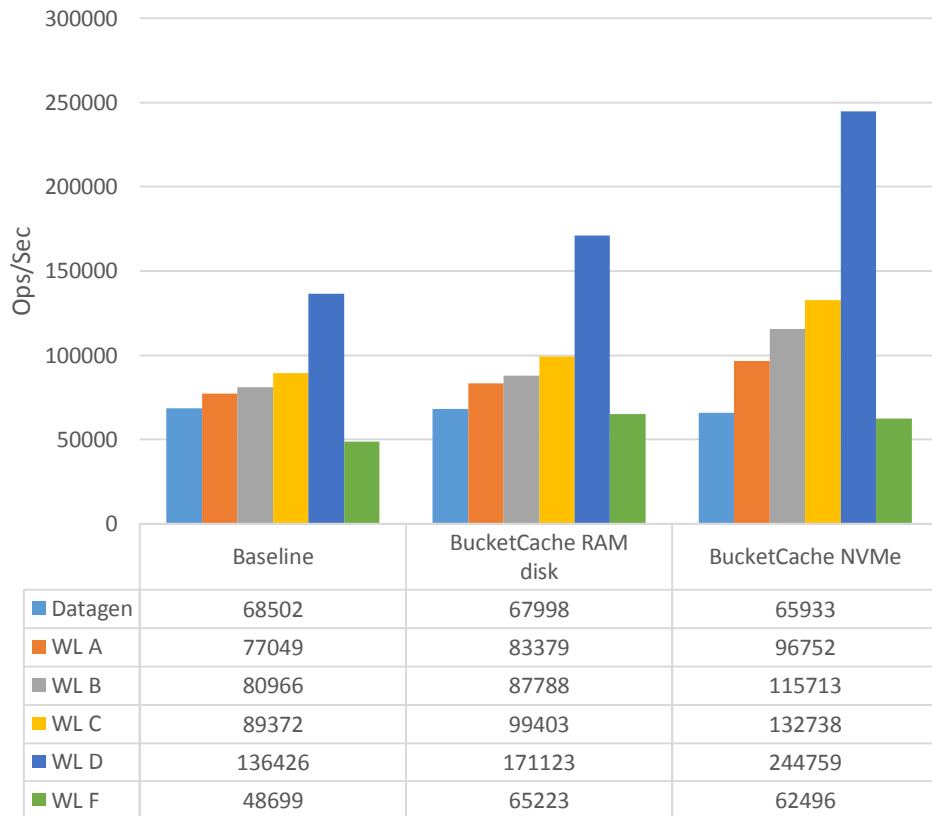
Benchmark suite: The Yahoo! Cloud Serving Benchmark (YCSB)

- Open source specification and kit, for comparing NoSQL DBs. (since 2010)
- Core workloads:
 - A: Update heavy workload
 - 50/50 R/W.
 - B: Read mostly workload
 - 95/5 R/W mix.
 - **C: Read only**
 - **100% read.**
 - D: Read latest workload
 - Inserts new records and reads them
 - **Workload E: Short ranges (Not used, takes too long to run SCAN type)**
 - **Short ranges of records are queried, instead of individual records**
 - F: Read-modify-write
 - read a record, modify it, and write back.

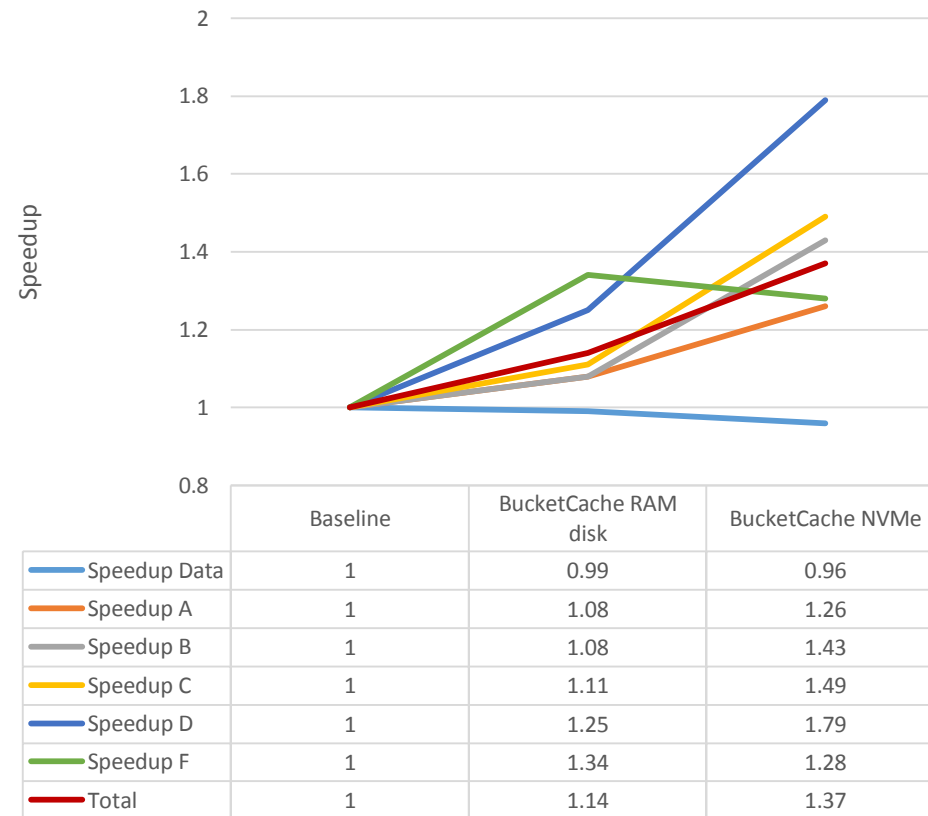


E2.1: Throughput and Speedup ALL (128GB RAM)

Throughput of workloads A-D,F (128GB, 1 iteration)



Speedup of workloads A-D,F (128GB, 1 iteration)

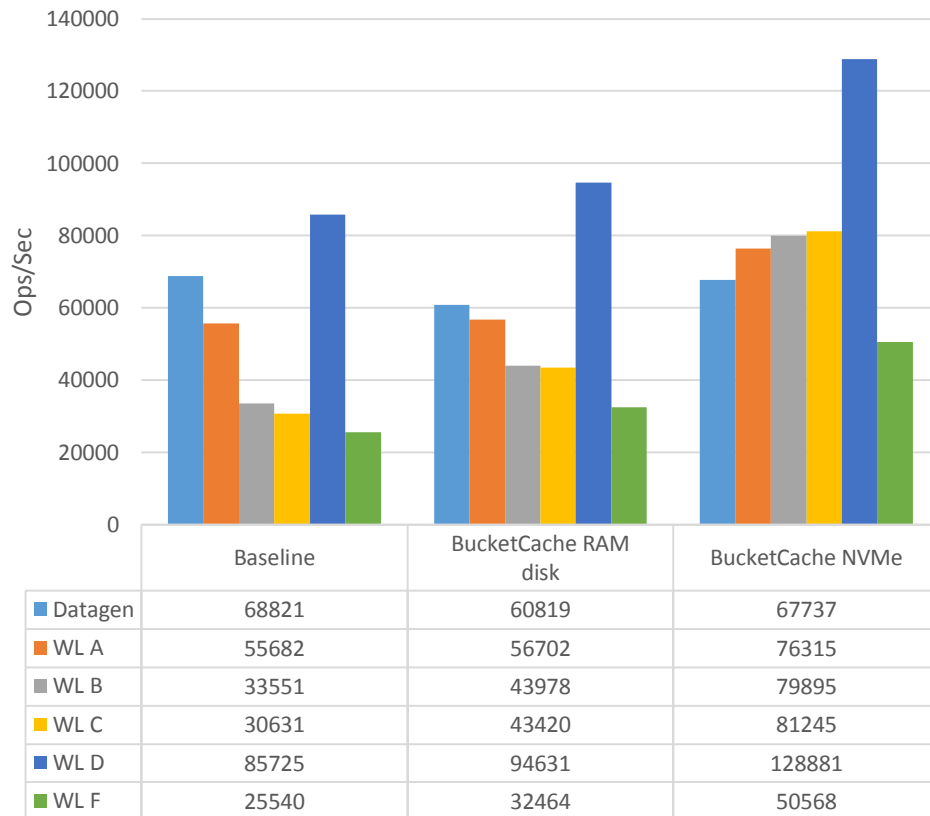


Results:

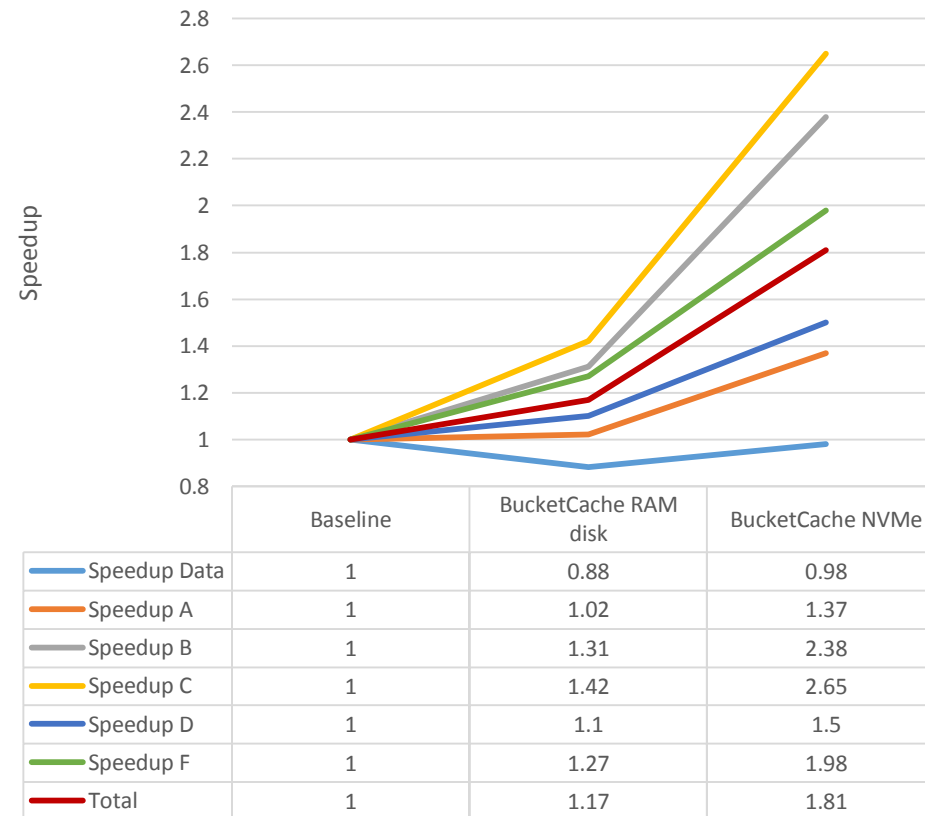
- Datagen same in all
 - (write-only)
- Overall: Ops/Sec improve with BC
 - RAMd 14%
 - NVMe 37%
- WL D gets higher speedup with NVMe
- WL F 6% faster on RAMd than in NVMe
- Need to run more iterations to see max improvement

E2.1: Throughput and Speedup ALL (32GB RAM)

Throughput of workloads A-D,F (128GB, 1 iteration)



Speedup of workloads A-D,F (128GB, 1 iteration)

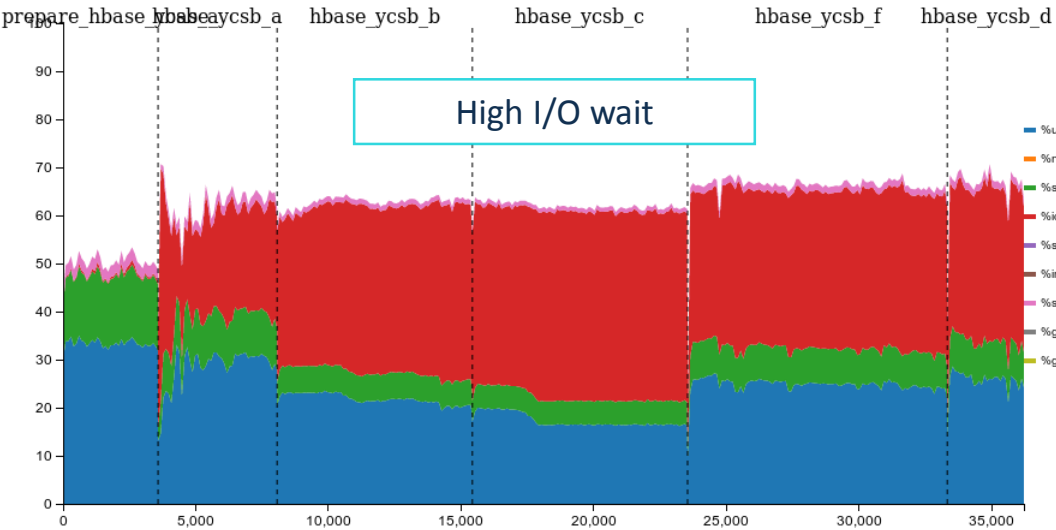


Results:

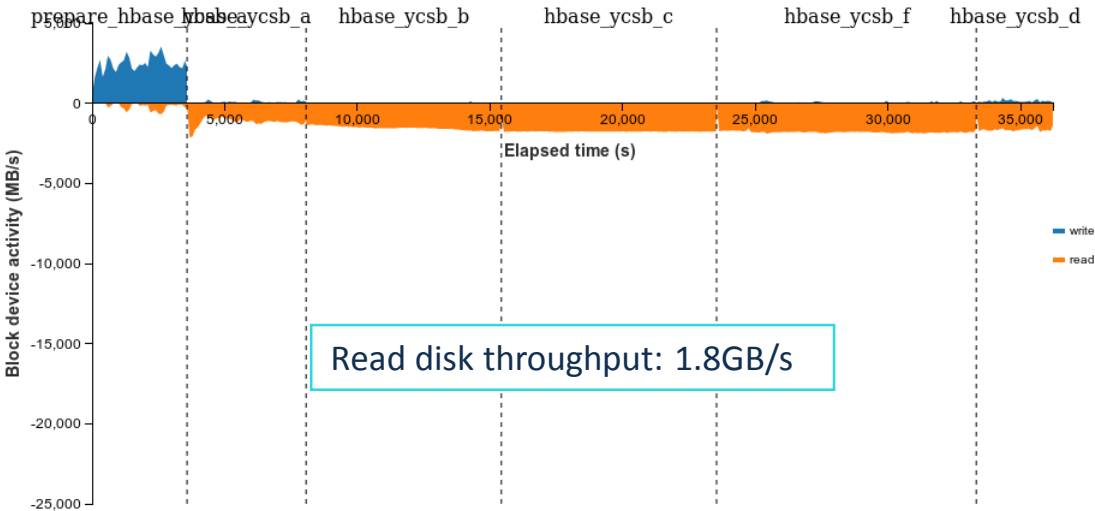
- Datagen slower with the RAMd (less OS RAM)
Overall: Ops/Sec improve with BC
 - RAMd 17%
 - NVMe 87%
- WL C gets higher speedup with NVMe
- WL F now faster with NVMe
- Need to run more iterations to see max improvement

E2.1: CPU and Disk for ALL (32GB RAM)

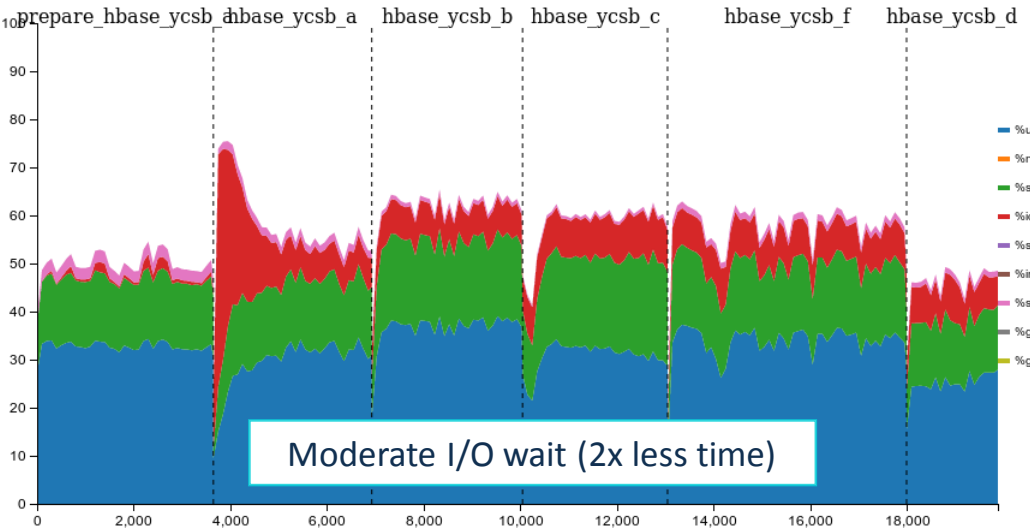
Baseline CPU %



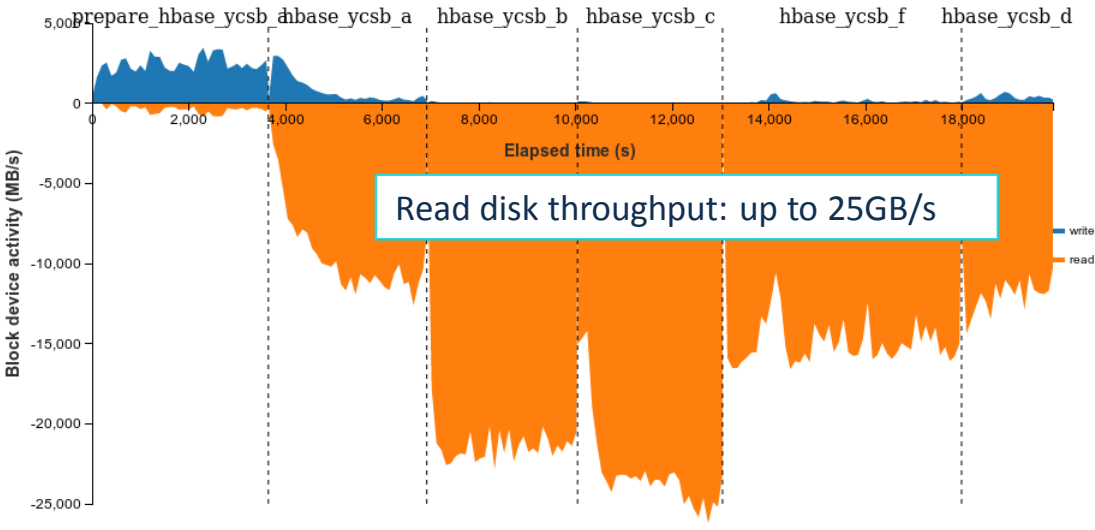
Baseline Disk R/w



NVMe CPU%

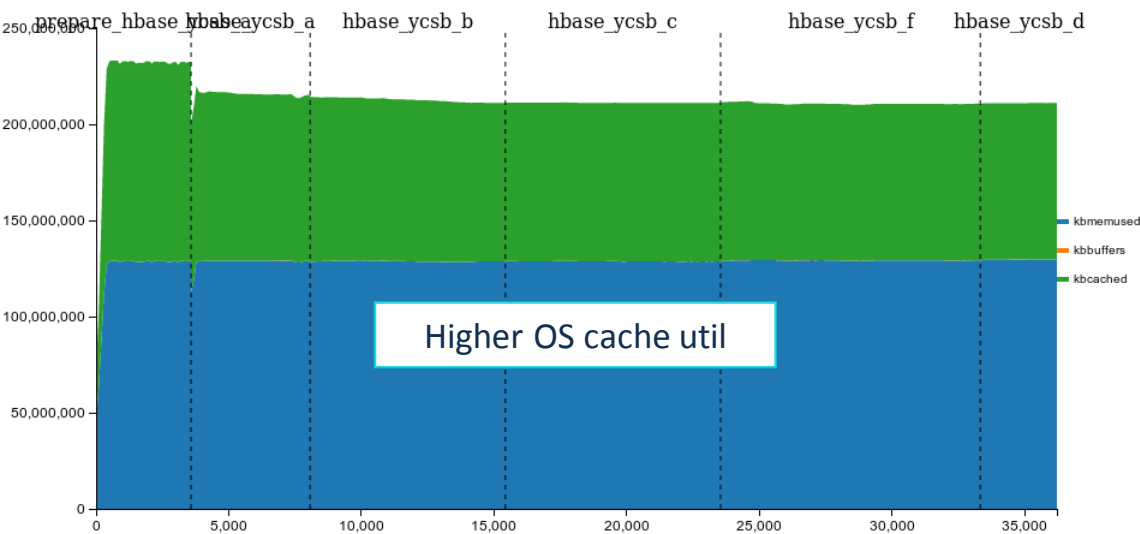


NVMe Disk R/W

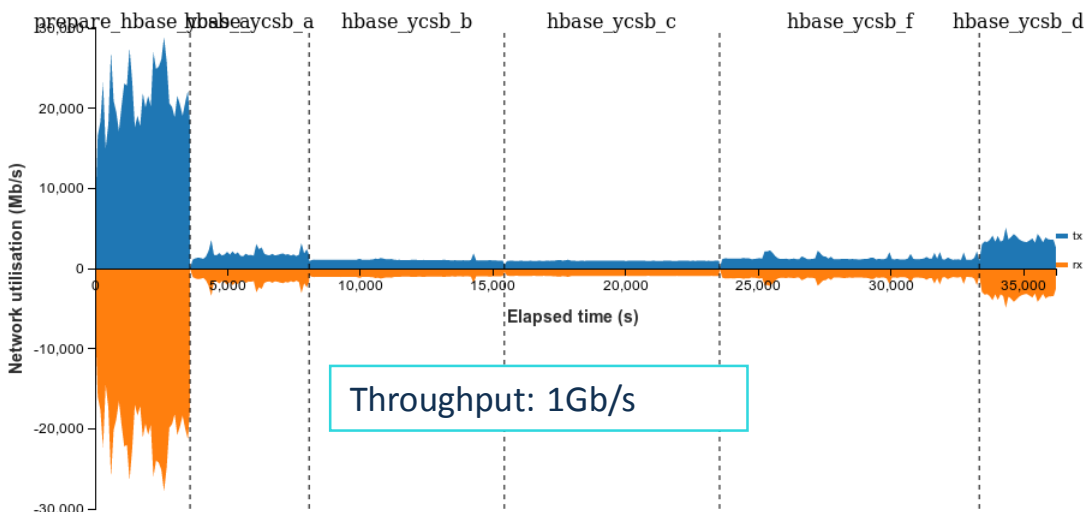


E2.1: NET and MEM ALL (32GB RAM)

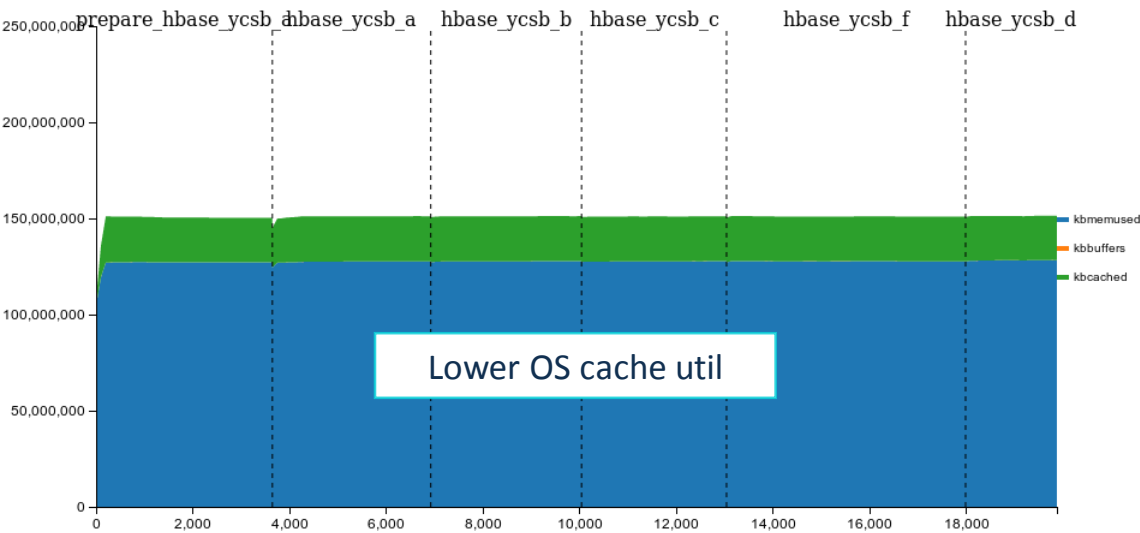
Baseline MEM



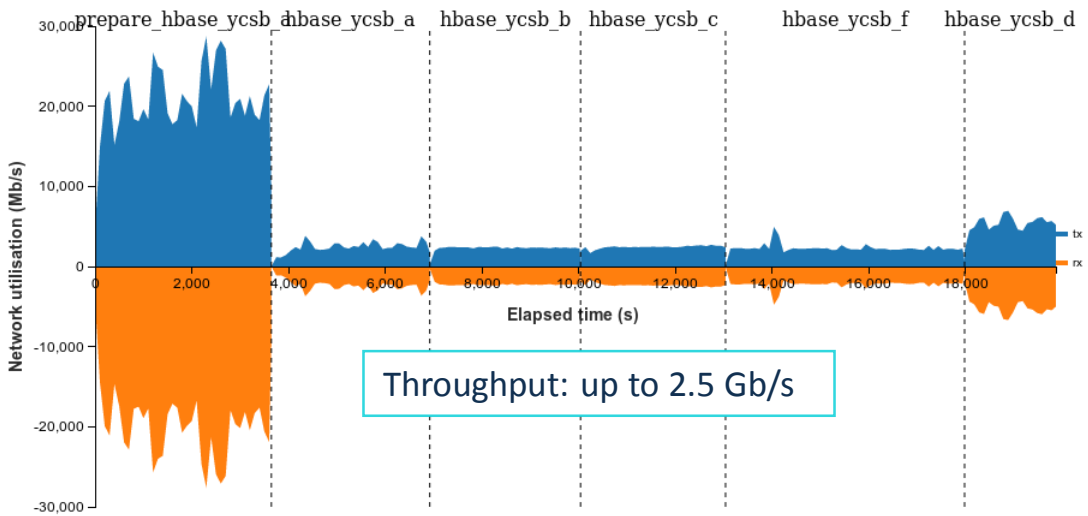
Baseline NET R/W



NVMe MEM

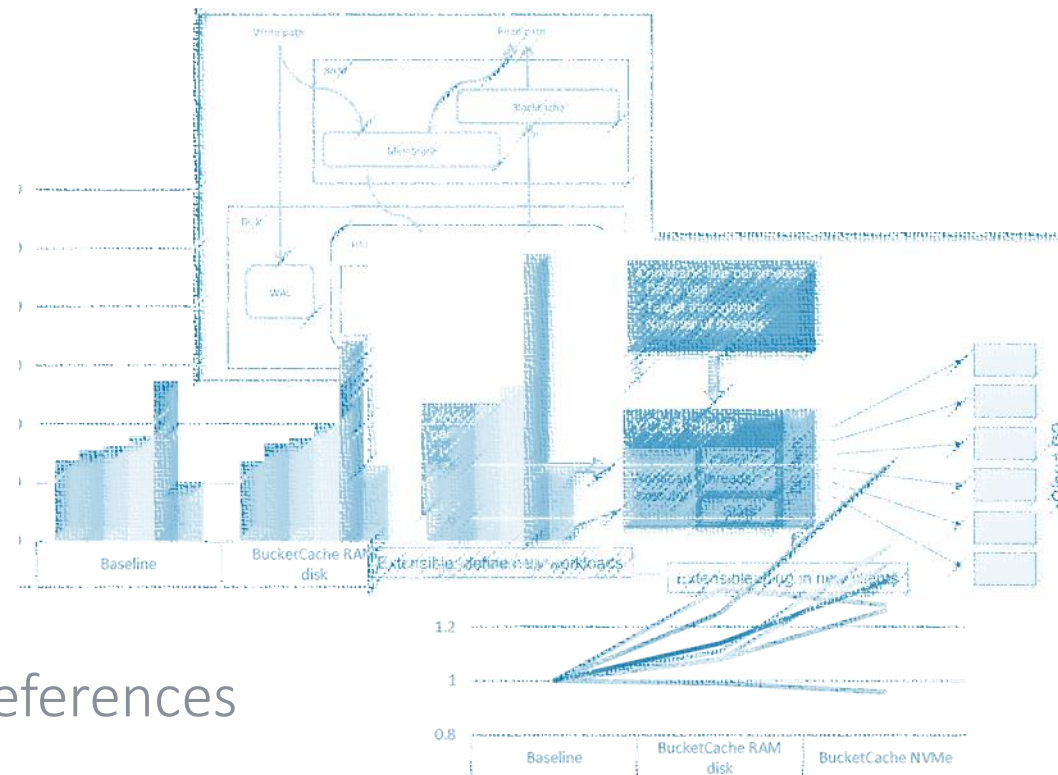


NVMe NET R/W



Summary

Lessons learned, findings, conclusions, references

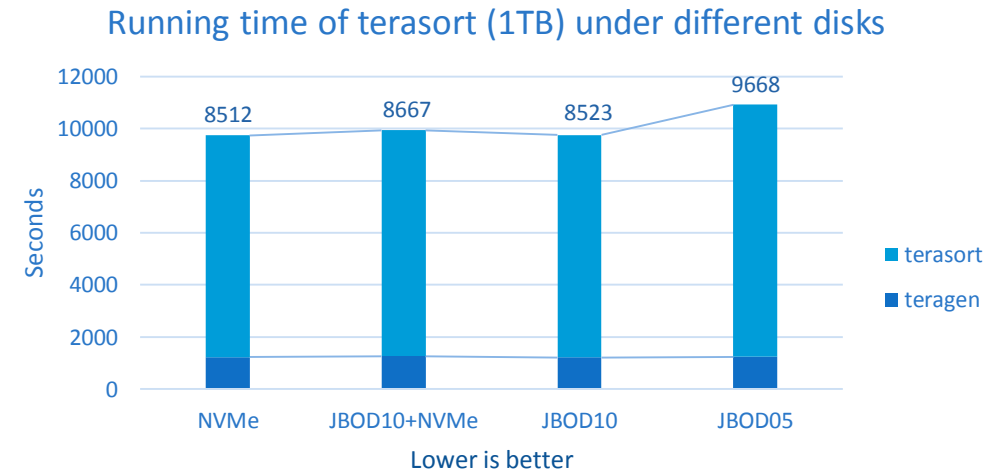


Bucket Cache results recap (medium sized WL) aloja

- Full cluster (128GB RAM / node)
 - WL-C up to **2.7x** speedup (warm cache)
 - Full benchmark (CRUD) from 0.3 to 0.9x speedup (cold cache)
- Limiting resources
 - 32GB RAM / node
 - WL-C NVMe gets up to **8X** improvement (warm cache)
 - Other techniques failed/poor results
 - Full benchmark between **0.4 and 2.7x** speedup (cold cache)
 - Drop OS cache WL-C
 - Up to **9x** with NVMe, **only < 0.5x** with other techniques (warm cache)
- Latency reduces significantly with cached results
- Onheap BC not recommended
 - just give more RAM to BlockCache

Open challenges / Lessons learned

- Generating app level workloads that **stresses newer HW**
 - At acceptable time / cost
- Still need to
 - Run micro-benchmarks
 - Per DEV, node, cluster
 - Large working sets
 - > RAM (128GB / Node)
 - > NMVe (1.6TB / Node)
- OS buffer cache **highly effective**
 - at least with HDFS and HBase
 - Still, a RAM app “L2 cache” is able to speedup
 - App level LRU more effective
 - YCSB **Zipfian distribution** (popular records)
 - The larger the WL, higher the gains
 - **Can be simulated by limiting resources or dropping caches effectively**



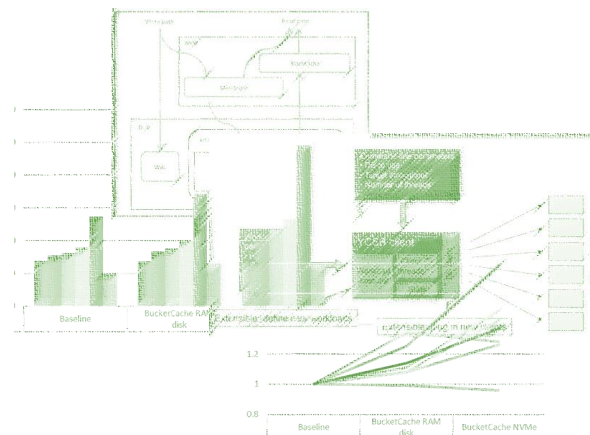
To conclude...

- NVMe offers significant BW and Latency improvement over SAS/SATA, but
 - JBODs still perform well for seq R/W
 - Also cheaper €/TB
 - Big Data apps still designed for rotational (avoid random I/O)
- Full tiered-storage support is missing by Big Data frameworks
 - Byte addressable vs. block access
 - Research shows improvements
 - Need to rely on external tools/file systems
 - Alluxio (Tachyon), Triple-H, New file systems (SSDFS), ...
- Fast devices speedup, but still caching is the simple use case...

References

- ALOJA
 - <http://aloja.bsc.es>
 - <https://github.com/aloja/aloja>
- Bucket cache and HBase
 - BlockCache (and Bucket Cache 101) <http://www.n10k.com/blog/blockcache-101/>
 - Intel brief on bucket cache:
<http://www.intel.com/content/dam/www/public/us/en/documents/solution-briefs/apache-hbase-block-cache-testing-brief.pdf>
 - <http://www.slideshare.net/larsgeorge/hbase-status-report-hadoop-summit-europe-2014>
 - HBase performance: <http://www.slideshare.net/bijugs/h-base-performance>
- Benchmarks
 - FIO <https://github.com/axboe/fio>
 - Brian F. Cooper, et. Al. 2010. Benchmarking cloud serving systems with YCSB.
<http://dx.doi.org/10.1145/1807128.1807152>

Evaluating NVMe drives for accelerating HBase



Thanks, questions?

Follow up / feedback : Nicolas.Poggi@bsc.es

Twitter: @ni_po

FIO commands:

- Sequential read
 - `fio --name=read --directory=${dir} --ioengine=libaio --direct=1 --bs=${bs} --rw=read --iodepth=${iodepth} --numjobs=1 --buffered=0 --size=2gb --runtime=30 --time_based --randrepeat=0 --norandommap --refill_buffers --output-format=json`
- Random read
 - `fio --name=randread --directory=${dir} --ioengine=libaio --direct=1 --bs=${bs} --rw=randread --iodepth=${iodepth} --numjobs=${numjobs} --buffered=0 --size=2gb --runtime=30 --time_based --randrepeat=0 --norandommap --refill_buffers --output-format=json`
- Sequential read on raw devices:
 - `fio --name=read --filename=${dev} --ioengine=libaio --direct=1 --bs=${bs} --rw=read --iodepth=${iodepth} --numjobs=1 --buffered=0 --size=2gb --runtime=30 --time_based --randrepeat=0 --norandommap --refill_buffers --output-format=json`
- Sequential write
 - `fio --name=write --directory=${dir} --ioengine=libaio --direct=1 --bs=${bs} --rw=write --iodepth=${iodepth} --numjobs=1 --buffered=0 --size=2gb --runtime=30 --time_based --randrepeat=0 --norandommap --refill_buffers --output-format=json`
- Random write
 - `fio --name=randwrite --directory=${dir} --ioengine=libaio --direct=1 --bs=${bs} --rw=randwrite --iodepth=${iodepth} --numjobs=${numjobs} --buffered=0 --size=2gb --runtime=30 --time_based --randrepeat=0 --norandommap --refill_buffers --output-format=json`
- Random read on raw devices:
 - `fio --name=randread --filename=${dev} --ioengine=libaio --direct=1 --bs=${bs} --rw=randread --iodepth=${iodepth} --numjobs=${numjobs} --buffered=0 --size=2gb --runtime=30 --time_based --randrepeat=0 --norandommap --refill_buffers --offset_increment=2gb --output-format=json`