

A Media-Aware Cloud Storage Acceleration Layer (CSAL) Cache Solution with Intel® Optane™ SSDs for Alibaba EBS Local Disk D3C Service

Authors

Yanbo Zhou, Alibaba Corporation;
Li Zhang, Alibaba Corporation;
Kapil Karkra, Intel Corporation;
Wayne Gao, Intel Corporation;
Chunhong Mao, Intel Corporation;
Mariusz Barczak, Intel Corporation

Table of Contents

1 Executive Summary	1
2 Background and Motivation	1
3 Architecture	4
4 Evaluations	5
4.1 Tests on FIO Baseline Workloads . . .	5
4.1.1 4K Write Performance for Uniform Random Workloads	5
4.1.2 4K Write Performance for zipf Random Workloads.....	5
4.1.3 64K Write Performance for zipf Random Workloads.....	6
4.1.4 Read/Write Mixed Workloads... .	6
4.1.5 Write Amplification Reduction with CSAL + ZNS.....	6
4.2 Performance for Cloud Production Workloads.....	7
4.2.1 Spark Workloads on Alibaba D3C Public Cloud with Local Disk Services.....	7
4.2.2 TPCx-HS Workloads.....	7
4.2.3 TPC-DS Workloads.....	7
5 Conclusion	7
6 Future of CSAL Software	7
7 Acronyms	8
8 Reference	8

1 Executive Summary

NAND-based solid-state disks (SSDs) are on a trajectory to replace hard disk drives (HDDs) as predominant storage media. Ongoing advances in NAND media are driving higher density and lower costs. However, with density increasing, the endurance and write performance of these media decrease. For example, PLC (5 bits per cell) SSDs, HLC (6 bits per cell) SSDs, or even QLC (4 bits per cell) SSDs are not able to efficiently sustain random workloads or match the line rates of higher-speed networks. So, the storage industry has developed a new Zoned Namespace (ZNS) interface to mitigate endurance and performance challenges (ref[2]).

While ZNS mitigates NAND flash challenges, it requires the host software to be changed because ZNS SSDs only accept sequential writes. For that, we propose a Cloud Storage Acceleration Layer (CSAL) that preserves the existing software interface while transforming write workloads of all forms, random or sequential, large or small, to large sequential writes. This greatly improves the performance and lifetime of SSDs.

The deployment of CSAL in the Alibaba ECS [D3C](#) local-disk product shows that for the Big Data workloads, CSAL with high-density QLC SSDs delivers 2x higher performance than the previous generation of the local-disk product—Alibaba ECS [D2C](#). In addition, storage density in D3C doubles compared to D2C. This whitepaper intends to explain in detail how Alibaba has been able to achieve 2x improvement in both performance and density by using CSAL with high-density flash media.

2 Background and Motivation

Alibaba's D2C instance family is equipped with high-capacity and high-throughput local SATA HDDs. With the exponential growth in data, it needs higher storage capacity and performance. However, there exist two challenges in upgrading this instance family with HDDs: the capacity scaling challenge and the challenge in performance per terabyte (TB). As there are constraints on space and power in the server, it is no longer an option to add more HDDs in the server for larger storage capacity. This is what we call "the capacity scaling challenge". In recent years, the compute (e.g. core counts) and IO interfaces are keeping up with the growth in data, but HDDs' performance per TB drops as indicated by the declining Perf/TB and Perf/vCPU lines as we scale from smaller instance size to larger size with the D2C family offerings (see Figure 1a and 1b). This means HDDs turn out to be a bottleneck at the system level, leading to under-utilization of server resources. We call it "the challenge in performance per TB".

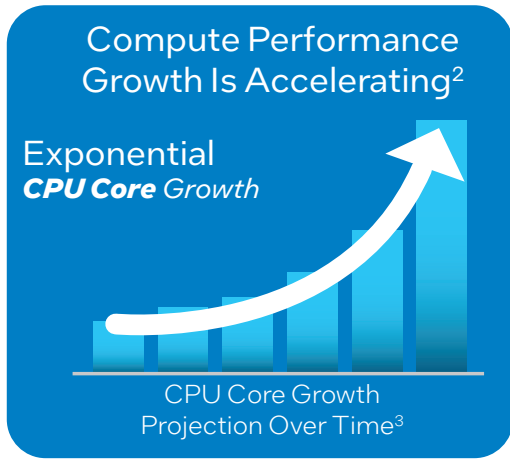


Figure 1a: CPU core growth

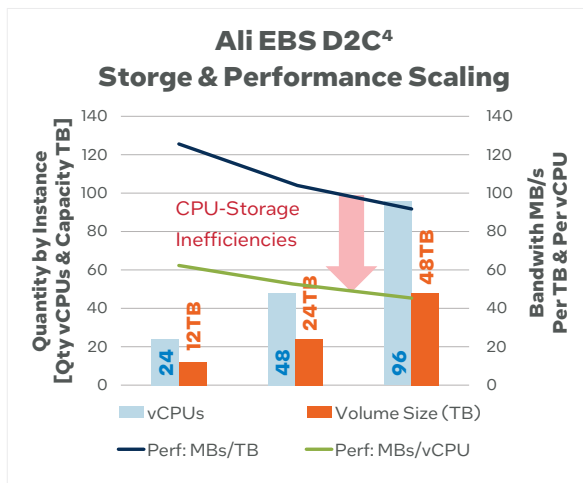


Figure 1b: Performance per TB challenge

To solve these challenges, high-density low-cost QLC NANDs seemed to be a no-brainer solution. However, when we tried it, it didn't work. For some workloads, the QLC media even performed worse than HDDs. Generally, there are four interconnected reasons for this failure of QLC. Firstly, much like HDDs, NAND SSDs' write performance per TB decreases with storage density going up. This is depicted in Figure 2a with the declining line as we transition from low-density SLC and MLC media to higher-density TLC and QLC media. It is expected to get even worse with PLC and HLC NAND media. Secondly, the available Program/Erase cycles decline with higher density NAND. This reduces the endurance/effective lifetime of a QLC SSD significantly as depicted in Figure 2b. Thirdly, high-density NAND SSDs use larger internal storage unit (i.e. indirection

unit, or IU) to reduce the DRAM BOM cost. For example, a standard 4TB 4K IU SSD needs 4GB of DRAM to hold the direct logical to physical (L2P) lookup table. Thus, for a 64TB SSD, 64K IU has to be used if no more on-SSD DRAM is added. However, when we use 64K IU sized SSDs (e.g. Solidigm's D5-P5316 SSDs), executing 4K random writes results in 16x write amplification due to read-modify-write penalty. As depicted in Figure 2c, any mis-sized or misaligned writes may incur significant write amplification. Fourthly, and also the primary reason, is the fragmentation due to multi-tenancy (i.e. multiple tenants sharing a single large SSD) as depicted in Figure 2d. For example, if 8 tenants share a single D5-P5316 16TB QLC SSD, each tenant will have a 2TB virtual disk technically. However, when these 8 tenants simultaneously write to the SSD with different velocities and data lifetime, the SSD, which delivers 3000MiB/s sequential write bandwidth, is fragmented so much that defragmentation makes the write bandwidth drop down to 400MiB/s. Each tenant, therefore, gets 50MiB/s of bandwidth. This is 20% of the bandwidth for a tenant in the D2C case where each tenant is allocated one physical 2TB HDD which delivers 250MiB/s sequential write performance.

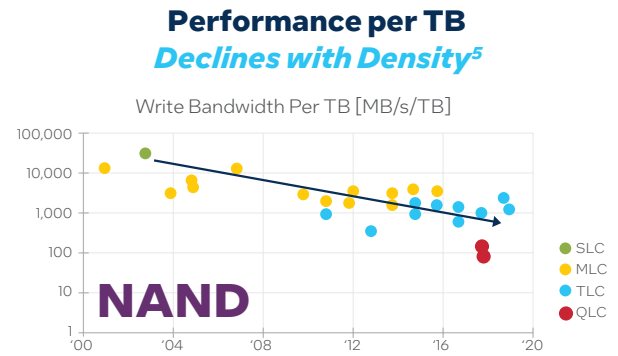


Figure 2a: Performance per TB

Endurance Declines with Density

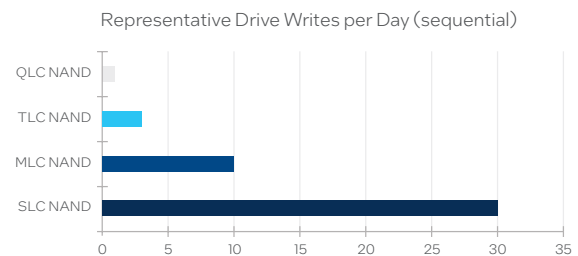


Figure 2b: Endurance of media

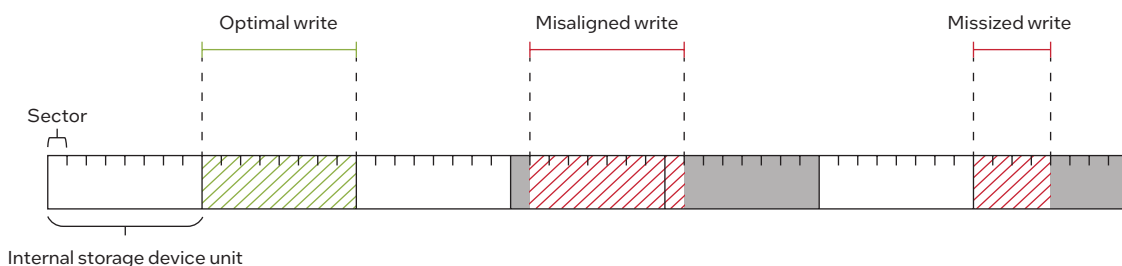


Figure 2c: Small random writes

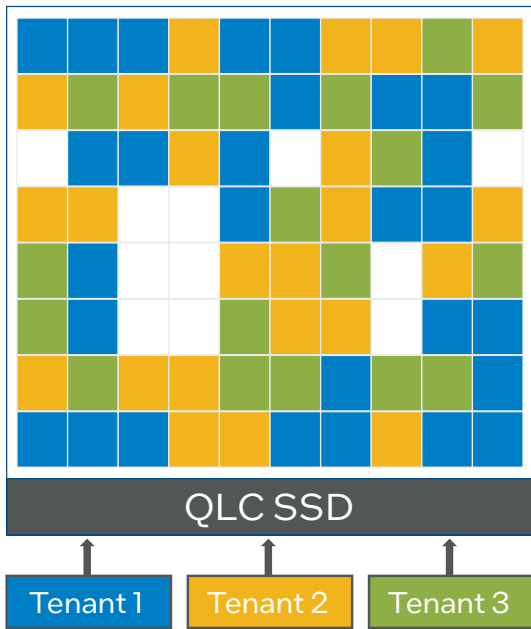


Figure 2d: Multiple tenancy

Fortunately, emerging Storage Class Memory (SCM) SSDs, such as Intel® Optane™ SSDs P5800X/P5810X/P5820X and SLC SSDs, have introduced a new tier into the memory and storage hierarchy to mitigate the performance per TB and endurance issues with high-density QLC NAND media, as shown in Figure 3 (ref[5]). This is a promising technology that uses caching to simultaneously reap the benefits of both performance (SCM) and capacity (QLC) tiers. In a traditional cache architecture, high performance storage, such as an SCM SSD, is put in front of primary storage, such as a QLC SSD. Instead of writing data to primary storage directly, writes are acknowledged to users or applications as soon as data is written to the cache tier. Then data is written back to the capacity tier.

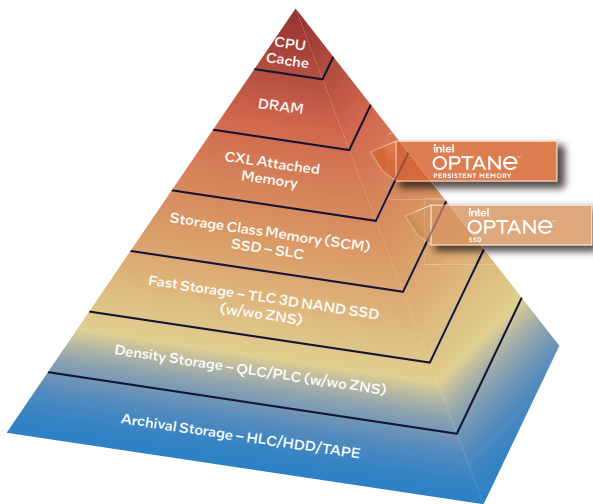


Figure 3: Optane SSD/PMem and CXL-attached memory/SCM SSD in a tiering hierarchy

Traditional caches can help high-density NAND media to mitigate performance per TB and endurance issues for high-temporal-locality workloads. For example, a high-performance, high-endurance SCM tier can absorb frequently updated writes without sending them to the QLC NAND tier. However, the design of a traditional set-associative cache is not suitable for SSDs with large IU size. For write-heavy workloads, when data is written from a set-associative cache to the QLC NAND tier, a mis-sized/misaligned random write workload will be generated, incurring significant read-modify-write penalty. Furthermore, in traditional caches, a sequential workload can be turned into a random workload over time as the cache free space becomes fragmented, and the traditional cache may spread the contiguous Logical Block Addresses (LBAs) to available cache lines that may not be contiguous. The cleaning policies may allow best-effort sequential writes to the QLC tier but cannot guarantee strictly sequential writes to QLC as required by the ZNS interface.

To address these problems, we propose CSAL, the Cloud Storage Acceleration Layer, an open-source cloud-scale share-nothing storage software layer (bdev, i.e. block device) in Storage Performance Development Kit (SPDK). The key strategy of CSAL is to leverage SCM as write cache to shape the write workloads in all forms and sizes to NAND friendly big IO size sequential writes. The criteria of IO size for sequential writes chosen here is that it must be no less than the IU granularity of a QLC drive. With that, write amplification factor (WAF) of a drive is reduced significantly, for example, 4KB random write WAF drops from over 70 (ref[3], Table 2) to 1.02. In addition, CSAL also outperforms traditional cache technologies in three ways:

1. CSAL uses an ultra-fast write buffer (SCM) to “sequentialize” I/O writes to the QLC device for higher performance and endurance at the system level.
2. CSAL absorbs and compacts large quantities of user writes in the cache tier, further improving endurance and lifespan of the capacity tier—QLC NAND SSDs.
3. CSAL guarantees that data in the cache tier can be written back to the capacity tier in predictable time.

Thus, CSAL can utilize Zoned Namespace (ZNS) to remove extra costs of capacity over-provision and DRAM for internal Flash Translation Layer (FTL) within generic SSDs, which decreases total costs of SSDs as primary storage as well.

We implement CSAL in SPDK for high performance storage systems. SPDK offers full stack storage system from a logic volume, a generic block layer to a NVMe driver. CSAL is implemented in the SPDK block layer and exposed as a virtual block device that consists of two physical block devices: SCM SSD as the cache tier and NAND SSD as the capacity tier. Storage applications, such as vhost-blk, iSCSI and NVMeoF, can use this virtual block device as a generic block device.

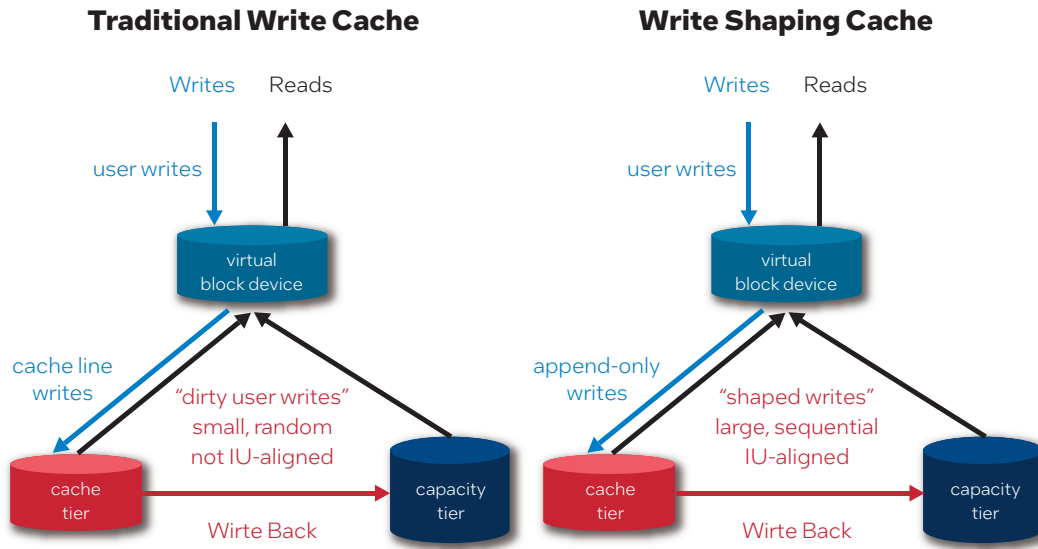


Figure 4: Differences between a traditional write cache and a CSAL design

Figure 4 illustrates the main differences between a traditional write cache and a write shaping cache proposed in this paper. The traditional write cache on the left is a **set-associative cache**. The Write Shaping Cache we propose is a **log-structured cache**.

The key strategy of CSAL is to leverage an SCM SSD as the cache to compact and shape user random writes to NAND friendly writes. The goal of a CSAL design is to minimize the system-level write amplification and the wear for NAND SSDs, hence improving overall performance and system endurance of NAND-based primary storage.

3 Architecture

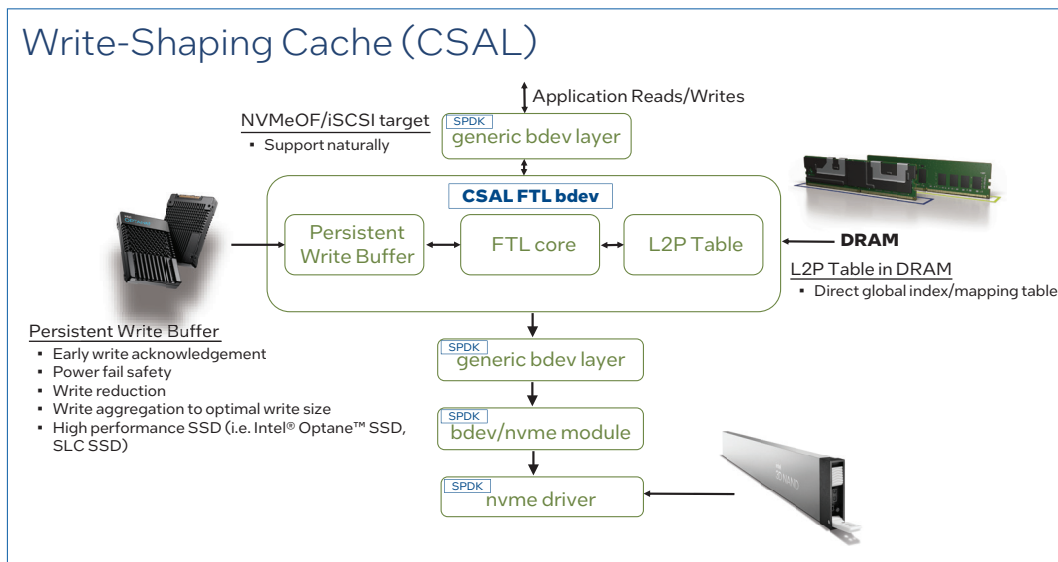


Figure 5: Diagram of a write-shaping cache block

Figure 5 shows the overall architecture of CSAL.

1. CSAL is a generic SPDK block device (bdev) that supports NVMeOF and iSCSI targets naturally.
2. Application Reads/Writes go through the SPDK generic bdev layer first, and then go into CSAL bdev.
3. The CSAL bdev layer is a virtualized FTL device that will shape a random workload into a sequential workload by leveraging the Persistent Write Buffer and the L2P table.

- a. FTL will record user write IOs to Persistent Write Buffer as FIFO logs on the Intel Optane SSD, and the L2P table is then updated to point to the Optane LBA.
- b. When the cache free capacity reaches a certain threshold, FTL background compaction process will kick in to:
 - i. Read FIFO logs from the Optane SSD
 - ii. Evict invalid logs
 - iii. Merge and write valid logs as large sequential IOs to the QLC SSD
 - iv. Update the L2P table to point to the QLC LBA.

4. Data is written to the QLC and Optane SSDs via the standard SPDK bdev again.
5. A FTL device is similar to an SSD device, and defragmentation is designed to do housekeeping jobs to maintain the free space for new writes.

To achieve the above data transition, CSAL manages four key components: logical to physical (L2P) address table, Persistent Write Buffer, compaction worker, and garbage collection (GC) worker.

4 Evaluations

4.1 Tests on FIO Baseline Workloads

To evaluate the performance of CSAL, we set up the environment with the following configuration for comparison:

Storage Server—Supermicro SYS-220U-TNR System Configuration	
BIOS Version	1.1
OS	Fedora 33 (Server Edition)
Kernel	5.12.15-200.fc33.x86_64
CPU Model	Intel® Xeon® Platinum 8375C CPU @ 2.90GHz
NUMA Node(s)	2
DRAM Installed	256GB (16x16GB DDR4 3200MT/s)
Huge Pages Size	2048 kB
NIC Summary	Ethernet Controller X710 for 10GBASE-T
Drive Summary	<ol style="list-style-type: none"> 1. O is Intel® Optane™ SSD P5800X 800GB SSDPF21Q016TB for a cache tier 2. Q is a Solidigm QLC SSD P5316 16TB SSDPF2NV153TZ for a capacity tier 3. TLC is a Solidigm TLC SSD P5510 8TB SSDPF2KX076TZO for a capacity tier 4. O+Q: O: 1x P5800 SSDPF21Q016TB; Q: 1x P5316 SSDPF2NV153TZ 5. ZNS SSD: WD ZN540 4TB 6. Regular SSD: (used for ZNS comparison)—SN640 7.68TB
SPDK	21.04
CSAL	1.0
FIO	3.20

Table 1: FIO server configuration for the tests in Section 4.1

Before the evaluation, QLC and TLC SSDs have been pre-conditioned to the steady state. FIO is used in all the tests for bandwidth and WAF on SPDK block layer. To be fair, O+Q, TLC and QLC SSDs are all 10% over provisioned. Our evaluation of CSAL is divided into three aspects:

- 8x jobs’ 4K Write performance for purely uniform random and zipf1.2 workloads on 8x virtual partitions. The updated workloads of SQL server or distributed Database indexing are mostly zipf random 4K/16K write with high locality.

Significant performance improvement is observed with CSAL for this kind of workloads, which is key for D3C to adopt cost-effective QLC SSDs while still meeting the end users’ requirements of the Big Data application SLA.

- 8x jobs’ 64K Write performance with specific distribution ranging from zipf0.8 to zipf1.2 on 8x virtual partitions. Most real-world workloads follow the 20/80 principle which is cache friendly due to the zipf locality.
- Performance of mixed reads and writes and mixed sequential and random writes. Both modern Big Data and legacy applications have reads and writes as well as sequential and random writes mixed together.

4.1.1 4K Write Performance for Uniform Random Workloads

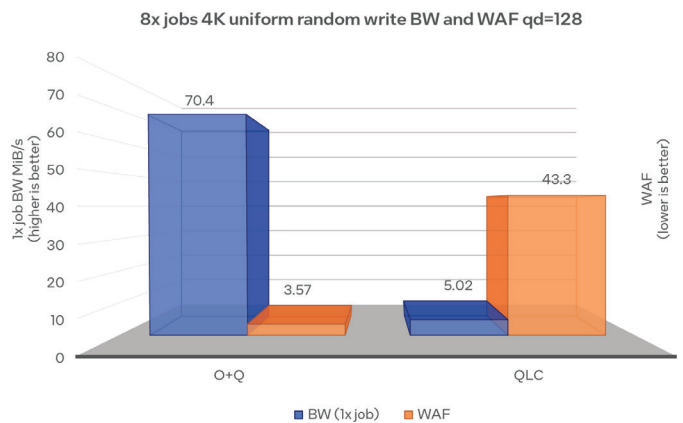


Figure 6: 4K uniform random write performance

Figure 6 shows that for 4K uniform random writes, the bandwidth (BW) that a O+Q combination provides for a single job is 70.4 MiB/s, about 14 times of the QLC BW (5.02 MiB/s), while the O+Q WAF is just 3.57, 8.2% of the QLC WAF (43.3).

4.1.2 4K Write Performance for zipf Random Workloads

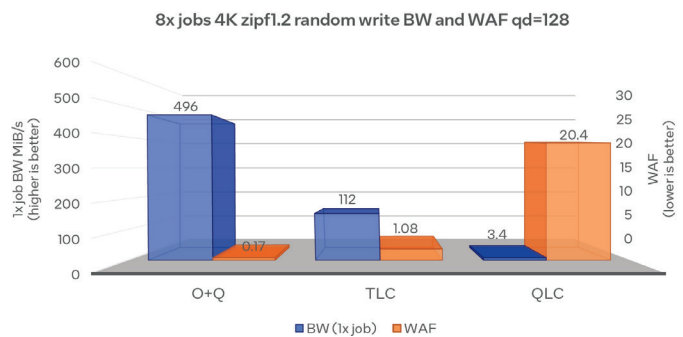


Figure 7: 4K zipf random write performance

Figure 7 shows that for workloads with high locality, such as zipf1.2 4K random write which follows the 20/80 real-world rule, the O+Q BW is 496 MiB/s, 2.5 times of the TLC BW (201 MiB/s) and 145 times of the QLC BW (3.4 MiB/s), while the O+Q WAF is only 0.17 (data written to QLC drive), 23% of the TLC WAF (0.735), and 0.8% of the QLC WAF (20.4).

4.1.3 64K Write Performance for zipf Random Workloads

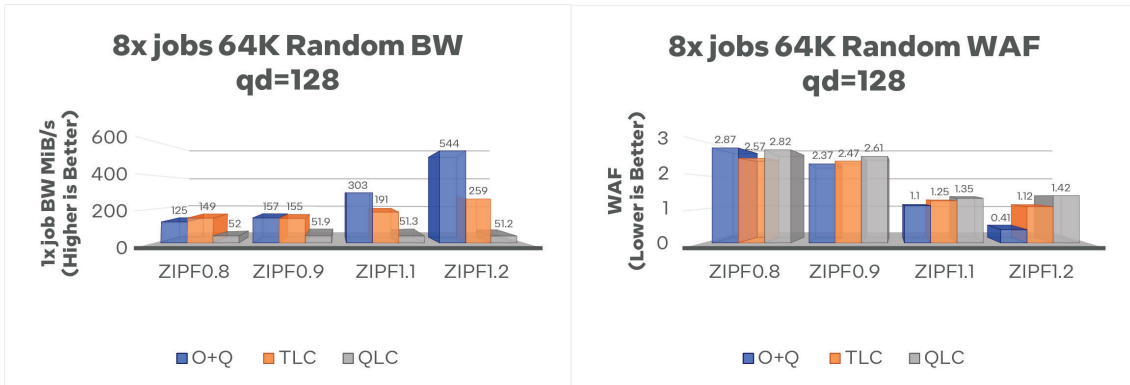


Figure 8: 64K zipf0.8/0.9/1.1/1.2 random write performance

Figure 8 shows that:

- As the workload features more locality, the O+Q solution improves the BW from 8x 125MiB/s for zipf0.8 to 8x 544MiB/s for zipf1.2, while the WAF drops from 2.87 for zipf0.8 to 0.41 for zipf1.2.
- When comparing the Q+Q configuration with the QLC only and TLC only configurations:
 - O+Q vs. QLC only BW for zipf0.8: $125 \div 52 = 2.4$
 - O+Q vs. QLC only BW for zipf1.2: $544 \div 51.2 = 10.6$
 - O+Q vs. QLC only WAF for zipf0.8: $(2.87 \div 2.82) \times 100\% = 102\%$
 - O+Q vs. QLC only WAF for zipf1.2: $(0.41 \div 1.42) \times 100\% = 29\%$
 - O+Q vs. TLC only BW for zipf0.8: $125 \div 149 = 0.84$
 - O+Q vs. TLC only BW for zipf1.2: $544 \div 259 = 2.1$
 - O+Q vs. TLC only WAF for zipf0.8: $(2.87 \div 2.57) \times 100\% = 112\%$
 - O+Q vs. TLC only WAF for zipf1.2: $(0.41 \div 1.12) \times 100\% = 36.6\%$

Overall, O+Q has bigger BW and WAF improvements over QLC and over TLC. The performance improvement on TLC is mostly observed on high locality workloads.

4.1.4 Read/Write Mixed Workloads

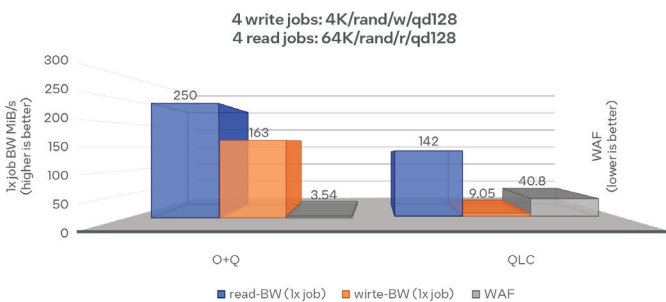


Figure 9: Mixed workloads with 4K random write and 64K random read

Figure 9 shows that for mixed workloads:

- The read BW for a single job on QLC is only 142MiB/s, which fails to meet the target BW—250MiB/s.
- The write BW for a single job on the O+Q combination is 163MiB/s, 18 times of the QLC BW (9.05MiB/s), while the WAF is 3.54, only 8.7% of the QLC WAF (40.8).

Key takeaways

Compared to a QLC only solution, CSAL delivers:

- up to 145 times BW and 0.8% of WAF for the 4K random write workload
- up to 10.6 times BW and 29% of WAF for the 64K random write workload
- up to 18 times BW and 8.7% of WAF for the mixed workload

4.1.5 Write Amplification Reduction with CSAL + ZNS

Figure 10 below shows write amplification reduced by using the CSAL solution with a ZNS drive as the capacity tier compared to the one with a regular SSD. In this benchmark test, several workloads were defined with different types of mixed write jobs (tenants)—each workload (pair of columns) include four jobs (tenants). CSAL isolates each job (tenant) into a separate set of zones, which results in write amplification reduction. The biggest WAF reduction is observed in multi-tenant workloads where each tenant has a different characteristic.

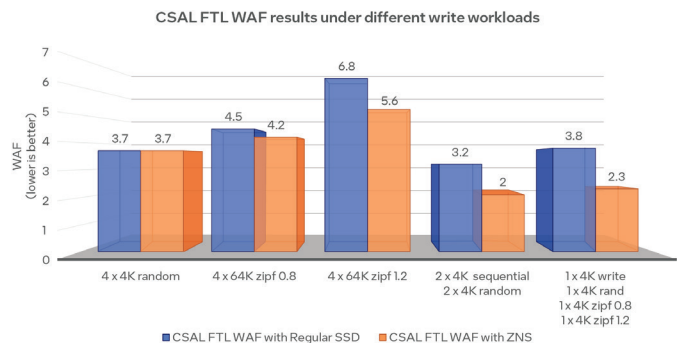


Figure 10: Write amplification reduction by CSAL with a ZNS drive

4.2 Performance for Cloud Production Workloads

With this new technology, Alibaba Cloud released a new D-series Big Data instance named ECS D3C instance in which Optane SSDs and QLC SSDs are used to replace HDDs for higher storage density and performance.

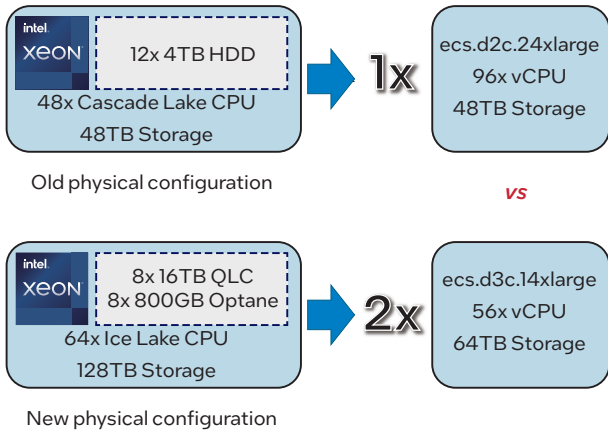


Figure 11: Alibaba Cloud D-Series

As shown in Figure 11, one physical server now provides 2x d3c.14xlarge instances (56x vCPUs and 64TB storage) with the new configuration while with the old configuration, one physical server just provides one d2c.24xlarge instance (96x vCPUs & 48TB storage). The following section describes performance comparison between ecs.d2c.24xlarge and ecs.d3c.14xlarge instances from the perspective of users in the cloud.

4.2.1 Spark Workloads on Alibaba D3C Public Cloud with Local Disk Services

There are two typical Big Data benchmarks: TPCx-HS and TPC-DS. We test d3c.14xlarge, the largest standard D-series instance on Alibaba Cloud, and the d2c.24xlarge instance using these two benchmarks with a 3TB dataset. Two Hadoop clusters have been set up using these Alibaba Cloud ECS Big Data instances, with three nodes for each cluster.

4.2.2 TPCx-HS Workloads

TPCx-HS (3TB dataset)	d2c.24xlarge	d3c.14xlarge	Improvement
HSGen (min)	7.11	4.16	70.91%
HSSort (min)	20.31	9.96	103.92%
HSValidate (min)	3.46	1.18	193.22%
Total Time (min)	31	15.25	103.28%
HSph@SF	1.9357	3.9354	

Table 2: Test results with TPCx-HS

TPCx-HS is the first industry standard Big Data benchmark designed to stress test on both hardware and software that is based on Apache HDFS API compatible distributions. Table 2 is the comparison between two D-series instances. You can see the time for each process and the total time in TPCx-HS. The D3C instance shows 103% improvement in total time for the 3TB data set. In each process, the D3C instance also outperforms the D2C instance by 70.91%, 103.92% and 193.22% improvement in HSGen, HSSort and HSValidate respectively.

4.2.3 TPC-DS Workloads

TPC-DS (3TB dataset)	d2c.24xlarge	d3c.14xlarge	Improvement
datagen (min)	40.8	41.93	-2.69%
SQL (min)	50.02	50.58	-1.11%
Total Time (min)	90.82	92.51	-1.83%

Table 3: Test results with TPC-DS

TPC-DS is a decision support benchmark that models several generally applicable aspects of a decision support system. The test indicators feature query response time, query throughput, and data maintenance performance with a given hardware, operating system, and data processing system configuration. Now TPC-DS enables emerging technologies, such as Big Data systems. We evaluate the Alibaba Cloud D-Series instances for TPC-DS workloads with a 3TB dataset. As shown in Table 3, in datagen and SQL processes, the D3C instance has 2.69% and 1.11% performance improvement compared with the D2C instance. For total time, the D3C instance shows 1.83% performance improvement over the D2C instance. TPC-DS is a compute-intensive workload. Storage performance is not a bottleneck for this workload. Therefore, the D3C instance achieves almost the same performance as the D2C instance does under TPC-DS with fewer CPU cores (96x vCPUs for d2c.24xlarge vs. 48x vCPUs for d3c.14xlarge). This is because the D3C instance is based on Intel® Xeon® CPU codenamed Ice Lake, which is more efficient than Cascade Lake that supports the D2C instance.

5 Conclusion

CSAL is a write shaping cache that unleashes the value of high-density NAND flash media. By leveraging the host-side FTL, CSAL preserves the existing software interface while transforming any write workload to a sequential write workload for ZNS flash storage. Furthermore, CSAL minimizes the frequency of writes by caching frequently updated or temporary data on Optane media or SCM. With these two strategies, CSAL mitigates endurance and performance challenges for modern flash media. Our tests demonstrate that CSAL brings substantial throughput improvement for write-intensive workloads compared to a NAND flash itself and a traditional cache (OpenCAS). CSAL is a software defined and flexible storage architecture for next gen media. It is easy to scale out in data centers. You can tune it to your different performance and TCO requirements.

6 Future of CSAL Software

While this paper showcases the performance data is primarily focused on using Optane SSD as the cache for CSAL, it's important to understand that the CSAL software component itself is a general-purpose storage shaping and caching software solution designed to unlock the benefits of media, such as SCM and future NAND media.

Intel is committed to working with its ecosystem, including server and storage OEMs as well as component suppliers, to support customers' businesses today and in the future. [Engage now](#) to lay the foundation for [CXL support](#) on Intel® processor platforms as the future standard of tiered-memory solutions. Please contact your Intel representative or send an email to wayne.gao@intel.com for any more information, for example, SLC SSD performance testing data.

7 Acronyms

Acronym	Description
Intel® 3D XPoint™	Intel's stable low latency persistent memory and storage media that supports 64 bytes level in-place write
AI	Artificial Intelligence
CSAL	Cloud Storage Acceleration Layer
DRAM	A physical memory in a computer
FTL	Flash Translation Layer
GC	Garbage Collection
HCI	Hyper Converged Infrastructure
HDFS	Apache's highly distributed file system
HPC	High performance computing
L2P	Logical to Physical
LBA	Logical Block Address
MASF	Media Aware Storage Framework
IU	Indirection Unit
NAND	Not AND, it is a circuit that builds a NAND flash-based SSD
OCF	Open CAS Framework
OpenCAS	Open Cache Acceleration Software
PLC	5 bits level cell
QLC	4 bits level cell
SLA	Service Level Agreement
SSD	Solid State Drive
WAF	Write Amplification Factor
WRF	Write Reduction Factor
SPDK	Storage Performance Development Kit
ZNS	Zoned Namespace

8 Reference

- [1] [IDC Global DataSphere Forecast, May 2022](#)
- [2] [Zoned Namespaces \(ZNS\) SSDs: Disrupting the Storage Industry](#)
- [3] [Achieving Optimal Performance and Endurance on Coarse-grained Indirection Unit SSDs](#)
- [4] [Media Aware Storage Framework \(MASF\), Solving the Challenges of Media Modernization](#)
- [5] [Intel® Optane™ Technology: Memory or Storage? Both](#)
- [6] [Intel Storage Analytics Standalone Linux IO tracer tool](#)
- [7] [Open Cache Acceleration Software](#)



Notices & Disclaimers

Performance varies by use, configuration and other factors. Learn more on the [Performance Index site](#).

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See Table 1 for configuration details. No product or component can be absolutely secure.

Your costs and results may vary.

Intel technologies may require enabled hardware, software or service activation.

Intel does not control or audit third party data. You should consult other sources for accuracy.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

Alibaba, the Alibaba logo, and other Alibaba marks are trademarks of Alibaba or its subsidiaries.