# Faster than Flash: An In-Depth Study of System Challenges for Emerging Ultra-Low Latency SSDs

Sungjoon Koh, Junhyeok Jang, Changrim Lee, Miryeong Kwon, Jie Zhang and Myoungsoo Jung

Computer Architecture and Memory Systems Laboratory, KAIST

{skoh, jhjang, crlee, mkwon, jie}@camelab.org, and mj@camelab.org

*Abstract*—Emerging storage systems with new flash exhibit ultra-low latency (ULL) that can address performance disparities between DRAM and conventional solid state drives (SSDs) in the memory hierarchy. Considering the advanced low-latency characteristics, different types of I/O completion methods (polling/hybrid) and storage stack architecture (SPDK) are proposed. While these new techniques are expected to take costly software interventions off the critical path in ULL-applied systems, unfortunately no study exists to quantitatively analyze system-level characteristics and challenges of combining such newly-introduced techniques with real ULL SSDs.

In this work, we comprehensively perform empirical evaluations with 800GB ULL SSD prototypes and characterize ULL behaviors by considering a wide range of I/O path parameters, such as different queues and access patterns. We then analyze the efficiencies and challenges of the polled-mode and hybrid polling I/O completion methods (added into Linux kernels 4.4 and 4.10, respectively) and compare them with the efficiencies of a conventional interrupt-based I/O path. In addition, we revisit the common expectations of SPDK by examining all the system resources and parameters. Finally, we demonstrate the challenges of ULL SSDs in a real SPDK-enabled server-client system. Based on the performance behaviors that this study uncovers, we also discuss several system implications, which are required to take a full advantage of ULL SSD in the future.

## I. INTRODUCTION

The state-of-the-art solid state drives (SSDs) begun to offer extremely high bandwidth by connecting with an on-chip processor controller, such as a memory controller hub or processor controller hub, via PCI Express (PCIe) buses [33, 36]. For example, NVMe SSDs (e.g., Intel 750 [30]) provide read and write bandwidths as high as 2.4 GB/s and 1.2 GB/s, respectively, which are approximately $4.3\times$ and $2.4\times$ higher than the bandwidth that a conventional SSD offers [31]. There also exist several industry prototypes that promise to deliver even much higher performance, ranging from 840 KIOPS to one million IOPS [26, 61]. Thanks to such performance superiority, high-end NVMe SSDs are widely used for diverse computing domains such as data-intensive applications [13, 15, 52, 53] and server systems as a disk drive cache [27, 58], burst buffer [47, 55], and in-memory computing storage [48, 52, 70].

Even though the bandwidth of such modern high-end SSDs almost reaches the maximum performance capability that PCIe buses can deliver, their system-level turn-around delays are still far different from those of other memory technologies or fast peripheral devices, such as DRAM and GPGPU. To bridge the latency disparity between the on-chip processor controller and underlying SSD, new flash based archives, called *Z-SSD*, have come to the attention of both industry and academia [46, 79]. This new type of SSDs can provide *ultra-low latency*

(ULL) behaviors, which exhibit great potential to get block devices close to computational components [33, 41, 63]. Specifically, the new flash medium that ULL SSDs employ is a revised version of vertically-stacked 3D NAND flash (48 layers) whose memory read latency for a page is 3 $\mu$s [24], which is $8\times$ faster than the fastest page access latency of modern multi-level cell (MLC) flash memory [18]. ULL SSDs are expected to satisfy a high level of service-level agreement (SLA) and quality of service (QoS) while providing a much higher storage capacity, compared with other types of new memory, such as resistive memory (ReRAM) [28].

To take the advantage of such fast storage, advanced types of I/O completion methods are proposed and implemented in modern Linux systems [74, 76, 79]. For example, [74] reveals that a *polled-mode* I/O completion method can be better than an interrupt-driven I/O completion. This claim is reasonably accepted in cases where the device-level latency of high-performance NVMe SSDs is shorter than the latency delays of context switches and interrupt service routine (ISR) management. To make polling more efficient, a hybrid polling is also proposed, which sleeps the polling processes at a certain period and then performs it after the period [23, 49]. These polled mode and hybrid I/O completion methods are newly implemented and published in Linux kernel 4.4 [22] and 4.10 [49], respectively. While Linux has simply applied aforementioned new types of I/O completion methods, the benefits of all those studies are evaluated on a DRAM emulation platform. Thus, the system performance of real ULL SSDs may be different than expected by the previous studies.

On the other hand, kernel-bypass storage architecture, such as storage performance development kit (SPDK) [76], can be employed to directly expose the latency superiority of ULL SSDs to user-level applications. Specifically, SPDK moves most drivers into the userspace, which aims to avoid multiple system calls and to eliminate redundant memory copies from the application accesses. However, because of this user-level implementation, SPDK cannot efficiently manage system interrupt handlers, which in turn makes SPDK employ the polled-mode I/O completion method instead of the interrupts. The real system therefore may not fully enjoy the benefits of SPDK, but, the system-level challenges imposed by such practical challenges of SPDK are not quantitatively analyzed and studied in the literature by far.

An industry article and gray literature uncover the low device-level latency characteristics of ULL SSDs [69]. However, it is unfortunately difficult to estimate their actual performance behaviors by considering diverse system execution parameters, such as different block sizes and queue depths. In addition, from the viewpoint of the advanced I/O completions

and storage architecture, it is non-trivial to analyze the system-level challenges that should be addressed to take full advantages of ULL behaviors since these new types of advanced SSD technologies are unavailable in a public market yet. In this work, we characterize the performance behaviors of a real 800GB Z-SSD prototype and describe a wide spectrum of challenges in integrating ULL SSDs into the current software storage stack.

The main questions that we want to address and observations can be categorized into three groups.

### A. System-Level Performance and ULL-Specific Behaviors

Regarding new flash, there are a very few studies to demonstrate low-level circuitries and peripherals [18] or high-level marketing data [59]. However, designing an I/O path with new technologies requires considering many combinations of system and storage parameters, such as diverse queue and block device configuration values. In addition, many system designers need to understand ULL-specific performance characteristics, such as garbage collection and read-and-write interference.

- How fast are ULL SSDs compared to NVMe SSDs?
- How much can ULL SSDs cut the long-tail latency?
- Is the traditional NVMe multi-queue mechanism also affordable for ULL SSDs?
- Do ULL SSDs also have a critical path that most flash suffers from?

### B. System Impacts with Polled-Mode Completion

Modern polled-mode I/O completion and hybrid polling mechanism are designed to a new storage stack for high-performance SSDs. Since high-end NVMe SSDs have potential to serve most read and write requests in several $\mu$s, [74] studies that the polled-mode I/O completion method can reduce the latency further in addressing the software overheads related to interrupt handling. Motivated by this, Linux already implemented those new I/O completion methods in its kernel and published them [22, 49].

- Is the polled-mode I/O completion faster than interrupts?
- What are the practical overheads when systems employ the polled-mode?
- Will be a hybrid polling better than the polled-mode?

### C. Advanced Storage Stack Analysis

Linux storage stack is considered as the performance bottleneck when ones use high-end NVMe SSDs [43, 76]. Intel SPDK reorganizes the storage stack such that users can bypass the kernel and directly access the underlying storage. A common expectation behind this kind of kernel-bypass schemes is to reduce software intervention and shorten the storage latency.

- Can SPDK really eliminate the latency overheads, imposed by the storage stack? What about ULL SSDs?
- Will SPDK be the best option for future low-latency storage? Would it have any side effects?
- How much can SPDK reduce the system-level latency under a real server execution environment?

| 3D NAND | BiCS | V-NAND | Z-NAND |
|---|---|---|---|
| # layer | 48 | 64 | 48 |
| $t_R$ | $45\mu s$ | $60\mu s$ | $3\mu s$ |
| $t_{PROG}$ | $660\mu s$ | $700\mu s$ | $100\mu s$ |
| Capacity | 256Gb | 512Gb | 64Gb |
| Page Size | 16KB/Page | 16KB/Page | 2KB/Page |

TABLE I: Analysis of 3D flash characteristics [18].

## II. BACKGROUND

### A. Ulta-Low Latency SSD

*1) New Flash:* Modern SSDs can satisfy a high bandwidth requirement of users by exploiting various architectural supports, such as internal parallelism, I/O queuing/scheduling and DRAM buffers. However, shortening the latency for a basic unit of I/O operation requires low-level memory design changes and device updates. As shown in the right most of Figure 1, new flash memory, called *Z-NAND*, leverages single-level cell (SLC) based 3D flash design, but optimizes several I/O circuitries, such as a default I/O page size and DDR interfaces, to offer the low flash latency and shorter data transfer delays, respectively [24]. Table I summarizes the device-level characteristics of three different state-of-the-art 3D flash technologies: i) Bit Cost Scaling (*BiCS*) 3D flash [73] ii) Vertically stacked (*V-NAND*) 3D flash [42] and iii) ULL-tailored flash (*Z-NAND*) [18]. Z-NAND uses 48 stacked word-line layer, which exhibits $3\mu$s and $100\mu$s for a read operation and a write operation, respectively. The write latency of Z-NAND is shorter than that of BiCS and V-NAND by $6.6\times$ and $7\times$, respectively, while its read latency is $15\sim20 \times$ shorter than those of such two modern 3D flash technologies. Even though the storage capacity and page size of Z-NAND are smaller than those of BiCS/V-NAND, ULL SSDs can offer a bigger storage volume with shorter latency by putting more Z-NAND packages into their device platform as a scale-out solution.

*2) Multiple Super-Channel Architecture:* The storage architecture of high-end SSDs consists of multiple system buses, referred to as *channels*, each employing many flash packages via multiple datapaths, called *ways* [25, 34, 40]. Similarly, ULL SSDs adopt this multi-channel and multi-way architecture, but further optimize the datapaths and its channel-level striping method. As described in Table I, the basic I/O unit of Z-NAND (i.e., page) is much smaller than that of other flash technologies, which in turn can introduce a higher level of device-level and flash-level parallelism by serving a host request with finer granular operations [34, 38, 40]. Specifically, ULL SSDs split a 4KB-sized host request into two operations and issue them to two different channels simultaneously, as shown in Figure 1 (left). These two channels always handle flash transactions together as a pair of system buses, which are referred to as a *super-channel*. To efficiently manage data transfers and flash transactions upon the super-channel architecture, ULL SSDs exploit an optimized circuit that automatically adjusts data-flow and manages a timing skew on individual channels [18]. This circuit, called *split-DMA* management engine, can reduce the read access time thereby tailoring ULL further. One of the concerns behind the super-channel architecture is to manage wear-out blocks (i.e., bad blocks) each potentially existing in different channels. Since two operations spread across a pair of system buses based on
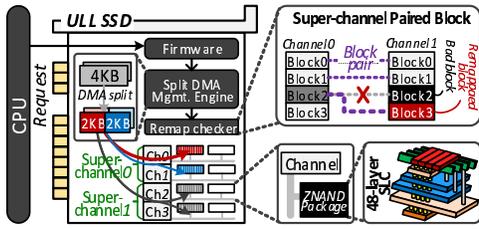
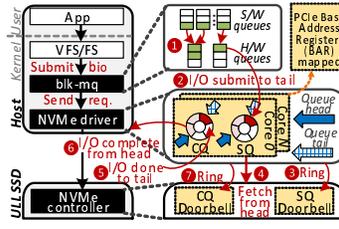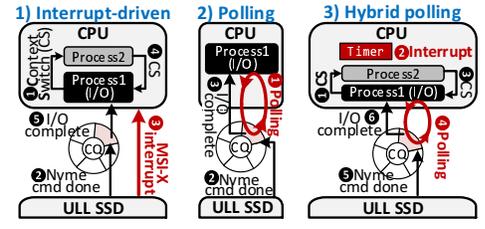Fig. 1: ULL SSD internals and Split DMA.



Fig. 2: NVMe storage stack.



Fig. 3: NVMe I/O completion methods.

their offsets (within a super-channel), flash firmware can waste the storage space if the bad blocks appear in one of the channel pairs. To address these challenges, the split-DMA management engine also employs a remap checker to automatically remap the physical block addresses between a bad block and a clean block. The remap checker exposes this semi-virtual address space to the flash firmware, such that the storage space can be fully utilized with the super-channel technique.

*3) Suspend and Resume Support:* To reduce the read latency more, Z-SSDs also apply a suspend/resume DMA technique [18]. While the split-DMA management engine can shorten the latency with a finer granular data access, it may not be able to immediately serve a read request if the target super-channel, associated with the data, is busy to transfer a write request, which was issued at an earlier time. This resource conflict can increase the latency of such reads, waiting for a service in a device-level pending queue. The suspend/resume DMA method of ULL SSDs pauses the write service in progress for the super-channel that the read targets by storing an exact point of the flash write operation into a controller-side small buffer. The suspend/resume DMA engine then issues the read command to its target flash package without a wait for resolving the resource conflict. Once the read has been successfully committed, which takes only a few cycles, the engine resumes the write by restoring the context (e.g., the stored write point). Thus, this suspend/resume mechanism can maximize the resource utilization of multi-way of super-channels and reduce the overall latency of ULL SSDs, thereby satisfying diverse levels of QoS and SLA.

### B. Advanced Storage Stack

*1) NVMe Storage Stack:* In conventional systems, most block I/O services should be requested to kernel-level I/O modules through system call(s). The requests are delivered at first to virtual file system (VFS) and native file system (e.g., ext4). They not only provide compatible user interfaces and file management, but also perform page caching and guarantee ACID (atomicity, consistency, isolation, durability) of the underlying storage. Underneath the file systems, a multi-queue block layer, `blk-mq` [11], schedules I/O requests while an NVMe driver manages block storage protocols. The `blk-mq` employs two queues, each corresponding to a software queue and hardware queue. The software queue exists per CPU core and is responsible for file system's block request (i.e., `bio` structure) handling. On the other hand, the number of existing hardware queues equals the number of queues provided by the underlying NVMe driver provides and is used for delivering `bio` requests to the NVMe driver.

*2) NVMe Multi-Queue Mechanism:* All NVMe SSDs, including ULL SSDs, expose a certain memory region, which is

used for multi-queue communications to the host-side NVMe driver through *PCIe base address registers* (BARs) [7]. When a host issues an NVMe command (as a form of NVMe queue entry) to the underlying NVMe SSD, it is enqueued into the target *submission queue* (SQ), and SQs are managed in a first-in first-out manner. To complete the corresponding I/O service, the NVMe SSD requires issuing a completion command to the corresponding *completion queue* (CQ). Each queue can accommodate $2^{16}$ NVMe command entries, and all the queues that the target SSD can support are mapped to the BARs. To synchronize the SQ/CQ states between the host and SSD, NVMe also specifies a pair of *doorbell* (DB) registers, each being allocated to an SQ and CQ, respectively. If a host module or SSD firmware module updates any entry of the SQ or CQ, it can simply inform such update by writing to the corresponding DB; as all the DBs are mapped to the BARs, they are visible for both the host and SSD. While this multi-queue mechanism is affordable to fully utilize the PCIe bandwidth capability, the queue structure itself is very rich and complicated, which can increase queue waiting time at some extent.

*3) I/O Completion Methods for NVMe:* When there is an update for a CQ, the SSD device informs its completion to the host through a *message signal interrupt* (MSI). If an MSI is reported, the host's NVMe driver is required to complete the request by checking the target CQ and handling ISR. This interrupt-driven I/O path is sufficient to manage slow block devices, but as NVMe SSDs provide 3.2× higher bandwidth and 1.8× shorter latency than conventional SSDs (e.g., SATA SSD), a polled-mode I/O completion method has attracted the attention of both academia [23, 74, 79] and industry [23, 49] instead of interrupts. Typically, `blk-mq` polls hardware queues to check the target entry by retrieving the corresponding queue and request information from a bitmap mask, called *cookie* (returned at the submission time). As the hardware queues are one-to-one matched with the NVMe queues (offered by the SSD device), the NVMe driver provides a polling interface, called `nvme_poll()`, to `blk-mq`. In this function, the corresponding CQ entry is iteratively checked. If an update of the CQ entry is detected during polling process, it stops polling and completes the corresponding I/O process by returning the results to the users.

This polling method is implemented in Linux 4.4 kernel, but from 4.10 kernel, Linux supports a hybrid polling [23, 49]. The hybrid polling sleeps for a while in an attempt to reduce the number of polling the queues. Specifically, the hybrid polling method checks the queue like the polled-mode I/O completion, but it adds up a certain time of sleep before jumping to the poll process. In the latest version of kernel that we tested (Linux 4.14), the hybrid polling calculates an average time for the
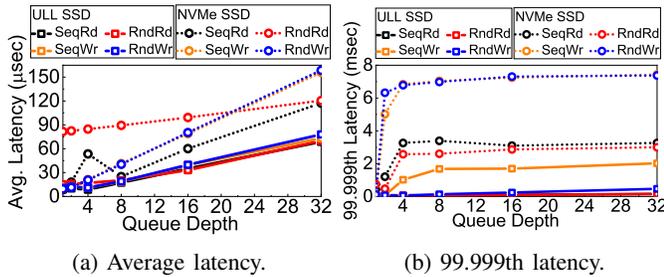
(a) Average latency.　　　　(b) 99.999th latency.

Fig. 4: Latency analysis of NVMe SSD and ULL SSD.

previous I/O completions and then sleeps the poll process as long as the half of the average time, monitored.

*4) Kernel-bypass (SPDK):* To reduce the overheads brought by the storage stack, kernel-bypass schemes are applied in modern NVMe SSDs [43, 44, 57, 60, 76]. Intel SPDK is a practical approach that implements the NVMe driver at user-level [76]. To this end, SPDK unbinds the underlying NVMe device and rebinds it to a *user-space I/O driver* (uio) or *virtual function driver* (vfio) driver. SPDK maps PCIe BARs to a huge page [2], which is allocated and managed by a kernel and network bypass-framework, called *DPDK* [6]. Since the huge page is mostly not swapped out, uio/vfio of SPDK can directly check NVMe queues (SQs/CQs) and doorbell registers from the user-level, which makes it bypass the storage stack and allows user applications directly to access the underlying NVMe SSD. However, ISR cannot be handled from the user level driver, only polling can be adopted as an I/O completion routine.

## III. EVALUATION SETUP

### A. Benchmark

To characterize NVMe and ULL SSDs, we use FIO (v3.13) as our microbenchmark suite [9]. We set an O_DIRECT flag for all evaluations to bypass page caches and directly serve the I/O requests to/from the underlying SSDs. In this test, we also use Linux native AIO (libaio) [10] as an I/O engine to generate asynchronous block I/O requests. Even though we test the SSDs with different block I/O sizes, ranging from 4KB to 32KB, for specific evaluations, such as performance analysis and CPU utilization, we configure the default block size as 4KB. On the other hand, synchronous preadv2/pwritev2 (pvsync2) is enabled by the I/O engine to analyze the system impacts brought by different types of I/O completion methods.

### B. Device Configuration and Profilers

We use a testbed that employs a 4.6 GHz, 6-core Intel i7 processor (i7-8700) and 16GB DDR4 DRAM. All OS modules and executables are stored and executed from an external 400GB SSD [30]. We evaluate an NVMe SSD by using Intel 750 SSD [30] as it is the only commercially available device that uses the PCIe interface (without a M.2 bridge) over a standard NVMe protocol in the market. To make the underlying NVMe SSD up-to-date, we use the latest firmware (8EV101H0). We evaluate a 800GB Z-SSD prototype as an ULL SSD. In this test, all the SSDs are connected to the host of the testbed via PCIe 4x 3.0 lanes. For this study, we use CentOS 7.6 and Linux kernel 4.14, which contains the most recent and stable version of NVMe storage stack. To make sure that CPU cores do not reside on the critical path, we configure the policy of CPUfreq governors (in the Linux

kernel) with "performance", which statically sets CPU to the highest frequency in the range of target CPU scaling. Lastly, we use the FIO report and Intel Vtune Amplifier 2019 [3] to characterize memory access patterns that the NVMe storage stack exhibits.

### C. Storage Stack

To evaluate the polled-mode I/O completion, we set the polling flag (queue_io_poll) of Linux pseudo file system (sysfs) as '1'; otherwise, it is configured with '0' for all the tests that we performed. In addition, the priority of I/O requests (i.e., hipri flag) is set with a value, higher than that of other tasks; we make sure that only one core, running at the maximum frequency, is utilized to manage incoming I/O requests. For kernel-bypass, we use SPDK 19.07 [5] and configure the size of a huge page with 2MB. The huge page is directly mapped to PCIe BARs by DPDK 19.02 [6] in our evaluation. For the user-level NVMe driver, we use uio 0.01.0 [45] and utilize fio_plugin() (provided by SPDK) for our SPDK-enabled benchmark executions.

## IV. SYSTEM-LEVEL ANALYSIS FOR ULL SSDs

### A. How fast are ULL SSDs compared to NVMe SSDs?

Figure 4a shows the overall latency characteristics of ULL SSD and NVMe SSD that we tested with varying I/O depths, ranging from 1 to 32. With a low I/O queue depth (1~4), the write latency of NVMe SSD is around $14.1\mu s$, which is slightly worse than that of ULL SSD; ULL SSD offers $12.6\mu s$ and $11.3\mu s$ for reads and writes, on average, respectively. The reason why NVMe SSD can exhibit much shorter write latency than its actual flash execution time is that it caches and/or buffers the data using their large size of internal DRAM. Nevertheless, NVMe SSD cannot hide the long latency, imposed by the low-level flash, in the case of random reads ($82.9\mu s$). Specifically it is $5.2\times$ slower than ULL SSD ($15.9\mu s$). This is because the low locality of random reads enforces the internal cache to frequently access the underlying flash media, which makes NVMe SSD expose the actual flash performance to the host. Similarly, as the queue depth increases, the execution time characteristics of NVMe SSD significantly get worse, and its latency increases as high as $121\mu s$ and $159\mu s$ for random writes and reads, respectively. In contrast, as shown in the figure, ULL SSD provides reasonably sustainable performance even in the test with the high I/O queue depth.

### B. How much can ULL SSDs cut the long-tail latency?

The performance difference between NVMe SSD and ULL SSD becomes more notable when we examine their long tail latency behaviors. Figure 4b analyzes *five-nines* (99.999%) latency for both NVMe SSD and ULL SSD. The evaluation results show the worst case performance brought by each SSD's low-level flash due to internal tasks, such as garbage collection, and/or internal DRAM cache misses. For example, even though the average latency of random/sequential writes is better than that of reads, in this long tail latency evaluation, the writes are worse than the random reads by $2.1\times$, on average. In overall, NVMe SSD increases the five-nines latency of reads and writes than the average latency of them by $17.9\times$ and $108\times$, respectively. The reason behind this significant performance degradation is that the the five-nines latency cannot take
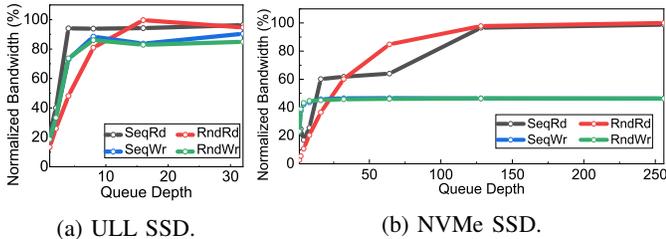
(a) ULL SSD.

(b) NVMe SSD.

Fig. 5: Bandwidth analysis (normalized to max B/W).

(a) Average power consumption.

(b) Latency degradation by GC.

Fig. 7: Power analysis and garbage collection.
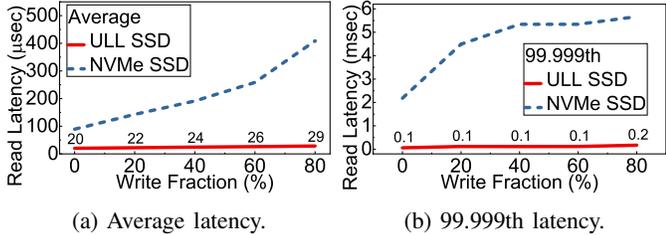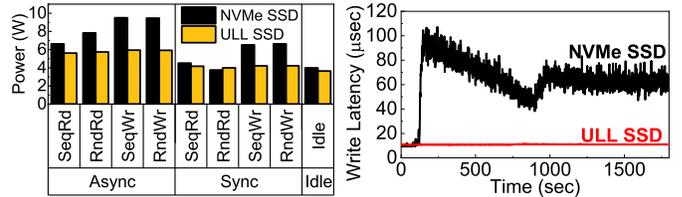
(a) Average latency.

(b) 99.999th latency.

Fig. 6: I/O interference analysis.

an advantage of NVMe SSD's architectural supports due to many systemic problems (e.g., resource conflicts, insufficient internal buffer size), which are insufficient to accommodate all the incoming I/O requests, and heavy internal tasks [37, 39]. In contrast, ULL SSD exhibits much shorter latencies, ranging from a few $\mu s$ to hundreds $\mu s$, for both the reads and writes. In contrast to NVMe SSD, the backend flash media of ULL SSD offers significantly low latency by new memory design and device-level optimizations. These low-level optimizations and Z-NAND characteristics not only can reduce memory-level latency further but also offer a better opportunity to fully reap the benefits of multi-channel and multi-way architectural support.

*C. Is the traditional NVMe multi-queue mechanism also affordable for ULL SSDs?*

Figure 5 shows the bandwidth utilizations of NVMe SSD and ULL SSD with varying I/O depths (1∼256). While ones can expect that the level of parallelism can increase as the queue depth increases (thereby higher bandwidth), unfortunately, NVMe SSD cannot reach the maximum bandwidth for the service of many 4KB-sized I/O requests. NVMe SSD only utilizes 40% of the total performance capacity for the 4KB-sized writes. Interestingly, in contrast to the previous read latency evaluations, the bandwidth of NVMe SSDs for random reads outperforms that of all other I/O patterns (that we tested) with higher queue depths (more than 128). This is because, with more I/O requests, scheduled in the queue, SSDs can easily find out a set of flash media that can simultaneously serve multiple I/O requests (by spreading the requests across different flash dies in parallel). Thus, random and sequential reads on NVMe SSD can offer the maximum bandwidth (1.8GB/s). In contrast, the bandwidth utilization of ULL SSD bumps against the maximum bandwidth for all the read test scenarios, and even for the sequential and random writes, ULL SSD can utilize 90% and 87% of the total bandwidth, on average, respectively.

It is worthwhile to report that, ULL SSD needs only "8 queue entries" for the sequential accesses; even in the worst case, 16 queue entries are sufficient to achieve the maximum bandwidth that ULL SSD offers. We believe that the rich

queue mechanism and software-based protocol management of NVMe (which are managed by an NVMe driver and/or Linux `blk-mq`) are very reasonable design strategies to maximize the benefits of modern high performance SSDs in the sense that such SSDs require securing many I/O resources for higher parallelism. However, once the latency becomes shorter by employing new flash or memory technologies (similar to ULL SSD), we believe that the rich queue and existing NVMe protocol specification are *overkilled*; a future ULL-enabled system may require to have a lighter queue mechanism and simpler protocol, such as NCQ [29] of SATA [62].

*D. Do ULL SSDs also have a critical path that most flash suffers from?*

*1) I/O Interference Impact:* Figure 6 analyzes the degree of I/O interference when reads and writes are intermixed. For this analysis, we randomly read data from NVMe SSD and ULL SSD by sporadically issuing the writes among many read requests. In addition, we increase the write fraction of the total I/O executions, ranging from 20% to 80%, in an attempt to analyze the different I/O interference behaviors. As shown in Figure 6a, the average read latency of NVMe SSD linearly increases as the write fraction (in the intermixed workloads) increases. Even when only 20% writes are interleaved with the reads, the writes seriously interfere I/O services of such reads, which can make the read latency worse than that of read-only workloads by $1.6\times$ ($54\mu s$, on average). There are two root causes. First, a write operation of a conventional flash (at memory-level) takes a significantly longer time than a read operation ($19\times$ at most), and the write blocks all other subsequent read services. Second, data transfers for the write (4KB) also occupy a specific channel for around $60\mu s$, which prevents an incoming read command from being issued to a target memory in the channel. In contrast, we observe that ULL SSD exhibits sustainable latency behaviors irrespective of the amount of write operations (interleaved with the reads). These anti-interference characteristics of ULL SSD are captured by our five-nines latency analysis as well. As shown in Figure 6b, while the five-nines read latency of NVMe SSD increases as high as 4.5ms even with 20% sporadic writes, ULL SSD maintains its latency under $118\mu s$. Since modern file systems and OS kernels are required to periodically write metadata or to perform journaling [65], we believe that ULL SSD can be a more desirable solution even in cases where many data-intensive applications intensively reads multiple data.

*2) Power and Garbage Collection:* Figure 7a analyzes power consumption characteristics of NVMe and ULL SSDs. ULL SSD consumes power less than NVMe SSD for asynchronous I/O operations by 30%, on average. In particular,
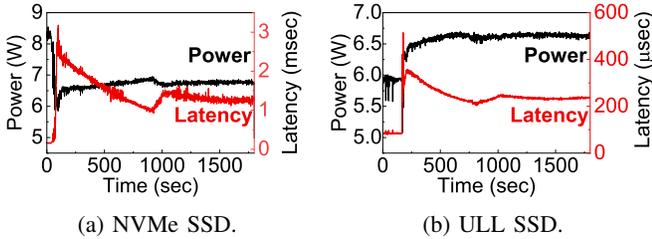
(a) NVMe SSD.          (b) ULL SSD.

Fig. 8: Power consumption during garbage collection.



(a) Seq. reads.   (b) Rnd. reads.   (c) Seq. writes.   (d) Rnd. writes.

Fig. 9: Latency comparison (interrupt vs. poll) in NVMe SSD.



(a) Seq. reads.   (b) Rnd. reads.   (c) Seq. writes.   (d) Rnd. writes.

Fig. 10: Average latency of interrupt and polling based I/O in ULL SSD.

disparities of the power consumption trend between NVMe SSD and ULL SSD become more notable for the write services (compared with reads). We conjecture that this is because of SLC-like-Z-NAND; SLC in practice employs much fewer steps of writes (per request) on the target flash cell, compared to MLC [17, 20, 67, 72], which exhibits lower write latency. On the other hand, the power consumptions in idle states and sequential/random reads are similar to each other; 3.8W and 4.1W are consumed for the idles and read services, respectively. Flash reads require enabling only internal peripherals to sense out the target data (whose dynamic power is much lower than write power). This makes the power of internal DRAM buffers and controllers dominant for the I/O services rather than the power of backend flash.

One of the performance bottlenecks that most SSDs suffer from is garbage collections [21, 35, 37, 39, 64]. Since flash does not allow to overwrite data without erasing a flash block, flash firmware [40, 50, 51] forwards all incoming overwrite requests to a new space (i.e., block), which was erased in a previous stage, and remaps the target addresses, associated with the request, to the new block. If there is no available page in the remaining block to serve the overwrites, the firmware should reclaim pages by erasing a block. It then reads the corresponding data and writes them to the erased new block, which is called a *garbage collection* (GC). Since it introduces multiple extra reads and writes, GCs exhibit long latency, which in turn can consume more power than a normal operation.

Figure 7b shows a time series analysis, which keep randomly writing 4KB data block after writing the entire address range of the underlying SSD. One can observe from this figure that the write latency of NVMe SSD sharply increases once GCs begin to reclaim flash blocks ($107\mu s$). The overall write latency with GCs is $6.3\times$ higher than the average latency of writes ($52\mu s$). In contrast, ULL SSD exhibits a very sustained latency behavior with the same evaluation scenario that we tested for NVMe SSD (less than $1\mu s$). We conjecture that there are three reasons behind ULL SSD does not suffer from GCs. First, as analyzed earlier, most of the flash reads and writes of ULL SSD are simply served faster than those of NVMe SSD. This flash-level performance superiority can shorten the long latency, imposed by GCs. Second, ULL SSD employs multiple super-channels, which can reclaim flash blocks with a higher parallelism by enabling more flash chips at a given time. Third, the suspend and resume operations (cf. Section II-A3) can make GCs interleaved with many incoming write requests, which makes the GC latency invisible to users at some extent. To be precise, we also analyze the power consumption behaviors of GCs. As shown in Figure 8a, power consumption of
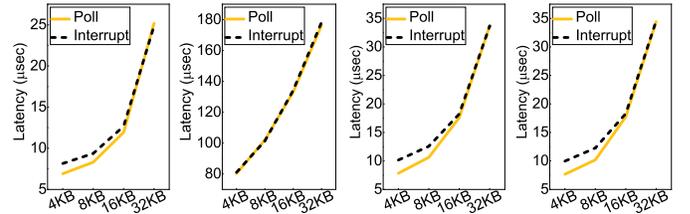
NVMe SSD decrease when it starts GCs; we believe that this is because a few flash chips are involved in GCs per I/O request arrival (thereby low GC performance). In addition, as the data migration between the old flash block and new flash block makes multiple SSD resources busy (e.g., channels, die and planes), NVMe SSD cannot serve I/O requests while GC is in progress. As shown in Figure 8b ULL SSD shows completely different story, compared with the power consumption trend of NVMe SSD. While the GC latency of NVMe SSD is as high as *3ms*, that of ULL SSD is around $500\mu s$, which is even faster than a write latency of conventional flash. Once a GC is invoked, ULL SSD consumes 12% more power than a non-GC workload execution, on average. As described earlier, this is because ULL SSD enables many flash chips and perform GCs in parallel, so incoming requests targeting a specific flash chip can still be serviced by leveraging the suspend/resume operations.

## V. I/O COMPLETION METHODS AND CHALLENGES

### A. Is the polled-mode I/O completion faster than interrupts?

*1) Overall Latency Comparison:* Figure 9 shows the latency difference between polling-based and interrupt-based I/O services with NVMe SSD. While polling gets attention from both industry and academia as a promising solution to expose the true latency of PCIe based SSDs to user-level applications, one can observe from the figure that polling unfortunately has no performance impact with the modern high-end SSD technology. Specifically, the latency difference of reads and writes brought by the interrupt-based and polling-based I/O services is less than 2.2% and 11.2%, on average, respectively. In contrast, as shown in Figure 10, polling improves the performance compared to the conventional interrupt-based I/O completion when ULL SSD is applied to the current NVMe stack. The average read and write latency of ULL SSD with polling are respectively $9.6\mu s$ and $9.2\mu s$ (for 4KB-sized requests), while those with interrupts are $11.8\mu s$ and $11.2\mu s$, respectively. Even though polling shortens the read and write

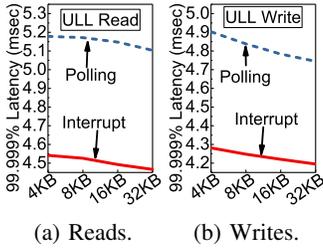(a) Reads.  (b) Writes.

Fig. 11: 99.999th latency of ULL SSD.



Fig. 12: CPU utilization of hybrid polling.



(a) Module.  (b) Function.

Fig. 14: CPU utilization breakdown.



(a) Reads.  (b) Writes.

Fig. 15: Normalized memory instruction count of poll.



(a) Seq. reads.  (b) Rnd. reads.  (c) Seq. writes.  (d) Rnd. writes.
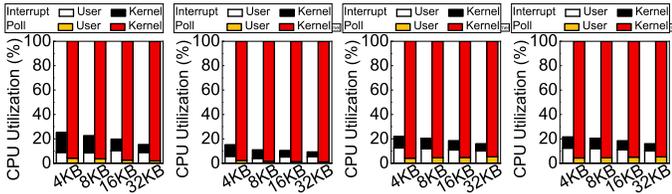
Fig. 13: CPU utilization of interrupt/polling based I/O.

latency by only 16.3% and 13.5% on average, respectively, we believe that the polling mechanism in the NVMe stack can bring an advantage for latency sensitive applications. The benefits can be more notable with future SSDs that employ faster NVM technologies such as resistive random access memory (ReRAM) [19, 32, 78, 80]. To appropriately exploit the polling method in the future systems, we believe that there exists several system-level challenges that the conventional NVMe storage stack should address; we will analyze the challenges shortly.

*2) Five-Nines Latency:* Figures 11a and 11b analyze five-nines latency for reads and writes with two different I/O completion methods, respectively. In contrast to our previous observations, the long tail latency of polling for reads and writes is worse than those of interrupts by 12.5% and 11.4%, on average, respectively. This should be addressed in the future NVMe storage stack, as well. The key functions of polling (cf. nvme_poll and blk_poll) require acquiring spin locks when they process CQ(s), which should be iterated until the target request is completed. The polling may not release and relax CPU as spin locks are used for polling NVMe queues. Thus, other incoming requests will be pending as polling does not allow a context switch while its I/O request is outstanding. CPU scheduling should be revised to address the aforementioned shortcomings, exposed by the current polled-mode operation [1].
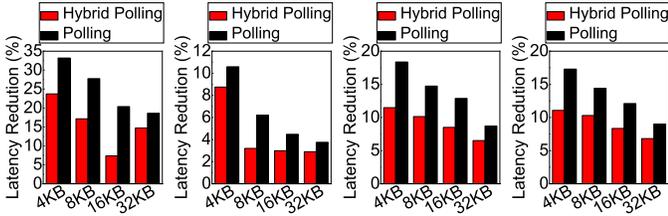
*B. What are the practical overheads when systems employ the polled-mode?*

*1) CPU Utilization Analysis:* Figure 13 compares the CPU utilization of two systems, each employing the interrupt-based and polling-based I/O services. In this evaluation, we measure the CPU usages by reading/writing 4KB-sized data from/to ULL SSD with those two different I/O completion methods. As shown in the figure, the interrupt-based I/O services only take 8.4% and 9.2% of the total CPU cycles for kernel mode and user mode, on average, respectively. However, the polling shows a completely different trend on the CPU utilization analysis, compared to that on interrupts. While the CPU cycles

consumed by polling at the user-level are similar to those of interrupts, such CPU cycles (of polling) at the kernel-level accounts for 96.4% of the entire application execution, on average. This is because the NVMe driver does not release the CPU and keeps iterating the completion procedure. Specifically, it looks up NVMe's completion queue (CQ) head/tail pointers and checks all phase tags within the CQ's entries when the driver polls the I/O completion.
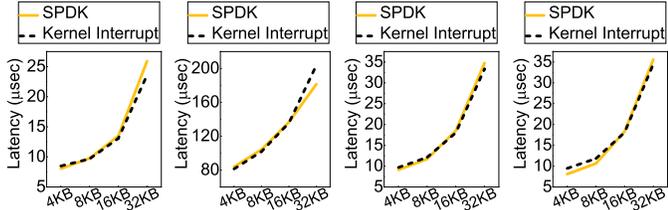
Figure 14a decomposes the kernel-side CPU cycles consumed by a single NVMe driver and other modules (such as blk-mq). We found out that the NVMe driver uses only 17.5% of the total CPU executions for all I/O access patterns. Then, we decomposes CPU utilization into function-level to study how the NVMe driver and I/O related modules spend the CPU cylces cycles. Figure 14b shows CPU utilization breakdown results, obtained by VTune amplifier performance profiler; the CPU cycles consumed by different stack functions are normalized to the total CPU cycles used for the I/O services. It shows two major functions: i) blk_mq_poll() and ii) nvme_poll(). These two functions are respectively associated to blk_mq of the storage stack and the NVMe driver. Surprisingly, only these two functions consume 84% of the total CPU cycles among all kernel modules in the storage stack. blk_mq_poll() takes around 67% CPU cycles to check if the current thread needs to be rescheduled and whether a pending request exists, coming from other kernel modules or processes. nvme_poll(), the callee of blk_mq_poll(), simply loads the target CQ entries and checks up if the corresponding I/O requests finish. Since the target entries should be written by the underlying SSD (over DMA), it is necessary to keep loading the CQ entries, repeatedly. This simple load and check up procedure takes 17% of entire CPU cycles, on average, and this accounts for 97% of cycles spent by NVMe module. We believe that, even though polling can shorten the device-level latency, allocating an entire core to refer the I/O completions can hurt the overall system performance as it prevents several computational tasks from running on the host system in a timely manner.

*2) Memory Requirements:* Figure 15 normalizes the number of memory instructions (load/store) for the polled-mode system executions to that of a conventional interrupt-driven system. As polling requires checking up all target CQ entries, it exhibits 137% more load instructions (to compare and check the I/O completion) than what the interrupt-driven system introduces. Note that the target CQ entries are managed by the underlying SSD devices, and thus, loads should go through all CPU caches and bring them back to the caches again; the caches always hold stale data for the CQ entries' information. Thus, actual execution of the load instructions is significantly

(a) Seq. reads. (b) Rnd. reads. (c) Seq. writes. (d) Rnd. writes.

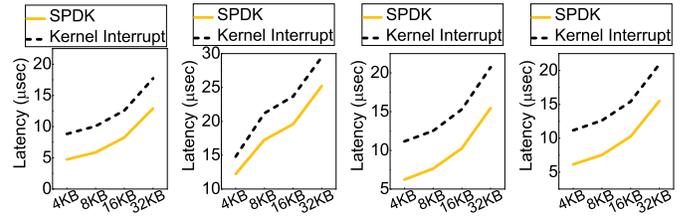Fig. 16: Latency reduction of hybrid polling.



(a) Seq. reads. (b) Rnd. reads. (c) Seq. writes. (d) Rnd. writes.

Fig. 17: Avg. latency of interrupt and SPDK in NVMe SSD.



(a) Seq. reads. (b) Rnd. reads. (c) Seq. writes. (d) Rnd. writes.

Fig. 18: Avg. latency of interrupt and SPDK in ULL SSD.



(a) Seq. reads. (b) Rnd. reads. (c) Seq. writes. (d) Rnd. writes.

Fig. 19: Average latency of interrupt and SPDK with big size requests in ULL SSD.

inefficient compared to the instructions which are usually performed to manage conventional memory subsystems. Similarly, polling exhibits 78% more store instructions to serve I/O operations than interrupts, on average. Even though polling spends most times to examine the CQ entries' states, polling also keeps metadata, indicating how many loop iterations are performed for a poll invocation (at `blk-mq`). Note that these activities of polling is not observed in the interrupt-driven method. We believe that this kind of store instructions can be reduced by optimizing memory management routine of `blk-mq` and NVMe driver for the polled-mode I/O completion in the near future.
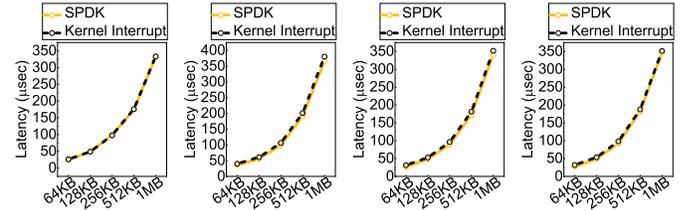
### C. Will be a hybrid polling better than the polled-mode?

The hybrid polling is expected to reduce the CPU cycles, wasted by intensively polling CQs. Figure 16 shows how much the latency can be reduced by the polled-mode I/O completion method and hybrid polling, compared to the interrupt-driven I/O completion. In this evaluation, we tested on ULL SSD with the block sizes varying from 4KB to 32KB. One can observe from this figure that compared to the interrupts, the hybrid polling only can reduce the I/O latency 8.2%, at most, while the polled-mode completion can reduce the I/O latency as high as 33%, on average. Importantly, the reason why the hybrid polling method exhibits 5% longer latency, compared to the polled-mode completion, is that the expected time to sleep before jumping to its poll process is highly inaccurate. Even though the hybrid polling calculates the sleep epoch by monitoring the average latency of each request, the device-level latency of SSDs varies due to many different parameters such as internal DRAM caching, different I/O operation types, garbage collection, and execution delays of error correction codes. Thus, the hybrid polling can oversleep or wake up the poll process early, which in turn can increase the latency or reduce benefits of the sleep. Figure 12 analyzes the CPU utilization that the hybrid polling exhibits. While it reduces the CPU cycles, the hybrid polling still consumes CPU cycles to handle sequential and random I/O requests by 58% and 56%, on average, respectively; in other words, the hybrid polling

still consumes many CPU cycles, which are 2.2× more than the conventional interrupt-driven method consumes.
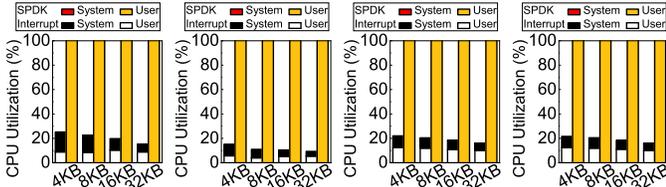
## VI. ADVANCED STORAGE STACK ANALYSIS

Emerging kernel bypass methods are already used in many computing domains, especially network (DPDK) and storage (SPDK). These methods are expected to reduce latency by removing all kernel module interventions while serving the target requests.

### A. Can SPDK really eliminate the latency overheads, imposed by the storage stack? What about ULL SSDs?
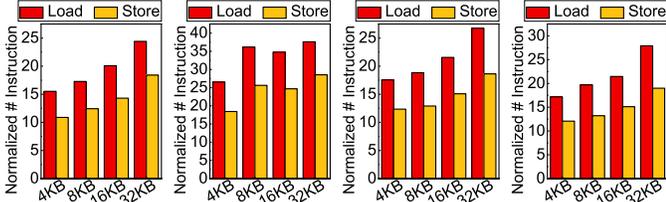
Figure 17 compares NVMe SSD's latency, exhibited by SPDK and a conventional interrupt-based method that requires going through the entire storage stack. In contrast to the common expectations of SPDK [43, 76], the latency difference between reads and writes of NVMe SSD are only 4.3% and 11.1%, on average, respectively, which are almost similar to each other and negligible. We agree that the CPU burst used for a complicated software stack can be a performance bottleneck, but one of the main roles that operating systems (OSs) need to perform is to schedule the CPU burst and I/O burst by nicely overlapping them. While the current NVMe SSDs override PCIe and provide high performance by having the internal DRAM cache, the latency of underlying flash media is not short enough yet to reduce the latency more than that of host's CPU burst.

Figure 18 compares the latency of ULL SSD with the same scenario tested for NVMe SSD. The latency using SPDK can reduce by 25.2%, 6.3%, 13.7% and 13.3% for sequential reads, random reads, sequential writes, and random writes, on average, respectively. We can learn from this evaluation that current storage system can benefit from removing the complicated software stack from the I/O path only when the latency becomes shorter than that of most conventional flash exhibits.

(a) Seq. reads.    (b) Rnd. reads.    (c) Seq. writes.    (d) Rnd. writes.

Fig. 20: CPU utilization in SPDK.



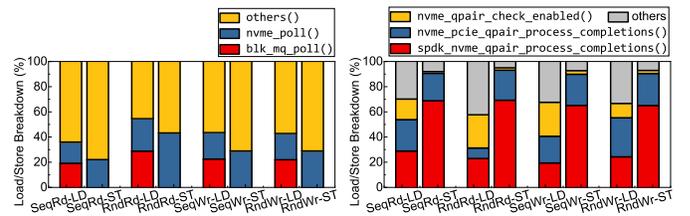(a) Seq. reads.    (b) Rnd. reads.    (c) Seq. writes.    (d) Rnd. writes.

Fig. 21: Normalized memory instruction count of SPDK.



(a) Polling.           (b) SPDK.

Fig. 22: Memory instruction breakdown.

*B. Will SPDK be the best option for future low-latency storage? Would it have any side effects?*

While SPDK can reduce the latency on ULL SSDs, we observe several issues that system designers or application users need to pay attention. First, the benefits of SPDK with ULL SSD become negligible when the request size increases. Figure 19 shows the latency comparison with SPDK and the conventional system under the execution of I/O requests whose block size is greater than 64KB (64KB~1MB). One can observe from the figure, the latency trend of SPDK is completely overlapped with that of the conventional system. Based on the results, SPDK is only meaningful with ULL, and even with ULL, users may selectively use SPDK in particular for small-sized I/O requests, as SPDK itself exhibits poor efficiency of host system resource management.

Figure 20 studies the CPU utilization of SPDK, which decomposes it into the CPU usages of userland and kernel spaces by comparing SPDK with the conventional system. Even though SPDK removes all kernel module involvements from its datapath, SPDK consumes all the CPU resources while the conventional system only uses, on average, 10% and 15% for userland and kernel executions, respectively. These excessive CPU usages are coming from polling process of the user-level NVMe driver. Since SPDK takes the NVMe driver from the kernel and puts it into `uio`, SPDK cannot handle ISR and therefore it cannot use interrupts as the I/O completion method. Thus, SPDK wastes CPU cycles to poll the huge page, which is mapped to BARs, and it increases power consumption. As multiple user applications need to communicate with `uio`, 100% CPU occupancies (took by SPDK) may be insufficient to process data from the userland directly (as there is no room for such applications).

Similar to the overheads, imposed by polling-based I/O completion methods, SPDK also exhibits heavy memory accesses. Figure 21 shows the number of load/store instructions for SPDK, which are normalized to the ones for the conventional system. SPDK generates load and store instructions more than those of the conventional system by 23× and 16.2×, on average, respectively. Note that one of the reasons why the number of loads/stores with random reads is bigger than that
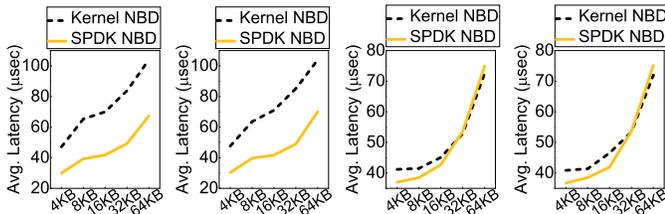
of sequential reads or writes is that random read latency of ULL SSD is relatively slower than other patterns; the longer time it waits for I/O completion, the more polling is executed.

Figures 22a and 22b analyze the contributor functions that generate most load/store instructions for the polled mode and SPDK, respectively, while we perform the previous memory access evaluations. In the conventional system, its polling takes 39% of the total load/store instructions whereas SPDK's NVMe driver generates 45% and 91% of total load and store instructions, on average, respectively. Polling of the conventional system spends most of the execution time at `blk_mq_poll` and `nvme_poll`, as discussed in Section V-B1. Surprisingly, even though SPDK also uses the polling as it cannot handle ISR, the number of load and store instructions is greater than that of the kernel side polled mode I/O completion by dozens of times (cf. Figure 15). Specifically, `spdk_nvme_qpair_process_completions()` and `nvme_pcie_qpair_process_completions()` generates 37% and 22% of total memory instructions (loads/stores), on average, respectively. Note that the role of `spdk_nvme_qpair_process_completions()` and `nvme_pcie_qpair_process_completions()` are same with those of `blk_mq_poll` and `nvme_poll`, respectively. One of the reasons behind two functions of SPDK generate more memory instructions is that SPDK's `uio` does not need to traverse the software and hardware queues (managed by `blk-mq`) to check up the CQ entries, and therefore more frequent memory accesses can be introduced. Interestingly, in contrast to kernel-side polling, SPDK uses `nvme_qpair_check_enabled()`, which is implemented by an inline function. This inline function checks if the queue pairs of CQ and SQ are valid whenever an I/O request is issued. This in turn introduces more loads, which account for 20% of the total load instructions. `nvme_qpair_check_enabled()` is used for SPDK to avoid memory accesses during a controller reset, which is triggered under a certain conditions, such as PCIe reset [4]. Note that a hybrid polling may be applied to the future storage system in an attempt to reduce CPU and memory access overheads. However, it is in the practical difficulty to be applied since the hybrid polling only supports synchronous operations, while SPDK supports both synchronous and asynchronous operations.

*C. How much can SPDK reduce the system-level latency under a real server execution environment?*

Many prior studies [8, 12, 56, 71, 75, 77] use SPDK to directly access the underlying SSD by bypassing all kernel modules in the storage stack. Even though the kernel-bypass is a very promising and attractive approach, it is limited to apply

(a) Seq. reads.  (b) Rnd. reads.  (c) Seq. writes.  (d) Rnd. writes.

Fig. 23: Analysis of SPDK for a server-client scenarios.

the kernel-bypass to a system that does not have a file system or ACID support. However, in the server-client model, a file system should be employed by the client side node, which cannot be bypassed. We configured a server-client system to evaluate the performance improvement of SPDK when the file system exists at the client-side. Note that, in Section VI-A, we demonstrated that the current high-end NVMe SSD has almost no benefit of using SPDK, and therefore, we only test this server-client model with ULL SSD. To be precise, we configure two network block devices: *Kernel NBD* and *SPDK NBD*. Kernel NBD uses a conventional network block device implemented in Linux kernel, while SPDK NBD is a server-client system that enables ULL SSD as a network block device by employing SPDK and DPDK. In this evaluation, both kernel NBD and SPDK NBD employ ext4 as their file system. ACID of I/O services for both kernel NBD and SPDK NBD is always guaranteed since every I/O request heading to the server side SSD passes through the file system. In other words, client side kernel (i.e., storage stack) cannot be bypassed even in case of SPDK NBD; SPDK and DPDK only bypass the server-side kernel modules. Figure 23 compares the latency brought by those two kernel NBD and SPDK NBD when FIO accesses a set of files, connected to the network block devices. In this test, we build 10 million files whose size varies, ranging from 4KB to 64KB. For read services, SPDK NBD can reduce the read latency, compared to kernel NBD by 39% (sequential) and 38% (random), on average, even though the execution time of client-side file system is included in its I/O path. However, for writes, SPDK cannot reduce the latency as much as what we observed in Figure 18. Specifically, the latency of sequential and random writes exhibited by SPDK are respectively 3.7% and 4.6% shorter than those of kernel NBD, on average. We believe that this is because the writes need to handle more file system level operations than the reads. The writes should perform creating or modifying multiple metadata (i.e., inodes and bitmaps) and manage system journaling to guarantee the system consistency, while the reads only exhibit minor changes on the metadata (i,e., access time). These file system level operations in writes in turn can unfortunately introduce significant kernel module involvement in accessing ULL, which lose the benefits of SPDK NBD on the writes.

## VII. RELATED WORK

While there is no new flash archive (i.e., ULL SSD) publicly available in the market, it is essential to revamp the current storage stack for such future storage. [74] pointed out that the existing interrupt-driven method cannot make systems utilize the time that CPUs are released since the latency of NVMe devices is shorter than that of the existing

storage. The study proposed to employ a polling-based synchronous I/O service in high performance systems. Similarly, [14, 23, 43, 49, 54, 68, 76, 79] advocate to employ the polled-mode system for fast NVMe devices. To use the polling in an efficient manner, [23, 49, 79] proposed different types of a hybrid polling, which combines the current interrupt-driven method with the polling by being aware of application characteristics or runtime environment. While these studies are mostly performed by simulation or DRAM emulation, of which the actual datapath exists on the memory controller side rather than block storage controller side. In this work, we performed an in-depth analysis of the entire storage stack, including NVMe driver and `blk-mq`, and revealed several system-level characteristics and implications, which are not observed by the previous work.

On the other hand, the kernel-bypass and user-level I/O management are also popularly studied in both academia and industry. Specifically, [16, 43, 66] consider using user-space I/O framework to directly access the underlying NVMe devices and employing a dedicated polling thread to reduce the system-level latency, which are similar to the industry-driven I/O frameworks such as SPDK [71, 76], OpenMPDK [60] and uNVMe [44, 57, 60]. As an application of this advanced storage stack, [76] leverages SPDK to reduce the latency, imposed by multiple virtual machines; this study uses a userland driver to connect hypervisor to the underlying SSD and distributes queue pairs across individual virtual machine instances. As this work pointed out, the current NVMe device is not fast enough to take the full benefits of employing SPDK. However, in the near future, SPDK-applied systems can reduce the overheads imposed by software kernel involvement from their I/O path. To maximize the advantage of SPDK, we also provide several in-depth studies such as file system overheads, which should be reconsidered in a real server-client usage.

## VIII. CONCLUSION

We analyzed the performance behaviors of ULL SSDs and brought several system-level challenges of the current and future storage stack to take full advantages of ULL devices. Specifically, we observed that ULL SSDs can offer their maximum bandwidth with only a few queue entries, which contradicts with the design direction of the current rich NVMe queue. While it is beneficial to employ a polling-based I/O completion routine for ULL SSDs, system-level overheads delivered by polling, such as high CPU cycles and frequent memory accesses, incur frequent CPU stalls and increase overall system power consumption.

REFERENCES

[1] blk: optimization for classic polling. https://lore.kernel.org/patchwork/patch/885868/.

[2] Direct memory access (dma) from user space. https://spdk.io/doc/memory.html.

[3] Intel vtune amplifier, 2019. DocumentationattheURL:https://software.intel.com/en-us/intel-vtune-amplifier-xe.

[4] NVM Express. http://nvmexpress.org/wp-content/uploads/NVM-Express-1_3a-20171024_ratified.pdf.

[5] Storage performance development kit. https://spdk.io/development/.

[6] Data plane development kit. https://www.dpdk.org/, 2014.

[7] Nvm express revision 1.2.1. https://nvmexpress.org/wp-content/uploads/NVM_Express_1_2_1_Gold_20160603.pdf, 2016.

[8] AN, Z., ZHANG, Z., LI, Q., XING, J., DU, H., WANG, Z., HUO, Z., AND MA, J. Optimizing the datapath for key-value middleware with nvme ssds over rdma interconnects. In *2017 IEEE International Conference on Cluster Computing (CLUSTER)* (2017), IEEE, pp. 582–586.

[9] AXBOE, J. Flexible io tester. *https://github.com/axboe/fio*.

[10] BHATTACHARYA, S. Linux asynchronous i/o design: Evolution & challenges.

[11] BJØRLING, M., AXBOE, J., NELLANS, D., AND BONNET, P. Linux block io: introducing multi-queue ssd access on multi-core systems. In *Proceedings of the 6th international systems and storage conference* (2013), ACM, p. 22.

[12] CAO, W., LIU, Z., WANG, P., CHEN, S., ZHU, C., ZHENG, S., WANG, Y., AND MA, G. Polarfs: an ultra-low latency and failure resilient distributed file system for shared storage cloud database. *Proceedings of the VLDB Endowment 11*, 12 (2018), 1849–1862.

[13] CAULFIELD, A. M., COBURN, J., MOLLOV, T., DE, A., AKEL, A., HE, J., JAGATHEESAN, A., GUPTA, R. K., SNAVELY, A., AND SWANSON, S. Understanding the impact of emerging non-volatile memories on high-performance, io-intensive computing. In *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis* (2010), IEEE Computer Society, pp. 1–11.

[14] CAULFIELD, A. M., DE, A., COBURN, J., MOLLOW, T. I., GUPTA, R. K., AND SWANSON, S. Moneta: A high-performance storage array architecture for next-generation, non-volatile memories. In *Proceedings of the 2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture* (2010), IEEE Computer Society, pp. 385–395.

[15] CAULFIELD, A. M., GRUPP, L. M., AND SWANSON, S. Gordon: using flash memory to build fast, power-efficient clusters for data-intensive applications. *ACM Sigplan Notices 44*, 3 (2009), 217–228.

[16] CAULFIELD, A. M., MOLLOV, T. I., EISNER, L. A., DE, A., COBURN, J., AND SWANSON, S. Providing safe, user space access to fast, solid state disks. *ACM SIGARCH Computer Architecture News 40*, 1 (2012), 387–400.

[17] CHANG, Y.-M., CHANG, Y.-H., KUO, T.-W., LI, Y.-C., AND LI, H.-P. Achieving slc performance with mlc flash memory. In *Proceedings of the 52nd Annual Design Automation Conference* (2015), ACM, p. 192.

[18] CHEONG, W., YOON, C., WOO, S., HAN, K., KIM, D., LEE, C., CHOI, Y., KIM, S., KANG, D., YU, G., KIM, J., PARK, J., SONG, K.-W., PARK, K.-T., CHO, S., OH, H., LEE, D. D., CHOI, J.-H., AND JEONG, J. A flash memory controller for $15\mu s$ ultra-low-latency ssd using high-speed 3d nand flash with $3\mu s$ read time. In *Solid-State Circuits Conference-(ISSCC), 2018 IEEE International* (2018), IEEE, pp. 338–340.

[19] CHI, P., LI, S., XU, C., ZHANG, T., ZHAO, J., LIU, Y., WANG, Y., AND XIE, Y. Prime: A novel processing-in-memory architecture for neural network computation in reram-based main memory. In *ACM SIGARCH Computer Architecture News* (2016), vol. 44, IEEE Press, pp. 27–39.

[20] CHOI, W., JUNG, M., AND KANDEMIR, M. Invalid data-aware coding to enhance the read performance of high-density flash memories. In *2018 51st Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)* (2018), IEEE, pp. 482–493.

[21] CHOI, W., JUNG, M., KANDEMIR, M., AND DAS, C. Parallelizing garbage collection with i/o to improve flash resource utilization. In *Proceedings of the 27th International Symposium on High-Performance Parallel and Distributed Computing* (2018), ACM, pp. 243–254.

[22] CORBET, J. Block-layer i/o polling. https://lwn.net/Articles/663879.

[23] EISENMAN, A., GARDNER, D., ABDELRAHMAN, I., AXBOE, J., DONG, S., HAZELWOOD, K., PETERSEN, C., CIDON, A., AND KATTI, S. Reducing dram footprint with nvm in facebook. In *Proceedings of the Thirteenth EuroSys Conference* (2018), ACM, p. 42.

[24] ELLIOTT, J., AND JEONG, J. Advancements in ssds and 3d nand reshaping storage market. *Keynote presentation at Flash Memory Summit 3* (2017), 1.

[25] ESHGHI, K., AND MICHELONI, R. Ssd architecture and pci express interface. In *Inside solid state drives (SSDs)*. Springer, 2013, pp. 19–45.

[26] FADU, INC. FADU. http://www.fadu.io/, 2017.

[27] HICKEN, M. S., HOWE, S. M., SOKOLOV, D. J., SWATOSH, T., AND WILLIAMS, J. L. Disk drive cache system using a dynamic priority sequential stream of data segments continuously adapted according to prefetched sequential random, and repeating types of accesses, July 18 2000. US Patent 6,092,149.

[28] HSU, S. T., AND ZHUANG, W.-W. Electrically programmable resistance cross point memory, Mar. 11 2003. US Patent 6,531,371.

[29] HUFFMAN, A., AND CLARK, J. Serial ata native command queuing. *Intel Corporation and Seagate Technology, Whitepaper* (2003).

[30] INTEL CORPORATION. Intel SSD 750 Series. https://www.intel.com/content/www/us/en/products/memory-storage/solid-state-drives/gaming-enthusiast-ssds/750-series.html.

[31] INTEL CORPORATION. Intel SSD D3-S4610 Series. https://www.intel.com/content/www/us/en/products/memory-storage/solid-state-drives/data-center-ssds/d3-series/d3-s4610-series/d3-s4610-3-84tb-2-5inch-3d2.html, 2018.

[32] JI, Y., ZHANG, Y., XIE, X., LI, S., WANG, P., HU, X., ZHANG, Y., AND XIE, Y. Fpsa: A full system stack solution for reconfigurable reram-based nn accelerator architecture. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems* (2019), ACM, pp. 733–747.

[33] JUNG, M. Exploring design challenges in getting solid state drives closer to CPU. *IEEE Transactions on Computers 65*, 4 (2014), 1103–1115.

[34] JUNG, M. Exploring parallel data access methods in emerging non-volatile memory systems. *IEEE Transactions on Parallel and Distributed Systems 28*, 3 (2017), 746–759.

[35] JUNG, M., CHOI, W., KWON, M., SRIKANTAIAH, S., YOO, J., AND KANDEMIR, M. Design of a host interface logic for gc-free ssds. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* (2019).

[36] JUNG, M., AND KANDEMIR, M. Challenges in getting flash drives closer to CPU. In *Presented as part of the 5th USENIX Workshop on Hot Topics in Storage and File Systems* (2013).

[37] JUNG, M., AND KANDEMIR, M. Revisiting widely-held expectations of ssd and rethinking implications for systems. *SIGMETRICS13* (2013), 203–216.

[38] JUNG, M., AND KANDEMIR, M. T. An evaluation of different page allocation strategies on high-speed ssds. In *HotStorage* (2012).

[39] JUNG, M., PRABHAKAR, R., AND KANDEMIR, M. T. Taking garbage collection overheads off the critical path in SSDs. In *ACM/IFIP/USENIX International Conference on Distributed Systems Platforms and Open Distributed Processing* (2012), Springer, pp. 164–186.

[40] JUNG, M., WILSON III, E. H., AND KANDEMIR, M. Physically addressed queueing (paq): improving parallelism in solid state disks. In *ACM SIGARCH Computer Architecture News* (2012), vol. 40, IEEE Computer Society, pp. 404–415.

[41] KANG, Y., KEE, Y.-S., MILLER, E. L., AND PARK, C. Enabling cost-effective data processing with smart ssd. In *2013 IEEE 29th symposium on mass storage systems and technologies (MSST)* (2013), IEEE, pp. 1–12.

[42] KIM, C., KIM, D.-H., JEONG, W., KIM, H.-J., PARK, I. H., PARK, H.-W., LEE, J., PARK, J., AHN, Y.-L., LEE, J. Y., KIM, S.-B., YOON, H., YU, J. D., CHOI, N., KIM, N., JANG, H., PARK, J., SONG, S., PARK, Y., BANG, J., HONG, S., CHOI, Y., KIM, M.-S., KIM, H., KWAK, P., IHM, J.-D., BYEON, D. S., LEE, J.-Y., PARK, K.-T., AND KYUNG, K.-H. A 512-gb 3-b/cell 64-stacked wl 3-d-nand flash memory. *IEEE Journal of Solid-State Circuits 53*, 1 (2017), 124–133.

[43] KIM, H.-J., LEE, Y.-S., AND KIM, J.-S. Nvmedirect: A user-space i/o framework for application-specific optimization on nvme ssds. In *8th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 16)* (2016).

[44] KIM, K. KV and LBA SSD userspace NVMe driver. https://github.com/OpenMPDK/uNVMe.

[45] KOCH, H. J., AND GMB, H. L. Userspace i/o drivers in a realtime context. In *The 13th Realtime Linux Workshop* (2011).

[46] KOH, S., LEE, C., KWON, M., AND JUNG, M. Exploring system challenges of ultra-low latency solid state drives. In *10th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 18)* (2018).

[47] KOUGKAS, A., DEVARAJAN, H., AND SUN, X.-H. Hermes: a heterogeneous-aware multi-tiered distributed I/O buffering system. In *Proceedings of the 27th International Symposium on High-Performance Parallel and Distributed Computing* (2018), ACM, pp. 219–230.

[48] KUMAR, P., AND HUANG, H. H. Graphone: A data store for real-time analytics on evolving graphs. In *17th USENIX Conference on File and Storage Technologies (FAST 19)* (2019), pp. 249–263.

[49] Le Moal, D. I/o latency optimization with polling. *Vault–Linux Storage and Filesystem* (2017).

[50] Lee, S., Shin, D., Kim, Y.-J., and Kim, J. Last: locality-aware sector translation for nand flash memory-based storage systems. *ACM SIGOPS Operating Systems Review 42*, 6 (2008), 36–42.

[51] Lee, S.-W., Choi, W.-K., and Park, D.-J. Fast: An efficient flash translation layer for flash memory. In *International Conference on Embedded and Ubiquitous Computing* (2006), Springer, pp. 879–887.

[52] Li, J., Tseng, H.-W., Lin, C., Papakonstantinou, Y., and Swanson, S. Hippogriffdb: Balancing i/o and gpu bandwidth in big data analytics. *Proceedings of the VLDB Endowment 9*, 14 (2016), 1647–1658.

[53] Liu, H., and Huang, H. H. Graphene: Fine-grained IO management for graph computing. In *15th USENIX Conference on File and Storage Technologies (FAST 17)* (2017), pp. 285–300.

[54] Liu, J., and Abali, B. Virtualization polling engine (vpe): using dedicated cpu cores to accelerate i/o virtualization. In *Proceedings of the 23rd international conference on Supercomputing* (2009), ACM, pp. 225–234.

[55] Liu, N., Cope, J., Carns, P., Carothers, C., Ross, R., Grider, G., Crume, A., and Maltzahn, C. On the role of burst buffers in leadership-class storage systems. In *IEEE 28th Symposium on Mass Storage Systems and Technologies (MSST)* (2012), IEEE, pp. 1–11.

[56] Luse, P. Spdk blobstore: A look inside the NVM optimized allocator. https://www.snia.org/sites/default/files/SDC/2017/presentations/Solid_State_Stor_NVM_PM_NVDIMM/Luse_Paul_Verma_Vishal_SPDK_Blobstore_A_Look_Inside_the_NVM_Optimized_Allocator.pdf.

[57] Micron. User space nvme driver. https://github.com/MicronSSD/unvme.

[58] Oh, Y., Lee, E., Hyun, C., Choi, J., Lee, D., and Noh, S. H. Enabling cost-effective flash based caching with an array of commodity SSDs. In *Proceedings of the 16th Annual Middleware Conference* (2015), ACM, pp. 63–74.

[59] Samsung. Ultra-Low Latency with Samsung Z-NAND SSD. https://www.samsung.com/semiconductor/global.semi.static/Ultra-Low_Latency_with_Samsung_Z-NAND_SSD-0.pdf, 2017.

[60] Samsung Electronics Co., L. Openmpdk and uNVMe user space device driver for server and data center. https://www.samsung.com/semiconductor/global.semi.static/White_Paper_OpenMPDK_1807.pdf.

[61] Samsung Electronics Co., L. PM1725 NVMe PCIe SSD. http://www.samsung.com/semiconductor/insights/tech-leadership/pm1725-nvme-pcie-ssd, 2015.

[62] Serial, A. International organization. *Serial ATA Revision 3* (2007).

[63] Seshadri, S., Gahagan, M., Bhaskaran, S., Bunker, T., De, A., Jin, Y., Liu, Y., and Swanson, S. Willow: A user-programmable SSD. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)* (2014), pp. 67–80.

[64] Shahidi, N., Arjomand, M., Jung, M., Kandemir, M. T., Das, C. R., and Sivasubramaniam, A. Exploring the potentials of parallel garbage collection in ssds for enterprise storage systems. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* (2016), IEEE Press, p. 48.

[65] Shen, K., Park, S., and Zhu, M. Journaling of journal is (almost) free. In *Proceedings of the 12th USENIX Conference on File and Storage Technologies ({FAST} 14)* (2014), pp. 287–293.

[66] Shin, W., Chen, Q., Oh, M., Eom, H., and Yeom, H. Y. OS I/O path optimizations for flash solid-state drives. In *2014 USENIX Annual Technical Conference (USENIXATC 14)* (2014), pp. 483–488.

[67] Suh, K.-D., Suh, B.-H., Lim, Y.-H., Kim, J.-K., Choi, Y.-J., Koh, Y.-N., Lee, S.-S., Kwon, S.-C., Choi, B.-S., Yum, J.-S., et al. A 3.3 v 32 mb nand flash memory with incremental step pulse programming scheme. *IEEE Journal of Solid-State Circuits 30*, 11 (1995), 1149–1156.

[68] Swanson, S., and Caulfield, A. M. Refactor, reduce, recycle: Restructuring the i/o stack for the future of storage. *Computer 46*, 8 (2013), 52–59.

[69] Tallis, B. Samsung launches Z-SSD SZ985: Up to 800gb of Z-NAND. https://www.anandtech.com/show/12376/samsung-launches-zssd-sz985-up-to-800gb-of-znand.

[70] Tseng, H.-W., Zhao, Q., Zhou, Y., Gahagan, M., and Swanson, S. Morpheus: Exploring the potential of near-data processing for creating application objects in heterogeneous computing. *ACM SIGOPS Operating Systems Review 51*, 1 (2018), 71–83.

[71] Walker, B. Spdk: Building blocks for scalable, high performance storage applications. In *Storage Developer Conference. SNIA* (2016).

[72] Wu, G., and He, X. Reducing ssd read latency via nand flash program and erase suspension. In *FAST* (2012), vol. 2, p. 3.

[73] Yamashita, R., Magia, S., Higuchi, T., Yoneya, K., Yamamura, T., Mizukoshi, H., Zaitsu, S., Yamashita, M., Toyama, S., Kamae, N., Lee, J., Chen, S., Tao, J., Mak, W., Zhang, X., Yu, Y., Utsunomiya, Y., Kato, Y., Sakai, M., Matsumoto, M., Chibvongodze, H., Ookuma, N., Yabe, H., Taigor, S., Samineni, R., Kodama, T., Kamata, Y., Namai, Y., Huynh, J., Wang, S.-E., He, Y., Pham, T., Saraf, V., Petkar, A., Watanabe, M., Hayashi, K., Swarnkar, P., Miwa, H., Pradhan, A., Dey, S., Dwibedy, D., Xavier, T., Balaga, M., Agarwal, S., Kulkarni, S., Papasaheb, Z., Deora, S., Hong, P., Wei, M., Balakrishnan, G., Ariki, T., Verma, K., Siau, C., Dong, Y., Lu, C.-H., Miwa, T., and Moogat, F. 11.1 a 512gb 3b/cell flash memory on 64-word-line-layer bics technology. In *2017 IEEE International Solid-State Circuits Conference (ISSCC)* (2017), IEEE, pp. 196–197.

[74] Yang, J., Minturn, D. B., and Hady, F. When poll is better than interrupt. In *FAST* (2012), vol. 12, pp. 3–3.

[75] Yang, Z. Accelerate block service built on ceph via SPDK. https://www.snia.org/sites/default/files/SDC/2017/presentations/Storage_Architecture/Yang_Zie_Accelerate_block_service_built_on_Ceph_via_SPDK.pdf.

[76] Yang, Z., Harris, J. R., Walker, B., Verkamp, D., Liu, C., Chang, C., Cao, G., Stern, J., Verma, V., and Paul, L. E. Spdk: A development kit to build high performance storage applications. In *2017 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)* (2017), IEEE, pp. 154–161.

[77] Yang, Z., Liu, C., Zhou, Y., Liu, X., and Cao, G. SPDK vhost-NVMe: Accelerating I/Os in virtual machines on NVMe SSDs via user space vhost target. In *2018 IEEE 8th International Symposium on Cloud and Service Computing (SC2)* (2018), IEEE, pp. 67–76.

[78] Zhang, H., Xiao, N., Liu, F., and Chen, Z. Leader: Accelerating reram-based main memory by leveraging access latency discrepancy in crossbar arrays. In *Proceedings of the 2016 Conference on Design, Automation & Test in Europe* (2016), EDA Consortium, pp. 756–761.

[79] Zhang, J., Kwon, M., Gouk, D., Koh, S., Lee, C., Alian, M., Chun, M., Kandemir, M. T., Kim, N. S., Kim, J., et al. Flashshare: punching through server storage stack from kernel to firmware for ultra-low latency ssds. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)* (2018), pp. 477–492.

[80] Zhang, Y., Feng, D., Liu, J., Tong, W., Wu, B., and Fang, C. A novel reram-based main memory structure for optimizing access latency and reliability. In *2017 54th ACM/EDAC/IEEE Design Automation Conference (DAC)* (2017), IEEE, pp. 1–6.