

H-NVMe: A Hybrid Framework of NVMe-based Storage System in Cloud Computing Environment

Zhengyu Yang^{*}, Morteza Hoseinzadeh[‡], Ping Wong[†], John Artoux[†], Clay Mayers[†],
David (Thomas) Evans[†], Rory (Thomas) Bolt[†], Janki Bhimani^{*}, Ningfang Mi^{*}, and Steven Swanson[‡]

^{*}Dept. of Electrical and Computer Engineering, Northeastern University, Boston, MA 02115

[‡] Dept. of Computer Science and Engineering, University of California San Diego, San Diego, CA 92093

[†] Samsung Semiconductor Inc., Memory Solution Research Lab, Software Group, San Diego, CA 92121

Abstract—In the year of 2017, more and more datacenters have started to replace traditional SATA and SAS SSDs with NVMe SSDs due to NVMe’s outstanding performance [1]. However, for historical reasons, current popular deployments of NVMe in VM-hypervisor-based platforms (such as VMware ESXi [2]) have numbers of intermediate queues along the I/O stack. As a result, performance is bottlenecked by synchronization locks in these queues, cross-VM interference induces I/O latency, and most importantly, up-to-64K-queue capability of NVMe SSDs cannot be fully utilized. In this paper, we developed a hybrid framework of NVMe-based storage system called “H-NVMe”, which provides two VM I/O stack deployment modes “Parallel Queue Mode” and “Direct Access Mode”. The first mode increases parallelism and enables lock-free operations by implementing local lightweight queues in the NVMe driver. The second mode further bypasses the entire I/O stack in the hypervisor layer and allows trusted user applications whose hosting VMDK (Virtual Machine Disk) files are attached with our customized vSphere IOFilters [3] to directly access NVMe SSDs to improve the performance isolation. This suits premium users who have higher priorities and the permission to attach IOFilter to their VMDKs. H-NVMe is implemented on VMware ESXi 6.0.0, and our evaluation results show that the proposed H-NVMe framework can significantly improve throughputs and bandwidths compared to the original inbox NVMe solution.

Index Terms—NVMe Driver, Lock-free I/O Queue, User Mode Polling, Datacenter Storage, Resource Management, Tiering and Caching Algorithm, Big Data, Cloud Computing, Virtualization

I. INTRODUCTION

Ever decreasing price per gigabyte of SSDs (Solid-state Drive) and their capability of fast operation make them indispensable for large-scale cloud services. They are overtaking HDDs (Hard Disk Drive) and leaving them far behind by providing orders of magnitude more IOPS and lowering I/O latency. Thanks to the rapid improvement of the process technology, SSDs can possess several chips. Consequently, they may boost up their bandwidth and capacity by obtaining higher degrees of parallelism. With this regards, Non-Volatile Memory Express (NVMe) [1], [4] interface has been introduced to better utilize the parallelism of the emerging storage technologies.

This work was completed during Zhengyu Yang and Morteza Hoseinzadeh’s internship at Samsung Semiconductor Inc. This project is partially supported by NSF grant CNS-1452751.

However, the limitation induced by software overheads prevents users to perfectly perceive this performance advancement. The overhead of the legacy kernel I/O stack, which has been optimized for slow HDDs, is more noticeable as the storage devices and the connection interfaces get faster. There is a tremendous number of researches trying to reduce the kernel overhead by eliminating unnecessary context processing [5], employing a polling mechanism instead of interrupts [6], [7], and performance isolation [8]–[10].

In addition to the new emerging storage technologies which provide high performance operation through parallelism, container-based virtualization [11] has been a key cloud computing platform which can perfectly take advantages of this parallelism through allowing multiple isolated instances (i.e., containers) of the storage resources. Because of the independence of containers running on top of a single host operating system (i.e., hypervisor), the management of resources through isolation and sharing becomes more important in container-based virtualization. Obviously, having advanced guest kernels is not enough for efficiently sharing system resources. VMware ESXi [2] is one of the most common commercial virtualization platform which can potentially use the NVMe SSDs in the most efficient way by leveraging the massive parallelism and isolation characteristics of NVMe SSDs. Accordingly, VMware ESXi provides VMDK (Virtual Machine Disk) files for each VM, which are completely independent of each other. But, in the case of using NVMe SSDs, they access the storage resource through a single submission and completion queue in NVMe driver, regardless of high levels of parallelism provided by NVMe. This inefficiency originates from the NVMe driver in the hypervisor, and has become a bottleneck in the storage I/O stack.

In this paper, we propose H-NVMe, a novel NVMe framework on VMware ESXi. To best utilize NVMe SSDs, H-NVMe provides two different working modes: “*Parallel Queue Mode*” and “*Direct Access Mode*”. In the former working mode, H-NVMe circumvents the built-in *Adapter Queue* of ESXi by emptying it and spreading out its entities between multiple lightweight subqueues in our customized NVMe driver in order to use the parallelism of the device more efficiently. The latter working mode bypasses all the hypervisor queues and directly connects the trusted user application

threads to the NVMe *Driver Queue*, to achieve better performance isolation. H-NVMe can work in either of these two modes in whole or in partial (details see in Sec. III). It means that both parallelism and isolation can be provided at the same time. We evaluate the performance of H-NVMe with a set of representative applications, and the experimental results show that H-NVMe can significantly improve the I/O performance.

The rest of this paper is organized as follows. Sec. II presents the background and motivation. Sec. III proposes our hybrid framework algorithm. Experimental evaluation results and analysis are presented in Sec. IV. Sec. V discusses the related work. We finally present the conclusions in Sec. VI.

II. BACKGROUND AND MOTIVATION

A. Hardware Development

Hardware technology development trends can be categorized into two major paths: storage technology and hardware I/O interface.

1) *Storage Technology*: The storage technology is the means of maintaining digital data in form of binary information, and plays a key role in the latency of a device. In HDD devices, data is stored in form of sequential changes in the direction of magnetization in a thin film of ferromagnetic material on a disk. Reading and writing to such a device has mechanical nature of the rotating disks and moving heads [12]. Later on, flash memory cells came into play, and a new trend of storage technology formed by introducing Flash-based SSDs. Since flash memory cells need to be erased before storing new data, a Flash Translation Layer (FTL) [13] makes SSDs to present HDD-like I/O interface to an operating system. Even though the latency of an SSD varies under different workloads because of FTL, it is orders of magnitude faster than HDDs. As a result, the new trend in developing SSD devices is rapidly going on to provide cheap storage devices with large storage capacity. There have been a large number of studies on other Storage Class Memory (SCM) including Phase-Change Memory (PCM), Resistive RAM, and 3D XPoint [14], which use these storage devices either as primary storage in form of Non-Volatile Dual In-line Memory Module (NVDIMM) [15], [16], or as secondary storage [17], [18] in form of SSDs.

The ultimate goal of developing SCM is to provide near-DRAM latency and also to maintain data persistent. Up to 2017, there is no matured SCM in the market other than flash-based SSDs which will to replace HDDs very soon. Modern SSDs tend to exploit the parallelism of multiple flash chip, and reorder I/O requests for better scheduling [19]–[23], in order to amortize per-request latency overhead and achieve high bandwidth by performing concurrent data transfers.

2) *Hardware I/O Interface*: With the emergence of new and fast storage technologies, hardware I/O bus interface also has experienced significant upgrades to facilitate the use of these modern technologies. Parallel ATA interface (i.e., IDE) has become outdated as it was standardized for slow HDDs with a few megabytes of data transfer per second. Recent fast storage devices usually support Serial ATA (SATA) with a few gigabytes per second transfer rates. However, new storage

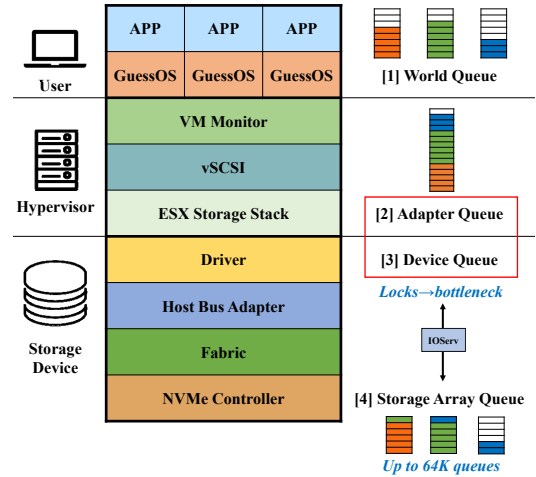


Fig. 1: Architecture of current NVMe-based VM Hypervisor

devices demand much higher transfer rates. The maximum I/O bandwidth is determined by the I/O bus interface which is migrating from IDE and SATA toward PCI Express and is expected to support up to 128 GB/s by 2019.

B. Software Development

Despite the mature infrastructures of recent hardware technology for using fast storage devices, the development of kernel storage stack still needs to pay more attention to these trends.

1) *NVMe Driver*: In the past, SATA, SAS or Fibre Channel buses were most common interfaces for SSDs. According to the enterprise storage system marketplace, SATA is the most typical way for interfacing SSDs and OS. Also, there was a few number of PCI Express based bus interfaces for high-end SSDs, but they use non-standard specification interfaces. After standardizing a universal interface of SSDs, the OS may only need one driver to interact with all SSDs. It is no longer needed to use additional resources to develop specific interface drivers. Non-Volatile Memory Express (NVMe) [1] is a scalable host controller interface designed for both enterprise and client systems to use SSDs over PCI Express. Recent NVM-Express standard [24] abridges the I/O path with several deeper queues. It provides a device driver which bypasses the block layer and the SCSI I/O subsystem. This driver directly issues the requests to a deeper hardware queue (up to 64K in depth) which enervates the need of background queue running context in the I/O completion path. Therefore, host hardware and software can fully exploit the highest levels of parallelism offered by modern SSDs which form a single unit that plugs into the PCIe bus.

2) *I/O Stack in Virtualization Environments*: Storage I/O stack in virtualization environments requires low latency. The main challenges are from (1) the presence of additional software layer such as guest OS; (2) context switching between VM and hypervisor; and (3) queuing delay for I/O operations. These challenges do not cause serious problems and lead to high latency when HDDs are dominating the datacenter [25].

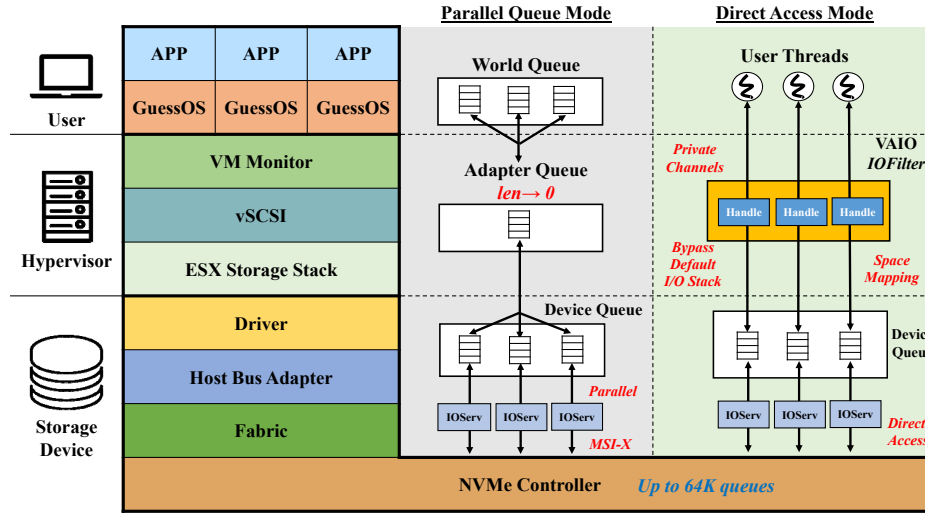


Fig. 2: I/O paths of two proposed modes.

However, in the year of 2017, more and more cloud service vendors started to adopt NVMe SSDs into their storage systems, and for some historical reasons, the current popular deployment of NVMe SSDs in cloud computing hypervisor cannot best utilize the NVMe. For example, Fig. 1 depicts a simplified architecture of NVMe I/O stack in VMware’s ESXi hypervisor solution (i.e., VMware ESXi Pluggable Storage Architecture). We can see that there are multiple intermediate queues in the I/O stack. In detail, when an I/O comes to the current original NVMe framework the following steps are needed:

- 1) *World Queue* is responsible for each VM I/O requests.
- 2) *Adapter Queue* in the hypervisor layer gathers jobs from *World Queues*.
- 3) Hypervisor translates the request to `ScsiCommand()` [26] and sends it to the *Device Queue* in the driver.
- 4) *IOServ* (I/O server) in the Driver acquires an internal completion *lock*.
- 5) *Driver* sends I/O to the device, waiting for the device to finish the I/O request, while holding the *lock*.
- 6) Once the I/O is done, *Driver* releases the *lock* and completes the I/O asynchronously when NVMe controller I/O completion interrupt occurs.

This architecture has many issues. Specifically, VMware’s PSA layer virtually completes all I/O asynchronously via interrupts. Historically, interrupt driven I/O is efficient for high latency devices such as hard disks, but it also induces a lot of overheads including many context switches. However, it is not efficient for NVMe SSDs. Most importantly, this interrupt driven approach is not lock-free and limits the parallel I/O capacity, since *IOServs* in the upper layer will not send another I/O to the NVMe controller until previous submitted I/O returns pending status. Thus, the NVMe’s multiple cores/multiple queues mechanisms are not fully utilized.

To further find the evidence of the disadvantage of having these non-lock-free queues, we conduct a preliminary test where we break down the five locks initialized in the sub-

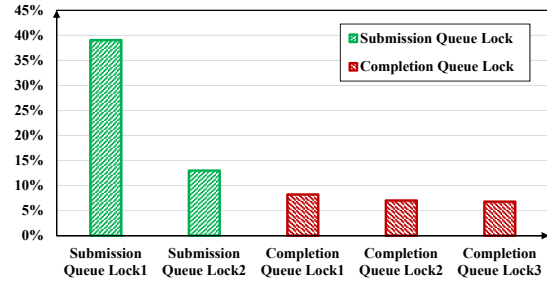


Fig. 3: Normalized time spent of locks in each queue.

mission and completion queues in the *Adapter Queue*. The statistic results in Fig. 3 show that temporal wastes on these locks are not negligible, especially the *lock1* in submission queue [27] takes near 40% of the total queuing time. As a result, the huge bottleneck between the *Adapter Queue* and the *Device Queue* makes the system not able to best utilize up to 64K queues in the NVMe controller, in addition to the cross-VM interference issues.

In summary, these issues motivate us to develop a novel solution that is highly parallelism with lock-free queues and supports performance isolation cross VMs.

III. ALGORITHM DESIGN

In this section, we propose a hybrid NVMe utilization approach called “H-NVMe” to support parallelism and performance isolation by introducing two modes, namely “*Parallel Queue Mode (PQM)*” and “*Direct Access Mode (DAM)*”. The PQM is to enforce the *Adapter Queue* to be empty and use our enhanced subqueues in the driver to speed up the I/O stack. This is relatively straightforward to implement and is a more generalized solution. On the other hand, the DAM can achieve better performance by directly connecting those trusted user application threads to the NVMe driver queues through our customized VAIO (vSphere APIs for IOFiltering) IOFilters [3] attached to their corresponding hosting VMDKs. In contrast, this approach breaks the encapsulation of NVMe

resources and thus needs higher permission control processes. Based on SLA (Service Level Agreement) of VMDKs owned by different users, and the security audit requirement (such as attaching IOFilter permission), the cloud manager can select from these two modes for VM deployment in the NVMe-based storage system.

A. Parallel Queue Mode

The challenge of solving the major bottleneck located in the hypervisor layer (observed in Sec. II-B2) is that unlike the “*Device Queue*” encapsulated in the driver layer which we have full control, the “*World Queue*” and “*Adapter Queue*” are much engineering-unfriendly and any changes on them will affect the system capability. Thus, our solution is to “force” the *Adapter Queue* to be emptied and forwarded to the driver layer (through some “*fake completion*” signals from the driver), where the enhanced subqueues can help to increase the parallelism and speed up the I/Os, see in the middle column of Fig. 2 and Alg. 1.

Unlike the original NVMe design where the single *IOServ* has the “lock” issue, H-NVMe allows to create more than one *IOServ* threads and subqueues in the “*Device Queue*” of the NVMe driver to handle the I/O with the NVMe controller. Therefore, the lock-free goal can be achieved by *IOServ* threads to only focusing on their own queues, see in Alg. 1 lines 11 to 13. In our implementation, the *callIOServ* function selects an idle *IOServ* to assign jobs in the round robin order.

Apparently, having more *IOServs* will indeed improve performance, but it is not free to have infinite *IOServs*, since more CPU and memory resources will be consumed by these threads. As a result, the service rate of each *IOServ* decreases, too. Therefore, the next question that PQM has to address is “*how to dynamically assign subqueue and IOServs number?*”, which motivates us to find an adaptive algorithm to automatically selecting a appropriate number of *IOServs*.

Since the the length of *Adapter Queue* is forced to be 0, we can model the problem based on the *M/M/c* queuing theory [28]. Let the total arrival rate (from the *Adapter Queue*) be λ , and let vector $\vec{\mu}$ denote each *IOServ*’s service rate. This vector has c dimensions, where $c \in [1, c_{max}]$ is the number of *IOServs*, and c_{max} is the preset maximal *IOServ* number. If we increase the number of c (e.g., creating and destroying *IOServs*), each server’s service rate will change, and this change can be estimated by a regression function (see in Alg. 1 line 24) based on periodically measurement. Once λ and $\vec{\mu}$ are calculated, H-NVMe calls the *optimizeServNum* function to decide the best number of c . The optimization objective of this function is to minimize the total latency. Specifically, different combinations of c and corresponding service rates $\vec{\mu}_c$ are tested. This procedure also uses the *ErlangC* function [29] to calculate the probability that an arriving job will need to queue (as opposed to immediately being served). Lastly, to further reduce the cost of changing c , we need to limit the frequency of updating c to a preset update epoch window *epochLen* in Alg. 1 line 22.

Algorithm 1: Main Procedures of Parallel Queue Mode.

```

1 Procedure IOSubmissionNVMeHypervisor ()
2   while OSLibIOSubmit(ctrlr, vmkCmd, devData) do
3     submissionQueue = getSubmissionQueue(ctrlr);
4     IORequest = getIORequest(ctrlr, vmkCmd);
5     if !IORequest then
6       return VMK_NO_MEMORY;
7     IORequest.vmkCmd = vmkCmd;
8     IORequest.devData = devData;
9     OrigLink = submissionQueue.link;
10    IORequest.link = OrigLink;
11    while !update(submissionQueue.link, OrigLink,
12      IORequest) do
13      | OrigLink = submissionQueue.link;
14      | IORequest.link = OrigLink;
15      callIOServ(submissionQueue);
16    return VMK_OK;
17 Procedure callIOServ (submissionQueue)
18 | curIOServ=selectIOServ();
19 | curIOServ.takeJob(submissionQueue);
20 return;
21 Procedure subQueueNum ()
22 while True do
23 | if curTime MOD epochLen = 0 then
24 |    $\lambda = \text{updateArrvRate}()$ ;
25 |    $\vec{\mu} = \text{regressSubQueueServRate}(C_{max})$ ;
26 |    $c = \text{optimizeServNum}(\lambda, \vec{\mu})$ ;
27 Procedure optimizeServNum ( $\lambda, \vec{\mu}$ )
28 return  $\underset{c \in [1, c_{max}]}{\text{argmin}} \left[ \frac{\text{ErlangC}(\lambda, \vec{\mu}_c)}{c\vec{\mu}_c - \lambda} + \frac{1}{\vec{\mu}_c} \right]$ ;
29 Procedure ErlangC ( $\lambda, \vec{\mu}_c$ )
30  $\rho = \frac{\lambda}{\vec{\mu}_c}$ ;
return  $\frac{1}{1 + (1 - \rho) \left[ \frac{c!}{(c\rho)^c} \sum_{k=0}^{c-1} \frac{(c\rho)^k}{k!} \right]}$ ;

```

B. Direct Access Mode

Although PQM can improve the queuing performance by moving all jobs from the *Adapter Queue* to our customized multiple subqueues in the driver, it still cannot simplify the complex VMware I/O stack and thus cannot fully utilize the low latency NVMe SSDs and avoid cross-VM interference. To thoroughly reduce the I/O path complexity and support performance isolation, we develop the “*Direct Access Mode*”, which allows *trusted* applications whose hosting VMDKs are attached with our customized VAIO IOFilters to bypass the entire VMware I/O stacks and directly use polling I/Os to the NVMe resource shared by multiple users. DAM is transparent to users since it does not require applications to be re-compiled (e.g., linking with a third-party library). This user-space I/O strives to fully utilize the performance of NVMe SSDs while meeting the diverse requirements from user applications and achieving performance isolation.

As illustrated in the third column of Fig. 2, in this mode, DAM bypasses the entire VMware I/O stack, and provides *Handles* to grant each *trusted* application in the VM the *root privilege* to manage the access permission of I/O queues in NVMe controller via *Private Channels*. Each user thread is assigned to each application inside VMs. The I/O submissions and completions do not require any driver interventions. As

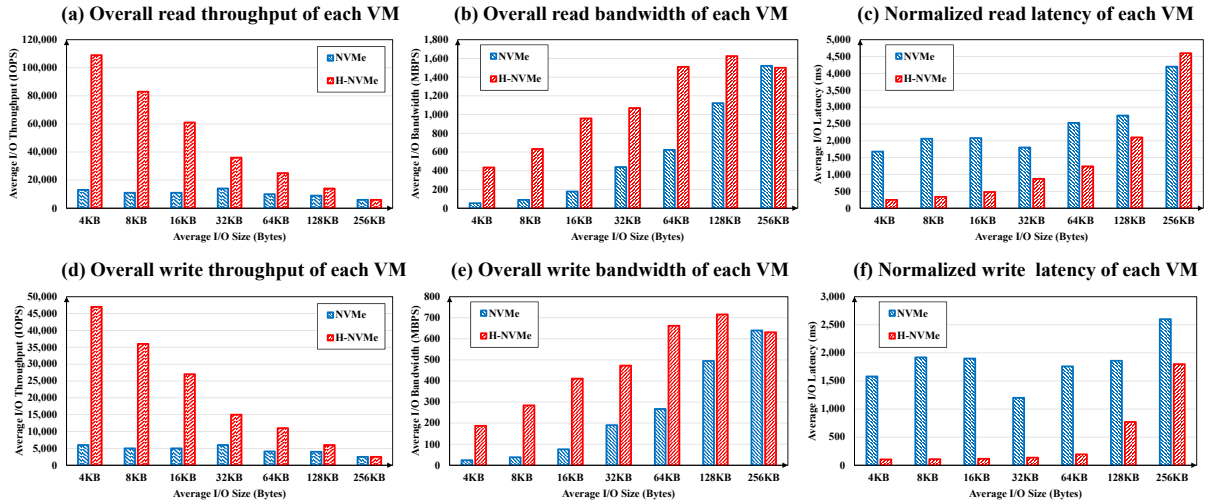


Fig. 4: Throughput, bandwidth, and latency results of 7 different workloads under original NVMe and proposed H-NVMe frameworks.

shown in the yellow box, *Handles* will take over all I/O stack of the hypervisor, and map the user space to the NVMe space. Next, the *Device Queue* will be directly assigned to each *Handle*, without being forwarded by the *World Queue* and *Adapter Queue*.

Algorithm 2: Main Procedures of Direct Access Mode.

```

1 Procedure IOSubmissionNVMeHypervisor()
2   userLib  $\rightarrow$  kernelDriver ;
3   /* Initialization */ ;
4   if curApp  $\in$  userLevelWhitelist then
5     new subQueue;
6     new compQueue;
7     subQueue.memoryRegion.map(curApp.memRegion);
8     compQueue.memoryRegion.map(curApp.memRegion);
9     curHandle=new handle(submissionQueue,
10      completionQueue, doorbellReg);
11     curHandleList+=curHandle;
12   forwardIOCmd(curHandleList, NVMeCtrlr);

```

Alg. 2 further explains the detail of DAM. When an application requests a single I/O queue to DAM, H-NVMe checks whether the application (as well as its hosting VM) is allowed to perform user-level I/Os. If it is in the whitelist (preset based on SLA and security audit requirements), the customized VAIO IOFilter creates a required submission queue and a completion queue (i.e., “*Private Channels*”), see Alg. 2 lines 4-6. H-NVMe then maps their memory regions (including those associated doorbell registers) to the user-space memory region of the application, in lines 7-8. After this initialization process, the application can issue I/O commands directly to the NVMe SSD without any hardware modification or help from the kernel I/O stack, see Alg. 2 lines 9-11.

IV. PERFORMANCE EVALUATION

A. Evaluation Methodology

We implement the proposed H-NVMe framework on VMware ESXi hypervisor 6.0.0 [2]. Table I summarizes our

TABLE I: Host server configuration.

Component	Specs
Host Server	HPE ProLiant DL380 G9
Host Processor	Intel Xeon CPU E5-2360 v3
Host Processor Speed	2.40GHz
Host Processor Cores	12 Cores
Host Memory Capacity	64GB DIMM DDR4
Host Memory Data Rate	2133 MHz
Host Hypervisor	VMware ESXi 6.0.0
Storage Device	Samsung NVMe 1725
Form Factor	HHHL
Per Disk Size	3.2TB
Max Seq Read Bandwidth	6,000 MBPS
Max Seq Write Bandwidth	2,000 MBPS

server configuration. We compare the I/O performance of H-NVMe with the original kernel-based I/O with asynchronous I/O support (i.e., *Kernel I/O*) using the Flexible IO Tester (FIO) benchmark [30]. We allocate around 10% of VMs under Direct Access Mode (DAM) to reflect the premium paid rate in enterprise market. Other VMs are deployed under Parallel Queue Mode (PQM). VMDKs of those DAM VMs are attached with a modified IO Filter [3] driver to support the user polling direct access feature. The IO Filter driver is responsible for collecting and forwarding I/O to *Handles* in Fig. 2, which thus breaks the encapsulation of the VMware and lets VMs directly operate on the NVMe device.

To compare solutions, we evaluate I/O speed results using average I/O throughputs (IOPS), bandwidth (MBPS) and latency for both read and write. To evaluate the cost of different algorithms, we focus on two metrics, i.e., (1) total runtime of a fixed amount of workload, and (2) the amount of data involved in context switching.

B. Throughput, Bandwidth and Latency of Mixed Workloads

We run 7 workloads with different average I/O sizes on 7 different VMs hosted by the same hypervisor, configuring the overall sequential ratio of each workloads varies from 30% to 70%. Fig. 4 illustrates the average throughput, bandwidth, and normalized latency for both read and write I/Os of these

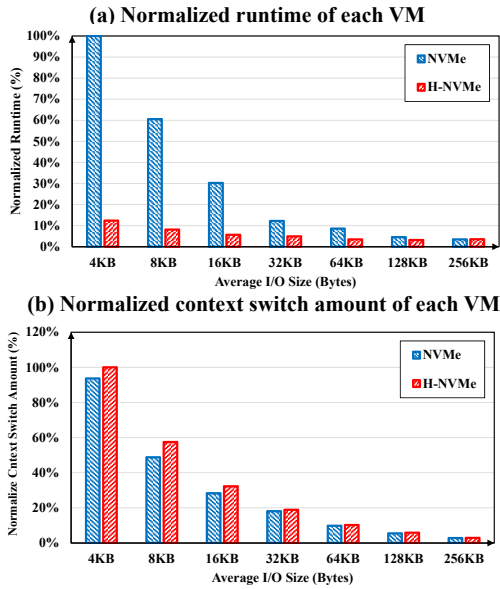


Fig. 5: Normalized runtime and context switch amount of 7 different workloads under original NVMe and proposed H-NVMe frameworks.

workloads under the original NVMe and proposed H-NVMe frameworks. Generally speaking, H-NVMe achieves better performance in all subfigures. We also notice that for small block size workloads, H-NVMe performs much better than the original NVMe framework. However, when workload I/O size increases to 256KB (e.g., the “256KB” bars in Fig. 4(b) and (e)), performance improvement becomes less. This is because larger average I/O sizes have better locality so that there is less room for performance improvement. Similarly, once I/O size reaches 256KB, the read latency of the original NVMe is even slightly better, see in Fig. 4(c).

C. Temporal and Spatial Overheads of Mixed Workloads

We further investigate the overhead of H-NVMe, using Fig. 5 to show both the temporal and spatial overhead of H-NVMe. In detail, from Fig. 5(a), we observe that H-NVMe significantly reduces the total runtime for the same amount of workload compared to the original NVMe framework, especially for those workloads with small I/O sizes. The reason is that unlike workloads with larger I/O sizes having high locality even with the original NVMe framework, workloads with small I/O sizes are more beneficial from PQM’s parallel processing subqueues. Additionally, according to our sensitivity analysis, we find that the “sweet spot” for the average number of subqueues in the NVMe driver is 4, which validates our analysis in Sec. III-A that simply deploying more subqueues will not always improve performance.

On the other hand, we evaluate the spatial overhead by measuring the amount of data in the context switching [4]. As shown in Fig. 5(b), H-NVMe has slightly higher spatial overhead, because it needs to maintain multiple subqueues in the driver layer which requires few extra context switches. However, the difference of their spatial overheads gets smaller when the average I/O size increases. In fact, the penalty of this extra spatial overhead is much lower than the time saved

by H-NVMe. Therefore, we conclude that by eliminating the synchronization locks contention from multiple queues, H-NVMe can significantly reduce the performance overhead of the I/O stack.

D. Stress Tests on Sequential and Random Workloads

To further investigate the boundary performance improvement brought by H-NVMe under extreme cases, we conduct several stress tests on read- and write-intensive, pure sequential and random workloads in this subsection.

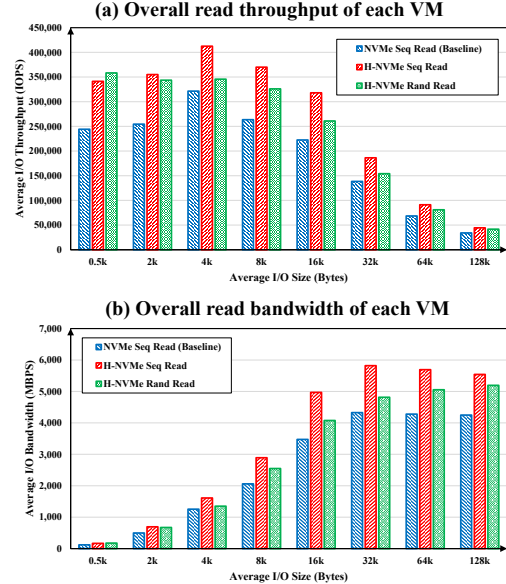


Fig. 6: Read throughput and bandwidth of each VM.

We use the original NVMe’s sequential read and write as the benchmark baseline, which reflects the read and write *steady states* (i.e., the “upper bound”) under the original NVMe framework. We first present the study on the read-intensive workloads. From Fig. 6(a), we observe that H-NVMe has higher throughput and bandwidth in both sequential and random read than the original NVMe upper bound, which validates the effectiveness of H-NVMe. Furthermore, once the average I/O size is greater than 16KB, H-NVMe’s sequential read bandwidth is close to the device capacity (i.e., 6,000 MBPS), and meanwhile, the original NVMe cannot reach that point even at the datapoint of 128KB, see Fig. 6(b).

We next examine the write-intensive sequential and random workloads. As shown in Fig. 7(a), the throughputs of H-NVMe under those sequential write workloads are close to those of NVMe. In fact, even with multiple queues in the I/O stack, these sequential writes still have relatively low latency, and thus there is no so much room to improve for sequential writes. The only exception from our result is when the average I/O size is 4KB, H-NVMe has an outstanding performance (more than 350,000 IOPS). This ascribes to the fact that this 4KB I/O size is aligned to the cache line size in our implementation. It is worth to mention that the random write throughput of H-NVMe is getting close to the sequential write throughput of the original NVMe, once average I/O sizes exceed 32 KB.

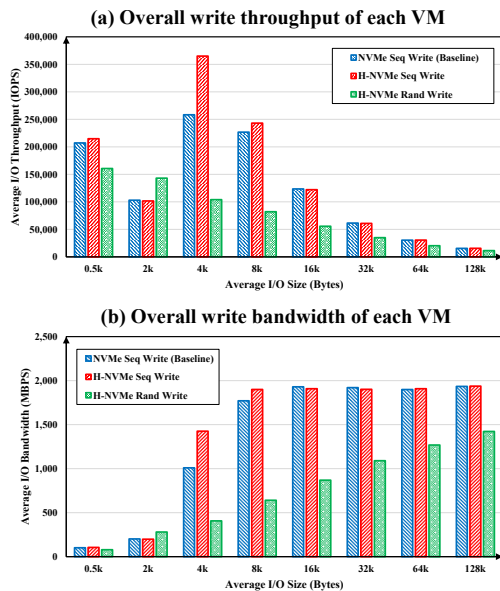


Fig. 7: Write throughput and bandwidth of each VM.

Similarly, Fig. 7(b) also shows that H-NVMe works slightly better than the original NVMe, and the *steady state* is reached after 8KB, where both NVMe and H-NVMe’s bandwidths are close to the device capacity (i.e., 2,000 MBPS) for sequential write. Another observation from Fig. 7(b) is that the random write bandwidth (see the green bar) is increasing linearly with increasing I/O sizes, and thus reduces the gap from the sequential write, which ascribes to H-NVMe’s ability to widen the queuing bottlenecks.

V. RELATED WORK

As the hardware latency constantly decreases, many studies have been conducted to diminish the storage stack latency, along the way. Shin et al. [31] present a low level hardware abstraction layer interface which curtails scheduling delays caused by extra contexts to optimize the I/O path. Studies [32]–[36] focus on how to fully utilize SSD resources in big data and cloud computing platforms such as Apache Hadoop and Spark [37]–[41]. The proposed schemes in [42] rely on a hardware support to expand parallelism inside the SSD. Similar to [31], they also eliminate context switches in the I/O path. Additionally, they exploit polling I/O completion, merging I/O, and double buffering. Likewise, Yang et al. [6] compare polling-based and interrupt-based I/O paths, and eventually come up with the fact that synchronously polling model for I/O completion is much faster in the era of non-volatile memories and very fast SSDs.

To maximize parallelism further, P. Kumar and H. Huang [43] propose Falcon which is a single flush thread per drive, instead of per volume, and separates I/O batching and I/O serving in the storage stack. Studies [44]–[46] investigate data deduplication and replication techniques to improve the performance and reliability of modern datacenters. Besides the above design towards optimizing the I/O path, the kernel also imposes overhead to the storage stack. To address this

issue, researchers have suggested to grant direct access to the user applications without involving the kernel. Caulfield et al. [7] present a new hardware/software architecture which achieves high performance by skipping the kernel involvement and leaving the file-system permission checking phase to the hardware (i.e. their special storage device, Moneta [47]). Aerie [48], another flexible file-system, is introduced by Volos et al. to expose the storage devices to user applications for accessing without kernel interaction. HJ Kim et al. propose NVMeDirect [5] which is a user I/O framework allowing user applications directly access commercial NVMe SSDs by associating NVMe I/O queues to them upon request. Furthermore, the parallelism provided by SSD brings an opportunity for performance isolation for multiple tenants sharing the device.

On one hand, many of these approaches require specialized hardware and/or running complex software frameworks. On the other hand, their focus is on stand alone computer systems running a kernel and several jobs in form of processing threads. However, in data center clusters running hypervisors and several virtual machines, there still exist many unnecessary latency overheads. In contrast, H-NVMe targets the needs in big-data applications which results in a dramatic improvement in large data centers and better performance isolation for VMs by modifying the NVMe driver in VMware ESXi hypervisor.

VI. CONCLUSION

In this paper, we present a hybrid framework called “H-NVMe” to better utilize NVMe SSDs in the modern super-scale cloud computing datacenters. Our work is motivated by the bottlenecks caused by multiple intermediate queues in the current deployment of NVMe in the VM-hypervisor environment, which cannot fully utilize the maximum performance throughput (up to 64K in depth) of NVMe resources. To solve this bottleneck issue, H-NVMe offers two modes to deploy VMs, i.e., “Parallel Queue Mode” and “Direct Access Mode”. The first mode is to increase parallelism and enable lock-free operations by forwarding jobs from the *Adapter Queue* to the our proposed enhanced subqueues in the driver. This mode is a generalized solution and is relatively straightforward to implement. The second mode allows trusted applications with VAIO IOFilter attached to their user VMDBs to directly access NVMe SSDs and bypass the entire I/O stack in the hypervisor layer to further ensure performance isolation, which suits premium users who have higher priorities and the permission to attach IOFilter to their VMDBs. We implement H-NVMe on VMware ESXi, and our evaluation results show that the proposed framework outperforms the original NVMe solution under multiple benchmarks.

REFERENCES

- [1] “NVM-EXPRESS,” <http://www.nvmexpress.org/>.
- [2] “VMware ESXi,” www.vmware.com/products/vsphere-hypervisor.html.
- [3] “vsphere apis for i/o filtering (vaio) program,” <https://code.vmware.com/programs/vsphere-apis-for-io-filtering>.
- [4] J. Bhimani, J. Yang, Z. Yang, N. Mi, Q. Xu, M. Awasthi, R. Pandurangan, and V. Balakrishnan, “Understanding Performance of I/O Intensive Containerized Applications for NVMe SSDs,” in *35th IEEE International Performance Computing and Communications Conference*. IEEE, 2016.

- [5] H.-J. Kim, Y.-S. Lee, and J.-S. Kim, "Nvmedirect: A user-space i/o framework for application-specific optimization on nvme ssds." in *Hot-Storage*, 2016.
- [6] J. Yang, D. B. Minturn, and F. Hady, "When poll is better than interrupt." in *FAST*, vol. 12, 2012, pp. 3–3.
- [7] A. M. Caulfield, T. I. Mollov, L. A. Eisner, A. De, J. Coburn, and S. Swanson, "Providing safe, user space access to fast, solid state disks," *ACM SIGARCH Computer Architecture News*, vol. 40, no. 1, pp. 387–400, 2012.
- [8] J. Huang, A. Badam, L. Caulfield, S. Nath, S. Sengupta, B. Sharma, and M. K. Qureshi, "Flashblox: Achieving both performance isolation and uniform lifetime for virtualized ssds." in *FAST*, 2017, pp. 375–390.
- [9] Z. Yang, M. Hoseinzadeh, A. Andrews, C. Mayers, D. T. Evans, R. T. Bolt, J. Bhimani, N. Mi, and S. Swanson, "AutoTiering: Automatic Data Placement Manager in Multi-Tier All-Flash Datacenter," in *36th IEEE International Performance Computing and Communications Conference*. IEEE, 2017.
- [10] Z. Yang and D. Evans, "Automatic Data Placement Manager in Multi-Tier All-Flash Datacenter," Patent US62/534 647, 2017.
- [11] J. Bhimani, Z. Yang, M. Leeser, and N. Mi, "Accelerating Big Data Applications Using Lightweight Virtualization Framework on Enterprise Cloud," in *21st IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE, 2017.
- [12] D. I. Shin, Y. J. Yu, and H. Y. Yeom, "Shedding light in the black-box: Structural modeling of modern disk drives," in *Modeling, Analysis, and Simulation of Computer and Telecommunication Systems, 2007. MASCOTS'07. 15th International Symposium on*. IEEE, 2007, pp. 410–417.
- [13] Z. Yang, S. Hassani, and M. Awasthi, "Memory Device Having a Translation Layer with Multiple Associative Sectors," Patent US15/093 682, US20 170 242 583A1, 2015.
- [14] "3D XPoint." [Online]. Available: https://en.wikipedia.org/wiki/3D_XPoint
- [15] M. Hoseinzadeh, M. Arjomand, and H. Sarbazi-Azad, "Reducing access latency of mlc pcms through line striping," *ACM SIGARCH Computer Architecture News*, vol. 42, no. 3, pp. 277–288, 2014.
- [16] —, "Spcm: The striped phase change memory," *ACM Transactions on Architecture and Code Optimization (TACO)*, vol. 12, no. 4, p. 38, 2016.
- [17] J. Tai, D. Liu, Z. Yang, X. Zhu, J. Lo, and N. Mi, "Improving Flash Resource Utilization at Minimal Management Cost in Virtualized Flash-based Storage Systems," *Cloud Computing, IEEE Transactions on*, no. 99, p. 1, 2015.
- [18] Z. Yang, J. Tai, J. Bhimani, J. Wang, N. Mi, and B. Sheng, "GREM: Dynamic SSD Resource Allocation In Virtualized Storage Systems With Heterogeneous IO Workloads," in *35th IEEE International Performance Computing and Communications Conference*. IEEE, 2016.
- [19] E. H. Nam, B. S. J. Kim, H. Eom, and S. L. Min, "Ozone (o3): An out-of-order flash memory controller architecture," *IEEE Transactions on Computers*, vol. 60, no. 5, pp. 653–666, 2011.
- [20] Z. Yang, M. Awasthi, M. Ghosh, and N. Mi, "A Fresh Perspective on Total Cost of Ownership Models for Flash Storage in Datacenters," in *2016 IEEE 8th International Conference on Cloud Computing Technology and Science*. IEEE, 2016.
- [21] Z. Yang, M. Ghosh, M. Awasthi, and V. Balakrishnan, "Online Flash Resource Allocation Manager Based on TCO Model," Patent US15/092 156, US20 170 046 089A1, 2016.
- [22] Z. Yang and M. Awasthi, "I/O Workload Scheduling Manager for RAID/non-RAID Flash Based Storage Systems for TCO and WAF Optimizations," Patent US15/396 186, 2017.
- [23] Z. Yang, M. Ghosh, M. Awasthi, and V. Balakrishnan, "Online Flash Resource Migration, Allocation, Retire and Replacement Manager Based on a Cost of Ownership Model," Patent US15/094 971, US20 170 046 098A1, 2016.
- [24] A. Huffman, "NVM Express specification 1.1a," http://www.nvmeexpress.org/wp-content/uploads/NVM-Express-1_1a.pdf, Sep 2013.
- [25] M. Oh, H. Eom, and H. Y. Yeom, "Enhancing the i/o system for virtual machines using high performance ssds," in *Performance Computing and Communications Conference, 2014 IEEE International*. IEEE, 2014, pp. 1–8.
- [26] "SCSI Command." [Online]. Available: https://en.wikipedia.org/wiki/SCSI_command
- [27] "Performance Best Practices for VMware vSphere 6.5, VMware ESXi 6.5, vCenter Server 6.5," <http://www.vmware.com/support/2017>.
- [28] "MMC Queue." [Online]. Available: https://en.wikipedia.org/wiki/M/M_c_queue
- [29] "Erlang C Formula." [Online]. Available: [https://en.wikipedia.org/wiki/Erlang_\(unit\)#Erlang_C_formula](https://en.wikipedia.org/wiki/Erlang_(unit)#Erlang_C_formula)
- [30] "FIO: Flexible I/O Tester," <http://linux.die.net/man/1/fio>.
- [31] W. Shin, Q. Chen, M. Oh, H. Eom, and H. Y. Yeom, "Os i/o path optimizations for flash solid-state drives." in *USENIX Annual Technical Conference*, 2014, pp. 483–488.
- [32] Z. Yang, M. Hoseinzadeh, P. Wong, J. Artoux, and D. Evans, "A Hybrid Framework Design of NVMe-based Storage System in Cloud Computing Storage System," Patent US 62/540 555, 2017.
- [33] J. Wang, Z. Yang, and D. Evans, "Efficient Data Caching Management in Scalable Multi-stage Data Processing Systems," Patent US15/423 384, 2017.
- [34] Z. Yang, J. Wang, and D. Evans, "A Duplicate In-memory Shared-intermediate Data Detection and Reuse Module in Spark Framework," Patent US15/404 100, 2017.
- [35] T. Wang, J. Wang, N. Nguyen, Z. Yang, N. Mi, and B. Sheng, "EA2S2: An Efficient Application-Aware Storage System for Big Data Processing in Heterogeneous Clusters," in *26th International Conference on Computer Communications and Networks (ICCCN)*. IEEE, 2017.
- [36] Z. Yang, J. Wang, and D. Evans, "Adaptive Caching Replacement Manager with Dynamic Updating Granulates and Partitions for Shared Flash-Based Storage System," Patent US15/400 835, 2017.
- [37] J. Wang, T. Wang, Z. Yang, N. Mi, and S. Bo, "eSplash: Efficient Speculation in Large Scale Heterogeneous Computing Systems," in *35th IEEE International Performance Computing and Communications Conference*. IEEE, 2016.
- [38] J. Wang, T. Wang, Z. Yang, Y. Mao, N. Mi, and B. Sheng, "SEINA: A Stealthy and Effective Internal Attack in Hadoop Systems," in *International Conference on Computing, Networking and Communications (ICNC 2017)*. IEEE, 2017.
- [39] H. Gao, Z. Yang, J. Bhimani, T. Wang, J. Wang, B. Sheng, and N. Mi, "AutoPath: Harnessing Parallel Execution Paths for Efficient Resource Allocation in Multi-Stage Big Data Frameworks," in *26th International Conference on Computer Communications and Networks (ICCCN)*. IEEE, 2017.
- [40] H. Harvey, Y. Mao, Y. Hou, and B. Sheng, "Edos: Edge assisted offloading system for mobile devices," *arXiv preprint arXiv:1705.07524*, 2017.
- [41] J. Bhimani, N. Mi, M. Leeser, and Z. Yang, "FiM: Performance Prediction Model for Parallel Computation in Iterative Data Processing Applications," in *10th IEEE International Conference on Cloud Computing (CLOUD)*. IEEE, 2017.
- [42] Y. J. Yu, D. I. Shin, W. Shin, N. Y. Song, J. W. Choi, H. S. Kim, H. Eom, and H. Y. Yeom, "Optimizing the block i/o subsystem for fast storage devices," *ACM Transactions on Computer Systems (TOCS)*, vol. 32, no. 2, p. 6, 2014.
- [43] P. Kumar and H. H. Huang, "Falcon: Scaling io performance in multi-ssd volumes," in *2017 USENIX Annual Technical Conference (USENIX ATC 17)*. USENIX Association, 2017, pp. 41–53.
- [44] J. Roemer, M. Groman, Z. Yang, Y. Wang, C. C. Tan, and N. Mi, "Improving Virtual Machine Migration via Deduplication," in *11th IEEE International Conference on Mobile Ad Hoc and Sensor Systems (MASS 2014)*. IEEE, 2014, pp. 702–707.
- [45] Z. Yang, J. Wang, D. Evans, and N. Mi, "AutoReplica: Automatic Data Replica Manager in Distributed Caching and Data Processing Systems," in *1st International Workshop on Communication, Computing, and Networking in Cyber Physical Systems (CCNCPs)*. IEEE, 2016.
- [46] Z. Yang, J. Wang, and D. Evans, "Automatic Data Replica Manager in Distributed Caching and Data Processing Systems," Patent US15/408 328, 2017.
- [47] A. M. Caulfield, A. De, J. Coburn, T. I. Mollov, R. K. Gupta, and S. Swanson, "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*. IEEE Computer Society, 2010, pp. 385–395.
- [48] H. Volos, S. Nalli, S. Panneerselvam, V. Varadarajan, P. Saxena, and M. M. Swift, "Aerie: Flexible file-system interfaces to storage-class memory," in *Proceedings of the Ninth European Conference on Computer Systems*. ACM, 2014, p. 14.