



Enhancing the Performance of Next-Generation SSD Arrays: A Holistic Approach

JIE ZHANG, Peking University, Beijing, China

SHUSHU YI, Peking University, Beijing, China

XIURUI PAN, Peking University, Beijing, China

YIMING XU, Peking University, Beijing, China

QIAO LI, City University of Hong Kong Department of Computer Science, Kowloon, Hong Kong

QIANG LI, Alibaba Group, Hangzhou, China

CHENXI WANG, University of the Chinese Academy of Sciences, Beijing, China

BO MAO, Software School, Xiamen University, Xiamen, China

MYOUNGSOO JUNG, KAIST, Daejeon, Korea (the Republic of)

All-flash array (AFA) is a popular approach to aggregate the capacity of multiple solid-state drives (SSDs) while guaranteeing fault tolerance. Unfortunately, existing AFA engines inflict substantial software overheads on the I/O path, such as the user-kernel context switches and AFA internal tasks (e.g., parity preparation), thereby failing to adopt next-generation high-performance SSDs. Tackling this challenge, we propose ScalaAFA, a unique holistic design of AFA engine that can extend the throughput of next-generation SSD arrays in scale with low CPU costs. We incorporate ScalaAFA into user space to avoid user-kernel context switches while harnessing SSD built-in resources for handling AFA internal tasks. Specifically, in adherence to the lock-free principle of existing user-space storage framework, ScalaAFA substitutes the traditional locks with an efficient message-passing-based permission management scheme to facilitate inter-thread synchronization. Considering the CPU burden imposed by background I/O and parity computation, ScalaAFA proposes to offload these tasks to SSDs. To mitigate host-SSD communication overheads in offloading, ScalaAFA takes a novel data placement policy that enables transparent data gathering and in-situ parity computation. ScalaAFA also addresses the write amplification and metadata persistence issues and avoids GC-caused latency spikes, by thoroughly exploiting SSD architectural innovations. Comprehensive evaluation results indicate that ScalaAFA can achieve $2.5\times$ write throughput and reduce average write latency by a significant 52.7%, compared to the state-of-the-art AFA engines.

CCS Concepts: • **Information systems** → RAID; Flash memory; • **Software and its engineering** → Secondary storage.

Additional Key Words and Phrases: All Flash Array, User Space, In-Storage Processing, SSD GC

Authors' Contact Information: Jie Zhang, Peking University, Beijing, China; e-mail: jiezh@pku.edu.cn; Shushu Yi, Peking University, Beijing, China; e-mail: shusyi@stu.pku.edu.cn; Xiurui Pan, Peking University, Beijing, China; e-mail: panxr23@stu.pku.edu.cn; Yiming Xu, Peking University, Beijing, China; e-mail: teddyxym@outlook.com; Qiao Li, City University of Hong Kong Department of Computer Science, Kowloon, Hong Kong; e-mail: qiaoli045@gmail.com; Qiang Li, Alibaba Group, Hangzhou, Zhejiang, China; e-mail: flyandsing@126.com; Chenxi Wang, University of the Chinese Academy of Sciences, Beijing, China; e-mail: wangchenxi@ict.ac.cn; Bo Mao, Software School, Xiamen University, Xiamen, Fujian, China; e-mail: maobo@xmu.edu.cn; Myoungsoo Jung, KAIST, Daejeon, Korea (the Republic of); e-mail: mj@camelab.org.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 1553-3093/2025/5-ART

<https://doi.org/10.1145/3736588>

1 Introduction

The last decade has witnessed all-flash arrays (AFA) [39, 44, 51, 65, 69, 79] increasingly adopted as buffer layers in high-performance computing systems and datacenters [19, 52]. These systems have proven critical in speeding up numerous I/O-intensive scenarios, including big data analysis, scientific computing, and machine learning [1, 18, 22, 42, 58, 63, 68, 72]. In comparison to traditional storage media like hard disks, AFAs capitalize on the advantages of solid-state drives (SSDs) such as enhanced throughput, latency, and power efficiency. AFAs bundle multiple SSDs into an array to improve storage capacity at scale, thereby providing a large and uniform storage space. Additionally, AFAs tackle flash errors by integrating data redundancy mechanisms.

Considerable efforts have been expended in academia and industry to make AFAs more practical [5, 28, 39, 71, 74]. For instance, Linux software RAID, known as *mdraid* [5], has been developed to exploit multi-core processors for concurrent parity preparation. Building on *mdraid*, two-phase write schemes [28, 39, 71, 74] have revolutionized the write path. These systems use *replication* (e.g., RAID 10) as a stepping stone to *striping* (e.g., RAID 5 or 6) for fast I/O processing and update data out-of-place to accelerate small writes.

However, with the continual advancement in SSD technology, current AFA implementations risk becoming a bottleneck for future storage systems that will leverage next-generation SSDs. For instance, Samsung PM1743 PCIe 5.0 SSDs can deliver up to 13 GB/s I/O bandwidth [11], a stark contrast to the majority of AFA engines [5, 39] that were originally designed for slower storage interfaces (i.e., SATA) and top out at a maximum throughput of 500 MB/s. To elucidate the performance issue in AFA, we conduct an experiment with a state-of-the-art two-phase write AFA engine, *FusionRAID* [39] (cf. § 3.1 for details). During replication, *FusionRAID* only achieved 4.8 GB/s write throughput, a mere 36.9% of the ideal performance. This is attributed to the storage software stack (e.g., user-kernel context switches and tedious block layers) consuming 54.4% of the CPU cycles, whilst the CPU stall time for SSD I/O only accounts for 20.8%.

SPDK [78], one of the most popular user-space storage frameworks, presents an encouraging approach to lighten the storage software stack. However, directly incorporating existing AFA solutions into SPDK poses significant challenges. First, existing AFA designs rely on locks to facilitate concurrent multi-thread access [5, 79], whereas SPDK operates based on a lock-free principle. Secondly, while SPDK is latency-sensitive [76], AFA frequently experiences latency spikes due to SSD garbage collection (GC). Furthermore, SPDK cannot mitigate the intrinsic shortcomings of existing AFA schemes (e.g., two-phase write). For instance, the “conversion” process from replication to striping involves reading replicated data from SSDs, computing parities, and storing them back in a space-efficient layout, leading to intensive background I/O and computation, which results in notable performance degradation in I/O-intensive scenarios (evidenced by an 88.5% throughput reduction as shown in Figure 3(b)). In addition, out-of-place updates require extra mapping tables to track data locations. The persistence of this metadata is CPU-intensive. Lastly, the write amplification triggered by replication can significantly shorten SSD lifetime.

In response to these challenges, we introduce *ScalaAFA*, an innovative user-space AFA engine designed to maximize the performance of future SSD arrays while maintaining low CPU overhead. In accordance with SPDK’s lock-free principle, *ScalaAFA* substitutes the traditional locks with an SPDK-compatible message-passing-based permission management scheme for concurrent multi-thread accesses. It also minimizes the synchronization overhead with a lightweight storage space abstraction and batch processing method. To maintain consistently low latency, *ScalaAFA* synergistically schedules I/O requests and SSD GC to effectively mitigate interference between them. Recognizing the CPU load imposed by conversion, we propose an inventive data placement policy that curtails background I/O and offloads parity computation to the embedded resources of the SSDs transparently. In the face of frequently updated mapping tables, *ScalaAFA* leverages a hardware-based crash consistency mechanism, significantly reducing software overhead. Finally, to temper the damage of write amplification on SSD lifetime, *ScalaAFA* suggests harnessing the SSD-internal high-endurance write buffer to accommodate

redundant writes. Evaluation results demonstrate that ScalaAFA surpasses leading AFA engines, delivering 2.5× write throughput and reducing average write latency by 52.7%.

Our core **contributions** can be summarized as follows:

- *Constructing a high-performance AFA engine in user space:* An in-depth analysis of the existing AFA engines reveals that the primary hurdles for integrating AFA engines into user space stem from the complex lock mechanisms for concurrent multi-thread accesses and latency spikes induced by GC. To overcome these impediments, we propose to manage write permission with a lightweight storage space abstraction and efficiently grant/retrieve these permissions in batches with a message-passing scheme. This solution conforms to the lock-free principle of user space designs and facilitates thread-level parallelism at a low cost. Furthermore, we employ a comprehensive design to circumvent the latency spikes caused by SSD GC. These solutions decrease the average and 99.99th percentile write latency by 45.2% and 44.1%. To the best of our knowledge, this is the first study to design an AFA engine in the user space.
- *Offloading conversion to SSDs:* In two-phase write AFAs, data chunks belonging to the same stripe are dispersed across different SSDs after replication. Offloading conversion to SSDs is non-trivial, as it requires the host to transfer data chunks to the target SSD manually, which imposes extra I/O and huge CPU burdens. We propose an innovative data placement policy that can transparently gather data chunks into the target SSDs. Specifically, when the AFA engine generates redundant writes for fault tolerance, we redirect the redundant replicas of the same stripe to the SSD where the parity will be stored. By leveraging the SSD built-in XOR engine to calculate the parity codes from the local replicas, we eliminate the need of host involvement (i.e., computation and data transfer). This design further improves the write throughput by 36.9% and reduces average latency by 36.1%.
- *Optimizing two-phase write with holistic designs:* Prior work has revealed that two-phase write AFA can outperform traditional AFA (e.g., RAID 5) especially when serving small write requests. Nevertheless, how to conceal the inherent drawbacks of two-phase write (i.e., metadata persistence and write amplification) is still unsolved. Notably, we propose to tackle these challenges by thoroughly exploiting SSD architectural innovations. We suggest a hardware-based crash consistency mechanism that employs the out-of-band areas of SSDs to persist the metadata with minor overhead. Moreover, we take a comprehensive design that accommodates transient data within the SSD-internal durable buffers and avoids flushing them to vulnerable flash cells. This design succeeds in reducing the impact of write amplification by 38.6%.

2 Background

2.1 SSD Internal

Baseline architecture. Figure 1 depicts a common architecture of modern SSDs [4, 81]. The SSD is comprised of multiple embedded processors (e.g., ARM), a DDR controller, and specialized processing elements, such as DMA and XOR engines [16, 17, 73]. The XOR engine bolsters the device’s reliability by calculating parity codes for data stored in the SSD. These processors are linked to a flash backbone through the flash physical layer (PHY). The flash backbone consists of 4 to 16 channels, each connecting to several flash packages. Each flash package encloses multiple flash dies, each comprised of 2 to 4 flash planes. A single plane can be divided into thousands of flash blocks. The flash blocks can be categorized as *single-level-cell (SLC)* and *multiple-level-cell (MLC)* based on the number of data bits stored in a flash cell [12, 30]. The SLC blocks offer shorter I/O latency, extended endurance, and reduced capacity, while the MLC blocks provide larger capacity but exhibit longer I/O latency and reduced endurance. To optimize SSD lifetime, capacity, and performance, SSD manufacturers use SLC blocks as the write buffer to accommodate small writes and employ MLC blocks as storage backend [35, 37]. A flash block contains hundreds of flash pages, each divided into data and *out-of-band (OOB)* areas. The data area stores data whilst the OOB area stores metadata (e.g., error-correction codes [36]). Note that the size of OOB area in NAND flash is typically greater than what metadata requires (e.g., tens of bytes reserved per flash page [17, 31]).

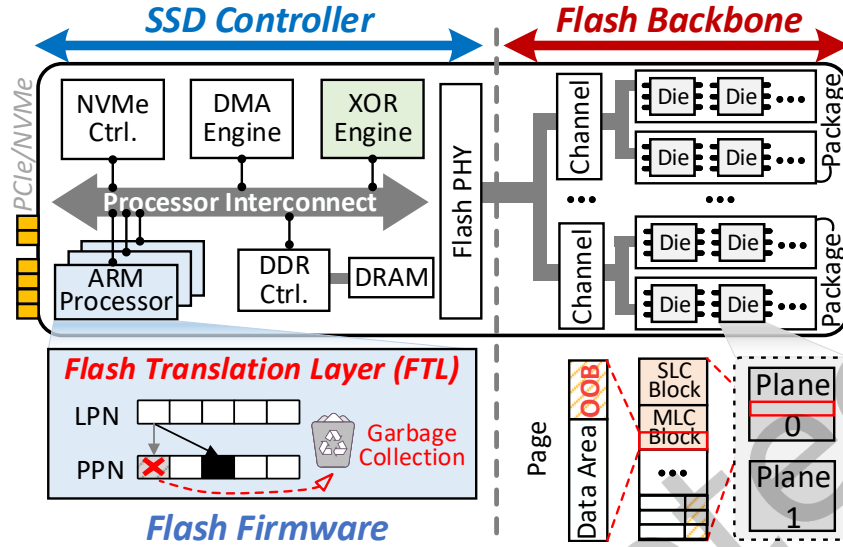


Fig. 1. Details of SSD internal.

Flash firmware. NAND flash only supports out-of-place updates due to its physical attributes [59]. To shield users from flash intrinsic, the flash firmware internally constructs an indirection layer, known as *flash translation layer (FTL)* [21, 64], to remap incoming write requests to new flash pages and invalidate the stale flash pages (cf. Figure 1). FTL maintains the mapping information between the request logical address (i.e., logical page number, LPN) and the flash physical address (i.e., physical page number, PPN) in the SSD-internal DRAM. The flash firmware also persists this mapping information in the OOB area, enabling it to be recovered after system corruption or reboot. Flash page invalidation due to write requests can significantly diminish the available SSD capacity for users. To tackle this problem, the firmware performs garbage collection (GC) to reclaim invalid pages. Specifically, it selects a victim block and moves all valid pages in it to an erased block. Subsequently, the victim block is erased and reused. Note that GC can significantly degrade the performance of SSD [39, 44, 51].

2.2 All-Flash Array

All-flash array (AFA) is a storage organization that groups multiple SSDs as a single logical unit to aggregate their capacity and throughput while providing fault tolerance. Existing AFA designs can be classified into two types, stripe write and two-phase write, based on how data is written to SSDs.

Stripe write AFA. Stripe write AFA orchestrates data as *stripes*. Each stripe consists of k data chunks and m parity chunks with a fixed size (e.g., 64 KB). The parity chunks are calculated with a specific error-correction code (ECC) algorithm (e.g., Reed-Solomon codes [70]). The $k + m$ chunks are distributed and stored in $k + m$ SSDs (we call this *striping layout*). If a few data chunks are lost due to SSD failures, the same number of parity chunks can be used to recover the lost data chunks. Therefore, a $k + m$ AFA can tolerate up to m SSD corruptions. Figure 2(a) shows how a stripe write AFA serves write requests. When a write request arrives, it first slices the request into multiple data chunks. If the number of data chunks is less than k , the write request is called *partial write*. Partial write, unfortunately, is unfriendly to stripe write AFA. The AFA engine must fetch the missing data chunks belonging to the same stripe from SSDs (①). Then, the host CPU calculates the parity chunks for this stripe (②). Finally, both the data and parity chunks will be written to the storage devices (③). This procedure is called

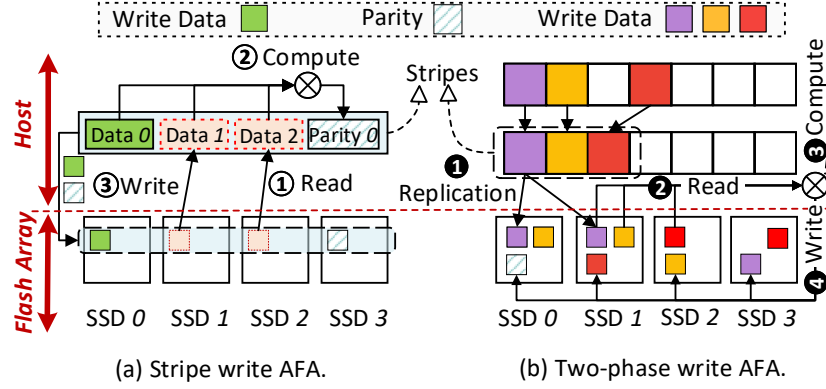


Fig. 2. Write paths of all-flash arrays ($k=3$ and $m=1$).

read-construct-write. Partial writes introduce extra data reads and intensive parity updates, which significantly delay the I/O completion time.

Two-phase write AFA. To tackle partial write issue, prior work proposes two-phase write AFA (or log-structured write AFA) [28, 39, 71, 74], which employs replication as the prelude of striping to absorb small write requests. Two-phase write AFA [39] exhibits superiority over stripe write AFA [5] in terms of both throughput and latency (e.g., 51.2% throughput improvement and 33.9 % average latency reduction for 64 KB sequential write, cf. § 6.2), two of the most crucial metrics when employing AFAs as the buffer layers in datacenters [18, 42, 63, 68]. Figure 2(b) shows the procedure of two-phase writes. Specifically, two-phase write AFA writes data chunks in two phases: *replication phase* and *conversion phase*. In the replication phase, the data chunks are replicated into $m + 1$ copies and distributed across $m + 1$ SSDs (①), which provides the same fault tolerance as stripe write AFA with $k + m$ SSDs. When the size of replicated data exceeds a given watermark (e.g., 5% of the AFA capacity), user I/O requests are delayed and the conversion phase starts. In particular, the host reads all data chunks of the same stripe from SSDs (②) and computes the parity chunks (③). Afterward, these data and parity chunks are written back to SSDs and stored in the space-efficient striping layout (④). Finally, the space previously used to store replicated data chunks will be recycled. Note that all the updates in two-phase write AFA are out-of-place. Therefore, when updating data chunks, two-phase write AFA only needs to replicate the updated data to new spaces and invalidate the stale data, which avoids the tedious read-construct-write procedure.

3 Preliminary Study

3.1 Challenges

While two-phase write has demonstrated its superiority over stripe write in terms of both throughput and latency, we observe that the designs of existing two-phase write AFA engines are still the bottleneck when adopting high-performance SSDs. To illustrate this, we reproduce FusionRAID [39], one of the state-of-the-art two-phase write AFA engines, and set up an experiment with 4+1 ($k + m$) Samsung 980 Pro SSDs [3] to analyze its write performance. We use fio [14] and perf [25] to evaluate the I/O performance and capture CPU cycles of the key functions in the storage stack.

Challenge in replication phase. Figure 3(a) shows the 64 KB sequential write throughput and CPU overhead breakdown with different numbers of I/O threads in the replication phase of FusionRAID (*FSR-Rep*). We categorize the overhead into five parts: preAFA, AFA, postAFA, I/O, and Other. preAFA is the time of submitting I/O requests from user space to FSR-Rep through context switches and block layers, while AFA represents the time consumed by

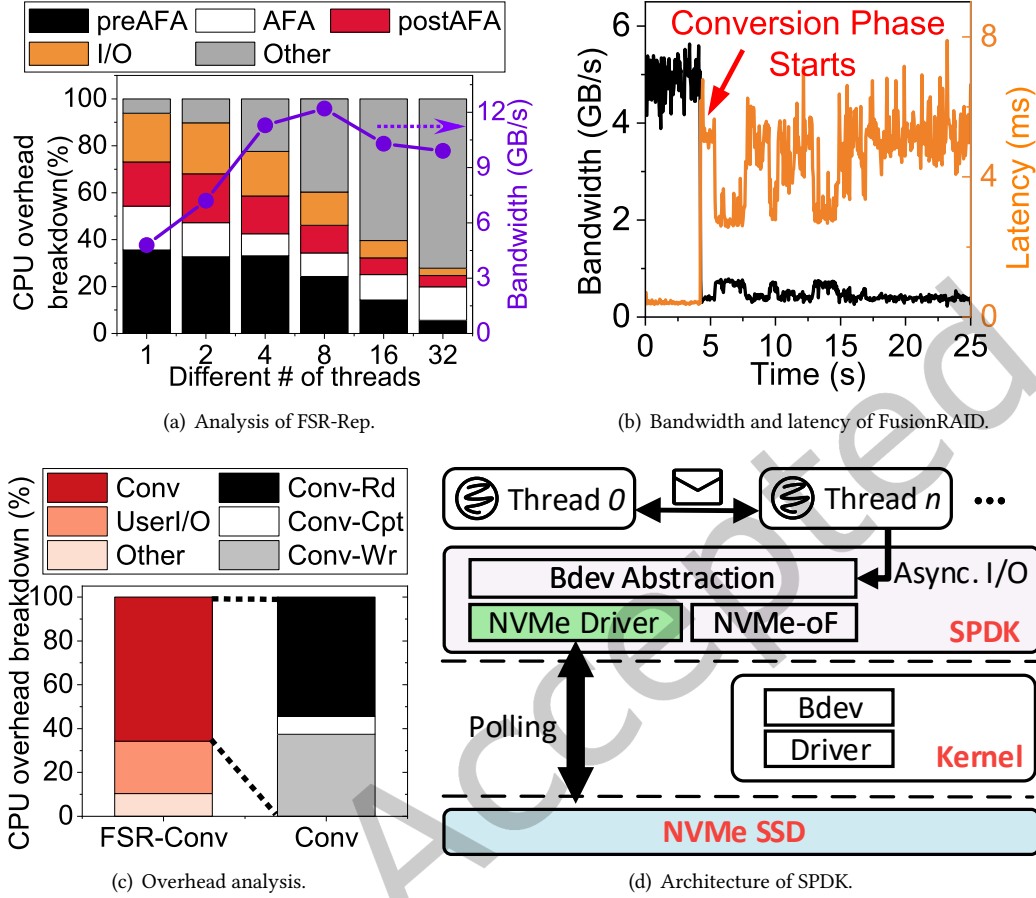


Fig. 3. Deep analysis of two-phase write AFA and key insights.

FSR-Rep; postAFA is the time of sending I/O requests to SSDs through the NVMe driver; I/O is the time of serving I/O requests in SSDs; Lastly, Other summarizes other software overheads (e.g., CPU spins for synchronization among the I/O threads and kernel worker threads [69, 79]).

FSR-Rep achieves 4.8 GB/s and 7.2 GB/s write throughput with 1 and 2 threads, which are only 36.9% and 55.4% of the ideal case (i.e., $(k + m)/(m + 1) * S$, where S is the peak throughput of a single SSD, that is, 5.2 GB/s for Samsung 980 Pro SSD). This is because the frequent user-kernel context switches and the tedious storage stack (i.e., preAFA and postAFA) consume tremendous CPU ticks. For example, with one thread, I/O only accounts for 20.8%, while preAFA and postAFA consume 54.4%. One way to improve the throughput is employing more CPU resources. For example, FSR-Rep achieves 12.2 GB/s with 8 threads. More than 8 threads instead degrades the throughput because inter-thread synchronization (i.e., Other) dominates the CPU overhead [69, 79]. However, allocating 8 threads per 4+1 AFA is still infeasible for most existing storage servers, which equip hundreds of SSDs but limited CPU resources. For instance, PowerStore 500T [10] holds up to 97 SSDs while only equipping two

24-core CPUs. To summarize, *the software overhead imposed by the tedious storage stack has become a hindrance to achieving high performance with limited CPU resources (Challenge 1).*

Challenge in conversion phase. The performance becomes worse when the conversion phase starts (*FSR-Conv*). We continuously send 64 KB sequential write requests to FusionRAID with one I/O thread. Figure 3(b) shows the write throughput and average latency over time. When the conversion phase starts, write throughput drops from about 4.8 GB/s to 550 MB/s, and the latency increases from 0.4 ms to 5.9 ms. We break down the CPU overhead of *FSR-Conv* into three parts, which are shown in Figure 3(c). *Conv* is the time consumed by the conversion in the background. *UserI/O* is the time of serving I/O requests sent by the user (i.e., replication of 64 KB sequential write). Lastly, *Other* includes other software overheads in *FSR-Conv* (e.g., CPU spins for synchronization among the I/O thread and kernel worker threads). *Conv* accounts for 65.7% while *UserI/O* only consumes 23.9% of the CPU resources. We further categorize *Conv* into *Conv-Rd*, *Conv-Cpt*, and *Conv-Wr*, which represent the overheads of the aforementioned read, compute, and write operations (i.e., ②, ③, ④ in Figure 2(b)), respectively. As shown in Figure 3(c), the read and write (i.e., *Conv-Rd* and *Conv-Wr*) overheads dominate the conversion phase, while *Conv-Cpt* only accounts for 8.3%. To sum up, *in the conversion phase, the background tasks (i.e., I/O and computation) significantly degrade the performance of user I/O (Challenge 2).*

Intrinsic issues of two-phase write. Apart from the aforementioned performance penalties, two-phase write AFA engines impose three more challenges.

First, *two-phase write introduces extra metadata that need persisted, which increases the crash consistency cost (Challenge 3).* Specifically, to support out-of-place updates, the AFA engine maintains extra mapping tables in host memory to record where data is actually stored. As these mapping tables are updated frequently, persisting them imposes a huge burden on the host (e.g., by writing an undo or redo log [29] before every update).

Second, *two-phase write causes significant write amplification (Challenge 4).* Assume that $1\times$ data needs writing to a 4+1 two-phase write AFA. In the replication phase, the AFA engine duplicates the data by $2\times$ to provide guaranteed fault tolerance. In the conversion phase, it reads $1\times$ data out for parity computation and writes $1\times$ data and $0.25\times$ parity back to SSDs. The total write amplification is $3.25\times$. These extra writes significantly shorten the lifetime of SSDs.

Lastly, *AFA is GC-sensitive (Challenge 5).* Prior work [57, 75, 77] has revealed that SSD GC can drastically hamper its throughput and latency. This problem occurs more frequently in AFA. Because in AFA, an I/O request needs collaborative services from multiple SSDs (e.g., $m + 1$ SSDs for a chunk write request in the replication phase). If any SSD becomes straggler due to GC, AFA cannot complete the write on time.

3.2 Key Insights

User-space storage stack. The user-space storage framework, such as Storage Performance Development Kit (SPDK) [78], is a promising solution to address **Challenge 1**, that is, the software overheads in storage stack. Figure 3(d) shows the architecture of SPDK. SPDK employs a block device abstraction called *Bdev* to perform the same functions as the block device layer of kernel. SPDK also implements a user-space, asynchronous, polling-based NVMe driver. Therefore, users can directly access NVMe SSDs in user space without trapping into the intricate kernel. These designs achieve a significant acceleration of the storage stack. However, it is non-trivial to directly integrate existing AFA designs into SPDK. This is because traditional AFA engines rely on multiple complicated locks to facilitate thread-level parallelism [5, 69, 79], which violates the lock-free principle of SPDK. *Our key insight is that the SPDK-compatible message-passing mechanism can be the alternative to locks for supporting multi-thread access.*

SSD-internal hardware resources. It is non-trivial to solve **Challenges 2~5** with purely software solutions. For example, the read, compute, and write operations in the conversion phase heavily rely on the host CPU. Even if we can optimize parity computation with sophisticated algorithms, the major overheads (i.e., *Conv-Rd*

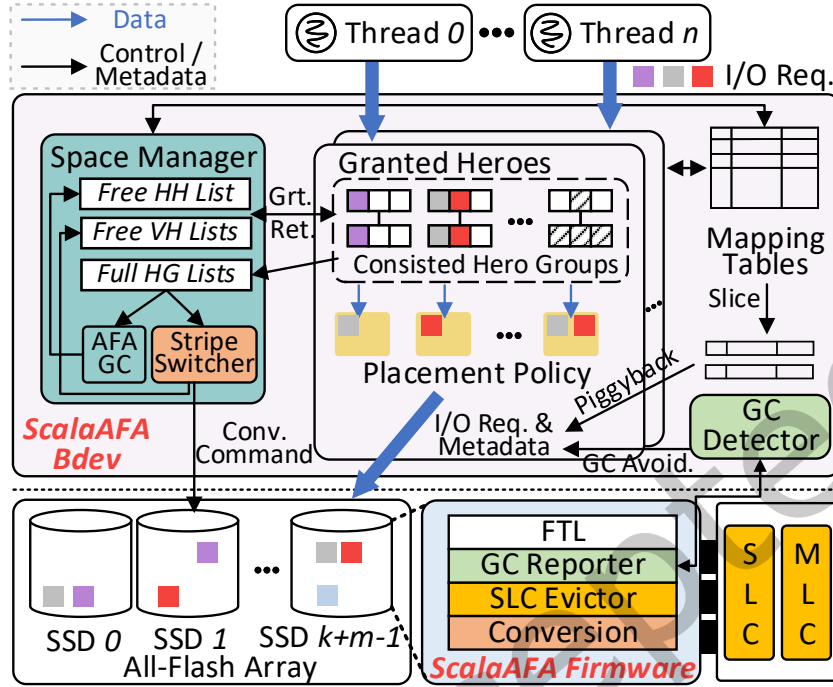


Fig. 4. Architecture of ScalaAFA.

and Conv-Wr in Figure 3(c)) remain unavoidable. In addition, although we can predict SSD GC based on the latency spikes and mitigate its impact by serving I/O requests with SSDs that are not in GC [39], the issued I/O requests are already postponed by GC, which leads to a poor user experience. Fortunately, holistic designs can be a promising solution to address these issues. *Our key insight is that the available resources in SSDs can overcome the software constraints of existing AFA designs, making them fit for next-generation storage.*

4 ScalaAFA Overview

Inspired by the above analysis and key insights, we propose *ScalaAFA*, a high-performance AFA engine built from holistic designs, which overcomes all the aforementioned challenges by **embracing user-space storage stack** and **exploiting SSD-internal hardware resources**. Figure 4 illustrates the architecture of ScalaAFA. Note that, ScalaAFA focuses on improving the write performance of AFA engines and follows the conventional designs of read path [39]. Prior work [69, 79] has proved the scalability and high performance of read operations in available AFA solutions.

For **Challenge 1**, we adopt SPDK [78] to take advantage of its high-performance NVMe driver and lightweight storage stack. Considering the lock-free principle of SPDK, ScalaAFA replaces the conventional lock mechanism with an SPDK-compatible message-passing scheme to avoid write collision among threads. Specifically, ScalaAFA first abstracts storage space with a novel data structure named *hero*¹. Afterward, ScalaAFA employs a *space manager* to manage the write permissions of all the storage spaces (i.e., heroes). These permissions are granted to

¹Hero is the name of the I-shape piece in Tetris.

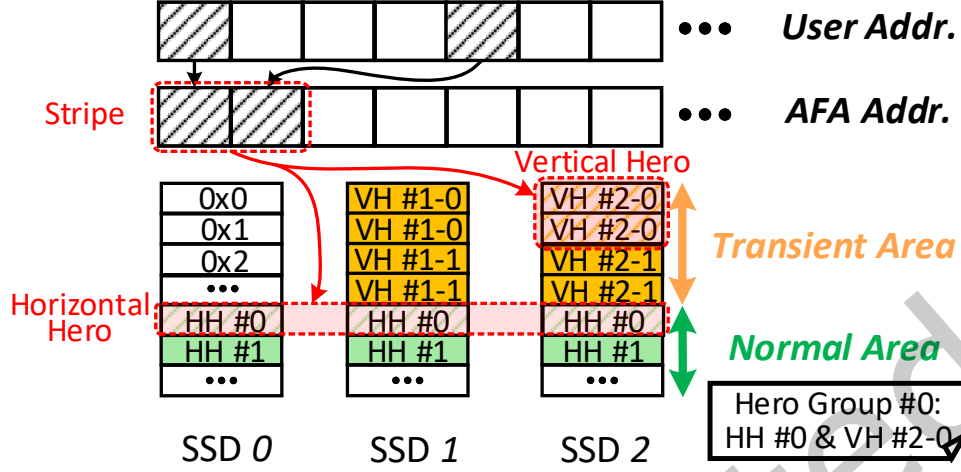
threads via message passing in batches. Threads can only serve write requests with the heroes where permissions are already granted.

For **Challenge 2**, a promising solution to relieve the CPU burden is to offload the parity computation to the storage, similar to the prior in-storage processing approaches [40, 46, 56, 62]. Note that parity computation only consumes minor CPU resources (cf. Conv-Cpt in Figure 3(c)). Such computation can be easily taken with available on-device hardware (e.g., the SSD-embedded XOR engines), which is originally designed to enable SSD-internal parity computation. For instance, our prototype takes only 20 μ s and 16 mW dynamic power to calculate a 64 KB parity chunk from 6 data chunks (cf. § 6.1 for details). However, this method cannot address the background read and write issues in the conversion phase, which have become the major bottleneck in two-phase write AFA (cf. Conv-Rd and Conv-Wr in Figure 3(c)). This is because two-phase write AFA scatters data chunks across different SSDs in the replication phase. The conversion phase requires the host CPU to copy all the corresponding data chunks to the target SSD before offloading parity computation to that SSD. Tackling this issue, we propose a novel *placement policy*, which can transparently gather data chunks of the same stripe to the SSD where parity chunks will be stored. Afterward, in the conversion phase, ScalaAFA only needs to send a command through the *stripe switcher* to the SSD. Once receiving these commands, the *conversion* module in each SSD is in charge of computing parity chunks on device (without collecting data from other SSDs) and storing the parity locally (without writing to other SSDs).

For **Challenge 3**, one feasible solution is to persist the mapping tables in host-side battery-backed DRAM [39]. However, it increases monetary costs and cannot provide enough fault tolerance (i.e., data cannot be located if the DRAM fails). It is also impractical to store them in SSD internal battery-backed DRAM as its limited capacity only keeps up to tens of MB of data persistent [15]. To address this, we propose to persist these mapping tables in the OOB area of SSDs. In particular, considering OOB is scattered across different physical pages of SSDs (cf. Figure 1), ScalaAFA first reconstructs the mapping tables into a segmentable data structure. It then slices and piggybacks the mappings to SSDs via write requests. Finally, ScalaAFA persists the sliced mappings in OOB when data is written to the data area. Note that writing data and its metadata in the data and OOB areas can be done by a single flash program operation, the cost of which is negligible.

For **Challenge 4**, while our novel placement policy has eliminated extra I/O in the conversion phase, the replication phase still causes $m + 1$ times write amplification. We alleviate its damage to SSD endurance by fully utilizing the durable SLC buffer within SSDs. Specifically, ScalaAFA avoids flushing redundant replicas from the SLC buffer to the vulnerable MLC blocks. This is because these replicas will be invalidated soon in the conversion phase (cf. § 2.2). To this end, ScalaAFA deploys *SLC evictor*, which gives a low priority to these replicas when selecting victims to evict.

For **Challenge 5**, we employ a *GC reporter* and a *GC detector* in flash firmware and the host to collaboratively schedule every GC event by being aware of the busy time of both the host and the SSDs. GC reporter and GC detector are leveraged to report and handle GC states, respectively. To be specific, when SSD is going to perform GC, the GC reporter reports this state to the GC detector by piggybacking the information in the I/O completion message. Afterward, the GC detector executes the *GC avoidance mechanism*. Specifically, in the replication phase, write requests are redirected to be served with $m + 1$ SSDs that are not in GC. However, the GC issue is a different story in the conversion phase as user I/O can be blocked by the background conversion, which originally needs to collect data from k SSDs and write data and parity back to all $k + m$ SSDs. If any SSD is in GC, the user I/O will be delayed. Fortunately, thanks to our placement policy that can transparently gather data chunks of the same stripe into the target SSDs beforehand, the conversion phase can now be completed with only m SSDs (i.e., the SSDs that are responsible for parity computation). Thus, by executing conversion and handling I/O requests with SSDs that are not in GC, ScalaAFA can also avoid latency spikes in the conversion phase.

Fig. 5. Space abstraction of ScalaAFA ($k=2$ and $m=1$).

5 Design Details of ScalaAFA

5.1 Storage Space Abstraction

Figure 5 shows the storage space abstraction of ScalaAFA. From the user's perspective, ScalaAFA functions as a standard block device, which exposes a continuous space, referred to as *user address space*, and enables random writes as well as in-situ updates. The user-written data will be sequentially logged into an intermediary space, named *AFA address space*. Contiguous k chunks (also named k slots) in AFA address space are orchestrated as a *stripe*. Each stripe documents the storage locations of the k data chunks and their corresponding m parity chunks in the SSDs (i.e., in which SSD and the SSD logical address in that SSD). ScalaAFA partitions the storage space of SSDs into two virtual areas: *normal area* and *transient area*. The former is used to store long-term data (e.g., data chunks after conversion) while the latter accommodates transient data (e.g., replicated data chunks). Note that the size of the transient area is variable and determined by the total capacity of the SLC buffers (cf. § 2.1) within SSDs. ScalaAFA organizes the two areas as sets of *horizontal hero* (HH) and *vertical hero* (VH), respectively. In the normal area, $k + m$ chunks with the same SSD logical address from $k + m$ SSDs are grouped as an HH (e.g., HH #0 in Figure 5), while contiguous k chunks in the same SSD make up a VH in the transient area (e.g., VH #2-0 in SSD 2). Further, ScalaAFA composes one HH and m VHs from m different SSDs as a *hero group* (HG), such as the HG #0 that consists of HH #0 and VH #2-0 (we will describe the rationale of HG in § 5.3). With the abstraction of HG, ScalaAFA can track the storage locations of data and parities by associating stripes with HGs.

5.2 Enable Lock-free Multi-Thread Access

ScalaAFA supports multiple I/O threads to access it concurrently. Similar to conventional block devices, ScalaAFA does not provide sequential consistency [13]. Thus, there is no need to prevent I/O threads from accessing the same user address simultaneously. However, to prohibit these threads from mapping different user addresses to the same AFA address or mapping different AFA addresses to the same SSD storage location (i.e., HG), ScalaAFA has to manage the write permission of AFA address space and SSD storage space. ScalaAFA achieves this with a user-space-friendly (i.e., lock-free) message-passing scheme.

Considering the communication overhead, we first simplify address mappings to reduce the need for communicating. Specifically, we fix the mappings between stripes and HHs. One stripe is bound with one HH that has

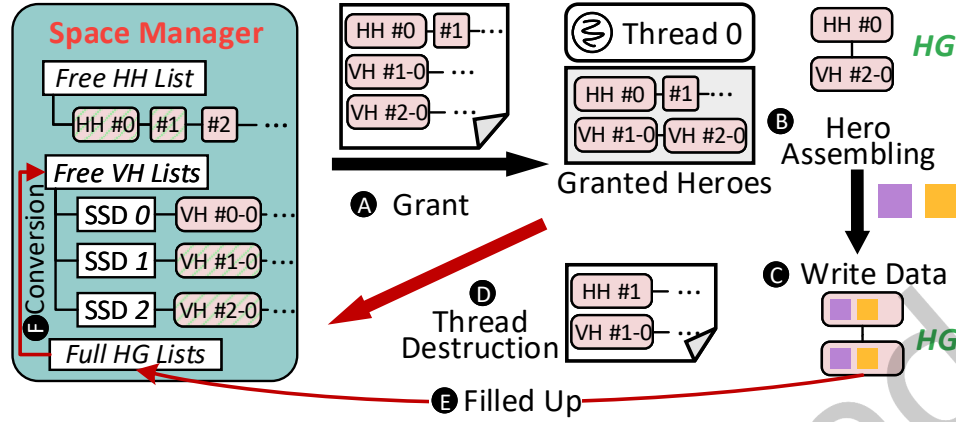


Fig. 6. An example of write permission management.

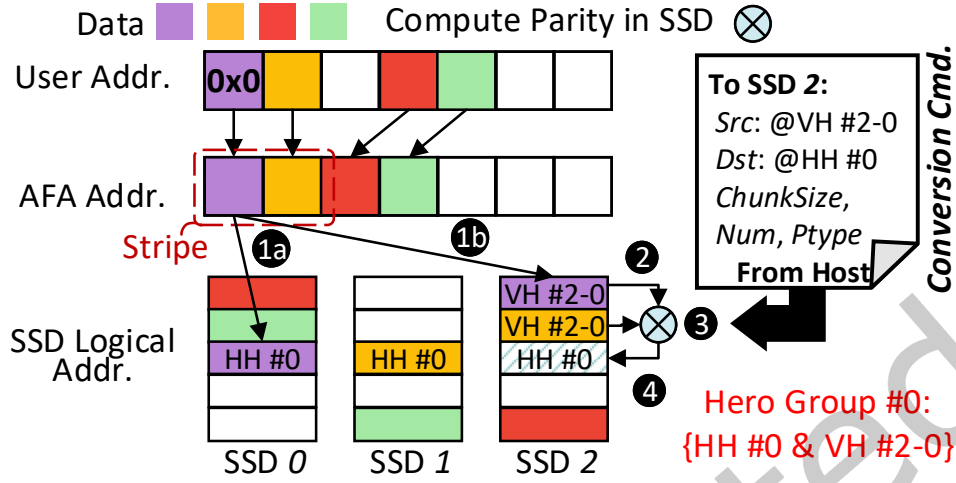
the same offset. For example, the first stripe in the AFA address space is bound with the first HH in the normal area (i.e., HH #0 in Figure 5). Therefore, we only need to assure that different user addresses will not be mapped to the same HH and VHs (i.e., the same HG) without considering the intermediary AFA address. To this end, the *space manager* is responsible for granting and retrieving the write permissions of heroes. Figure 6 gives an example of our solution.

Grant. Initially, all write permissions of heroes belong to the space manager (i.e., *Free HH List* and *Free VH Lists*). When an SPDK I/O thread is set, it asks the space manager for VHs and HHs. Then, the space manager grants heroes in batches (e.g., $1024 * m$ VHs from m SSDs and 1024 HHs at once, **A**). Afterward, if an I/O thread wants to write data chunks to a certain user address, it first assembles granted HH and VHs as an HG (**B**). Then, the data will be written to the HG (**C**), and the user address will be mapped to the corresponding stripe that is bound with the HH. If heroes are run out, the I/O thread will ask the space manager for more grants.

Retrieve. The space manager can retrieve heroes in four ways: (1) When an I/O thread is being destructed, it returns unused heroes to the space manager (**D**); (2) When an HG is filled up by an I/O thread, it will be inserted into a list (i.e., *Full HG Lists*) associated with this thread (**E**). When the CPU is idle or heroes are running out (i.e., the conversion phase starts), the space manager scans the list of each I/O thread and conducts conversion for HGs in the list. Note that, after conversion phase, data and parity chunks are stored in the more space-efficient striping layout (cf. § 2.2). Therefore, the saved space (heroes) can be recycled (**F**); (3) If no more heroes can be granted, the space manager broadcasts a message to recall unused heroes from all I/O threads. The more often a thread requests heroes over a period of time, the fewer heroes will be recalled from that thread; (4) Out-of-place write invalidates stale data and diminishes available storage space. The space manager executes AFA-level GC [39] to recycle these spaces and reuse heroes.

5.3 Evolve the Write Path

Replication phase with placement policy. With the aforementioned space abstraction, we can place data chunks in heroes to transparently gather them in the SSDs where conversion will be executed. Figure 7 illustrates this process. In the replication phase, if one data chunk is written to the stripe, it will be replicated $m + 1$ times. One of them (named *primary copy*) will be stored in HH while the other m copies (named *backup copies*) are written to m VHs. Taking Figure 7 as an example, the user sends a write request to 0x0 and the data chunk is logged in the first slot of the stripe, which is mapped with HG #0. Therefore, in the replication phase, the primary

Fig. 7. Procedure of conversion offloading ($k=2$ and $m=1$).

copy is written to the first slot of HH #0 (1a), and the backup copy is placed in the first slot of VH #2-0 (1b). Repeating this process, data chunks belonging to the same stripe will be transparently gathered in VHs (e.g., the chunks in VH #2-0).

Conversion phase with offloading. Afterward, when the host is idle or VHs are running out, ScalaAFA starts conversion phase and offloads the conversion tasks to SSDs where the VHs locate (e.g., SSD 2). To make the SSDs aware of this, ScalaAFA extends the NVMe command set with a *conversion command*. We include the following information in the new command. *Src*: offset of the VH in SSD; *Dst*: the logical address where parity chunk will be stored; *ChunkSize*: the size of chunks; *Num*: the number of data chunks in one stripe (i.e., k); *Ptype*: the type of parity chunk that will be computed. *Ptype* is used to determine the parity computing method. For example, in Reed-Solomon codes, *Ptype* is the row number of the encoding vector in Vandermonde matrix [70]. Once receiving the command, the SSD executes conversion locally. Specifically, the SSD controller reads *Num* chunks with *ChunkSize* from *Src* (2) to its internal DRAM. It then computes the parity chunk according to *Ptype* (3) and writes the parity chunk to *Dst* (4). Finally, the SSD cleans the data in VH (e.g., VH #2-0) by marking the flash pages as invalid. Note that, after 4, the HH (e.g., HH #0) has stored data and parity chunks in striping layout. Thus, ScalaAFA avoids extra writes for scattering chunks (cf. § 2.2). By doing these, ScalaAFA successfully reduces the host-SSD I/O thereby mitigating the performance degradation in the conversion phase.

5.4 Persist the Metadata

Mapping table. As shown in Figure 8, ScalaAFA maintains two mapping tables in host memory: *AFA Mapping Table (AMT)* and *Hero Group Mapping Table (HGMT)*. AMT maps user addresses to AFA addresses, which records a 32-bit *AFA chunk number* for each chunk in user address space. HGMT translates AFA addresses (i.e., stripes) to SSD storage locations (i.e., HGs). Since stripes and HHs are one-to-one mapped (cf. § 5.2), ScalaAFA only needs to record the mapping information between stripes and VHs. To this end, each row in HGMT represents a stripe and maintains a pointer to a VH list that records the m VHs of this HG. Each item in the VH list contains two components: an 8-bit *SSD ID* to record which SSD the VH locates in and a 32-bit *VH number*, which is the offset of the VH in that SSD. If a stripe is in striping layout (i.e., after conversion), its VHs have been recycled (cf. § 5.3).

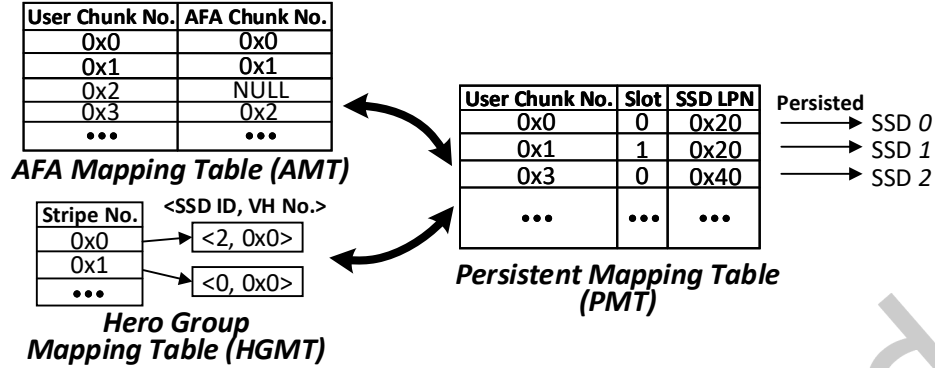


Fig. 8. Mapping tables in ScalaAFA.

and the VH list is NULL. Note that, in ScalaAFA, the chunk size is 64 KB by default. Therefore, to create a 6+3 AFA with 2TB SSDs, ScalaAFA needs only about 928 MB (0.005%) memory to maintain AMT and HGMT.

Persisting the mapping table. As two-phase write AFA has to update its mapping tables frequently, persisting such metadata, in turn, introduces huge software overheads (e.g., by persisting an undo/redo log [29] before every update). One possible solution to tackling this challenge is to maintain the metadata in SSD-internal OOB. However, OOB is scattered across different physical pages and the size of OOB in each flash page is tens of bytes (cf. § 2.1). It is difficult to store the entire AMT and HGMT in one OOB.

To better utilize the constrained OOB, we first convert AMT and HGMT to a segmentable mapping table, called *persistent mapping table (PMT)*. As shown in Figure 8, each entry in PMT corresponds to a chunk in the user address space (i.e., *user chunk number*). It contains an 8-bit *slot number* and a 32-bit *SSD LPN*. The former represents the slot in the stripe where the chunk is placed. The latter records the SSD logical address of the chunk's primary copy (i.e., the copy placed in HH, cf. § 5.3). For example, user chunk 0x0 is placed in slot 0 of the first stripe in AFA address space, and its primary copy is stored in 0x20 of SSD 0 (cf. Figures 7 and 8). Since logical pages are one-to-one mapped to physical pages in SSD, we slice PMT based on *SSD LPN* and store it in the corresponding OOB. Specifically, we piggyback the user chunk number and slot number in write requests and send them to SSDs with a 32-bit timestamp. When writing data to flash physical pages, the SSD controller persists the metadata to OOB via the same program operations (i.e., without extra write). To guarantee fault tolerance of PMT, we also compute parity codes for metadata. Specifically, in the conversion phase, the conversion module reads both data and metadata from the same physical page with a single read operation. It then computes parity codes for them simultaneously. Finally, the parity codes of data and metadata will be written back together and stored in the same page (cf. § 5.3). Note that our solution only utilizes 72 bits *spare* space of OOB (cf. § 2.1) per chunk to persist our customized metadata, which does not affect other OOB-based functions (e.g., ECC).

Recovery of mapping table. Take Figure 8 as an example to depict the recovery AMT and HGMT. For simplicity, we assume the chunk size, flash page size, and transient area size of each SSD are 64, 4, and 128 KB. First, we scan OOB to recover PMT. However, this process can be time-consuming as it needs to scan all flash pages. We accelerate this process with periodical (e.g., every 2 seconds) checkpoints [80] and only scanning the un-checkpointed metadata. From PMT, we know that the primary copy of user chunk 0 is stored at 0x20 of SSD 0, which is the first chunk of the normal area (i.e., HH #0 in Figure 7). Therefore, user chunk 0 is stored in the first stripe of AFA address space. Considering that the slot number of user chunk 0 is 0, the chunk is stored in the first chunk of AFA address space, i.e., its AFA chunk number is 0. Repeating this, we restore AMT. As user chunks 0 and 1 are stored in SSD 0 and 1, we know that their parity chunk is stored in SSD 2 (i.e., SSD ID). We

then read the OOB from 0x20 of SSD 2 and verify if it is the parity. If so, this stripe is in the striping layout and the corresponding item in HGMT is NULL. Otherwise, we need to scan the transient area of SSD 2 and match the OOB with user chunks 0 and 1 to locate their VH (i.e., VH number). Thereby, we recover HGMT.

Recovery of space manager and I/O threads. Another critical metadata in ScalaAFA are the states of space manager (i.e., Free HH List, Free VG Lists, and Full HG Lists, cf. Figure 6) and I/O threads (i.e., granted heroes). We consider the recovery of them in three corruption conditions. First, both the space manager and the I/O threads are corrupted (e.g., system reboot). In this condition, all granted heroes are nullified. After reboot, the space manager can restore its states by scanning the mapping tables to find out: (a) Whether VHs/HHs are free. If so, their corresponding AFA chunk numbers in AMT are NULL; (b) Whether HGs need conversion. If so, the corresponding items in HGMT are not NULL and all chunks in the HG are not free. Note that, partially-used HG can be reclaimed in AFA-level GC. Second, only the space manager is corrupted. In this condition, I/O threads can continue to serve requests with the already granted heroes but fall asleep after heroes run out. When the space manager restarts, it restores the states following the procedure in the first condition and communicates with I/O threads to confirm which heroes have been granted; Third, only the I/O thread corrupts. The heroes that have been granted to this thread become unavailable. Afterward, when AFA is idle, the space manager reclaims the space by following the procedure in the second condition (i.e., scanning mapping tables and communicating with active I/O threads).

5.5 Reduce the Impact of Write Amplification

While our conversion offloading has eliminated extra I/Os in the conversion phase (cf. § 5.3), the replication phase still causes $m + 1$ times write amplification, which significantly shortens the SSD lifetime. Our key insight is that backup copies (i.e., data in VHs) in two-phase write are transient and will be invalidated soon in the conversion phase. Thus, there is no need to flush them from the durable SLC buffer to the vulnerable MLC blocks. Unfortunately, SSDs are unaware of which data are transient. Tackling this issue, we propose a holistic design. Specifically, when writing backup copies to SSDs, ScalaAFA marks them as transient data by flagging one bit in the write requests. Once receiving these requests, SSD records the transient tags. Afterward, when the SLC buffer is full, the SLC evictor selectively evicts SLC blocks by considering the transient tags. To be specific, the SLC evictor will evict the SLC block with the highest evicting score, which can be calculated with the following formula:

$$score = factor * num_invalid - num_transient$$

$num_invalid$ and $num_transient$ are the numbers of invalid and transient pages in this block, while $factor$ is a constant. We set it as 1 empirically in this paper. If one block has more invalid pages, we evict it in higher priority, because it generates fewer writes in this eviction. On the other hand, if one block contains more transient pages, we evict it in lower priority. Since these pages will be invalidated soon, there is no need to flush them to MLC blocks immediately.

5.6 Mitigate the Latency Spikes Caused by GC

To mitigate the performance degradation caused by SSD GC, we employ two designs: *GC avoidance* and *GC detection*, as described in Algorithms 1 and 2, respectively.

GC avoidance. In traditional AFAs, GC avoidance is non-trivial as a write request needs to be served with multiple SSDs. If any SSD starts GC, AFA cannot complete write requests in time. Fortunately, ScalaAFA can complete writes (both in replication and conversion phases) with only part of the SSDs by leveraging its unique data placement policy (cf. § 5.3). Thus, ScalaAFA always chooses the SSDs that are not in GC to serve I/Os. To be specific, in the replication phase (cf. Lines 1-13 in Algorithm 1), I/O threads stop writing data to the HGs, where their VHs are located at the SSD in GC. Instead, they assemble new HGs that only include VHs coming from

Algorithm 1: GC avoidance.

```

// Get HG for writes in the replication phase (cf. § 5.2).
1 Function getHG():
2    $HG_{tgt} \leftarrow NULL$ ;
3   foreach Assembled  $HG_i$  do
4     if all VHS in  $HG_i$  are not located at the SSD in GC then
5        $HG_{tgt} \leftarrow HG_i$ ;
6       break;
7     end
8   end
9   if  $HG_{tgt} = NULL$  then
10     $HG_{tgt} \leftarrow$  assemble new HGs that only include VHS coming from normal SSDs;
11  end
12  return  $HG_{tgt}$ ;
13 end

// Choose a HG to convert in the conversion phase (cf. § 5.2).
14 Function chooseConvert():
15    $HG_{tgt} \leftarrow NULL$ ;
16   foreach  $HG_i$  in full HG list do
17     if all VHS in  $HG_i$  are not located at the SSD in GC then
18        $HG_{tgt} \leftarrow HG_i$ ;
19       break;
20     end
21   end
22   if  $HG_{tgt} = NULL$  then
23      $HG_{tgt} \leftarrow$  head of the full HG list;
24   end
25   return  $HG_{tgt}$ ;
26 end

```

normal SSDs (i.e., not in GC). ScalaAFA then writes data chunks to these new HGs and uses the slots that are mapped with normal SSDs. In the conversion phase (cf. Lines 14-26), ScalaAFA will convert the VHS placed in normal SSDs first, thereby freeing the transient area for replication and continuing the blocked user I/O quickly.

GC detection. Since flash firmware schedules every GC event, we adopt a holistic design that leverages this feature to detect GC in both passive and proactive ways. Passively, when the flash firmware schedules GC, the GC reporter (cf. Figure 4) encapsulates this information into the completion queue entry (CQE) of the following NVMe requests (cf. Lines 1-6 in Algorithm 2). Therefore, the GC detector (cf. Figure 4) in the host is informed that the SSD will start GC. It then broadcasts this information to all I/O threads (cf. Lines 9-13), which thereafter redirects the I/O requests to normal SSDs via GC avoidance. However, since I/O threads stop sending I/O requests to this SSD, the SSD cannot piggyback signals to GC detector after GC completion, which causes under-utilization of this SSD. To address this, GC detector periodically sends *GC probe commands* to the SSD in GC, which checks if GC is done (cf. Line 14). If so, I/O threads can resume accessing this SSD.

Algorithm 2: GC detection.

```

// Codes of the GC reporter (in the flash firmware).
1 Function reporter():
2   while true do
3      $State_{GC} \leftarrow$  if flash firmware schedules GC then 1 else 0;
4     Encapsulates  $State_{GC}$  into the CQE of the next NVMe request;
5   end
6 end

// Codes of the GC detector (in the host).
7 Function detector():
8   while true do
9     Unwrap  $State_{GC}$  from the CQE of NVMe request  $Q$ ;
10    if  $State_{GC} = 1$  then
11       $N \leftarrow$  the unique number of the target SSD for the request  $Q$ ;
12      Inform all I/O threads that SSD  $N$  is in GC;
13    end
14    Periodically send GC probe commands to SSD in GC;
15  end
16 end

```

Host System		Femu		Software	
CPU	Intel Xeon 5320	Virtual machine	32 CPU threads	Linux kernel	v5.11.0
	1×26 Core / 2.2 GHz with hyper-threading		32 GB DRAM		
Mem.	8×64 GB / DDR4	Flash	8 Channel / 12 Die / 1 Plane / 352 Block / 512 Page / 4 KB	fio	v3.30
SSD	Samsung 980 Pro			perf	v5.11
	R/W: 7000/5200 MB/s	Bw.	R/W: 7500/4890 MB/s	mdadm	v4.1
				SPDK	v20.05

Table 1. System configurations.

5.7 Implementation

ScalaAFA employs a daemon thread to play the role of space manager. When creating an AFA, we initialize the daemon thread and bind it to a fixed CPU core by `spdk_env_thread_launch_pinned()`. All of our modifications to NVMe protocols are based on NVMe Base Specification 2.0c [6] and NVMe Command Set Specification 1.0c [7]. The conversion and GC probe commands are implemented as an NVMe IO command and an NVMe admin command, respectively. The transient tag uses one reserved bit in NVMe write command. The GC-related information consumes one reserved bit in CQE. We piggyback the user chunk number, lot number, and timestamp to SSDs and program them in OOB with the support of NVMe Protection Information feature [7]. We develop ScalaAFA in SPDK v22.05 [9] with 6K LOC. Modification of the SSD firmware is implemented in a popular SSD emulator [50] with 1K LOC.

Trace	Wr. cnt. (Kops)	Rd. cnt. (Kops)	Avg. wr. len. (KB)	Avg. rd. len. (KB)	Data wr. (GB)	Data rd. (GB)
proxy0	12135.4	383.5	4.6	8.3	53.8	3.1
prn	7753.0	9066.3	10.4	22.5	76.8	194.5
src2	2201.1	1171.3	23.2	54.8	48.8	61.2
CFS	1173.0	3304.1	12.6	8.7	14.1	27.3
DAP	475.3	610.4	97.2	62.1	44.1	36.2
webmail	6381.9	1413.8	4.0	4.0	24.3	5.4

Table 2. Characteristics of real workloads.

6 Evaluation

6.1 Experimental Setup

Methodology. We use Femu [50], a QEMU-based SSD emulator, to evaluate our holistic designs. We set up the QEMU virtual machine to run on Linux v5.11.0 with 32 CPU threads, 32 GB DRAM, and up to 8 NVMe SSDs. We configure the emulated SSDs as high-performance storage devices with 7500 MB/s and 4890 MB/s peak read and write bandwidths, respectively. We set the size of the SLC buffer to 4 GB. These configurations match with commercial high-end SSD products [3]. In addition, our simulator employs an XOR engine in the SSD controller, which is the same as the state-of-the-art SSD devices [16, 17, 73]. We get configurations of the XOR engine from Xilinx xc7a200t FPGA with 200 MHz clock. It takes 20 μ s and 16 mW dynamic power to calculate a 64 KB parity chunk from 6 data chunks. The conversion of a stripe in SSD costs 103 μ s in total. The key configurations in our experiments are listed in Table 1.

AFA platforms. We compare ScalaAFA with five other popular AFA engines. (1) mdraid [5]: the default stripe write AFA engine implemented in Linux kernel; (2) ScalaRAID [79]: a state-of-the-art stripe write AFA engine, which mitigates the software overheads of mdraid with fine-grained locks and improves its performance to some extent; (3) stRAID [69]: another stripe write AFA engine, which alleviates the software overheads in mdraid with a run-to-complete I/O processing scheme; (4) RAID5F [8]: an *incomplete* stripe write AFA engine in SPDK, which has no crash consistency and *can only serve RAID 5 full-stripe I/O* (i.e., the I/O size is equal to the stripe size). We consider the performance of RAID5F is close to the ideal case of software-only stripe write AFA designs; (5) FusionRAID [39]: a state-of-the-art two-phase write AFA engine; (6) ScalaAFA: the user-space AFA engine that includes all the designs proposed in this paper. We set the chunk size in all the AFA platforms as 64 KB.

Workloads. We evaluate ScalaAFA with various benchmark suites and applications. Specifically, we measure the performance of different AFA engines by employing fio v3.30 [14] to execute microbenchmarks. By default, we use a single I/O thread to generate asynchronous I/O requests. We also consider the multi-thread scenarios in scalability testing (cf. § 6.2). To reflect the impacts of block devices on system performance, we evaluate ScalaAFA with block I/O traces [48]. We select six representative workloads from both industries (Microsoft MSRC [41] and MSPC [61]) and academia (FIU [67]) including both write-dominated (e.g., proxy0) and read-write mixed (e.g., DAP) traces with varied I/O sizes. Table 2 summarizes the key characteristics of the selected workloads. We also conduct a comparison on RocksDB [26], a popular KV store, with db_bench [27], which demonstrates the end-to-end performance improvement brought by our designs. Unless specified, we configure the GC threshold to a high value (i.e., 95% of SSD capacity) so that SSD GC is not triggered in these tests. We will evaluate the impacts of GC in § 6.5.

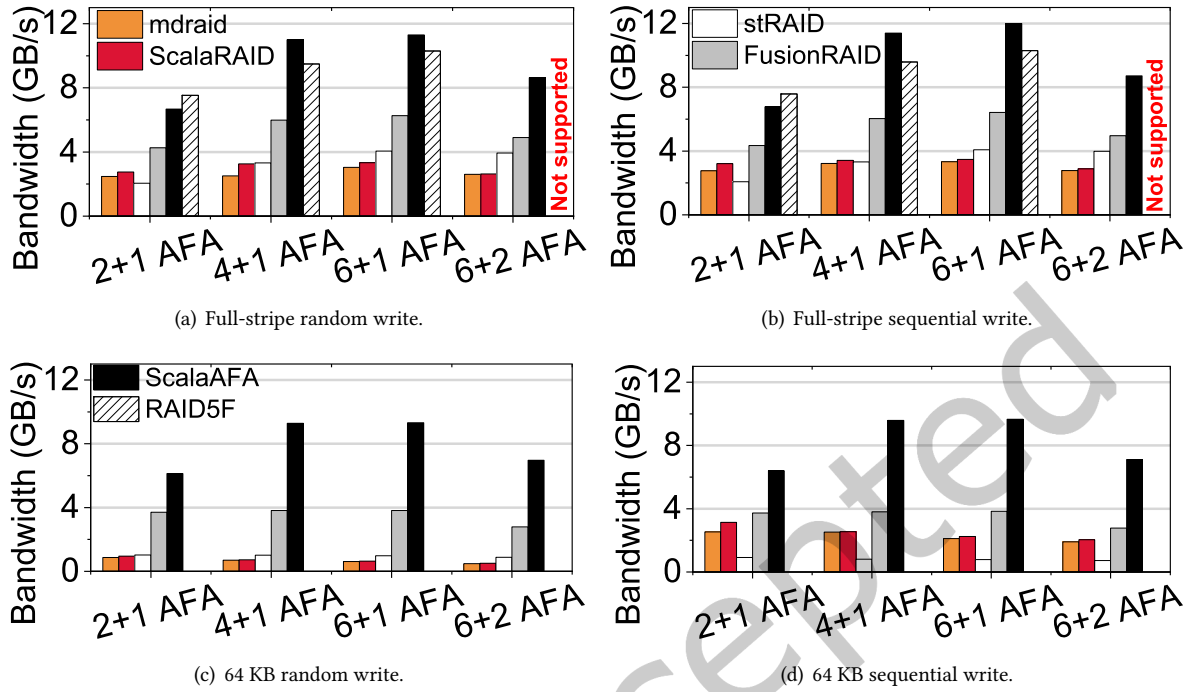


Fig. 9. Comparison of write throughput on microbenchmarks.

6.2 Overall Performance

Throughput. We compare the write throughput of different AFA engines, which is shown in Figure 9. We set the I/O depth to 32. ScalaRAID slightly outperforms mdraid. This is because ScalaRAID aims at mitigating the overheads of lock mechanism. However, such overheads are not severe when employing only one I/O thread (cf. § 3.1). stRAID also achieves 17.9% higher full-stripe write throughput than mdraid, on average, thanks to its run-to-complete I/O processing scheme. As shown in Figure 9(c), mdraid, ScalaRAID and stRAID all have poor performance in 64 KB random write accesses because stripe write AFA are unfriendly to partial write. Since mdraid and ScalaRAID opportunistically aggregate multiple contiguous write requests as a full-stripe write by employing a DRAM cache [5, 69], they achieve higher performance for 64 KB sequential write than stRAID (cf. Figure 9(d)). FusionRAID outperforms stRAID by 3.6× in 64 KB random write, on average. This is because it absorbs partial write requests with replication. Also, it updates data chunks in an out-of-place way, which avoids the tedious read-construct-write procedure. ScalaAFA outperforms FusionRAID in all types of I/O access patterns. For example, in 64 KB sequential write, ScalaAFA further improves the write throughput by 1.6×, 2.4×, 2.4×, and 2.5× in 2+1, 4+1, 6+1, and 6+2 AFAs, respectively. This is because the key techniques in ScalaAFA, such as the user-space design and conversion offloading, significantly reduce the CPU burdens. RAID5F outperforms ScalaAFA by 12.4% in 2+1 AFA, but has worse performance (17.4% and 13.1%) in 4+1 and 6+1 AFAs. This is because conversion overhead is minor in 2+1 AFA whereas becoming severe as the number of SSDs increases. With our hardware/software co-designs, ScalaAFA not only offloads parity computation to SSDs but also eliminates the penalty of data preparation (cf. §5.3).

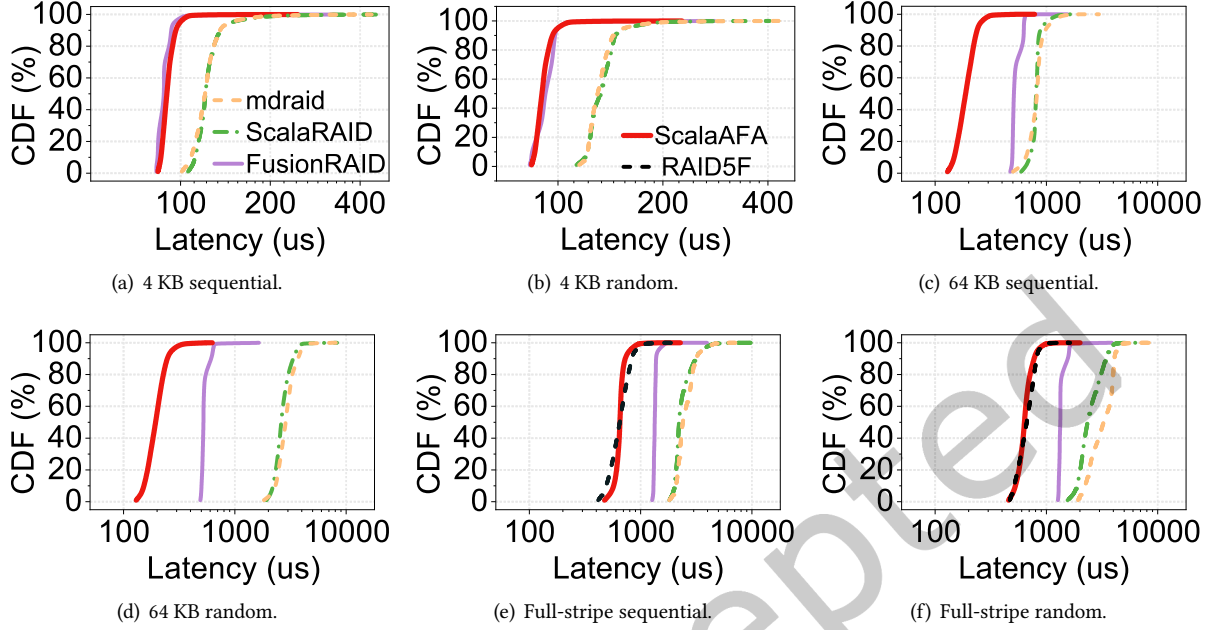


Fig. 10. Comparison of write latency CDF on microbenchmarks.

Latency. Figure 10 shows the cumulative distribution function (CDF) of all the tested AFA engines in terms of latency. We omit stRAID in this comparison for the implementation bugs in its open-sourced codes. We select 4+1 AFAs in this test. We set the I/O depth to 1 for 4 KB write while it is 32 for both 64 KB and full-stripe write. There is no visible difference between ScalaRAID and mdraid in the 4 KB write scenario. However, ScalaRAID reduces the average and 99th percentile latency by 22.9% and 15.0% for full-stripe random write, where the lock overhead is not negligible. FusionRAID outperforms mdraid in terms of both average and tail latencies. For example, it reduces the average and 99th percentile latency to 66.1% and 48.2%, respectively, for 64 KB sequential write. ScalaAFA shows significantly better CDF profiles in all scenarios. Compared with FusionRAID, ScalaAFA further reduces the average and 99th percentile latency by 52.7% and 41.9% for full-stripe write, respectively. This is because ScalaAFA not only benefits from the high-performance user-space storage stack but also mitigates performance degradation caused by conversion with SSD architectural innovations. Although RAID5F is considered an ideal software-only AFA solution, ScalaAFA even achieves 4.8% shorter 99th percentile tail latency than RAID5F thanks to our exploration of hardware resources.

Scalability. We further measure the scalability of different AFA engines by varying the number of I/O threads from 1 to 12, which is shown in Figure 11. For mdraid, ScalaRAID, and stRAID, we employ the same number of worker threads as I/O threads (as suggested by prior work [69, 79]). We also plot the ideal throughput of two-phase write AFA engines (i.e., $(k + m)/(m + 1) * S$, where S is the peak bandwidth of a single SSD).

mdraid gains limited benefits from extra threads. It achieves peak performance when employing 8 I/O threads. This is because lock overhead caused by multi-thread access dominates the time cost of the storage stack [79]. Benefited from the lower software overhead (i.e., fine-grained locks [79] and run-to-complete I/O processing scheme [69]), ScalaRAID and stRAID are more scalable than mdraid. They achieve 11.0 GB/s and 11.3 GB/s, respectively, in 4+1 AFA with the cost of 12 I/O threads. In comparison, ScalaAFA can achieve 11.4 GB/s with

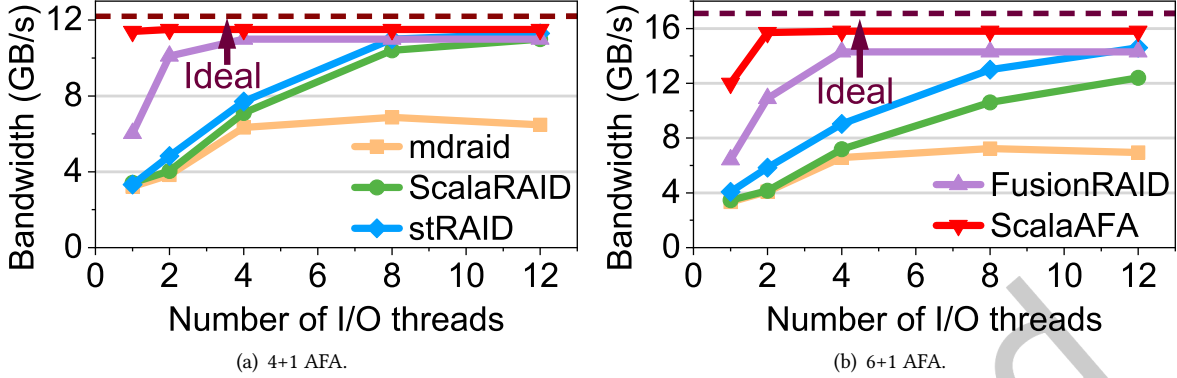


Fig. 11. Comparison of scalability.

only one I/O thread, which is almost the same as the ideal case for 4+1 AFA. The minor performance difference between ScalaAFA and the ideal case comes from the background tasks within SSDs (e.g., conversion). For 6+1 AFA, FusionRAID achieves only 6.4 GB/s with one I/O thread, 37.5% of the ideal performance. This is because FusionRAID suffers from the huge CPU burden caused by the tedious storage stack (cf. § 3.1). ScalaAFA improves the throughput to 12.0 GB/s and 15.7 GB/s when employing 1 and 2 I/O threads, respectively. This is because ScalaAFA imposes minor overhead on the host CPU, thereby allowing I/O threads to fully exploit the bandwidth of SSDs.

6.3 Analysis of Real Workloads

Throughput. Figure 12 illustrates the throughput of different AFA engines in different real workloads. We set the I/O depth to 32 in this test. Compared with mdraid, FusionRAID achieves 34.1% higher throughput in all workloads, on average. Based on FusionRAID, ScalaAFA further improves the average throughput by 2.9×, 2.2×, and 1.2× in workloads prn, src2, and DAP, respectively. Moreover, it shortens the completion time of all workloads to 30.0%, on average, compared with FusionRAID. Workload DAP has the largest average write request size, in which ScalaAFA achieves the highest average throughput (i.e., 4.8 GB/s) among all workloads. All AFA solutions exhibit the lowest throughput in webmail. This is because webmail has the smallest write request size, which cannot fully utilize the SSD-internal parallelism. However, the bandwidth of ScalaAFA still exceeds that of FusionRAID by 2.8×.

Latency. Figure 13 shows the latency CDFs in different real workloads. Compared with mdraid, FusionRAID reduces the average and 99th percentile latency by 13.8% and 59.7%. In all workloads, ScalaAFA achieves the best CDF profiles. For example, ScalaAFA reduces the average latency to 24.2% and 26.9% in proxy0 and prn, compared to FusionRAID. ScalaAFA reduces the 99th percentile latency of FusionRAID by 59.1%, 65.1%, and 72.5% in the src2, CFS, and webmail. This is because ScalaAFA not only offloads conversion to SSDs that eliminates background I/O but also benefits from the lock-free user-space designs.

6.4 End-to-end Evaluation

To demonstrate the superiority of ScalaAFA on applications, we conduct an end-to-end evaluation on RocksDB [26], a popular KV store. We run RocksDB on Ext4 and BlobFS [66] file systems for kernel-space (i.e., mdraid and FusionRAID) and user-space (i.e., ScalaAFA) AFA engines, respectively. We use five representative benchmarks from db_bench including both write-dominated (e.g., fillrandom) and read-write-mixed (e.g., fillseekseq)

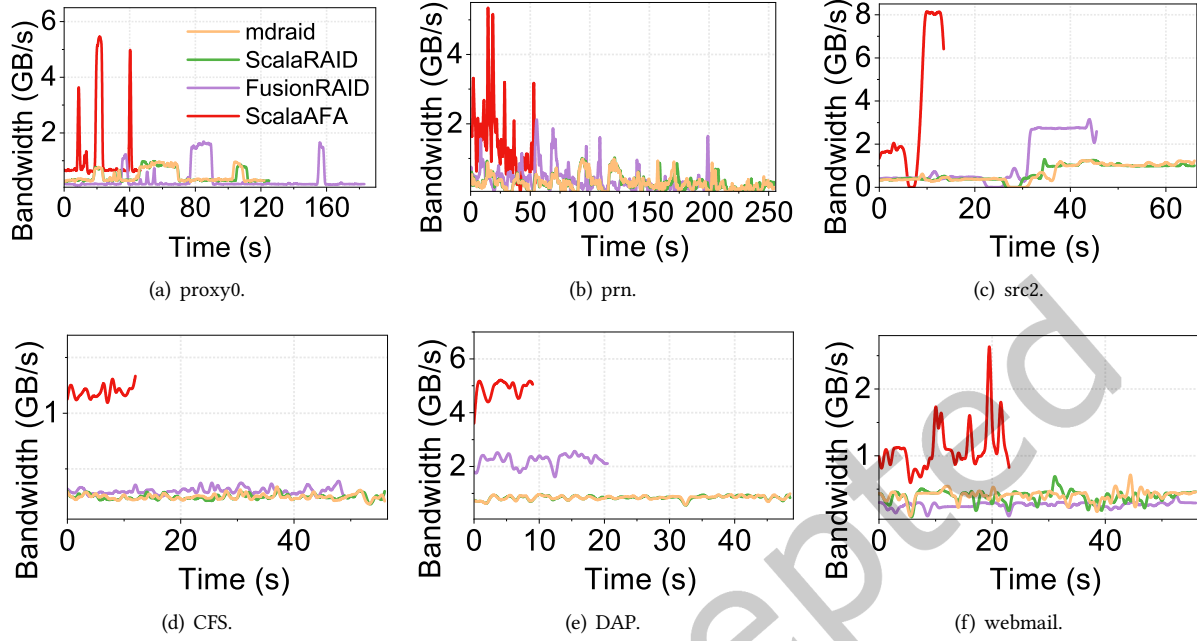


Fig. 12. Comparison of write throughput on real workloads.

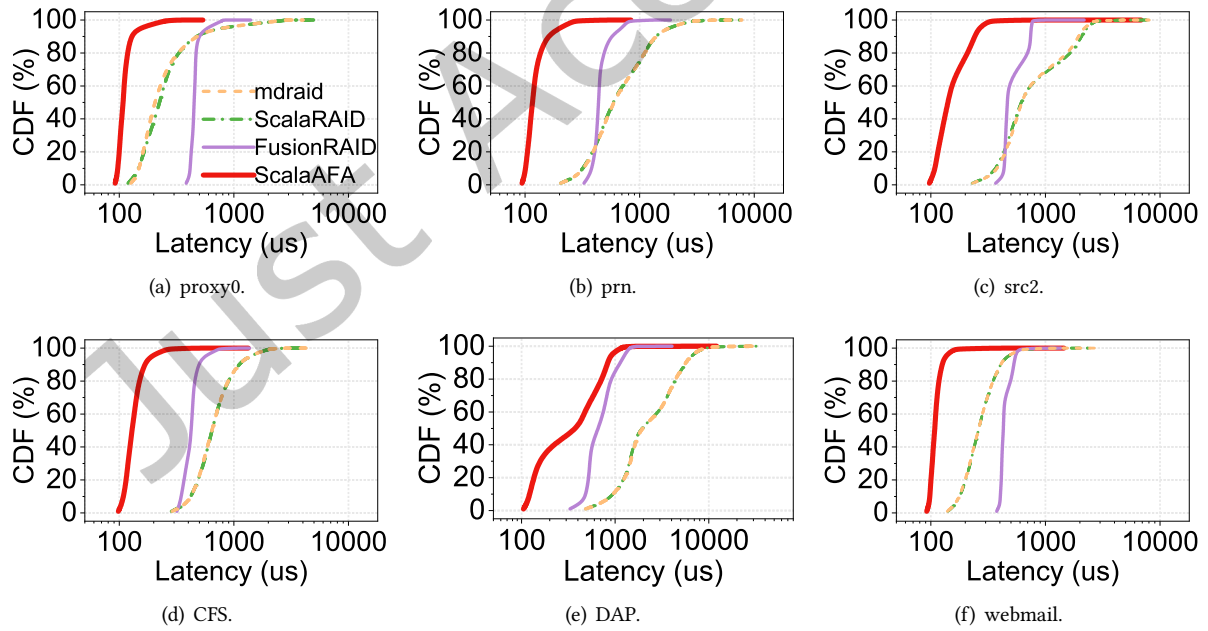


Fig. 13. Comparison of write latency CDF on real workloads.

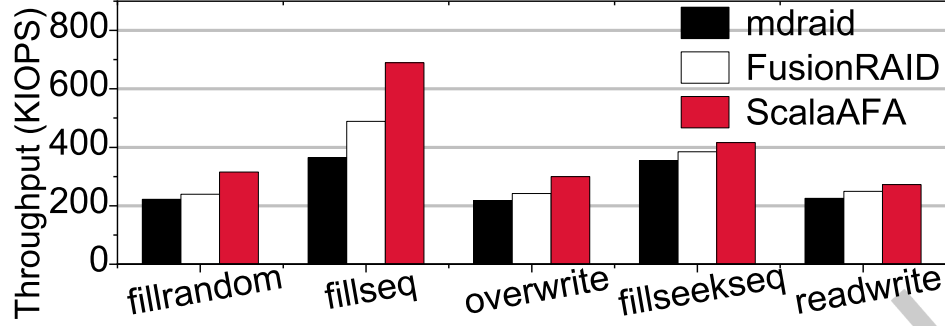


Fig. 14. Throughput comparison in RocksDB.

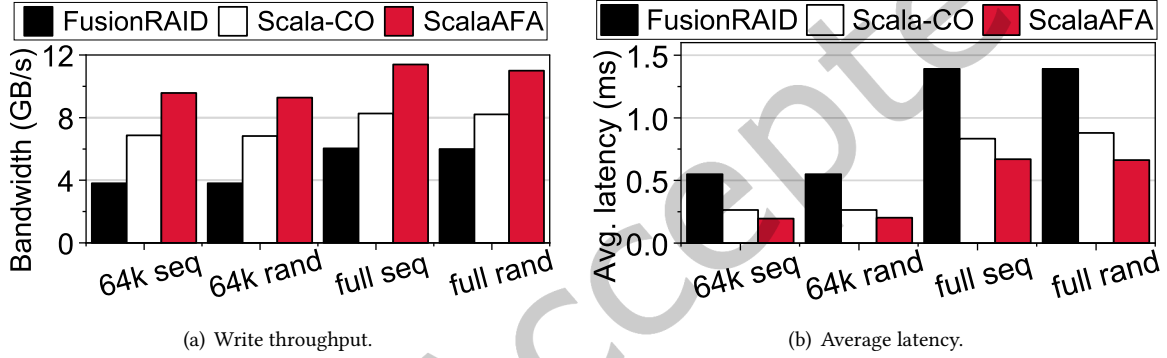


Fig. 15. Sources of performance improvement when no GC happens.

workloads. We set the key and value sizes as 16 B and 1 KB, respectively. Figure 14 illustrates the throughput comparison of different AFA engines. Thanks to the utilization of SSD-internal resources and the high-performance storage software stack, ScalaAFA achieves the highest throughput in all workloads. Specifically, ScalaAFA outperforms mdraid by 41.4% in all tests, on average. Compared with FusionRAID, ScalaAFA improves the average throughput by 31.9%, 41.1%, and 23.8% in fillrandom, fillseq, and overwrite, respectively. The improvement decreases to 8.1% in fillseekseq. This is because, in this workload, lots of I/O requests are absorbed by the DRAM cache [20] of RocksDB without entering AFAs.

6.5 Benefits of Individual Techniques

In this section, we evaluate the benefits brought by each key design in ScalaAFA. We implement three new AFA engines, which incorporate only parts of ScalaAFA techniques. (1) Scala-CO: based on ScalaAFA, we execute the conversion in the host. Specifically, in the conversion phase, the daemon thread reads the data chunks from the SSDs, computes the parity chunks, and finally writes them back to SSDs; (2) Scala-WR: based on ScalaAFA, we do not set the replicated chunks as low priority; (3) Scala-GC: based on ScalaAFA, we remove the GC detection and avoidance mechanism. We use 4+1 AFA in this evaluation.

Lock-free user-space design. Figure 15 shows the write throughput and latency of different AFA solutions when no GC happens. Compared with FusionRAID, Scala-CO reduces the average latency by 45.2% and improves

Metatdata type	Mapping tables	Space manager	I/O thread
Recovery time (ms)	1091	201	227

Table 3. Time consumption of metadata recovery.

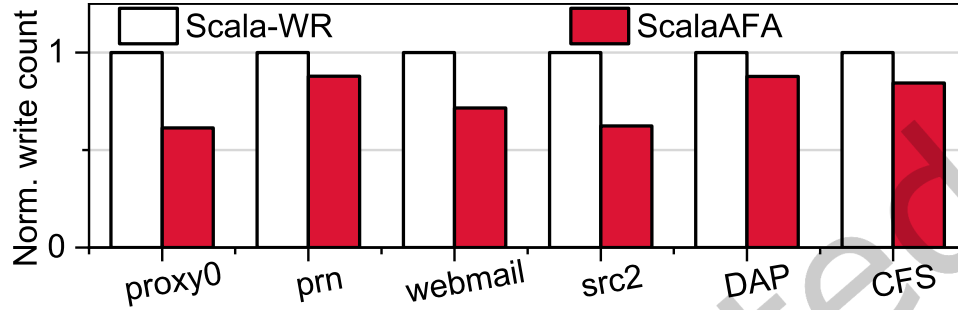


Fig. 16. Analysis of write amplification.

the write throughput by 58.4%, for all scenarios on average. Note that, FusionRAID [39] recommends persisting metadata with host-side battery-backed DRAM, which is impractical for productions considering its monetary cost and insufficient fault tolerance. In our reproduced version, FusionRAID employs the same metadata persistence scheme as ScalaAFA. Therefore, the performance differences between FusionRAID and Scala-CO mainly come from our lock-free user-space design (i.e., message-passing-based permission management mechanism).

Conversion offloading. As shown in Figure 15, ScalaAFA further reduces the average latency by 36.1%, in comparison to Scala-CO. In terms of write throughput, ScalaAFA improves the performance by 39.6%, 36.1%, 37.9%, and 34.0% for 64 KB sequential and random write, full-stripe sequential and random write, respectively. This is because the conversion takes up 77.6% of the CPU ticks of the daemon thread, which significantly disrupts the space manager running on the daemon thread thereby blocking user I/Os. In contrast, ScalaAFA achieves higher throughput in all testing scenarios thanks to the conversion offloading design.

OOB-based crash consistency. To mitigate the software overhead of metadata persistence, we suggest an OOB-based crash consistency mechanism (cf. § 5.4). Note that SSDs can write OOB area by hitchhiking the same flash programming operation of data writes. Therefore, this method consumes no extra time for metadata persistence. Moreover, the periodical checkpoint (cf. § 5.4) can significantly reduce the time consumption of metadata recovery. As listed in Table 3, ScalaAFA can restore its metadata in about 1 second.

Write amplification reduction on MLC blocks. To evaluate how much our design can mitigate the impact of write amplification, we set the size of the SLC buffer to 4 GB and count the number of writes on MLC blocks. The results are shown in Figure 16. Compared to Scala-WR, which performs the same as existing two-phase write engines [28, 39, 71], ScalaAFA can decrease write count in all workloads. For example, it alleviates the write amplification by 38.6%, 12.2%, 37.7%, and 28.4% in proxy0, prn, src2, and webmail, respectively. This is because, with our holistic write amplification reduction scheme, ScalaAFA evicts the replicated data to MLC blocks in a low priority as these data will be invalidated right after the conversion phase. Therefore, ScalaAFA reduces the amount of data written to the vulnerable MLC blocks.

GC detection and avoidance. We evaluate the impact of the GC detection and avoidance mechanism with both microbenchmark (i.e., 64 KB sequential write) and real workloads. Since CFS only writes 14 GB data and does not trigger SSD GC, we omit the result of CFS. Figures 17(a) and 17(b) show the comparison of 99.99th

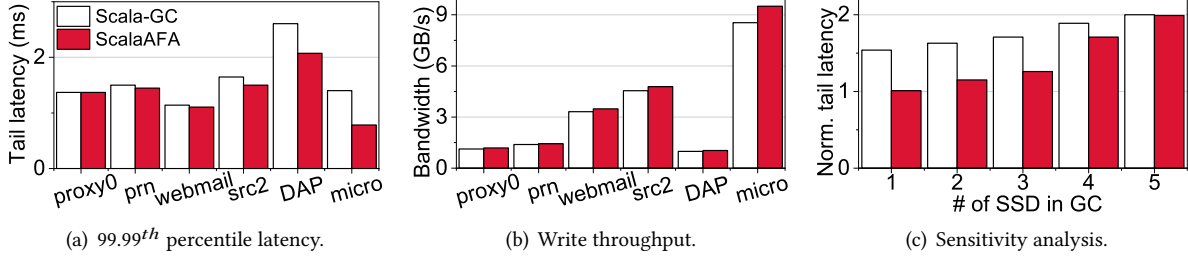


Fig. 17. Improvement brought by GC detection and avoidance.

Trace	proxy0	prn	webmail	src2	DAP	micro
Ratio (%)	0.003	0.05	0.001	0.40	0.74	4.00

Table 4. Fraction of I/O requests influenced by GC.

percentile latency and throughput, respectively. We also summarize the fractions of I/O requests influenced by SSD GC in Scala-GC, listed in Table 4. Compared with Scala-GC, ScalaAFA decreases the tail latency by 8.9% and 20.5% in workloads src2 and DAP (cf. Figure 17(a)). This improvement becomes 44.1% in microbenchmark. This is because microbenchmark generates more intensive write bursts, which trigger SSD GC frequently (cf. Table 4). Since our design can redirect I/O to SSD that is not in GC, it eliminates the penalties from stragglers and thereby reduces tail latency. This superiority also exists in throughput comparisons (cf. Figure 17(b)). For example, ScalaAFA outperforms Scala-GC by 11.3% in microbenchmark. On the contrary, in proxy0, prn, and webmail workloads, only small ratios of I/O requests are influenced by SSD GC (cf. Table 4). Therefore, our GC detection and avoidance mechanism can merely deliver minor improvement. Figure 17(c) presents the 64 KB write latency comparison when there are different numbers of SSDs in GC. We manipulate the number of SSDs in GC by pre-filling certain SSDs before each test, thereby triggering GC. We normalized the results to the ideal case (i.e., no GC occurs). Thanks to our GC-related designs, ScalaAFA surpasses Scala-GC in all the tests, achieving a 19.9% latency reduction on average. Interestingly, ScalaAFA succeeds in mitigating GC-caused latency spikes when the number of SSDs in GC is less than $(k + m) - (m + 1)$ (i.e., 3), while delivering minor benefits when an excessive number of SSDs are in GC (e.g., 4 and 5). This is because, although ScalaAFA can redirect write requests to normal SSDs (i.e., not in GC), it requires at least $m + 1$ SSDs for replication. Therefore, when the number of normal SSDs falls below $m + 1$, ScalaAFA cannot avoid GC-induced performance degradation.

Summary. To sum up, our user-space lock-free designs (cf. § 5.2) deliver up to 58.4% write throughput improvement. The conversion offloading design (cf. § 5.3) further increases write throughput by 39.6%. The OOB-based crash consistency mechanism (cf. § 5.4) can persist metadata with negligible overhead. In addition, our SLC evictor (cf. § 5.5) succeeds in mitigating write amplification by up to 38.6%. Finally, the GC detection and avoidance (cf. § 5.6) mechanism can reduce up to 44.1% tail latency.

7 Related Work and Discussion

Overheads of AFA engines. Multiple studies [33, 34, 38, 39, 53, 54, 69, 79] have been proposed to mitigate the software overheads of AFA engines. ScalaRAID [79] and stRAID [69] partially mitigate the overhead of multi-thread access with fine-grained lock schemes and distributed data structures. WAFL [34] integrates AFA engine in a specific file system and mitigates partial write overhead with its write-anywhere file layout. In comparison,

Challenge	Challenge 1		Challenge 2	Challenge 3	Challenge 4	Challenge 5
Overhead	Context switch	Lock	Conversion	Metadata persisting	Write amplification	SSD GC
mdraid	✗	✗	✗	✗	✗	✗
ScalaRAID	✗	✓	✗	✗	✗	✗
stRAID	✗	✓	✗	✗	✗	✗
EAR	✗	✗	✓	✗	✗	✗
FusionRAID	✗	✗	✗	✗	✗	✓
ScalaAFA	✓ (cf. § 5.2)	✓ (cf. § 5.2)	✓ (cf. § 5.3)	✓ (cf. § 5.4)	✓ (cf. § 5.5)	✓ (cf. § 5.6)

Table 5. Overall comparison across different AFA solutions.

ScalaAFA is a general AFA engine that functions as a standard block device and can adapt to more file systems. Moreover, WAFL relies on battery-backed DRAM to persist metadata, which increases monetary costs and cannot provide enough fault tolerance (i.e., data cannot be located after the DRAM fails). Instead, ScalaAFA leverages SSD-internal hardware resources (i.e., OOB) for crash consistency. FusionRAID [39] reduces the overhead of partial write with two-phase write scheme. However, the overheads of parity computation and background I/O (e.g., conversion) still exist in the prior work and impose a huge burden on the host CPU, especially in I/O-intensive scenarios. In contrast, the holistic designs in ScalaAFA not only alleviate the lock overheads but also eliminate the conversion penalty. Moreover, our hardware-based metadata persistence scheme can further accelerate AFA. EAR [53] tries to reduce the data reads/writes in conversion phase with a sophisticated flow graph matching algorithm. However, EAR still computes parities out of place and requires extra communications among storage nodes for forwarding computing results (i.e., map-reduce). Considering the lack of proactive communication capability between SSDs, it is hard to adopt EAR for AFA scenario. In comparison, ScalaAFA can generate parities locally, which completely eliminates the communications among SSDs.

Alleviating the impact of GC. Multiple prior work [17, 23, 32, 39, 43–45, 49, 51, 54, 75] has extensively studied the performance degradation caused by SSD GC in AFA. SWAN [44] classifies SSDs as foreground and background groups. It eliminates GC penalties by serving writes with only the foreground group where no SSD is in GC. SWAN suffers from low throughput since it underutilizes normal SSDs (i.e., not in GC) in background groups. In contrast, ScalaAFA handles I/O requests with every normal SSD to fully utilize the aggregated throughput of AFA. FusionRAID [39] mitigates the interference of SSD GC by employing a software GC detector and redirecting the I/O to other normal SSDs in the same storage pool. However, the software detector is passive. When the SSD GC can be detected at the software level, the SSD GC has already delayed I/O requests. ScalaAFA detects SSD GC with a holistic design, which is aware of the occurrence of GC in advance and avoids its impact. TolerAID [32] and IODA [51] reconstruct delayed read data from parities such that they can mitigate the GC-caused read latency spikes. In contrast, ScalaAFA focuses on optimizing the write path and is orthogonal to the approach.

ScalaAFA with emerging NICs. Some emerging NICs [2, 60] enable direct access to NVMe SSDs, bypassing the host CPU and saving precious PCIe bandwidth. ScalaAFA is orthogonal to this advanced hardware. In this condition, ScalaAFA can run on client servers and access remote NVMe SSDs via NVMe-oF protocols. The only difference is that we need to modify the NVMe driver integrated in NICs rather than in the SPDK or kernel. Moreover, in this scenario, the overhead of data read and write in the conversion phase can be further magnified, considering the extra latency of network communication. Therefore, our conversion offloading design can even achieve more improvements.

In-storage processing. As high-performance SSDs usually equip abundant computing resources (e.g., ARM processors and XOR engine), much prior work [24, 40, 46, 47, 55, 56, 62] proposes to offload tasks to SSDs. HolisticGNN [46] offloads GNN operators to SSDs. Cognitive SSD [56] accelerates deep learning with on-device resources. Willow [62] extends the same idea to a general-purpose framework, which allows applications to offload data-intensive operands to the underlying SSDs. ScalaAFA extends this idea by taking the system overheads of the AFA scenario into consideration.

Comparison summary. Taking the discussion of related work into consideration, we summarize the key differences between prior work and our proposed ScalaAFA in Table 5.

Trade-off in ScalaAFA. Compared to traditional software-only designs, the main obstacle of implementing ScalaAFA is the modification of SSD firmware. Limited by SSD vendors, it is difficult to achieve this in off-the-shelf SSDs. However, we believe in its prospects, given the significant benefits brought by our holistic designs (cf. § 6) and the emergence of in-storage processing.

8 Conclusion

Existing AFA engines fail to adopt next-generation high-performance SSDs because of the huge software overhead caused by the storage stack and AFA internal tasks. Tackling this issue, we propose a new user-space AFA engine with holistic designs, called *ScalaAFA*, which is tightly integrated into SPDK while harnessing SSD built-in resources to deliver scalable performance at low CPU costs. To be specific, ScalaAFA employs a message-passing-based permission management mechanism for concurrent access thereby conforming to the lock-free principle of SPDK. ScalaAFA offloads the conversion tasks to SSDs, which reduces the CPU burden. Lastly, ScalaAFA addresses the metadata persistence and write amplification issues by exploiting SSD architectural innovations. Evaluation results reveal that ScalaAFA improves the write throughput to 2.5× and reduces the average latency by 52.7% compared to the state-of-the-art AFA solutions.

Acknowledgement

We thank the anonymous reviewers for their constructive feedback. This work is mainly supported by the National Key Research and Development Program of China under Grant No. 2023YFB4502702, the Natural Science Foundation of China under Grant No. 62332021 and 62472007. Shushu Yi is supported in part by the National Natural Science Foundation of China under Grant No. 624B2004. Dr. Li is supported in part by the National Natural Science Foundation of China under Grant No. 62202396. Dr. Wang is supported in part by the Innovation Funding of ICT, CAS under Grant E261110. Dr. Mao is supported in part by the National Natural Science Foundation of China under Grant No. U22A2027. Dr. Jung is supported in part by NRF (RS-2021-NR059881), IITP (RS-2023-00256472), Samsung Electronics, KAIST IDEC, Technology development Program of MSS (RS-2023-00303967), Information & communications Technology Planning & Evaluation (IITP) grant (No.RS-2023-00221040, Enabling A Low-Power and High-Performance Computational SSD System that Supports the Execution of General-Purpose Applications, 10%), (No.RS-2024-00460762, Development of Core System Software for Fabric Memory-Based Computational Storage, 10%), (No.RS-2025-02214652, Development of SoC Technology for AI Semiconductor-Converged Pooled Storage/Memory, 10%), (No.RS-2025-02214654, Development of AI Semiconductor Converged Computational Memory/Storage SoC and Application Technology, 10%), (No.RS-2025-02263080, Development of PIM Memory Technology for a High-Speed Interface-Based Programmable Computing Architecture, 10%). The corresponding author is Jie Zhang.

References

- [1] 2010. NETAPP INC. Data ontap 8. <http://www.netapp.com/us/products/platform-os/data-ontap-8/>.
- [2] 2020. NVIDIA® Mellanox® ConnectX®-5 Adapters. <https://www.nvidia.com/en-sg/networking/ethernet/connectx-5/>.

- [3] 2020. Samsung 980Pro NVMe SSD. <https://www.samsung.com/us/computing/memory-storage/solid-state-drives/980-pro-pcie-4-0-nvme-ssd-1tb-mz-v8p1t0b-am/>.
- [4] 2022. Marvel Bravera SC5 SSD Controllers. <https://www.marvell.com/products/ssd-controllers/mv-ss1331-1333.html>.
- [5] 2022. mdraid layer. <https://github.com/torvalds/linux/tree/master/drivers/md>.
- [6] 2022. NVM Express Base Specification 2.0c. <https://nvmexpress.org/wp-content/uploads/NVM-Express-Base-Specification-2.0c-2022.10.04-Ratified.pdf>.
- [7] 2022. NVM Express NVM Command Set Specification 1.0c. <https://nvmexpress.org/wp-content/uploads/NVM-Express-NVM-Command-Set-Specification-1.0c-2022.10.03-Ratified.pdf>.
- [8] 2022. RAID5F. <https://github.com/spdk/spdk/tree/master/module/bdev/raid/raid5f.c>.
- [9] 2022. SPDK v22.05. <https://github.com/spdk/spdk/tree/v22.05.x>.
- [10] 2023. DELL Dell PowerStore 500T STORAGE ARRAY. <https://www.delltechnologies.com/asset/en-ca/products/storage/technical-support/dell-powerstore-gen2-spec-sheet.pdf>.
- [11] 2023. Samsung PM1743. <https://semiconductor.samsung.com/ssd/enterprise-ssd/pm1743/>.
- [12] Ahmed Izzat Alslibi, Sparsh Mittal, Mohammed Azmi Al-Betar, and Putra Bin Sumari. 2018. A survey of techniques for architecting SLC/MLC/TLC hybrid Flash memory-based SSDs. *Concurrency and Computation: Practice and Experience* 30, 13 (2018), e4420.
- [13] Hagit Attiya and Jennifer L Welch. 1994. Sequential consistency versus linearizability. *ACM Transactions on Computer Systems (TOCS)* 12, 2 (1994), 91–122.
- [14] Jens Axboe. 2019. Flexible I/O Tester. <https://github.com/axboe/fio>.
- [15] Duck-Ho Bae, Insoon Jo, Youra Adel Choi, Joo-Young Hwang, Sangyeun Cho, Dong-Gi Lee, and Jaeheon Jeong. 2018. 2B-SSD: the case for dual, byte- and block-addressable solid-state drives. In *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*. 425–438.
- [16] Matias Björling, Abutalib Aghayev, Hans Holmberg, Aravind Ramesh, Damien Le Moal, Gregory R Ganger, and George Amvrosiadis. 2021. ZNS: Avoiding the Block Interface Tax for Flash-based SSDs. In *2021 USENIX Annual Technical Conference (ATC)*. 689–703.
- [17] Matias Björling, Javier Gonzalez, and Philippe Bonnet. 2017. LightNVM: The Linux Open-ChannelSSD Subsystem. In *15th USENIX Conference on File and Storage Technologies (FAST)*. 359–374.
- [18] John Canny, Huasha Zhao, Bobby Jaros, Ye Chen, and Jiangchang Mao. 2015. Machine learning at the limit. In *2015 IEEE International Conference on Big Data (Big Data)*. 233–242.
- [19] Zhichao Cao. 2020. *High-Performance and Cost-Effective Storage Systems for Supporting Big Data Applications*. Ph.D. Dissertation. University of Minnesota.
- [20] Zhichao Cao, Siying Dong, Sagar Vemuri, and David HC Du. 2020. Characterizing, modeling, and benchmarking RocksDBKey-Value workloads at facebook. In *18th USENIX Conference on File and Storage Technologies (FAST)*. 209–223.
- [21] Feng Chen, Tian Luo, and Xiaodong Zhang. 2011. CAFTL: A Content-Aware Flash Translation Layer Enhancing the Lifespan of Flash Memory based Solid State Drives. In *9th USENIX Conference on File and Storage Technologies (FAST)*.
- [22] Man-Sheng Chen, Jia-Qi Lin, Xiang-Long Li, Bao-Yu Liu, Chang-Dong Wang, Dong Huang, and Jian-Huang Lai. 2022. Representation learning in multi-view clustering: A literature review. *Data Science and Engineering* 7, 3 (2022), 225–241.
- [23] Tzi-cker Chiueh, Weafon Tsao, Hou-Chiang Sun, Ting-Fang Chien, An-Nan Chang, and Cheng-Ding Chen. 2014. Software orchestrated flash array. In *Proceedings of International Conference on Systems and Storage*. 1–11.
- [24] Hyeokjun Choe, Seil Lee, Hyunha Nam, Seongsik Park, Seijoon Kim, Eui-Young Chung, and Sungroh Yoon. 2016. Near-data processing for differentiable machine learning models. *arXiv preprint arXiv:1610.02273* (2016).
- [25] Arnaldo Carvalho De Melo. 2010. The new linux ‘perf’ tools. In *Slides from Linux Kongress*, Vol. 18. 1–42.
- [26] Facebook. 2015. RocksDB. <http://rocksdb.org/>.
- [27] Facebook. 2021. Performance Benchmarks. <https://github.com/facebook/rocksdb/wiki/Benchmarking-tools>.
- [28] Bin Fan, Wittawat Tantisiriroj, Lin Xiao, and Garth Gibson. 2009. DiskReduce: RAID for data-intensive scalable computing. In *Proceedings of the 4th annual workshop on petascale data storage*. 6–10.
- [29] Jim Gray, Paul McJones, Mike Blasgen, Bruce Lindsay, Raymond Lorie, Tom Price, Franco Putzolu, and Irving Traiger. 1981. The recovery manager of the System R database manager. *ACM Computing Surveys (CSUR)* 13, 2 (1981), 223–242.
- [30] Laura M Grupp, John D Davis, and Steven Swanson. 2012. The bleak future of NAND flash memory.. In *FAST*, Vol. 7. 10–2.
- [31] Aayush Gupta, Raghav Pisolkar, Bhuvan Ugaonkar, and Anand Sivasubramaniam. 2011. Leveraging Value Locality in Optimizing NAND Flash-based SSDs. In *9th USENIX Conference on File and Storage Technologies (FAST)*.
- [32] Mingzhe Hao, Gokul Soundararajan, Deepak Kenchammana-Hosekote, Andrew A Chien, and Haryadi S Gunawi. 2016. The tail at store: A revelation from millions of hours of disk and SSD deployments. In *14th USENIX Conference on File and Storage Technologies (FAST)*. 263–276.
- [33] Brian Hickmann and Kynan Shook. 2007. ZFS and RAID-Z: The Über-FS? *University of Wisconsin–Madison* (2007).
- [34] Dave Hitz, James Lau, and Michael A Malcolm. 1994. File System Design for an NFS File Server Appliance.. In *USENIX winter*, Vol. 94. 10–5555.

- [35] Seongcheol Hong and Dongkun Shin. 2010. NAND flash-based disk cache using SLC/MLC combined flash memory. In *2010 International Workshop on Storage Network Architecture and Parallel I/Os*. 21–30.
- [36] Ping Huang, Pradeep Subedi, Xubin He, Shuang He, and Ke Zhou. 2014. FlexECC: Partially Relaxing ECC of MLCSSD for Better Cache Performance. In *2014 USENIX Annual Technical Conference (ATC)*. 489–500.
- [37] Soojun Im and Dongkun Shin. 2010. ComboFTL: Improving performance and lifespan of MLC flash memory using SLC flash buffer. *Journal of Systems Architecture* 56, 12 (2010), 641–653.
- [38] Nikolaus Jeremic, Gero Mühl, Anselm Busse, and Jan Richling. 2011. The pitfalls of deploying solid-state drive RAIDs. In *Proceedings of the 4th Annual International Conference on Systems and Storage*. 1–13.
- [39] Tianyang Jiang, Guangyan Zhang, Zican Huang, Xiaosong Ma, Junyu Wei, Zhiyue Li, and Weimin Zheng. 2021. FusionRAID: Achieving Consistent Low Latency for Commodity SSD Arrays. In *19th USENIX Conference on File and Storage Technologies (FAST)*. 355–370.
- [40] Yangwook Kang, Yang-suk Kee, Ethan L Miller, and Chanik Park. 2013. Enabling cost-effective data processing with smart SSD. In *2013 IEEE 29th symposium on mass storage systems and technologies (MSST)*. 1–12.
- [41] Swaroop Kavalanekar, Bruce Worthington, Qi Zhang, and Vishal Sharda. 2008. Characterization of storage workload traces from production windows servers. In *2008 IEEE International Symposium on Workload Characterization*. 119–128.
- [42] Aleksandr Khasymiski, M Mustafa Rafique, Ali R Butt, Sudharshan S Vazhkudai, and Dimitrios S Nikolopoulos. 2012. On the use of GPUs in realizing cost-effective distributed RAID. In *2012 IEEE 20th International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems*. 469–478.
- [43] Byungseok Kim, Jaeho Kim, and Sam H Noh. 2017. Managing array of SSDs when the storage device is no longer the performance bottleneck. In *Proceedings of the 9th USENIX Conference on Hot Topics in Storage and File Systems*. 20–20.
- [44] Jaeho Kim, Kwanghyun Lim, Youngdon Jung, Sungjin Lee, Changwoo Min, and Sam H Noh. 2019. Alleviating garbage collection interference through spatial separation in all flash arrays. In *2019 USENIX Annual Technical Conference (ATC)*. 799–812.
- [45] Youngjae Kim, Sarp Oral, Galen M Shipman, Junghee Lee, David A Dillow, and Feiyi Wang. 2011. Harmonia: A globally coordinated garbage collector for arrays of solid-state drives. In *2011 IEEE 27th Symposium on Mass Storage Systems and Technologies (MSST)*. 1–12.
- [46] Miryeong Kwon, Donghyun Gouk, Sangwon Lee, and Myoungsoo Jung. 2022. Hardware/SoftwareCo-Programmable Framework for Computational SSDs to Accelerate Deep Learning Service on Large-Scale Graphs. In *20th USENIX Conference on File and Storage Technologies (FAST)*. 147–164.
- [47] Miryeong Kwon, Donghyun Gouk, Sangwon Lee, and Myoungsoo Jung. 2022. Hardware/SoftwareCo-Programmable Framework for Computational SSDs to Accelerate Deep Learning Service on Large-Scale Graphs. In *20th USENIX Conference on File and Storage Technologies (FAST)*. 147–164.
- [48] Miryeong Kwon, Jie Zhang, Gyuyoung Park, Wonil Choi, David Donofrio, John Shalf, Mahmut Kandemir, and Myoungsoo Jung. 2017. Tracetracker: Hardware/software co-evaluation for large-scale i/o workload reconstruction. In *2017 IEEE International Symposium on Workload Characterization (IISWC)*. 87–96.
- [49] Sungjin Lee, Ming Liu, Sangwoo Jun, Shuotao Xu, Jihong Kim, et al. 2016. Application-Managed Flash. In *14th USENIX Conference on File and Storage Technologies (FAST)*. 339–353.
- [50] Huaicheng Li, Mingzhe Hao, Michael Hao Tong, Swaminathan Sundararaman, Matias Björling, and Haryadi S Gunawi. 2018. The CASE of FEMU: Cheap, accurate, scalable and extensible flash emulator. In *16th USENIX Conference on File and Storage Technologies (FAST)*. 83–90.
- [51] Huaicheng Li, Martin L Putra, Ronald Shi, Xing Lin, Gregory R Ganger, and Haryadi S Gunawi. 2021. IODA: A Host/Device Co-Design for Strong Predictability Contract on Modern Flash Storage. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*. 263–279.
- [52] Qiang Li, Qiao Xiang, Yuxin Wang, Haohao Song, Ridi Wen, Wenhui Yao, Yuanyuan Dong, Shuqi Zhao, Shuo Huang, Zhaosheng Zhu, et al. 2023. More Than Capacity: Performance-oriented Evolution of Pangu in Alibaba. In *21st USENIX Conference on File and Storage Technologies (FAST)*. 331–346.
- [53] Runhui Li, Yuchong Hu, and Patrick PC Lee. 2015. Enabling Efficient and Reliable Transition from Replication to Erasure Coding for Clustered File Systems. In *2015 45th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*. 148–159.
- [54] Yongkun Li, Helen HW Chan, Patrick PC Lee, and Yinlong Xu. 2016. Elastic parity logging for SSD RAID arrays. In *2016 46th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. 49–60.
- [55] Shengwen Liang, Ying Wang, Cheng Liu, Huawei Li, and Xiaowei Li. 2019. InS-DLA: An in-SSD deep learning accelerator for near-data processing. In *2019 29th International Conference on Field Programmable Logic and Applications (FPL)*. 173–179.
- [56] Shengwen Liang, Ying Wang, Youyou Lu, Zhe Yang, Huawei Li, and Xiaowei Li. 2019. Cognitive SSD: A Deep Learning Engine for In-Storage Data Retrieval. In *2019 USENIX Annual Technical Conference (ATC)*. 395–410.
- [57] Bo Mao, Suzhen Wu, and Lide Duan. 2017. Improving the SSD performance by exploiting request characteristics and internal parallelism. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 37, 2 (2017), 472–484.
- [58] Xiangfu Meng, Hongjin Huo, Xiaoyan Zhang, Wanchun Wang, and Jinxia Zhu. 2023. A survey of personalized news recommendation. *Data Science and Engineering* 8, 4 (2023), 396–416.

- [59] Rino Micheloni, Luca Crippa, and Alessia Marelli. 2010. *Inside NAND flash memories*. Springer Science & Business Media.
- [60] Jaehong Min, Ming Liu, Tapan Chugh, Chenxingyu Zhao, Andrew Wei, In Hwan Doh, and Arvind Krishnamurthy. 2021. Gimbal: enabling multi-tenant storage disaggregation on SmartNIC JBOFs. In *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*. 106–122.
- [61] Dushyanth Narayanan, Austin Donnelly, and Antony Rowstron. 2007. MSR Cambridge Traces (SNIA IOTTA Trace Set 388). In *SNIA IOTTA Trace Repository*, Geoff Kuenning (Ed.). Storage Networking Industry Association. <http://iota.snia.org/traces/block-io?only=388>
- [62] Sudharsan Seshadri, Mark Gahagan, Sundaram Bhaskaran, Trevor Bunker, Arup De, Yanqin Jin, Yang Liu, and Steven Swanson. 2014. Willow: A User-Programmable SSD. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*. 67–80.
- [63] Xuanhua Shi, Ming Li, Wei Liu, Hai Jin, Chen Yu, and Yong Chen. 2017. Ssdup: a traffic-aware ssd burst buffer for hpc systems. In *Proceedings of the international conference on supercomputing*. 1–10.
- [64] Ji-Yong Shin, Zeng-Lin Xia, Ning-Yi Xu, Rui Gao, Xiong-Fei Cai, Seungryoul Maeng, and Feng-Hsiung Hsu. 2009. FTL design exploration in reconfigurable high-performance SSD for server applications. In *Proceedings of the 23rd international conference on Supercomputing*. 338–349.
- [65] Junyi Shu, Ruidong Zhu, Yun Ma, Gang Huang, Hong Mei, Xuanzhe Liu, and Xin Jin. 2023. Disaggregated RAID Storage in Modern Datacenters. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*. 147–163.
- [66] SPDK. [n. d.]. Blobstore filesystem. <https://spdk.io/doc/blobfs.html>.
- [67] Akshat Verma, Ricardo Koller, Luis Useche, and Raju Rangaswami. 2009. FIU Traces (SNIA IOTTA Trace Set 390). In *SNIA IOTTA Trace Repository*, Geoff Kuenning (Ed.). Storage Networking Industry Association. <http://iota.snia.org/traces/block-io?only=390>
- [68] Rui Wang, Yongkun Li, Hong Xie, Yinlong Xu, and John CS Lui. 2020. Graphwalker: An i/o-efficient and resource-friendly graph analytic system for fast and scalable random walks. In *Proceedings of the 2020 USENIX Conference on Usenix Annual Technical Conference*. 559–571.
- [69] Shucheng Wang, Qiang Cao, Ziyi Lu, Hong Jiang, Jie Yao, and Yuanyuan Dong. 2022. StRAID: Stripe-threaded Architecture for Parity-based RAID with Ultra-fast SSDs. In *2022 USENIX Annual Technical Conference (ATC)*. 915–932.
- [70] Stephen B Wicker and Vijay K Bhargava. 1999. *Reed-Solomon codes and their applications*. John Wiley & Sons.
- [71] John Wilkes, Richard Golding, Carl Staelin, and Tim Sullivan. 1996. The HP AutoRAID hierarchical storage system. *ACM Transactions on Computer Systems (TOCS)* 14, 1 (1996), 108–136.
- [72] Greg Wong. 2013. SSD market overview. In *Inside Solid State Drives (SSDs)*. Springer, 1–17.
- [73] Suzhen Wu, Haijun Li, Bo Mao, Xiaoxi Chen, and Kuan-Ching Li. 2018. Overcome the GC-induced performance variability in SSD-based RAID with request redirection. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 38, 5 (2018), 822–833.
- [74] Suzhen Wu, Bo Mao, Xiaolan Chen, and Hong Jiang. 2016. LDM: Log disk mirroring with improved performance and reliability for SSD-based disk arrays. *ACM Transactions on Storage (TOS)* 12, 4 (2016), 1–21.
- [75] Shiqin Yan, Huaicheng Li, Mingzhe Hao, Michael Hao Tong, Swaminathan Sundararaman, Andrew A Chien, and Haryadi S Gunawi. 2017. Tiny-tail flash: Near-perfect elimination of garbage collection tail latencies in NAND SSDs. *ACM Transactions on Storage (TOS)* 13, 3 (2017), 1–26.
- [76] Jisoo Yang, Dave B Minturn, and Frank Hady. 2012. When poll is better than interrupt.. In *FAST*, Vol. 12. 3–3.
- [77] Pan Yang, Ni Xue, Yuqi Zhang, Yangxu Zhou, Li Sun, Wenwen Chen, Zhonggang Chen, Wei Xia, Junke Li, and Kihyoun Kwon. 2019. Reducing garbage collection overhead in SSD based on workload prediction. In *11th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 19)*.
- [78] Ziye Yang, James R Harris, Benjamin Walker, Daniel Verkamp, Changpeng Liu, Cunyin Chang, Gang Cao, Jonathan Stern, Vishal Verma, and Luse E Paul. 2017. SPDK: A development kit to build high performance storage applications. In *2017 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*. 154–161.
- [79] Shushu Yi, Yanning Yang, Yunxiao Tang, Zixuan Zhou, Junzhe Li, Chen Yue, Myoungsoo Jung, and Jie Zhang. 2022. ScalaRAID: optimizing linux software RAID system for next-generation storage. In *Proceedings of the 14th ACM Workshop on Hot Topics in Storage and File Systems*. 119–125.
- [80] Chi Zhang, Yi Wang, Tianzheng Wang, Renhai Chen, Duo Liu, and Zili Shao. 2014. Deterministic crash recovery for NAND flash based storage systems. In *Proceedings of the 51st Annual Design Automation Conference*. 1–6.
- [81] Jie Zhang, Miryeong Kwon, Michael Swift, and Myoungsoo Jung. 2020. Scalable parallel flash firmware for many-core architectures. In *18th USENIX Conference on File and Storage Technologies (FAST)*. 121–136.

Received 22 August 2024; revised 14 February 2025; accepted 30 April 2025