

Fastmove: A Comprehensive Study of On-Chip DMA and its Demonstration for Accelerating Data Movement in NVM-based Storage Systems

JIAHAO LI, Computer Science, University of Science and Technology of China, Hefei, China
JINGBO SU, Computer Science, University of Science and Technology of China, Hefei, China
LUOFAN CHEN, Computer Science, University of Science and Technology of China, Hefei, China
CHENG LI, CS. Dept, University of Science and Technology of China, Hefei, China
KAI ZHANG, SmartX, Beijing, China
LIANG YANG, SmartX, Beijing, China
SAM NOH, Virginia Tech, Blacksburg, United States
YINLONG XU, University of Science and Technology of China, Hefei, China

Data-intensive applications executing on NVM-based storage systems experience serious bottlenecks when moving data between DRAM and NVM. We advocate for the use of the long-existing but recently neglected on-chip DMA to expedite data movement with three contributions. First, we explore new latency-oriented optimization directions, driven by a comprehensive DMA study, to design a high-performance DMA module, which significantly lowers the I/O size threshold to observe benefits. Second, we propose a new data movement engine, Fastmove, that coordinates the use of the DMA along with the CPU with DDIO-aware strategies, judicious scheduling and load splitting such that the DMA's limitations are compensated, and the overall gains are maximized. Finally, with a general kernel-based design, simple APIs, and DAX file system integration, Fastmove allows applications to transparently exploit the DMA and its new features without code change. We run three data-intensive applications MySQL, GraphWalker, and Filebench atop NOVA, ext4-DAX, and XFS-DAX, with standard benchmarks like TPC-C, and popular graph algorithms like PageRank. Across single- and multi-socket settings, compared to the conventional CPU-only NVM accesses, Fastmove introduces to TPC-C with MySQL 1.13-2.16 \times speedups of peak throughput, reduces the average latency by 17.7-60.8%, and saves 37.1-68.9% CPU usage spent in data movement. It also shortens the execution time of graph algorithms with GraphWalker by 39.7-53.4%, and introduces 1.01-1.48 \times throughput speedups for Filebench.

CCS Concepts: • **Hardware** → **Non-volatile memory**; **Hardware accelerators**; • **Software and its engineering** → *Software usability*.

Additional Key Words and Phrases: persistent memory, direct memory access, direct cache access

"SmartX" is also known as Beijing Zhiling Haina Technology Co., Ltd.

Authors' Contact Information: Jiahao Li, Computer Science, University of Science and Technology of China, Hefei, Anhui, China; e-mail: lijh2015@mail.ustc.edu.cn; Jingbo Su, Computer Science, University of Science and Technology of China, Hefei, Anhui, China; e-mail: sujib2113@mail.ustc.edu.cn; Luofan Chen, Computer Science, University of Science and Technology of China, Hefei, Anhui, China; e-mail: clfbbn@mail.ustc.edu.cn; Cheng Li, CS. Dept, University of Science and Technology of China, Hefei, Anhui, China; e-mail: chengli7@ustc.edu.cn; Kai Zhang, SmartX, Beijing, Beijing, China; e-mail: kyle@smartx.com; Liang Yang, SmartX, Beijing, Beijing, China; e-mail: liang.yang@smartx.com; Sam Noh, Virginia Tech, Blacksburg, Virginia, United States; e-mail: samhnoh@vt.edu; Yinlong Xu, University of Science and Technology of China, Hefei, Anhui, China; e-mail: ylXu@ustc.edu.cn.

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, or post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM /2024/5-ART

<https://doi.org/10.1145/3656477>

1 INTRODUCTION

Emerging non-volatile memory (NVM) technologies such as STT-MRAM [50], PCM [44], ReRAM [6], and 3D-XPoint [17] offer byte-addressability and comparable latency as DRAM but with substantially larger capacity. In addition, it provides data durability with orders of magnitude higher performance than prior durable devices like SSDs [60]. Recently, numerous studies have been proposed to combine faster, volatile DRAM, for caching, with slightly slower, denser NVM, for persisting data, in storage systems to revolutionize I/O performance of data-intensive applications with persistence demands [12].

In NVM-based storage systems, data are often moved between the two types of memories, due to DRAM cache fill-up, logging, or flushing. However, recent studies [31, 60] highlight that the DRAM-NVM data movement is not efficient, mainly because of their performance gaps in latency and bandwidth [13]. Additionally, we further notice that such data movement leads to heavy CPU consumption since NVM chips are attached to the memory bus, and their accesses must make use of the load and store instructions. Such negative performance effects worsen with multiple sockets, which modern high-end servers often provide, because of the negative NUMA impact [31]. This data movement bottleneck severely impairs the overall performance of I/O intensive applications and consequently, undermines the benefits brought by incorporating NVM.

To address this bottleneck, the slowness of NVM motivates us to re-think the usage of the on-chip DMAs that still come with the CPU but have deteriorated in use with the advent of fast storage devices. In this paper, we seek to transparently expedite data movement in NVM-based storage systems by (partially) offloading data movement to DMA to improve overall performance. However, while exploiting the on-chip DMA is a natural optimization, there are a few obstacles to incorporating it into NVM-based storage systems.

First, we need to handle more complex I/O patterns and have significantly different optimization goals than existing work [45, 59], which have already applied DMA as a minor technique to free CPU cycles of page migration in tiered DRAM-NVM systems. They handle I/Os that are always large, i.e., 2MB, and run in the background. However, NVM-based storage systems face I/Os with much smaller and variable sizes that are often on the critical path of the foreground user requests. Thus, our primary optimization goal is to shorten the execution time of DMA requests. Second, latency-critical optimization requires an in-depth understanding of the strengths and limits of DMA, in conjunction with NVM and storage-facing I/Os, which is largely beyond existing studies [24].

To address the above challenges, first, we conduct a comprehensive study to understand the latency behaviors of using DMA for DRAM-NVM data movement on the Intel I/OAT and Optane PM combination, the only such pair in existence. This study suggests that the potential of DMA is heavily constrained by various factors, e.g., uneven advantages between reads and writes over the CPU, the non-negligible costs that grow with I/O size, Data Directed I/O (DDIO) configuration, bandwidth and concurrency limits, etc.

Second, we derive principles from the study to design Fastmove, a general data movement system that sits at the lower level of the software hierarchy. At the core, it includes a high-performance DMA module, which encapsulates the upper-level I/O requests into low-level hardware commands that comply with the workflows of data movement in NVM-based storage systems. We also maximize the benefits of DMA by introducing various optimizations such as batching the page pinning and descriptor submission activities for grouped DMA tasks and balancing and coordinating concurrent accesses to DMA channels. Furthermore, to compensate for the limitations of the stand-alone DMA solution such as the extra overhead and the concurrency and bandwidth constraints, we devise a lightweight Scheduler to prioritize bulk I/Os to go through DMA, while smaller I/Os are routed to the original CPU path. Scheduler additionally splits bulk read I/Os and balance loads between the DMA and CPU paths, adapting to real-time changes in DMA resource availability.

Finally, we incorporate Fastmove into the Linux kernel as an OS library with a limited number of simple APIs, which can be used to easily replace system functions that trigger data movements. To demonstrate its practicality, we adapt three NVM-based storage systems, NOVA, ext4-DAX, and XFS-DAX, to make use of Fastmove with

minimal (2 to 4 lines of code) change. Consequently, applications running atop these systems can transparently enjoy the data movement acceleration brought by Fastmove. Additionally, we enable such acceleration for the cross-socket setting by deploying file systems atop the Linux device mapper with 2 lines of code change. This design enables the POSIX `read()` and `write()` APIs to freely employ Fastmove. To prove this, we successfully run three I/O-intensive applications, one industry-adopted database, MySQL [3], one graph engine, GraphWalker [54], and one file system and storage benchmark, Filebench [1] atop the modified file systems without any modifications to the applications.

We conduct extensive evaluations with three standard benchmarks FIO [8], fileserv [1], and TPC-C [5], and three popular random walk algorithms GraphLet, PageRank and SimRank. The results highlight that, for workloads containing substantial I/Os with moderately large sizes and beyond, considerable performance improvements are attained, regardless of local or remote NVM access. For TPC-C in MySQL, Fastmove increases its peak throughput by 13-116% compared to the original ones that use only the CPU, reduces the average latency by 17.7-60.8%, and saves CPU cycles used for data movement by 37.1-68.9%. Also, Fastmove brings 1.65-2.14× speedups of execution time for the GraphWalker algorithms, and 1.01-1.48× speedups of throughput for the Filebench fileserv workload.

In summary, Fastmove makes the following contributions:

- We present a comprehensive and general study to understand the characteristics of on-chip DMA in conjunction with NVM far beyond earlier studies [24], which showed DMA use just as a minor optimization in limited experimental settings [7, 28, 45].
- We propose and implement a fast memory copy engine Fastmove that accelerates DRAM-NVM data movement in NVM-based storage systems. Driven by the study findings, it incorporates new latency-oriented optimizations to reduce associated DMA costs and coordinates the CPU-only and DMA paths to maximize overall performance. Fastmove’s design principles significantly differ from earlier studies that concentrated on movement of data in a tiered memory setting [45, 59], where optimizations are simple due to the large size of memory copy requests.
- We present transparent in-kernel system support with integration of Fastmove into three NVM-aware DAX file systems, while extending the device mapper to enable cross-socket NVM access. This allows unmodified applications to run atop Fastmove.

2 BACKGROUND AND MOTIVATION

2.1 NVM-based Storage Systems

NVM chips sit close to the CPU either by being placed on the memory bus and connected to CPU sockets via the processors’ integrated memory controller (iMC) or by being exposed via cache coherence interconnects like Compute Express Link (CXL) [11, 21, 43]. In 2019, Intel released Optane PM, the first commercial NVM chip based on the 3D XPoint technology [17]. Beyond Optane PM, multiple companies are developing new products based on technologies other than 3D XPoint [17] such as STT-MRAM [50], FRAM [26], Nano-RAM [46], and ReRAM [6].

Despite the different implementations, they are expected to offer memory interfaces with byte-addressability, data persistence, and large capacity. Therefore, there have been extensive research focusing on incorporating NVM to build scalable storage systems [9, 25, 27, 30, 38, 58] that accelerate the data access of latency-critical, data-intensive applications. These applications persist all their data on NVM, while caching the working set and metadata like indexes in DRAM. When accessing non-cached data, applications need to load them from storage, while upon modification, the dirty pages and log entries need to be flushed back to storage for data durability.

Table 1. I/O size (KB) distribution of various workloads

	size	TPC-C	fileserver	Graphlet/PPR/SR
read	[0,16)	-	80.2%	-
	[16,32)	100%	11.5%	-
	[32,∞)	-	8.3%	100%
write	[0,16)	6.5%	82.2%	-
	[16,32)	82.9%	10.2%	-
	[32,∞)	10.6%	7.6%	-

Typically, they make use of NVM-aware DAX file systems such as NOVA that retain the standard file system interfaces and provide strong consistency guarantees along with various NVM-oriented performance optimizations [7, 23, 58]. Therefore, the aforementioned data copies often involve memory allocated in user space, while requiring kernel memory copy module support.

2.2 The Data Movement Bottleneck

DRAM-NVM data movement can be a critical bottleneck in terms of performance in data-intensive applications. To understand this, we perform a study on the I/O size distributions of various applications, from domains ranging from traditional SQL databases to graph analytic frameworks, and their impact on performance and resource usage.

As shown in Table 1, driven by the standard database TPC-C workloads with 5000 warehouses and 16KB innodb page size in MySQL, more than 93% of write I/Os in MySQL are beyond 16KB, where a significant number of these bulk writes are sitting on the critical path of writing logs for foreground update transactions. In the fileserver workload of Filebench, 8.3% and 7.6% of the reads and writes are beyond 32KB, respectively. Though the number of bulk I/Os is relatively small in fileserver, they already account for 44.1% of the overall data movement volume. Finally, GraphWalker, a single-threaded graph processing system, periodically reads from NVM into DRAM, all in 128KB chunks, which it later consumes with its in-memory processing [54].

To assess the negative impacts of data movement, we run the msprr workload [54] in GraphWalker atop NOVA, an NVM-based file system, with Optane PM. Note that NOVA uses Linux memcpy to access data on Optane and does not make use of SIMD as SIMD cannot be used within the kernel [47]. We find that over 92% of the execution time is spent on reading data from NVM under a single socket setting, while, when cross-socket data movement is involved, this number increases to over 97%. While these numbers will vary depending on the application, our observation is that for many applications, the time consumed for data movement is a clear bottleneck.

The inefficiencies of CPU-directed data movement are mainly caused by the performance gap between DRAM and NVM. In particular, with 6 interleaved Optane DIMMs within a single socket, reading a 4K page from Optane takes 952ns, 2.9× longer than that of DRAM. Similar to latency, PM shows 74.4%/35.3% lower read/write throughput than DRAM. Even worse, it takes 18 CPU cores for Optane to reach its peak load throughput while it only takes 5 for DRAM to reach a similar load throughput [60]. Finally, when accessing remote memory across sockets, both DRAM and NVM suffer negative NUMA effects due to the extra writes introduced by the default directory-based cache coherence protocol [31]. However, the performance loss of remote NVM accesses is larger because of its lower write bandwidth. Our findings are consistent with recent studies [16, 31, 60].

2.3 On-Chip DMA and its Challenges

Modern processors have included on-chip DMA engines since as far as one can remember. For instance, Intel’s I/O Acceleration Technology (I/OAT) DMA engine [18] lies in the integrated I/O module of the CPU, which also

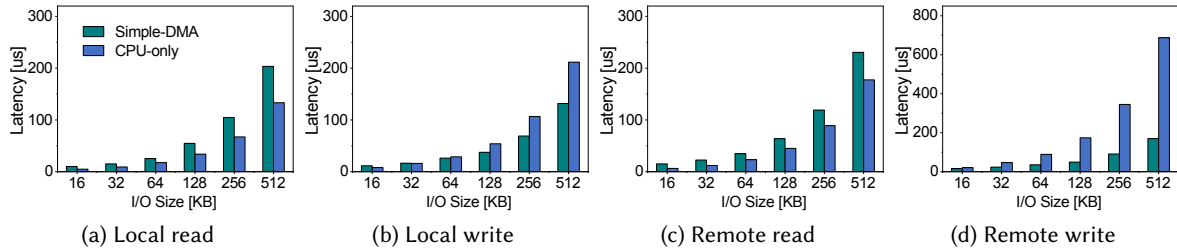


Fig. 1. Simple-DMA versus CPU-only read/write latency as request size is varied with FIO workloads.

connects to cores and iMCs through a mesh interconnect. Similarly, AMD’s second-generation EPYC processors are also equipped with on-chip DMA engines [41]. With the advent of high performance storage devices, however, they have deteriorated to a mostly unused component. The observations behind the data movement overhead problem motivate us to re-think the role of the on-chip DMA in NVM-based storage systems. We advocate that it will be beneficial to use on-chip DMAs to offload data copy jobs in NVM-based storage systems, thereby improving the copy performance itself as well as saving CPU cycles that could be used for other useful work.

To explore the latency improvement potential of DMA, we evaluate the speed of moving data between DRAM and NVM achieved by Intel I/OAT, in comparison with the CPU-only counterparts. Note that Data Direct I/O technology (DDIO) allows on-chip DMA to write directly into processor cache instead of through the memory controller. Initially, for all our experiments, we disable DDIO and use 64B aligned addresses. Then, separately in Section 3.2.4, we thoroughly examine the impact of DDIO and address alignment on data movement performance. Here, we refer to the I/OAT setting as Simple-DMA as we use it as-is without optimizations, which are explored later.

We use the FIO benchmark [8] to generate single-threaded read and write requests with I/O sizes ranging from 16KB to 512KB, where the former load data from NVM to DRAM while the latter store data in the opposite direction. These requests trigger kernel memory copy functions through NOVA to operate the underlying NVM–Optane PM [60], and we measure the time consumed for those functions.

Figure 1a and Figure 1c show that Simple-DMA performs consistently worse than CPU-only, and delivers 29.9-134.4% higher read latency, regardless of local and remote accesses. Contrary to reads, for local writes as shown in Figure 1b, Simple-DMA delivers comparable latency as CPU-only at 64KB, with meaningful differences expanding with I/O sizes from 128KB and beyond. For instance, when writing 256KB, the latency of Simple-DMA is only 64.7% of the CPU-only latency. Compared to the single-socket results, in Figure 1d, when considering two sockets, we observe that the performance of remote writes achieved by CPU-only and Simple-DMA both worsen. However, the request size threshold where Simple-DMA catches up with CPU-only becomes smaller at 16KB, which is only 25% of that observed for local writes.

The above latency comparison suggests that there is hardly any opportunity to allow reads within NVM-based storage systems to benefit from Simple-DMA; while for large writes, opportunities seem to exist. However, whether such large writes (not smaller than 128KB for local writes) are amply available in typical applications is questionable. For example, as shown in Table 1 in our evaluation, around 80% of the bulk writes for MySQL-TPC-C concentrate on the range of [16KB, 32KB], which is certainly below the benefit threshold of Simple-DMA. Our conclusion is that we need to explore whether there are optimization opportunities.

Moreover, we have witnessed initial adoptions [7, 24, 45, 59] of on-chip DMA to accelerate DRAM-NVM data movement. However, these early attempts mostly focus on tiered memory systems, and cannot be directly applied to NVM-based storage systems, which is our focus, due to the following reasons.

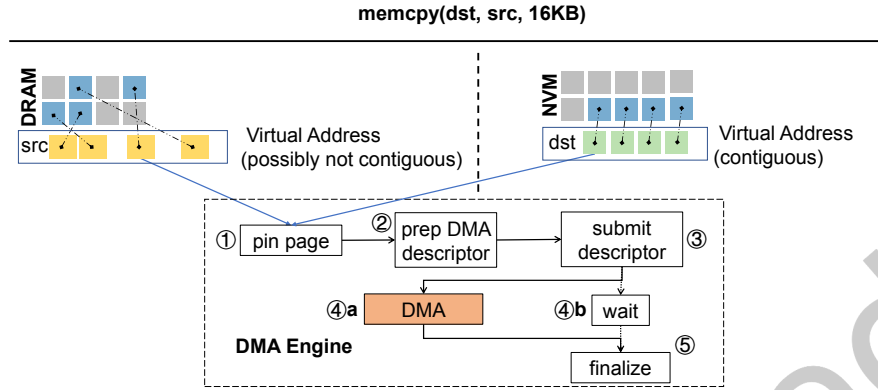


Fig. 2. Workflow of memory copy using Simple-DMA.

First, our optimization goal differs from using DMA in tiered memory systems, where data movements triggered by page migration run in the background, not on the critical path of user requests. Related works primarily focus their optimization goal on deriving advanced migration policies, and use DMA as a minor optimization to free CPU cycles [45]. In contrast, for NVM-based storage systems, data copy jobs such as user reads and log flushing are part of an end-to-end execution of foreground requests, which directly affect user experience. Thus, the key performance measure is latency.

Second, the I/O patterns and workflows differ significantly between NVM-based storage systems and tiered memory systems. The page migration workloads in tiered memory systems are quite simple and always happen at 2MB huge page granularity [45, 59]. In contrast, the sizes of bulk I/Os in NVM-based storage systems are much smaller and vary considerably. It is equally important that the workflow of handling memory copies via DMA in NVM-storage systems contains considerably more steps than that of tiered memory. These differences imply that the associated overhead of DMA is not negligible in NVM-based storage systems.

In summary, the Simple-DMA performance, the demand for reducing latency and the storage-specific I/O patterns present us with unique challenges in making use of the DMA in NVM-based storage systems. In this paper, through an in-depth study of the behavior of on-chip DMA, we explore avenues of optimization opportunities. In addition, through `Fastmove`, we develop the necessary abstractions and transparent latency-sensitive optimizations so that applications may reap the benefits of the DMA without any code change.

3 DMA OPTIMIZATION OPPORTUNITIES

Here we provide a comprehensive study on DMA in conjunction with NVM to derive the optimization directions for lowering the latency of DMA-enabled memory copies and for unleashing its potentials to (partially) alleviate the above DRAM-NVM data stall problem.

3.1 DMA-enabled Data Moving Workflow

To begin our study, we first illustrate in Figure 2 the workflow of handling memory copy requests issued by applications via DMA, which implements exactly the same logic as the Linux `memcpy`. Take a 16KB I/O as an example. The virtual addresses of data residing in DRAM for NVM-based storage systems are possibly not contiguous, which leads to this single memory copy operation at the application side being divided up into four

Table 2. Breakdown time costs of local read and write requests that use Simple-DMA

	size (KB)	cost (%)				#subtasks
		pin	submit	I/OAT	other	
read	16	4.7	12.7	80.4	2.2	8
	32	4.9	14.1	78.7	2.3	16
	64	5.1	14.8	77.9	2.3	32
write	16	6.2	15.7	75.3	2.9	8
	32	7.1	19.8	69.8	3.3	16
	64	7.9	23.3	65.3	3.5	32

DMA subtasks. Each subtask corresponds to a 4KB page and will go through the following steps. ① pins the target DRAM pages as we need to prevent those pages from being swapped out or modified during DMA execution. An alternative way to do so is to allocate a DMA buffer, but at the cost of imposing extra memory copies or giving up transparent support to applications. ② prepares the DMA descriptor, the required metadata for I/OAT, which is then submitted to the hardware at step ③. Meanwhile, the submitter waits (④b) until the completion of ④a and reaches the final step ⑤ to finalize the corresponding DMA subtask execution, e.g., unpinning the page and notifying the application. Note that all steps except ④a are managed by a CPU thread, often the I/O thread of the application.

3.2 I/OAT and Optane PM Demonstration

To make the study concrete, in this section, we focus on the combination of Optane PM and Intel’s I/OAT DMA.

3.2.1 Associated Time Costs. First, we investigate the latency breakdown results of Simple-DMA, which are summarized in Table 2, with the same setup as Figure 1a and Figure 1b. “pin”, “submit”, and “I/OAT” correspond to steps ①, ②-③, and ④a of Figure 2, respectively, while “others” denotes the remaining overhead.

The execution on the I/OAT hardware is the longest step of DMA-enabled memory copy requests across reads and writes. However, its ratio decreases from 80.4% to 77.9%, and 75.3% to 65.3% for reads and writes, respectively, when I/Os expand from 16KB to 64KB. In contrast, the associated overhead, excluding I/OAT, is also non-negligible and grows proportionally with request size, reaching to 34.7% for local 64KB writes. This is mainly because bulk I/Os within NVM-based storage systems trigger a series of I/OAT subtasks at 4KB granularity, as introduced in Section 3.1.

This growing overhead can be further doubled when the source and destination addresses of the corresponding I/O request are not aligned to 4KB boundary. Figure 3 illustrates such an example. The *src* of page#1 is not aligned with *dst* of page#a. As the DMA does not support cross-page copy when it cannot tell if the physical address is contiguous between pages, we have to split page#1 into two separate portions, namely ① and ②, where the former fits in the empty space of page#a, while the latter will have to fit on the lower part of page#b. Each of these portions will trigger a separate I/OAT subtask. Moreover, the remaining two pages #2 and #3 will go through the same effort. As the FIO workloads exhibit unaligned memory addresses, as shown in Table 2, bulk I/Os consist of 8-32 DMA subtasks and pay the associated time cost one more time. As this shows, in the case of transferring unaligned memory addresses, the overhead involved can turn out to be even more significant.

Trimming down the associated costs seems promising for improving the latency of writes. For instance, one can imagine that reducing them by 30.7% for 16KB writes will allow the DMA latency turning point to be reduced from 64KB to 16KB, enabling more applications, like MySQL, to gain performance benefits. However, this is not so with reads, since even completely eliminating these overheads still results in the DMA performing 11.1%-39.3% slower than CPU-only. Thus, other means to overcome this challenge must be conceived.

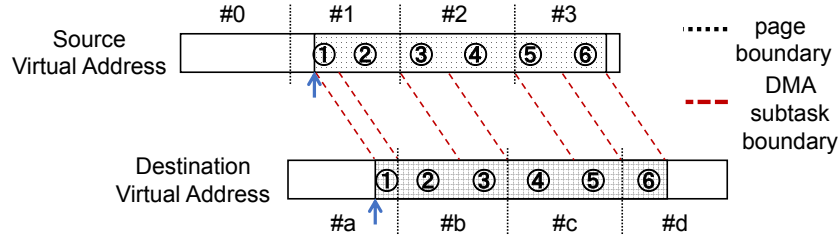


Fig. 3. Composition of DMA subtasks for a single application data movement job with unaligned source and destination addresses. Three source pages (#1-#3) are involved but six subtasks are generated (①-⑥).

In addition, using Transparent Huge Page (THP) in the kernel makes the addresses, with high probability, to be contiguous. For contiguous copies, the cost of I/OAT still dominates, but with the submission and unalignment cost significantly diminished, compared to the above non-contiguous ones. This is because under such setting, memory copy requests will no longer be divided into multiple DMA subtasks.

3.2.2 Intra-Request Parallel Copy. Each DMA device consists of M multiple channels that can process DMA subtasks in parallel. Therefore, we explore parallelizing hardware copy of a single request, where we split the request into N chunks ($N \leq M$) and thus, N DMA subtasks, each chunk making use of one channel. Here, we derive two different parallel execution modes, namely, para-A and para-B, where para-A uses a single submitter for channel submission, while para-B spawns N submitters, each of which manages its own channel independently.

For 64KB reads and writes, compared to Simple-DMA, para-A indeed reduces the I/OAT copy time, but the reduction is not proportional to the number of parallel chunks. In addition, we observe a significant increase in the submission overhead, which eventually offsets the benefits of intra-request, multi-channel parallel copy. In the end, para-A does not improve much on the end-to-end latency of Simple-DMA for reads, while even leading to performance loss for writes.

Para-B fares worse than para-A, worsening latency for both reads and writes. Our analysis shows that para-B sharply increases the hardware copy time by up to 68.7%. This is because of the heavy contention on DMA bandwidth driven by the parallel subtasks. This case differs from para-A, as the single submitter setting in para-A enables pipeline parallelism, which does not heavily stress the DMA. In addition, para-B introduces heavy CPU usage due to the multiple submitters.

Finally, as we cannot parallelize intra-request copies within DMA, we also explore the possibilities of balancing these copy subtasks between the CPU and DMA. Unfortunately, this is not applicable for writes, as using the DMA can easily saturate NVM's bandwidth. We find that the bandwidth competition can lead to amplified interference between the two tasks, resulting in 14.6% higher latency compared to the sole execution of using the DMA. In contrast, we find this solution works well for reads as the DMA cannot consume all of the NVM bandwidth, and thus, the joint use of the CPU and DMA leads to better bandwidth consumption. We take this last approach as part of our optimization.

3.2.3 Impacts of Inter-Request Parallelism. Next, we evaluate the impact of inter-request parallelism as in reality, application threads may concurrently execute data movement requests and make use of the DMA. First, we investigate whether using more DMA channels influences the performance of DMA operators. To this end, we use four concurrent threads to submit DMA requests to its multiple DMA channels on our two-socket NUMA machine. Here, we exercise up to 8 channels per DMA device/NUMA node. Figure 4 shows the latencies of DMA operators with varied I/O sizes. With the increasing number of channels, irrespective of local/remote reads/writes,

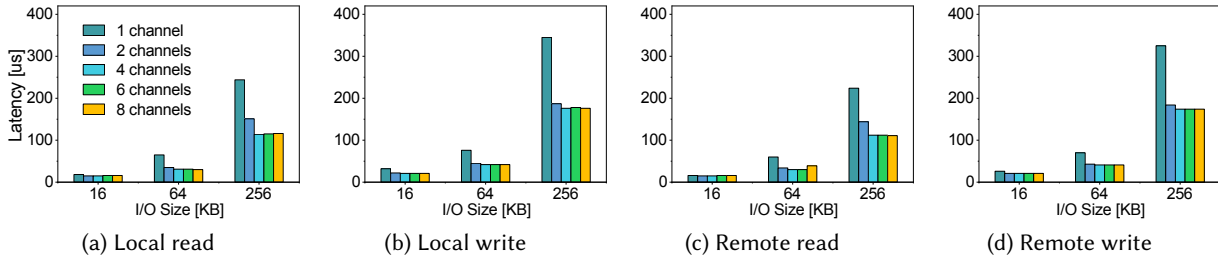


Fig. 4. I/OAT latency when 4 FIO threads doing read/write workloads as number of channels and as request size is varied.

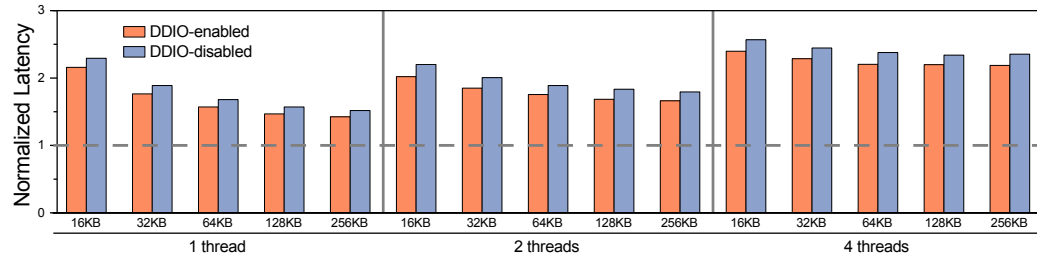
the DMA operators become faster. For instance, compared to the 1 channel setting, adding one more channel leads to 38.1%-53.3% latency reduction for the 256KB memory copy operators. Trends are similar with more concurrent threads and cross-socket NVM accesses.

Second, we explore the changes in read/write effective bandwidth with the increase in the number of concurrent threads submitting DMA requests with bulk I/Os. We find that Simple-DMA observes an increase in read/write effective bandwidth for up to four threads, but beyond this, it starts to decline sharply. (Results not shown due to space limit.) The key limiting factor here is not drive scalability but, instead, the I/OAT DMA bandwidth. This suggests that a limit on concurrent DMA access should be set to prevent the DMA resource from being over-used.

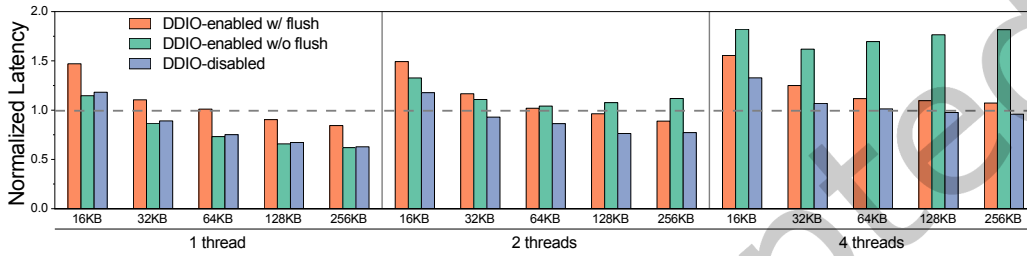
3.2.4 Impact of DDIO and Address Alignment. Finally, we examine how DDIO impacts the performance of DMA. When DDIO is enabled, DMA writes to PM are redirected to the CPU cache instead of being persisted. Subsequently, data is passively written back to NVM when the corresponding cache line is evicted. Kalia et al. [24] find that disabling DDIO can increase the peak PM write throughput of IOAT DMA. We replicate the experiment by running the FIO workloads on the NOVA file system equipped with Simple-DMA. We generate local read and local write workloads on 64B aligned addresses with various I/O sizes ranging from 16KB to 256KB that are issued by 1 to 4 concurrent threads. For DDIO-enabled writes, we have two different setups. In the first, flush does not occur after the request is done (denoted “DDIO-enabled w/o flush” or WOF). In the other, a manual flush is done after each write request (denoted as “DDIO-enabled w/ flush” or WF). Figure 5 shows the normalized average latency of Simple-DMA against CPU-only memory copying. We omit the remote copying results as they are similar to the local ones.

As shown in Figure 5a, disabling DDIO shows slightly worse read latency. This is because DDIO allows DMA to write from PM to CPU cache, shortening the data movement path, which originally involves DRAM. Overall, though, DDIO does not have a strong influence on read performance.

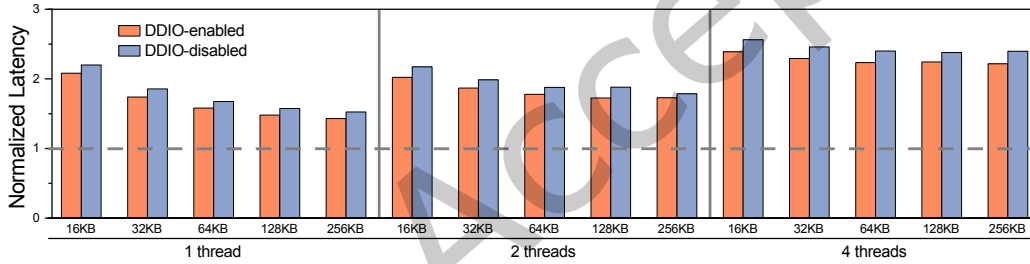
The results are different for writes that are shown Figure 5b. We first consider enabling DDIO. The performance of the two DDIO-enabled configurations are very close with 1 or 2 threads. WOF is slightly faster because it does not require waiting for flush operations. However, with 4 threads, WF shows significant advantage over WOF. To investigate the underlying reason, we further examine the write amplification using `ipmwatch`, a utility in Intel Vtune Profiler [19]. We find that for the experiments, WF has a write amplification of 1.0 while WOF has a write amplification of 2.1. With the huge write amplification, WOF saturates PM bandwidth early with only 2 threads. As explained in other studies [16, 60], data writes to PM are first stored in an internal buffer (XPbuffer) consisting of multiple XPLines (256B), which is the PM’s minimum access granularity. Under the WF mode, the modified cachelines are written back consecutively. In contrast, under the WOF mode, CPU only evicts modified cachelines when there is no enough space to place incoming new writes. However, eviction likely results in random writes to PM, leading to high write amplification and less effective PM bandwidth usage. Although



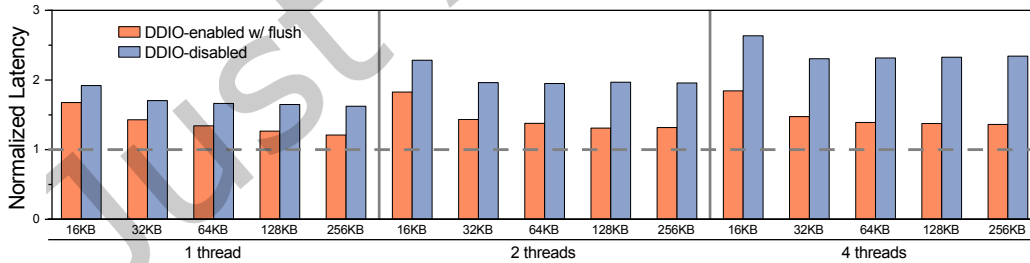
(a) Aligned Local read



(b) Aligned Local write



(c) Unaligned Local Read



(d) Unaligned Local Write

Fig. 5. The latency comparison between Simple-DMA under different DDIO settings, with 1,2,4-threaded FIO workloads, normalized to the latencies of CPU-only memory copying.

WOF shows lower latency without contention, i.e., one thread, it is harmful for multi-threading and sharing

cases, which are common in modern applications and file systems. In conclusion, we find that WF is a better configuration over WOF with additional persistence guarantee.

Second, also in Figure 5b, we explore the effects of disabling DDIO on writes. Unlike the above results, disabling DDIO constantly improves Simple-DMA's latency over enabling DDIO with flush, and it is the only configuration that can beat CPU-only with 4 threads.

We also observe that the alignment of addresses to the 64B boundary impacts DMA performance. To investigate how exactly it influences performance, we conduct an unaligned version of the previous experiments by modifying `mem_align` to 31. The results are shown in Figure 5c and Figure 5d. Unaligned addresses marginally increase the latency of read and WF write. In contrast, they noticeably increase the latency of write when DDIO is disabled. DMA performs poorly when data is moved directly into PM with unaligned addresses. We attempted to explain the phenomenon, but failed to witness any increase in write amplification or other anomalies. We also reproduce the experiment on DRAM (using emulated NVM [39]), but do not observe the influence of alignment. Therefore, we speculate the root cause lies in the implementation of I/OAT DMA and its interaction with PM. Based on the aforementioned observations, we conclude that disabling DDIO adversely impacts the performance of unaligned data movement.

3.3 Study Generalization

While the performance study above takes into consideration the performance characteristics of the underlying hardware, it also lays out the general study flow and key factors to be considered independent of particular NVM and DMA devices. With the advent of new hardware, the general study always needs to answer the following two questions:

First, *how can the DMA be best configured so that using it can be faster than CPU-only even for small I/O requests?* This part requires understanding the impact of DDIO, DMA subtask associated cost, the DMA parallel execution, and the effects of concurrency that drive the latency-oriented optimizations. Furthermore, it also requires exploring the effects of balancing loads among DMA channels and even between DMA and CPU.

Second, *how do we choose among the different copy paths?* We decide the best-effort path with the minimal time cost among three choices, namely, CPU, DMA, and DMA-CPU cooperation. Furthermore, we have to check if there are available DMA resources, i.e., the current DMA bandwidth usage, monitored during DMA execution, is below the profiled maximum bandwidths of DMA and NVM, respectively.

In summary, our general study framework will offer useful guidelines for accelerating data movement in storage systems that combine future DMA implementations and near-DRAM storage devices such as the upcoming CXL devices.

4 OVERVIEW OF FASTMOVE

Driven by the study in Section 3, we aim to let data-intensive applications transparently make the best use of DMA to alleviate the NVM data stall problem presented in Section 2.2. Done properly, this should lead to better performance and alleviate CPU involvement required for memory copies between DRAM and NVM. First, we need to improve the latency of DMA-enabled data movement by taking into consideration the access constraints of DMA such as extra overhead, resource allocation, and interference within DMA or with CPU. Second, to complement DMA's limitations, we need to judiciously determine when and how much to resort to the normal CPU data path. Finally, while DMA is supported by Linux kernel functions, applications should not be burdened by high development and optimization overhead to exploit the DMA. Thus, a clean abstraction that requires minimal changes to applications is imperative.

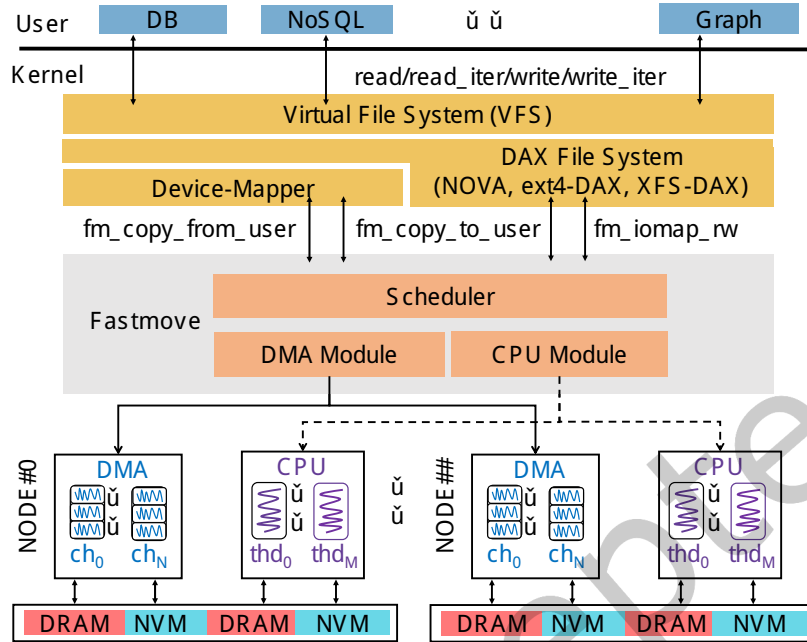


Fig. 6. The overall architecture of Fastmove, which manages both DMA and CPU resources. Each NUMA node has a DMA device (dashed line box), which has multiple channels.

4.1 Fastmove’s Architecture

Figure 6 shows the overall design of Fastmove, our efficient data movement engine. It sits below DAX file systems such as NOVA, ext4-DAX, and XFS-DAX, which are compatible with POSIX APIs and designed to use recent PM, as well as the Linux device mapper module, which allows file systems to use PMs across sockets. With this design, applications that run atop a POSIX file system should seamlessly be able to use our engine. Fastmove consists of three major system components, namely, Scheduler, DMA module, and CPU module. We retain the original design of the CPU module, where we let the corresponding I/O request execute the load and store instructions as usual. However, we introduce a new DMA module that manages DMA resource allocation and memory copy offloading, with various optimizations to alleviate DMA costs and improve DMA resource usage. (Details will be discussed in Section 5.1.)

As the core logic, Scheduler is responsible for making decisions on selecting either DMA or CPU to execute requests. Its detailed design will be discussed in Section 5.4. This decision-making procedure should be fast so as not to incur overhead on the end-to-end request latency. It should also be smart so as to prioritize the use of DMA to fully make use of its strengths, while resorting to the CPU-only path as needed to compensate for the limitations of DMA for overall enhanced performance (Section 5.3).

4.2 API Abstraction

To exploit DMA transparently at the application level, we introduce three APIs that are simple extensions to existing APIs used by DAX file systems. The key observation here is that DAX file systems universally make use of a limited number of APIs for data movement, namely, copy_from_user, copy_to_user, and dax_iomap_rw.

Table 3. Fastmove APIs

```

fm_copy_from_user(dst, src, len, bdev);
fm_copy_to_user(dst, src, len, bdev);
fm_iomap_rw(iocb, iov_iter, iomap_ops);

```

The first two are called by the read and write file system functions, while the last API is used by the `read_iter` and `write_iter` file system functions to perform memory copies in batches. These APIs are replaced by the APIs that we describe below.

As shown in Table 3, the three APIs that we introduce are `fm_copy_from_user`, `fm_copy_to_user`, and `fm_iomap_rw`. The first two new APIs have four arguments, `dst`, `src`, `len`, and `bdev`. `dst` and `src` specifies the destination and source of the copy (from PM to DRAM or vice versa), while `len` refers to the number of bytes to copy. The last argument `bdev` is the PM block device descriptor that includes rich information of the PM device such as the NUMA node id of the target PM. The last API `fm_iomap_rw` has three parameters, where `iocb` specifies the operational semantics such as read or write, `iov_iter` encodes parameters such as source and destination address vectors, and `iomap_ops` that is passed by file systems for I/O address mapping.

Finally, we only need to replace the old APIs with the new ones at the file system level. Thus, upper layer applications can take advantage of Fastmove without any code change. (Details are discussed in Section 5.5.)

5 DESIGN AND IMPLEMENTATION

5.1 High-Performance DMA Module

Under Fastmove, we offer a dedicated wrapper module to easily use the low-level primitives that DMA offers. This wrapper executes the I/O requests admitted by Scheduler. Here, we encapsulate the DMA requests by inheriting the values of parameters from the Fastmove APIs and the DMA channel assignment from Scheduler. Then, the wrapper executes DMA requests by going through all the steps in Figure 2 with the following major techniques and optimizations.

Batched DRAM page pinning. Memory addresses passed from user space are all virtual and need to be translated into physical ones that the DMA can consume. Furthermore, to satisfy DMA requirements, the virtual-to-physical address mapping must remain valid and unchanged during the execution of the corresponding DMA copy. This can be done by calling the `pin_user_page` and the `dma_map_page` kernel functions. However, pinning user pages one by one incurs high overhead for bulk I/Os, which span across multiple pages. To lighten this overhead, we leverage the `pin_user_pages` function available in the recent Linux kernel (version 5.9) that pins all the pages belonging to a single I/O. Similarly, we apply the same optimization for `unpin_user_page` via the new `unpin_user_pages` function.

src/dst page pairing. A bulk I/O will be mapped to a list of DMA subtasks at 4KB page granularity, each of which requires to pair the addresses of the source and destination pages for preparing the DMA descriptor. If the two addresses are not aligned, to ensure the correctness of DMA execution, which assumes that copies take place within page boundary, we have to carefully match the capacity of `dst` pages and the content size of `src` pages so that cross-page copies can be avoided. However, this leads to doubling the number of DMA subtasks, as described in Section 3.2.1.

Here, we make a key observation that NVM is managed contiguously in the kernel, and thus the cross-page copies can be tolerated. We exploit this finding as when preparing the DMA descriptor, we specify the length of the corresponding subtask in a page aligned manner on the DRAM side. For instance, take the situation in Figure 7 assuming that the source and destination are DRAM and NVM, respectively. We take the first portion of the source (① of page#1), which will always be smaller or equal to a page but aligned on the right end, as the size

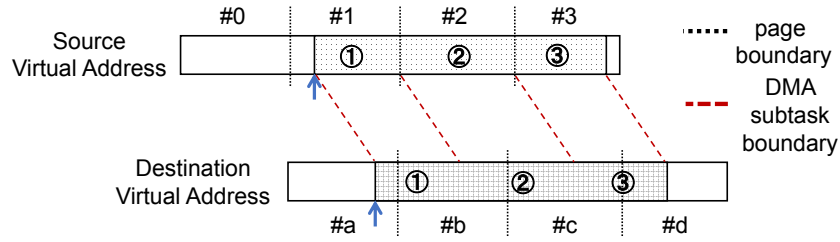


Fig. 7. Composition of DMA subtasks with halved numbers for a data movement job, enabled by the contiguous NVM address management, in comparison to Figure 3.

of the first DMA subtask. Thereafter, the size of the subsequent DMA subtasks will always be a page and aligned (② of page#2) except possibly for the last portion (③ of page#3), which will be page aligned on the left end. This enables us to reduce the number of DMA subtasks by half, in comparison to Figure 3.

Metadata buffer pre-allocation. DRAM space must be allocated with varying sizes to store the DMA request metadata, i.e., descriptors. The `scatterlist` structure is used to store the list of descriptors of DMA subtasks belonging to a single bulk I/O, where each item is typically 32 bytes. To accelerate memory allocation, we pre-allocate a fixed-size buffer to store this information prior to the execution of DMA copy. We set the buffer size to 4KB, which can accommodate DMA request metadata for 128 user pages (in total 512KB) at once.

Batch submission. Finally, to amortize the DMA subtask submission, considering that leveraging multiple channels performs no better than using a single channel (Section 3.2.2), we submit `scatterlist` in a batch to a single DMA channel assigned by Scheduler. This batched submission reduces the locking overhead for coordinating the concurrent accesses of the task queue associated with the DMA channel [20].

5.2 DDIO-aware Writes

As we discussed in Section 3.2.4, the configuration of DDIO has a strong influence on I/OAT write performance. Leveraging insights gained from the experiments, we devise strategies according to DDIO configurations. When DDIO is enabled, we perform flush operations after each write to minimize write amplification. In cases where addresses are not aligned to 64B, we first attempt to align them. To achieve this, we check the offsets of the source and destination addresses on 64B boundaries. If they are equal, we use the CPU to execute the copy of the left fringes. Subsequently, the remaining addresses are aligned and can be efficiently copied by the DMA. However, if the offsets are not equal and DDIO is disabled, we employ the CPU, instead of the DMA, to execute PM writes. We also recommend users to disable DDIO if they are certain that their addresses are (or can be) aligned to 64B boundaries. Conversely, if there is uncertainty regarding address alignment, enabling DDIO is advised.

5.3 DMA-CPU Cooperated Bulk Reads

With Simple-DMA, the application thread (CPU) submits requests to the DMA, which solely moves the data (see top part of Figure 8). However, as shown in Section 3.2.3, bulk reads could be made faster through DMA and CPU cooperation. Motivated by this, we design an optimized bulk read within Fastmove that is depicted by the lower part in Figure 8. Here, the application thread first splits the bulk read into two chunks, and then submits one chunk (#1) via the normal DMA path with optimizations mentioned in Section 5.1, followed by the other chunk (#2) being copied by the CPU. Upon completion of chunk#2, the corresponding CPU thread polls the status of the DMA. Finally, the execution of the target read completes when both the DMA and CPU finish their assigned chunks. This design not only improves the NVM read bandwidth but also hides the copy latency due to the CPU.

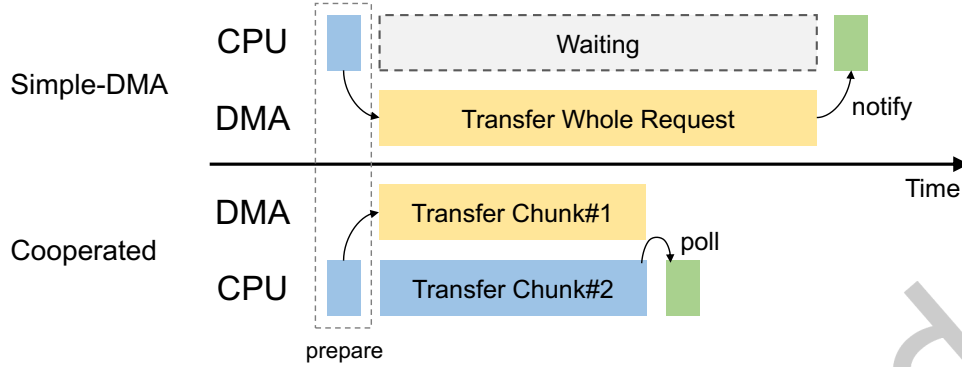


Fig. 8. The workflow of DMA-CPU cooperated reads.

While the optimized bulk read is a natural sharing of load, the challenge we face here is how to decide the loads that will go through the CPU and DMA. Chosen inappropriately, the gap between the execution time of CPU and DMA could lead to either waste of CPU cycles for polling the DMA status or lower DMA utilization. To balance their execution time, we set the chunk #1 and #2 size ratio to the ratio of the average single-threaded bandwidth on the CPU and DMA paths, which are monitored by our Scheduler.

5.4 Controlling and Scheduling

We design a light-weight Scheduler that outputs the proper memory copy path assignment plan for each I/O request going through the above Fastmove's APIs, distributes loads of bulk reads between CPU and DMA, and properly allocates DMA resources for offloaded tasks.

Initial configuration. Decision-making by Scheduler is driven by the DDIO configuration, the four pre-chosen I/O size thresholds for local/remote NVM reads/writes, beyond which DMA path should be involved for better performance, and the concurrency sweet spot M per DMA device, which corresponds to the maximal number of concurrent threads leading DMA to reach the peak bandwidth. In addition, Scheduler also monitors the following four variables: (1) C_i , which is used to keep track of the number of on-the-fly requests submitted to device i and that works as an indicator of the workload intensity level of that device; (2) S_i , which points to the next available DMA channel on the DMA device i ; and (3) B_C and B_D , that record the bandwidth dynamically consumed by the CPU and DMA, respectively.

Scheduling. Scheduler first inspects every I/O request to figure out the following parameters: the alignment of the request (A), the NUMA node id of the target NVM (N_p), the request type RW , the NUMA information LR , and I/O length L . A , RW and LR are both boolean values indicating *aligned/unaligned*, *read/write* and *local/remote*, respectively. Then, the path scheduling logic is straightforward as follows. Scheduler first checks if DDIO is disabled, A is true, and RW is false. If so, Scheduler chooses the CPU-only data path. Then, Scheduler compares the request length L to the DMA threshold, pointed by the pair of RW and LR , to identify bulk I/Os. For bulk I/Os, Scheduler chooses the DMA as long as the DMA device on node N_p is under its concurrent limit, i.e., $C_{N_p} < M$. If so, Scheduler chooses the next DMA channel associated with N_p 's DMA in a round-robin fashion (based on S_{N_p}) and updates the required resource variables, i.e., $C_{N_p} = C_{N_p} + 1$ and $S_{N_p} = (S_{N_p} + 1) \bmod G$. Otherwise, we fall back to the CPU-only data path. Additionally, we use B_C and B_D to derive the split ratio of bulk reads between the CPU and DMA by following the logic presented in Section 5.3.

Performance consideration. To minimize the overhead that may incur due to request processing, we make the following two design choices. First, instead of implementing the `Fastmove` logic as a centralized component for coordination, we provide the logic as a function, which runs at the `memcpy` caller side. This precludes inter-thread communication between I/O threads and Scheduler helping enhance performance. Second, coordination of concurrent access to globally shared variables like S_N and C_N adopt lightweight mechanisms such as atomic counters to further reduce overhead.

5.5 Implementation Details

We implement `Fastmove`¹ under the DMA framework [37] in Linux kernel 5.9 with 2417 lines of C code for its core logic.

Integration with NVM-based storage systems. We integrate `Fastmove` into three widely-adopted, DAX file systems, namely, NOVA [58], `ext4-DAX` [34], and `XFS-DAX` [35], where NOVA is tailored for hybrid DRAM-NVM settings, while the other two systems are more general and compatible with NVM. `Fastmove`'s transparent design leads to minimal changes to the above systems. Specifically, we introduce only 2 lines of code changes to both `ext4-DAX` and `XFS-DAX`, which simply replace the memory copy functions in `read_iter()` and `write_iter()` system calls with the APIs in Table 3. NOVA requires 2 additional changes to its `read()` and `write()` functions.

Though `Fastmove` enables NUMA NVM access by design, DAX file systems cannot naturally use NVM devices sitting across NUMA sockets. We address this problem by leveraging the Linux native device mapper [36], as shown in Figure 6. For the device mapper, similarly, only 2 lines in `dm_copy_from_iter` and `dm_copy_to_iter` functions need to be replaced. Note, however, that the current version of NOVA does not support the use of the device mapper. Therefore, we extend NOVA to work with the device mapper and its new code base can be found in `Fastmove`¹. With these minimal changes, `Fastmove` is able to transparently benefit many applications that run atop these three file systems.

Correctness guarantee. The use of DMA in `Fastmove` will not introduce any data inconsistencies compared to CPU-only data accesses. First, while not mentioned in any public documentation from the hardware vendor, Kalia et al. [24] experimentally show that I/OAT preserves ordering during execution. Second, `Fastmove` always monitors the execution status of parallel DMA subtasks and knows which set of pages failed to be copied even though these pages may not be consecutive. This slightly relaxed `memcpy` semantic is enough since (1) most applications including filesystems and databases have their own well-designed fault handling mechanism, which can leverage `Fastmove`'s fault reports to recover state correctly, and (2) in kernel, there are many strict checks to avoid copy failures, such as permission validation prior to copy execution. Thus, failures will be rare.

6 EVALUATION

6.1 Experimental Setup

We deploy our experiments on a physical server with two 20-core Xeon Gold 6248 processors and 192GB DRAM. This machine has two NUMA nodes, each connected with six Intel Optane PM chips (128GB each and 1.5TB in total). We evaluate `Fastmove` with both the Optane PM device and emulated NVM to demonstrate the generality of `Fastmove`. With Optane PM, we configure it to be interleaved within each NUMA node and under the App Direct mode, and use the Linux device mapper under its striped mode to enable cross-socket NVM accesses. For the NVM emulated experiments, we use 64GB DRAM to emulate an advanced NVM device with DRAM-like latency and bandwidth, which is significantly better than Optane PM, using a Linux built-in emulator [39]. Note that our evaluation primarily focuses on Optane PM, while the emulated NVM performance results are only presented in Section 6.3.7.

¹Publicly available at <https://github.com/fastmove-open/fastmove>

Table 4. Latencies (us) of CPU-only memory copying obtained by running FIO workloads

latency (us)		1 thread		2 threads		4 threads	
		local	remote	local	remote	local	remote
read	16K	5.0	6.5	5.2	6.6	5.4	6.9
	32K	9.2	12.2	9.5	12.4	10.0	12.9
	64K	17.6	23.2	18.1	23.6	19.1	24.6
	128K	34.1	45.3	35.3	46.1	37.3	48.0
	256K	67.1	89.1	69.6	90.7	73.6	94.8
write	16K	8.7	20.6	8.9	21.7	9.5	23.5
	32K	16.0	46.7	16.5	48.1	19.0	52.1
	64K	28.8	89.7	30.3	92.3	35.3	100.2
	128K	54.7	174.8	56.9	180.0	67.7	195.5
	256K	106.2	344.8	112.0	354.3	131.8	384.1

Baseline and configurations. We exercise NOVA, ext4-DAX, and XFS-DAX enhanced by Fastmove. Our natural baselines are these file systems with their memory copy operations going through the conventional CPU path, denoted by “CPU-only”. We use default configurations for both baselines.

Case study applications and workloads. We take three data-intensive applications MySQL, GraphWalker and Filebench, with no code changes, to transparently use Fastmove by simply running them atop the three slightly modified DAX file systems. To evaluate Fastmove’s benefits, we run experiments with the FIO microbenchmark [8] and a synthetic workload generated based on FIO, application workloads like the widely-adopted standard database workload TPC-C [5] and the file access workload fileserv [1], and three popular graph processing tasks, namely, Graphlet Concentration, Personalized PageRank and SimRank. The detailed configurations are presented in Section 6.3.

6.2 Microbenchmark Results

6.2.1 Latency Threshold Choices. To help figure out the read/write thresholds with different concurrency levels required to drive the memory copy path selection in Fastmove, we run the FIO workloads to evaluate both the original and modified NOVA file systems. Here, we generate read and write workloads with different I/O sizes ranging from 16KB to 256KB that are issued by 1 to 4 concurrent threads. We test both local and remote (cross-socket) NVM accesses, considering both alignment settings. We enable DDIO when address is not aligned, and disable it when address is aligned. In Figure 9, we show the normalized average latency of Fastmove against the CPU-only baseline. In addition, we also include the results of “Simple-DMA”, the baseline with DMA enabled but not highly optimized, to demonstrate the validity and effectiveness of Fastmove’s optimizations and our Fastmove. We also show the absolute latency numbers of the CPU-only baseline in Table 4, which helps recover the absolute latencies for all other system configurations from the normalized results in Figure 9. Note that, these results look exactly the same across the three file systems. Thus, we do not present the results of ext4-DAX and XFS-DAX.

As shown in Figures 9a, 9c, 9e, and 9g, across all exercised configurations, CPU-only delivers constantly lower read latency than Simple-DMA. In contrast, Fastmove visibly improves on the performance of Simple-DMA and introduces 1.20-3.09× speedups for various I/O sizes, even leading read requests with relatively small sizes to benefit from DMA. Compared to CPU-only, the turning points of Fastmove are uniformly 32KB across the 1, 2, and 4 threaded workloads. Including and beyond these turning points, Fastmove starts to observe a visible

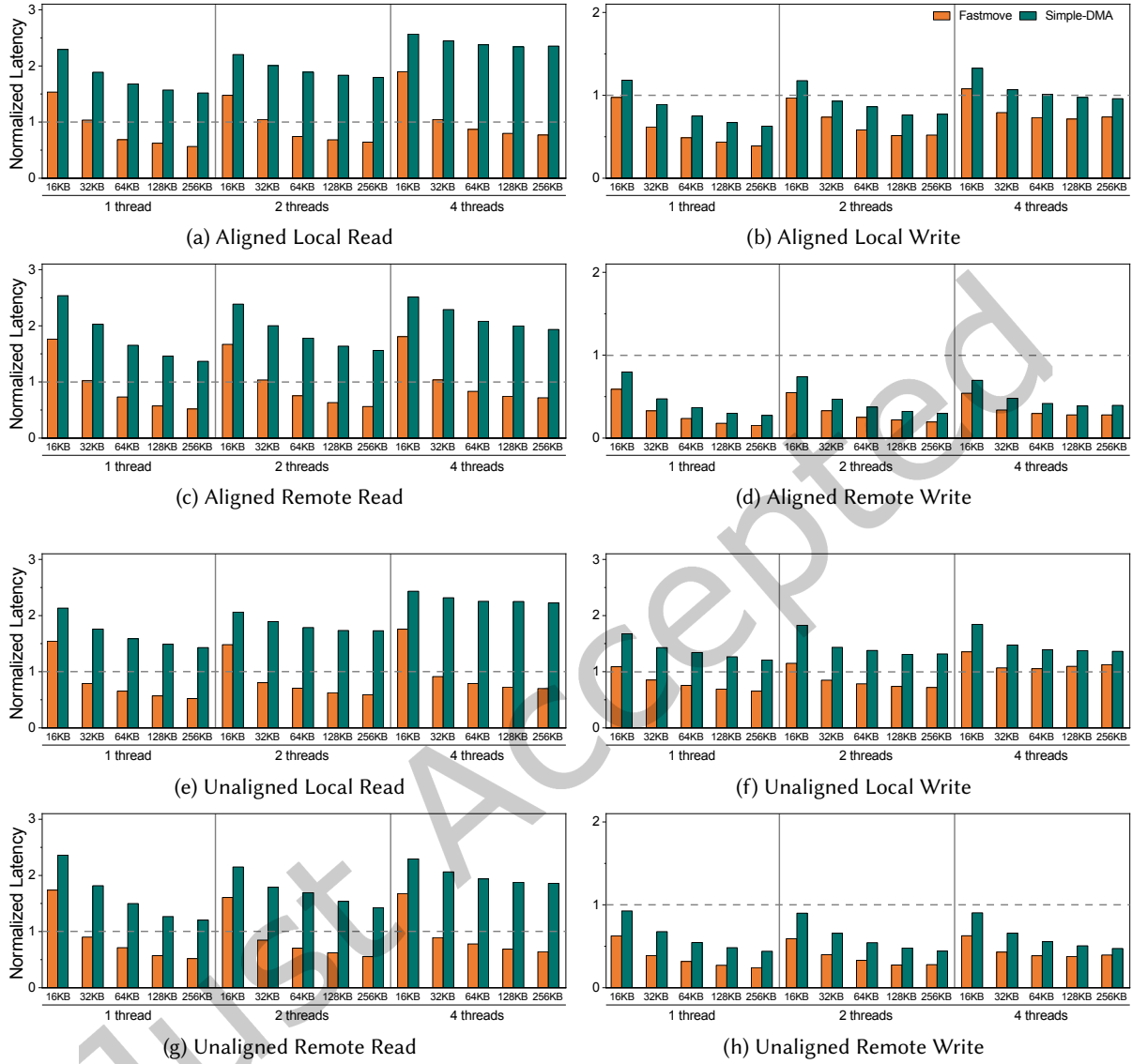


Fig. 9. The latency comparison between Fastmove and Simple-DMA, with 1,2,4-threaded FIO workloads, normalized to the latencies of CPU-only memory copying.

reduction in average request latency. For instance, Fastmove reduces the local read latency of CPU-only accesses by 13.0-25.6% for 64KB.

For writes, we observe larger improvements than reads. Figure 9b shows that for aligned local writes, Simple-DMA runs faster than CPU-only at 64KB, 128KB, and 128KB for the three concurrency settings, respectively. Fastmove dramatically improves Simple-DMA’s latency, dropping the turning points to 16KB, 16KB, and 32KB, respectively, for 1, 2, and 4 threads. With 2 threads, Fastmove achieves 36.9-49.0% and 26.3-48.6% reduction on

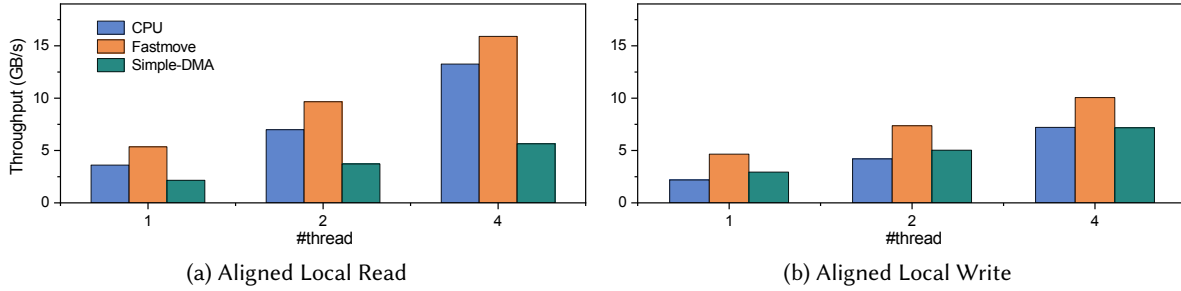


Fig. 10. The bandwidth comparison between CPU-only, Fastmove and Simple-DMA, with 1,2,4-threaded 64KB FIO workloads.

Table 5. P99 latency (us) comparison of aligned local read/write with batching enabled or disabled in Fastmove, corresponding to the same setting of 2-thread experiments in Figure 9.

size (KB)	read		write	
	batching	non-batching	batching	non-batching
16	8	13	10	11
32	9	11	14	19
64	16	21	23	30
128	27	37	46	55
256	49	72	80	87

average latency for I/Os at 32KB and beyond, compared to Simple-DMA and CPU-only, respectively. Figure 9f shows that in the case of unaligned local writes, Simple-DMA is unable to outperform CPU-only. Conversely, Fastmove drops the turning points to 32KB with 1 and 2 threads. With 2 threads, Fastmove achieves 40.7-45.3% and 14.9-28.1% reduction on average latency for I/Os at 32KB and beyond, compared to Simple-DMA and CPU-only, respectively. With 4 threads, Fastmove achieves comparable performance to CPU-only. The benefits of the two DMA variants further expand for remote writes, as shown in Figure 9d and 9h. First, they perform better than CPU-only for even 16KB. Second, the latency gap between the DMA usage and CPU-only becomes visibly larger, e.g., for 256KB aligned cross-socket I/O requests, Simple-DMA and Fastmove reduce latency by 75.3% and 86.1%, respectively, compared to CPU-only. Third, Fastmove significantly outperforms Simple-DMA by 40.7-48.1%, 65.9-73.7%, and 86.7-96.2% for 16KB, 32KB, and 64KB, respectively.

Finally, Table 5 illustrates the impact of the batched submission optimization on tail latency. We find that batching within Fastmove does not prolong, but rather, improves tail latency. For instance, with the same setting of 2-thread experiments in Figure 9, the P99 latency numbers in Table 5, indicating a 8.0-38.5% reduction, compared to the non-batching baseline. This is because Fastmove is not batching DMA subtasks across I/O requests from upper applications but is batching submissions of DMA subtasks that belong to a single request.

Bandwidth understanding We further investigate the data copy bandwidth consumption achieved by CPU-only, Fastmove and Simple-DMA. Here, Figure 10 shows the bandwidth statistics of the local read/write, aligned workload with 64KB I/O requests. With regard to read, Fastmove improves the total bandwidth over Simple-DMA and CPU-only by up to 159.5% and 48.1%, respectively. CPU-DMA cooperation plays a key role in the improvement. In more detail, DMA accounts for 40.2%, 44.5% and 52.8% of the total bandwidth for 1, 2 and 4 threads, respectively. For write, Fastmove improves the bandwidth over Simple-DMA and CPU-only by up to 58.0% and 112.6%, respectively. With 4 threads, Fastmove reaches the write bandwidth capacity of NVM for about 10GB/s.

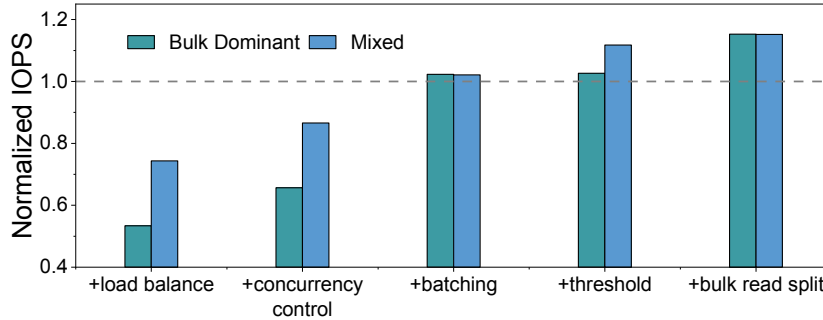


Fig. 11. Breakdown analysis of Fastmove with synthetic FIO workloads when gradually enabling optimizations. Throughput is normalized to the CPU-only baseline.

6.2.2 Breakdown Analysis. We use two synthetic FIO workloads to investigate the performance improvements introduced by each individual optimization within Fastmove. The bulk dominating workload contains I/Os with an average size of 256KB, while the mixed one has a mixture of bulk and small I/Os, ranging between 8KB and 256KB. For the two workloads, we use 6 concurrent threads to issue aligned local read or write requests to the underlying NOVA file system.

Figure 11 reports the normalized throughput numbers, which indicate that different workloads see different optimization sweet points. The direct usage of DMA with loads evenly distributed among channels leads to a 46.6% and 25.7% throughput drop for the bulk and mixed workloads, respectively, compared to CPU-only. This is because small I/Os do not benefit, yet still go through the DMA, and the associated DMA overheads have not yet been ameliorated. As we start to avoid overloading the DMA resources by adding the concurrent limit optimization (here, set to 4), Fastmove’s performance improves by 23.0% and 16.5% for the two workloads. The batching optimization makes Fastmove begin to outperform CPU-only, with a throughput increase of 55.8% and 17.9%. The latency threshold filtering further improves Fastmove’s performance by 0.3% and 9.5%, where the mixed workload observes larger improvements as this optimization avoids its small I/Os from paying the latency penalty of going through the DMA. Finally, the bulk read split design choice brings another 12.3% and 3.1% improvement. In the end, adding all these optimizations together brings a 1.15× improvement in throughput for the two workloads, compared to CPU-only.

6.3 Overall Performance

Next, we evaluate the positive impact of using Fastmove on the performance of real-world applications that introduce more complex characteristics than microbenchmarks such as non-uniformed I/O size distribution, computation-related cost, foreground and background processing division, etc.

6.3.1 Application Configurations. MySQL. We install MySQL version 5.7.33 with the default 16KB `innodb_page_size`. `innodb_buffer_pool_size` is set to half of the DRAM space, the recommended setting. We run the TPC-C workload with a read and write ratio of 1.78:1. For each run, we populate a 466GB database with 5000 Warehouses during the initialization phase and use 14 connections during the evaluation phase.

GraphWalker. GraphWalker [54] supports fast random walks on large graphs with a single machine. We exercise three common random walk algorithms, namely, Graphlet Concentration (Graphlet), Personalized PageRank (PPR) and SimRank (SR). We also follow GraphWalker to generate a Kron30 dataset using the Graph500 Kronecker [2], which consists of 1 billion vertices and 32 billion edges that take 638GB and 136GB of persistent

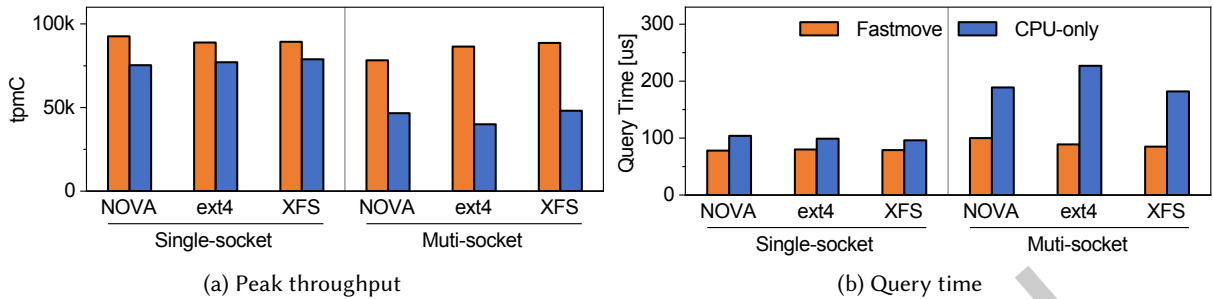


Fig. 12. Throughput (measured as tpmC) and query time achieved by running TPC-C against MySQL.

media space to store its original text data and the compressed CSR data, respectively. We use the GraphWalker default configurations.

Fileserver. We exercise the predefined workload, fileserver, within the Filebench framework [1]. It uses different numbers of concurrent threads (1, 2, 4, 8) to issue I/Os with various sizes presented in Table 1.

Enabling/disabling THP. We test MySQL and Fileserver without using transparent huge pages (THP), resulting in non-contiguous memory copies. This is recommended by the MySQL official site as THP introduces negative performance impacts on random memory accesses with small I/O sizes. Contrary, we enable THP for GraphWalker with contiguous copies, as its workloads are read-dominating and bulk-sized.

DDIO configuration. Among the applications above, Filebench is the only one that involves unaligned writes to NVM. We add a modified Filebench that only allocates aligned IO buffer as a comparison. We keep DDIO enabled for the unmodified Filebench and disable it for all other applications.

6.3.2 Fastmove Configurations. To run the above applications atop Fastmove with maximized performance, we need to set two parameters, namely, I/O size thresholds for prioritizing requests to go through the DMA path, and concurrent limit that avoids I/OAT bandwidth over-provision. Here, we give the detailed steps of the manual parameter setting process, which can be further automated in future.

I/O size threshold We begin with running the FIO benchmark to obtain the latency numbers of each configuration, combining local/remote and read/write, for both Fastmove and CPU-only with I/O size varying from 16KB to 256KB and the thread number varying from 1 to 4. The two ranges can be extended if the DMA hardware becomes more powerful. Then, we determine the thresholds for each configuration by simply finding the latency turning points, where Fastmove performs faster than CPU-only.

Concurrency limit We determine the bandwidth capacity of I/OAT by configuring Fastmove to use DMA only and running the FIO workload with the 64KB I/O size and varied thread numbers starting from 1. Then, we plot a bandwidth capacity figure to find the saturating thread number, which is the minimum number making the DMA bandwidth reach its maximum value.

6.3.3 MySQL Enhancement. Single-socket results. First, we consider the performance within a single socket, where application threads and PM are located under socket 0. Figure 12a shows the throughput comparison (officially measured as tpmC by TPC-C) between CPU-only and Fastmove execution of MySQL. Across all settings, Fastmove consistently delivers better performance than CPU-only, and the improvements associated with different underlying file systems look similar. For instance, Fastmove introduces 1.23 \times , 1.15 \times , and 1.13 \times speedups of peak throughput over the CPU-only baseline across NOVA, ext4-DAX, and XFS-DAX, respectively.

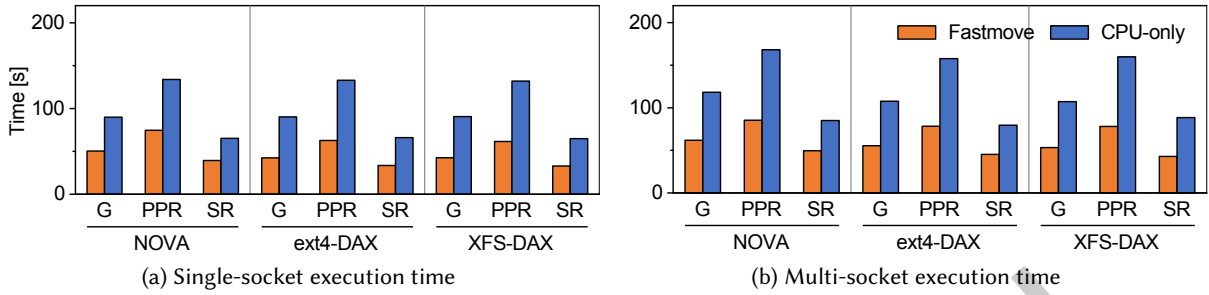


Fig. 13. Execution times running Graphlet (G), PPR and SR over GraphWalker with NOVA/ext4-DAX/XFS-DAX.

Figure 12b reports the corresponding average query time results. Consistent with the throughput results, Fastmove reduces the average latency of CPU-only by 17.7-25.0%.

To understand the source of improvements, we profile the I/O distribution of the TPC-C workload. As shown in Table 1, almost all of its read requests are smaller than 32KB. As this is below the 32KB threshold, the vast majority of read requests go through the ordinary CPU-only path in Fastmove. As a consequence, the performance improvements here are driven by the 90.9% of bulk writes beyond 16KB, which correspond to the logging activities handled by the 4 background flush threads. To conclude, Fastmove indeed choose proper memory copy paths for I/O with varied sizes, and I/OAT DMA does alleviate the NVM accessing data stalls.

Multi-socket results. Next, we explore the performance implications under two sockets, where we replicate the above experiments by evenly distributing application threads to two CPUs and spreading the data on all 12 PMs via Linux device mapper under its striped mode.

Figure 12a shows the absolute throughput numbers achieved by CPU-only with two sockets decrease by 38.1-48.1%, compared to the single-socket counterparts. This is because performance degrades for cross-socket memory copy operations as depicted in Figure 9. In contrast, Fastmove observes lighter negative impact of cross-socket NVM access with only 0.8-15.5% drop in peak throughput. Fastmove significantly outperforms the CPU-only baseline, introducing 1.68-2.16 \times tmpC improvements. Additionally, in Figure 12b, Fastmove brings a significant latency reduction of 47.1-60.8%. Contrary to the single-socket results, we see that Fastmove's improvements over CPU-only expand. This is because the threshold for remote reads drops to 16KB, which allows for cross-socket NVM reads to take advantage of the DMA if DMA usage is not full, and also the DMA benefits for remote reads and writes are larger than those for local ones.

6.3.4 GraphWalker Enhancement. Single-socket results. Figure 13a shows the execution times of three graph analytic tasks over GraphWalker. The performance of the CPU-only baseline looks similar across different file systems, and so does our Fastmove. However, we observe that Fastmove significantly reduces the execution time over CPU-only, despite the fact that the graph analytic jobs are read-only workloads towards the underlying data systems. More specifically, Fastmove introduces 1.78-2.13 \times , 1.79-2.14 \times , and 1.65-1.97 \times speedups for Graphlet, PPR, and SR, respectively. The significant improvements come from the dominating bulk read I/Os as shown in Table 1.

Multi-socket results. Consistent with the above TPC-C results, the improvements of Fastmove for the graph analytic workloads become larger compared to the single-socket counterparts. Figure 13b depicts that Fastmove brings 1.91-2.01 \times , 1.97-2.05 \times , and 1.71-2.06 \times execution time speedups for Graphlet, PPR and SR running in GraphWalker, respectively, across three different NVM-based file systems.

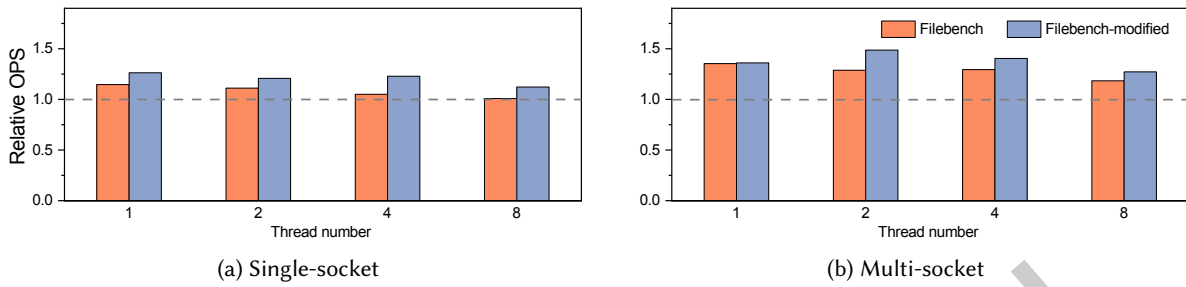


Fig. 14. Operations per second (OPS) of Filebench and Filebench-modified running the Fileserver workload with Fastmove-powered NOVA, normalized to the OPS of CPU-only.

6.3.5 Fileserver Enhancement. Single-Socket results. Figure 14a shows the normalized operations per second (OPS) achieved by different numbers of Filebench threads executing the fileserver workload atop of NOVA. Filebench-modified consistently performs better than vanilla Filebench across different thread numbers. Fastmove succeeds to speedup both Filebench and Filebench-modified by 1.01-1.15 \times and 1.12-1.26 \times , respectively. With 8 threads, Fastmove achieves parity with CPU-only on Filebench. This is consistent with the FIO results that Fastmove performs worse than CPU-only in unaligned local writes with 4 or more concurrent users. As shown in Table 1, the benefits of Fastmove comes from its acceleration to 8.3% reads (over 32KB) and 17.8% writes (over 16KB).

Multi-socket results. Multi-socket advantages can also be found in Filebench. As shown in Figure 14b, the improvements of Fastmove becomes larger with the multi-socket setting. Fastmove introduces 1.18-1.35 \times and 1.27-1.48 \times speedups for Filebench and Filebench-modified, respectively.

6.3.6 CPU Consumption Improvement. Finally, we explore another possible benefit of using Fastmove, which is the CPU consumption improvement. Here, we measure the CPU cycles spent in moving data between DRAM-NVM and processing the application logic. For MySQL TPC-C workload, Fastmove reduces its data movement CPU usage from 62% to 39% and from 90% to 28% for single-socket and multi-socket settings, respectively. We also observe a significant increase in its utime. This is because the saved CPU cycles from data movement are used to perform useful work, leading to improved throughput numbers (presented in Section 6.3.3). Unlike this, for GraphWalker, Fastmove’s CPU usage improvement seems little. For instance, Fastmove reduces its CPU usage for data movement by up to 5%. This is because workloads with GraphWalker benefits largely by the DMA-CPU cooperated bulk read optimization, which requires CPU involvement.

6.3.7 Emulated NVM Performance. We deploy NOVA on emulated NVM, replicate the experiments for Figure 9, and report the latency comparison results between CPU-only, Simple-DMA, and Fastmove in Figure 15. Fastmove outperforms CPU-only for local reads and writes with I/O sizes of 16KB and beyond. The benefits observed are larger than those corresponding to experiments with Optane PM (Figure 9). This is because emulated NVM is of DRAM-like read and write performance. Considering the association cost in Section 3.2.1, the dominating DMA copy execution step becomes faster, leading to visible end-to-end read/write latency improvements. Also, this implies that the time cost of NVM device access plays a key role in assigning DMA resources, i.e., the performance turning point based on I/O size decreases when NVM device performance improves, and vice versa. Furthermore, we find that the DMA bandwidth within Fastmove saturates when concurrency reaches 4 threads, exactly the same as the Optane PM experiments. This is because both the emulated NVM and Optane based experiments

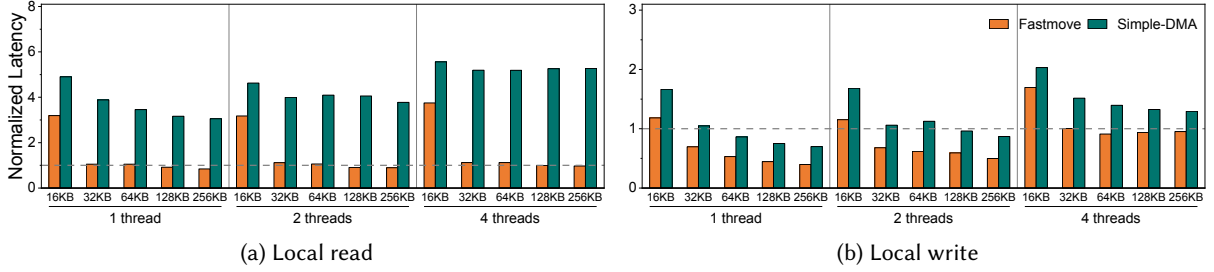


Fig. 15. The latency comparison between Fastmove and Simple-DMA, with 1,2,4-threaded FIO workloads, normalized to the latencies of CPU-only memory copying when deploying NOVA on in-kernel NVM-emulator.

make use of I/OAT DMA, and under both cases, DMA bandwidth capacity is lower than Optane PM and emulated NVM.

7 DISCUSSION

In this section, we explore potential enhancements in Fastmove’s design and contemplate its evolution to accommodate emerging, swifter DMA/memory devices.

Automatic profiling. Currently, manual profiling is necessary to determine the appropriate thresholds and concurrency limits. With each new machine setup, such as new CPUs and different number of NVM modules, these empirical parameters may vary. Thus, automated profiling would not only streamline this labor-intensive process but also render Fastmove more versatile and user-friendly. A practical design of the profiling procedure would follow the steps that were laid out in Section 6.2.1.

Resource sharing. In scenarios where multiple I/O intensive applications are deployed in a single machine, the existing design of Fastmove lacks provisions for intervention or fairness between these applications. However, due to our commitment to transparency, it is difficult for Fastmove alone to coordinate the applications behind the VFS interface. A possible solution is to integrate Fastmove into alternative memory bandwidth regulation systems such as MT² [61].

Embracing future DMA hardware. Hardware has been constantly evolving to become faster with more features. For instance, Intel Data Streaming Accelerator (DSA) [49], the successor to I/OAT DMA, is equipped with larger DMA bandwidth and many new features, including address sharing through IOMMU, and per-request DDIO control [32]. To cope with the new hardware trends, Fastmove may also need change to exploit their potentials. Referring to the results in Section 6.3.7, even with better hardware capacities, most of Fastmove’s design choices will remain valid. In addition, new features introduced to DMA devices will need extra consideration. First, let’s focus on the DMA bandwidth and parallelism. In I/OAT’s implementation, different channels share the bandwidth. However, if a future DMA provides channels with separate, independent bandwidth, Fastmove will need to design a channel selection algorithm for better load balancing and resource utilization. Second, we shift our attention to DMA performance trade-offs. As shown in Table 2, pin overhead is non-negligible in DMA copy cost. If DMA can share addresses with cores through IOMMU, this overhead can be eliminated. Furthermore, the support of this new IOMMU feature allows virtual address in descriptors, and imposes no restrictions on address alignment. Therefore, we will not need to process in page granularity, and each request can fit into a single DMA subtask, reducing the pairing and submitting overhead. Nevertheless, this benefit may come at the cost of extra overhead for across-core communication upon page faults. Third, If DDIO can be controlled on per-request basis, we can assign the best-fit DDIO configuration for each request based on the lessons learned from this paper rather than a global static configuration. It can potentially improve the performance for DMA shared usage. In

summary, new hardware advances likely make DMA overheads and benefits be completely different from what we present in this paper, potentially leading to the failure of some of Fastmove's current design. Thus, a novel approach would be required to utilize future DMA efficiently.

CXL implications. The CXL technology can extend the CPU's address space, enabling it to access external storage via memory interface [49]. A CXL-native accelerator will definitely help data movements in such environments. To fully exploit its potential, one still needs to study the characteristics of that accelerator, by largely follow the study we conducted against I/OAT and NVM in Section 3. To name a few points, the study may involve latency evaluation to determine I/O size thresholds for choosing hardware-assisted data path, understanding the bandwidth capacity and impact of parallelism to prevent resource over-provision. Additionally, the CPU-DMA cooperative data movement may be also an effective optimization. In summary, most of the study principles and optimization proposals remain valid. We also realize that the CXL accelerator can come with fundamentally new features that I/OAT does not have, e.g., pages no longer needs to be pinned, which can inspire new use cases, open door for new optimization opportunities, and introduce technical challenges for future research.

8 RELATED WORK

I/OAT usage. Previous studies have used I/OAT to offload memcpy operations that move data from DRAM to DRAM [51, 52] as well as to improve network bandwidth with lower CPU utilization in data center environments [28, 53]. Unlike these, our study stands to speed up data movement between DRAM and NVM, where the interaction between I/OAT and hybrid memory architectures is more complex and its acceleration demands careful system design. Most recent work have included I/OAT as a minor optimization for data movement in NVM-based systems with special purposes ranging from log replication [7] to memory migration [45].

Unlike them, Fastmove is a general system to make use of on-chip DMA to address the inefficiencies (e.g., lower bandwidth or extra CPU overhead) introduced by CPU-only accesses to NVM for bulk, storage-facing I/Os, which has been observed to be a critical performance limitation in combined use of Optane PM and Intel processors. These systems can take advantage of Fastmove with little effort. Kalia et al. [24] present a number of optimizations for efficient remote NVM accesses via network, which includes an initial attempt to use I/OAT to improve single-core RPC performance for bulk remote NVM writes. This work is orthogonal to ours.

DDIO study. Several studies have also delved into the characteristics of DDIO. Farshin et al. [14] investigate the implementation details of DDIO and analyze its attributes in the context of network packet transmission. This work is distinct from ours in that we concentrate on the interplay between I/OAT DMA, DDIO, and NVM. Kalia et al. [24] finds that disabling DDIO increases peak bandwidth of I/OAT DMA when writing to NVM. Our findings in Section 3.2.4 show that this conclusion is limited to large I/Os and the WOF configuration. Furthermore, we study the impact of address alignment, which has not been considered before.

NVM-related studies and systems. There is a large body of work focusing on the analysis of the basic performance characteristics of using NVM [12, 16, 56, 60, 61]. The rich findings from these studies have spawned numerous studies for re-designing scalable and high performance data structures[27, 31, 55], file systems[22, 33, 58, 62] and key-value stores[9, 30]. Our work extends the existing study by incorporating the interaction between NVM and DMA, and complements the prior NVM-based systems as they can benefit from either the general design or the real implementation of Fastmove to alleviate data stalls. OdinFS [62] decouples application threads from the background NVM access threads and additionally parallelizes NVM accesses across sockets. Its NVM threads can benefit from Fastmove and its integration will be explored in the future.

Tiered memory systems. Fastmove handles more complex I/O patterns than those in tiered memory. In addition, Fastmove is implemented in the kernel with simple APIs. Therefore, Fastmove could be directly used in tiered memory systems. In fact, we have successfully adapted Nimble [59] to transparently use Fastmove through simple API replacement. However, through preliminary evaluations, we find that the DMA, in particular I/OAT,

may not be a good option for improving page migration in tiered memory. This is because the DMA bandwidth is easily overwhelmed by the workload. Therefore, `Fastmove` does not deliver any significant improvement over `Nimble-DMA` [59], a Linux patch that adapts `Nimble` to use I/OAT.

Zero-copy technologies. Another line of work on PM attempts to move data management from kernel space to user space to eliminate data copies along the I/O path. For instance, the memory mapped file I/O (e.g., the `mmap` system call) is enabled such that users may access files in the same way as memory data [57]. However, `mmap`-based solutions may incur high overhead due to page faults [10, 29] and may have to have applications handle data persistence and reliability on their own [40, 42]. Yet another line of work leverages kernel by-pass I/O interfaces such as `SPDK` and `PMDK` [4] to avoid the use of the complicated OS I/O stack [48]. However, the performance gains come at the price of substantial effort for re-writing the I/O handling part of the target applications.

In contrast, our work demonstrates better applicability since there is no code change required to run existing applications atop `Fastmove`, as long as they use kernel file systems. Moreover, it is possible to extend our design to handle memory copy operations in user space, where these operations may have an even bigger impact on the overall performance compared to their counterparts in kernel space. This is because by bypassing the kernel, memory copying will contribute to a larger portion of the end-to-end access performance.

Misc. `DaeMon` [15] investigates the data movement bottleneck in disaggregated memory, and proposes adaptive migration granularity for shorter latency and increasing bandwidth utilization. There is a chance to combine `Fastmove` and memory disaggregation, which can spawn new research directions in future.

9 CONCLUSION

In this paper, we first study the DRAM-NVM data movement problem and then propose and implement `Fastmove`, a general engine that exploits the on-chip DMA technology. With a clean abstraction and transparent design, applications can use `Fastmove` via slightly-modified file systems with no further changes. Experimental results with industry-standard workloads on `MySQL` and popular random walk algorithms on `GraphWalker` highlight that `Fastmove` brings significant benefits such as peak throughput increase, execution time reduction, and CPU consumption savings.

ACKNOWLEDGMENTS

We sincerely thank all anonymous reviewers for their insightful feedback. This work was supported in part by National Nature Science Foundation of China under the grant No. 62141216 and U.S. National Science Foundation under the grant No. 2312785. Cheng Li is the corresponding author.

REFERENCES

- [1] 2023. Filebench. <https://github.com/filebench/filebench>. [Online; accessed Jan-2023].
- [2] 2023. Graph500. <https://graph500.org/>. [Online; accessed Jan-2023].
- [3] 2023. MySQL. <https://github.com/mysql>. [Online; accessed Jan-2023].
- [4] 2023. PMDK. <https://github.com/pmemp/pmdk>. [Online; accessed Jan-2023].
- [5] 2023. TPC Benchamrk C. <http://tpc.org/tpcc/>. [Online; accessed Jan-2023].
- [6] Hiroyuki Akinaga and Hisashi Shima. 2010. Resistive random access memory (ReRAM) based on metal oxides. *Proc. IEEE* 98, 12 (2010), 2237–2251.
- [7] Thomas E Anderson, Marco Canini, Jongyul Kim, Dejan Kostić, Youngjin Kwon, Simon Peter, Waleed Reda, Henry N Schuh, and Emmett Witchel. 2020. Assise: Performance and Availability via Client-local NVM in a Distributed File System. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. 1011–1027.
- [8] Jens Axboe. 2023. FIO. <https://github.com/axboe/fio>. [Online; accessed Jan-2023].
- [9] Lawrence Benson, Hendrik Makait, and Tilmann Rabl. 2021. Viper: An Efficient Hybrid PMem-DRAM Key-Value Store. *Proc. VLDB Endow.* 14, 9 (may 2021), 1544–1556. <https://doi.org/10.14778/3461535.3461543>

- [10] Jungsik Choi, Jiwon Kim, and Hwansoo Han. 2017. Efficient Memory Mapped File I/O for In-Memory File Systems. In *9th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 17)*. Santa Clara, CA. <https://www.usenix.org/conference/hotstorage17/program/presentation/choi>
- [11] CXL Consortium. 2022. Compute Express Link: The Breakthrough CPU-to-Device Interconnect. <https://www.computeexpresslink.org/>. [Online; accessed Jan-2023].
- [12] Björn Daase, Lars Jonas Bollmeier, Lawrence Benson, and Tilmann Rabl. 2021. Maximizing Persistent Memory Bandwidth Utilization for OLAP Workloads. In *Proceedings of the 2021 International Conference on Management of Data (Virtual Event, China) (SIGMOD/PODS '21)*. New York, NY, USA, 339–351. <https://doi.org/10.1145/3448016.3457292>
- [13] Subramanya R. Dulloor, Amitabha Roy, Zheguang Zhao, Narayanan Sundaram, Nadathur Satish, Rajesh Sankaran, Jeff Jackson, and Karsten Schwan. 2016. Data Tiering in Heterogeneous Memory Systems. In *Proceedings of the Eleventh European Conference on Computer Systems (London, United Kingdom) (EuroSys '16)*. Association for Computing Machinery, New York, NY, USA, Article 15, 16 pages. <https://doi.org/10.1145/2901318.2901344>
- [14] Alireza Farshin, Amir Roozbeh, Gerald Q Maguire Jr, and Dejan Kostić. 2020. Reexamining Direct Cache Access to Optimize {I/O} Intensive Applications for Multi-hundred-gigabit Networks. In *2020 USENIX Annual Technical Conference (USENIX ATC 20)*. 673–689.
- [15] Christina Giannoula, Kailong Huang, Jonathan Tang, Nectarios Koziris, Georgios Goumas, Zeshan Chishti, and Nandita Vijaykumar. 2023. DaeMon: Architectural Support for Efficient Data Movement in Fully Disaggregated Systems. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* 7, 1 (2023), 1–36.
- [16] Shashank Gugnani, Arjun Kashyap, and Xiaoyi Lu. 2020. Understanding the Idiosyncrasies of Real Persistent Memory. *Proc. VLDB Endow.* 14, 4 (Dec. 2020), 626–639. <https://doi.org/10.14778/3436905.3436921>
- [17] Frank T Hady, Annie Foong, Bryan Veal, and Dan Williams. 2017. Platform storage performance with 3D XPoint technology. *Proc. IEEE* 105, 9 (2017), 1822–1833.
- [18] Intel. 2023. Intel I/O Acceleration Technology. <https://www.intel.com/content/www/us/en/wireless-network/accel-technology.html>. [Online; accessed Jan-2023].
- [19] Intel. 2023. Intel Vtune Profiler. <https://www.intel.com/content/www/us/en/developer/tools/oneapi/vtune-profiler.html#gs.3f5fmb>. [Online; accessed Jun-2023].
- [20] Dave Jiang. [n. d.]. libnvdimm: add DMA supported blk-mq pmem driver. <https://lore.kernel.org/linux-nvdimm/150412628764.69288.12074115435918322858.stgit@djiang5-desk3.ch.intel.com/#r>. [Online; accessed Jan-2023].
- [21] Myoungsoo Jung. 2022. Hello Bytes, Bye Blocks: PCIe Storage Meets Compute Express Link for Memory Expansion (CXL-SSD). In *Proceedings of the 14th ACM Workshop on Hot Topics in Storage and File Systems (Virtual Event) (HotStorage '22)*. Association for Computing Machinery, New York, NY, USA, 45–51. <https://doi.org/10.1145/3538643.3539745>
- [22] Rohan Kadekodi, Saurabh Kadekodi, Soujanya Ponnappalli, Harshad Shirwadkar, Gregory R. Ganger, Aasheesh Kolli, and Vijay Chidambaram. 2021. WineFS: A Hugepage-Aware File System for Persistent Memory That Ages Gracefully. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles (Virtual Event, Germany) (SOSP '21)*. Association for Computing Machinery, New York, NY, USA, 804–818. <https://doi.org/10.1145/3477132.3483567>
- [23] Rohan Kadekodi, Se Kwon Lee, Sanidhya Kashyap, Taesoo Kim, Aasheesh Kolli, and Vijay Chidambaram. 2019. SplitFS: Reducing Software Overhead in File Systems for Persistent Memory. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles (Huntsville, Ontario, Canada) (SOSP '19)*. New York, NY, USA, 494–508. <https://doi.org/10.1145/3341301.3359631>
- [24] Anuj Kalia, David Andersen, and Michael Kaminsky. 2020. Challenges and Solutions for Fast Remote Persistent Memory Access. In *Proceedings of the 11th ACM Symposium on Cloud Computing (Virtual Event, USA) (SoCC '20)*. New York, NY, USA, 105–119. <https://doi.org/10.1145/3419111.3421294>
- [25] Sudarsun Kannan, Nitish Bhat, Ada Gavrilovska, Andrea Arpaci-Dusseau, and Remzi Arpaci-Dusseau. 2018. Redesigning LSMs for Nonvolatile Memory with NoveLSM. In *2018 USENIX Annual Technical Conference (USENIX ATC 18)*. 993–1005.
- [26] Yoshihisa Kato, Yukihiko Kaneko, Hiroyuki Tanaka, Kazuhiro Kaibara, Shinzo Koyama, Kazunori Isogai, Takayoshi Yamada, and Yasuhiro Shimada. 2007. Overview and future challenge of ferroelectric random access memory technologies. *Japanese Journal of Applied Physics* 46, 4S (2007), 2157.
- [27] Ana Khorguani, Thomas Ropars, and Noel De Palma. 2022. ResPCT: Fast Checkpointing in Non-Volatile Memory for Multi-Threaded Applications. In *Proceedings of the Seventeenth European Conference on Computer Systems (Rennes, France) (EuroSys '22)*. Association for Computing Machinery, New York, NY, USA, 525–540. <https://doi.org/10.1145/3492321.3519590>
- [28] Jongyul Kim, Insu Jang, Waleed Reda, Jaeseong Im, Marco Canini, Dejan Kostić, Youngjin Kwon, Simon Peter, and Emmett Witchel. 2021. LineFS: Efficient SmartNIC Offload of a Distributed File System with Pipeline Parallelism. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles (Virtual Event, Germany) (SOSP '21)*. New York, NY, USA, 756–771. <https://doi.org/10.1145/3477132.3483565>
- [29] Juno Kim, Yun Joon Soh, Joseph Izraelevitz, Jishen Zhao, and Steven Swanson. 2020. SubZero: Zero-Copy IO for Persistent Main Memory File Systems. In *Proceedings of the 11th ACM SIGOPS Asia-Pacific Workshop on Systems (Tsukuba, Japan) (APSys '20)*. New York, NY, USA, 1–8. <https://doi.org/10.1145/3409963.3410489>

- [30] Wonbae Kim, Chanyeol Park, Dongui Kim, Hyeongjun Park, Young ri Choi, Alan Sussman, and Beomseok Nam. 2022. ListDB: Union of Write-Ahead Logs and Persistent SkipLists for Incremental Checkpointing on Persistent Memory. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. USENIX Association, Carlsbad, CA, 161–177. <https://www.usenix.org/conference/osdi22/presentation/kim>
- [31] Wook-Hee Kim, R. Madhava Krishnan, Xinwei Fu, Sanidhya Kashyap, and Changwoo Min. 2021. PACTree: A High Performance Persistent Range Index Using PAC Guidelines. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles (Virtual Event, Germany) (SOSP '21)*. New York, NY, USA, 424–439. <https://doi.org/10.1145/3477132.3483589>
- [32] Reese Kuper, Ipoom Jeong, Yifan Yuan, Jiayu Hu, Ren Wang, Narayan Ranganathan, and Nam Sung Kim. 2023. A Quantitative Analysis and Guideline of Data Streaming Accelerator in Intel 4th Gen Xeon Scalable Processors. *arXiv preprint arXiv:2305.02480* (2023).
- [33] Ruibin Li, Xiang Ren, Xu Zhao, Siwei He, Michael Stumm, and Ding Yuan. 2022. ctFS: Replacing File Indexing with Hardware Memory Translation through Contiguous File Allocation for Persistent Memory. In *20th USENIX Conference on File and Storage Technologies (FAST 22)*. USENIX Association, Santa Clara, CA, 35–50. <https://www.usenix.org/conference/fast22/presentation/li>
- [34] Linux. 2014. Add support for NV-DIMMs to ext4. <https://lwn.net/Articles/613384/>. [Online; accessed Jan-2023].
- [35] Linux. 2015. xfs: DAX support. <https://lwn.net/Articles/635514/>. [Online; accessed Jan-2023].
- [36] Linux. 2023. Device Mapper. <https://www.kernel.org/doc/Documentation/device-mapper/>. [Online; accessed Jan-2023].
- [37] Linux. 2023. DMAEngine framework. <https://www.kernel.org/doc/Documentation/driver-api/dmaengine/>. [Online; accessed Jan-2023].
- [38] Youyou Lu, Jiwu Shu, Youmin Chen, and Tao Li. 2017. Octopus: an RDMA-enabled Distributed Persistent Memory File System. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*. USENIX Association, Santa Clara, CA, 773–785. <https://www.usenix.org/conference/atc17/technical-sessions/presentation/lu>
- [39] Maciej Maciejewski. 2016. How to emulate Persistent Memory. <https://pmem.io/blog/2016/02/how-to-emulate-persistent-memory/>. [Online; accessed Jan-2023].
- [40] Ian Neal, Gefei Zuo, Eric Shiple, Tanvir Ahmed Khan, Youngjin Kwon, Simon Peter, and Baris Kasikci. 2021. Rethinking File Mapping for Persistent Memory. In *19th USENIX Conference on File and Storage Technologies (FAST 21)*. 97–111. <https://www.usenix.org/conference/fast21/presentation/neal>
- [41] Philip Ng. 2019. Accelerating Intra-Host PVRDMA Storage Traffic in a Future Dell AMD Server. Talk at VMWorld 2019. [Online; accessed Jan-2023].
- [42] Anastasios Papagiannis, Manolis Marazakis, and Angelos Bilas. 2021. Memory-Mapped I/O on Steroids. In *Proceedings of the Sixteenth European Conference on Computer Systems*. New York, NY, USA, 277–293. <https://doi.org/10.1145/3447786.3456242>
- [43] Jonathan Prout. 2022. Expanding Beyond Limits With CXL-based Memory. [Online; accessed Jan-2023].
- [44] Simone Raoux, Geoffrey W Burr, Matthew J Breitwisch, Charles T Rettner, Y-C Chen, Robert M Shelby, Martin Salinga, Daniel Krebs, S-H Chen, H-L Lung, et al. 2008. Phase-change random access memory: A scalable technology. *IBM Journal of Research and Development* 52, 4.5 (2008), 465–479.
- [45] Amanda Raybuck, Tim Stamler, Wei Zhang, Mattan Erez, and Simon Peter. 2021. HeMem: Scalable Tiered Memory Management for Big Data Applications and Real NVM. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles (Virtual Event, Germany) (SOSP '21)*. New York, NY, USA, 392–407. <https://doi.org/10.1145/3477132.3483550>
- [46] Thomas Rueckes. 2011. High density, high reliability carbon nanotube NRAM. In *Flash Memory Summit*.
- [47] Stackoverflow. 2022. Why are SIMD instructions not used in kernel? <https://stackoverflow.com/questions/46677676/why-are-simd-instructions-not-used-in-kernel>. [Online; accessed Jan-2023].
- [48] Timothy Stamler, Deukyeon Hwang, Amanda Raybuck, Wei Zhang, and Simon Peter. 2022. zIO: Accelerating IO-Intensive Applications with Transparent Zero-Copy IO. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. USENIX Association, Carlsbad, CA, 431–445. <https://www.usenix.org/conference/osdi22/presentation/stamler>
- [49] Yan Sun, Yifan Yuan, Zeduo Yu, Reese Kuper, Chihun Song, Jinghan Huang, Houxiang Ji, Siddharth Agarwal, Jiaqi Lou, Ipoom Jeong, Ren Wang, Jung Ho Ahn, Tianyin Xu, and Nam Sung Kim. 2023. Demystifying CXL Memory with Genuine CXL-Ready Systems and Devices. In *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture (, Toronto, ON, Canada,) (MICRO '23)*. Association for Computing Machinery, New York, NY, USA, 105–121. <https://doi.org/10.1145/3613424.3614256>
- [50] AA Tulapurkar, Y Suzuki, A Fukushima, H Kubota, H Maehara, K Tsunekawa, DD Djayaprawira, N Watanabe, and S Yuasa. 2005. Spin-torque diode effect in magnetic tunnel junctions. *Nature* 438, 7066 (2005), 339–342.
- [51] K. Vaidyanathan, L. Chai, W. Huang, and D. K. Panda. 2007. Efficient asynchronous memory copy operations on multi-core systems and I/OAT. In *2007 IEEE International Conference on Cluster Computing*. 159–168. <https://doi.org/10.1109/CLUSTER.2007.4629228>
- [52] K. Vaidyanathan, W. Huang, L. Chai, and D. K. Panda. 2007. Designing Efficient Asynchronous Memory Operations Using Hardware Copy Engine: A Case Study with I/OAT. In *2007 IEEE International Parallel and Distributed Processing Symposium*. 1–8. <https://doi.org/10.1109/IPDPS.2007.370479>
- [53] Karthikeyan Vaidyanathan and Dhableswar K Panda. 2007. Benefits of I/O acceleration technology (I/OAT) in clusters. In *2007 IEEE International Symposium on Performance Analysis of Systems & Software*. IEEE, 220–229.

- [54] 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 *2020 USENIX Annual Technical Conference (USENIX ATC 20)*. 559–571.
- [55] Kan Wu, Kaiwei Tu, Yuvraj Patel, Rathijit Sen, Kwanghyun Park, Andrea Arpaci-Dusseau, and Remzi Arpaci-Dusseau. 2022. NyxCache: Flexible and Efficient Multi-tenant Persistent Memory Caching. In *20th USENIX Conference on File and Storage Technologies (FAST 22)*. USENIX Association, Santa Clara, CA, 1–16. <https://www.usenix.org/conference/fast22/presentation/wu>
- [56] Lingfeng Xiang, Xingsheng Zhao, Jia Rao, Song Jiang, and Hong Jiang. 2022. Characterizing the Performance of Intel Optane Persistent Memory: A Close Look at Its on-DIMM Buffering. In *Proceedings of the Seventeenth European Conference on Computer Systems (Rennes, France) (EuroSys '22)*. Association for Computing Machinery, New York, NY, USA, 488–505. <https://doi.org/10.1145/3492321.3519556>
- [57] Jian Xu, Juno Kim, Amirsaman Memaripour, and Steven Swanson. 2019. Finding and Fixing Performance Pathologies in Persistent Memory Software Stacks. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems (Providence, RI, USA) (ASPLOS '19)*. New York, NY, USA, 427–439. <https://doi.org/10.1145/3297858.3304077>
- [58] Jian Xu and Steven Swanson. 2016. NOVA: A Log-structured File System for Hybrid Volatile/Non-volatile Main Memories. In *14th USENIX Conference on File and Storage Technologies (FAST 16)*. Santa Clara, CA, 323–338. <https://www.usenix.org/conference/fast16/technical-sessions/presentation/xu>
- [59] Zi Yan. 2019. Accelerate page migration and use memcg for PMEM management. <https://lwn.net/Articles/784925/>. [Online; accessed Jan-2023].
- [60] Jian Yang, Juno Kim, Morteza Hoseinzadeh, Joseph Izraelevitz, and Steve Swanson. 2020. An Empirical Guide to the Behavior and Use of Scalable Persistent Memory. In *18th USENIX Conference on File and Storage Technologies (FAST 20)*. Santa Clara, CA, 169–182. <https://www.usenix.org/conference/fast20/presentation/yang>
- [61] Jifei Yi, Benchao Dong, Mingkai Dong, Ruizhe Tong, and Haibo Chen. 2022. MT²: Memory Bandwidth Regulation on Hybrid NVM/DRAM Platforms. In *20th USENIX Conference on File and Storage Technologies (FAST 22)*. USENIX Association, Santa Clara, CA, 199–216. <https://www.usenix.org/conference/fast22/presentation/yi-mt2>
- [62] Diyu Zhou, Yuchen Qian, Vishal Gupta, Zhifei Yang, Changwoo Min, and Sanidhya Kashyap. 2022. ODINFS: Scaling PM Performance with Opportunistic Delegation. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. USENIX Association, Carlsbad, CA, 179–193. <https://www.usenix.org/conference/osdi22/presentation/zhou-diyu>

Received 16 October 2023; revised 14 January 2024; accepted 4 March 2024