

HiNFS: A Persistent Memory File System with Both Buffering and Direct-Access

YOUMIN CHEN, JIWU SHU, JIAXIN OU, and YOUYOU LU, Tsinghua University

Persistent memory provides data persistence at main memory with emerging non-volatile main memories (NVMMs). Recent persistent memory file systems aggressively use *direct access*, which directly copy data between user buffer and the storage layer, to avoid the double-copy overheads through the OS page cache. However, we observe they all suffer from slow writes due to NVMMs' asymmetric read-write performance and much slower performance than DRAM.

In this article, we propose HiNFS, a high-performance file system for non-volatile main memory, to combine both *buffering* and *direct access* for fine-grained file system operations. HiNFS uses an *NVMM-aware Write Buffer* to buffer the lazy-persistent file writes in DRAM, while performing direct access to NVMM for eager-persistent file writes. It directly reads file data from both DRAM and NVMM, by ensuring read consistency with a combination of the *DRAM Block Index* and *Cacheline Bitmap* to track the latest data between DRAM and NVMM. HiNFS also employs a *Buffer Benefit Model* to identify the eager-persistent file writes before issuing I/Os. Evaluations show that HiNFS significantly improves throughput by up to 184% and reduces execution time by up to 64% comparing with state-of-the-art persistent memory file systems PMFS and EXT4-DAX.

Categories and Subject Descriptors: D.4.2 [Operating Systems]: Storage Management–Persistent Memory

General Terms: Design, Performance

Additional Key Words and Phrases: Persistent memory, file system, direct access, buffering

ACM Reference format:

Youmin Chen, Jiwu Shu, Jiaxin Ou, and Youyou Lu. 2018. HiNFS: A Persistent Memory File System with Both Buffering and Direct-Access. *ACM Trans. Storage* 14, 1, Article 4 (March 2018), 30 pages. <https://doi.org/10.1145/3204454>

1 INTRODUCTION

Emerging fast, byte-addressable non-volatile memories (NVMs), such as phase change memory (PCM) (Doller 2009; Lee et al. 2010; Burr et al. 2010), resistive RAM (ReRAM), and memristor (Yang and Williams 2013), are promised to be employed to build fast, cheap, and persistent memory

An earlier version of this article appeared in *Proceedings of the 11th European Conference on Computer Systems (EuroSys 2016)* (Ou et al. 2016).

This work is supported by the National Natural Science Foundation of China (Grant No. 61772300, 61232003), the Beijing Municipal Science and Technology Commission of China (Grant No. D151100000815003), China Postdoctoral Science Foundation (Grant No. 2016T90094, 2015M580098). Youyou Lu is also supported by the Young Elite Scientists Sponsorship Program of China Association for Science and Technology (CAST).

Authors' addresses: Y. Chen, J. Shu, J. Ou, and Y. Lu (corresponding author), Department of Computer Science and Technology, Tsinghua University, Beijing, China; emails: {chenym16, ojx11}@mails.tsinghua.edu.cn, {shujw, luyouyou}@tsinghua.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 ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 ACM 1553-3077/2018/03-ART4 \$15.00

<https://doi.org/10.1145/3204454>

systems. Attaching NVMMs directly to processors produces non-volatile main memories (NVMMs), exposing the performance, flexibility, and persistence of these memories to applications (Zhang and Swanson 2015; Zhang et al. 2015). Moreover, these devices are expected to become a common component of the memory/storage hierarchy for laptops, PCs, and servers in the near future (Condit et al. 2009; Qureshi et al. 2009; Lee et al. 2009; Zhou et al. 2009; Chen et al. 2011; Jiang et al. 2012; Jung et al. 2013).

Given the anticipated high-performance characteristics of emerging NVMMs, recent research (Dulloor et al. 2014; Condit et al. 2009; Wu and Reddy 2011; DAX 2014) shows that the overheads from the generic block layer and copying data between the OS page cache and the NVMM storage significantly degrade the system performance. To avoid these overheads, state-of-the-art NVMM-aware file systems, such as BPFS (Condit et al. 2009), PMFS (Dulloor et al. 2014), EXT4-DAX (DAX 2014; EXT 2014), and so on, bypass the OS page cache and the generic block layer. Specifically, all of them directly copy data between the user buffer and the NVMM storage without going through the OS page cache, implying that all requests incur prompt access to NVMM.

Unfortunately, one major drawback of NVMM is the slow writes (Volos et al. 2011; Chen et al. 2011; Huang et al. 2014). The asymmetric read-write performance of NVMM indicates that, while DRAM and NVMM have similar read performance, the write operations of existing NVMM technologies, such as PCM and ReRAM, incur longer latency and lower bandwidth compared to DRAM (Suzuki and Swanson 2015; Zhang and Swanson 2015). Therefore, direct access to NVMM can lead to suboptimal system performance as it exposes the long write latency of NVMM to the critical path. Furthermore, our experiments of running existing NVMM-aware file systems on a simulated NVMM device show that the overhead from the direct write access can dominate the system performance degradation.

The relatively large write performance gap between DRAM and NVMM indicates that buffering writes in DRAM is important for improving the NVMM system performance, because (1) writes to the same block may be coalesced, since many I/O workloads have access locality (Min et al. 2012; Roselli et al. 2000; Ruemmler and Wilkes 1993; Ou et al. 2014), and (2) writes to files that are later deleted do not need to be performed. In addition, writes in file systems typically involve a trade-off between performance and persistence, and applications usually have alternative approaches to persisting their data (Nightingale et al. 2006; Harter et al. 2011).

However, simply using DRAM as a cache of NVMM is inefficient due to the double-copy overheads in the critical path among the user buffer, the DRAM cache, and the NVMM storage (Dulloor et al. 2014; DAX 2014). On one hand, reading data to a block not present in the DRAM cache causes the double-copy overhead in the read path, because the operating system needs to first copy the data from the storage layer to the DRAM cache, and then copy it from the DRAM cache to the user buffer. On the other hand, synchronous writes or synchronization operations, such as `fsync`, also lead to the double-copy overheads in the write path. For instance, if an application issues a write operation to block *A* followed by a `fsync` operation to persist block *A*, it incurs double data copies for block *A*. (The operating system first copies it to the DRAM cache at the write operation, and then copies it to the storage layer at the `fsync` operation.) The double-copy overheads can substantially impact the system performance when the storage device is attached directly to the memory bus and can be accessed at memory speeds (Dulloor et al. 2014; DAX 2014; Wu and Reddy 2011; Condit et al. 2009).

To address these problems, we propose HiNFS, a high-performance file system for non-volatile main memory. The **goal** of HiNFS is to hide the long write latency of NVMM whenever possible but without incurring extra overheads, such as the double-copy or software stack overheads, thereby improving the system performance. Specifically, HiNFS buffers the lazy-persistent file

writes (i.e., write operations that are allowed to be persisted lazily by file systems) in DRAM temporarily to hide the long write latency of NVMM. To improve the fetch/writeback performance of a buffer block, HiNFS manages the DRAM buffer at a fine-grained granularity by leveraging the byte-addressable property of NVMM. In addition, HiNFS interacts between the DRAM buffer and the NVMM storage using a memory interface, rather than going through the generic block layer, to avoid the high software stack overhead. To eliminate the double-copy overheads from the critical path, HiNFS performs direct access to NVMM for the eager-persistent file writes (i.e., write operations that are required to be persisted immediately), and directly reads file data from both DRAM and NVMM as they have similar read performance. However, writing data to DRAM and NVMM alternatively imposes a challenge for ensuring read consistency. Meanwhile, it also requires the file system to identify the eager-persistent writes before issuing the write operations.

This article makes four contributions:

- We reveal the problem of the direct access overheads by quantifying the copy overheads of state-of-the-art NVMM-aware file systems on a simulated NVMM device. Based on our experimental results, we find that the overhead from the direct write access dominates the system performance degradation in most cases.
- We propose an *NVMM-aware Write Buffer* policy to hide the long write latency of NVMM by buffering the lazy-persistent file writes in DRAM temporarily. To eliminate the double-copy overheads, we use direct access for file reads and eager-persistent file writes.
- We ensure read consistency by using a combination of the *DRAM Block Index* and *Cacheline Bitmap* to track the latest data between DRAM and NVMM. We also design a *Buffer Benefit Model* to identify the eager-persistent file writes before issuing the write operations.
- We implement HiNFS as a kernel module in Linux kernel 3.11.0 and evaluate it to demonstrate the benefits gains from fine-grained combination of buffering and direct access.

The remainder of this article is organized as follows. Section 2 discusses the problem in state-of-the-art NVMM-aware file systems and analyzes their direct access overheads. We present the design and implementation of HiNFS in Sections 3 and 4, respectively. We then present the evaluation results of HiNFS in Section 5. Finally, we discuss related work in Section 6 and conclude in Section 7.

2 BACKGROUND AND MOTIVATION

2.1 Problem in NVMM-aware File Systems

State-of-the-art NVMM-aware file systems, like BPFS (Condit et al. 2009), SCMFS (Wu and Reddy 2011), PMFS (Dulloor et al. 2014), and EXT4-DAX (EXT 2014), eliminate the OS page cache, which access the byte-addressable NVMM storage device directly. As an example, a `write()` syscall copies the written data from the user buffer to the NVMM device directly without going through the OS page cache and the generic block layer.

While this approach avoids the double-copy overheads, direct access to NVMM also exposes its long write latency to the critical path, leading to suboptimal system performance. In addition, to ensure data persistence and consistency, file systems either employ a cache bypass write interface¹

¹Different from the DRAM buffer cache, the CPU cache is hardware controlled, which is cumbersome for the file system to track the state of the written data. As a result, existing NVMM-aware file systems, such as PMFS, use a cache bypass interface (e.g., `copy_from_user_inatomic_nocache()`) to enforce that the written data becomes persistent before the associated file system metadata does, because they wouldn't be able to control the writeback from the processor caches to the NVMM storage without using an expensive `clflush` operation.

or use a combination of the `clflush` and `mfence` instructions behind write operations to explicitly flush data from the CPU caches to the NVMM device to enforce ordering (Dulloor et al. 2014; Wu and Reddy 2011), because existing cache hierarchies that were designed for volatile memory may reorder writes to improve the performance. For this reason, write latency is usually in the critical path, which cannot be tolerated by the CPU caches when NVMM is used as a persistent storage device rather than a volatile memory device (Pelley et al. 2014; Lu et al. 2014, 2015). Although BPFS’s epoch-based caching architecture offers an elegant solution, it requires complex hardware modifications, which involve non-trivial changes to cache and memory controllers (Condit et al. 2009). In our work, we would therefore like to design an NVMM system without any hardware modifications.

In this article, we mainly investigate how to design a high-performance file system for NVMM by hiding the long write latency of NVMM but without introducing extra overheads. Our work is based on several assumptions shown as follows.

- First, we assume that NVMM devices are attached directly to the memory bus alongside DRAM, and the operating system is able to distinguish the NVMM devices from the DRAM ones (Cooperation 2015).
- Second, we use the `clflush/mfence` instructions to enforce ordering and persistence, and assume that the `clflush` instruction guarantees that the flushing data actually reaches the persistent point (i.e., NVMM device). While Intel has proposed new instructions (CLWB/CLFLUSHOPT/PCOMMIT) to improve the cacheline flush performance and the CPU cache efficiency (Cooperation 2016), these approaches are still unavailable in existing hardware. This article, therefore, does not take them into consideration.
- Finally, HiNFS is mainly optimized for file-based I/O (i.e., *read* and *write* system calls) rather than memory-mapped I/O, as many important applications rely on traditional file I/O interfaces to access file data. However, HiNFS still supports direct access for memory-mapped I/O similar to existing NVMM-aware file systems (e.g., PMFS), which means that it does not sacrifice the performance of memory-mapped I/O. For the remainder of the article, we refer to file write simply as *write* and file read simply as *read*.

2.2 The Direct Access Overheads of NVMM-aware File Systems

In this section, we will show that the overhead from the direct write access in existing NVMM-aware file systems can dominate the system performance degradation, and hence it is essential to reduce such overhead whenever possible.

To quantify the direct access overheads of existing NVMM-aware file systems, we run the `fio` (FIO) microbenchmark on PMFS (Dulloor et al. 2014),² and use the `perf` profiling utility to obtain a breakdown of the time spent on running the benchmark. We use DRAM to emulate NVMM by introducing an extra configurable delay to NVMM writes to emulate NVMM’s slower writes relative to DRAM. More technical details about our experimental setup are given in Section 5.1.

Each test is run for 60s, and the results are shown in Figure 1. In all tests, we set the read/write ratio to 1:2 by default. In this figure, the time breakdown is organized into three categories: (1) *Read Access* refers to the overhead of copying data from the NVMM storage to the user buffer for read requests; (2) *Write Access* represents the overhead of copying data from the user buffer to the NVMM storage for write requests; and (3) *Others* is the overhead excluding the *Read Access* and

²We choose PMFS (Dulloor et al. 2014) as a case study of the baseline system, because it along with EXT4-DAX (EXT 2014) are the only available open-source NVMM-aware file systems at present. We also perform the same tests on EXT4-DAX, and it shows similar results. While BPFS (Condit et al. 2009) and SCMFS (Wu and Reddy 2011) are not open-source, we believe our observations also apply to them as they both perform direct access to NVMM.

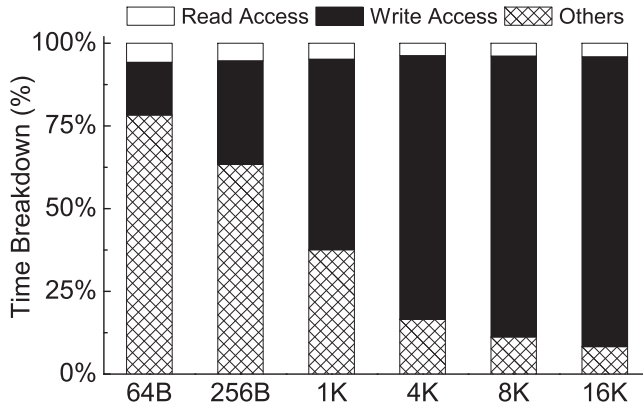


Fig. 1. Time breakdown of running the Fio Benchmark on PMFS.

Write Access overheads, which mainly includes overheads from user-kernel mode switch, file abstraction, and so on. From this figure, we observe that the direct write access is a major source of overhead in most cases. We also notice that the proportion of write access time increases as the I/O size becomes larger. One possible reason is that the data blocks are written back to non-volatile main memory in the granularity of cache lines, and larger blocks require higher overhead in book-keeping (i.e., the metadata to index which parts need to be persisted) and cache flush. While the overhead of read operations are less likely to be affected by the I/O sizes, because they don't involve cache flushing and metadata update. When the I/O size is no less than 4KB, the direct write access overhead can account for over 80% of the total overheads, which substantially degrades the system performance. When the I/O size becomes smaller, such as 64B, the direct write access overhead becomes relatively less significant than others, but still accounts for at least 16% of the total overheads.

While file systems can optimize the performance of the write operations that are not required to be persisted immediately, others, such as write operations enforced by synchronization operations, must enter the stable storage instantly to guarantee the data persistence required by user applications. Thus, their NVMM access overheads cannot be avoided. To see if there is enough room for optimizing those lazy-persistent writes, we perform another experiment that collects the `fsync` bytes across various workloads. Figure 2 shows the results of the percentage of `fsync` bytes with different workloads. More detailed descriptions of these workloads are given in Section 5. In this figure, we observe that different workloads have different persistence requirements. For example, TPC-C has over 90% `fsync` writes whereas LASR has no `fsync` writes. To conclude, a large number of applications have a significant portion of lazy-persistent writes, which are consistent with prior research results (Harter et al. 2011).

The above observations have interesting implications for the design of the file system for fast NVMM. On one hand, the revealed direct write access overhead strongly suggests that we need to reduce prompt writes to NVMM to improve the performance. On the other hand, we believe that an elegant design should be flexible. In other words, it should not improve the performance in some particular cases, while sacrificing the performance in other cases. For example, simply using DRAM as a cache of NVMM may improve the performance for workloads having many lazy-persistent writes, but this simple design will significantly degrade the system performance for workloads containing many eager-persistent writes due to the double-copy overheads.

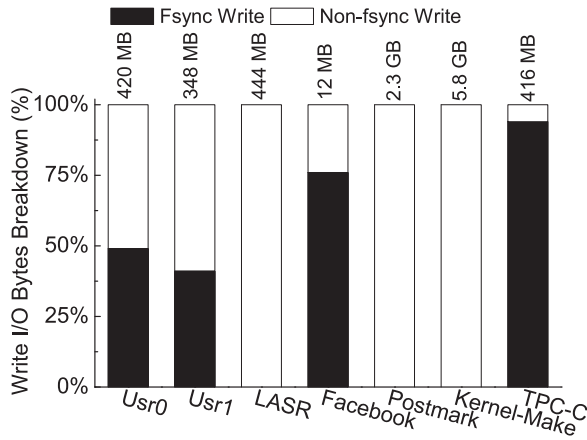


Fig. 2. Percentage of Fsync Bytes with different workloads. *The value atop each bar shows total bytes written.*

3 HINFS DESIGN

In this section, we first describe the high-level system architecture comparison of existing file systems and HiNFS. We then present an *NVMM-aware Write Buffer* policy to reduce prompt writes to NVMM by buffering the lazy-persistent writes in DRAM temporarily. Finally, we discuss how to eliminate the double-copy overheads resulted from conventional buffer management.

3.1 System Architecture

Figure 3(a) shows the system architecture of traditional block-based file systems on a RAMDISK-like NVMM block device. This is the most straightforward way to use NVMM as a persistent storage in which legacy file systems, such as ext2/ext4, can directly work on NVMM without extra modifications by emulating it as a block device. In a block-based file system, each file I/O usually requires two data copies, one between the block device and the OS page cache through the generic block layer, and one between the OS page cache and the user buffer through the memory interface. However, it has been recently reported that the overheads from the double-copy and the generic block layer can significantly impact the NVMM system performance (Condit et al. 2009; Dulloor et al. 2014; Wu and Reddy 2011; DAX 2014). As a result, state-of-the-art NVMM-aware file systems, such as BPFS (Condit et al. 2009), PMFS (Dulloor et al. 2014), and so on, access the NVMM device directly as shown in Figure 3(b). In these NVMM-aware file systems, each file I/O requires only a single data copy, directly between the NVMM and the user buffer (a.k.a., direct access). Unfortunately, the major drawback of this approach is that it does not consider NVMM’s relatively longer write latency compared to DRAM. Specifically, each write operation leads to prompt access to NVMM, which always expose the long write latency of NVMM to the critical path, leading to suboptimal system performance. Therefore, to get the best system performance, we propose another system architecture for the NVMM storage as shown in Figure 3(c). The design objectives of HiNFS are twofold:

- (1) *Buffering* to hide the long write latency of NVMM behind the critical path. HiNFS uses an NVMM-aware *Write Buffer* policy to buffer the lazy-persistent writes in DRAM temporarily. HiNFS design, including fine-grained buffer management and using a memory interface to interact between DRAM and NVMM, is optimized for the NVMM storage (Section 3.3).

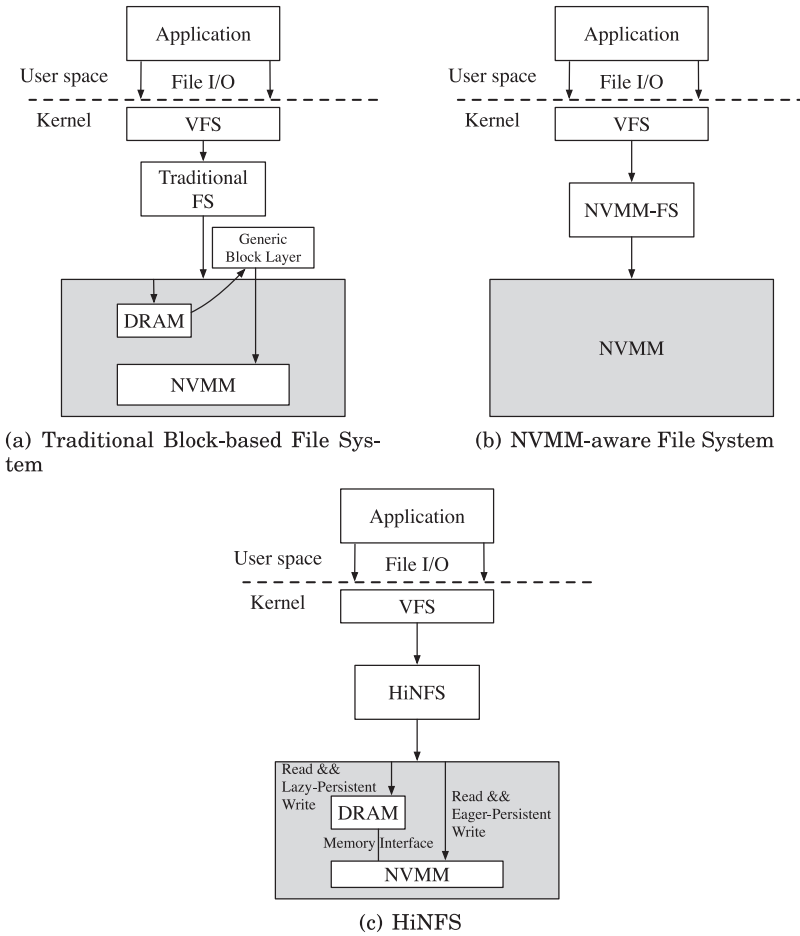


Fig. 3. Architecture comparison of different file systems for NVMM.

- (2) *Direct Access* to eliminate the double-copy overheads. Although buffering can help hide the long write latency of NVMM, it may introduce the double-copy overheads. For this reason, HiNFS optimizes read and eager-persistent write by avoiding unnecessary data copies. Read or eager-persistent write, in HiNFS, requires only a single data copy between DRAM/NVMM and the user buffer (Section 3.4).

Moreover, the combination of *Buffering* and *Direct Access* is also beneficial to extending the lifetime of NVMM devices. This is because multiple asynchronous updates to the same physical block are more likely to be merged to a single write, hence, the write traffic to the NVMM devices can be dramatically reduced.

3.2 Buffer Benefit Model

In a file system with both buffering and direct access, one challenge is how to determine whether a I/O request is eager-persistent or lazy-persistent during I/O issues. The eager-persistent writes include both synchronous writes that are explicitly tagged with sync (e.g., writes in a file opened with `O_SYNC` flag) and asynchronous writes followed by sync operations (e.g., `fsync`, `fdatasync`).

To predict whether a write request is eager-persistent HiNFS designs a *Buffer Benefit Model* for prediction when a request arrives.

The Prediction Model. The Buffer Benefit Model performs prediction using history information of the recent synchronization information. It is designed based on our observations from various workloads that synchronization operations remain nearly the same within a short time period in most cases. For example, one file marked as Eager-Persistent is either opened with O_SYNC flag, or is always updated with a following sync operations in later time, and we call this behavior as good locality.

With this policy, the model associates each data block with one bit, namely Eager-Persistent, to indicate the synchronization state. In HiNFS, each 4KB data block needs one extra bit to indicate its current state, implying that this overhead is small and can be acceptable. Moreover, we store the block states in DRAM rather than in slow NVMM. If a data block is decided to be in the Eager-Persistent state, then all the subsequent asynchronous writes to this data block are considered as the eager-persistent writes. Otherwise, they are considered as the lazy-persistent writes, which are issued to the DRAM buffer first.

In the *Buffer Benefit Model*, the DRAM write latency is denoted as L_{dram} , and the NVMM write latency is expressed as L_{nvmm} . N_{cw} indicates the total number of cacheline writes between the previous and current synchronization operation of a data block, while N_{cf} is the total number of cacheline flushes from DRAM to NVMM of a data block, which are performed by the current synchronization process rather than the background writeback threads. Then, buffering is more efficient than non-buffering for this block only if it satisfies the following inequality:

$$N_{cw} * L_{dram} + N_{cf} * L_{nvmm} < N_{cw} * L_{nvmm}. \quad (1)$$

This inequality means that the total execution time if writing to DRAM first is less than that if writing to NVMM directly for a data block. If a block satisfies this inequality, then it will be set to the Lazy-Persistent state. Otherwise, it would be set to the Eager-Persistent state.

When the file system is mounted, all the existing or newly created data blocks are initialized to the Lazy-Persistent state before the arrival of their first synchronization operations. After that, we dynamically decide the data block states at each file operation. At each synchronization operation,³ we calculate to see if the related data blocks, which are required to be persisted to NVMM in the current synchronization operation, satisfy the above inequality. If a data block cannot satisfy this inequality, then the state of this block is set to Eager-Persistent, which means that any subsequent asynchronous writes to this data block go directly to NVMM. Otherwise, we set the block state to Lazy-Persistent. Moreover, the state of a data block is switched from Eager-Persistent to Lazy-Persistent if it has not met a synchronization operation for a certain period of time, which is set to 5s by default and can be adjusted. It is worth noting that we achieve this by deciding the data block state at the time of writing this block using the last synchronization time of its dependent file,⁴ rather than scanning all the data blocks at each fixed time, as it is lightweight to record the file synchronization time.

To get the value of N_{cf} of a buffer block, we maintain a *ghost buffer* to measure the total number of cacheline flushes from DRAM to NVMM of a buffer block during each synchronization operation. Ghost buffer assumes that every write goes to the DRAM buffer first but maintains only the buffer index metadata rather than the actual data. This leads to low memory overhead, which requires less than 1% of the total DRAM buffer space.

³In the current implementation, HiNFS only regards the fsync system call as the synchronization operation. While the msync operation is also a synchronization point in HiNFS, it is related to mmap I/O rather than file I/O.

⁴As the synchronization operation, such as fsync, is based on the file granularity, HiNFS adds a new field to the file metadata structure to record the last synchronization time of its related data blocks.

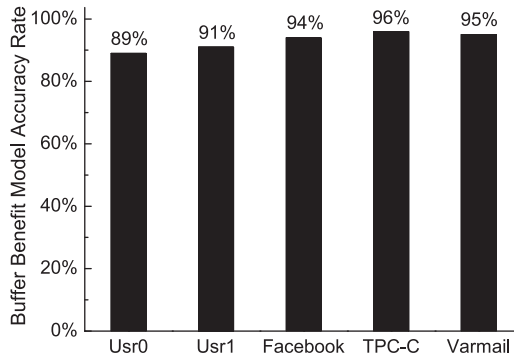


Fig. 4. The accuracy rate of the Buffer Benefit Model using the most recent synchronization information for different workloads.

Prediction Accuracy. To see whether using the most recent synchronization information of a block to predict the state of its next synchronization operation is accurate, we measure the accuracy rate of our model using various workloads. The results are shown in Figure 4. We select five workloads that contain the synchronization operations and the descriptions of these workloads are shown in Section 5. Moreover, we measure it during the synchronization operations for each block. That is, if both the current and previous synchronization operation for a block satisfy or violate Inequality Equation (1), it is accurate; Otherwise, it is inaccurate. In this figure, we can see that the accuracy ratio is close to 90% even in the worst case (i.e., Ustr0). These results demonstrate that the synchronization information of a block remains nearly the same within a short time period, and thus our *Buffer Benefit Model* is effective in most cases.

Note that the *Buffer Benefit Model* is only a heuristic prediction model, but it is enough to provide satisfiable prediction accuracy. Many unimportant factors are not taken into consideration, such as the overhead of CPU cache misses when data in DRAM buffer is written back to NVMM. This is because (1) we assume that the NVMM access latency is much higher than that of DRAM, thus the read back latency is insignificant compared to other parts, and (2) the asynchronous I/Os that followed by synchronization operations are more likely to be reside in the CPU cache, so the cache miss ratio is relatively low. According to the prediction accuracy in Figure 4, we conclude that our *Buffer Benefit Model* works both concisely and effectively.

3.3 NVMM-aware Write Buffer Policy

To buffer the lazy-persistent writes, we propose an *NVMM-aware Write Buffer* policy to store them in DRAM temporarily, to hide the relatively long write latency of NVMM behind the critical path. Figure 5 shows an overview of HiNFS. When a write request is serviced, the *Eager-Persistent Write Checker* module would decide whether the current write operation is a lazy-persistent or eager-persistent write, based on the Buffer Benefit Model discussed in the previous section. If it is a lazy-persistent write, then HiNFS would like to issue this write request to the fast DRAM buffer, thereby eliminating the overhead of writing the NVMM. In the buffering, HiNFS only buffers file data blocks. Metadata blocks are not buffered currently to ease consistency implementations.

3.3.1 DRAM Block Index. HiNFS builds per-file B-tree in DRAM to index the DRAM blocks. The per-file indexing structure is feasible, because eager-persistent writes have good locality according to discussions in Section 3.2. Figure 6 shows the details of the *DRAM Block Index* in HiNFS. In the *DRAM Block Index*, the key of the index is the logic file offset, which is aligned to the DRAM

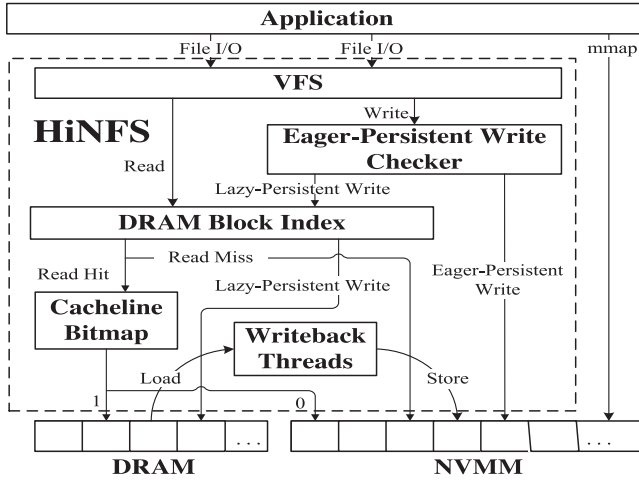


Fig. 5. HiNFS overview.

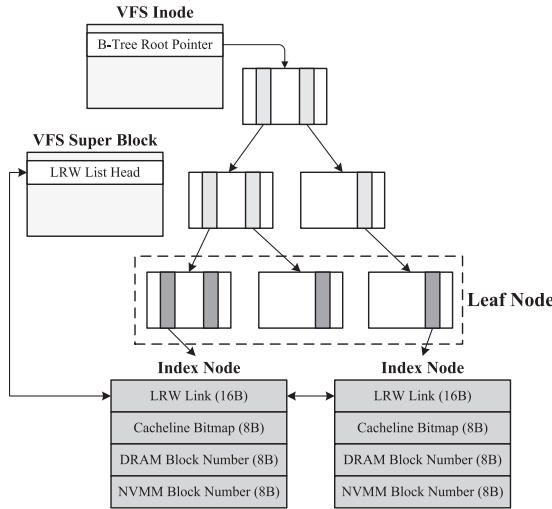


Fig. 6. DRAM Block Index.

block size, and the value field is the Index Node shown in Figure 6. The Index Node is a 40byte value, which contains a 16byte LRW Link pointers, a 8byte Cacheline Bitmap (for the *Cacheline Level Fetch/Writeback (CLFW)* in the following section), a 8byte DRAM block number, and its corresponding 8byte NVMM block number. The NVMM block number in the value field enables the background writeback threads to flush the DRAM block to the corresponding NVMM block address. The root pointer of the B-tree is stored in the kernel’s VFS inode structure. Moreover, all the index nodes are allocated from DRAM and linked to a global LRW (Least Recently Written) list, the head of which is located in the kernel’s VFS super block structure.

Note that the *DRAM Block Index* structure is located in DRAM entirely rather than in NVMM, to enable fast index operations and effectively track the status of each file data block. However, the memory space consumed by the *DRAM Block Index* is insignificant, because (1) Each Index

Node in the leaf nodes has the size of 40bytes, which is less than 1% of a 4KB data block; (2) Only Lazy-Persistent files have their DRAM Block Index kept in DRAM, thus further reduces the memory consumption.

We use the B-tree structure for the *DRAM Block Index*, because we would like to reuse the B-tree data structure from the PMFS (Dulloor et al. 2014) implementation as HiNFS is implemented based on it. While other index structures, such as hash table, can also be employed by HiNFS, the difference between B-tree and them may be only several bytes of DRAM access for each 4KB block access, the overhead of which is far less than that of the data copy operations. Therefore, we believe that the data structure selection for the *DRAM Block Index* is not a critical issue, and thus there will be little performance difference between the index implementations of B-tree and other structures for HiNFS.

3.3.2 Fine-Grained Buffer Block Fetch and Writeback. Conventional buffer management in the OS page cache maintains the DRAM buffer space at the block granularity (i.e., 4KB). This coarse-grained buffer management is inefficient for HiNFS. On one hand, an unaligned lazy-persistent write to a block not present in the DRAM buffer causes the operating system to synchronously fetch the block from the NVMM storage into the DRAM buffer before the write is applied. Such *fetch-before-write* requirement impacts the system performance, because the fetching process can block the writing process (Campello et al. 2015). On the other hand, a whole buffer block would be flushed to storage even though only a few bytes of data are written to this block, causing a significant impact on the foreground application performance for two main reasons. First, when the DRAM buffer has no free blocks, the foreground lazy-persistent writes may stall until the background writeback threads reclaim enough free DRAM buffer space. Second, the background writeback threads can also compete the limited NVMM write bandwidth with the foreground eager-persistent writes. As a result, it is essential to improve the fetch/writeback performance of a buffer block to achieve higher system performance.

To address the above issue, we propose *Cacheline Level Fetch/Writeback (CLFW)*, which tracks the writes to the DRAM blocks on the basis of processor's cache lines. In *CLFW*, data is fetched from or flushed to NVMM in a fine-grained way rather than the block level. To do so, we use a *Cacheline Bitmap* (as shown in Figure 5) to track the state of each cacheline within a DRAM block. In this scheme, when a dirty DRAM block is selected for eviction, the writeback thread will check the *Cacheline Bitmap* of this block. Only if the P bit is 1 (i.e., the Pth cacheline is dirty), the cacheline should be written back to the NVMM. For an unaligned lazy-persistent write to a block not present in the DRAM buffer, we only need to fetch the corresponding cachelines instead of the whole block into the DRAM buffer. For example, for the baseline system with 4KB DRAM block size and 64B cacheline size, if a user writes to the 0~112B region of a block, traditional system needs to fetch the whole block (0~4096B) into the DRAM buffer, while *CLFW* only needs to fetch the second cacheline of this block (64~128B) into the DRAM buffer. In summary, *CLFW* significantly reduces the wasteful data-fetch and data-flush for workloads containing many small block-unaligned lazy-persistent writes, thereby improving the performance in these cases.

3.4 Elimination of the Double-Copy Overheads

As fast NVMM is attached directly to the processor's memory bus and can be accessed at memory speeds, extra data copies would be inefficient for NVMM systems, which can substantially degrade their performance (Dulloor et al. 2014; Wu and Reddy 2011; Condit et al. 2009; DAX 2014). As a result, it is essential to avoid such overheads whenever possible. To this end, we find two key reasons to cause the double-copy overheads resulted from conventional buffer management. This section describes them and discusses how we overcome them separately. It is worth noting that

all the double-copy overheads, we pay attention to in this article, mainly refer to those that occur in the critical I/O path, as they are the key factors of affecting the system performance.

3.4.1 Direct Read. In conventional buffer management, reading data to a block not present in the DRAM buffer causes the operating system to fetch the block into the DRAM buffer first and then copy the data from the DRAM buffer to the user buffer, thereby leading to the double-copy overhead in the read path. To address this issue, HiNFS directly read data from both DRAM and NVMM to the user buffer, as they have similar read performance. Such direct copy policy is more efficient than conventional two-step copy policy as it eliminates unnecessary data copies.

The **read consistency** needs to be ensured for read operations to find the latest data blocks, when writing data to either DRAM or NVMM. To find the up-to-date data for a read operation, HiNFS first checks the *DRAM Block Index* to see if the corresponding block is in DRAM. If not, then it uses the file system block index to get the corresponding NVMM block address, and then performs this read operation to NVMM directly. Otherwise, it further checks the *Cacheline Bitmap* of the corresponding DRAM block to see which parts of data are in the DRAM block and which parts of data are in the NVMM block, and then copies the corresponding parts of data to the user buffer from both the DRAM and NVMM blocks on the basis of the *Cacheline Bitmap*. To minimize the number of memory copy (i.e., `memcpy`) operations, a single `memcpy` operation is used to copy the data in the consecutive cachelines, the corresponding bits of which in the *Cacheline Bitmap* have the same value, to the user buffer.

3.4.2 Direct Eager-Persistent Write. To further avoid the double-copy overhead in the write path, we issue the eager-persistent writes to NVMM directly rather than copying them to DRAM first. This is because writing them to DRAM not only causes unnecessary copy overheads, but also pollutes the buffer space, which may evict other valuable buffer blocks. In HiNFS, the *eager-persistent writes* are defined as the following two cases:

- (1) *Synchronous writes.* This happens when the file system is mounted with the sync option or the written file is opened with the `O_SYNC` flag.
- (2) *Asynchronous writes followed by explicit synchronization operations.* We divide this scenario into two cases. If enough asynchronous writes can be coalesced before the arrival of the next explicit synchronization operation, in which case buffering is more efficient than direct access, then we still regard them as the lazy-persistent writes. Otherwise, they are considered as the eager-persistent writes.

As HiNFS needs to choose either direct or buffer write mode for a write request, it is important to identify the eager-persistent writes before issuing the write operations. It is straightforward to identify case (1), because we can check the file system state by reading the file system super block and the file opening state by reading the file inode. To identify case (2), HiNFS uses the Buffer Benefit Model to predict the eager-persistent writes in advance.

The **write consistency** needs to be ensured between DRAM and NVMM. In other words, updates to a same block in either DRAM or NVMM need to be ordered to NVMM. To ensure the consistency, when a write operation is identified as the eager-persistent write, if it is in case (1), then we further check if the written block is present in the DRAM buffer before directly accessing the NVMM. If so, then we still write the data to the corresponding DRAM block, and explicitly evict it from the DRAM buffer before returning to users. Fortunately, this case rarely happens, unless the file opening or file system state is altered frequently. If it is in case (2), then we can always perform direct access to NVMM as long as the written block is in the Eager-Persistent

state, because the latest data of this block is guaranteed to be persisted to NVMM, since the last synchronization operation of this block.

4 IMPLEMENTATION

HiNFS is implemented based on the PMFS (Dulloor et al. 2014) file system in Linux kernel 3.11.0. HiNFS shares the file system data structures of PMFS but adds a new DRAM buffer layer and modifies the file I/O execution paths. In this section, we mainly discuss some details related to the implementation.

4.1 Buffering and Persistence

In HiNFS, allocation and replacement for the DRAM buffer are block-oriented. By default, the DRAM block size is 4KB, which equals to the default block size of the NVMM storage. Currently, we use the LRW (Least Recently Written) policy, a variant of the LRU (Least Recently Used) algorithm, for the replacement of the DRAM buffer blocks due to the simplicity and efficiency of the LRU policy over decades (Coffman and Denning 1973; Denning 1968). Specifically, we maintain the LRW list to keep track of the recency of write references of blocks in the DRAM buffer. That is, all the DRAM blocks are sorted by their last written time. When a DRAM block is written, it would be moved to the MRW (Most Recently Written) position. It is worth noting that this does not limit HiNFS of using other sophisticated buffer replacement policies, such as Least Frequently Used (LFU) (Willick et al. 1993), Adaptive Replacement Cache (ARC) (Megiddo and Modha 2003), 2Q (Johnson and Shasha 1994), and so on. Different buffer replacement policies have different buffer write hit ratios, which decide how many writes can be coalesced before a buffer block is written back to the NVMM. However, these policies also increase the complexity of the buffer design, and the adding software overhead is non-trivial for the NVMM system. For this reason, we believe that the LRW-based policy is a good candidate to help us improve the performance, as a large majority of file system workloads show strong locality and high I/O skewness (Min et al. 2012; Roselli et al. 2000; Ruemmler and Wilkes 1993; Ou et al. 2014). We leave the research of using different buffer replacement policies in the future.

To ensure data persistence, HiNFS creates multiple independent kernel threads at mount time to flush the dirty DRAM blocks to the NVMM periodically in background. The flushed DRAM blocks can be released to secure free DRAM blocks for further buffering. There are two different cases of waking up the background writeback threads:

- (1) The first case occurs when there are less than Low_f free DRAM blocks, where Low_f is a pre-defined threshold. In HiNFS, Low_f is set to 5% of the total DRAM blocks by default and is configurable.
- (2) The second case is that the background thread wakes up every 5s and periodically writes the updated data from the DRAM buffer to the NVMM storage.

When a writeback thread is woken up, it first selects the victim DRAM blocks from the LRW position of the LRW list. These victim DRAM blocks are then written back to the corresponding NVMM block addresses via a memory interface (e.g., `memcpy()`), rather than going through the generic block layer. After that, these DRAM blocks can be reclaimed for future write operations. The writeback thread reclaims several DRAM blocks at a time until the number of free DRAM blocks surpasses the $High_f$ threshold, which is set to 20% of the total DRAM blocks by default and can be adjusted. Then, the background writeback thread continues to scan the rest LRW list to write back any dirty DRAM blocks that were updated more than 30s ago. In addition, HiNFS flushes all the DRAM blocks to the NVMM when unmounting the file system.

4.2 System Consistency

To maintain file system consistency, traditional journaling file systems provide multiple levels of consistency using different journaling modes (e.g., *writeback*, *ordered data*, or *journal data* mode). However, the current implementation of HiNFS only provides *ordered data* mode, which means that it only guarantees the data updates become persistent before the related metadata updates. To achieve this, HiNFS reuses the PMFS's journaling mechanism, which only journals the file system metadata at the cacheline granularity (Dulloor et al. 2014). Note that HiNFS does not buffer any file system metadata (e.g., inode or directory entry).

Different from the journaling mechanism in PMFS, HiNFS needs to keep the persistence ordering of the lazy-persistent writes. To do this, each lazy-persistent write operation will create a new transaction. The file system data blocks in the lazy-persistent write operation are buffered to DRAM first without being journaled to NVMM. These data blocks in DRAM are tracked using a transaction handler. In contrast, the file system metadata and its undo log entries are written to NVMM directly using the PMFS's logging scheme. To guarantee the ordered mode journaling invariant, HiNFS does not write the commit log entry to the NVMM log space until the related DRAM data blocks are persisted to NVMM. Additionally, HiNFS ensures ordering and persistence using the `clflush` and `mfence` instructions. Each writeback operation of a data block is followed by the `clflush`/`mfence` instructions so that the subsequent commit log entry will not be persisted to NVMM before this data block.

To be able to identify the partially written log entries during recovery, HiNFS includes a valid flag in each cacheline size log entry and leverages the architectural guarantee in the processor caching hierarchy that writes to the same cacheline are never reordered, to indicate the integrity of a log entry, the approach of which is similar to that of PMFS (Dulloor et al. 2014). To achieve this, the valid flag is written last when writing a log entry so that it will not become persistent before the data of this log entry.

4.3 Direct Memory-mapped I/O (mmap) Support

One of the key features of state-of-the-art NVMM-aware file systems (e.g., PMFS) is that they can support direct memory-mapped I/O, thus removing unnecessary data copies. HiNFS also supports this feature. When *mmap* a file, HiNFS first flushes all the dirty DRAM blocks of this file to NVMM and then sets the states of all its related data blocks to Eager-Persistent, which remain unchanged until this file is *munmapped*. Then, it directly maps the file data into the application's virtual address space so that users can access NVMM directly. However, the *mmap* write operations are not guaranteed to be persistent until the arrival of the next *msync* operation, as they are performed to the CPU caches first before being persisted to the NVMM storage.

5 EVALUATION

In this section, we evaluate HiNFS to address the following questions:

- (1) How does HiNFS perform against existing file systems?
- (2) What are the benefits of eliminating the double-copy overheads?
- (3) How is the scalability of HiNFS compared to other file systems?
- (4) How is HiNFS sensitive to the variation of the I/O size of the workload, the DRAM buffer size, and the NVMM write latency?

We use the Filebench microbenchmark (Fil) to address questions (1), (2), (3), and (4). We use a variety of data-intensive traces and macrobenchmarks to further analyze questions (1) and (2). Table 1 provides a description of all the workloads we evaluate.

Table 1. Workloads and Descriptions

Type	Workload	Description
Micro	Fileserver	Emulates a simple file server, which consists of creates, deletes, appends, reads, and writes.
	Webserver	Emulates a web server, which performs file reads and log appends.
	Webproxy	Emulates a simple web proxy server with a mix of create-write-close, open-read-close, and delete operations, as well as log appends.
	Varmail	Emulates a mail server comprised of create-append-sync, read-append-sync, read, and delete operations.
Macro	Postmark (Katcher 1997)	Measures the performance of a file system used for e-mail and web-based services.
	TPC-C	Emulates the activity of a wholesale supplier where a population of users execute transactions against a database, we execute DBT2 workload (DBT) on PostgreSQL 8.4.10 database system with three warehouses.
	Kernel-Grep	Searching for an absent pattern under the Linux 3.11.0 kernel source directory.
	Kernel-Make	Running make inside the Linux 3.11.0 kernel source tree.
Traces	Usr0	System call trace collected from research desktop by FIU (Use).
	Usr1	System call trace collected from research desktop by FIU (Use) at different time from Usr0.
	LASR (LAS)	System call trace collected from computers used for software development by CS researchers.
	Facebook	MobiBench (LABORATORY 2013) facebook system call trace.

5.1 Experimental Setup

NVMM Emulator. As real NVMM devices are not available for us yet, we develop a simple performance emulator based on the NVMM emulator used in the Mnemosyne (Volos et al. 2011) project to evaluate HiNFS’s performance. Similar to prior projects (Volos et al. 2011; Huang et al. 2014; Volos et al. 2014), our NVMM emulator introduces an extra latency for each NVMM store operation to emulate the slower writes of NVMM relative to DRAM, while introducing no extra latency on the NVMM load operations. We have two considerations in assuming that NVMM and DRAM have the same read latency. First, we focus on the asymmetry of the read and write operations of NVMMs in HiNFS, and our evaluations focus on showing the benefits of the write performance rather than the read performance of HiNFS compared to state-of-the-art NVMM-aware file systems. Second, emulating the NVMM read latency is complicated due to CPU features such as speculative execution, memory parallelism, prefetching, and so on, which is hard to make it accurate (Dulloor et al. 2014).

NVMM Latency Emulation: Our emulator emulates NVMM using DRAM. To account for NVMM’s slower writes relative to DRAM, we introduce an extra configurable delay when writing to NVMM. We create delays using a software spin loop that uses the x86 RDTSCP instruction to read the processor timestamp counter and spins until the counter reaches the intended delay. Moreover, we add these delays after executing the *clflush* instruction. By default, we set the NVMM write latency to 200ns (Volos et al. 2011).

Table 2. Server Configurations for Quartz Measurement

CPU	Intel Xeon E5-2680 v3, 2.5GHz, x2
CPU cores	48
DRAM	384GB
Operating system	Ubuntu 12.04, linux 4.4.16, linux 3.11.0

NVMM Bandwidth Emulation: NVMM has significantly lower write bandwidth than DRAM (Zhang and Swanson 2015; Suzuki and Swanson 2015). Assume that B_{NVMM} indicates NVMM's write bandwidth and L_{NVMM} is NVMM's write latency. Then, we emulate the NVMM write bandwidth by limiting the maximum number of the concurrent NVMM writing threads (denoted as N_w), where N_w equals to $(B_{NVMM}/(1/L_{NVMM}))$. An NVMM writing thread would be queued if the number of the current NVMM writing threads reaches N_w , and the waiting queue will be woken up when one of the current NVMM writing threads completes. By default, the maximum sustained write bandwidth of NVMM is set to 1GB/s, about 1/8 of the available DRAM bandwidth on the unmodified system (Dulloor et al. 2014).

NVMMBD Emulator. To compare HiNFS against traditional block-based file systems, we construct another emulator, *NVMMBD*, to emulate the NVMM-based block device. We modify Linux's RAM disk module (brd device driver) and use the above NVMM performance model to emulate the NVMM latency and bandwidth.

Quartz Emulator. In latency emulation, both NVMM Emulator and NVMMBD Emulator works with cache line granularity, and thereby introduce extra latency with larger I/O size. So we also deploy Quartz (Volos et al. 2015) to run these benchmarks as comparison. Quartz is a NVM emulator that is able to emulate a wide range of NVM latencies and bandwidth characteristics, and can also be used to emulate the application execution on heterogeneous systems with both fast, regular volatile DRAM and slower persistent memory. Quartz achieves this by calculating the hardware performance monitoring counters (PMC) and add extra latencies with external signal interrupt. Quartz partitions the application execution time into continuous epoches, and dynamically injects software created delays at boundaries of each epoch. Quartz can also support multi-thread applications by capturing inter-thread events and adjusting the length of each epoch to keep accuracy. HiNFS is designed to run on systems with both fast DRAM and persistent memory, so we choose to deploy HiNFS on NUMA-enabled server, by reserving both local memory and remote memory with `mmap` command when the system is powered on. Then, HiNFS is initialized to use local memory as buffer space and remote memory as data storage. When executing applications, Quartz can automatically differentiate the local memory and remote memory and only add extra latency to any access to remote memory.

We evaluate the performance of HiNFS against five existing file systems listed in Table 3. NOVA, PMFS, and EXT4+DAX are the three available open-source NVMM-aware file systems that access NVMM directly. EXT2/EXT4+NVMMBD are traditional block-based file systems, which are built on the NVMMBD block device emulator. Both of them are mounted with default settings. All the experiments are conducted on a x86 server with NVMM and NVMMBD emulators. The configurations of the server are listed in Table 4. For all the experiments, each data-point is calculated using the average of at least five executions. Besides, we measure the performance of these file systems with Quartz for comparison and the configurations of the server are listed in Table 2. Quartz emulator is only used to run filebench, as Quartz will disable some system calls and many benchmarks will not work on top of it.

Table 3. Existing File Systems for Comparison

NOVA (Xu and Swanson 2016)	a log-structured file system with direct access to NVMM
PMFS (Dulloor et al. 2014)	a NVMM-aware file system with direct access to NVMM
EXT4+DAX (EXT 2014)	DAX is a kernel patch that supports EXT4 for bypassing the OS page cache
EXT2+NVMMBD	a traditional file system without journaling
EXT4+NVMMBD	a traditional journaling file system

Table 4. Server Configurations

CPU	Intel Xeon E5-2620, 2.1GHz
CPU cores	12
Processor cache	384KB 8-way L1, 1.5MB 8-way L2, 15MB 20-way L3
DRAM	16GB
NVMM	Emulated with slowdown, the write latency is 200ns, the write bandwidth is 1GB/s
Operating system	RHEL 6.3, kernel version 3.11.0

5.2 Microbenchmarks

In this section, we run four types of workloads from the Filebench benchmark. Each workload is run for 60s using 5GB pre-allocated files after clearing the contents of the OS page cache. Unless otherwise specified, all the experiments are run with multiple threads and the mean I/O size is set to 1MB by default.⁵ Moreover, HiNFS is mounted with 2GB DRAM buffer size, while EXT2/EXT4+NVMMBD are run with the available memory size being set to 8GB (5GB for storing the dataset on the NVMMBD and 3GB for the system memory). We use the number of operations per second, which is reported by the Filebench benchmark, as the performance metric.

We first evaluate the overall performance. Figure 7 shows the throughput normalized to that of PMFS. As shown in the figure, HiNFS achieves (one of) the best performance among the five file systems for all the evaluated workloads.

5.2.1 Overall Performance.

Fileserver. Comparing HiNFS with NOVA, PMFS, and EXT4+DAX, HiNFS gains performance improvement by up to 127%, this is because almost all the writes in the Fileserver workload are lazy-persistent, and HiNFS asynchronously persists them to NVMM, thereby hiding the long write latency of NVMM behind the critical path. EXT2 and EXT4 significantly underperform PMFS, NOVA and EXT4+DAX in this workload, as the benefits of the DRAM buffer are offset by the overheads from the double-copy and the generic block layer.

Webproxy. This workload has strong access locality, and EXT2+NVMMBD and EXT4+NVMMBD use the OS page cache to buffer the writes, we can see that only in this case can they outperform PMFS, NOVA, and EXT4+DAX.

⁵We choose 1MB as the mean I/O size for two reasons. First, this is the default configuration of the Filebench benchmark. Second, we adopt this configuration from the Aerie article (Volos et al. 2014). Sensitivity to different I/O sizes is also evaluated in Section 5.2.4 and Figure 10.

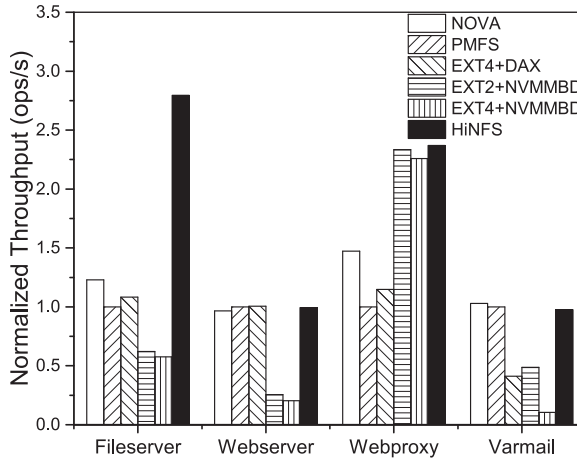


Fig. 7. Overall Performance.

Webserver. Webserver is a read-intensive workload, EXT2 and EXT4 with NVMDB show $3\times$ lower performance than PMFS due to the unnecessary read copies between the DRAM buffer and the NVMM storage. Comparatively, we can see that HiNFS, NOVA and PMFS achieve almost the same performance for the Webserver workload, demonstrating the benefits of eliminating the double-copy overheads.

Varmail. For the Varmail workload, we find that it contains a large part of synchronization operations. Moreover, all the writes in this workload are append operations, which cannot be coalesced in the DRAM buffer before the arrival of a synchronization operation. Therefore, we can see that HiNFS performs at par with PMFS and NOVA due to that HiNFS bypasses the buffer for these eager-persistent writes. The eager-persistent writes also cause the double-copy overheads, which account for the bad performance of EXT2 and EXT4. However, EXT4+DAX shows much lower performance than PMFS and NOVA in this case. This is because the Varmail workload contains many metadata operations, and EXT4+DAX still follows the cache-oriented methods for them, while PMFS follows direct access for both data and metadata. Different from PMFS, NOVA optimize metadata operations by placing index tree in DRAM and directly accessing metadata in NVMM, as a consequence, NOVA has the best performance among the six file systems in Varmail workload.

NOVA has slightly higher throughput than PMFS but still underperform HiNFS in most cases, because NOVA access NVMM directly like PMFS and EXT4+DAX, the result in our experiments shows closer throughput between NOVA and PMFS compared to the data in NOVA article, as we set the I/O size to 1MB and worker threads are varied between 1 to 10, however, in NOVA's experiments, write size is set to 16KB, and worker threads are set to 50, which lead to the performance variation.

5.2.2 System Scalability. We also evaluate the system scalability of HiNFS and other file systems. Figure 8 shows the throughput for the four filebench workloads as we vary the number of threads in a single client process. Surprisingly, HiNFS achieves the best scalability for all the evaluated workloads.

Fileserver. For the Fileserver workload, the performance of PMFS, NOVA and EXT4+DAX are gradually limited by the NVMM write bandwidth when going from 1 to 10 threads, while the performance of EXT2/EXT4+NVMDB is constrained by the overheads from the double-copy and

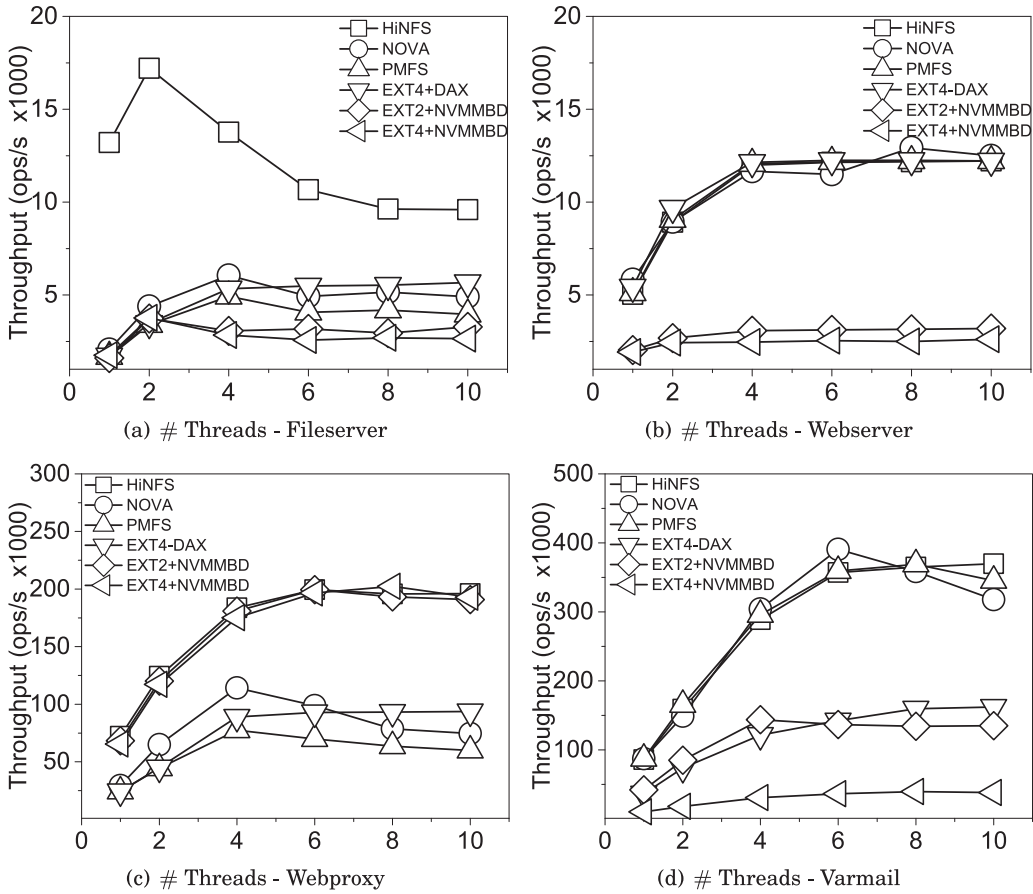


Fig. 8. Throughput (operations per second) for 1–10 threads (NVMM emulator).

the generic block layer. Therefore, HiNFS scales better than the other five file systems as it buffers and coalesces the writes before writing to NVMM. However, we find that HiNFS’s throughput drops when the thread count goes from 2 to 8, this is because the buffer write hit ratio decreases as the number of threads increases. Fortunately, the performance becomes stable beyond 8 threads, and HiNFS still achieves nearly $1.5\times$ higher performance than PMFS when going to 10 threads.

Webproxy. In fact, the performance of HiNFS basically depends on the write locality of the workloads. With better write locality, such as Webproxy, we can see that HiNFS always scales well and its performance never decreases as the thread count increases.

Webserver. For read-intensive workloads and workloads containing many eager-persistent writes, such as the Webserver workloads, HiNFS achieves almost the same scalability with NOVA and PMFS, both of which are much better than EXT2/EXT4+NVMMBD.

Varmail. HiNFS shows almost the same performance as NOVA and PMFS as they both access file data and metadata in NVMM directly, and in Varmail workloads, NOVA can achieve the best performance among the six file systems with 6 threads, as NOVA has many optimizations on metadata processing. But the throughput drops when the number of worker threads increased to 8 and 10, for our system is equipped with only one CPU (6 physical cores), and cannot fully exploit the scalability of NOVA.

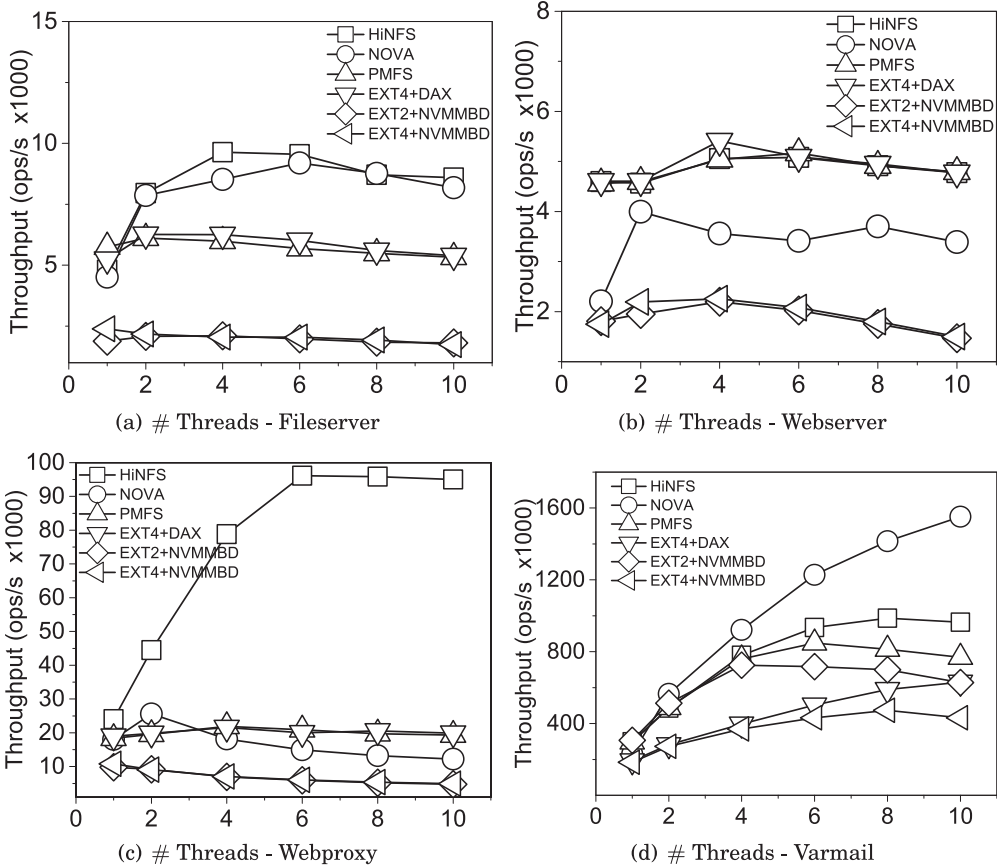


Fig. 9. Throughput (operations per second) for 1–10 threads (Quartz emulator).

5.2.3 System Scalability Measured by Quartz. Moreover, we use Quartz to evaluate the system scalability of these file systems for comparison (as shown in Figure 9). HiNFS can outperform other file systems in general, which shows the same trend as in Figure 8, and can be explained by above analysis. However, the throughput between NOVA, PMFS and EXT4+DAX shows different behavior.

Fileserver. In this workload, NOVA has relatively higher throughput than PMFS and EXT4+DAX, which is close to that of HiNFS, and this is different from the data in Figure 8 (NOVA has the same throughput as that of PMFS and EXT4+DAX). However, the experiment result with Quartz shows more similar trend as in NOVA article, where NOVA outperform PMFS and EXT4+DAX with Fileserver workload, and this indicating that Quartz emulator may work more accurately.

Webserver and Webproxy. For these two workloads, NOVA, PMFS, and EXT4+DAX exhibit extremely different performance behavior, as these two workloads are both read-dominated and Quartz cannot inject accurate latency in each epoch. Quartz can only support symmetrical read-write latency model and inject latency by calculating PMC; however, these counters cannot be separated to represent read operations and write operations, respectively. In our experiment assumption, persistent memory provides asymmetrical read-write performance and read latency is

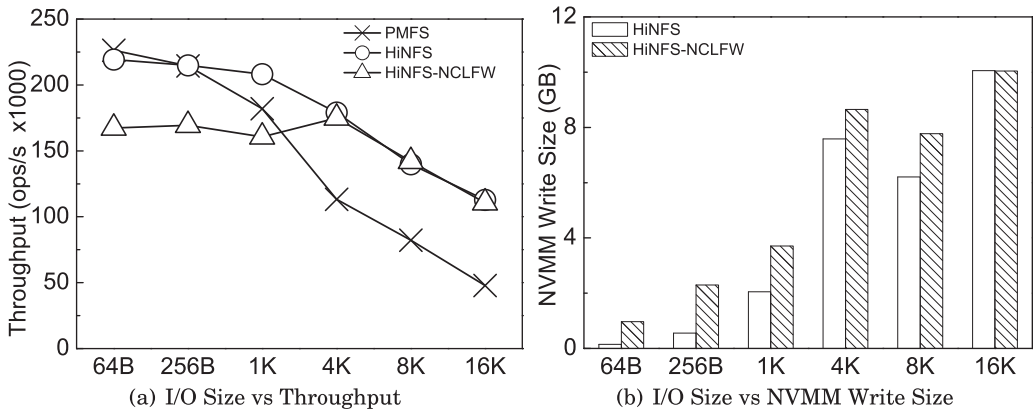


Fig. 10. Throughput (operations per second) and NVMM write size with different I/O sizes for fileserver workload.

close to that of DRAM. So under Quartz emulation, file systems (like PMFS, EXT4+DAX, and NOVA) that running read-dominated workloads will be interrupted more frequently by Quartz and show inaccurate throughput.

Varmail. NOVA can outperform HiNFS in Varmail workload, and when threads increase, the gap between NOVA and other file systems become larger, this is because Varmail contains a large part of synchronization operations and the buffer cache in HiNFS has little effects, so HiNFS has no advantages. While evaluating with Quartz, there are 48 CPU cores available, and NOVA is designed with consideration of NUMA architecture. NOVA also add optimization on metadata processing, which account for the high performance of NOVA when running Varmail, for Varmail contains many metadata operations. Besides, Quartz can only emulate symmetrical read-write latency model and this can also lead to the deviation in this experiment. As a result, experiment result here shows different characteristics from that on NVMM emulator.

5.2.4 Sensitivity Analysis. As the I/O size of the workload, the DRAM buffer size, and the NVMM write latency can affect the system performance, we measure their impacts on HiNFS's performance in this section.

Sensitivity to the I/O Size. The I/O size of the workload can affect the performance. Figure 10(a) presents the throughput performance with different I/O sizes for the Fileserver workload. For brevity, we omit the other three workloads. Webserver is a read-intensive workload while Varmail includes a large portion of eager-persistent writes, both of which cannot benefit from the DRAM buffer, thus HiNFS always yields performance similar to PMFS with different I/O sizes. We omit the Webproxy workload, because it shows similar results with the Fileserver workload. To investigate the benefits of the *CLFW* scheme, we compare the performance and NVMM write sizes (i.e., total bytes that are written to NVMM) of HiNFS and HiNFS-NCLFW. HiNFS-NCLFW is a version of HiNFS that does not implement the *CLFW* scheme.

From Figure 10(a), we observe that HiNFS and HiNFS-NCLFW show a great difference in throughput when the I/O size is less than the DRAM block size (i.e., 4KB), and HiNFS shows up to nearly 30% performance improvement over HiNFS-NCLFW. From Figure 10(b), we can see that HiNFS shows a remarkable drop in NVMM write size compared to HiNFS-NCLFW when the I/O size is less than the DRAM block size. The reason is that the background NVMM write traffic can also impact the system performance, because when the DRAM buffer is full, the normal writing

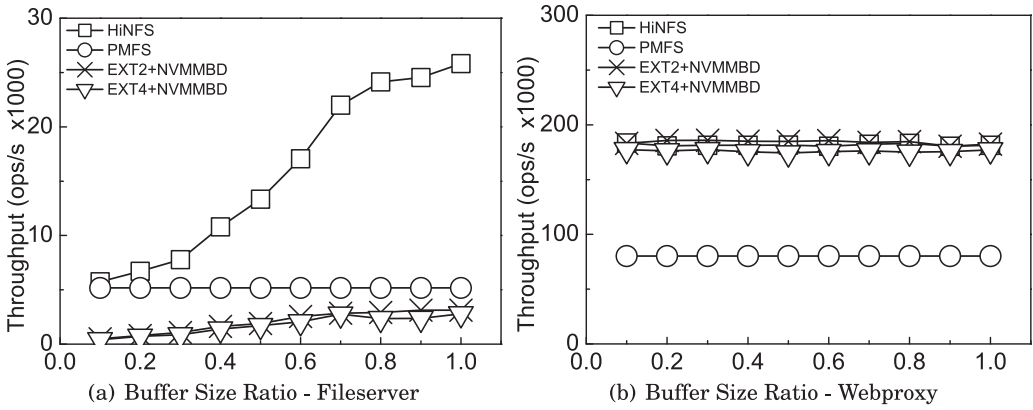


Fig. 11. Throughput (operations per second) as a function of the DRAM buffer size.

threads may need to wait for the background writeback threads to clean out free buffer blocks. HiNFS significantly reduces the NVMM write traffic when the I/O size is unaligned to the DRAM block size, thereby improving the system performance. In contrast, the performance gap between them is bridged when the I/O size is larger than and aligned to the DRAM block size.

We also make another observation from Figure 10(a) that the performance gap between HiNFS and PMFS grows as the I/O size increases. For example, HiNFS outperforms PMFS by 58% when the I/O size is 4KB, while improves the performance by 136% over PMFS when the I/O size is 16 KB. This is mainly due to that the copy overheads gradually become relatively more significant than other parts as the I/O size increases. When the I/O size is small (e.g., 64B), the overheads from other parts, such as system call, user-kernel mode switch, and so on, become dominant, thus hiding the benefits of reducing the copy overheads.

Sensitivity to the DRAM Buffer Size. The DRAM buffer size also has a strong impact on HiNFS's performance. Figure 11 shows the throughput performance as we vary the buffer size from 0.1 (10%) to 1.0 (100%) relative to the workload size. In Figure 11, we observe that the performance of HiNFS exhibits great improvement as the buffer size increases for the Fileserver workload, because more write operations will hit in the buffer when the buffer size increases. However, HiNFS's throughput remains nearly unchanged for the Webproxy workload when the buffer size ratio goes from 0.1 to 1.0 due to that the Webproxy workload has strong locality. Moreover, we find that the Webproxy workload exhibits many short-lived files, which would be deleted before the written data is flushed to NVMM. Therefore, the Webproxy workload is insensitive to the buffer size, and this is the only case where EXT2/EXT4+NVMMBD and HiNFS show nearly the same performance. For the Fileserver workload, EXT2/EXT4+NVMMBD have much lower performance than PMFS even when the buffer size ratio is 1.0, this is due to that the read copy overhead degrades the overall performance. Before running the benchmark, we clear the contents of the OS page cache, so the read operations should first fetch the data from the NVMM storage into the DRAM buffer through the generic block layer. The overheads from the double-copy and the generic block layer significantly degrade their performance.

Sensitivity to the NVMM Write Latency. Another aspect that can affect the system performance is the NVMM write latency. Figure 12 shows the throughput performance when we vary the NVMM write latency from 50 to 800ns using a single thread. In this figure, we can observe that the performance benefits of HiNFS become more obvious with longer NVMM write latency. For instance, HiNFS outperforms PMFS by only 53% when the NVMM write latency is 100ns, but

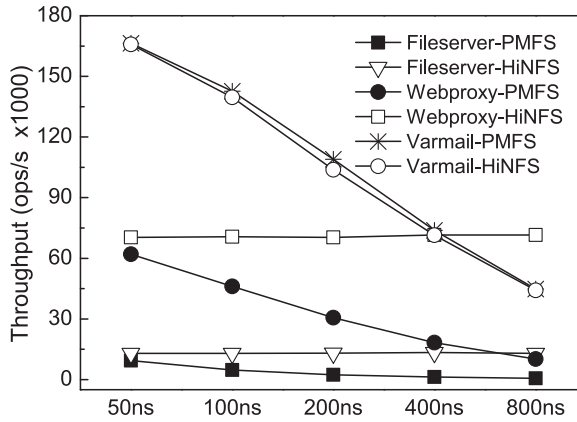


Fig. 12. Throughput (operations per second) for different NVMM write latencies.

improves the performance by nearly $6\times$ over PMFS when the NVMM write latency is 800ns for the Webproxy workload. This is attributed to the fact that the system can get more performance benefits from the DRAM buffer as the speed gap between DRAM and NVMM increases. Even when the write latency of NVMM is close to that of DRAM (e.g., 50ns), HiNFS still performs no worse than PMFS. This is because most of the write operations, in this case, will bypass the DRAM buffer with the *Buffer Benefit Model*, thereby eliminating the high double-copy overheads.

5.3 Data-Intensive Traces and Macrobenchmarks

To further investigate the performance of HiNFS and other file systems on real workloads, we replay a series of traces and run a set of macrobenchmarks on these file systems. In these experiments, the DRAM buffer size is set to 1/10 of the workload size by default. To demonstrate the benefits of bypassing the buffer for the eager-persistent writes, we also compare HiNFS with HiNFS-WB. HiNFS-WB refers to a system that simply uses DRAM as a write buffer of NVMM, which is implemented by closing the function of the *Eager-Persistent Write Checker* in HiNFS. In HiNFS-WB, all the writes are buffered in DRAM first before being persisted to NVMM. For the traces replay, all the traces are system call level I/O traces, and we extract the *read*, *write*, *unlink*, and *fsync* operations from the traces, and replay them on the five different file systems. Moreover, we collect the time spent on these four different types of I/O operations respectively, and report a breakdown of the execution time in Figure 13. For the macrobenchmarks, we report the normalized runtime of all the benchmarks and show the results in Figure 14.

Lazy-persistent Workload. In Figure 13, we observe that HiNFS exhibits a reduction in execution time when comparing with PMFS by 37%, 35%, and 38% for the Ustr0, Ustr1, and LASR traces, respectively, and outperform NOVA by 24%, 26%, and 53%. As we can see in the figure, this is mainly attributed to the reduction of the write time of HiNFS compared to PMFS and NOVA.

Eager-persistent Workload. HiNFS significantly outperforms PMFS except the Facebook trace, in which they yield similar performance. When we analyze this trace, we find that it contains a significant amount of sync operations. Moreover, we observe that HiNFS sets most of the related data blocks to the Eager-Persistent state with the *Buffer Benefit Model* in this case. Thus, it bypasses the DRAM buffer for most writes that are directly performed to NVMM, because the sync operations in this workload appear too frequent to coalesce enough writes in the DRAM buffer.

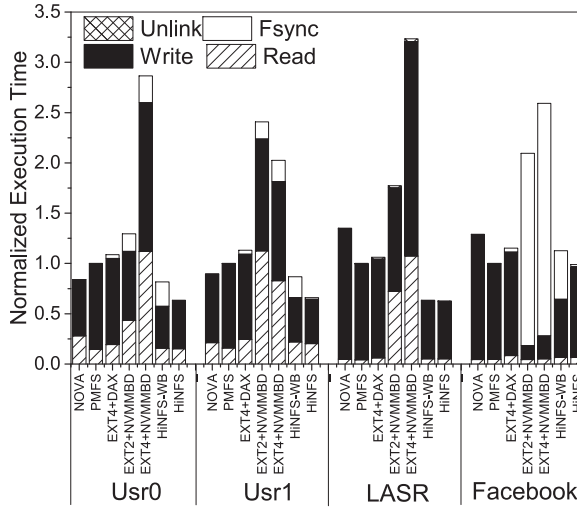


Fig. 13. Breakdown of the time spent on replaying traces. *Normalized to PMFS's execution time.*

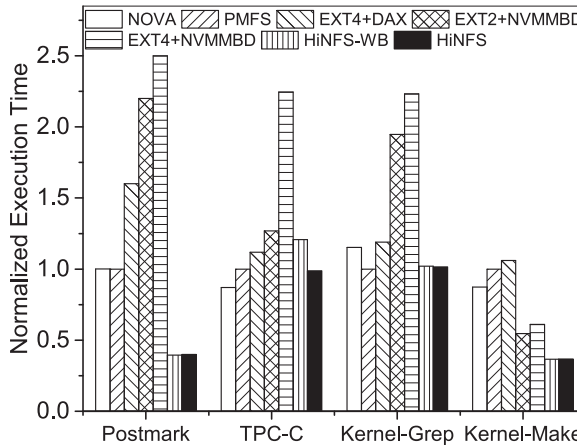


Fig. 14. The elapsed time of running macrobenchmarks. *Normalized to PMFS's execution time.*

Worth noticing, NOVA shows lower execution time than PMFS in Usr0 and Usr1 traces, and exhibits higher execution time in LASR and Facebook traces, but NOVA shows higher execution time on all four traces compared to HiNFS.

In Figure 14, HiNFS reduces the execution time of running the Postmark and Kernel-Make benchmarks by 60% and 64%, respectively, when comparing with PMFS, and 60% and 58%, when comparing with NOVA. We find that the Postmark workload contains many short-lived files, where many lazy-persistent writes in this workload can benefit from the DRAM buffer for HiNFS, as writes to these files that are later deleted do not need to be performed to NVMM. In the remaining two cases (i.e., TPC-C and Kernel-Grep), we can see that HiNFS and NOVA/(PMFS, EXT4+DAX) show nearly the same performance, all of which exhibit a remarkable drop in execution time when comparing with EXT2/EXT4+NVMDB. We find that Kernel-Grep is a read-intensive workload while TPC-C contains many sync operations. In these cases, HiNFS bypasses the DRAM buffer for

most I/O operations. This set of experiments also demonstrate the notable benefits of eliminating the double-copy overheads. In this figure, we also observe that EXT2+NVMMBD is much faster than EXT4+NVMMBD due to the absence of the journaling-related overheads.

Comparing HiNFS with HiNFS-WB in the two figures, we can see that HiNFS-WB increases the execution time over HiNFS by 28%, 32%, 14%, and 22% for the Ustr0, Ustr1, Facebook, and TPC-C workloads, respectively. As buffering the eager-persistent writes not only increases the system copy overheads, but also may evict other valuable buffer blocks, which in turn decreases the ratio of write coalescing and increases the buffer writeback traffic, this performance improvement with HiNFS is due to that it effectively identifies the eager-persistent write operations, and then performs them to NVMM directly, demonstrating the benefits of the direct eager-persistent write policy of HiNFS. In other workloads, these two systems yield similar performance due to the absence of the synchronization operations in these workloads. However, because of the small mean I/O size (less than 1KB) exhibited in the Facebook workload, we observe that it shows less difference between HiNFS and HiNFS-WB than that in the Ustr0, Ustr1, and TPC-C workloads.

6 RELATED WORK

In this section, we discuss and draw connections to classes of previous works that are closely related.

6.1 Buffering and Caching in Non-Volatile Memories

In storage systems, buffering and caching are commonly used to improve performance by increasing cache hit ratio. However, with the advent of non-volatile memories, including both flash memory and persistent memory, buffering and caching are used with different optimization goals.

In flash-based storage systems, buffers are introduced mostly because flash memory has slower writes than reads and much slower erases. Considering the relatively poor write performance of the flash memory, some cache studies (Jo et al. 2006; Kim and Ahn 2008; Kang et al. 2009) have investigated how to increase its write performance using a RAM write buffer. Flash-Aware Buffer (FAB) management (Jo et al. 2006) groups pages in the same flash block and evicts the group that has the largest number of pages when the buffer is full. However, FAB only considers the group size while overlooking the recency. To accommodate both the temporal locality and group size, the Cold and Largest Cluster (CLC) policy (Kang et al. 2009) combines the FAB and LRU algorithms. Both the FAB and CLC schemes aim to reduce the number of write and erase operations of the flash memory. In contrast, the Block Padding Least Recently Used (BPLRU) strategy (Kim and Ahn 2008) focuses on optimizing the random write performance of the flash memory by establishing a desirable write pattern with RAM buffering.

However, these flash-aware write buffer policies are not suitable for the NVMM storage due to the following reasons: First, they manage the buffer space at the page granularity rather than the cacheline level, which will generate a large amount of wasteful fetching and flushing data. Second, their designs are based on the unique characteristics of the flash storage, such as reducing the random write or erase operations, most of which are not applicable to the NVMM storage, as the random and sequential access of existing NVMM technologies are nearly identical and they have no erase operations. In contrast, the relatively high performance of existing NVMM technologies indicates that the system designers should carefully deal with the copy overheads among the user buffer, the file system buffer, and the NVMM storage (DAX 2014; Dulloor et al. 2014). Therefore, HiNFS's write buffer policy is highly optimized for the NVMM storage, which focuses on reducing unnecessary data-fetch and data-flush by leveraging the unique characteristics of NVMM's byte addressability, and eliminating the double-copy overheads resulted from conventional buffer management from the critical path, thereby improving the NVMM system performance.

With high-speed storage medias, like PCM, have emerged recently, the performance gap between the main memory and the storage device drops dramatically. To figure out whether the buffer cache is still effective for them, Lee and Bahn (2014) propose a new buffer cache management scheme appropriately designed for the system where the speed gap between cache and storage is narrow. To our knowledge, this is the only work that analyzes the effectiveness of the buffer cache under the fast NVM storage. Our work differs from them in the following aspects: First, their work is based on the assumption that NVM sits behind the I/O bus, while our work assumes that NVM is attached directly to the memory bus. Second, they aim to optimize the OS page cache and focus on improving the hit ratio of the buffer cache. HiNFS, in contrast, completely replaces the OS page cache with a new DRAM write buffer using a novel NVMM-aware buffer policy, which is cacheline-oriented and eliminates the software stack overhead of the block device layer altogether. Finally, their algorithm copies data to the buffer cache first for all file operations, which will incur the double-copy overheads. Based on our observation, these overheads are non-trivial for NVM storage system. HiNFS, therefore, buffers only the lazy-persistent writes, while uses direct access for reads and eager-persistent writes to eliminate the double-copy overheads from the critical path.

6.2 File Systems on Non-Volatile Memories

Flash File Systems. File systems need to be redesigned for non-volatile memories, which differ a lot from hard disk drives in many aspects. Direct File System (DFS) (Josephson et al. 2010) is designed on top of FusionIO's ioDrive to leverage the Flash Translation Layer (FTL) for data allocation, in order to avoid duplicated data allocations in both file system and FTL layers. Objected Flash Storage System (OFSS) (Lu et al. 2013) proposes to re-architect the storage stack for flash-based storage systems. OFSS directly managed raw-flash devices (a.k.a., open-channel SSDs) via software in an object-based way. With the co-design of both software and hardware, OFSS introduces endurance-aware design to file systems to significantly reduce the write amplification inside file systems. ReconFS (Lu et al. 2014) then reorganizes the namespace management of file systems by trading low metadata write cost off recovery performance, which is acceptable due to fast read performance. F2FS (Lee et al. 2015) takes a less aggressive approach. It keeps the same read/write interface to flash SSDs, but optimize data layout to flash memory characteristics. ParaFS (Zhang et al. 2016) further studies the collaborations between software and hardware in open-channel SSDs, and proposes to simplify the FTL and export the internal physical layout to file systems for parallelism exploration in the system software. The above-mentioned designs successfully demonstrate the benefits of removing or reorganizing duplicated layers in storage systems based non-volatile memories. These design concepts are partially absorbed in persistent memory file system designs.

Persistent Memory File Systems. For byte-addressable NVMM, a number of file systems have been proposed. BPFS (Condit et al. 2009) uses shadow paging techniques and 8byte atomic updates to provide fast and consistent updates. However, BPFS doesn't support mmap and relies on a hardware approach (epochs) to support data persistence and ordering. While HiNFS is not optimized for mmap I/O, it still supports direct mmap access. PMFS (Dulloor et al. 2014) is a light-weight file system that is optimized for persistent memory, it avoids the block layer and eliminates the copy overheads by enabling applications to access persistent memory directly. Similar to PMFS's direct access policy, DAX (DAX 2014; EXT 2014) is a kernel patch that can support traditional ext4 file system for bypassing the OS page cache and direct access to memory-like storage. However, all three file systems discussed above do not take into account NVMM's slow write operations, and direct access to NVMM for all file operations leads to suboptimal system performance. In contrast, HiNFS buffers the lazy-persistent writes in the DRAM buffer, which can hide the long NVMM write latency, thereby improving the performance.

SCMFS (Wu and Reddy 2011) leverages the OS VMM to reduce the complexity of the file system. Aerie (Volos et al. 2014) provides flexible file system interfaces to reduce the hierarchical file system abstraction. Both SCMFS and Aerie focus on reducing the software overheads. However, based on our analysis, only in cases of metadata-intensive workloads or workloads with a small mean I/O size can the software overheads become relatively more significant than the storage access overheads. HiNFS, in contrast, focuses on reducing the storage access overheads (i.e., copy overheads) for data-intensive workloads.

NOVA (Xu and Swanson 2016) is a recently proposed log-structured file system for persistent memory and is optimized with per-core log for currency exploration. NOVA also optimizes the indexing in DRAM while keeping data persistent in NVMM for better performance. Different from the index buffering in NOVA, HiNFS performs buffering for file system operations and keeps the lazy-persistent writes in the DRAM for lazy persistence.

In addition to above-mentioned local file systems, Octopus (Lu et al. 2017) is a distributed persistent memory file system that is built from scratch based on NVMM and RDMA. Octopus further reduces memory copies in remote I/Os leveraging features of both NVMM and RDMA. It also re-designs a number of data or metadata mechanisms in distributed file systems. While Octopus is designed for distributed storage, the combination usage of buffering and direct access in HiNFS can be integrated into Octopus.

6.3 Persistent Memory Techniques

Hybrid NVMM/DRAM Architecture. To take advantage of both DRAM and NVMM, hybrid PCM/DRAM memory systems have been discussed (Qureshi et al. 2009; Ramos et al. 2011). Qureshi et al. (2009) use a DRAM device as a cache of PCM in the hierarchy, while Ramos et al. (2011) present a page placement policy on memory controller to implement PCM-DRAM hybrid memory systems. Different from these architectural level designs, our proposed HiNFS uses buffering in the system level, which can take advantages of file system semantics and is more effective in exploring the buffering benefits.

Consistency and Reliability in Persistent Memory. Consistency and reliability are also important design issues in persistent memory. Consistency and reliability techniques are proposed in the NVMM-based programming models (Coburn et al. 2011; Volos et al. 2011; Hwang et al. 2015; Lu et al. 2015, 2016), using persistent data structures (Venkataraman et al. 2011; Yang et al. 2015), or with hardware support (Moraru et al. 2013; Zhao et al. 2013; Sun et al. 2015; Lu et al. 2014, 2017), and in distributed systems (Zhang et al. 2015). However, it requires significant modifications or hardware support for legacy applications to use the above-mentioned consistency or reliability mechanisms. Rather than using these techniques, HiNFS provides consistency in the file system by ordering the writes and reusing PMFS's journaling mechanism. We also propose a combined redo and undo logging for the journaling optimization in FCFS (Ou and Shu 2016), which journaling technique can also be incorporated into HiNFS for consistency.

7 CONCLUSION

Direct Access, which is now an overwhelming approach in persistent memory file system designs, is not a panacea. This is the major reason that existing persistent memory file systems suffer from slow writes. However, buffering can hide long write latency for lazy-persistent writes due to higher I/O performance in DRAM against NVMM. Our proposed HiNFS combines buffering and direct access for fine-grained file system operations. It buffers the lazy-persistent writes in DRAM temporarily to hide the long write latency of NVMM, while eliminating the double-copy overheads resulted from conventional buffer management by using direct access for reads and

eager-persistent writes. Extensive evaluations demonstrate that HiNFS significantly outperforms both traditional block-based file systems and state-of-the-art NVMM-aware file systems.

REFERENCES

- Dbt2 test suite. Retrieved from <http://sourceforge.net/apps/mediawiki/osdl/dbt>.
- Filebench 1.4.9.1. Retrieved from <http://sourceforge.net/projects/filebench/>.
- FIU system call io trace. Retrieved from <http://syllab-srv.cs.fiu.edu/dokuwiki/doku.php?id=projects:nbw:start>.
- Flexible IO (fio) Tester. Retrieved from <http://freecode.com/projects/fio>.
- Lasr system call io trace. Retrieved from http://iotta.snia.org/historical_section?tracetype_id=1.
2014. Support ext4 on NV-DIMMs. Retrieved from <http://lwn.net/Articles/588218/>.
2014. Supporting file systems in persistent memory. Retrieved from <https://lwn.net/Articles/610174/>.
- Geoffrey W. Burr, Matthew J. Breitwisch, Michele Franceschini, Davide Garetto, Kailash Gopalakrishnan, Bryan Jackson, Bülent Kurdi, Chung Lam, Luis A. Lastras, Alvaro Padilla, Bipin Rajendran, Simone Raoux, and Rohit S. Shenoy. 2010. Phase change memory technology. *J. Vacuum Sci. Technol. B* 28, 2 (2010), 223–262.
- Daniel Campello, Hector Lopez, Ricardo Koller, Raju Rangaswami, and Luis Useche. 2015. Non-blocking writes to files. In *Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'15)*. 151–165.
- Shimin Chen, Phillip B. Gibbons, and Suman Nath. 2011. Rethinking database algorithms for phase change memory. In *Proceedings of the 5th Biennial Conference on Innovative Data Systems Research (CIDR'11)*. 21–31.
- Joel Coburn, Adrian M. Caulfield, Ameen Akel, Laura M. Grupp, Rajesh K. Gupta, Ranjit Jhala, and Steven Swanson. 2011. NV-heaps: Making persistent objects fast and safe with next-generation, non-volatile memories. In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'11)*. 105–118.
- Edward Grady Coffman and Peter J. Denning. 1973. *Operating Systems Theory*. Vol. 973.
- Jeremy Condit, Edmund B. Nightingale, Christopher Frost, Engin Ipek, Benjamin Lee, Doug Burger, and Derrick Coetzee. 2009. Better I/O through byte-addressable, persistent memory. In *Proceedings of the ACM SIGOPS 22nd Symposium on Operating Systems Principles (SOSP'09)*. 133–146.
- Intel Cooperation. 2015. NVDIMM Namespace Specification. Retrieved from http://pmem.io/documents/NVDIMM_Namespace_Spec.pdf.
- Intel Cooperation. 2016. Intel Architecture Instruction Set Extensions Programming Reference. Retrieved from <https://software.intel.com/sites/default/files/managed/b4/3a/319433-024.pdf>.
- Peter J. Denning. 1968. The working set model for program behavior. *Commun. ACM* 11, 5 (1968), 323–333.
- Edward Doller. 2009. Phase change memory and its impacts on memory hierarchy. Retrieved from <http://www.pdl.cmu.edu/SDI/2009/slides/Numonyx.pdf>.
- Subramanya R. Dulloor, Sanjay Kumar, Anil Keshavamurthy, Philip Lantz, Dheeraj Reddy, Rajesh Sankaran, and Jeff Jackson. 2014. System software for persistent memory. In *Proceedings of the 9th European Conference on Computer Systems (EuroSys'14)*. 15:1–15:15.
- Tyler Harter, Chris Dragga, Michael Vaughn, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. 2011. A file is not a file: Understanding the I/O behavior of apple desktop applications. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP'11)*. 71–83.
- Jian Huang, Karsten Schwan, and Moinuddin K. Qureshi. 2014. NVRAM-aware logging in transaction systems. *Proc. VLDB Endow.* 8, 4 (Dec. 2014), 389–400.
- Taeho Hwang, Jaemin Jung, and Youjip Won. 2015. Heapo: Heap-based persistent object store. *ACM Trans. Storage (TOS)* 11, 1 (2015), 3.
- Lei Jiang, Bo Zhao, Youtao Zhang, Jun Yang, and B. R. Childers. 2012. Improving write operations in MLC phase change memory. In *Proceedings of the 18th International Symposium on High Performance Computer Architecture (HPCA'12)*. 1–10.
- Heeseung Jo, Jeong-Uk Kang, Seon-Yeong Park, Jin-Soo Kim, and Joonwon Lee. 2006. FAB: Flash-aware buffer management policy for portable media players. *IEEE Trans. Consum. Electron.* 52, 2 (May 2006), 485–493.
- Theodore Johnson and Dennis Shasha. 1994. 2Q: A low-overhead high-performance buffer management replacement algorithm. In *Proceedings of the 20th International Conference on Very Large Data Bases (VLDB'94)*. 439–450.
- William K. Josephson, Lars A. Bongo, David Flynn, and Kai Li. 2010. DFS: A file system for virtualized flash storage. In *Proceedings of the 8th USENIX Conference on File and Storage Technologies (FAST'10)*. USENIX, Berkeley, CA.
- Myoungsoo Jung, John Shalf, and Mahmut Kandemir. 2013. Design of a large-scale storage-class RRAM system. In *Proceedings of the 27th International ACM Conference on International Conference on Supercomputing (ICS'13)*. 103–114.
- Sooyong Kang, Sungmin Park, Hoyoung Jung, Hyoki Shim, and Jaehyuk Cha. 2009. Performance trade-offs in using NVRAM write buffer for flash memory-based storage devices. *IEEE Trans. Comput. (TC)* 58, 6 (June 2009), 744–758.
- Jeffrey Katcher. 1997. *Postmark: A New File System Benchmark*. Technical Report TR3022, Network Appliance.

- Hyojun Kim and Seongjun Ahn. 2008. BPLRU: A buffer management scheme for improving random writes in flash storage. In *Proceedings of the 6th USENIX Conference on File and Storage Technologies (FAST'08)*. 16:1–16:14.
- ESOS LABORATORY. 2013. Mobibench benchmark tool. Retrieved from <http://www.mobibench.co.kr/>.
- B. C. Lee, Ping Zhou, Jun Yang, Youtao Zhang, Bo Zhao, E. Ipek, O. Mutlu, and D. Burger. 2010. Phase-change technology and the future of main memory. *Micro, IEEE* 30, 1 (Jan 2010), 143–143.
- Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger. 2009. Architecting phase change memory as a scalable dram alternative. In *Proceedings of the 36th Annual International Symposium on Computer Architecture (ISCA'09)*. 2–13.
- Changman Lee, Dongho Sim, Jooyoung Hwang, and Sangyeun Cho. 2015. F2FS: A new file system for flash storage. In *Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'15)*. USENIX, Santa Clara, CA. Retrieved from <https://www.usenix.org/conference/fast15/technical-sessions/presentation/lee>.
- Eunji Lee and Hyokyung Bahn. 2014. Caching strategies for high-performance storage media. *ACM Trans. Storage (TOS)* 10, 3 (Aug. 2014), 11:1–11:22.
- Youyou Lu, Jiwu Shu, Youmin Chen, and Tao Li. 2017. Octopus: An RDMA-enabled distributed persistent memory file system. In *Proceedings of the 2017 USENIX Annual Technical Conference (USENIX ATC'17)*. USENIX Association, 773–785.
- Y. Lu, J. Shu, and L. Sun. 2015. Blurred persistence in transactional persistent memory. In *Proceedings of the 31st Symposium on Mass Storage Systems and Technologies (MSST'15)*. 1–13.
- Youyou Lu, Jiwu Shu, and Long Sun. 2016. Blurred persistence: Efficient transactions in persistent memory. *ACM Trans. Storage (TOS)* 12, 1 (2016), 3.
- Y. Lu, J. Shu, L. Sun, and O. Mutlu. 2014. Loose-ordering consistency for persistent memory. In *Proceedings of the 32nd IEEE International Conference on Computer Design (ICCD'14)*. 216–223.
- Youyou Lu, Jiwu Shu, Long Sun, and Onur Mutlu. 2017. Improving performance and endurance of persistent memory with loose-ordering consistency. *IEEE Trans. Parallel Distrib. Syst.* PP, 99 (2017), 1–1.
- Youyou Lu, Jiwu Shu, and Wei Wang. 2014. ReconFS: A reconstructable file system on flash storage. In *Proceedings of the 12th USENIX Conference on File and Storage Technologies (FAST'14)*. USENIX, Berkeley, CA, 75–88.
- Youyou Lu, Jiwu Shu, and Weimin Zheng. 2013. Extending the lifetime of flash-based storage through reducing write amplification from file systems. In *Proceedings of the 11th USENIX Conference on File and Storage Technologies (FAST'13)*. USENIX, Berkeley, CA.
- Nimrod Megiddo and Dharmendra S. Modha. 2003. ARC: A self-tuning, low overhead replacement cache. In *Proceedings of the 2nd USENIX Conference on File and Storage Technologies (FAST'03)*. 115–130.
- Changwoo Min, Kangnyeon Kim, Hyunjin Cho, Sang-Won Lee, and Young Ik Eom. 2012. SFS: Random write considered harmful in solid state drives. In *Proceedings of the 10th USENIX Conference on File and Storage Technologies (FAST'12)*. San Jose, CA, 139–154.
- Iulian Moraru, David G. Andersen, Michael Kaminsky, Niraj Tolia, Parthasarathy Ranganathan, and Nathan Binkert. 2013. Consistent, durable, and safe memory management for byte-addressable nonvolatile main memory. In *Proceedings of the First ACM SIGOPS Conference on Timely Results in Operating Systems (TRIOS'13)*. 1:1–1:17.
- Edmund B. Nightingale, Kaushik Veeraraghavan, Peter M. Chen, and Jason Flinn. 2006. Rethink the sync. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation (OSDI'06)*. 1–14.
- Jiaxin Ou and Jiwu Shu. 2016. Fast and failure-consistent updates of application data in non-volatile main memory file system. In *Proceedings of the 32st Symposium on Mass Storage Systems and Technologies (MSST'16)*.
- Jiaxin Ou, Jiwu Shu, and Youyou Lu. 2016. A high-performance file system for non-volatile main memory. In *Proceedings of the 11th European Conference on Computer Systems*. ACM, 12.
- Jiaxin Ou, Jiwu Shu, Youyou Lu, Letian Yi, and Wei Wang. 2014. EDM: An endurance-aware data migration scheme for load balancing in SSD storage clusters. In *Proceedings of the 28th International Parallel and Distributed Processing Symposium (IPDPS'14)*. 787–796.
- Steven Pelley, Peter M. Chen, and Thomas F. Wenisch. 2014. Memory persistency. In *Proceedings of the 41st Annual International Symposium on Computer Architecture (ISCA'14)*. 265–276.
- Moinuddin K. Qureshi, Vijayalakshmi Srinivasan, and Jude A. Rivers. 2009. Scalable high-performance main memory system using phase-change memory technology. In *Proceedings of the 36th Annual International Symposium on Computer Architecture (ISCA'09)*. 24–33.
- Luiz E. Ramos, Eugene Gorbatov, and Ricardo Bianchini. 2011. Page placement in hybrid memory systems. In *Proceedings of the International Conference on Supercomputing (ICS'11)*. 85–95.
- Drew Roselli, Jacob R. Lorch, and Thomas E. Anderson. 2000. A comparison of file system workloads. In *Proceedings of the Annual Conference on USENIX Annual Technical Conference (ATEC'00)*. 41–54.
- Chris Ruemmler and John Wilkes. 1993. UNIX disk access patterns. In *Proceedings of the USENIX Winter Conference*, Vol. 93. 405–420.

- Long Sun, Youyou Lu, and Jiwu Shu. 2015. DP2: Reducing transaction overhead with differential and dual persistency in persistent memory. In *Proceedings of the ACM International Conference on Computing Frontiers (CF'15)*. ACM.
- Kosuke Suzuki and Steven Swanson. 2015. *The Non-Volatile Memory Technology Database (NVMDB)*. Technical Report CS2015-1011. Department of Computer Science & Engineering, University of California, San Diego.
- Shivaram Venkataraman, Niraj Tolia, Parthasarathy Ranganathan, and Roy H. Campbell. 2011. Consistent and durable data structures for non-volatile byte-addressable memory. In *Proceedings of the 9th USENIX Conference on File and Storage Technologies (FAST'11)*. San Jose, CA, 61–75.
- Haris Volos, Guilherme Magalhaes, Ludmila Cherkasova, and Jun Li. 2015. Quartz: A lightweight performance emulator for persistent memory software. In *Proceedings of the 16th Annual Middleware Conference*. ACM, 37–49.
- Haris Volos, Sanketh Nalli, Sankarlingam Panneerselvam, Venkatanathan Varadarajan, Prashant Saxena, and Michael M. Swift. 2014. Aerie: Flexible file-system interfaces to storage-class memory. In *Proceedings of the 9th European Conference on Computer Systems (EuroSys'14)*. 14:1–14:14.
- Haris Volos, Andres Jaan Tack, and Michael M. Swift. 2011. Mnemosyne: Lightweight persistent memory. In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'11)*. 91–104.
- D. L. Willick, D. L. Eager, and R. B. Bunt. 1993. Disk cache replacement policies for network file servers. In *Proceedings the 13th International Conference on Distributed Computing Systems (ICDCS'93)*. 2–11.
- Xiaojuan Wu and A. L. Narasimha Reddy. 2011. SCMFs: A file system for storage class memory. In *Proceedings of the 2011 International Conference for High Performance Computing, Networking, Storage, and Analysis (SC'11)*. 39:1–39:11.
- Jian Xu and Steven Swanson. 2016. NOVA: A log-structured file system for hybrid volatile/non-volatile main memories. In *Proceedings of the 14th USENIX Conference on File and Storage Technologies (FAST'16)*. 323–338.
- Jun Yang, Qingsong Wei, Cheng Chen, Chundong Wang, Khai Leong Yong, and Bingsheng He. 2015. NV-tree: Reducing consistency cost for NVM-based single level systems. In *Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'15)*. 167–181.
- J. Joshua Yang and R. Stanley Williams. 2013. Memristive devices in computing system: Promises and challenges. *J. Emerg. Technol. Comput. Syst. (JETC'13)* 9, 2 (May 2013), 11:1–11:20.
- Jiacheng Zhang, Jiwu Shu, and Youyou Lu. 2016. ParaFS: A log-structured file system to exploit the internal parallelism of flash devices. In *Proceedings of the 2016 USENIX Annual Technical Conference (USENIX ATC'16)*.
- Yiyang Zhang and Steven Swanson. 2015. A study of application performance with non-volatile main memory. In *Proceedings of the 31st Symposium on Mass Storage Systems and Technologies (MSST'15)*. 1–10.
- Yiyang Zhang, Jian Yang, Amirsaman Memaripour, and Steven Swanson. 2015. Mojim: A reliable and highly-available non-volatile memory system. In *Proceedings of the 20th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'15)*. 3–18.
- Jishen Zhao, Sheng Li, Doe Hyun Yoon, Yuan Xie, and Norman P. Jouppi. 2013. Kiln: Closing the performance gap between systems with and without persistence support. In *Proceedings of the 46th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO'13)*. ACM, 421–432.
- Ping Zhou, Bo Zhao, Jun Yang, and Youtao Zhang. 2009. A durable and energy efficient main memory using phase change memory technology. In *Proceedings of the 36th Annual International Symposium on Computer Architecture (ISCA'09)*. 14–23.

Received February 2017; accepted July 2017