I/O Transit Caching for PMem-based Block Device

Qing Xu, Qisheng Jiang and Chundong Wang*

School of Information Science and Technology, ShanghaiTech University, Shanghai, China

ARTICLE INFO

Keywords: Persistent Memory Block Translation Table I/O Transit Caching

ABSTRACT

Byte-addressable non-volatile memory (NVM) sitting on the memory bus is employed to make *persistent memory* (PMem) in general-purpose computing systems and embedded systems for data storage. Researchers develop software drivers such as the block translation table (BTT) to build block devices on PMem, so programmers can keep using mature and reliable conventional storage stack while expecting high performance by exploiting fast PMem. However, our quantitative study shows that BTT underutilizes PMem and yields inferior performance, due to the absence of the imperative in-device cache. We add a conventional I/O staging cache made of DRAM space to BTT. As DRAM and PMem have comparable access latency, I/O staging cache is likely to be fully filled over time. Continual cache evictions and fsyncs thus cause on-demand flushes with severe stalls, such that the I/O staging cache is concretely unappealing for PMem-based block devices. We accordingly propose an algorithm named Caiti with novel *I/O transit caching*. Caiti eagerly evicts buffered data to PMem through CPU's multi-cores. It also conditionally bypasses a full cache and directly writes data into PMem to further alleviate I/O stalls. Experiments confirm that Caiti significantly boosts the performance with BTT by up to 3.6x, without loss of block-level write atomicity.

1. Introduction

Non-volatile memory (NVM) technologies bring about the availability of *persistent memory* (PMem) that is placed on the memory bus alongside DRAM, for CPU to load and store data. As NVM generally has shorter access latency and higher write endurance than flash memory, researchers have considered using it for data storage in generalpurpose computing systems, cloud and virtual machines (VMs), internet-of-things (IoT) endpoints, and broad embedded systems [2, 22, 29, 50]. Although Intel has discontinued its Optane DC memory business [3, 77], the exploration of NVM technologies continues. The other types of NVM, such as STT-RAM [16], are still being developed and deployed in embedded systems. A number of new file systems for PMem have been developed [15, 12, 11, 81, 33]. These file systems mainly follow the direct access (DAX) fashion to directly write and read files with PMem, bypassing DRAM page cache of operating system (OS). In spite of exploiting the performance potential of PMem, they fail to satisfy a few fundamental requirements raised by applications with respect to reliability, compatibility, deployment cost, and so on. One of them is block-level write atomicity, which means that writing a block (e.g., 512B or 4KB) of data shall be done in an atomic (all-or-nothing) manner. Typical applications and system softwares, such databases, rely on the atomic write of a block, because non-atomic write might leave a mix of up-to-date and stale data in one block and in turn cause the breach of data integrity. However, many of PMem-oriented file systems do not support this feature, while conventional block devices have it [35, 6, 44, 15, 33, 80].

As of today, PMem-oriented file systems have not been merged into Linux kernel or employed in production environments, mainly due to the issues of reliability, compatibility, and deployment cost concerned by applications. By contrast, conventional block-based file system mounted on a block device has been the de facto storage stack for decades. To be compatible with this storage stack, researchers have tried to build block devices on top of PMem, e.g., PMBD [10]. Intel's developers, regarding the need of block-level write atomicity with Optane memory, proposed the block translation table (BTT) [69]. BTT is actively maintained as a software driver in Linux kernel [65] for programmers to make block devices from raw PMem and mount conventional file systems like Ext4 or XFS. Compared to their DAX variants without block-level write atomicity (e.g., Ext4-DAX and XFS-DAX) that show 'WARNING: use at your own risk' upon being mounted, Ext4 and XFS are of maturity, robustness, and reliability after undergoing evolutions in past years. Programmers can continue using their tactics tuned for conventional storage stack with block I/O (bio) interfaces that BTT provides atop fast PMem.

As a software driver, BTT employs a combination of copy-on-write (CoW) and logging to achieve block-level write atomicity. However, the cost of doing so is non-trivial. We have done a comparative test with mature Ext4 on PMem formatted with BTT, mature Ext4 with raw PMem, and Ext4-DAX. Ext4 on PMem formatted with BTT yields inferior performance. For example, when randomly writing 64GB data with 4KB per I/O request, Ext4 with BTT spends 37.4% and 16.6% more time than Ext4 with raw PMem and Ext4-DAX, respectively. Whereas, neither of the latter two guarantees the block-level write atomicity.

When using fast PMem, we aim to yield high performance like raw PMem or Ext4-DAX while preserving the

^{*}Corresponding author at: School of Information Science and Technology, ShanghaiTech University, 393 Middle Huaxia Road, Pudong, Shanghai, China 201210.

^{🖄 (}cd_wang@outlook.com) (C. Wang)

ORCID(s): 0000-0001-9069-2650 (C. Wang)

block-level write atomicity for critical system and application softwares. To attain this aim, we thoroughly analyze the source code of BTT and compare it against traditional block devices. We find that a critical component is absent in it, i.e., the internal device cache¹ commonly installed in hard disk drives (HDD) or solid state drives (SSD) for I/O staging and acceleration [34, 74, 78, 73, 56]. The reason behind the absence of device cache is that BTT has been designed in line with DAX, by which data is directly written and read with PMem. Assuming that we complement BTT with a device cache, we could use it to absorb write requests. This is likely to conceal the performance overhead caused by maintaining block-level write atomicity at the BTT driver and in turn promote the overall performance. Inspired by this potential gain, we consider adding and managing a DRAM cache within BTT in the I/O staging fashion. We follow two common polices to manage the cache when the cache space is used up. One is a PMBD-like caching that flushes a batch of buffered blocks when the cache is filled to an extent (watermark), e.g., 70% or 100% full [10]. The other one evicts the least-recently-used (LRU) buffered block to make space. We expect them to improve the performance with BTT. However, with the foregoing test, both algorithms decrease the performance by 6.0% and 15.1%, respectively. When no free space is left in the cache, I/O requests have to stall for the drain of buffered data in one or multiple cache slots. Moreover, the upper-level file system such as Ext4 periodically (five seconds by default) issues a bio request with a flag named REQ_PREFLUSH being set to asynchronously flush the internal cache of block device [64, 66]. Frequent ondemand flushes are hence occurring over time. An explicit call of fsync is another source of cache flushes [45, 76]. Our further test confirms that I/O stalls also emerge when I/O staging caches serve fsyncs.

The PMBD-like and LRU caching policies represent the conventional I/O staging strategy, which, however, is practically ineffectual for BTT-like PMem-based block device with fast access speed. Aforementioned tests motivate us to consider what a gainful device cache shall be for PMem managed by BTT driver. Firstly, it shall be simple and efficient to incur minimal performance penalty, since PMem is relatively fast while BTT itself suffers from the software cost of preserving block-level write atomicity. Secondly, cache is likely to be fully filled at runtime. Using fsyncs to flush data is also common for applications. We need to ensure that flushing buffered data does not cause stalls on the critical path of serving I/O requests. Thirdly, as both DRAM and PMem are operated by CPU through the memory bus, we shall leverage multi-cores to handle I/O requests for high parallelism and efficiency. Last but not the least, as mentioned, the caching strategy must not impair blocklevel write atomicity but simultaneously retain all features required by applications and file systems, such as the support

for bio flags and fsyncs [45, 76]. Taking into account these concerns, we design a new caching algorithm named Caiti (\underline{ca} ching with \underline{I} /O \underline{t} rans \underline{i} t) for PMem-based block device. The main ideas of Caiti are summarized as follows.

- *Caiti eagerly evicts buffered data.* Once the cache receives a data block, Caiti promptly launches a writeback for eviction, instead of waiting for a flush or replacement. It places buffered blocks in multi-queues and engages a pool of background threads in concurrently moving them to PMem.
- *Caiti bypasses a fully filled cache.* When no space is left in DRAM cache and a cache miss occurs at the arriving write request, Caiti directly writes data to PMem.This avoids I/O congestion at the cache and further reduces response time on the critical path of serving I/O request.
- *Caiti exploits multi-core CPU for high concurrency and scalability.* Each bio request carries an *lba* and runs on a CPU core. Caiti logically partitions cache space into multiple sets and hashes each *lba* to find an appropriate set without maintaining a mapping table. With a core handling a bio request in one cache set, Caiti manages to simultaneously proceed concurrent bio requests. The multi-queues and background thread pool also accelerate concurrent write-backs of buffered data blocks to PMem and help to achieve scalability.

In contrast to I/O staging strategy intentionally buffers data for sufficiently long time, Caiti exploits scores of CPU cores to place data into cache in the foreground and swiftly *transits* data to PMem in the background to avoid I/O stalls. Extensive experiments on benchmarks and applications show that I/O transit caching enables Caiti to boost the performance with BTT by up to $3.6\times$. Besides PMBD and LRU, we further port and implement a state-of-the-art caching algorithm named Co-Active with PMem-based block device [61]. Caiti significantly outperforms them with up to $3.6\times$, $3.6\times$, and $2.9\times$ higher throughput, respectively.

The rest of this paper is organized as follows. In Section 2 we briefly present PMem and BTT. We show a motivational study in Section 3. We detail the design of Caiti in Section 4 and evaluate it in Section 5. We discuss related works in Section 6. We conclude the paper in Section 7.

2. Background

2.1. The Usage of Persistent Memory

PMem products are available in DRAM backed by flash memory (referred to as NVDIMM or NVDIMM-N) or NVM technologies such as spin-transfer torque RAM (STT-RAM) [16], phase-change memory (PCM) [23, 1, 42, 53, 25], and Intel Optane memory [28]. PCM and Optane memory technologies generally hold inferior access speeds compared to DRAM and are the main focus of this paper. General-purpose computing systems, cloud and VMs, and

¹We interchangeably refer to the cache inside a block device as *DRAM cache*, *internal cache*, or *device cache*. However, for distinguishing, we always refer to the kernel-space page cache of OS with the combinational terminology *page cache*, without calling it *DRAM cache* or *cache*.

embedded systems are using PMem for storage [2, 22, 29, 50]. As CPU loads and stores data with PMem through the memory bus, researchers have developed new file systems upon using PMem as memory device (e.g., [80, 33, 15]). These file systems mainly bypass OS's page cache with the DAX feature and directly operate with files stored in PMem. Without the involvement of OS's page cache, DAX helps to reduce the cost of traversing software stack and alleviate the overhead of memory copying. With regards to space efficiency, the DAX feature saves DRAM space that can be used for OS and applications.

Though, most of file systems developed with DAX are not merged in the mainline of Linux kernel. For Ext4and XFS-DAX that are already contained in Linux kernel because of being based on stable Ext4 and XFS, a mount of them shows a warning that says 'use at your own risk'. To deploy them in product environments has to take into account reliability, compatibility, and cost efficiency, particularly regarding the absence of sector- to block-level (e.g., 512B or 4KB) write atomicity within them. These require considerable efforts. In addition, the aforementioned Intel Optane memory was only one commercial NVM product shipped in a large capacity (up to terabytes) to build scalable PMem. Whereas, Optane memory is physically not byteaddressable. Its access unit for write and read operations is 256B [82], which deviates from the assumption of DAX that PMem shall be accessible at the same byte granularity as DRAM. Using DAX on Optane memory-based PMem causes unnecessary write and read amplifications, which in turn hinder the use of DAX on real-world NVM products.

Meanwhile, how to use PMem depends on the need of systems and applications. Mainstream file systems are built on the assumption of sector- to block-level atomic I/O that eases development, maintenance, and extension [14, 6, 44, 52]. Many applications, such as databases that have gained wide popularity for decades, explicitly or implicitly demand the support of block-level write atomicity, because a mix of up-to-date and stale data in one block leaves ambiguity [35, 6]. As a result, **configuring PMem in the form of block device is a promising and necessary alternative to use it**.

To facilitate the use of real PMem products, Intel provides two configuration modes. One is AppDirect mode and the other one is Memory mode [82, 70]. On one hand, the Memory mode uses PMem to expand main memory capacity without persistence [82]. When PMem is configured in the Memory mode, it presents capacious but volatile memory space, because the Memory mode employs DRAM as a cache to hide PMem's higher latency, with hardware-controlled caching policy between DRAM and PMem [48]. On the other hand, if programmers intend to explicitly utilize a separate, visible PMem space, they must choose the AppDirect mode [82, 48, 5]. In the AppDirect mode, programmers can create and use a namespace in different usage modes for different purposes [31, 30, 32]. For example, they can set up a namespace in the "fsdax" mode that exposes the raw PMem space without sector- or blocklevel write atomicity [30]. In practice, with such an exposed

PMem space, programmers can mount a file system with aforementioned DAX feature (e.g., Ext4-DAX), and they can also mount a mature file system (e.g., Ext4) in spite of no support for block-level write atomicity. Also, programmers can create a namespace for a "sector" mode and mount a mature file system with block-level write atomicity through the *block translation table*.

2.2. Block Translation Table

Researchers have developed software-based block devices on PMem, such as PMBD [10] and Block Translation Table (BTT) [65, 32]. Both of them are kernel-space device drivers that achieve block-level write atomicity in the software approach. As BTT has been actively maintained in Linux kernel since version 4.2, we focus on it in this paper. In fact, the aforementioned "sector" mode for practically creating a PMem space with block-level write atomicity is also known as "btt" mode or "safe" mode [31]. BTT is named after its main component, i.e., the block translation table recording the mapping from logical block number (*lba*) to physical block number (pba). However, it is much more than address translation. To upper-level file system, BTT provides standard bio interfaces and allows programmers to configure the block size (e.g., 4KB). Programmers thus can make and mount mature file systems such as Ext4 on BTT. Inside the device, it achieves block-level write atomicity by a hybrid scheme of CoW and logging for data and metadata, respectively. In addition, programmers cannot mount a file system with DAX feature on a PMem space created in the "btt" mode to bypass the OS's page cache [31]. However, they can still perform direct I/Os in the "btt" mode by using the O_DIRECT flag that has been used for conventional block devices such as HDDs and SSDs.

Figure 1 illustrates the layout of PMem space formatted with BTT and how BTT manages the PMem-based block device. BTT driver has no data or metadata kept in DRAM. It divides the PMem space into multiple *arenas*. Each arena can have a maximum capacity of 512GB. Two identical *Info blocks* are placed at the head and tail of each arena for redundant backup. BTT divides an arena's PMem space into *data blocks*, each of which is indexed with a *pba*. As mentioned, BTT utilizes a table to record address mapping from *lbas* carried in bio requests to *pbas* of PMem space.

BTT introduces the concept of *lanes* to support simultaneous writes to PMem. The number of lanes is configured as the number of CPU cores or 256, whichever is smaller. Among all data blocks per arena, BTT keeps a dynamic set of up to 256 *free blocks*. Each lane is associated to a free block. The *BTT Flog* of an arena is composed of multiple log entries, used to track changes of address mapping from *lbas* to *pbas* when BTT is serving write requests.

Each bio request is running on a CPU core with an ID for high concurrency. This core ID determines in which lane BTT serves that bio request. On a read request with a target *lba*, BTT driver locates the corresponding arena before going to the lane indicated by core ID. It then looks up the mapping table and gets the *pba*. Next, BTT loads



Figure 1: An Illustration of Block Translation Table (BTT)

data indexed by *pba* to complete the read request. On a write request, as shown in Figure 1, BTT still takes the corresponding lane (1). Instead of in-place updating, BTT writes data carried in the request into the lane's free block for out-ofplace updating (2), i.e., the CoW way. Next, BTT initiates redo logging to record the change of address mapping for the lba. BTT records lba, old pba, new pba of the lane's free block in the log entry for crash recoverability (3). After doing so, BTT modifies the mapping indexed by *lba* in the mapping table (4). The previously-mapped block indexed by old *pba* is factually swapped out to be free and BTT employs this free block to supplement the lane. As a result, each lane is always with an active free block at runtime. In summary, CoW and logging jointly ensure block-level write atomicity for BTT with the capability of rolling back on failed write for an *lba* if system crashes (e.g., power outage or kernel panic).

3. Motivation

We have conducted a study on a real PMem device configured in the AppDirect mode to observe the performance of BTT. More details of the platform are presented in Section 5. We consider three variants based on different usages of PMem. Firstly, we create a namespace in the "btt" mode and mount Ext4 on the PMem-based block device. Secondly, we mount mature Ext4 on bare-metal raw PMem created in the default "fsdax" mode. Thirdly, we mount Ext4-DAX on raw PMem reinitialized and created in the default "fsdax" mode. They are denoted as BTT, PMem, and DAX, respectively. As shown by the green arrows in Figure 1, both BTT and PMem utilize block I/O interfaces, while DAX bypasses the page cache of OS. Note that among them only BTT supports blocklevel write atomicity. PMem and DAX may leave data corrupted if a power outage occurs.

To quantitatively test three utilization ways, we use Fio [4] to generate a write-intensive workload that randomly writes a total volume of 64GB data in 4KB per I/O request under the direct I/O mode in order to exclude the impact of OS's page cache. Figure 2a compares their results measured in execution time. Although BTT provides the strongest crash consistency and write atomicity, it spends 37.4% and 16.6% more time than PMem and DAX, respectively. DAX is faster than BTT because Ext4-DAX, as tuned with the DAX feature, has been optimized in the software stack for PMem. BTT creates a block device on raw PMem that is just used by PMem. The 37.4% gap between PMem and BTT indicates that, although BTT has employed CPU's multi-cores to concurrently process I/O requests, the cost of maintaining block-level atomic writes by BTT driver. In short, **The performance with BTT suffers from such non-trivial cost of enforcing block-level write atomicity**.

When using fast PMem, we intend to make use of it for high performance like raw PMem or Ext4-DAX. In the meantime, we shall retain the block-level write atomicity for critical system and application softwares. In order to improve the performance with BTT, we analyze the source code of it. We find that an internal device cache, a key component that is widely built in modern HDDs and SSDs, is not contained in BTT. I/O staging cache is widely used to accelerate performance for storage [34, 74, 78, 73, 56]. For example, PMBD employs a DRAM cache to temporarily store dirty pages for I/O staging [10]. It has maintained a syncer daemon thread that flushes the buffer to PMem when the buffer is filled to an extent (watermark). Though, the developers of PMBD did not fully consider how to utilize the multi-cores of CPU to accelerate caching for high concurrency when in designing PMBD, because they still followed the classic caching strategy used for a conventional block device. In short, PMBD divides the DRAM cache into multiple sub-buffers. Once one sub-buffer is filled to the watermark (e.g., 70% or 100%), PMBD drains it by writing back buffered data to PMem. We have added multibuffers in an overall capacity of 512MB to BTT with PMBDlike caching (referred to as PMBD). However, in contrast to an expectation that I/O staging cache should have boosted the performance of BTT, the execution time spent by PMBD, as shown in Figure 2a, is even 6.3% longer than that of BTT. The flush of a sub-buffer upon preset watermark is likely to cause high I/O congestion on the critical path. Moreover, Ext4 periodically performs a journal commit for crash consistency and data durability every five seconds by default, through issuing a bio request with the REQ_PREFLUSH flag being set to flush the volatile internal cache of storage device [64, 66]. However, Ext4 does not synchronously wait for the completion of such a periodical flush with the SYNC flag unset per request [51, 76]. Whereas, a user request that encounters one such asynchronous flush still suffers from additional overhead. Besides PMBD, we introduce the conventional LRU strategy that evicts the least recently used (LRU) cached block when the DRAM cache is full. Yet LRU still periodically flushes all data that it caches per bio request with REQ_PREFLUSH set. As demonstrated by Figure 2a, the execution time of LRU is even a bit longer than that of PMBD.

To figure out why the performance with BTT is not improved with the employment of I/O staging cache, we have recorded the individual response time of serving every write



Figure 2: A Comparison on BTT, Ext4-DAX, PMem, and PMBD

request. Figure 2c, Figure 2d, and Figure 2e capture a part of these records in a contiguous window of 50,000 requests for BTT, PMBD, and LRU, respectively. The bottom margins in Figure 2d and Figure 2e are evidently narrower than that in Figure 2c, which indicates that writing data into cache truly helps to cause shorter response time ($<5\mu$ s) in serving a part of I/O requests. However, a comparison among three diagrams clearly shows that PMBD and LRU have completed numerous I/O requests with more than $20\mu s$, as depicted in Figure 2d and Figure 2e. Moreover, the continuous and frequent spikes exhibited in Figure 2d and Figure 2e are not commonly present in Figure 2c. At runtime, PMBD is likely to use up DRAM space under I/O pressure and flushes buffered data to make space for incoming requests. Similarly, LRU lays a 2-step write where LRU firstly evicts a buffered block to PMem and then writes arriving data to DRAM cache.

Furthermore, the periodical flushes that Ext4 launches also cause PMBD and LRU to asynchronously flush cached data. To observe the impact of all flushes, we have captured the runtime response time in a larger time window with one million (1,000,000) requests. Figure 3 captures each request's response time for two caching algorithms. Both of them exhibit almost uniformly distributed long tail latency over time. Particularly, in Figure 3a, we can distinguish for PMBD what tail latencies are caused by synchronous flushes (higher ones) due to a fulfilled cache or periodical asynchronous flushes (relatively lower ones) demanded by Ext4. **Those flushes occur on the critical path of serving I/O requests, thereby generating frequent latency spikes**.

In the meantime, applications explicitly call fsyncs to persistently store data and impose I/O ordering with the REQ_PREFLUSH, REQ_FUA, and SYNC flags set, like what LevelDB does with SSTable files for compactions [19]. On receiving an fsync, the device management firmware or software has to forcefully flush buffered data to persistent storage [76]. Since fsyncs are the other source of on-demand flushes, we have modified the foregoing test case by inserting one fsync after every 128 write requests with 512KB (128×4 KB). As shown by the right part of Figure 2a, fsyncs evidently increase the execution time for all five. Note that PMBD and LRU still take 6.1% and 10.8% more time than BTT, respectively.

In order to further explore how I/O staging cache influences the fsync performance, we invoke an fsync in different frequencies, i.e, after randomly writing 512KB to 128MB data. Figure 2b shows two close curves for PMBD and LRU, respectively. The fsync time of them rises sharply as more data is written between two consecutive fsyncs. More written data means more data is buffered in the cache. As fsync drains the cache, the impact of flushing buffered data becomes more and more detrimental. In practical, the cache is likely to be fully filled over time since it continuously and uniformly receives data. On-demand flushes are hence inevitable and severely impair performance of I/O staging cache.

To sum up, the performance with BTT suffers from the cost of preserving block-level write atomicity, while the fashion of I/O staging cache with frequent on-demand flushes cannot improve its performance. The exploitation of CPU's multi-cores is also promising for caching data in a PMem-based block device which is directly handled by the CPU. We thus reconsider how to manage an effectual cache for BTT-like block device made of fast PMem. The cache management must avoid flushing data on the critical path of serving I/O requests. It shall be simple, without introducing dramatic performance penalty. It may apply non-uniform caching policy under different conditions. For example, in case of a fully filled cache that is undergoing I/O congestion, it should handle arriving requests in a different but efficient way. Notably, BTT *emulates* a block device sitting on the memory bus. Both DRAM and PMem are operated by multicore CPU. If we exploit scores of CPU cores to place data into cache in the foreground and swiftly transit data to PMem in the background, we may effectively boost the



Figure 3: The response time for PMBD and LRU in a window of one million requests

performance with BTT. This summarizes the essence of Caiti.

4. Design of Caiti

4.1. Overview

Caiti manages a DRAM cache that significantly differs from conventional I/O staging caches. Firstly, it decouples cache space management from address mapping. By logically organizing cache space into sets and hashing *lbas* to localize a set, Caiti places a block freely in the cache and finds a block swiftly (Section 4.2). Secondly, in order to avoid stalls that prevent I/O requests from proceeding, Caiti has two writing policies to enable I/O transit caching (Section 4.3). On one hand, Caiti immediately initiates an eager eviction to write back a block into PMem-based block device once the block is put down into a cache slot. On the other hand, if no free space is available in the DRAM cache, Caiti conditionally bypasses the full cache but directly writes a block to PMem, so as to circumvent I/O congestion. Caiti bases both policies on CPU's multi-cores, which in turn enables scalability and concurrency for it. Thirdly, without losing block-level write atomicity, Caiti comprehensively supports all standard bio flags and facilitates order preserving needed by applications and system softwares (Section 4.4).

In the storage stack, Caiti's cache is positioned under the OS's page cache. They complement each other and share the functionality of buffering data for application and system softwares. The OS's page cache helps to absorb frequently accessed data for Caiti. In turn Caiti's cache helps to further reduce the time cost on the critical path of serving write requests. In addition, Caiti is software-only solution and incurs no change to hardware like memory controllers for PMem or DRAM. It has been practicable on a platform with real DRAM and PMem devices.

4.2. Caiti's Cache Space Management

Organization. Caiti demands and reserves a contiguous DRAM space, e.g., in 512MB, from the OS. It freely places blocks across this cache space and partitions the space into cache *slots*. For example, the cache illustrated in Figure 4 has six slots. All slots share a configurable uniform size, which by default is equivalent to the block size of BTT. To track slots and handle I/O requests, Caiti logically maintains

a number of cache *sets* holding valid data indexed by the hash value of *lba*. Caiti organizes each cache set as a write-back queue (WBQ). As shown in Figure 4, Caiti also manages a global *free set* to group unoccupied slots over time.

Slot header. As illustrated by Figure 4, Caiti tracks and manipulates a slot with slot header, in which there are a slot number, an *lba*, a slot state, a WBQ pointer, and a lock. The lock is used for supporting concurrent accesses between multiple threads. Slot number is the identity of a slot and points to where a block is stored in the cache. The *lba* indicates where a block is stored in the PMem-based block device. Although the mapping from an *lba* to a slot number is fundamental for caching algorithms, it only holds transiently for Caiti, because Caiti employs the I/O transit strategy that swiftly sends a cached block down to device (see Section 4.3). Any slot staying in the free set is assigned an outlier lba, e.g., -1. Given a normal lba, Caiti hashes the *lba* to get the corresponding cache set number. The hash function exemplified in Figure 4 is a modulo operation with four. In practical, we do a similar modulo operation with *lba* over the number of cache sets. A collision may happen if two *lbas* are hashed to the same set. However, we allocate multiple cache slots per set and, more important, the eager eviction to be presented later helps to swiftly vacate cache slots. These jointly alleviate the impact of hash collisions.

The WBQ pointer in each slot header is used to link slots in the cache set's WBQ as well as ones in the free set. As to the slot state, Caiti defines four legal ones, i.e., *Free*, *Pending*, *Valid*, and *Evicting*. Initially all slots are placed on standby in the free set with the *Free* states. At the *Pending* state, data is being written into the slot. The *Valid* state means data has been put in the slot while *Evicting* means Caiti is doing write-back to PMem-based block device. After a successful write-back (eviction), Caiti recycles the slot to be *Free* and puts it back to the free set for use in the near future.

4.3. Caiti's Writing and Reading Policies *4.3.1.* Caiti's Writing Policies

The way Caiti handles a write request significantly differs from ordinary storage device caches. Conventionally, since SRAM or DRAM embraces much shorter write latency than HDD and SSD, researchers employ a cache to keep blocks buffered, expecting a sufficiently long period of I/O



Figure 4: Main Components and Write Procedure of Caiti

staging to hide the slow access speed of underlying storage media. When the device cache is full or an fsync is received, the device management flushes partial or all buffered blocks for orderly persistence, which incurs drastically long response time on the critical path. Even though BTT builds a block device on top of relatively fast PMem, it demands nontrivial cost to enforce block-level write atomicity (see Figure 2a and Section 3), so waiting for data to be flushed and persisted is still costly for PMem-based block device. As a result, we attempt to avoid the DRAM cache from being fully packed and minimize I/O stalls at runtime. This intention leads to two main policies Caiti utilizes in handling a write request.

Eager eviction is a policy that Caiti employs for prompt write-backs. In contrast to buffering a block for long-term I/O staging, Caiti immediately launches a write-back and delivers received blocks to underlying PMem-based block device through a background thread. By eager eviction, Caiti aims to free up cache slots in a timely fashion, so as not to engage incoming write requests in waiting for the completion of flushing buffered blocks due to no available cache space. Concretely, Caiti manages to write arriving data to a free slot with one DRAM write on the critical path, regardless of cache hit or miss, which, interestingly, renders the same effect of a write hit in I/O staging cache. The latter even occupies cache space and affects serving other write requests. In addition, when receiving an explicit fsync call that shall trigger the flush of device cache, Caiti is able to quickly complete the fsync since continuous eager evictions have drained most of the buffered data to PMem.

Conditional bypass is the other writing policy of Caiti. Under high I/O pressure, it is possible that no free slot is available in the DRAM cache. Waiting for a vacant slot by evicting a buffered block to make room would stall the current I/O request with a long response latency. Whereas, Caiti bypasses the full cache and directly stores arriving block to the PMem-based block device. The reason why Caiti does so is twofold. Firstly, a fully occupied cache means a tense I/O congestion is ongoing. Waiting for a slot to be vacated by flushing the slot's buffered data down further worsens the congestion. Secondly, vacating a slot to make space implies that the current write request cannot be finished until Caiti Algorithm 1 Caiti's Write Procedure (caiti_write(*lba*, *d*))

Require: A write request that carries <i>lba</i> and data <i>d</i>					
1: $\eta := \text{find_set_by_hash}(lba); \triangleright \text{Hash to find cache s}$	et				
2: $\tilde{q} := \text{get}_WBQ(\eta)$; \triangleright Get the WBQ pointer \tilde{q} of set	η				
3: $sh := \text{dequeue}(\tilde{q}, lba); \triangleright \text{ Get slot header } sh \text{ by dequeue}$	ıe				
4: if $(sh \neq \text{NULL})$ then \triangleright Cache h	it				
5: \triangleright We make Caiti first test and set the hit slot's sta	te				
6: test_and_set_state(<i>sh</i> , <i>Pending</i>); \triangleright state \rightarrow <i>Pendin</i>	ıg				
7: \triangleright Caiti writes data to the hit slot and set <i>Valid</i> sta	te				
8: $ret := write_slot(\eta, sh, lba, d, Valid); \triangleright state -$	\rightarrow				
Valid					
9: $ret := enqueue(sh, \eta, \tilde{q}); \triangleright For eager eviction$	n				
10: else \triangleright Cache miss, <i>lba</i> is not matched in any sle	ot				
11: if (is_cache_full() == False) then \triangleright Cache is not	ot				
full					
12: \triangleright Let Caiti allocate a free slot <i>sh</i> from the free	e				
set					
13: $sh := \text{allocate_slot_from_free_set}(\eta);$					
14: set_state(<i>sh</i> , <i>Pending</i>); \triangleright state \rightarrow <i>Pendin</i>	g				
15: \triangleright Caiti writes data to slot <i>sh</i> and sets <i>Valid</i> sta	te				
$ret := write_slot(\eta, sh, lba, d, Valid);$					
17: \triangleright Let Caiti get the WBQ of set η for lat	er				
eviction					
18: $\tilde{q} := \operatorname{get}_{WBQ}(\eta);$					
19: $ret := enqueue(sh, \eta, \tilde{q}); \triangleright For eager eviction$	n				
20: else \triangleright Cache is currently fu	III				
21: $ret := btt_write(lba, d); \triangleright Conditional bypa.$	55				
22: return ret; \triangleright SUCCESS or -EIC	J				
24: end II	1				
25: If $(ret \neq raise)$ then \triangleright writing cache slot not range	20				
20: ret := notify_eager_eviction(sn , η); \triangleright Edge	zr				
Eviction 27. and if					
$27: \text{ chu ii} \qquad \qquad$	\cap				
$20.$ return ret, \square	J.				

performs one PMem write and one DRAM write for putting down the arriving data into the vacated slot. The conjunction of these operations, however, costs even more time than a direct single write to PMem.

Write procedure. Eager eviction and conditional bypass jointly make Caiti function as I/O transit rather than staging. Algorithm 1 shows main steps of Caiti's write procedure. On receiving a bio request with a target *lba* and data, it firstly hashes the *lba* to decide the corresponding cache set (Line 1). Hashing helps to balance I/O loads among the entire cache space as well as quickly identify a buffered block. The hash function may incur conflicts with different *lbas*. As mentioned, Caiti links the slot headers of conflicting lbas in a cache set's WBQ. Over time, the eager eviction policy of Caiti ensures that not many slots stay in one set. Caiti gains access to the set's WBQ (Line 2). It performs a scan to dequeue and grab the slot header for the target *lba* (Line 3). It expects a cache hit with a non-null slot header (Line 4). If so, Caiti tests the slot's state (Line 6). Given a Valid state, it transitions the slot to the Pending state in order to prevent the slot from being written back as the eager eviction is continually working. In case of the Evicting state, which means BTT is writing down the block to PMem, Caiti waits for the completion from BTT so as not to interrupt a persist operation. This makes Caiti retain block-level write atomicity. If Caiti finds the slot in the Pending state, which means a previous cache write is still ongoing, it waits for the completion of that DRAM write. Next Caiti updates data in the slot and sets the slot's state as Valid (Line 8). It then puts the slot to the set's WBQ for eager eviction (Line 9).

If a cache miss happens to target *lba* while the cache is not full (Line 11), Caiti allocates a free slot and places the slot into proper cache set η (Line 13), with the slot's state set as *Pending* (Line 14). Next, Caiti puts down data into allocated slot for caching and fills metadata in the slot's header, such as the *lba* and *Valid* state (Line 16). Note that the state transition from *Pending* to *Valid* is essential and useful for Caiti. During the duration of writing a block of data to DRAM, a read request that hits at the cache slot with matched *lba* would not see incomplete data because of the *Pending* state. Also, Caiti would not do eager eviction on a slot unless the slot has accommodated an integral block with a *Valid* state stationed. Then, Caiti fetches the entry pointer of WBQ (Line 18) and prepares for eager eviction by enqueuing the slot (Line 19).

In case of a full cache, Caiti directly persists the data to PMem-based block device and returns a success or not (Lines 20 to 22). This corresponds to conditional bypass. On the other hand, every time a slot is placed in the set's WBQ, Caiti instantly notifies the submodule of eager eviction (Line 26). Upon evicting a slot, Caiti changes the slot's state from *Valid* to *Evicting*. A background thread now initiates the write-back. In the meantime, Caiti returns a signal of success or fail for serving the write request (Line 28).

Figure 4 briefly illustrates three cases for a write hit (**1**), a write miss with a slot allocation (**234**), and a concurrent write-back in the background for eager eviction (**5**). As mentioned, the hash function used in Figure 4 is the modulo operation with four. On a write request with *lba* #0060 that is executing on CPU core 0 (**1**), Caiti hashes the *lba* and gets set 0. Then, through core 0, Caiti finds that a cache slot with

the same *lba* already stays in set 0's WBQ, thereby indicating a write hit. Thus, Caiti writes data into the cache slot directly. After finishing the write, Caiti transitions the slot's state to Valid. As to the next write request with lba #2022 running on CPU core 2, after hashing, Caiti does not find any slot holding #2022 among cache set 2. Therefore, Caiti needs to handle this request as a write miss (2). Since the cache is not full, it allocates a free slot from the free set (3). In the end, Caiti uses core 2 to set the slot as Valid as well and add it into set 2's WBQ for eager eviction (4). In Figure 4, Caiti is leveraging a background thread on core 3 to conduct eager eviction (**5**). Caiti finds a slot in the *Valid* state staying in set 3's WBQ. It transitions the slot state to *Evicting* and starts writing back the data of Iba #9527. In these procedures, readwrite locks of involved slots shall be acquired and released carefully in order to rule out any disorders on writing and reading data for upper-level software layers.

Scalable and concurrent I/Os. The efficiency of Caiti is supported by CPU's multi-cores. An ongoing bio request naturally executes on a CPU core that is responsible for proceeding the write to a cache slot or in-PMem location. BTT functions as the backend of Caiti and embraces multiple lanes that are ready to take and handle multiple I/O streams simultaneously [65]. As shown in Figure 4, Caiti maintains a thread pool. At runtime, it grabs a background thread from the pool and runs the thread on an idle CPU core. The thread checks WBQs and writes buffered blocks within PMembased block device. In addition, Caiti leverages the lock in each slot header to coordinate race conditions and resolve write-write or write-read conflicts between concurrent bio requests on one *lba*. Additionally, the use of a global free set neither posits a bottleneck for multi-threads nor impairs the scalability of Caiti. On one hand, when allocating or deallocating a free cache slot, Caiti deals with the lock in the slot header through efficient *compare-and-swap* (CAS) operations. On the other hand, a global free set, rather than distributed local ones that are respectively dedicated to multiple CPU cores, is more friendly to workloads that are with write or read skewness. Assume that a workload keeps writing data via few particular cores with multiple threads. Local free sets, if employed, are likely to be used up for those working cores while the free sets of other cores stay idle. Caiti's globally shared free set has no such inefficiency in utilizing cache slots. In summary, these above-mentioned components concretely guarantee high scalability and concurrency for the collaboration between Caiti and BTT.

4.3.2. Caiti's Reading Policy

When Caiti receives a read request with an *lba*, Caiti firstly checks whether the *lba* exists in a cache slot with *Valid* or *Evicting* state. This ensures that applications and file systems always perceive the latest valid and complete data. If a cache hit occurs, Caiti returns the buffered data copy from the matched slot. Otherwise, it redirects the read request to underlying PMem-based block device. Additionally, Caiti

prioritizes write I/Os and a read miss in the DRAM cache does not entail loading a block for buffering.

4.4. Important Aspects of Caiti

Cache mapping and replacement. As mentioned, address mapping from storage space to cache space is basic for cache management. Employing a table for address mapping was widely utilized in block devices [49]. Comparatively, rather than keeping a mapping table, Caiti performs the mapping from an *lba* to a cache set by hash calculation. After hashing a target *lba*, Caiti avoids checking lots of slot headers in the calculated cache set because of continual eager evictions in the background, so it gains promising efficiency in satisfying on-demand I/O requests. The conditional bypass also helps to avoid cache replacement for Caiti. Waiving the mapping table and replacement structurally distinguishes Caiti from conventional caching designs.

bio and ordering. Caiti supports all bio flags. The aforementioned REQ_PREFLUSH is used to flush the entire cache. Applications such as databases and mail service frequently call fsyncs to orderly flush file data to persistent storage. An fsync is eventually translated to a bio request with two flags, namely REQ_PREFLUSH and REQ_FUA, turned on. REQ_FUA ensures that storage device signals I/O completion only after the data has been persistently committed [64]. To support these flags for order preservation, Caiti flushes all WBQ entries on receiving both REQ_PREFLUSH and REQ_FUA and waits for I/O completion signals from underlying PMem-based block device. Then Caiti continues to serve subsequent bio requests. The alignment of bio standard secures compatibility and viability for Caiti. As to recent research works that proposed new methods to restructure emerging storage devices by adjusting store ordering [76, 8, 44], Caiti can be easily adapted to suit their flags, primitives, or commands.

There is another ordering case in which an application keeps writing the same *lba* for consecutive updating. Given a high update frequency, the OS's page cache is very likely to filter and absorb repeated write I/Os. As to a low frequency, Caiti handles them as ordinary requests since previous data versions should have been evicted to PMem. Neither scenario affects Caiti's work, which in turn addresses its robustness.

A volatile internal cache has been widely employed for decades in HDDs and SSDs. The aforementioned periodical flush that Ext4 issues is used to ensure that, once a power fail or kernel panic happens, at most modified data put in the past five seconds would be lost in the volatile device cache [66]. If an application such as a database demands a higher guarantee of crash consistency and data durability, it shall call fsync or fdatasync, with the REQ_PREFLUSH set alongside write bio requests to forcefully flush cached data [76, 64].

Positioning. Caiti is positioned in the latest developments of software and hardware for storage stack. For example, current Linux kernel has included multi-queues (blk-mq [67] for modern storage devices that BTT emulates with high parallelism. Caiti's WBQs align with blk-mq. Also, manufacturers are shipping emerging block devices (e.g.,

Samsung CXL SSDs [54]) with large DRAM cache. The ideas of Caiti are applicable to enhance them.

Wear leveling. Write endurance is crucial for NVM [9, 46, 39, 26, 82, 24, 53, 25]. Previous works have proposed wear leveling algorithms at both hardware and software levels. For example, researchers have reverse-engineered the hardware wear leveling algorithm that Intel deploys for Optane memory [75, 47]. On the other hand, softwarebased wear leveling is promising to prolong the lifetime of PMem with wear-aware space allocation and recycling through the software stack [23, 1, 42, 18, 21]. During its procedures of eager eviction and conditional bypass, Caiti can take into account the write endurance issue of NVM and collaborate with a software-based wear leveling scheme customized with the consideration of I/O transit caching. By doing so, Caiti is supposed to simultaneously gain both high performance and enhanced lifetime. We leave this as one of our works for exploration in the near future.

5. Evaluation

Setup. All experiments have been conducted on a machine with real PMem products, i.e., Intel Optane DC memory in 768GB configured in the *AppDirect* mode. There are also 384GB DRAM installed in the machine. The CPU is 36-core Intel Xeon Gold 6240. The OS is Ubuntu 20.04.5 with Linux kernel 5.1. The compiler is GCC/G++ 8.4.0. The vanilla BTT has 2,001 lines of code (LOC). We add or change 1,066 LOC to implement the core functions and structures of Caiti.

We compare BTT with Caiti to Ext4-DAX, Ext4 mounted on raw PMem, original BTT, and NOVA in CoW mode [80]. They represent common ways of using PMem and are referred to as DAX, PMem, BTT, and NOVA, respectively. Note that DAX and PMem do not provide block-level write atomicty. As to caching with the I/O staging strategy, we consider algorithms including aforementioned PMBD and LRU, as well as the state-of-the-art Co-Active [61]. For PMBD-like caching, we have implemented two variants. One is denoted as PMBD. It flushes the entire cache if the cache is 100% full (see Section 3). In other words, when there is no free cache slot available upon the arrival of a write request, PMBD does a cache flush. The other one, referred to as PMBD-70, strictly follows the literature and source code of PMBD. PMBD-70 flushes the buffer when the buffer is 70% full. PMBD-70 further differs from PMBD in that the former flushes the cache through a syncer daemon thread [10]. Unlike PMBD and PMBD-70, LRU evicts the LRU slot rather than an entire cache if cache is full. Co-Active is a collaborative active write-back cache management approach. We port it from NVMe SSD to PMem-based block device managed by BTT. It employs a cold/hot separation module to distinguish cold and hot data via a Bloom Filter. Meanwhile, it maintains two linked lists for dirty and clean blocks in DRAM cache, respectively. When PMem stays idle, Co-Active proactively evicts data from the dirty list to PMem. We refer to it as COA. We configure the block size as 4KB for BTT. Caiti, PMBD, PMBD-70,

I/O Transit Caching for PMem-based Block Device



Figure 5: A Comparison with Fio on Average/Runtime Response Time, Tail Latency, and Multi-threads

LRU and COA are four designs that manage a DRAM cache on top of PMem-based block device and for each of them, we set the cache slot size and default cache capacity as 4KB and 512MB, respectively. We conduct experiments with micro-benchmarks (e.g., Fio) and real-world applications (e.g., LevelDB [19] and QEMU [68]) to thoroughly evaluate Caiti. The main metrics

to measure performance is response (execution) time and throughput (bandwidth in MB/s).

5.1. Micro-benchmark (Fio)

Fio is a powerful and capacious benchmark [4]. The reason why we choose it for evaluation is multifold. Firstly, it can continuously generate I/O requests with varying I/O depth to launch different pressure tests. Secondly, it reports meaningful results that help to observe and interpret the performance of Caiti. Thirdly, it allows us to configure multiple threads, enabling us to test the scalability of Caiti. Fourthly, it further supports a configurable composition of I/O engine, access pattern, I/O depth, and multiple jobs (threads) to comprehensively evaluate caching algorithms under comparison.

With Fio, we firstly test Caiti with substantial I/O pressure. As we have done in the motivational study, we configure Fio in the direct I/O mode, with libaio engine continuously issuing random writes for 30 minutes in a 64GB file with 4KB I/O size. We vary the I/O depth from 32 to 1024 so that the storage stack undergoes increasingly more orderless I/O requests in flight [44]. By doing so, we aim to impose stressful workloads and thus deeply explore the capability of Caiti with ample I/O pressure. As shown in Figure 5a, Caiti consistently outperforms other designs by taking much less response time. For example, with 128 I/O depth, it boosts the performance with BTT by up to 3.6x, while PMBD, PMBD-70, LRU, COA and NOVA that enforce block-level write atomicity incur 3.6×, 3.6×, 3.5×, 2.9× and 4.1× execution time compared to Caiti, respectively. The superior performance of Caiti justifies the efficacy of its strategy of I/O transit caching with timeliness and concurrency. On massive write requests arriving, Caiti handles them with multiple cache sets supported by multi-cores in the foreground and leverages a pool of background threads and WBQs to promptly move buffered data to PMem. As a result, I/O stall and congestion hardly occur to Caiti. Meanwhile, BTT and NOVA directly write data to PMem with efforts for achieving block-level write atomicity. PMBD's cache, albeit employing multi-buffers, is likely to be overfilled at runtime and stall for the drain of buffers. Although PMBD-70 flushes the cache when the cache space is used as much as 70%, it cannot timely handle burst requests with the syncer daemon thread, resulting in limited improvement. LRU also frequently stalls due to its 2-step write upon no free cache space (see Section 3) and yields similar performance compared to PMBD. Compared to other I/O staging algorithms, COA reduces response time by taking advantage of cold/hot data separation. However, COA is still inferior to Caiti. Upon the arrival of continuous I/O requests, PMem is unlikely to be idle. COA has to evict buffered data and stall incoming requests. As a result, many cache replacements occur on the critical path for COA. This is a common issue shared by caching algorithms with the I/O staging strategy. Separating cold and hot data also incurs cost for it. The performance gap between these I/O staging algorithms and Caiti is concretely significant as regards I/O stall and congestion the formers are encountering.

DAX and PMem spend 115.7% and 120.2% more time than Caiti, respectively. They directly persist data to PMem, while Caiti puts data to DRAM cache. Caiti's hashing on *lbas* distributes data in multiple sets to leverage multi-cores and eases its eager eviction that concurrently vacates filled cache slots. Conditional bypasses also occasionally alleviate I/O congestion at DRAM cache. In addition, Yang et al. [82] obtained a similar observation as they found that writing data to PMem has lower throughput than doing so with DRAM.

Secondly, we have captured the runtime response time per request for Caiti. We present 50,000 points of response time in Figure 5b and a comparison with Figure 2c to Figure 2e conveys that Caiti spends much shorter time in serving each request than BTT, PMBD, and LRU. At runtime, the number of flushes does not change significantly for Caiti, since Ext4 consistently issues a bio request with the REQ_PREFLUSH set for flush every five seconds. Though, the eager eviction has moved data aggressively to PMem and results in much more lightweight flushes for Caiti compared to PMBD and LRU. This explains why the majority of runtime latencies in Figure 5b for Caiti is below 20µs. In fact, we also record the runtime response time in the large window of one million requests and shown their points in Figure 5c. A comparison between Figure 5c and Figure 3 indicates that in a long run, Caiti yields much shorter runtime response time at the presence of periodical flushes since eager eviction and conditional bypass help to reduce the volume of cached data that is supposed to be flushed every five seconds. Though, foreground and background threads of Caiti may contend for the same cache slot and incur lock/unlock operations. Meanwhile, many processes are running simultaneously and numerous events (e.g., I/O interrupts or exceptions) are concurrently happening. When serving an I/O request, the CPU core on which Caiti is running might be scheduled to run for the other process or handle a sudden event. These two factors are likely to bring latency spikes to Caiti over time. To grab a more quantitative and comprehensive view on the critical path of serving write requests, we have recorded the 99.99 percentile (99.99P) tail latency that Fio reports after finishing all requests. The tail latencies shown in Figure 5d with Y axis being in the logarithmic scale complement our observation with Figure 5a. Caiti generally achieves shorter tail latency by I/O transit caching. With greater I/O depth, all nine algorithms experience higher pressure and the tail latency dramatically increases. Caiti is likely to trigger more conditional bypasses to avoid I/O congestion while, for instance, PMBD at 1024 I/O depth makes 14.0× 99.99P latency than Caiti because PMBD insists on waiting for the flush of buffered data.

Thirdly, we test the scalability of Caiti by varying the number of Fio's jobs (threads) under 32 I/O depth. In Figure 5e, with more jobs (1 to 32) issuing write requests, Caiti is more performant than others. It works with multiple sets by dynamically utilizing multi-cores over its foreground and background threads. Caiti thus yields high performance in a scalable and balanced way with multi-threading workloads. However, COA shows a clear distinction between a single

Table 1			
The Impact	of Cache S	ize (Avg. Resp.	Time in μ s)

Capacity	64MB	128MB	256MB	512MB	1GB	2GB
PMBD	16.81	16.50	16.67	16.39	16.53	16.23
PMBD-70	16.84	16.75	16.58	16.47	16.47	16.29
LRU	16.72	16.95	16.42	16.85	16.55	16.53
COA	12.89	12.81	12.85	12.95	12.55	12.30
Caiti	4.44	4.42	4.29	4.41	4.44	4.37

job and multiple jobs. This is because COA has difficulty in handling concurrent and continuous I/O requests from multiple threads, since its cold/hot data separation approach cannot efficiently identify overwhelming data.

Fourthly, we observe the impact of cache capacity on Caiti, PMBD, PMBD-70, LRU and COA with one job and 32 I/O depth. We set six cache sizes which, as shown in Table 1, hardly affect performance for all. This phenomenon is mainly due to the large volume of data received at the cache. For example, in the starting first minute, no less than 15GB data has been written. With massive data continuously arriving, a cache in tens of or even more gigabytes stays overloaded. This again explains why I/O transit caching is effective and useful for PMem.

Fifthly, we measure the spatial cost of metadata Caiti and other algorithms need for caching. For every 4KB block (cache slot), the spatial cost for each caching algorithm is as follows.

- Caiti with 102B in all: 8B for *lba*, 4B for slot_number, 1B for state, 40B for lock, 33B for work_struct, and 16B for two pointers (WBQ and free list).
- PMBD, PMBD-70 and LRU with 84B in all: 8B for *lba*, 4B for slot_number, 40B for lock, and 32B for lists.
- COA with 102B in all: 8B for *lba*, 4B for slot_number, 40B for lock, 48B for lists, and 2B for Bloom Filter.

For Caiti, 512MB cache costs about 12.75MB and a 4KB slot demands 102B on average. The ratio of 2.5% $(\frac{102}{4096})$ indicates high space efficiency for Caiti.

5.2. The Breakdown for Caiti

In order to comprehensively investigate what factors contribute to the performance of Caiti and other caching algorithms, we have conducted a breakdown test. To obtain relatively pure time trajectory, we choose the most fundamental POSIX I/O calls for the test. We overhaul the critical paths of calling POSIX write and read (pwrite and pread), respectively, when each caching algorithm serves them. The main steps of Caiti include 'cache metadata management', 'cache eviction and write' (a stalled write), 'conditional bypass', 'cache write only' (due to cache hit or a vacant slot), 'WBQ enqueue', 'cache flush' (due to Ext4 committing its journal every five seconds [66] or fsync [45]), and others (e.g., software overhead across kernel- and user-spaces). Moreover, in order to measure the respective impact of

eager eviction or conditional bypass, we also include two variants of Caiti with either eager eviction or conditional bypass disabled. They are denoted as 'w/o EE' and 'w/o BP', respectively. In our test program, we continuously send write requests to a file stored in PMem. The file is opened in the O_DIRECT mode to rule out the impact of OS's page cache. Each write operation is performed at a granularity of 4KB. Overall 1024×1024 write operations are done. The target offset in the file for each write is randomly generated in advance by following a uniform distribution. The contributions of aforementioned factors are captured in Figure 6a. With these results, we have obtained following key insights.

Firstly, the occurrence of 'cache eviction and write' is extremely rare for Caiti, which means that almost no stall has been detected for Caiti. This is because the eager eviction timely makes space through multiple background threads and an arriving write request does not need to wait for a free cache slot. As shown by Figure 6a, PMBD, PMBD-70, LRU, and COA all suffer from stalls caused by insufficient cache space and spend 40.5%, 25.3%, 40.0%, and 32.7% on 'cache eviction and write', respectively. Comparatively, the time taken for Caiti's all writes on the critical path of pwrite, including both 'cache write only' and 'cache eviction and write', is just 11.5% of all the time cost.

Secondly, without either eager eviction or conditional bypass, the time cost spent on corresponding operations dramatically rises up and impacts the overall performance of Caiti. For example, as shown in Figure 6a, with the bar of 'w/o EE', the absence of eager eviction accumulates data in the DRAM cache and conditional bypasses are triggered more frequently, thereby incurring more percentage of conditional bypasses. In the meanwhile, without conditional bypass, all data blocks must be written into cache slots, which helps to increase cache hits (i.e., 'Cache Write Only'). Though, the pwrite workload studied here is quite simple with regard to simply writing 1024×1024 data blocks and I/O stalls hardly happen. Consequently, with both eager eviction and conditional bypass enabled, Caiti achieves close performance to the variant with eager eviction only.

Thirdly, as to the critical path for pread, each caching algorithm achieves a similar breakdown (see Figure 6b). We note that our test program writes and reads data across a 4GB space, which is eight times the cache capacity of 512MB. Meanwhile, in the Linux kernel, Ext4 commits its journal every five seconds (periodical write-back) and in turn enforces a cache flush with REQ_PREFLUSH flag [64]. Consequently, the limited cache capacity and periodical cache flushes jointly entail cache misses. Yet the penalty of loading target data from PMem is equivalent among all caching algorithms.

Fourthly, Figure 6c illustrates the percentages of 'writing cache only' due to cache hits or available free slots, 'cache eviction and write', as well as 'conditional bypass' over all writes. It is evident that Caiti almost handles all write requests by writing data to cache slots. This again justifies why Caiti yields superior write performance. The syncer daemon thread helps PMBD-70 to vacate cache slots and in turn exhibits different percentages on 'cache eviction and



(c) The Percentages of Cache and Non-cache Writes

(d) Average Time under Different Write Conditions

Figure 6: A Detailed Breakdown of Different Cache Management Categories

write' compared with PMBD. In addition, COA writes more data directly to cache than LRU and PMBD. However, the more complex management of COA demands substantially more processing time, especially upon a REQ_PREFLUSH flag due to Ext4's periodical write-back every five seconds. As shown in Figure 6a, COA spends $1.9\times$ time compared to that of LRU and PMBD on 'cache flush'. Comparatively, Caiti eagerly evicts cached data to underlying PMem-based block device through multiple concurrent threads in the background. As mentioned, upon a periodical flush every five seconds, Caiti is very likely to write down only a handful of data to PMem. This in turn justifies why the time Caiti uses to flush the cache is almost negligible as shown in Figure 6a.

Fifthly, Figure 6d displays the average time needed for 'writing cache only', 'cache eviction and write', and 'conditional bypass'. For PMBD-70, the time for 'cache eviction and write' increases due to contention between the daemon thread and the working thread (e.g., for the list lock). Conversely, Caiti's 'conditional bypass' costs shorter time than other algorithms' 'cache eviction and write', highlighting the efficiency of conditional bypass when the cache is under congestion.

Sixthly, as told by the bars in Figure 6a, the remaining time ('others') accounts for a large proportion in the overall

time cost. To gain a deep insight, we capture the detailed timeline of 100 consecutive write requests and present a contiguous period in Figure 7 for both Caiti and BTT. In the upper half, we mark the timestamps when the test program has initiated the write requests, when these requests reach Caiti, when Caiti places data into a WBQ, when Caiti completes processing each write request, and when Caiti's background thread finishes evicting the cached block. The lower half of Figure 7 is for original BTT. The latency between a request's issuance and the request's arrival at Caiti is significant, taking 54.0% of the time cost per write request on average. This latency is mainly caused by the software penalty due to handling a system call between userand kernel-spaces [59, 84, 7, 40, 37]. The latency, however, provides a sufficient time window for Caiti's background threads to eagerly evict cached data to PMem. Consequently, given that a program consecutively writes data to the same file location twice without the involvement of OS's page cache, it is impossible that Caiti's eager eviction stalls the second write request from being served since the cache slot would be timely vacated when the request is still traversing the software stack. In the buffered I/O mode, the OS's page cache would accumulate two consecutive file writes before they are sent to underlying device driver.

I/O Transit Caching for PMem-based Block Device



Figure 7: Runtime Time Fragments of pwrite for Caiti and BTT

Last but not the least, the metadata management accounts for a very small proportion in the overall time cost for Caiti. As shown by the bottom segment in Figure 6a, it only takes 2.9% for Caiti. This low percentage addresses the cost efficiency of Caiti's management tactics. Take slot allocation/deallocation for example. Allocating or deallocating a free slot with the free set is done through efficient CAS operation. Furthermore, as Caiti makes the free set globally shared, each core is able to quickly fetch a vacant slot and avoids being blocked upon serving highly-skewed workloads.

5.3. Real-world Application

5.3.1. LevelDB

We illustrate with LevelDB for several reasons. Firstly, LevelDB is a key-value store gaining wide popularity. Secondly, LevelDB frequently calls fsyncs to persistently store data through storage stack and causes on-demand flushes not covered in Fio tests. Thirdly, Fio performs random I/Os in block size while LevelDB generates bulky I/Os batched in megabytes with SSTable files. Fourthly, we run Fio under direct I/O mode but with LevelDB, OS's page cache and LevelDB's application-level caches are involved. Fifthly, db_bench built in LevelDB has various write- and readintensive workloads. While we focus on optimizing write performance, we fully test Caiti's capability in handling read requests with LevelDB.

Figure 8 shows the average response time of all designs on serving ten million requests issued by db_bench on fillrandom, overwrite, readrandom, and readhot with value size varied from 128B to 4KB. Take 2KB values for example. Caiti spends 40.6% and 38.2% less time than BTT on fillrandom and overwrite, respectively. Still with 2KB values and fillrandom, Caiti takes 41.9%, 46.2%, 20.0%, 17.8%, 19.6% and 48.8% less time than DAX, NOVA, PMBD, PMBD-70, LRU and COA, respectively. LevelDB writes an SSTable in 4MB or 2MB with an fsync followed. Once receiving such a bulk of data, Caiti absorbs with cache slots and background threads launch eager evictions. Ideally, when a foreground thread finishes writing a slot, a background thread promptly initiates write-back. When fsync drains the cache, Caiti leaves only a handful of slots being in Valid states. If many SSTables arrive and congest at the cache, Caiti forwards them to PMem for avoidance of stalling. Comparatively, PMBD, PMBD-70, and LRU forcefully drain buffered data upon fsync or full cache while DAX's and NOVA's DAX does not use OS's page cache for buffering. As to PMBD-70, overwhelming bulky I/Os continually fill the 70% cache capacity and many cache evictions thus have to happen on the critical path. COA struggles to find an opportune moment to move out dirty data to PMem. However, the cold/hot prediction strategy is rendered ineffective for COA, since all SSTable files are written only once to become immutable. Worse, fsyncs continually drain the cache over time and hinder COA from accurately identifying cold or hot data. For the same reasons, the fsync time of Caiti is also much lower than that of PMBD and LRU, as shown in Figure 2b. PMem yet achieves comparable or a bit higher write performance because it benefits from no cost for block-level write atomicity and the use of OS's page cache for I/O staging. Figure 9a depicts the cumulative distribution of response time for LevelDB on fillrandom with 4KB values. Evidently Caiti climbs to the end of 100 percent at a much steeper rate. A comparison on these curves also tells that Caiti incurs shorter tail latency, thereby enabling higher quality of service.

As to read-intensive workloads, both readrandom and readhot follow a uniform distribution to fetch data, except that readhot's accesses are limited to a small range to simulate the hot spots commonly found in real-world situations. A comparison between Figure 8c and Figure 8d shows that readhot generally demands much less response time. This is because it is easier for LevelDB's and OS's caches to buffer and hit hot data than all data. DAX and NOVA are two that read data with DAX bypassing OS's page cache, so they spend increasingly more time in loading larger values

I/O Transit Caching for PMem-based Block Device



Figure 8: A Comparison with LevelDB on Fillrandom, Overwrite, Readrandom, and Readhot

from PMem. By taking advantage of upper-level caches, Caiti yields comparable read performance than others, which further justifies the soundness of its write buffering strategy.

In addition, we also use aforementioned four workloads issued by db_bench to test the two variants with either eager eviction or conditional bypass disabled, i.e., 'w/o EE' and 'w/o BP', respectively. As shown in Figure 8a and Figure 8b with two write-intensive workloads, the full Caiti spends less time than either variant lacking one of two components. Interestingly, as shown in Figure 8c and Figure 8d for two read-intensive workloads, i.e., readrandom and readhot, respectively, the variant without eager eviction generally achieves a bit shorter average response time than two other Caiti variants with eager eviction. The reason is that, the eager eviction strategy aggressively evicts cached data down to PMem-based block device and is likely to result in more cache misses, leading to longer response latency.

To showcase the versatility of Caiti in a more rigorous manner, we subject LevelDB to typical workloads generated by YCSB (Yahoo! Cloud Serving Benchmark) [13]. Besides loading data into LevelDB, we evaluate Caiti with YCSB's A (update-heavy workload with 50% point queries and 50% updates) and F (workload with 50% read-modify-write and 50% read) that are both mixed with write and read requests. We configure uniform, zipfian, and latest distributions to cover various data access patterns that mimic real-world



(b) YCSB Benchmark with Uniform Distribution on LevelDB Fillrandom with 4KB Values DAX Pmem DOVA BTT PMBD PMBD-70 BLRU COA Caiti DAX Pmem NOVA BTT PMBD PMBD-70 BLRU COA Caiti 80 40 Throughput (MB/s) Throughput (MB/s) 30 60 40 20 10 20 0 0 Run A Load A Run A Run F Load A Run F Zipfian Latest

(c) YCSB Benchmark with Zipfian Distribution on LevelDB

(d) YCSB Benchmark with Latest Distribution on LevelDB

Figure 9: Fillrandom Latency and YCSB Throughput on LevelDB



Figure 10: A Comparison with Filebench on Webserver, Webproxy, Varmail, and Fileserver in VM

scenarios. As shown by Figure 9b to Figure 9d, Caiti still achieves higher or comparable performance than other caching algorithms. For example, as depicted in Figure 9c with data following the skewed Zipfian distribution, the throughput of Caiti is 65.8%, 40.0%, 66.3% and 51.8% higher than PMBD, PMBD-70, LRU and COA, respectively, upon loading data. Caiti also outperforms other caching algorithms under the latest distribution with workloads A and F which alternately write and read data over time. We note that the highly-skewed latest distribution always chooses the most recent data for operation. As indicated by Figure 9d, Caiti efficiently exploits the locality of cached data. These results with varying data access distributions justify the rationality of Caiti's I/O transit strategy with eager eviction and conditional bypass.

5.3.2. Virtual Machine

Virtual machines (VMs) are widely used in multi-tenant cloud and practitioners have considered deploying PMem to enhance VMs [22, 29, 50]. To evaluate the efficacy of Caiti for VM, we utilize Filebench [17] to generate different access patterns with Webserver, Webproxy, Varmail, and Fileserver workloads in a VM with disk images atop PMem. The guest OS is Ubuntu 21.04 with Ext4 as file system, running within QEMU 6.2.0 and KVM [68, 38]. VM disk images are set in the RAW format with virtio enabled. We use the default writeback cache policy for QEMU and the VM is with 32GB DRAM and eight vCPUs.

Figure 10a depicts average response time for Caiti and other algorithms with aforementioned workloads running in VM. Caiti reduces the average response time by 12.0%, 13.3%, 62.0%, 43.9%, 52.4%, 48.6%, 53.1% and 43.2% than DAX, PMem, NOVA, BTT, PMBD, PMBD-70, LRU and COA respectively, on Fileserver. Fileserver is write-intensive and emulates I/O activities for a multi-threading file-server that generates numerous write operations for multiple files. Caiti is effectual because of concurrent eager evictions and flexible conditional bypasses. As to Varmail that simulates a mailserver and performs mandatory synchronization operations (i.e., fsync) after writing files, Caiti takes 21.7%, 22.4%, 26.8%, 31.8%, 8.2% and 30.9% less time than NOVA, BTT, PMBD, PMBD-70, LRU and COA. Though, it is slightly inferior to DAX and

PMem since the latter two have no cost of guaranteeing blocklevel atomic writes and flushing buffered data.

As for two read-intensive workloads, i.e., Webserver and Webproxy, Caiti still achieves comparable performance with other algorithms. Webserver produces a sequence of read operations on files and simulates a web-server's I/O activity, while Webproxy emulates a web-proxy server. Interestingly, on Webproxy, Caiti takes 12.7%, 14.3%, 11.9%, 15.3%, 9.6% and 13.4% less time than NOVA, BTT, PMBD, PMBD-70, LRU and COA. The reason is that Webproxy consists of a mix of write and read operations while Webserver is dominated by read operations. In particular, we collect the read latency with Filebench's workloads and present them in Figure 10b. It is evident that compared to other caching algorithms, Caiti's eager eviction hardly affects its read latency.

6. Related Work

In this section, we discuss related works mainly on cache management and storage systems built on PMem.

6.1. Cache Management

Caching has been proven to be useful for storage systems. Researchers explored various novel cache management policies tailored with regard to specific access patterns and the usage status of storage systems [61, 10, 12, 71, 62, 63, 36, 58, 60].

SSD cache management. Leveraging the device cache to buffer metadata for address translation or data has been thoroughly studied in the development of SSDs [20, 57, 27, 74]. For example, For example, Wang and Wong [71] designed TreeFTL that organizes address translation and data pages in a tree-like structure in the device cache of SSD. TreeFTL dynamically adjusts the partitions for address mapping and data buffering in order to adapt to a workload's access patterns. Sun et al. [62] proposed dynamic active and collaborative cache management, namely DAC, with a cache composed of cold/hot caches and ghost cold/hot caches. DAC adjusts the real cache size based on observing cold/hot data in I/O requests with ghost caches. The idea of Caiti is portable with SSDs, especially NVMe ones with fast access speed and multiple hardware queues [40, 37, 41, 61].

Cache management in file system. Researchers have considered effectual cache management in developing file systems for the evolving storage stack. HiNFS is a promising design that uses DRAM buffer to accelerate direct accesses with NVM in the file system [12]. In short, HiNFS employs a write buffer in DRAM to cache lazy-persistent writes while performing direct access to PMem for eager-persistent writes. Kannan et al. [36] proposed DevFS that is a devicelevel file system, providing performant direct-access capabilities within a storage device. In particular, DevFS employs reverse caching to move inactive data structures of the file system off the device to host memory and coordinates with the OS to ensure secured file access. Caiti shares similarities with HiNFS and DevFS in caching data by using a part of host OS's DRAM space, instead of a device's internal cache, for underlying storage. Whereas, Caiti takes effects as I/O transit rather than I/O staging.

Machine learning-based cache management. Recently, a few researchers have taken machine learning (ML) based approaches for cache management [55, 58, 60]. For example, at the CPU cache level, Sethumurugan et al. [55] used machine learning as an offline tool to design a new replacement policy for CPU cache. TCRM [58] addresses the trade-off between thermal and cache contention-induced slowdowns. It uses a neural network-based model to predict slowdown caused by cache contention. At the storage cache level, NCache [60] uses an ML model to predict data reaccess before eviction, preferentially evicting data unlikely to be accessed again and conserving cache space for frequently accessed data. The idea of Caiti and such ML-based techniques complement each other. The way Caiti separates cold from hot data through LRU evictions is simple but a bit coarsegrained. These ML-based techniques can help Caiti to make a fine-grained separation. However, Caiti has to schedule an offline training and profiling, as it is non-trivial to build an ML framework in Linux kernel. Moreover, the workloads Caiti serves vary significantly in different environments and over time. This is the second concrete challenge that Caiti shall consider when employing ML-based techniques. We leave this work for future exploration.

In all, these works aim to advance cache management policies, with a focus on improving performance, reducing access latency, and enhancing load and store efficiency in storage systems. Comparatively, Caiti boosts the performance of PMem-based block devices by leveraging a DRAM buffer with I/O transit strategy to accelerate writes.

6.2. Storage Systems on PMem

The development of NVM technologies has motivated researchers to develop various approaches in effectively utilizing NVM as PMem [80, 11, 81, 24, 43, 72].

File system on PMem. Developing new file systems for PMem has drawn wide attention, such as PMFS [15], HiNFS, and NOVA. They mainly follow the DAX fashion that does not take PMem as block device but memory. Chen et al. [11] considered improving file access speed by minimizing kernel overhead. They expose files into user-space in constant time independent of file sizes. They also implement efficient user-space journaling for consistency. Yang et al. [81] addressed problems caused by physical superpages in PMem file system. Their design utilizes virtual superpages and includes Multi-grained Copy-on-Write (MCoW) and Zero-copy File Data Migration (ZFDM) mechanisms to reduce write amplification and improve space utilization efficiency. The concept of Caiti can be applied to DAXbased file systems developed for PMem. Although such file systems do not format PMem with BTT to be block devices, they mainly manage PMem space in the unit of pages and tend to perform data updates in the COW fashion for data consistency. These leave an opportunity for Caiti.

The use of DRAM-PMem system. The slower access speed of PMem has motivated researchers to build a hybrid DRAM-PMem system [24, 83, 85, 72]. Sorting, indexing, and KV stores have been studied with DRAM-PMem. Hua et al. [24] conducted extensive experiments on various sorting methods adapted for PMem and further considered designing PMem-friendly sorting techniques on DRAM-PMem system. The PMSort they proposed adaptively selects optimal algorithms and reduces failure recovery overhead. Li et al. [43] designed MuHash that is a novel persistent and concurrent hashing index for DRAM-PMem. MuHash employs a multi-hash function scheme to solve the cascading write problem in PMem-based open-addressed hash-based indexes. Wang et al. [72] presented a server-bypass architecture for KV stores on DRAM-PMem system. Their design incorporates hopscotch hashing for latch-enabled append operations and a fully server-bypass read/write paradigm, so as to eliminate network round trips and reduce hash conflicts for efficient client access to KV stores. Caiti can also be viewed as a design built on DRAM-PMem system. However, Caiti works at a lower level as part of device driver. In addition, indexing, sorting, and KV stores need applicationlevel caches [79, 85], which can be managed with Caiti's I/O transit strategy.

7. Conclusion

In this paper, we revisit the use of PMem as block device and consider adding a cache to BTT that achieves block-level write atomicity. We develop Caiti which, in contrast to I/O staging caches, promptly transits buffered data into PMem in order to avoid I/O stalls caused by full cache or fsyncs. To further alleviate I/O congestion, it conditionally bypasses full cache without waiting for the drain of cache slots. Leveraging multi-core CPU, Caiti achieves high concurrency and scalability. It also retains block-level write atomicity. We have thoroughly evaluated Caiti. Caiti substantially boosts performance for BTT by as much as $3.6 \times$ in extensive experiments conducted with both micro-benchmarks and real-world applications.

Acknowledgment

This work was jointly supported by National Key R&D Program of China No. 2022YFB4401700, Natural Science Foundation of Shanghai under Grants No. 22ZR1442000 and 23ZR1442300, and ShanghaiTech Startup Funding. We are very grateful to Mr. Meng Chen for his valuable help.

In addition, this paper was eventually accepted for production by the Journal of Systems Architecture: Embedded Software Design (JSA) ² after 10 months (12 May 2023 to 10 March 2024) of processing. The Associate Editor of JSA who was in charge of this paper invited overall fourteen (14) reviewers and the paper underwent five (5) rounds of revisions. To our best knowledge, such a number of reviewers should be very rare for a research paper in the domain of electrical engineering and computer science.

References

- Aghaei Khouzani, H., Xue, Y., Yang, C., Pandurangi, A., 2014. Prolonging PCM lifetime through energy-efficient, segment-aware, and wear-resistant page allocation, in: Proceedings of the 2014 International Symposium on Low Power Electronics and Design, Association for Computing Machinery, New York, NY, USA. p. 327–330. URL: https://doi.org/10.1145/2627369.2627667, doi: 10.1145/2627369. 2627667.
- [2] Ahmed, S., Bhatti, N.A., Alizai, M.H., Siddiqui, J.H., Mottola, L., 2019. Efficient intermittent computing with differential checkpointing, in: Proceedings of the 20th ACM SIGPLAN/SIGBED International Conference on Languages, Compilers, and Tools for Embedded Systems, Association for Computing Machinery, New York, NY, USA. p. 70–81. URL: https://doi.org/10.1145/3316482.3326357, doi: 10.1145/3316482.3326357.
- [3] Alcorn, P., 2022. Intel kills Optane memory business, pays \$559 million inventory write-off. https://www.tomshardware.com/news/in tel-kills-optane-memory-business-for-good. [Online; accessed 02-June-2023].
- [4] Axboe, J., Fio flexible I/O tester. https://github.com/axboe/fio.[Online; accessed 02-June-2023].
- [5] Benson, L., Makait, H., Rabl, T., 2021. Viper: An efficient hybrid PMem-DRAM key-value store. Proc. VLDB Endow. 14, 1544–1556. URL: https://doi.org/10.14778/3461535.3461543, doi: 10.14778/3 461535.3461543.
- [6] BTRFS, 2022. Hardware considerations. https://btrfs.readthedocs.io/en/latest/Hardware.html#when-thingsgo-wrong. [Online: accessed 02-June-2023].
- [7] Caulfield, A.M., Mollov, T.I., Eisner, L.A., De, A., Coburn, J., Swanson, S., 2012. Providing safe, user space access to fast, solid state disks, in: Proceedings of the Seventeenth International Conference on Architectural Support for Programming Languages and Operating Systems, Association for Computing Machinery, New York, NY, USA. p. 387–400. URL: https://doi.org/10.1145/2150976.2151017, doi: 10.1145/2150976.2151017.
- [8] Chang, Y.S., Liu, R.S., 2019. OPTR: Order-Preserving translation and recovery design for SSDs with a standard block device interface, in: 2019 USENIX Annual Technical Conference (USENIX ATC 19), USENIX Association. pp. 1009–1024.
- [9] Chen, C.H., Hsiu, P.C., Kuo, T.W., Yang, C.L., Wang, C.Y.M., 2012. Age-based PCM wear leveling with nearly zero search cost, in: Proceedings of the 49th Annual Design Automation Conference (DAC '12), ACM. pp. 453–458. doi: 10.1145/2228360.2228439.
- [10] Chen, F., Mesnier, M.P., Hahn, S., 2014. A protected block device for persistent memory, in: 2014 30th Symposium on Mass Storage Systems and Technologies (MSST), pp. 1–12. doi: 10.1109/MSST.2 014.6855541.
- [11] Chen, X., Sha, E.H.M., Zhuge, Q., Wu, T., Jiang, W., Zeng, X., Wu, L., 2018a. UMFS: An efficient user-space file system for nonvolatile memory. Journal of Systems Architecture 89, 18–29. URL: https://www.sciencedirect.com/science/article/pii/S13837621173 05064, doi: https://doi.org/10.1016/j.sysarc.2018.04.004.
- [12] Chen, Y., Shu, J., Ou, J., Lu, Y., 2018b. HiNFS: A persistent memory file system with both buffering and direct-access. ACM Trans. Storage 14. URL: https://doi.org/10.1145/3204454, doi: 10.1145/3204454.
- [13] Cooper, B.F., Silberstein, A., Tam, E., Ramakrishnan, R., Sears, R., 2010. Benchmarking cloud serving systems with YCSB, in: Proceedings of the 1st ACM Symposium on Cloud Computing, Association for Computing Machinery, New York, NY, USA. p. 143–154. URL: https://doi.org/10.1145/1807128.1807152, doi: 10.1145/1807128. 1807152.
- [14] Corbet, J., . Atomic I/O operations. https://lwn.net/Articles/552 095/. [Online; accessed 02-June-2023].

 $^{^2} The website of JSA: https://www.sciencedirect.com/journal/journal-of-systems-architecture$

- [15] Dulloor, S.R., Kumar, S., Keshavamurthy, A., Lantz, P., Reddy, D., Sankaran, R., Jackson, J., 2014. System software for persistent memory, in: Proceedings of the Ninth European Conference on Computer Systems, ACM. pp. 1–15. doi: 10.1145/2592798.2592814.
- [16] Everspin Technologies, Spin-transfer torque MRAM technology. https://www.everspin.com/spin-transfer-torque-mram-technology. [Online; accessed 02-June-2023].
- [17] Filebench, 2020. Filebench: File system and storage benchmark that uses a custom language to generate a large variety of workloads. https://github.com/filebench/filebench. [Online; accessed 02-June-2023].
- [18] Gogte, V., Wang, W., Diestelhorst, S., Kolli, A., Chen, P.M., Narayanasamy, S., Wenisch, T.F., 2019. Software wear management for persistent memories, in: 17th USENIX Conference on File and Storage Technologies (FAST 19), USENIX Association, Boston, MA. pp. 45–63. URL: https://www.usenix.org/conference/fast19/prese ntation/gogte.
- [19] Google, 2024. LevelDB. https://github.com/google/leveldb; [Online; accessed 16-February-2023].
- [20] Gupta, A., Kim, Y., Urgaonkar, B., 2009. DFTL: A flash translation layer employing demand-based selective caching of page-level address mappings, in: Proceedings of the 14th International Conference on Architectural Support for Programming Languages and Operating Systems, Association for Computing Machinery, New York, NY, USA. p. 229–240. URL: https://doi.org/10.1145/1508244.1508271, doi: 10.1145/1508244.1508271.
- [21] Hakert, C., Chen, K.H., Schirmeier, H., Bauer, L., Genssler, P.R., von der Brüggen, G., Amrouch, H., Henkel, J., Chen, J.J., 2022. Software-managed read and write wear-leveling for non-volatile main memory. ACM Trans. Embed. Comput. Syst. 21. URL: https: //doi.org/10.1145/3483839, doi: 10.1145/3483839.
- [22] Hansen, A., Downie, K., . Understand and deploy persistent memory. https://learn.microsoft.com/en-us/azure-stack/hci/concepts/dep loy-persistent-memory. [Online; accessed 02-June-2023].
- [23] Hu, J., Zhuge, Q., Xue, C.J., Tseng, W.C., Sha, E.H.M., 2013. Software enabled wear-leveling for hybrid PCM main memory on embedded systems, in: 2013 Design, Automation & Test in Europe Conference & Exhibition (DATE), pp. 599–602. doi: 10.7873/DATE .2013.131.
- [24] Hua, Y., Huang, K., Zheng, S., Huang, L., 2021. PMSort: An adaptive sorting engine for persistent memory. J. Syst. Archit. 120. URL: https://doi.org/10.1016/j.sysarc.2021.102279, doi: 10.1016/j.sysa rc.2021.102279.
- [25] Huang, F., Feng, D., Xia, W., Zhou, W., Zhang, Y., Fu, M., Jiang, C., Zhou, Y., 2016. Security RBSG: Protecting phase change memory with security-level adjustable dynamic mapping, in: 2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS), pp. 1081–1090. doi: 10.1109/IPDPS.2016.22.
- [26] Huang, K., Mei, Y., Huang, L., 2020. Quail: Using NVM write monitor to enable transparent wear-leveling. Journal of Systems Architecture 102, 101658. URL: https://www.sciencedirect.com/ science/article/pii/S1383762119304655, doi: https://doi.org/10.1 016/j.sysarc.2019.101658.
- [27] Huang, P.C., Chang, Y.H., Kuo, T.W., 2012. Joint management of RAM and flash memory with access pattern considerations, in: Proceedings of the 49th Annual Design Automation Conference, Association for Computing Machinery, New York, NY, USA. p. 882–887. URL: https://doi.org/10.1145/2228360.2228518, doi: 10.1145/2228360.2228518.
- [28] Intel, 2024. Intel optane memory responsive memory, accelerated performance. https://www.intel.com/content/www/us/en/products /details/memory-storage/optane-memory.html. [Online; accessed 16-February-2024].
- [29] Intel Coporation, . Scaling MySQL in the cloud with Intel Optane persistent memory. https://www.intel.com/content/dam/www/ public/us/en/documents/white-papers/scaling-mysql-in-the-clo ud-with-optane-persistent-memory-paper.pdf. [Online; accessed 02-June-2023].

- [30] Intel Corpration, a. ipmctl-create-namespace creates a namespace from a persistent memory region. https: //github.com/intel/ipmctl/blob/master/Documentation/ipmctl/Persi stent_Memory_Provisioning/ipmctl-create-namespace.txt. [Online; accessed 15-Dec-2023].
- [31] Intel Corpration, b. ndctl-create-namespace provision or reconfigure a namespace. https://github.com/pmem/ndctl/blob/main/Document ation/ndctl/ndctl-create-namespace.txt. [Online; accessed 15-Dec-2023].
- [32] Intel Corpration, c. Speeding up I/O workloads with intel Optane persistent memory modules. https://www.intel.com/content/ww w/us/en/developer/articles/technical/speeding-up-io-workloadswith-intel-optane-dc-persistent-memory-modules.html. [Online; accessed 14-Dec-2023].
- [33] Kadekodi, R., Lee, S.K., Kashyap, S., Kim, T., Kolli, A., Chidambaram, V., 2019. SplitFS: Reducing software overhead in file systems for persistent memory, in: Proceedings of the 27th ACM Symposium on Operating Systems Principles (SOSP '19), ACM. p. 494–508. doi: 10.1145/3341301.3359631.
- [34] Kang, W.H., Lee, S.W., Moon, B., Kee, Y.S., Oh, M., 2014. Durable write cache in flash memory SSD for relational and NoSQL databases, in: Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data (SIGMOD '14), ACM. p. 529–540. doi: 10.1145/2588555.2595632.
- [35] Kang, W.H., Lee, S.W., Moon, B., Oh, G.H., Min, C., 2013. X-FTL: Transactional FTL for SQLite databases, in: Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data (SIGMOD '13), ACM. p. 97–108. doi: 10.1145/2463676.2465326.
- [36] Kannan, S., Arpaci-Dusseau, A.C., Arpaci-Dusseau, R.H., Wang, Y., Xu, J., Palani, G., 2018. Designing a true Direct-Access file system with DevFS, in: 16th USENIX Conference on File and Storage Technologies (FAST 18), USENIX Association. pp. 241–256.
- [37] Kim, S., Lee, G., Woo, J., Jeong, J., 2021. Zero-copying I/O stack for low-latency SSDs. IEEE Computer Architecture Letters 20, 50–53. doi: 10.1109/LCA.2021.3064876.
- [38] KVM, 2019. Kernel virtual machine. https://www.linux-kvm.org/ page/Main_Page. [Online; accessed 02-June-2023].
- [39] Lee, C., Song, Y., Shin, Y., 2019a. Endurance enhancement of multilevel cell phase change memory, in: 2019 IEEE/ACM International Conference on Computer-Aided Design (ICCAD), pp. 1–8. doi: 10.1109/ICCAD45719.2019.8942175.
- [40] Lee, G., Shin, S., Song, W., Ham, T.J., Lee, J.W., Jeong, J., 2019b. Asynchronous I/O stack: A low-latency kernel I/O stack for ultra-low latency SSDs, in: Proceedings of the 2019 USENIX Conference on Usenix Annual Technical Conference, USENIX Association, USA. p. 603–616.
- [41] Li, H., Putra, M.L., Shi, R., Lin, X., Ganger, G.R., Gunawi, H.S., 2021. LODA: A host/device co-design for strong predictability contract on modern flash storage, in: Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles, Association for Computing Machinery, New York, NY, USA. p. 263–279. URL: https://doi.org/10.1145/3477132.3483573, doi: 10.1145/3477132.34 83573.
- [42] Li, W., Shuai, Z., Xue, C.J., Yuan, M., Li, Q., 2019. A wear leveling aware memory allocator for both stack and heap management in PCMbased main memory systems, in: 2019 Design, Automation & Test in Europe Conference & Exhibition (DATE), pp. 228–233. doi: 10.23919/DATE.2019.8715132.
- [43] Li, Y., Zeng, L., Chen, G., Gu, C., Luo, F., Ding, W., Shi, Z., Fuentes, J., 2022. A multi-hashing index for hybrid DRAM-NVM memory systems. Journal of Systems Architecture 128, 102547. URL: https://www.sciencedirect.com/science/article/pii/S13837621220 01047, doi: https://doi.org/10.1016/j.sysarc.2022.102547.
- [44] Liao, X., Lu, Y., Xu, E., Shu, J., 2020. Write dependency disentanglement with HORAE, in: 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20), USENIX Association. pp. 549–565.

- [45] Linux manual page, . fsync, fdatasync synchronize a file's in-core state with storage device. https://man7.org/linux/man-pages/man2/ fsync.2.html; [Online; accessed 16-February-2024].
- [46] Liu, D., Wang, T., Wang, Y., Shao, Z., Zhuge, Q., Sha, E.H.M., 2014. Application-specific wear leveling for extending lifetime of phase change memory in embedded systems. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems 33, 1450– 1462. doi: 10.1109/TCAD.2014.2341922.
- [47] Liu, S., Kanniwadi, S., Schwarzl, M., Kogler, A., Gruss, D., Khan, S., 2023. Side-channel attacks on Optane persistent memory, in: 32nd USENIX Security Symposium (USENIX Security 23), USENIX Association, Anaheim, CA. pp. 6807–6824. URL: https://www. usenix.org/conference/usenixsecurity23/presentation/liu-sihang.
- [48] Lu, B., Hao, X., Wang, T., Lo, E., 2020. Dash: Scalable hashing on persistent memory. Proc. VLDB Endow. 13, 1147–1161. URL: https://doi.org/10.14778/3389133.3389134, doi: 10.14778/3389133 .3389134.
- [49] Ma, D., Feng, J., Li, G., 2014. A survey of address translation technologies for flash memories. ACM Comput. Surv. 46. URL: https://doi.org/10.1145/2512961, doi: 10.1145/2512961.
- [50] Morera, D., . Understanding persistent memory (pmem) in vSphere. https://core.vmware.com/blog/understanding-persistent-memory-p mem-vsphere. [Online; accessed 02-June-2023].
- [51] Park, D., Shin, D., 2017. iJournaling: Fine-Grained journaling for improving the latency of fsync system call, in: 2017 USENIX Annual Technical Conference (USENIX ATC 17), USENIX Association, Santa Clara, CA. pp. 787–798. URL: https://www.usenix.org/c onference/atc17/technical-sessions/presentation/park.
- [52] Qin, H., Feng, D., Tong, W., Zhao, Y., Qiu, S., Liu, F., Li, S., 2021. Better atomic writes by exposing the flash out-of-band area to file systems, in: Proceedings of the 22nd ACM SIGPLAN/SIGBED International Conference on Languages, Compilers, and Tools for Embedded Systems, Association for Computing Machinery, New York, NY, USA. pp. 12–23. URL: https://doi.org/10.1145/3461 648.3463843, doi: 10.1145/3461648.3463843.
- [53] Qureshi, M.K., Karidis, J., Franceschini, M., Srinivasan, V., Lastras, L., Abali, B., 2009. Enhancing lifetime and security of PCMbased main memory with start-gap wear leveling, in: 2009 42nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO), pp. 14–23. doi: 10.1145/1669112.1669117.
- [54] Samsung Semiconductor, 2022. Samsung electronics unveils far-reaching, next-generation memory solutions at flash memory summit 2022. https://news.samsung.com/global/samsung-electron ics-unveils-far-reaching-next-generation-memory-solutions-at-f lash-memory-summit-2022. [Online; accessed 02-June-2023].
- [55] Sethumurugan, S., Yin, J., Sartori, J., 2021. Designing a cost-effective cache replacement policy using machine learning, in: 2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA), pp. 291–303. doi: 10.1109/HPCA51647.2021.00033.
- [56] Shen, Z., Chen, F., Jia, Y., Shao, Z., 2018. DIDACache: An integration of device and application for flash-based key-value caching. ACM Trans. Storage 14. URL: https://doi.org/10.1145/3203410, doi: 10.1145/3203410.
- [57] Shim, H., Seo, B.K., Kim, J.S., Maeng, S., 2010. An adaptive partitioning scheme for DRAM-based cache in solid state drives, in: 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST), pp. 1–12. doi: 10.1109/MSST.2010.5496995.
- [58] Sikal, M.B., Khdr, H., Rapp, M., Henkel, J., 2022. Thermal- and cache-aware resource management based on ML-driven cache contention prediction, in: 2022 Design, Automation & Test in Europe Conference & Exhibition (DATE), pp. 1384–1388. doi: 10.23919 /DATE54114.2022.9774776.
- [59] Silberschatz, A., Galvin, P.B., Gagne, G., 2018. Operating System Concepts, 10th Edition. Wiley. URL: http://os-book.com/0S10/ind ex.html.
- [60] Sun, H., Cui, Q., Huang, J., Qin, X., 2023a. NCache: A machinelearning cache management scheme for computational SSDs. IEEE Transactions on Computer-Aided Design of Integrated Circuits and

Systems 42, 1810–1823. doi: 10.1109/TCAD.2022.3208769.

- [61] Sun, H., Dai, S., Huang, J., Qin, X., 2021. Co-active: A workloadaware collaborative cache management scheme for NVMe SSDs. IEEE Transactions on Parallel and Distributed Systems 32, 1437– 1451. doi: 10.1109/TPDS.2021.3052028.
- [62] Sun, H., Dai, S., Huang, J., Yue, Y., Qin, X., 2023b. DAC: A dynamic active and collaborative cache management scheme for solid state disks. Journal of Systems Architecture 140, 102896. URL: https://www.sciencedirect.com/science/article/pii/S13837621230 00759, doi: https://doi.org/10.1016/j.sysarc.2023.102896.
- [63] Tang, C., Sha, Z., Li, J., Lin, H., Chen, L., Cai, Z., Liao, J., 2023. Cache eviction for SSD-HDD hybrid storage based on sequential packing. Journal of Systems Architecture 141, 102930. URL: https://www.sciencedirect.com/science/article/pii/S13837621230 01091, doi: https://doi.org/10.1016/j.sysarc.2023.102930.
- [64] The kernel development community, . Explicit volatile write back cache control. https://docs.kernel.org/block/writeback_cache_con trol.html. [Online; accessed 02-June-2023].
- [65] The kernel development community, 2022. BTT block translation table. https://www.kernel.org/doc/html/latest/driver-api/nvdimm /btt.html.[Online;accessed02-June-2023].
- [66] The kernel development community, 2023a. Ext4 general information. https://docs.kernel.org/admin-guide/ext4.html. [Online; accessed 17-Dec-2023].
- [67] The kernel development community, 2023b. Multiqueue block IO queueing mechanism (blk-mq). https://www.kernel.org/doc/html/latest/block/blk-mq.html#multiqueue-block-io-queueing-mechanism-blk-mq. [Online; accessed 02-June-2023].
- [68] The QEMU Project Developers, . QEMU user documentation QEMU documentation. https://www.qemu.org/docs/master/syste m/qemu-manpage.html. [Online; accessed 02-June-2023].
- [69] Verma, V., 2014. Using the block translation table for sector atomicity. https://pmem.io/blog/2014/09/using-the-block-translation-table -for-sector-atomicity/. [Online; accessed 02-June-2023].
- [70] VMware Inc., Intel optane DC persistent memory "memory mode" virtualized performance study. https://www.vmware.com/content/d am/digitalmarketing/vmware/en/pdf/techpaper/performance/IntelOpt aneDC-PMEM-memory-mode-perf.pdf. [Online; accessed 09-Dec-2023].
- [71] Wang, C., Wong, W.F., 2016. TreeFTL: An efficient workloadadaptive algorithm for RAM buffer management of NAND flashbased devices. IEEE Transactions on Computers 65, 2618–2630. doi: 10.1109/TC.2015.2485221.
- [72] Wang, J., Huang, R., Huang, K., Chen, Y., 2023a. A server bypass architecture for hopscotch hashing key-value store on DRAM-NVM memories. Journal of Systems Architecture 134, 102777. URL: https://www.sciencedirect.com/science/article/pii/S13837621220 02624, doi: https://doi.org/10.1016/j.sysarc.2022.102777.
- [73] Wang, T., Liu, D., Wang, Y., Shao, Z., 2013. FTL²: A hybrid flash translation layer with logging for write reduction in flash memory, in: Proceedings of the 14th ACM SIGPLAN/SIGBED Conference on Languages, Compilers and Tools for Embedded Systems, Association for Computing Machinery, New York, NY, USA. p. 91–100. URL: https://doi.org/10.1145/2491899.2465563, doi: 10.1145/2491899.24 65563.
- [74] Wang, Y., Qin, Z., Chen, R., Shao, Z., Yang, L.T., 2016. An adaptive demand-based caching mechanism for NAND flash memory storage systems. ACM Trans. Des. Autom. Electron. Syst. 22. URL: https://doi.org/10.1145/2947658, doi: 10.1145/2947658.
- [75] Wang, Z., Taram, M., Moghimi, D., Swanson, S., Tullsen, D., Zhao, J., 2023b. NVLeak: Off-Chip Side-Channel attacks via Non-Volatile memory systems, in: 32nd USENIX Security Symposium (USENIX Security 23), USENIX Association, Anaheim, CA. pp. 6771–6788. URL: https://www.usenix.org/conference/usenixsecurity23/prese ntation/wang-zixuan.
- [76] Won, Y., Jung, J., Choi, G., Oh, J., Son, S., Hwang, J., Cho, S., 2018. Barrier-Enabled IO stack for flash storage, in: 16th USENIX Conference on File and Storage Technologies (FAST 18), USENIX

Association. pp. 211-226.

- [77] Woo, H., Han, D., Ha, S., Noh, S.H., Nam, B., 2023. On stacking a persistent memory file system on legacy file systems, in: 21st USENIX Conference on File and Storage Technologies (FAST 23), USENIX Association, Santa Clara, CA. pp. 281–296. URL: https: //www.usenix.org/conference/fast23/presentation/woo.
- [78] Wu, C., Li, Q., Ji, C., Kuo, T.W., Xue, C.J., 2020a. Boosting user experience via foreground-aware cache management in UFS mobile devices. IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems 39, 3263–3275. doi: 10.1109/TCAD.2020.30 13078.
- [79] Wu, F., Yang, M.H., Zhang, B., Du, D.H., 2020b. AC-Key: Adaptive caching for LSM-based key-value stores, in: 2020 USENIX Annual Technical Conference (USENIX ATC 20), USENIX Association. pp. 603–615. URL: https://www.usenix.org/conference/atc20/presen tation/wu-fenggang.
- [80] Xu, J., Swanson, S., 2016. NOVA: A log-structured file system for hybrid Volatile/Non-volatile main memories, in: 14th USENIX Conference on File and Storage Technologies (FAST 16), USENIX Association. pp. 323–338.
- [81] Yang, C., Yu, Z., Zhang, R., Nie, S., Li, H., Chen, X., Long, L., Liu, D., 2022. Efficient persistent memory file systems using virtual superpages with multi-level allocator. Journal of Systems Architecture 130, 102629. URL: https://www.sciencedirect.com/science/articl e/pii/S1383762122001552, doi: https://doi.org/10.1016/j.sysarc.2 022.102629.
- [82] Yang, J., Kim, J., Hoseinzadeh, M., Izraelevitz, J., Swanson, S., 2020. An empirical guide to the behavior and use of scalable persistent memory, in: 18th USENIX Conference on File and Storage Technologies (FAST 20), USENIX. pp. 169–182.
- [83] Zhan, J., Zhang, Y., Jiang, W., Yang, J., Li, L., Li, Y., 2018. Energy-aware page replacement and consistency guarantee for hybrid NVM–DRAM memory systems. Journal of Systems Architecture 89, 60–72. URL: https://www.sciencedirect.com/science/article/pii/ S1383762118300596, doi: https://doi.org/10.1016/j.sysarc.2018.07 .004.
- [84] Zhong, Y., Li, H., Wu, Y.J., Zarkadas, I., Tao, J., Mesterhazy, E., Makris, M., Yang, J., Tai, A., Stutsman, R., Cidon, A., 2022. XRP: In-Kernel storage functions with eBPF, in: 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22), USENIX Association, Carlsbad, CA. pp. 375–393. URL: https://www.usenix .org/conference/osdi22/presentation/zhong.
- [85] Zou, X., Wang, F., Feng, D., Zhu, J., Xiao, R., Su, N., 2022. A write-optimal and concurrent persistent dynamic hashing with radix tree assistance. Journal of Systems Architecture 125, 102462. URL: https://www.sciencedirect.com/science/article/pii/S13837621220 00522, doi: https://doi.org/10.1016/j.sysarc.2022.102462.



Qing Xu received her B.Eng. degree in software engineering from Hunan Normal University in 2021. She is currently a Master's student majoring in computer science in ShanghaiTech University. Her research interests include file systems, data storage, and persistent memory.



Qisheng Jiang obtained his B.Eng. degree in software engineering from Tongji University in 2021. He is currently a Master's student majoring in computer science in ShanghaiTech University. Qisheng's research interests include systems for AI, persistent memory, and key-value store.



Chundong Wang received the Bachelor's degree in computer science and technology from Xi'an Jiaotong University in 2008, and the Ph.D. degree in computer science from National University of Singapore in 2013. Currently he works in ShanghaiTech University as a tenure-track assistant

professor. Before joining ShanghaiTech, he successively worked in Data Storage Institute, A*STAR, Singapore and Singapore University of Technology and Design (SUTD). He has published more than forty research papers in IEEE TC, IEEE TDSC, ACM TOS, ACM TECS, SC, DAC, USENIX Security, USENIX ATC, USENIX FAST, etc. His research interests include data storage and computer architecture.