

High-Performance Remote Data Persisting for Key–Value Stores via Persistent Memory Region

Yongping Luo¹, Peiquan Jin¹, *Member, IEEE*, Xiaoliang Wang², *Member, IEEE*,
Zhaole Chu¹, Kuankuan Guo¹, Peng Xu, Jinhui Guo, and Fei Liu

Abstract—Key–value stores (KVStores), such as LevelDB and Redis, have been widely used in real-world production environments. To guarantee data durability and availability, traditional KVStores suffer from high write latency, mainly caused by the long network and data-persisting time. To solve this problem, this article presents a novel data-persisting path for KVStores, allowing remote clients to persist data to the KVStore server with μ s-level latency. The novelty of this study is threefold. First, we propose PMRDirect, which utilizes a persistent memory region (PMR) in the NVM express standard to construct a direct data-persisting path from the RDMA networking card (NIC) to the PMR region inside an SSD. Second, to showcase PMRDirect in KVStores, we developed a new accessing stack called PMRAccess, enabling remote clients to access existing KVStores and providing durability for each write request. Specifically, we present a low-latency RDMA-based messaging mode and a chunk-based PMR management in PMRAccess to reduce write latency and improve system throughput. Finally, we conducted extensive experiments to evaluate the performance of our proposals. We first compared PMRDirect with a few remote data-persisting paths to show its effectiveness. Then, we evaluated PMRAccess upon two KVStores, including LibCuckoo (an in-memory KVStore) and LevelDB (an in-storage KVStore). The results showed that PMRAccess outperformed the SSD-based accessing stack by up to 6.1 \times in write throughput and 36 \times in write tail latency, and it achieved 1.7 \times higher write throughput and 0.59 \times lower write tail latency over the PMEM-based accessing stack. Further, we conducted a system-to-system comparison between the PMRAccess-integrated LibCuckoo and Redis, and the results showed our proposal achieved up to 13 \times higher throughputs and 40 \times lower write latency than Redis.

Index Terms—Key–value stores (KVStores), persistent memory region (PMR), RDMA, write latency.

I. INTRODUCTION

KEY–VALUE stores (KVStores) have been extensively used to store and manage unstructured data generated from a wide range of data-intensive applications. Large

Internet companies like Google, Meta, and ByteDance all employ KVStores at scale. In these KVStores, new data or log entries from clients to the KVStore server are persisted into storage, and the associated write latency mainly includes networking and data-persisting latency. Under traditional TCP/IP-based accessing stacks and SSD-based storage, the write latency per request is at least hundreds of microseconds [1]. For latency-sensitive scenarios such as financial transaction services, it is necessary to reduce the write latency for higher performance.

Recently, some research and industry practices introduced RDMA networking card (NIC) to improve the networking latency [2], [3]. Commodity RDMA NICs can deliver $\sim 3\text{-}\mu$ s networking latency and up to 200 GbE networking bandwidth. Because RDMA only supports accessing main memory, RDMA-based KVStores can only store all data in memory and, therefore, do not support data durability. In addition, storing too much data in memory will increase the DRAM cost (about \$5.1 per GB). Some other researchers proposed to use PMEM to reduce persistent latency as PMEM offers $\sim 120\text{-ns}$ persistent latency [4]. Combined with RDMA NICs, PMEM-based KVStores [5], [6] can reduce the write latency to μ s-level while providing data durability. However, commercial PMEM products such as Optane PMEM (about \$4.2 per GB) are much more expensive than NVMe SSDs (about \$0.2 per GB) [7] and rely on specific costly CPUs. When users deploy PMEM-based KVStores at scale, the total cost is usually unacceptable for many companies. On the other hand, Intel has announced that the original Optane production line is closed, meaning that PMEM-based KVStores will become impractical because of the shortage of PMEM.

In this article, we propose to leverage an overlooked feature in the NVMe standard called persistent memory region (PMR) [10] to reduce the latency of write requests from remote clients. PMR is a PMR inside NVMe SSDs. Many commodity NVMe SSDs, such as Starblaze SSD [11] and DapuStor Haishen5 [12], have already supported the PMR technology, e.g., by adding power protection for its control memory buffer (CMB) at low cost. Unfortunately, PMR is not typically exposed and manipulated from userspace with portable interfaces. Although a few works exploited CPU’s functions for accessing PMR [13], [14], no study proposed RDMA-based access to PMR. Unlike existing studies, this article proposes to build a new data-persisting path from RDMA NICs to PMR directly and develop a low-latency accessing stack for KVStores to ensure data durability.

Manuscript received 30 July 2024; accepted 30 July 2024. This work was supported in part by the National Science Foundation of China under Grant 62072419, and in part by ByteDance Inc. This article was presented at the International Conference on Hardware/Software Codesign and System Synthesis (CODES + ISSS) 2024 and appeared as part of the ESWEK-TCAD Special Issue. This article was recommended by Associate Editor S. Dailey. (*Corresponding authors: Peiquan Jin; Fei Liu.*)

Yongping Luo, Peiquan Jin, Xiaoliang Wang, and Zhaole Chu are with the School of Computer Science and Technology, University of Science and Technology of China, Hefei 230027, China (e-mail: ypluo@mail.ustc.edu.cn; jpq@ustc.edu.cn; wx1147@mail.ustc.edu.cn; czle@mail.ustc.edu.cn).

Kuankuan Guo, Peng Xu, Jinhui Guo, and Fei Liu are with the Applied Research Center (Storage), ByteDance Inc., Beijing 100800, China (e-mail: guokuankuan@bytedance.com; peng.xu@bytedance.com; guojinhui.liam@bytedance.com; fei.liu@bytedance.com).

Digital Object Identifier 10.1109/TCAD.2024.3442992

85 On the other hand, it is not trivial to establish a fast data-
 86 persisting path from RDMA to PMR. The main challenges
 87 are: 1) no existing mechanism supports accessing PMR from
 88 local and remote CPUs simultaneously; 2) it lacks a high-
 89 performance PMR-based accessing stack for KVStores, which
 90 can ensure low latency, high throughput, and data durability
 91 when serving multiple clients simultaneously; and 3) PMR
 92 is capacity-constrained in commodity NVMe SSDs, and an
 93 efficient management scheme for PMR is needed.

94 To address the above challenges, we modify the Linux
 95 NVMe driver to expose PMR to the kernel space and
 96 userspace. Users can access PMR from remote CPUs through
 97 RDMA *Write* verbs and local CPUs through memory mapped
 98 I/O (MMIO). To the best of our knowledge, this is the first
 99 effort to build a direct high-performance data-persisting path
 100 from RDMA NICs to the PMR region inside NVMe SSDs.
 101 Briefly, we make the following contributions in this article.

- 102 1) We develop a new high-performance data-persisting path
 103 called PMRDirect for remote clients to directly write
 104 data from an RDMA NIC to the PMR region inside an
 105 NVMe SSD. We demonstrate that such a data-persisting
 106 path can offer much lower latency than existing data
 107 paths (Section III).
- 108 2) We present a new accessing stack called PMRAccess
 109 for client/server scenarios that utilizes PMRDirect.
 110 PMRAccess provides low latency and data durability for
 111 write requests from remote clients to the KVStore server.
 112 Specifically, we implemented an RDMA-based messag-
 113 ing mode and a new chunk-based PMR management
 114 scheme in PMRAccess (Section IV).
- 115 3) We conduct extensive experiments with real RDMA
 116 NICs and a PMR-enabled NVMe SSD to evalu-
 117 ate the performance of PMRDirect and PMRAccess.
 118 PMRAccess offers $0.59\times\text{--}36\times$ lower latency and
 119 $1.7\times\text{--}6.8\times$ higher write throughput than alterna-
 120 tive solutions, including SSD-based and PMEM-based
 121 accessing stacks. We further conduct a system-to-
 122 system comparison between the PMRAccess-integrated
 123 KVStore and Redis, and the results also suggest the
 124 efficiency of our proposal (Section V).

125 II. BACKGROUND AND MOTIVATION

126 A. KVStores and Its Write Latency

127 Traditional TCP/IP-based KVStores can offer 20–60- μ s
 128 access latency between a client and a KVStore server [15]. The
 129 high networking latency of TCP/IP-based KVStore accessing
 130 stacks is owing to three aspects: 1) TCP/IP-based networking
 131 stacks rely on the CPU to perform networking computation;
 132 2) KVStores may require several networking round trips to
 133 finish one request; and 3) the capability of traditional Ethernet
 134 cards is limited.

135 The latency of persisting a write to storage depends mainly
 136 on devices; the latency value is several milliseconds for HDDs,
 137 $\sim 100\ \mu$ s for SATA SSDs, and 10–60 μ s for NVMe SSDs.
 138 Except for device type, the persistent latency differs slightly
 139 according to how we achieve durability. Users can invoke
 140 `fsync` or `fdatasync` system calls to flush data to storage
 141 after the `write` system call that moves data to the page cache.

TABLE I
 DIFFERENT ACCESSING TECHNOLOGIES FOR KVSTORES REGARDING
 DURABILITY, TAIL LATENCY (P90), AND STORAGE COST. THE DATA
 IS FROM THE EXPERIMENTAL RESULTS IN SECTION V-C
 AND PREVIOUS LITERATURE [7]

| Technology | Durability | Tail Latency (P90) | Storage Cost |
|--------------------------|------------|--------------------|--------------|
| Traditional [8] | ✓ | 74.37 μ s | \$0.2 per GB |
| RDMA+DRAM [2], [3], [9] | ✗ | 2.75 μ s | \$5.1 per GB |
| RDMA+PMEM [5], [6] | ✓ | 6.55 μ s | \$4.2 per GB |
| PMRAccess (Our proposal) | ✓ | 2.75 μ s | \$0.2 per GB |

Users can also write data to storage directly by `DIRECT_IO` 142
 in POSIX interfaces. `DIRECT_IO` incurs less data copy which 143
 makes it faster than `fsync`, at the cost of page alignment 144
 constraint. `fdatasync` reduces the persistent cost of meta- 145
 data in the file system and offers lower latency than `fsync`. 146

B. RDMA NIC and PMEM-Based KVStores 147

To reduce the high write latency in KVStores, researchers 148
 propose many new hardware-oriented techniques that offer 149
 tremendous opportunities for low latency and higher through- 150
 put. Table I summarizes these techniques in terms of 151
 durability, latency (P90), and storage cost. 152

RDMA-Based KVStores: One line of research is RDMA- 153
 based networking latency reduction. They leverage the RDMA 154
 technique, namely, accessing remote memory directly, to 155
 enhance in-memory KVStore accessing. Distributed clients 156
 can send requests to the server through RDMA verbs (such as 157
 two-sided *Send/Recv* verbs and one-sided *Read/Write/Atomic* 158
 verbs). Owing to the high networking bandwidth (10–200 159
 GbE) and ultralow accessing latency ($\sim 3\ \mu$ s) of RDMA 160
 NICs, they can reduce networking latency for two orders 161
 of magnitude and improve CPU consumption for $20\times$ [9]. 162
 However, RDMA-based in-memory KVStores do not preserve 163
 data durability. Since all data are stored in main memory 164
 and accessed through RDMA, they offer high performance 165
 at the risk of data loss. Nowadays, many real-world storage 166
 scenarios, such as financial and healthcare applications, require 167
 no data loss after a system crash or power loss event. Besides, 168
 storing all data in memory exaggerates the total cost of 169
 ownership (TCO), which is crucial for commercial companies. 170

PMEM-Based KVStores: Another line of research is 171
 PMEM-based persistent latency reduction. They leverage 172
 PMEM as storage-class memory to employ KVStores. 173
 Commodity Optane PMEM modules sit on the memory bus. 174
 Users can issue memory load and store toward PMEM, similar 175
 to DRAM. When users explicitly issue `CLFLUSH` instructions 176
 on a dirty cache line associated with the PMEM address, 177
 it can force data into PMEM persistently. The persistent 178
 latency of PMEM is roughly 120 ns, which brings KVStore 179
 memory-level performance while preserving data durability. 180
 Besides, some distributed systems also combine PMEM and 181
 RDMA, which can reduce networking latency and persistent 182
 latency. However, it is uneconomical to employ PMEM-based 183
 KVStores at scale. Building a system equipped with Optane 184
 DC PMEM needs a vendor-specific CPU and motherboard, 185
 which costs \$547 for one 128-GB Optane DIMM and \$2517 186
 for one Intel Xeon Gold 6242R CPU. As the first pioneer of 187

188 PMEM, the Intel Optane production line will be closed soon
189 due to cost-benefit considerations.

190 C. NVMe PMR and Its Opportunity

191 The NVMe express (NVMe) 1.4 introduces the PMR concept
192 in 2019 [10]. PMR represents a piece of nonvolatile memory
193 located on the NVMe SSD devices. Users can map PMR
194 into virtual address space, and then the CPU can issue
195 memory read and write requests to PMR through MMIO.
196 So far, many commodity NVMe SSDs, such as Starblaze
197 SSD [11] and DapuStor SSD [12], have already supported
198 the PMR technology, which has received much attention from
199 the academia and industry [13], [14], [16], [17]. We can also
200 implement PMR by adding power protection to the DRAM
201 buffer inside SSDs. Introducing a 128-MB PMR to an 8-TB
202 NVMe SSD costs only 0.1% of the SSD price, according to
203 an internal estimation at ByteDance. As CMB and PMR are
204 functionally equivalent and we can implement PMR by CMB,
205 we refer to them only by PMR, as with previous works on
206 PMR [13], [14].

207 In this article, we propose to leverage PMR to realize
208 a low latency, low cost, and portable data-persisting path
209 named PMRDirect. In PMRDirect, clients write data directly
210 from RDMA NIC to the server-side PMR. As RDMA offers
211 low networking latency and PMR guarantees data durability,
212 PMRDirect can fulfill our goal.

213 III. PMRDIRECT: NEW REMOTE PERSISTING PATH

214 A. Overview of PMRDirect

215 Enabling the direct path from RDMA NIC to PMR can
216 reduce the networking latency and persistent latency of dis-
217 tributed write accesses, but it is a nontrivial task. Though
218 SPDK [18] provides a userspace driver to access PMR from
219 local CPUs and RDMA provides verbs to access remote
220 memory directly, users cannot issue RDMA verbs to the PMR
221 region due to lacking driver support so far. Specifically, the
222 NVMe driver needs to expose the PMR region to the kernel
223 space and provide support for pinning the PMR address and
224 preparing DMA mappings upon DMA operations; the RDMA
225 NIC driver should be able to identify the bus address of PMR.
226 In a word, Remote DMA operations toward PMR require
227 driver support from both the initiator and target sides.

228 We address the above challenges with the following con-
229 tributions. First, we resort to a long-standing framework in
230 Linux kernel called *dma-buf* [19]. *Dma-buf* enables sharing
231 buffers for hardware DMA access across multiple devices and
232 synchronizing asynchronous hardware accesses. Second, we
233 add a patch to the Linux NVMe driver (438 LoC) to expose
234 PMR as *dma-buf* objects. Afterward, we leverage an interface
235 in the RDMA library (*ibv_reg_dmabuf_mr*) to register
236 *dma-buf* objects as remote memory for clients (no LoC to the
237 RDMA NIC driver). Finally, remote clients can post RDMA
238 verbs to registered PMR buffers like normal RDMA. Besides,
239 users can map a *dma-buf* object into virtual address space
240 through *mmap*, allowing userspace access to PMR.

241 The NVMe driver acts as an exporter of *dma-buf* object
242 related to PMR. To operate on a *dma-buf* object, the

NVMe driver implements three pairs of kernel functions: 243
{*pmr_attach*, *pmr_detach*}, {*pmr_map_dmabuf*, 244
pmr_unmap_dmabuf}, and {*pmr_dmabuf_mmap*, 245
pmr_dmabuf_release}. The *pmr_attach* function is 246
called when other devices (such as RDMA NICs) want to attach 247
the *dma-buf* object. It exposes PMR memory as a DMA memory 248
pool for later DMA operations. The *pmr_map_dmabuf* 249
function is called upon DMA operations to pin the buffer for 250
RDMA NIC accessing and prepares the DMA mapping. The 251
pmr_dmabuf_mmap is called when users map this *dma-buf* 252
object to virtual address space, so it executes corresponding 253
virtual memory mappings. The other three functions are reverse 254
operations, respectively. By adding a new control operation 255
into struct *nvme_dev_ioctl*, users can allocate a specific size of 256
PMR as a *dma-buf* object through *ioctl* system call. Users 257
pass two arguments to *ioctl*, i.e., a file descriptor (related 258
to the NVMe SSD device) and an allocation size. The control 259
operation allocates the new *dma-buf* object and returns a file 260
descriptor (related to the *dma-buf* object) to the users. The 261
latter can be used to register the PMR region for RDMA or 262
mmap it to virtual address space. 263

264 B. Comparison With Other Remote Persisting Paths

265 Fig. 1 shows other remote data-persisting paths. Fig. 1(a)
266 shows a typical data-persisting path in classical KVStores,
267 which is adopted by NVMe-oF, Redis, and Memcached. A
268 remote client sends data writes to server-side main memory
269 through networking interfaces, such as TCP/IP or RDMA-
270 based networking. The server is responsible for persisting data
271 from memory to SSD. If the SSD supports PMR, an alternative
272 destination of the persisting can be PMR [13], which results
273 in Fig. 1(b). We refer to the remote data path in Fig. 1(a)
274 and (b) as SSDSync and PMRSync. After introducing PMEM
275 as the storage device, we can exploit a data-persisting path
276 from RDMA NIC to PMEM as in Fig. 1(c), which we refer to
277 as PMEMDirect [5], [6], [20]. With DDIO disable,¹ RDMA
278 *Write* verbs upon PMEM address can guarantee durability.
279 In this article, we present PMRDirect as a new remote data-
280 persisting path that offers μ s-level latency, as in Fig. 1(d).

281 *Advantage of PMRDirect:* The advantage of PMRDirect
282 over other remote data-persisting paths lies in three aspects.

- 283 1) Compared to SSDSync, it offers μ s-level low latency
284 and provides data durability. 284
- 285 2) Compared to PMEMSync, it reduces the high cost
286 of building a storage cluster that supports PMEM.
287 Evaluation in Section 5 shows that the latency of
288 PMRDirect is lower than PMEMSync under larger than
289 256-byte writes. 289
- 290 3) PMRDirect exploits peer-to-peer DMA between RDMA
291 NIC and PMR-enabled NVMe SSD, which liber-
292 ates memory bandwidth to applications. If a node is
293 equipped with multiple NVMe SSDs and RDMA ports,
294 PMRDirect is also advantageous for extending multiple
295 highly isolated direct data paths. 295

¹Intel DDIO is a feature that makes the processor cache the primary destination and source of I/O data rather than access main memory.

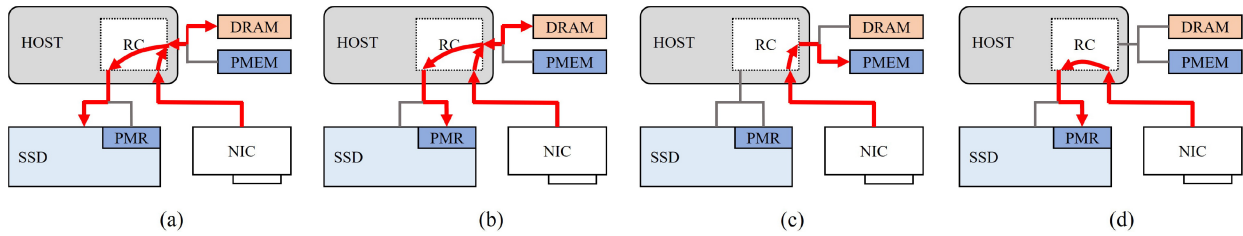


Fig. 1. Illustration of PMRDirect and other remote persisting paths. (a) SSDSync. (b) PMRSync. (c) PMEMDirect. (d) PMRDirect.

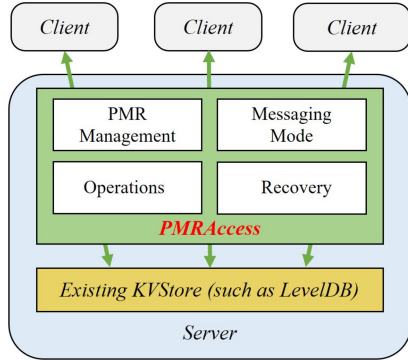


Fig. 2. Structure of PMRAccess. PMRAccess provides basic key-value operations for KVStores under client/server architecture.

296 *Potential of PMRDirect:* We can utilize PMRDirect to
 297 reduce the write latency of distributed systems that serve
 298 latency-sensitive applications such as financial transactions.
 299 Taking the chunk-based storage system at ByteDance as an
 300 example, when a client writes a piece of data to a single-
 301 node ChunkServer, the ChunkServer needs to persist data into
 302 storage before acknowledging it to the client. PMRDirect is
 303 promising for reducing this latency to μs -level, improving
 304 upstream service quality. A second use case of PMRDirect can
 305 be optimizing the well-known Raft consensus protocol [21].
 306 Some KVStores require the Raft consensus protocol to ensure
 307 consistency. A log entry in the Raft protocol is deemed
 308 committed only if most nodes have persisted it in their local
 309 log file. Therefore, the latency of remote data-persisting paths
 310 reflects the write latency of committing a new state in a
 311 distributed system. PMRDirect is promising to reduce the write
 312 latency of Raft-based distributed systems.

313 IV. PMRACCESS: PMRDIRECT-BASED ACCESSING 314 STACK FOR KVSTORES

315 To showcase the capability of PMRDirect, we propose to
 316 use PMRDirect to improve the accessing latency of existing
 317 KVStores, which results in our PMRAccess. Consider a
 318 KVStore server that serves a certain number of clients on
 319 different nodes. Clients send PUT/UPDATE/GET/DELETE
 320 requests to the server regarding any individual key-value pair,
 321 and the server then processes requests and sends results back
 322 to clients. If clients require write requests to be durable upon
 323 acknowledgment, the long write latency is always unavoidable.
 324 We present a PMRDirect-enabled KVStore accessing stack
 325 called PMRAccess to address this problem. The overview
 326 structure of PMRAccess is shown in Fig. 2. PMRAccess

is KVStore-independent, which means it can provide data 327
 durability and low access latency for both in-memory KVStore 328
 and persistent KVStore without durability requirements from 329
 the back-end KVStore. 330

331 A. Challenges of PMRAccess

PMRAccess leverages the PMRDirect remote data- 332
 persisting path to transport data directly from client-side 333
 RDMA NIC to server-side PMR, avoiding the long persistent 334
 latency exposed by NVMe SSD. However, several challenges 335
 should be overcome before PMRAccess becomes a full- 336
 fledged KVStore accessing stack. 337

*Challenge 1: How to reduce latency involved with KVStore 338
 requests with more efficient one-sided RDMA Write verbs?* 339
 Supporting requests in accessing persistent KVStore requires 340
 two-sided primitives. A complete KVStore request can be 341
 divided into three phases: 1) the client posts a request to the 342
 server; 2) the server accesses KVStore (might generate I/O 343
 to retrieve or persist data); and 3) the server acknowledges 344
 the client. The three phases are strictly ordered and require 345
 the involvement of both sides. Thus, it is straightforward 346
 to implement two-sided RDMA verbs. However, two-sided 347
 RDMA verbs exhibit higher latency than one-sided RDMA 348
 verbs. In addition, using two-sided RDMA verbs requires two 349
 round trips for each request from the client's perspective, 350
 one for sending the request and the other for receiving its 351
 acknowledgment. Therefore, the first challenge is how to use 352
 one-sided RDMA Write verbs to reduce write request latency. 353

*Challenge 2: How to manage memory effectively for PMR 354
 with limited capacity under multiple clients?* Currently, the 355
 PMR capacity is constrained. Therefore, the PMR space left 356
 for storing persistent data is limited to MB-level. If multiple 357
 clients write data to PMR simultaneously, one should be 358
 concerned about write stall due to high write contention of 359
 RDMA writes and PMR capacity shortage. Thus, PMRAccess 360
 should overcome this challenge. First, it has to persist data in 361
 PMR to SSD asynchronously and reclaim them for later usage. 362
 Second, it has to reduce the write contention of multiple clients 363
 while maximizing PMR utilization. Finally, it can recover 364
 from system crashes or power loss events and guarantee the 365
 durability of data writes. 366

367 B. RDMA Write-Based Messaging Mode

Previous works have proposed customized messaging 368
 modes for client/server scenarios, such as FaSST [22] and 369
 CatFish [23]. FaSST relies on the RDMA Send-based mes- 370
 saging mode, which incurs two network round trips per 371

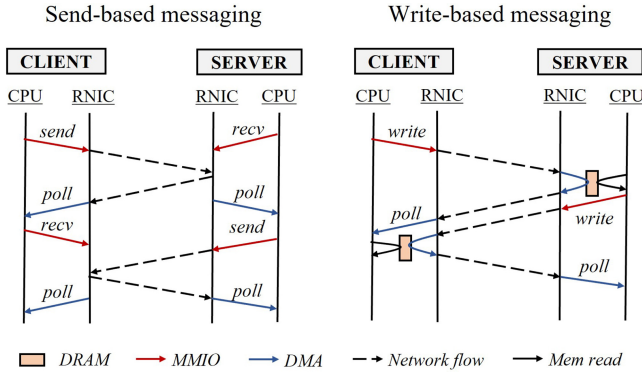


Fig. 3. Send-based messaging versus write-based messaging in RDMA.

372 request. CatFish proposes using one-sided RDMA verbs, but
 373 it maintains a central messaging queue to serve all clients and
 374 processes each request asynchronously, which increases the
 375 request latency. In contrast, we propose using RDMA Write
 376 verbs to reduce the network round trips to one per request,
 377 and the server processes request synchronously.

378 In the RDMA Send-based messaging mode [Fig. 3(a)], the
 379 client posts an RDMA *Send* verb to the server-side NIC, and
 380 then the client has to block and wait for the verb to complete.
 381 After the *Send* verb, the client posts an RDMA *Recv* verb
 382 and waits for completion. The *Send* and *Recv* verbs on each
 383 side cannot overlap as they will block the CPU and poll for
 384 a completion event. In addition, the receiver RDMA NIC has
 385 to ensure that a *Recv* verb has been posted before a *Send* verb
 386 arrives. Therefore, one request takes two network round trips.

387 In our RDMA Write-based messaging mode [Fig. 3(b)], the
 388 client sends two RDMA writes to the server-side memory
 389 addresses: one for the data request and the other to set a
 390 server-side flag that marks a data request that reaches the
 391 server. The two writes are combined in one RDMA *Write*
 392 verb and signal once to reduce network latency. Afterward,
 393 the client starts waiting for the server to set its local flag (as
 394 shown in *SendReq* in Algorithm 2). When the server finds
 395 its local flag is set, it parses and executes the new request.
 396 The server returns the result to the client and sets a client-side
 397 flag with a single RDMA *Write* verb. Our RDMA Write-based
 398 messaging mode overlaps the latency of waiting for a *Write*
 399 verb to complete with the server-side acknowledgment to come
 400 back. Therefore, it can reduce both the network round trip per
 401 request and round trip latency.

402 C. PMR Management

403 In PMRAccess, PMR is used to buffer writes for multiple
 404 clients. As the capacity of PMR is limited (e.g., only 8
 405 MB in *DapuStor Haishen5* [12]), we need to design a PMR
 406 management mechanism to reduce write contentions between
 407 multiple clients and batch-write data to the SSD efficiently,
 408 ensuring high PMR utilization efficiency.

409 RDMA-based locks are much slower than memory atomic
 410 operations [24]. Thus, we need to use a private PMR region
 411 for each client to reduce coordination costs. When the private
 412 PMR region becomes full, we must flush all its data to

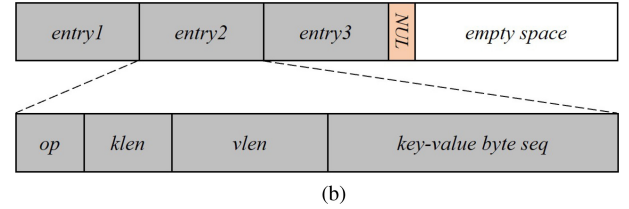
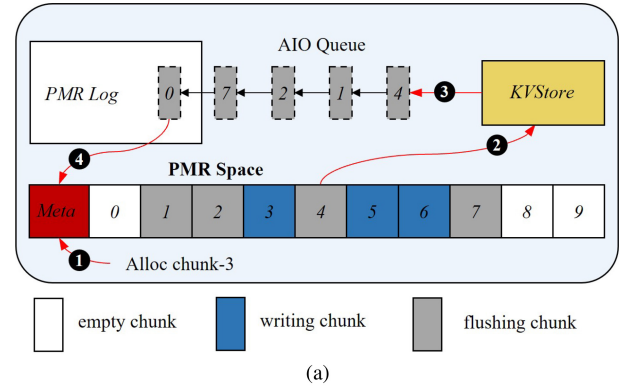


Fig. 4. PMR management scheme in PMRAccess. (a) Chunk-based PMR management. (b) Structure of a writing chunk.

413 SSD and reclaim it, blocking the client until it vacuums
 414 its PMR region and incurring high tail latency. Besides, a
 415 large private PMR region will limit the maximum number of
 416 connected clients, and a small private PMR region will limit
 417 the throughput of clients. To avoid this problem, we propose a
 418 chunk-based PMR management for multiple clients, as shown
 419 in Fig. 4(a). We divide the PMR space into identical small
 420 chunks. Initially, all chunks are empty chunks. When a new
 421 client connects to the server, it allocates an empty chunk (1)
 422 as a writing chunk. Then, it sends data write requests to this
 423 chunk and preserves durability by it. The client maintains the
 424 usage of a writing chunk until it contains no more space.
 425 Then, the client sends a *CHUNKALLOC* request to the server
 426 to transfer ownership of its writing chunk to the server. The
 427 server adds the old chunk into a queue and allocates a new
 428 empty one for the client (1). We store states of all chunks in
 429 a concurrent bitmap [Meta in Fig. 4(a)], where 0 represents
 430 an empty chunk, and 1 represents a writing chunk or a flush
 431 chunk.

432 The server processes flushing chunks with a background
 433 thread. A flushing chunk is first passed to the KVStore for
 434 potential usage (some KVStores support writing in a batch,
 435 and they can utilize the flushing chunk 2). Afterward, it is
 436 added to an asynchronous I/O (AIO) Queue to persist to an
 437 in-SSD log file (3). When a flushing chunk reaches SSD, the
 438 server modifies its state back to an empty chunk (4), and then
 439 it can be allocated again. As we persist chunk to SSD in a
 440 background manner, the long persisting latency is hidden from
 441 clients. When all empty chunks are used up, clients with no
 442 available PMR chunks will wait for the server to release new
 443 empty chunks (line 17 in Algorithm 2). In such situations,
 444 the throughput of PMRAccess will be limited by the SSD
 445 write bandwidth. This behavior is common in many existing
 446 systems, such as *RocksDB*, which also uses rate-limit and
 447 write-stall policies to postpone the client's requests.

Algorithm 1: Operations for Client

```

1 Function GET(key):
2   req.op = GET; req.key = key; req.val = NULL;
3   res = SendReq(req, pmrchunk + off);
4   return res;
5 Function PUT(key, val):
6   req.op = PUT; req.key = key; req.val = val;
7   res = SendReq(req, pmrchunk + off);
8   off = off + 8 + key.size + val.size;
9   return;
10 Function CHUNKALLOC():
11   req.op = CHUNKALLOC;
12   req.key = pmrchunk; req.val = NULL;
13   pmrchunk = SendReq(req, pmrchunk + off);
14   off = 0;
15   return pmrchunk;
16 Function SendReq(req, pmr_addr):
17   rdma_post_write(req, req.len, pmr_addr);
18   rdma_post_write_signaled(1, 1, &server_flag);
19   // wait for write completion
20   rdma_poll_complete();
21   while client_flag == 0 do
22     | nop; // wait for server
23   end
24   client_flag = 0;
25   return parse(client_buf);

```

Algorithm 2: Operations for Server (Serve One Client)

```

1 Function ProcessRequest(db):
2   while True do
3     while server_flag == 0 do
4       | nop; // wait for client
5     end
6     server_flag = 0;
7     req = parse(pmrchunk);
8     switch req.op do
9       case GET do
10        | res = db.get(req.key);
11        end
12        case PUT do
13        | res = db.put(req.key, req.val);
14        end
15        case CHUNKALLOC do
16        | add req.key to flushing AIO queue;
17        | pmrchunk = find an empty chunk;
18        | res = pmrchunk;
19        end
20      end
21      rdma_post_write(res, res.len, &client_buf);
22      rdma_post_write_signaled(1, 1, &client_flag);
23      // wait for write completion
24      rdma_poll_complete();
25    end

```

448 We depict the structure of a writing chunk in Fig. 4(b).
449 The write requests in a writing chunk are represented as an
450 entry. Each entry comprises a 1-byte operation field, a 3-byte
451 key-length field, a 4-byte value-length field, and a variable-
452 length byte sequence containing the key-value pair. The last
453 valid write request in a writing chunk has a NUL terminator
454 character followed, which separates the entry area and empty
455 area for the server upon recovery. For all write operations, the
456 client will append the request with a terminator right at the
457 position to override the last terminator. For nonwrite requests,
458 the client will not proceed to the entry area to reuse it later.

459 **D. Operations**

460 PMRAccess provides basic database access operations
461 (PUT/GET/UPDATE/DELETE) for clients to access server-
462 side KVStore (Algorithm 1). When a client connects to
463 the server, the server spawns a new thread to serve this
464 client exclusively. This thread runs the ProcessRequest
465 function as listed in Algorithm 2. ProcessRequest is
466 a dead loop that processes new requests iteratively. Basic
467 KVStore operations invoke interfaces provided by the back-
468 end KVStore. Clients send operations to the server by RDMA
469 Write-based messaging mode in Section IV-B. For write
470 requests, PMRAccess requires no durability assurance from
471 the KVStore itself, as write requests are either stored in
472 PMR chunks or flushed to the in-SSD log file. Therefore, the
473 KVStore is unnecessary to persist log entries or data writes to
474 SSD immediately.

475 PMRAccess also provides a CHUNKALLOC operation to
476 allocate a free chunk from the PMR address space. Before a
477 client can send basic operations to the server, it should check
478 whether the client owns a writing chunk with enough free
479 space. If not, it sends a CHUNKALLOC request to the server.

The old writing chunk address is embedded in the request. 480
The corresponding serving thread on the server processes this 481
request according to Algorithm 2. It works as follows: 1) pass 482
the old chunk to the back-end KVStore for potential usage; 483
2) add the old chunk to the AIO queue; and 3) allocate a new 484
chunk by changing its state in the Meta array from 0 to 1. 485

486 **E. Recovery**

Because only server-side threads can modify the Meta 487
bitmap, and these writes are all atomic, writes to PMR are 488
failure-atomic. In addition, we leverage *dma-buf* framework 489
to handle the coherent issue of CPU accessing PMR with 490
RDMA NIC, which guarantees that local writes reach PMR 491
space before RDMA NIC accesses it. Writing chunks are only 492
modified by RDMA Write verbs, where data goes from RDMA 493
NIC to PMR directly. The semantic of RDMA Write verb 494
guarantees that when the terminator character reaches PMR, 495
the data requests are valid. 496

When encountering power loss or system crash, we need 497
to recover the PMR area. It proceeds as follows: 1) scan the 498
Meta bitmap to gather all writing chunks and flushing chunks; 499
2) write each flushing chunk to the in-SSD log file and clear 500
their Meta state; 3) scan each writing chunk to locate the first 501
terminator, write all valid data requests to the log file, and clear 502
their bits; and 4) replay the log to the KVStore. To reduce the 503
log file size, we checkpoint the KVStore and then reclaim the 504
in-SSD log file. Note that the log file in PMRAccess differs 505
from the internal WAL of the KVStore. If we use PMRAccess 506
on a persistent KVStore, it may contain such a WAL file 507
that persists logs periodically, but when serving a purely in- 508
memory KVStore, the log file in PMRAccess is all we have 509
for durability. 510

V. PERFORMANCE EVALUATION

In this section, we conduct experiments to answer the following questions about our design.

- 1) How does PMRDirect perform when compared with other remote data-persisting paths? Does PMRDirect compare against PMEMDirect? The results are reported in Section V-B.
- 2) Can PMRAccess reduce the write latency of remote writes on KVStores? Can PMRAccess work efficiently on both in-memory KVStores and in-storage KVStores while providing data durability? Can PMRAccess help reduce the tail latency? The results are reported in Section V-C.
- 3) Can PMRAccess still work well in a system-to-system comparison? For example, can a PMRAccess-enabled KVStore outperform the well-known Redis? The results are shown in Section V-D.
- 4) Does the RDMA Write-based messaging mode outperform the classical RDMA Send-based messaging in terms of average latency and write bandwidth? The results are shown in Section V-E.
- 5) Can PMRAccess fully exploit the PMR capacity, and how does it impact the overall bandwidth? The results are shown in Section V-F.

A. Settings

We conduct all experiments on two physical machines: one node serves as the storage server node, and the other is the client node. Each node has two Intel Xeon Gold 6240 CPUs. Each CPU has 18 cores and a shared 24-MB L3 Cache. There are 32-kB L1I cache, a 32-kB L1D cache, and a 1-MB L2 cache per core. For networking, each node is equipped with a single-port 100 GbE Mellanox MCX515A-CCAT RDMA NIC to connect the server and client node back to back. The server node also has four 128-GB Optane PMEM DIMMs on each socket and one 3.20-TB DapuStor NVMe SSD [12] with an 8-MB on-chip PMR region.

Each node is running Linux with a 5.14.0 kernel. The Optane PMEM modules are configured into the `app-direct` mode and exposed as particular Device DAX (`devdax`) character devices. The DapuStor NVMe SSD is exposed to the users as an NVMe block device, and we install the `ext4` filesystem on it. For accessing PMEM from remote RDMA NICs, we map the `devdax` devices into virtual memory and register it for RDMA. For accessing PMR, we obtain the file descriptor through `ioctl` to the NVMe device. We map the file descriptor into virtual memory using `mmap` to support local accesses. For RDMA writes to PMR, we register PMR memory by the `ibv_reg_dmabuf_mr` interface and then issue RDMA verbs.

B. Performance of PMRDirect

This experiment aims to demonstrate the performance upper bound of our PMRDirect data path. We compare four remote data-persisting paths as elaborated in Section III-B. As we need not manipulate data, we reuse target address space as we care more about the latency and bandwidth metrics.

Fig. 5(a) shows the client-side request latency. Under small requests, PMRDirect and PMEMDirect exhibit similar write latency, which is $1.63 \mu\text{s}$ per 64-byte request. PMRSync shows 56% more write latency than PMRDirect because it needs to persist data from host memory to PMR, twice as much as PCIe transactions as in PMRDirect. SSDSync exhibits $16\text{-}\mu\text{s}$ write latency to persist data into SSD, $10\times$ that of PMRDirect. As the request size increases, PMRDirect achieves the lowest write latency. As the write request size increases 64 times, the write latency only increases for 94%, which is $3.1 \mu\text{s}$ at a 4-kB write request. PMRSync shows roughly $1 \mu\text{s}$ more write latency than PMRDirect. The latency of PMEMDirect starts rising after 256 B. As the write requests increase to 4 kB, it needs $13 \mu\text{s}$ to directly finish a write request from RDMA to PMEM. The write latency does not increase because write requests to NVMe SSD are amplified to 4 kB due to filesystem block size.

We measure the maximum write bandwidth of four competitors under 64 B, 256 B, 1 kB, and 4 kB report results in Fig. 5(b). In practice, it needs 32/16/12/32 clients to saturate the max write bandwidth for PMRDirect/PMEMDirect/SSDSync/PMRSync, respectively. PMRDirect achieves the highest write bandwidth under different request sizes. As the request size increases to 256 B, the maximum write bandwidth increases by $2.7\times$. The value reaches 2.6 GB for requests larger than 1 kB. Following PMRDirect, PMRSync exhibits a 33% lower write bandwidth than PMRDirect as a penalty for copying data up and down in terms of PCIe transactions. The max write bandwidth in PMEMDirect is only half of PMRDirect because of the low write bandwidth of commodity Optane PMEM. While for SSDSync, its max write bandwidth increases linearly as the request size. It exhibits 30-MB write bandwidth at 64-B write requests. This bandwidth reaches 1.2 GB at 4 kB requests, which is almost the peak write bandwidth of our NVMe SSD.

PMRDirect is a direct data path from RDMA NIC to PMR and will not occupy the memory bus. Therefore, it cannot be affected by other processes that occupy the memory bus bandwidth. To verify that, we measure the max write bandwidth of four data paths in the shadow of a bandwidth-stealing process that occupies server-side memory bandwidth (16 threads do random writes concurrently). We report the bandwidth loss compared with Fig. 5(b) in (c). It reveals that remote data-persisting paths, except PMRDirect, are more or less affected by the bandwidth-stealing process. The max write bandwidth of PMEMDirect/SSDSync/PMRSync drops by 23%/3%/13%, respectively. Fig. 5(c) shows that PMRDirect occupies no memory bandwidth and can scale the overall bandwidth linearly with more PMR-enabled NVMe SSDs.

C. Performance of PMRAccess

PMRAccess is a KVStore accessing stack that provides low latency and data durability for accessing any existing KVStores from distributed clients. We use PMRAccess to provide primary accesses for two types of KVStores, including

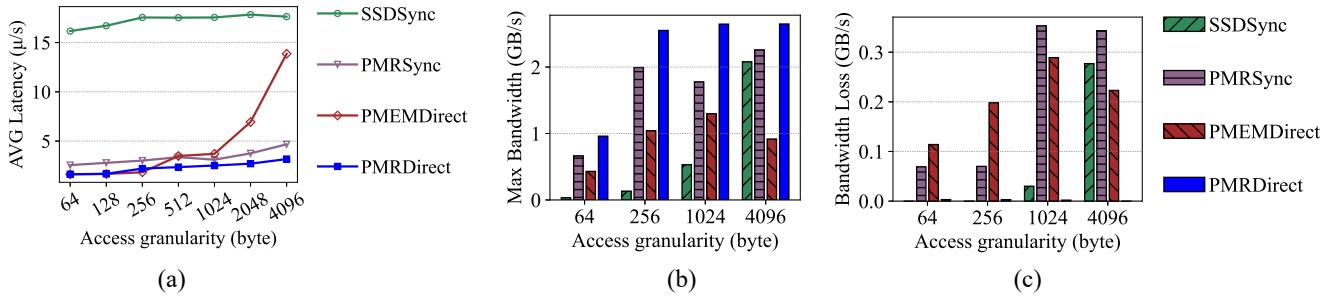


Fig. 5. Performance comparison between PMRDirect and three remote data-persisting paths: SSDSync, PMRSync, and PMEMDirect. (a) Average latency. (b) Max bandwidth. (c) Bandwidth loss.

622 LibCuckoo and LevelDB. We aim to demonstrate that our
623 PMRDirect data path is effective in optimizing KVStores.

624 *Competitors:* We compare the following KVStore accessing
625 stacks, including SyncAccess, GroupAccess, PMEMAccess,
626 and PMRAccess, all providing data durability for each write
627 request. Though they differ in how to guarantee durability
628 when sending write requests, these accessing stacks adopt
629 identical messaging modes and provide identical functional-
630 ities for accessing KVStores. In a word, the comparison is fair
631 in terms of network latency.

- 632 1) *PMRAccess:* Using RDMA Write-based messaging
633 mode and PMRDirect data path, clients send write
634 requests to the server while preserving WAL durability
635 by PMR. WAL records in PMR are persisted to SSD at
636 the unit of a Chunk.
- 637 2) *SyncAccess:* Using RDMA Write-based messaging
638 mode and SSDSync data path, clients send write requests
639 to the server while preserving WAL durability by persist-
640 ing data to SSD before acknowledging clients (adopted
641 by Redis [8]).
- 642 3) *GroupAccess:* WAL records from different clients are
643 grouped and persist to SSD in a batch as in LevelDB.
644 GroupAccess can reduce write amplification under
645 multiple-client scenarios (adopted by LevelDB [25]).
- 646 4) *PMEMAccess:* Using RDMA Write-based messaging
647 mode and PMEMDirect data path, clients send write
648 requests to the server while preserving WAL durability
649 by PMEM (adopted by FlatStore [5]).

650 *Back-End KVStores:* We compare these four accessing
651 stacks under two types of KVStores: 1) LibCuckoo and
652 2) LevelDB. It includes both in-memory KVStore and in-
653 storage KVStore to make the comparison diverse. LibCuckoo²
654 is a high-performance, concurrent hash table. LibCuckoo
655 stores key-value pairs in memory and guarantees no durability.
656 LevelDB³ is a fast key-value storage library written at Google
657 that provides an ordered mapping from string keys to string
658 values. It stores primary data in storage devices and guarantees
659 durability through log files. LevelDB can control the durability
660 of each write operation by setting its write option to persist
661 a WAL record. However, the performance is poor when
662 providing data durability for each write.

663 *Workloads:* Because PMRAccess focuses on write oper-
664 ations, we choose the Load Only, YCSB-A, and YCSB-B

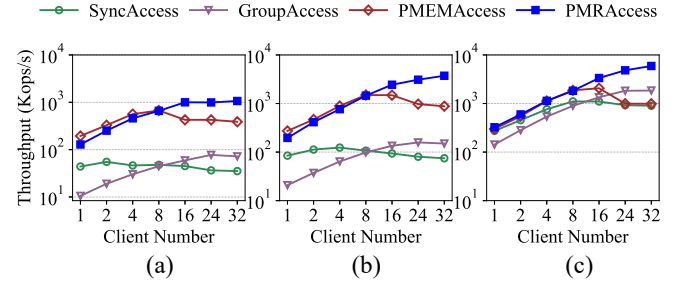


Fig. 6. Throughputs of different accessing stacks on LibCuckoo. (a) Load only. (b) YCSB-A. (c) YCSB-B.

665 query workload, which stands for write-intensive, read-write
666 balanced, and read-intensive workloads. Key-value pairs are
667 64-byte, composed of 16-byte string key and 48-byte value.
668 For larger value sizes, readers can refer to Section V-D for
669 detailed comparison under various value sizes. We populate
670 each KVStore with 128M key-value pairs. Then, we measure
671 the overall throughput and latency for each workload.

- 672 1) *Load Only:* This workload consists of 100% PUT
673 operations with randomized new keys.
- 674 2) *YCSB-A:* This workload consists of 50% GET operations
675 and 50% UPDATE operations. All keys follow the
676 Zipfian distribution. The skewness value of Zipfian
677 is 0.9.
- 678 3) *YCSB-B:* This workload consists of 95% GET opera-
679 tions and 5% UPDATE operations. Similarly, all keys
680 followed the Zipfian distribution.

681 *Performance on LibCuckoo:* This experiment evaluates
682 PMRAccess in LibCuckoo, a purely in-memory Data Store
683 that supports no durability. All accesses to KVStore come to
684 memory, but WAL records for each write request are persisted
685 by each accessing stacks with their strategies.

686 The throughputs of different accessing stacks on LibCuckoo
687 are shown in Fig. 6. We report the Load Only result in
688 Fig. 6(a). Under a single client, PMEMAccess outperforms
689 PMRAccess by 34%. As we increase the client number,
690 PMRAccess outperforms PMEMAccess and reaches 1.059
691 Mops/s, $2.3 \times / 14 \times / 27 \times$ that of PMEMAccess, SyncAccess,
692 and GroupAccess. Under the YCSB-A workload, the through-
693 puts are higher than the load-only workload. Under 32 clients,
694 PMRAccess is $3.2 \times / 47 \times / 23 \times$ better than PMEMAccess,
695 SyncAccess, and GroupAccess. Under YCSB-B workloads,
696 the gap between different competitors starts to narrow down.

²LibCuckoo: <https://github.com/efficient/libcuckoo>

³LevelDB: <https://github.com/google/leveldb>

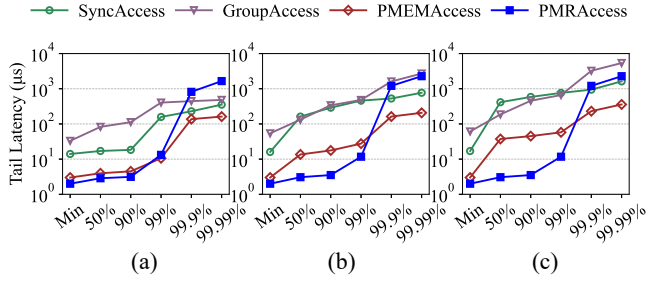


Fig. 7. Tail latency of different accessing stacks on LibCuckoo. (a) One client. (b) Eight clients. (c) Sixteen clients.

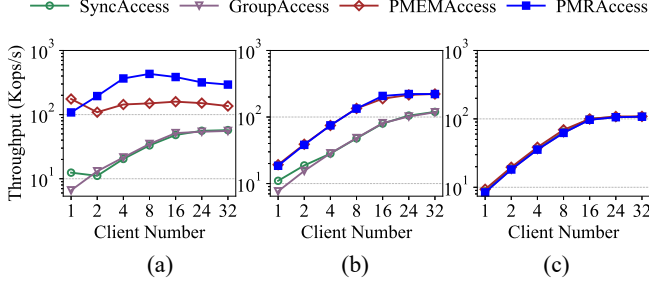


Fig. 8. Throughput of different accessing stacks on LevelDB. (a) Load only. (b) YCSB-A. (c) YCSB-B.

Both PMEMAccess and SyncAccess suffer from multiple-client contention on hardware resources. Under 32 clients, PMRAccess achieves $5.9 \times / 6.5 \times / 3.8 \times$ total throughput of PMEMAccess, SyncAccess, and GroupAccess, respectively.

The tail latency of PUT operations with different accessing stacks on LibCuckoo is in Fig. 7. We report the tail latency of four accessing stacks under 1 client [Fig. 7(a)], 8 clients [Fig. 7(b)], and 16 clients [Fig. 7(c)]. PMRAccess maintains the best Min latency, 50%, 90%, and 99% tail latency. Under one client, PMRAccess reduces latency to $3 \mu\text{s}$, followed by PMEMAccess, which is $4 \mu\text{s}$. Meanwhile, the PUT latency of SyncAccess is roughly 16 and $80 \mu\text{s}$ for GroupAccess. PMEMAccess maintains the best 99.9% tail latency at about $136 \mu\text{s}$, which is 80% less than PMRAccess, by storing all WAL records in PMEM. As we increase the client number, PMRAccess maintains its advantage at 99% tail latency. Under 16 clients, the 50% latency of PMRDirect is only 8%, 0.7%, and 1.6% of PMEMAccess, SyncAccess, and GroupAccess, respectively.

Performance on LevelDB: To evaluate PMRAccess in LevelDB, we choose the following system configurations. All competitors follow identical configurations to make the comparison fair. Each competitor is allocated with fixed memory and set with CGroups. The memory size is 4 GB, a third of the dataset, to ensure that requests are served from memory and persistent storage. The block cache size is 256 MB, and the Memtable size is 64 MB. In addition, for each write operation, the *sync* option is set to false, which means that LevelDB does not explicitly guarantee durability.

The throughputs of different accessing stacks on LevelDB are shown in Fig. 8. The overall throughput of accessing LevelDB is far less than that of LibCuckoo, as most of LevelDB’s data resides in SSD. In the *Load Only* workload,

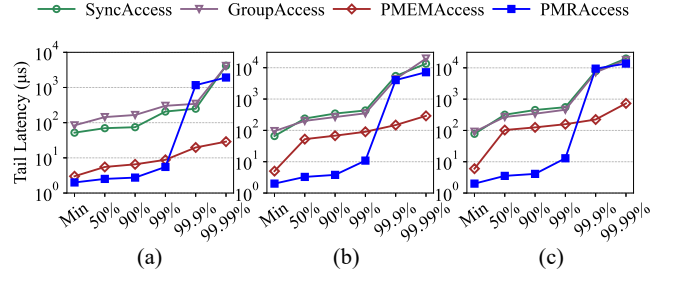


Fig. 9. Tail latency of different accessing stacks on LevelDB. (a) One client. (b) Eight clients. (c) Sixteen clients.

PMRAccess achieves the highest throughput, which is 120 Kops/s under one client. When we increase the client number, PMRDirect’s throughput increases to 430 Kops/s. The peak throughput of PMRAccess is $2.68 \times / 7.1 \times / 7.8 \times$ that of PMEMAccess, SyncAccess, and GroupAccess, respectively.

The tail latency of PUT operations with different accessing stacks on LevelDB is shown in Fig. 9. Similar to the tail-latency results on LibCuckoo, PMRAccess is the best among all competitors in terms of the Min, 50%, 90%, and 99% tail latency. Specifically, under one client, the 50% tail latency of PMRAccess is $2.5 \mu\text{s}$, which is only 45%/3.8%/1.7% of PMEMAccess, SyncAccess, and GroupAccess, respectively. In addition, PMRAccess keeps the 90% latency less than $3 \mu\text{s}$ and 99% latency less than $6 \mu\text{s}$, which is $36 \times$ faster than SyncAccess and even $0.5 \times$ faster than PMEMAccess. When we increase the client number to 16, PMRAccess’s write latency increases slightly within $1 \mu\text{s}$, which proves that PMRAccess has good performance isolation between multiple clients. SyncAccess and GroupAccess are bad in LevelDB as they waste CPU resources and SSD bandwidth on persisting WAL.

To sum up, PMRAccess provides extremely low write latency to $3 \mu\text{s}$ per remote write request, which is up to $36 \times$ better than SyncAccess/GroupAccess and $4 \times$ better than PMEMAccess. It can also maintain the best 50%, 90%, and 99% tail latency for both LibCuckoo and LevelDB. Regarding the throughput, PMRAccess achieves a high 1-Mops/s throughput on LibCuckoo and 430-Kops/s throughput on LevelDB and guarantees each remote write-request is durable. It achieves up to $18 \times$ higher throughput than SyncAccess/GroupAccess and a $1.5 \times$ higher throughput than PMEMAccess. Moreover, PMRAccess is also advantageous for performance isolation between clients and can keep high scalability when the number of remote clients increases.

D. System-to-System Comparison

Further, we conduct a system-to-system comparison to demonstrate the advantage of our proposal. We integrate PMRAccess into LibCuckoo and compare it with Redis [8] with two durability levels, including Redis-Strong and Redis-Everysec. Redis-Strong represents Redis with a durability guarantee for each write request; Redis-Everysec represents Redis with a durability guarantee for every second. We allocate identical hardware resources for each system, such as thread number and NIC resources.

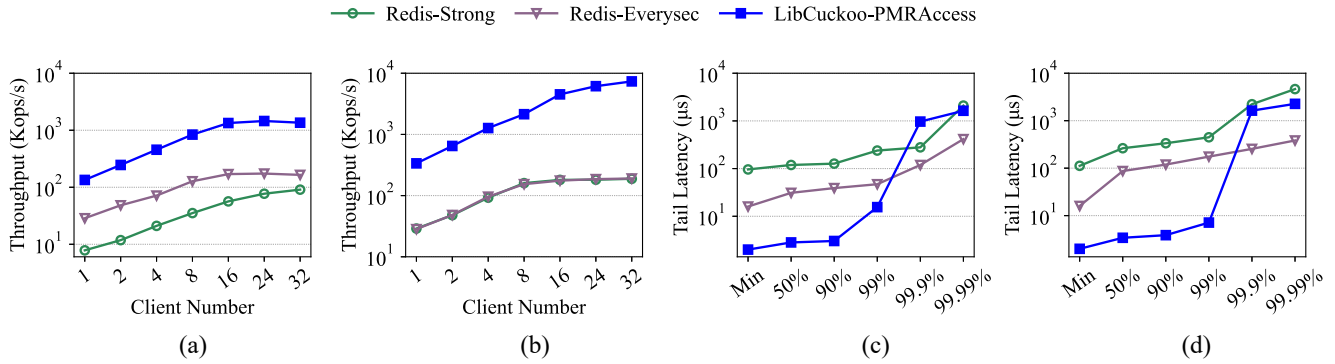


Fig. 10. Throughput and latency comparison between our system (LibCuckoo-PMRAccess) and Redis under different clients. (a) PUT. (b) GET. (c) One client. (d) Sixteen clients.

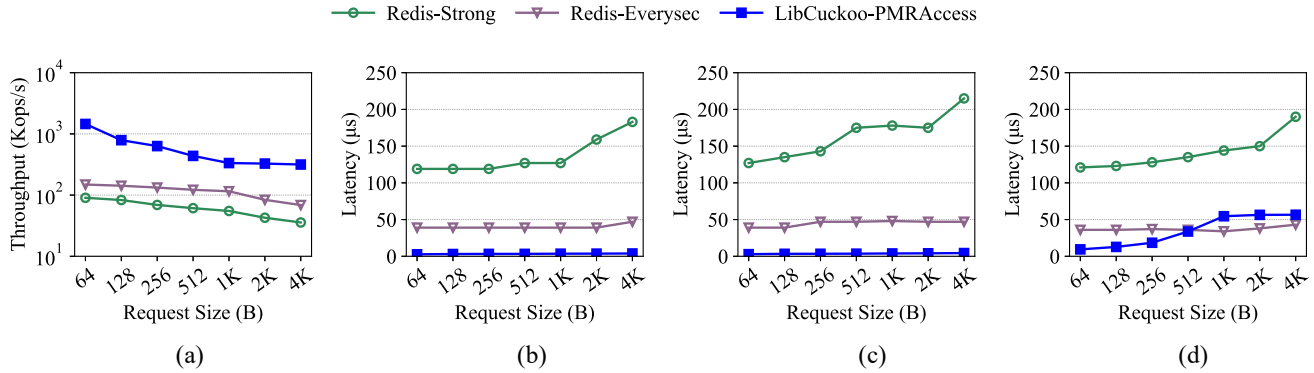


Fig. 11. Throughput and latency comparison between our system (LibCuckoo-PMRAccess) and Redis with different request sizes. (a) PUT IOPS. (b) 50% tail latency. (c) 90% tail latency. (d) AVG latency.

774 Fig. 10 shows the throughput and latency of three persistent KVStores under multiple clients. As Fig. 10(a) shows, 775 under the PUT workload, LibCuckoo-PMRAccess achieves 776 a 13.9 \times higher throughput than Redis-Strong and a 7.2 \times 777 higher throughput than Redis-Everysec with 32 clients. For the 778 GET workload, LibCuckoo-PMRAccess achieves 7 Mops/s, 779 about 38 \times higher than Redis-Strong/Redis-Everysec. The tail 780 latency of LibCuckoo-PMRAccess is reported in Fig. 10(c) 781 and (d). LibCuckoo-PMRAccess achieves the best Min latency 782 and the 50%/90%/99% tail latency. Specifically, the 90% tail 783 latency of LibCuckoo-PMRAccess is 2.3%/7.7% of Redis- 784 Strong/Redis-Everysec. LibCuckoo-PMRAccess exhibits 3- μ s 785 90% tail latency, almost as fast as one RDMA Write verb. 786 LibCuckoo-PMRAccess's 99.9% tail latency is worse than 787 Redis-Everysec as it puts more pressure on SSDs. 788

789 Fig. 11 shows the throughput and latency of three persistent 790 KVStore under request size varied from 64 B to 4 kB. At 791 a small request size, LibCuckoo-PMRAccess is 15 \times better 792 than Redis. As the request size increases to 1 kB, the PUT 793 throughput of LibCuckoo-PMRAccess is 2.89 \times and 6.06 \times 794 of Redis-Everysec and Redis-Strong throughput. Under 4 795 kB requests, LibCuckoo-PMRAccess's PUT throughput is 796 4.5 \times and 8.8 \times that of Redis-Everysec and Redis-Strong. 797 For Tail latency, LibCuckoo-PMRAccess maintains lower than 798 3.8- μ s 50% latency and 4.3- μ s 90% tail latency, which is 799 48 \times lower than that of Redis-Strong and 9 \times lower than 800 that of Redis-Everysec. In addition, LibCuckoo-PMRAccess's 801 50%/90% tail latency increases slightly with the request size,

and the 50% tail latency of Redis-Strong increases from 127 802 to 234 μ s. Fig. 11(d) reports the average latency of the three 803 KVStores. The average latency of Redis-Strong is 3.36 \times and 804 12.8 \times higher than that of Redis-Everysec and LibCuckoo- 805 PMRAccess under small request sizes. 806

In summary, the PMRAccess-enabled KVStore, LibCuckoo- 807 PMRAccess, achieves up to 13 \times higher throughputs and 40 \times 808 lower 50%/90% tail latency than Redis under small request 809 sizes. Although the throughput improvement of LibCuckoo- 810 PMRAccess decreases with the increase of the request size, 811 LibCuckoo-PMRAccess still outperforms Redis. 812

E. Benefits of the RDMA Write-Based Messaging Mode 813

In this experiment, we compare our RDMA Write-based 814 client/server messaging mode with the RDMA Send-based 815 messaging mode to demonstrate its efficiency in reducing the 816 request latency for remote clients. 817

Our RDMA Write-based messaging mode reduces the average 818 latency of client/server requests to 2.1–3.5 μ s, which is 819 only 26% of the latency observed with the RDMA Send-based 820 messaging mode. Additionally, it can improve the multiclient 821 write bandwidth by at least 3.1 \times . As the size of the write 822 request increases, our RDMA Write-based messaging can 823 maintain low request latency under 3.6 μ s and achieve up to 824 9-GB/s write bandwidth. Given that one RDMA Send verb 825 takes about 4 μ s and one RDMA Write verb takes about 2 μ s, 826 the results suggest that our messaging mode can reduce both 827

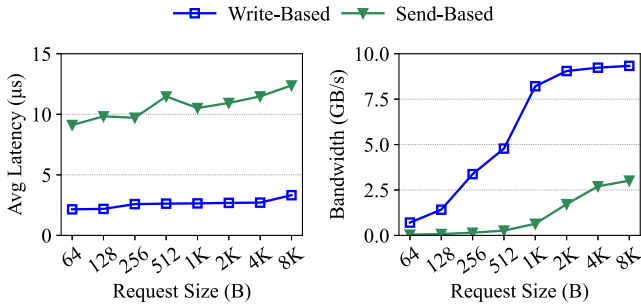


Fig. 12. Performance of different messaging modes: (a) average write latency and (b) total write bandwidth.

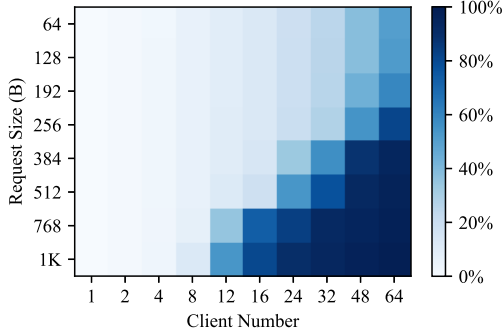


Fig. 13. Average PMR utilization ratio of PMRAccess under different configurations (darker colors mean higher utilization ratios).

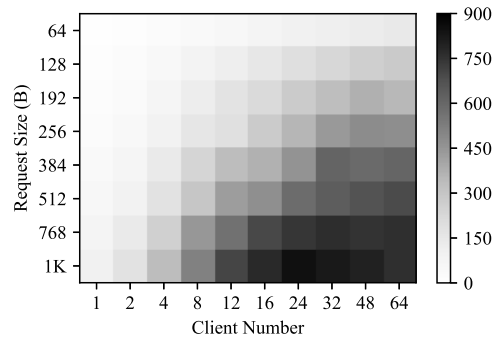


Fig. 14. Total write bandwidth of PMRAccess under different configurations (darker colors mean higher write bandwidths).

the number of network round trips and the latency of each round trip between the clients and the server.

F. PMR Utilization and Write Bandwidth

To uncover the PMR utilization of PMRAccess, we measure the real-time utilization ratio of PMR chunks under write-intensive workloads. We use LibCuckoo as the back-end KVStore engine and report the PMR utilization on the server side and the total write bandwidth on the client side. We evaluated various configurations, including the number of parallel clients and key-value request size, with the number of clients ranging from 1 to 64 and the request size ranging from 64 B to 1 kB.

Fig. 13 shows the average PMR utilization ratio under different configurations. With a small request size and a low number of clients, the PMR utilization ratio remains below 10%. As the number of clients increases to 64, the PMR utilization rises to between 50% and 80%, even with a small request size. When the number of clients exceeds 24, and the request size is larger than 512 B, the real-time PMR utilization ratio reaches 100% most of the time. Fig. 14 shows the write bandwidth of all clients under various configurations. The write bandwidth increases accordingly as the client number and request size grow to a certain extent. However, once the average PMR utilization reaches about 80%, the write bandwidth either plateaus or slightly decreases. To sum up, large write requests or excessive parallel clients tend to exhaust the PMR and eventually bottleneck the overall bandwidth.

VI. RELATED WORK

Traditional KVStores, such as Memcached, Redis [8], and Cassandra, are extensively used by commercial companies. With TCP/IP-based networking stacks and HDD/SSDs as storage devices, the write latency to KVStores for remote clients is usually hundreds of microseconds.

As the development of high RDMA NICs, many pioneers exploit the potential of reducing the networking latency in KVStores [3], [9] and transaction execution [22], [26]. However, RDMA supports accessing memory only; therefore, those RDMA-based KVStores store all data in DRAM and ignore data durability. Besides, the advent of PMEM brings new strategies for the durability of KVStores. Recently, many researchers propose PMEM-based Indexes [27], [28] and PMEM-based KVStore [4], [5]. Most of them exploit the characteristics of commodity PMEM, such as low latency and high bandwidth, to achieve high performance while ensuring data durability. Unfortunately, PMEM is too expensive compared to SSDs, and the Optane production line is closed.

So far, few studies have used PMR to improve storage and file systems. SineKV [13] first used PMR in KVStores to boost the durability of local WAL entries. Horae [14] used PMR to reduce the write order constraints of traditional file systems. X-SSD [16] proposed to use PMR for SSDs and integrated the database logging replication service into SSDs by Nontransparent Bridging networking. Unlike all the above techniques, our proposal opens up a new path that enables remote clients to write data to NVMe PMR directly, ensuring μs -level low latency while preserving data durability with low storage cost.

VII. CONCLUSION

In this article, we proposed to use NVMe PMR to build a new data-persisting path called PMRDirect to enable low-latency data writes from distributed clients. Further, we designed an accessing stack called PMRAccess for KVStores to reduce the latency of remote data writes. The experiments on real RDMA NICs and a PMR-enabled NVMe SSD showed that PMRDirect achieved the lowest write latency and the highest write bandwidth. Moreover, when evaluated on LevelDB, PMRAccess outperforms the SSD-based accessing stack by up to $6.1\times$ in write throughput and $36\times$ in write latency.

REFERENCES

- 898
- 899 [1] A. Dragojevic, D. Narayanan, M. Castro, and O. Hodson, “FaRM: Fast remote memory,” in *Proc. 11th USENIX Symp. Netw. Syst. Design Implement. (NSDI)*, 2014, pp. 401–414. 937
- 900 [2] B. Li et al., “KV-direct: High-performance in-memory key-value store with programmable NIC,” in *Proc. 26th SOSP*, 2017, pp. 137–152. 938
- 901 [3] H. Lim, D. Han, D. G. Andersen, and M. Kaminsky, “MICA: A holistic approach to fast in-memory key-value storage,” in *Proc. 11th USENIX Symp. Netw. Syst. Design Implement. (NSDI)*, 2014, pp. 429–444. 939
- 902 [4] L. Vogel, A. van Renen, S. Imamura, J. Giceva, T. Neumann, and A. Kemper, “Plush: A write-optimized persistent log-structured hash-table,” *Proc. VLDB Endowment*, vol. 15, no. 11, pp. 2895–2907, 2022. 940
- 903 [5] Y. Chen, Y. Lu, F. Yang, Q. Wang, Y. Wang, and J. Shu, “FlatStore: An efficient log-structured key-value storage engine for persistent memory,” in *Proc. 25th Int. Conf. ASPLOS*, 2020, pp. 1077–1091. 941
- 904 [6] T. Li, D. Shankar, S. Gugnani, and X. Lu, “RDMP-KV: Designing remote direct memory persistence based key-value stores with PMEM,” in *Proc. Int. Conf. High Perform. Comput., Netw., Storage Anal.*, 2020, pp. 1–14. 942
- 905 [7] K. Huang, D. Imai, T. Wang, and D. Xie, “SSDs striking back: The storage jungle and its implications to persistent indexes,” in *Proc. CIDR*, 2022, pp. 9–12. 943
- 906 [8] “Redis.” 2023. [Online]. Available: <https://redis.io/> 944
- 907 [9] C. Mitchell, Y. Geng, and J. Li, “Using one-sided RDMA reads to build a fast, CPU-efficient key-value store,” in *Proc. USENIX ATC*, 2013, pp. 103–114. 945
- 908 [10] “Persistent memory region.” NVMe. 2023. [Online]. Available: <https://nvmexpress.org/specification/nvm-express-base-specification/> 946
- 909 [11] “Starblaze OC SSD.” Starblaze. 2023. [Online]. Available: <https://www.starblaze-tech.com/en/lists/content/id/137.html> 947
- 910 [12] (DapuStor Corp., Shenzhen, China). *Dapustor Haishen5 PCIe Gen5 Enterprise SSDs*. (2023). [Online]. Available: <https://www.storagereview.com/news/dapustor-haishen5-pcie-gen5-enterprise-ssds-announced> 948
- 911 [13] F. Li, Y. Lu, Z. Yang, and J. Shu, “SineKV: Decoupled secondary indexing for LSM-based key-value stores,” in *Proc. IEEE 40th ICDCS*, 2020, pp. 1112–1122. 949
- 912 [14] X. Liao, Y. Lu, E. Xu, and J. Shu, “Write dependency disentanglement with HORAE,” in *Proc. 14th USENIX Symp. OSDI*, 2020, pp. 549–565. 950
- 913 [15] I. Zhang et al., “The Demikernel datapath OS architecture for microsecond-scale datacenter systems,” in *Proc. ACM SOSP*, 2021, pp. 195–211. 951
- 914 [16] S. Lee et al., “X-SSD: A storage system with native support for database logging and replication,” in *Proc. Int. Conf. Manag. Data*, 2022, pp. 988–1002. 952
- 915 [17] X. Liao, Z. Yang, and J. Shu, “RIO: Order-preserving and CPU-efficient remote storage access,” in *Proc. 18th Eur. Conf. Comput. Syst. (EuroSys)*, 2023, pp. 703–717. 953
- 916 [18] “Interfaces to access PMR through SPDK.” SPDK. 2023. [Online]. Available: https://spdk.io/doc/nvme_8h.html 954
- 917 [19] “Buffer sharing and synchronization.” Linux. 2023. [Online]. Available: <https://docs.kernel.org/driver-api/dma-buf.html> 955
- 918 [20] X. Wei, X. Xie, R. Chen, H. Chen, and B. Zang, “Characterizing and optimizing remote persistent memory with RDMA and NVM,” in *Proc. USENIX ATC*, 2021, pp. 523–536. 956
- 919 [21] D. Ongaro and J. Ousterhout, “In search of an understandable consensus algorithm,” in *Proc. USENIX ATC*, 2014, pp. 305–319. 957
- 920 [22] A. Kalia, M. Kaminsky, and D. G. Andersen, “FaSST: Fast, scalable and simple distributed transactions with two-sided RDMA datagram RPCs,” in *Proc. 12th USENIX Symp. OSDI*, 2016, pp. 185–201. 958
- 921 [23] M. Xiao, H. Wang, L. Geng, R. Lee, and X. Zhang, “Catfish: Adaptive RDMA-enabled R-tree for low latency and high throughput,” in *Proc. IEEE 39th ICDCS*, 2019, pp. 164–175. 959
- 922 [24] Q. Wang, Y. Lu, and J. Shu, “Sherman: A write-optimized distributed B+tree index on disaggregated memory,” in *Proc. SIGMOD*, 2022, pp. 1033–1048. 960
- 923 [25] “LevelDB.” 2023. [Online]. Available: <https://github.com/google/leveldb> 961
- 924 [26] A. Kalia, M. Kaminsky, and D. G. Andersen, “Design guidelines for high performance RDMA systems,” in *Proc. USENIX ATC*, 2016, pp. 437–450. 962
- 925 [27] Y. Luo, P. Jin, Q. Zhang, and B. Cheng, “TLBtree: A read/write-optimized tree index for non-volatile memory,” in *Proc. IEEE 37th ICDE*, 2021, pp. 1889–1894. 963
- 926 [28] Y. Luo, P. Jin, Z. Zhang, J. Zhang, B. Cheng, and Q. Zhang, “Two birds with one stone: Boosting both search and write performance for tree indices on persistent memory,” *ACM Trans. Embed. Comput. Syst.*, vol. 20, no. 5s, pp. 1–25, 2021. 964
- 927 965
- 928 966
- 929 967
- 930 968
- 931 969
- 932 970
- 933 971
- 934 972
- 935 973
- 936 974
- 975