Replicating Persistent Memory Key-Value Stores with Efficient RDMA Abstraction

Qing Wang, Youyou Lu, Jing Wang, and Jiwu Shu

Tsinghua University

Abstract

Combining persistent memory (PM) with RDMA is a promising approach to performant replicated distributed key-value stores (KVSs). However, existing replication approaches do not work well when applied to PM KVSs: 1) Using RPC induces software queueing and execution at backups, increasing request latency; 2) Using one-sided RDMA WRITE causes many streams of small PM writes, leading to severe device-level write amplification (DLWA) on PM.

In this paper, we propose Rowan, an efficient RDMA abstraction to handle replication writes in PM KVSs; it aggregates concurrent remote writes from different servers, and lands these writes to PM in a sequential (thus low DLWA) and one-sided (thus low latency) manner. We realize Rowan with off-the-shelf RDMA NICs. Further, we build Rowan-KV, a log-structured PM KVS using Rowan for replication. Evaluation shows that under write-intensive workloads, compared with PM KVSs using RPC and RDMA WRITE for replication, Rowan-KV boosts throughput by 1.22× and 1.39× as well as lowers median PUT latency by 1.77× and 2.11×, respectively, while largely eliminating DLWA.

1 Introduction

Replicated distributed key-value stores (KVSs) support many applications by providing durability and high availability [28,56,76]. The recent commercialization of persistent memory (PM), e.g., Intel's Optane DIMMs, enables local storage with extremely low latency (e.g., ~100ns when persisting small data [73]). When building replicated distributed KVSs with such fast storage media, network and CPU will become determinants of request latency, since replicating an object (i.e., key-value pair) involves multiple times of network communication and request queueing/execution.

RDMA, a widely-deployed network technology [34,37,53], is promising to mitigate the network and CPU overhead. First, RDMA delivers low latency ($\sim 2\mu$ s) due to protocol-offload RDMA NICs (RNICs) and kernel-bypass software. Second, RDMA provides one-sided WRITE and READ, allowing remote memory accesses without involvement of remote CPUs. Recent work have leveraged WRITE to replicate data in DRAM (i.e., WRITE-enabled replication) [17, 30, 31, 69]. This eliminates software queueing/execution of backups in the critical

This is the pre-print version of our OSDI'23 paper, which has been accepted after experiencing the revise-and-submit process of OSDI'22.

path, thus significantly cutting the replication latency compared with RPC-enabled replication.

Yet, in the context of PM KVSs, WRITE-enabled replication approach does not work well: it induces severe device-level write amplification (DLWA) on PM. Specifically, a KVS is typically finely sharded for load balancing and fast recovery, so every server acts as backups for many shards, receiving numerous concurrent replication writes from many remote threads; besides, these replication writes are typically small (~100B) due to prevalent tiny objects in real-world workloads [24, 52]. In WRITE-enabled replication approaches (e.g., FaRM [31]), each server allocates an exclusive backup log for every remote thread, to accommodate remote WRITE from primaries. When adopting WRITE-enabled replication to PM KVSs, these backup logs generate a huge number of PM write streams¹, which contain lots of small-sized writes. These numerous write streams lead to severe DLWA, since PM has block access granularity at media level (e.g., 256B in Optane DIMMs) and its hardware combining capacity is bounded. In our experiments, with 128B RDMA WRITE, 144 remote PM write streams cause 1.58× DLWA (§2.4). DLWA wastes limited PM write bandwidth, shortens PM lifetime, and harms PM's persistence efficiency.

In this paper, we propose Rowan, an efficient RDMA abstraction to handle replication writes on PM KVSs. Rowan can aggregate numerous concurrent remote writes from different servers, and land these writes to PM sequentially, so as to largely eliminate DLWA. Besides, it is one-sided as RDMA WRITE, enabling backup-passive replication with low latency and high CPU efficiency. We realize Rowan with offthe-shelf RNICs based on two observations: 1) RDMA SEND is two-sided on the control path but one-sided on the data path; 2) RNICs consume receive buffers in order. Thus, we let a control thread at the receiver side push PM-resident buffers into receive queues in increasing address order. Senders only need to issue SEND for remote PM writes and wait for ACKs generated by receiver-side RNICs. We leverage two RNIC hardware features, shared receive queue (SRQ) [11] and multipacket receive queue (MP RQ) [7,9], to merge writes from different connections and support variable-sized writes, respectively. We also streamline Rowan's control path by minimizing the control thread's tasks. A Rowan instance can

¹A write stream is a group of writes targeting contiguous addresses, e.g., writes that perform log appending.

achieve 54.5Mops/s for highly concurrent 64B remote PM writes, with almost no DLWA.

Further, we build Rowan-KV, a PM KVS leveraging Rowan for primary-backup replication. It adopts a log-structured approach to manage both local PM writes and remote PM writes. Specifically, each server maintains per-thread primary logs and a single backup log on PM. For a PUT request, a worker thread in servers generates a log entry containing the targeted object; then, it persists the log entry into its local primary log via CPU instructions and every backup's backup log via one-sided Rowan. For a GET request, the thread searches DRAM-resident indexes which point to objects in logs. In this way, Rowan-KV features high performance and low DLWA: 1) Replication bypasses CPUs of backups, ensuring low latency and saving CPU cycles for foreground operations; 2) The number of PM write streams in a server is small (i.e., n primary logs + 1 backup log, where n is local thread count), enabling efficient write combining in PM hardware and thus largely eliminating DLWA. Rowan-KV also introduces a failover mechanism for fault tolerance and a dynamic resharding mechanism for load balancing.

We evaluate Rowan-KV on Optane DIMMs under a cluster of 14 machines (8 clients and 6 servers). Our evaluation focuses on YCSB benchmarks [26] with object sizes from three typical Facebook KVSs workloads [24] (i.e., ZippyDB, UP2X and UDB). Compared with KVSs using RPC and WRITE for replication, Rowan-KV boosts throughput by 1.22× and 1.39×, lowers median PUT latency by 1.77× and 2.11×, and lowers 99% latency by 1.26× and 2.06×, respectively, under write-intensive workloads. In addition, the DLWA is less than 1.032× in Rowan-KV, while 1.54× in the WRITE-enabled KVSs. Under read-intensive workloads, they have similar performance. We also compare Rowan-KV with two software techniques mitigating DLWA, i.e., batching and log sharing; Rowan-KV still outperforms them.

In summary, this paper makes the following contributions:

- It demonstrates that WRITE-enabled replication can lead to severe device-level write amplification on PM KVSs.
- It introduces Rowan abstraction and Rowan-KV with goals of low latency and low device-level write amplification.
- It uses experiments to confirm the efficacy of Rowan-KV.

2 Background and Motivation

In this section, we first provide the background on PM (§2.1) and RDMA (§2.2). Then, we show that characteristics of typical KVSs architecture and workloads together lead to high fan-in small writes for replication (§2.3). Finally, with experiments, we demonstrate that when handling these writes, WRITE-enabled replication causes severe DLWA (§2.4).

2.1 Persistent Memory (PM)

PM is a new kind of storage device that sits on the memory bus. Thus, PM is byte-addressable and can be accessed by CPUs via load/store instructions. In this paper, we focus on Intel's Optane DIMM, the only available PM product.

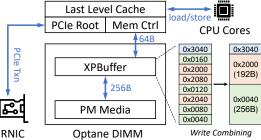


Figure 1: Architecture of Optane DIMMs and RNICs.

PM performance. Optane DIMMs have unique performance characteristics. In terms of bandwidth, an Optane DIMM offers about 2GB/s for writes and 6GB/s for reads, which are 1/6 and 1/3 of DRAM, respectively. In terms of latency, compared to DRAM, Optane DIMMs have similar write latency but $3 \times$ higher read latency [73]. The limited write bandwidth and high read latency of Optane DIMMs are the main design considerations for many PM systems [20,25,48,49,59,67,74]. PM architecture. Figure 1 presents the architecture of Optane DIMMs. The memory controller generates cache-line granularity (i.e., 64B) read/write requests to Optane DIMMs, but the internal PM media has a 256B access granularity (referred as XPLine in this paper). Such a granularity mismatch will trigger read-modify-write events, thus leading to devicelevel write amplification (DLWA). To mitigate DLWA, each Optane DIMM features an XPBuffer [73], which performs write combining for adjacent 64B writes, as shown in the right part of Figure 1. Yang et al. estimated that the XPBuffer in an Optane DIMM is approximately 16KB in size [73].

Persistent modes. There are two persistent modes for PM: ADR and eADR [36]. In ADR mode, once a store reaches the memory controller, it can survive power failure; but the CPU cache is volatile, so programmers must explicitly flush data from the CPU cache (using clwb or clflushopt instructions) or bypass the CPU cache (using ntstore instructions). In eADR mode, the CPU cache also belongs to the persistence domain: its data will be flushed to PM upon power failure.

2.2 Remote Direct Memory Access (RDMA)

RDMA is a network technology that offers high bandwidth (e.g., 100 Gbps) and low latency ($\sim 2\mu$ s).

Verb types. RDMA provides two types of verbs for network communication: *message verbs* and *memory verbs*. Message verbs, i.e., SEND and RECV, are the same as Linux socket interfaces: a SEND emits a message to a remote server that prepares receive buffers via RECV. Memory verbs include WRITE, READ and ATOMIC. These verbs can operate receivers' memory without involving receivers' CPUs. Due to the *one-sided* feature, memory verbs enjoy low latency and high CPU efficiency.

Queue pair. RDMA servers use queue pairs (QPs) for communication. A QP contains a *send queue* (SQ) and a *receive queue* (RQ). A server posts requests, including SEND, WRITE, READ, and ATOMIC, to the send queue, and posts RECV to the receive queue for accommodating incoming SEND messages. A send/receive queue is associated with a *completion queue*

	max shard size	# of backup shards (stored by one PM server)	
CosmosDB	20GB [10]	200	
DynamoDB	10GB [2]	400	
FoundationDB	500MB [3]	8,400	
Cassandra	100MB [1]	42,000	
TiKV	96MB [38]	43,000	

Table 1: A PM server hosts many backup shards for popular KVSs. We assume 3-way replication and a typical configuration of PM servers: 2 sockets, each with 3TB Optane DIMMs (6TB in total).

(CQ), which generates completion signals for posted verbs. **Remote persistence.** When issuing a WRITE to remote PM, to ensure the data persistence, we should take two extra actions. ① Since receiver-side RNICs return acknowledgements before data in WRITE is DMA-ed to PM, we should send a READ (1B in arbitrary addresses) to flush RNIC and PCIe buffers at the receiver side [42]. These two verbs (i.e., WRITE followed by READ) can be posted in one request according to the ordering guarantee of RDMA [70]. ② We should disable *Data Direct I/O (DDIO)* [5, 32], a technology of Intel CPUs that lets RNICs directly DMA data to last level cache (LLC). In ADR mode, disabling DDIO ensures that DMA-ed data can reach persistence domain. In eADR mode, it avoids PM write amplification resulting from LLC's near-random eviction (64B cache line vs. 256B XPLine) [42, 70].

2.3 High Fan-in Small Writes in KVSs

In KVSs, replication makes *high fan-in small writes* a dominant access pattern due to the following two reasons.

1) Data sharding. Distributed storage systems (including KVSs) typically split the entire data set into a large number of shards, and then distribute these shards across many servers [16, 50]. Each shard has multiple replicas, with one selected as *primary* and the others as *backups*. Data sharding has two advantages. First, it can improve load balancing and support dynamic data migration in a fine-grained manner. Second, it can improve availability: when a server fails, since replicas of its data are distributed to many servers, the system can perform recovery and re-replication in parallel. For example, FaRM [30] maps each server into 100 consistent hashing rings by default; in Facebook's RocksDB clusters, each server typically hosts tens or hundreds of shards [29].

With data sharding, each server acts as backups for tens or hundreds of shards, and their primaries are distributed to many servers. This makes every server receive messages for data replication, i.e., replication writes, from many primaries residing in many other servers. We call it *high fan-in writes*.

To solidify the argument of high fan-in writes in KVSs, we analyze five widely-used replicated KVSs. As shown in Table 1, these KVSs all have a maximum shard size, from tens of megabytes (i.e., Cassandra [1] and TiKV [38]) to several gigabytes (i.e., DynamoDB [2] and CosmosDB [10]). When we deploy these KVSs on servers having terabytes of PM, each server will host a considerable number of backup shards which ranges from 200 (CosmosDB) to 43,000 (TiKV),

generating high fan-in replication writes.

69]. To achieve multicore-scalable and squeeze out the raw performance of NICs, these systems run multiple threads, each independently processing requests using exclusive network connections. For example, in DrTM+H [69], every worker thread independently issues RDMA WRITE for replication. With this threading model, the degree of fan-in increases from the number of remote servers to the number of remote threads. 2) Numerous small-sized objects. Many important applications relying on KVSs generate numerous small objects, whose size is much smaller than the access granularity of PM media (e.g., 256B XPLine in Optane DIMMs). For example, in ZippyDB, the largest KVS at Facebook [15], the average size of objects is only 90.8B [24]. Moreover, the other two typical KVSs at Facebook — UP2X (a KVS for AI services) and UDB (a KVS for social graph) — have average object size of 57.25B and 153.8B, respectively [24]. Twitter exhibits a similar workload feature: the most common length of a tweet is only 33 characters [14, 52]. This paper focuses on these small objects because of their prevalence and importance.

The degree of fan-in is even higher in systems equipped with fast network hardware (e.g., RNIC) [30, 31, 44, 45, 61,

When a KVS handles PUT requests (from clients) for these small objects, primaries emit replication writes to associated backups. These writes are small, since they typically only contain replicated objects with tiny metadata [56]. These writes are also high fan-in due to data sharding, as explained before. As a result, we can conclude that high fan-in small writes are a dominant access pattern in the cluster of KVSs.

2.4 DLWA from WRITE-enabled Replication

Recent research demonstrates that for in-memory DRAM systems, compared with RPCs, leveraging RDMA WRITE for replication can obtain significant performance gain [17, 30, 31, 69]. In such WRITE-enabled replication, primaries issue replication writes to backups' logs via one-sided WRITE, and only need to wait for acknowledgements (ACKs) from the RNIC hardware of backups. This eliminates software queueing/execution of backups in the critical path, thus enjoying low latency (e.g., Mu [17] cuts the latency by 61%). Further, the saved CPU cycles in backups can serve requests (e.g., GET) from clients, thus improving system throughput.

In systems using WRITE-enabled replication, to handle high fan-in replication writes from many remote threads (recall §2.3), each server maintains lots of backup logs, each accommodating WRITE from an individual remote thread (which can act as primary) [31,69]. For example, in FaRM's evaluation with 90 machines (each running 30 worker threads) [31], there are thousands of backup logs (i.e., 89×30) in each server. Yet, when we apply WRITE-enabled replication to PM KVSs, these backup logs (which are placed in PM for durability) will cause a huge number of PM write streams, which contain lots of small writes, thus inducing severe DLWA. We conduct an experiment to demonstrate it.

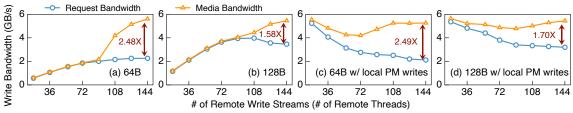


Figure 2: DLWA with varying remote write streams. *DLWA = media bandwidth/request bandwidth. Threads access PM on a remote server;* each thread generates a remote write stream. In (c) and (d), 18 CPU cores in the remote server perform local sequential PM writes.

In the experiment, we launch a number of threads (on four servers), each issuing sequential RDMA WRITE to an exclusive PM-resident log in a remote server and thus generating a PM write stream. We disable DDIO in the remote server and post a READ after each WRITE to guarantee persistence. The remote server is equipped with three 256GB Optane DIMMs and a 100Gbps RNIC. We use ipmctl [8] to periodically read hardware counters of Optane DIMMs, calculating request bandwidth and media bandwidth, which means write bandwidth received from memory bus and write bandwidth issued to PM media, respectively. Figure 2(a) and (b) show results with 64B and 128B WRITE size (representing small replication writes, §2.3), respectively. When remote write stream count is lower than 90, DLWA is negligible. This is because the XPBuffer on Optane DIMMs can combine adjacent small writes from the same write streams into 256B internal writes (§2.1). However, the capacity of combining is bounded due to the limited size of XPBuffer. Consequently, as the number of remote write streams continues to increase, severe DLWA appears. Specifically, when remote write stream count is 144, the DLWA is 2.48× and 1.58× in case of 64B WRITE and 128B WRITE, respectively.

Next, we consider a more practical scenario where local PM writes exist. In the remote server, we run 18 CPU cores, each performing sequential 128B PM writes using ntstore. We repeat the above experiment; Figure 2(c) and (d) show the results. Without remote RDMA WRITE, local PM writes can deliver high request bandwidth (i.e., available bandwidth). As the remote write stream count increases, DLWA in Optane DIMMs reaches 2.49× and 1.70× in case of 64B WRITE and 128B WRITE, respectively. In addition, the available bandwidth drops from 5.2GB/s to 2.1GB/s (60%) for 64B WRITE, and from 5.4GB/s to 3.2GB/s (41%) for 128B WRITE.

DLWA on PM leads to three issues. First, it reduces available PM write bandwidth, thus degrading system performance. The wasted bandwidth could also have been used for colocated applications [33,54,55]. Second, it shortens the lifetime of PM which has limited write endurance [6]. Third, severe DLWA consumes a considerable number of hardware resources (e.g., XPBuffer), harming persistence efficiency.

To efficiently handle high fan-in small writes, we need a new RDMA abstraction (rather than WRITE) for PM KVSs. This abstraction should *mitigate DLWA*, *while achieving benefits of one-sided verbs*—*low latency and high CPU efficiency*.

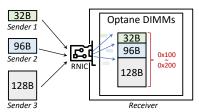


Figure 3: An instance of Rowan abstraction.

3 Rowan Abstraction

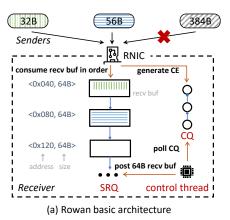
We propose Rowan, a new RDMA abstraction to handle high fan-in small writes in PM KVSs. In this section, we first describe Rowan's semantic and characteristics. Then, we present how to realize Rowan using off-the-shelf RNICs.

3.1 Rowan Semantic

Figure 3 presents a Rowan instance. A Rowan instance is associated with one receiver and a set of senders. Senders concurrently issue writes to the receiver which has registered a large PM area. The receiver-side RNIC lands these writes to the PM area *sequentially*, and finally returns ACKs to senders.

Rowan abstraction has the following advantages. First, by translating concurrent remote small writes into a single write stream, the XPBuffer in Optane DIMMs can easily combine them into 256B XPLine writes, largely eliminating DLWA. Second, since all the data operations are performed by the receiver-side RNIC without involving receiver-side CPUs, Rowan enjoys benefits of low latency and high CPU efficiency like RDMA WRITE. In addition, compared with CPUs, RNIC ASICs can deliver extremely high throughput.

Comparison with batching. Batching is also an approach that can mitigate DLWA on PM: it opportunistically accumulates multiple small writes at the sender side, and then emits the batched writes to the receiver via one RDMA WRITE. However, batching induces extra latency, sapping the benefits of extremely low-latency hardware (i.e., RNICs and PM). In contrast, Rowan does not delay any write and thus ensures low latency: senders immediately issue writes and receiver-side RNICs immediately land received writes to PM. In addition, as we will show in §6, batching frequently fails to accumulate enough small writes within a short time interval in KVSs, and Rowan outperforms batching in both latency and throughput. Our view of batching has been echoed by authors of RAMCloud — ". . . batching requires some operations to be delayed until a full batch has been collected, and this is not acceptable in a low-latency system such as RAMCloud" [56].



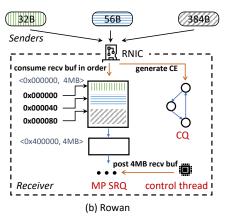


Figure 4: Realizing Rowan with off-the-shelf RNICs. (a) Rowan basic architecture using shared receive queue (SRQ). In this subfigure, the 32B write and 56B write are placed in the same XPLine. Yet, the 384B write fails to be received due to 64B receive buffers. (b) Rowan using multi-packet shared receive queue (MP SRQ) with 64B stride. In this subfigure, three writes are placed in two XPLines of the first receive buffer. We use a completion queue (CQ) ring to eliminate CQ polling in the control thread.

3.2 High-Performance Rowan

Rowan is conceptually simple but challenging to realize using off-the-shelf RNICs. We do not want to modify RNIC hardware like StRoM [60] and PRISM [23], so as to enable Rowan to be deployed immediately in datacenters today that are equipped with RNICs. Before describing our solution, we present a straightforward solution that has poor performance.

3.2.1 Straightforward Solution

A straightforward solution to realize Rowan abstraction is combining RDMA WRITE and atomic verb FETCH_AND_ADD. Specifically, there is a 64-bit sequencer stored in the receiver's memory. When performing a write, the sender first issues a FETCH_AND_ADD to the sequencer, reserving a PM address; then, it issues a WRITE to this address. This solution has two limitations. First, it needs two round trips, increasing the latency. Second, the poor performance of atomic verbs bottlenecks throughput: even storing the sequencer in RNICs' device memory [68], the throughput is less than 10Mops/s.

3.2.2 Our Solution

Counter-intuitively, we use RDMA SEND and RECV to realize Rowan. This is based on our two observations.

- RDMA SEND is two-sided on the control path but onesided on the data path. In the control path, the receiver's CPUs prepare receive buffers via RECV; however, in the data path, when handling SEND requests, the receiver-side RNIC performs all tasks, including landing SEND's data to receive buffers and returning ACKs.
- In a receive queue, receive buffers are consumed in order. Every time, the receiver-side RNIC pops the *first* buffer in the associated receive queue and lands data to it.

Key idea. On the control path, CPUs push PM buffers into the receive queue in increasing address order; on the data path, the receiver-side RNIC consumes them in order.

Basic architecture. Figure 4(a) shows the basic architecture of Rowan implementation. Rowan uses reliable connection (RC) QPs to delegate transmission reliability to RNICs. We create a shared receive queue (SRQ) [11] which is associated

with all QPs; thus, RNICs can land data of SEND from *different* remote QPs to the same receive queue. In the receiver, we reserve a dedicated thread, namely control thread, to perform control-path tasks; the RNIC performs data-path tasks.

Specifically, the control thread splits the PM area into fixed-sized (e.g., 64B in Figure 4)(a)) buffers, and posts these buffers (using RECV) into the SRQ in *increasing address order*. Senders encapsulate writes into SEND requests, and emit them to the receiver; each SEND is followed by a READ for persistence. When receiving a SEND (followed by a READ), the receiver-side RNIC pops the first buffer in SRQ, DMAs the SEND's data into the buffer, generates a completion entry (CE) to the SRQ's CQ, and finally returns an ACK to the sender. In this way, writes from different senders can be combined into the same XPLines on PM, mitigating DLWA.

Handling variable-sized writes. When the size of a SEND's data is larger than the first buffer in the SRQ, the RNIC cannot accommodate it and will trigger an error CE. For example, in Figure 4(a), with 64B receive buffers, the 384B write cannot be handled. If we use a buffer size larger than 256B for the SRQ to support relatively large writes, small writes from different senders will not be combined into the same XPLines, destroying the benefits of Rowan abstraction.

Fortunately, current RNICs (e.g., ConnectX-4/5/6) support a new type of RQ, called *multi-packet receive queue* [7,9] (MP RQ). In an MP RQ, each receive buffer can accommodate *multiple* SEND requests. We need to define a *stride* (e.g., 64B) for an MP RQ. When receiving a SEND, the RNIC appends the data to the receive buffer that is being used, and the start address is stride-aligned. If there is no enough space left, the RNIC pops a new receive buffer from the MP RQ to use.

Figure 4 shows Rowan that uses MP SRQ, where we set the stride to 64B and receive buffer size to 4MB. In the figure, three writes are placed in two XPLines (i.e., 512B area) in the first receive buffer, each having a 64B-aligned start address. By using MP SRQ, Rowan can support variable-sized writes, while combining small writes to mitigate DLWA.

There are two points worth noting when using MP SRQ:

- In Rowan, the stride is a fixed value of 64B. We do not choose a smaller value (e.g., 32B) for two reasons. First, in the RNIC we use (i.e., ConnectX-5), the minimum supported stride value is 64B. Second, recent studies suggest that senders should pad small writes to PCIe data word (64B) granularity [70], to avoid expensive read-modify-write operations on receivers' PM. Thus, we assume the incoming small writes are already 64B granularity.
- If a SEND is larger than maximum transmission unit (MTU), it is comprised of multiple packets. The RNIC may land these packets to non-contiguous addresses. We let the upper applications (e.g., KVSs) to handle this case.

Minimizing control-path tasks. On Rowan's data path, the receiver-side RNIC can deliver extremely high throughput (> 50Mops/s). On the control path, for CPU efficiency, we only want to use *one* control thread; thus, we minimize control-path tasks to make them can be easily handled by one thread.

There are two tasks performed by the control thread: posting receive buffers into the MP SRQ and polling the CQ to consume CEs. For the former, since we use large receive buffers (e.g., 4MB) by leveraging the multi-packet feature and post a batch of receive buffers at a time, this task is lightweight. For the latter, unfortunately, unlike other verbs, RECV can not be marked as unsignaled, so every SEND will generate a CE at the receiver side. The control thread cannot timely consume these CEs (considering > 50Mops/s throughput), making the CQ fill and thus causing QPs in an error state. We get inspiration from eRPC [43] to address this problem. Like eRPC, we create a CQ that forms a *ring structure*, so that the RNIC can overwrite entries in the CQ ring in a round-robin manner. In this way, the control thread does not need to poll the CQ.

4 Rowan-KV Design

We build Rowan-KV, a PM KVS that uses Rowan for primary-backup replication. It has two main design goals.

- Low latency. Rowan-KV exploits one-sided Rowan to eliminate software overhead at backups during replication.
- Low DLWA. Rowan-KV adopts a log-structured approach to manage PM writes from both local CPUs and remote CPUs. For the former, every thread appends data in its local log. For the latter, Rowan merges replication writes into a *single* backup log. Hence, Optane DIMMs only receive a small number of write streams and can efficiently combine adjacent small writes into XPLines, thus mitigating DLWA.

4.1 Overview

Figure 5 shows the architecture of Rowan-KV. Servers persistently store objects (i.e., key-value pairs) in PM and use RDMA for network communication. Rowan-KV divides the entire data set into many shards and distributes them across servers. Each shard is replicated for high availability: with the replication factor of k, it has one server as *primary* and k-1 servers as *backups*. Clients issue KV requests via RPCs. **Sharding mechanism.** Rowan-KV hashes each object's key

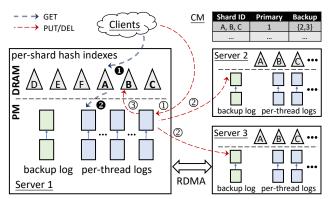


Figure 5: Architecture of Rowan-KV. *The per-thread logs (t-logs) and backup log (b-log) are divided into 4MB segments.*

into a 64-bit number and lets a shard manage a continuous range in the hashed keyspace. Shard distribution is maintained by a *configuration manager* (CM) and is cached in servers and clients. Rowan-KV uses a dynamic resharding mechanism to mitigate load imbalancing from overloaded servers (§4.6).

Log-structured approach. Rowan-KV adopts a log-structured approach, where each server has three components:

- Per-thread logs. Each server launches a number of worker threads to handle requests from clients. Each worker thread maintains a per-thread log (t-log) in PM, which stores objects of PUT/DEL requests. We do not allocate independent logs for each shard, to reduce random PM writes.
- Backup log. Each server has a single backup log (b-log) in PM, which receives replication writes from primaries using a Rowan instance. By doing so, Rowan-KV can largely eliminate DLWA from high fan-in small writes.
- Per-shard hash indexes. Each server builds a DRAM-resident hash table for every shard it manages, to index objects in t-logs or the b-log. Putting indexes in DRAM can avoid random PM writes and expensive PM reads [22, 25]. The t-logs and b-log are divided into 4MB segments.

Handling KV requests. When issuing a KV request for an object, the client sends an RPC to a worker thread residing in the server that is the targeted shard's primary.

For a PUT/DEL request, the worker thread generates a log entry containing the object (only the object's key for DEL), and persistently appends the log entry to its local t-log using ntstore instructions (① in Figure 5). Then, the worker thread issues replication write for every backup via one-sided Rowan, persistently appending the log entry to every backup's b-log (②). Upon receiving all ACKs from backups' RNICs, the worker thread updates the associated index to make the object (in t-logs) visible (③), and finally returns a response to the client. Rowan-KV has a strong durability guarantee: when a client receives the response of a PUT/DEL request, its effects have been persisted on all replicas.

For a GET request, the worker thread first locates the object by searching the associated index $(\mathbf{0})$. Then, it copies the object's value from t-logs $(\mathbf{2})$ and replies to the client.

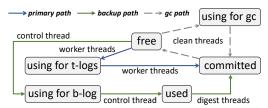


Figure 6: Life cycle of segments.

Background operations. Rowan-KV uses three types of threads to perform background operations.

- Control thread. One control thread performs control-path tasks for the Rowan instance (§3). In Rowan-KV, it pushes free segments to the b-log via RDMA RECV, and hands used segments over to digest threads.
- *Digest threads*. There are multiple digest threads. They digest used segments from the b-log. Specifically, they parse log entries and update associated indexes.
- Clean threads. There are multiple clean threads. They garbage collect stale objects in segments (from worker threads or digest threads) to reclaim free PM space.

4.2 Log Metadata

In Rowan-KV, t-logs and b-log are comprised of multiple segments, each storing a number of log entries. We describe segment metadata and log entry metadata, respectively.

4.2.1 Segment Metadata

A segment's metadata mainly includes its *state*. At any given time, each segment is in one of four states:

- *Free*. The segment can be allocated to t-logs by worker threads, the b-log by the control thread, or clean threads.
- *Using*. The segment is being used by t-logs, the b-log, or clean threads; it has space to store new log entries.
- *Used*. It has no space to store new log entries, and some of its log entries have *not* been persisted on all replicas.
- *Committed*. It has no space to store new log entries, and all of its log entries have been persisted on all replicas.

In addition to the state, a segment has an extra metadata called *owner*, indicating which type of thread allocates it (e.g., worker threads). Each server maintains a PM array called *segment meta table* to record metadata for all its segments.

Figure 6 presents the life cycle of segments. The path for primaries is simple: a worker thread allocates a *free* segment for its t-log, and the segment becomes *using* state. Once the segment has no space, it transitions into *committed*, since the worker thread can easily ensure that all of the segment's log entries have been persisted on all replicas. The path for backups is fairly complicated, where we should accurately distinguish between *used* segments and *committed* segments (§4.3 and §4.4). Such a distinguishment is essential for failover (§4.5).

4.2.2 Log Entry Metadata

A log entry contains the request type (i.e., PUT/DEL) and the targeted object (only the object's key for DEL). It also includes three metadata fields:

• 32-bit checksum. The checksum covers the whole log entry. Checksums eliminate persistent tails for logs: upon



Figure 7: A 2-MTU-sized log entry in the b-log.

recovery, we can identify the end of each log by calculating checksums. Besides, backups can use checksums to independently check the integrity of log entries in the b-log.

- 48-bit version. Each shard has a version, namely shard version, which is maintained by its primary. Upon a PUT/DEL request, the worker thread atomically increments the associated shard version, and stores the obtained version into the log entry. Upon recovery, the version allows us to identify the most recent objects from multiple t-logs.
- 16-bit shard ID. It indicates which shard the targeted object belongs to.

Handling larger-than-MTU log entries. For a log entry that is larger than MTU, backup-side RNICs may divide it into multiple packets and place them in non-contiguous addresses of the b-log (recall §3.2). To enable backups check the integrity of such a log entry, we design a simple counterbased metadata. Specifically, if a log entry is larger than MTU, we logically divide it into multiple MTU-sized blocks, and duplicate log entry metadata at the start of each block (each *checksum* field protects the individual block). Besides, we add two extra metadata to each block: 1) *cnt*: block count of the log entry, and 2) *seq*: the sequence number of the block.

Figure 7 shows a 2-MTU-sized log entry in the b-log, where its two blocks are not adjacent. The pair of $\langle shard\ ID : A, version : 64 \rangle$ uniquely identifies the log entry. When scanning the two blocks (checksums match) with their *cnt* and *seq*, backups can determine the log entry's integrity.

4.3 Managing the Backup Log

The control thread manages the b-log by communicating with the RNIC and digest threads. To minimize the communication overhead, the control thread performs tasks in a *batch* manner.

Specifically, when the system starts up, the control thread allocates a considerable number of free segments (e.g., 512) for the b-log, and pushes them into Rowan's MP SRQ via RECV. Then, it enters into a loop: 1) identifies a batch of segments (e.g., 128) that is in *used* state; 2) hands these segments over to digest threads; 3) allocates a batch of free segments and pushes them into the b-log via one RECV call. Note that a free segment transitions into the *using* state after it is allocated by the control thread (recall backup path in Figure 6).

Identifying used segments. The control thread adopts a simple method to identify used segments in the b-log. For every segment pushed into the b-log, its first 64 bits are set to zeros. Meanwhile, the first 64 bits in a log entry include the request type, which is non-zero. Thus, when the control thread finds that a segment has non-zero first 64 bits, it can ensure that all *previous* segments in the b-log (we call the set of segments S here) have been allocated by the RNIC for accommodating

log entries. However, this does not mean that segments in S are *used*, since maybe some DMA operations writing log entries in S are outstanding. Hence, we wait 2ms for all these DMA operations to complete, to guarantee that all segments in S have transitioned into *used* state. At the primary side, worker threads measure the time of each replication write: if it is more than 1ms, worker threads retry the replication write.

4.4 Digest and Garbage Collection

Digest. Multiple digest threads process used segments in the b-log in parallel. Each digest thread manages an exclusive set of shards: it extracts log entries from used segments in order and only processes shards it manages. For a log entry, digest threads update the index of the associated shard. Besides, digest threads identify *committed* segments, and hand these segments over to clean threads.

Identifying committed segments. To help digest threads identify committed segments in the b-log, primaries disseminate the information of log entries to backups. Specifically, for a shard, worker threads in its primary maintain a CommitVer; any log entry containing a version ≤ CommitVer has been persisted on all replicas. Every 15ms, worker threads write the ⟨shard ID, CommitVer⟩ pair into backups' b-logs via Rowan.

At the backup side, digest threads maintain an array CommitVerArray, which contains associated CommitVer for each shard. When encountering a \(\shard ID, CommitVer \) during parsing segments of the b-log, digest threads update CommitVerArray. Meanwhile, when processing a segment, digest threads generate an array MaxVerArray for it; for each shard, this array records the maximum version that digest threads have encountered in log entries. A used segment can transition into committed one, if its MaxVerArray \(\subseteq CommitVerArray \) (i.e., for every shard, the maximum version in MaxVerArray \(\subseteq CommitVer in CommitVerArray \)).

Garbage collection. Multiple clean threads garbage collect stale objects in committed segments. When memory utilization of a committed segment, i.e., the percentage of valid bytes, is lower than a pre-defined threshold (e.g., 75% in our evaluation), a clean thread cleans it. Specifically, the clean thread scans the committed segment and checks the liveness of objects in log entries (by searching indexes). For live objects, the clean thread copies associated log entries to a *using* segment and updates indexes. Finally, the committed segment transitions into *free* state for future usages.

4.5 Failover

We adopt FaRM's reconfiguration-style approach [31] to handle failover but tailor it for Rowan-KV. A *configuration* in Rowan-KV contains 1) 64-bit term, 2) membership, i.e., the set of live servers, and 3) shard distribution. The configuration is persistently stored in a Zookeeper instance [40], and is cached in the CM, clients, and servers. Rowan-KV uses leases to detect failures for servers and CM [31]. When the CM fails, Rowan-KV activates a new CM using the same mechanism as FaRM [31]. When a server fails, Rowan-KV performs failover

with the following three phases.

1) Generating and committing a new configuration. The CM generates a new configuration, where the term is incremented and the membership excludes the failed server. In the new shard distribution, the CM reassigns shards managed by the failed server to live servers, and promotes a backup to the new primary for each shard losing its primary.

Then, the CM stores the new configuration in Zookeeper and sends it to all servers. Servers cache the configuration, destroy QPs used for communicating with the failed server, and respond. From this point, servers block all requests from clients. Once the CM receives all responses, after ensuring that the lease for the failed server has expired, it sends a commit message to all servers. Now, servers can unblock requests. A server rejects requests containing terms that are lower than the one it caches. Clients will fetch the new configuration from CM upon receiving rejected responses.

- 2) Promoting backup to primary. When a backup of a shard (we call the shard A here) is promoted to the new primary, its worker threads block requests to A until digest threads build indexes for all objects of A. The new primary and backups should reach a consensus on the committed log entries. Hence, the new primary and backups process *using* and *used* segments in the b-log, collecting log entries belonging to A. These collected log entries are gathered to the new primary and then are scattered to backups. The new primary and backups store these log entries into segments. In this way, all replicas will own the same set of log entries for A. During digest, the new primary constructs a valid shard version for A, which is larger than versions in any A's log entry.
- **3) Re-replication.** The CM adds a new backup for the shard having replicas in the failed server. The new backup performs re-replication asynchronously. It first initializes an index for the shard, and then sends a message to the primary. Upon receiving the message, the primary traverses the shard's index and transmits associated log entries to the new backup.

4.6 Dynamic Resharding

ROWAN-KV introduces a dynamic resharding mechanism to migrate hotspot shards for improving load balancing.

CM detects overloaded servers and produces new shard distribution. Specifically, for each shard, each worker thread records the number of received requests during a fixed period (i.e., 500ms), and sends the statistic data to CM. Since Rowan is one-sided and thus backups are unaware of replication writes, we let worker threads in primaries record the number of received replication writes for backup shards. CM calculates the load of each server according to these statistics. If a server has a load that is higher than the average load by a threshold (i.e., 30%), CM determines that the server is overloaded. CM produces a new shard distribution, where the hottest shards in overloaded servers are moved to underloaded servers, with a goal of making the load of every server within 5% of the average. Then, it saves a migration list in the config-

uration, which contains a triple (source server, target server, shard ID) for each migration task. Finally, CM increments the term, writes the new configuration (including the new shard distribution) to Zookeeper, and sends it to all servers.

Next, we describe how Rowan-KV migrates a primary shard from a *source* server to a *target* server (migrating a backup shard is much easier since it does not serve client requests). Upon receiving the new configuration, servers cache it to local memory. From this point, the source server rejects client requests for the migrated shard. Clients will fetch the new configuration from CM when receiving rejected responses, so subsequent requests to the migrated shard will be sent to the target server. Then, the source server sends a message to the target server; the message contains the shard version of the migrated shard. Upon receiving both the message and the new configuration, the targeted server starts to serve requests for the migrated shard. In this way, Rowan-KV guarantees that only one server can serve the shard at any given time. Then, the process of data migration starts:

- In the source server, a migration thread requests free PM segments from the target server via RPCs, traverses the index of the migrated shard, and stores the associated log entries to remote segments via RDMA WRITE.
- In the target server, a migration thread scans segments written by the source server and installs log entries in the shard's index. Upon a PUT request to the migrated shard, the target server handles it as normal. Upon a GET request, the target server searches the index; if the corresponding key is not found, the target server routes the GET request to the source server since some objects have not been migrated yet. Of note, the versions in log entries resolve the conflicts between the migration thread and concurrent PUT requests.

The target server informs CM when it finishes data migration. Then, CM deletes the migration task from the migration list and writes the new configuration to Zookeeper. Finally, CM sends a message to the source server to inform it to free the index of migrated shard; the associated log entries in the source server will be removed by garbage collection.

If the migration is interrupted due to failures of the source/target server, the CM first rolls back the shard distribution in the configuration to the state before migration. Then, the CM deletes the associated task in the migration list and performs the normal failover process. In addition, the CM informs the target server (if alive) to release resources allocated for the interrupted migration task (e.g., migration thread and index).

4.7 Cold Start

When the entire cluster experiences a power failure, Rowan-KV can guarantee durability of data. Upon recovery, the CM fetches the configuration from Zookeeper, and disseminates it to all servers. Each server obtains the metadata for all its segments via the segment meta table (recall §4.2.1). For a shard, its primary extracts associated log entries from *using* segments whose *owner* is worker threads; then, the primary

sends these log entries to backups, to make all replicas own the same set of log entries. Each primary builds indexes for shards it manages by processing segments, and constructs valid shard versions. If two log entries have the same targeted key, the one with the larger version is more recent. Finally, Rowan-KV resumes unfinished migration tasks according to the migration list stored in the configuration.

5 Implementation

We implement Rowan-KV in Linux hosts. Rowan-KV is a fully user-space system: it uses *libibverbs* for RDMA operations and CPU memory instructions for accessing PM.

5.1 Threading Model

Rowan-KV binds each thread (i.e., worker threads, clean threads, digest threads, and control thread) to an exclusive CPU core. Rowan-KV follows two principles:

Minimizing inter-thread communication. First, each worker thread handles both network I/O and KV logic; this avoids request dispatch in systems that have dedicated threads to poll network requests [56], thus enjoying high multicore scalability. Second, a thread hands over segments to other threads in a batch manner (§4.3) using thread-safe queues. Avoiding thread blocking. To avoid blocking due to waiting for network events, worker threads adopt a coroutine-like approach to interleave work; after issuing Power operations

for network events, worker threads adopt a coroutine-like approach to interleave work: after issuing Rowan operations for a PUT, a worker thread saves the context of the PUT request (e.g., the targeted key); then, it polls the RDMA completion queue, getting new requests to execute. Upon receiving ACKs from backups, the worker thread restores the PUT's context and continues the remaining logic. In this way, a worker thread can concurrently handle multiple PUT requests.

5.2 Network Components

RPC. Rowan-KV uses an RPC framework for client-server and inter-server communication (not include replication). We build the RPC framework with RDMA SEND and RECV verbs using unreliable datagram (UD) QPs. Specifically, each worker thread creates a UD QP to receive requests and send responses. When a client joins the Rowan-KV cluster, it establishes RPC connections with a worker thread in every server. Like FaSST [44], our RPC framework currently does not support messages larger than an MTU. To reduce CPU consumption on PM reads: the RPC framework leverages RNICs' scatter-gather DMA to gather RPC headers and PM-resident objects, generating responses of GET requests.

Rowan. To realize Rowan, every worker thread builds a reliable connection (RC) QP with every remote control thread.

At the sender side, a worker thread uses the associated send queue in QPs to issue Rowan operations. A Rowan operation contains a SEND followed by a 1B READ for persistence (§3). SEND and READ are sent in one ibv_post_send call. For a worker thread, all its Rowan QPs and RPC QP share the same CQ, so that it can be aware of Rowan ACKs and new RPC messages by polling the CQ. We mark SEND as unsignaled to eliminate a completion event. For READ, we store the context

id of the associated PUT request (§5.1) into the wr_id field, so that worker threads can distinguish Rowan ACKs belonging to different PUT requests when polling the CQ.

At the receiver side, a control thread manages all Rowan QPs connected to remote worker threads; these QPs share an MP SRQ. The control thread pushes PM segments to the MP SRQ via RECV. We register PM to RNICs using physical addresses [64], to remove virtual-to-physical translation tables in RNICs and thus reduce cache thrash of RNICs.

Mitigating the impact of disabled DDIO. We disable DDIO to ensure the RNICs can land data to PM (rather than CPU cache). However, disabling DDIO will ① cause CPU cache miss when handling RPCs and ② degrade performance of DMA operations between RNICs and memory. For ①, worker threads poll multiple RPC messages at a time, and issue prefetch instructions to them. For ②, for RDMA READ used for persistence, we set its source address to RNICs' device memory [4, 68], to eliminate a DMA write at senders. We expect that DDIO does not need to be disabled, with nextgeneration RNICs supporting RDMA flush extensions [12].

5.3 Storage Components

PM management. We configure Optane DIMMs in App-Direct mode, which exposes PM as a range of physical memory. Rowan-KV splits the PM space into 4MB segments and stores the segment meta table in a predefined PM area (recall §4.2.1). A DRAM-resident free list records free segments, to serve segment allocation. We add padding for each log entry, making it 64B-aligned; it can ① avoid expensive PM readmodify-writes on receiver-side RNICs [70] when performing Rowan operations, and ② avoid slow repeated writes to the same cache lines [25, 42] in logs.

DRAM indexes. Each per-shard index is implemented with a concurrent bucket hash table [51]. The hash table is organized into a bucket array, where each bucket contains multiple 64-bit items. An item is composed of a 16-bit tag and a 48-bit PM address: the tag is a part of a key's hash value, to filter out mismatched searches and thus reduce PM reads; the PM address points to log entries. For a key, its targeted bucket is calculated by *hash(key)* % *sizeof(bucket array)*. If the targeted bucket is full when inserting a key, threads create a new free bucket and link it to the targeted bucket, forming a bucket chain. Indexes support conditional update to resolve conflicts between threads: indexes omit an update if its log entry has version that is smaller than the one indexes are pointing to.

6 Evaluation

6.1 Experimental Setup

Environment. We use 6 machines as servers and 8 machines as clients. Each machine is equipped with the Intel Xeon Gold 6240M CPU (18 physical/36 logical cores), 96GB DRAM, and one 100Gbps Mellanox ConnectX-5 RNIC. All machines are connected to a 100Gbps Mellanox IB switch. Each server machine owns three 256GB Optane DIMMs (ADR mode).

Unless otherwise specified, we run Rowan-KV on 6 servers.

In each server, we use 24 cores for worker threads, 5 cores for digest threads, 6 cores for clean threads, and 1 core for control thread. The control thread also manages leases, with a lease time of 10ms. The CM and Zookeeper instance (3-way replication) run on client machines. Each client machine runs multiple client threads to issue requests to servers. We set the replication factor to 3. Each server holds 48 shards.

Workloads. We evaluate Rowan-KV using YCSB [26] with different PUT:GET ratios: Load A — 100% PUT (write-only); A — 50% PUT and 50% GET (write-intensive); B — 5% PUT and 95% GET (read-intensive); C — 100% GET (read-only). Key distribution follows Zipfian with parameter 0.99 (default parameter in YCSB). We populate 200 million objects into KVSs before each experiment. We use three Facebook workloads [24] to generate object size: ZippyDB (for general data store) — 90.8B average object size; UP2X (for AI/ML services) — 57.25B average object size; UDB (for social graph) — 153.8B average object size.

Comparing targets. We compare Rowan-KV with four KVSs, each using a specific replication approach:

- RPC-KV. It uses RPC to perform replication. Each server
 maintains per-thread b-logs, and primaries issue replication
 writes via RPC. Upon receiving a replication RPC, the
 worker thread appends the log entry into its local b-log.
- **RWrite-KV**. It uses FaRM's approach [31] to perform replication. Each worker thread has an exclusive remote b-log at every remote server. During replication, the worker thread issues WRITE for appending log entries to its b-logs at backups. Each server stores (m-1)*n b-logs, where m is the number of servers and n is the worker thread count.
- Batch-KV. Batch-KV is a variant of RWrite-KV and uses WRITE for replication. Each worker thread generates large-sized WRITE requests to its remote b-logs by batching log entries, to mitigate DLWA. To reduce latency, worker threads immediately send batched log entries to backups once 1) the total size is larger than an XPLine, i.e., 256B, or 2) 5µs timeout is triggered.
- Share-KV. It is another variant of RWrite-KV and uses WRITE for replication. Worker threads in a server share the same remote b-log at a remote server, to reduce b-log count and thus mitigate DLWA. Worker threads use local atomic increment to obtain contiguous addresses in remote b-logs.

All systems are implemented in the same codebase (including optimizations in §5), to allow us to focus on the effects of replication approaches. By default, we disable DDIO to provide one-sided persistence. For RPC-KV, DDIO is enabled. We compare Rowan-KV with two off-the-shelf KVSs in §6.7.

6.2 Rowan Performance

We repeat the experiment in §2.4, to show performance of Rowan abstraction. Figure 8 presents the result of one Rowan instance. Rowan can largely eliminate DLWA in case of numerous concurrent remote small writes. The DLWA is less than 1.029× when no local PM writes exist (Figure 8(a) and

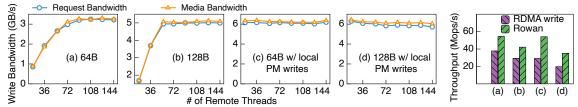


Figure 8: Rowan performance. *DLWA = media bandwidth/request bandwidth. A number of threads issue 64B/128B writes to a remote server's PM via a* Rowan instance. *In (c) and (d), 18 CPU cores in the remote server perform local sequential PM writes.*

(b)), and less than 1.056× when local PM writes exist (Figure 8(c) and (d))). This is because Rowan can merge remote small writes into a single write stream, enabling efficient hardware combining in Optane DIMMs' XPBuffer.

Further, we report the peak throughput of Rowan and RDMA WRITE under these four cases, as shown in the rightmost subfigure of Figure 8. When no local PM writes exist, Rowan can deliver 54.5 Mops/s for 64B remote PM writes and 42.2 Mops/s for 128B one, outperforming WRITE by 1.44× and 1.43×, respectively. When local PM writes appear, Rowan outperforms WRITE by 1.85×/1.78× for 64B/128B writes. Three causes make Rowan performant. First, Rowan largely eliminates DLWA, improving the available PM bandwidth. Second, on the data path of Rowan, all PM writes are performed by the receiver-side RNIC, ensuring high throughput. Finally, on the control path, by leveraging ring CQ and MP SRQ, the control thread only performs very lightweight tasks, so it does not become the bottleneck. Of note, the bottleneck of Rowan performance is 6GB/s PM write bandwidth in Figure 8(b)-(c), but processing capacity of RNICs in Figure 8(a). Rowan does not achieve 75Mops/s (a maximal message rate that a 100Gbps RNIC can provide), since we disable DDIO and send an extra RDMA READ for each Rowan operation.

6.3 Rowan-KV Performance

Figure 9 shows median latency and throughput (6 servers) under YCSB workloads with ZippyDB object size. Since Rowan-KV aims to accelerate replication, we report latency of PUT and GET separately. We increase the load generated by clients, and ensure that KVSs reach their peak throughput. We make two observations.

First, under read-only workloads (Figure 9(d)), RPC-KV has 5% higher throughput against other KVSs. This is because for KVSs using WRITE or Rowan, DDIO is disabled, lowering RPC performance. Such performance gap can be eliminated with next-generation RNICs supporting RDMA flush extensions [12]. Under read-intensive workloads (Figure 9(e) and (f)), RPC-KV and Rowan-KV have the similar throughput, since RPC-KV consumes CPU cycles of backups for 5% PUT requests, offsetting the benefits of DDIO. Compared with RPC-KV, Rowan-KV has 1.09× lower median PUT latency due to elimination of backups' software queueing, and 1.27× higher median GET latency due to disabled DDIO.

Second, under write-only and write-intensive (i.e., 50% PUT) workloads (Figure 9(a)-(c)), Rowan-KV has the highest throughput with the lowest median latency. We compare

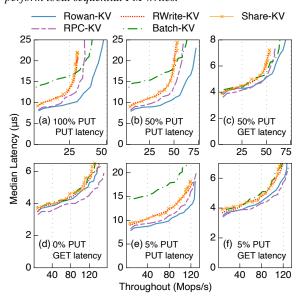


Figure 9: Median latency vs. throughput. *ZippyDB object size. We report PUT latency and GET latency separately.*

ROWAN-KV with the other four KVSs in turn.

With RPC-KV. Rowan-KV achieves peak throughput of 72.7/48.2Mops/s under write-intensive/write-only workloads, outperforming RPC-KV by 1.22×/1.37×. This is because Rowan-KV replicates log entries via one-sided Rowan, saving CPU cycles that handle replication RPCs. The saved CPU cycles can be used for primaries to handle RPCs from clients. At the peak throughput of RPC-KV, Rowan-KV has 1.77×/1.61× lower median PUT latency under write-intensive/write-only workloads. This is because compared with RPCs, one-sided Rowan eliminates backup-side software queueing/execution on the critical path of replication. Avoiding replication RPCs also makes Rowan-KV reduce median GET latency by 23%. Figure 10 shows DLWA of write-only and write-intensive workloads (6 servers). For Rowan-KV and RPC-KV, the DLWA is less than 1.032×. This is because they generate a small number of PM write streams: in each server, Rowan-KV has 24 t-logs and 1 b-log; RPC-KV has 24 t-logs and 24 b-logs (recall we use 24 worker threads in experiments). Optane DIMMs can efficiently combine adjacent small writes of the same logs into XPLine writes, when write stream count is not high (recall Figure 2(c) and (d)).

<u>With RWrite-KV.</u> Compared to RWrite-KV, Rowan-KV yields $1.39 \times /1.61 \times$ higher throughput and $2.06 \times /2.1 \times$ lower median PUT latency under write-intensive/write-only work-

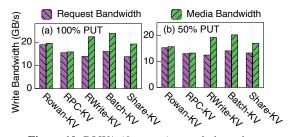


Figure 10: DLWA (6 servers) at peak throughput.

loads. The main culprit of RWrite-KV's low performance is DLWA: as shown in Figure 10(a), it suffers 1.54× DLWA.

This is because RWrite-KV owns lots of logs (i.e., 24×6 in experiments) in a server to accommodate small writes, exceeding the combining capacity of Optane DIMMs: a large number of write streams are equivalent to random writes. In RWrite-

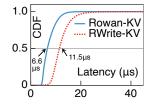


Figure 11: Latency CDF.

KV, Optane DIMMs trigger lots of read-modify-write events, which squander a considerable number of hardware resources (e.g., XPBuffer), degrading performance of PM accesses. To demonstrate it, we measure the latency of remote persistence operations of Rowan-KV and RWrite-KV under writeintensive workloads. Figure 11 shows the latency distribution. Remote persistence in RWrite-KV is slow (against Rowan-KV), with 11.5 μ s median latency and 24 μ s 99% latency. Of note, although RNICs are ideally capable of providing an RTT of $\sim 2\mu s$, the 6.6 μs median latency of Rowan is reasonable, since 1) we disable DDIO and each Rowan operation contains a synchronous RDMA READ, and 2) the latency is measured under high loads where RNICs suffer from DMA queueing. With Batch-KV. Batch-KV boosts the throughput of RWrite-KV by 1.23×/1.35× under write-intensive/write-only workloads, since it reduces the number of WRITE and mitigates DLWA (by 12%) via batching. However, batching makes Batch-KV suffer the highest PUT latency among all KVSs: even under low loads, Batch-KV has more than 50% higher PUT latency compared with Rowan-KV. In terms of throughput, Rowan-KV outperforms Batch-KV by 1.13×/1.19× under write-intensive/write-only workloads. This is because Batch-KV still experiences DLWA: it frequently fails to accumulate enough small writes within 5μ s timeout for two reasons: 1) All GET requests do not generate writes but consume CPU time; 2) Only writes to the same destination can be batched; yet, due to sharding of KVSs, for a server acting as primaries, the backups of its shards are distributed to *multiple* servers, greatly decreasing the batching opportunity. We also change the timeout value to 20µs, and Batch-KV delivers 9% lower throughput against Rowan-KV, with unacceptable latency. With Share-KV. Share-KV reduces DLWA of RWrite-KV by 26%/22% under write-intensive/write-only workloads, since it lets worker threads share the same b-logs. However, it still

	Rowan-KV	RPC-KV	RWrite-KV	Batch-KV	Share-KV
UP2X	73.9Mops/s	61.5Mops/s	56.2Mops/s	70.3Mops/s	56.0Mops/s
UDB	62.5Mops/s	50.4Mops/s	49.9Mops/s	57.1Mops/s	50.6Mops/s

Table 2: Throughput under write-intensive workloads.

suffers sizable DLWA (1.28×~1.39×), resulting in lower performance against Rowan-KV. This is because although worker threads in a Share-KV server generate contiguous remote addresses for WRITE, the asynchronous network makes receiverside RNICs receive and perform these writes in an out-of-order manner. In contrast, for Rowan-KV, leveraging Rowan, receiver-side RNICs decide destination addresses of writes. Besides, Rowan can merge writes from *different servers*.

Tail latency. Under write-intensive workloads with 50Mops/s throughput, Rowan-KV's 99% latency is 20.5μ s, which is $1.26\times$, $2.11\times$, $1.53\times$, and $1.87\times$ lower than that of RPC-KV, RWrite-KV, Batch-KV, and Share-KV, respectively.

Performance under uniform workloads. We evaluate Rowan-KV using uniform key distribution. Rowan-KV delivers 67.86Mops/s and 108.19Mops/s in cases of 50% PUT and 5% PUT, respectively, which are 6.6% and 15.5% slower than throughput of Zipfian skewed workloads (see Figure 9). Rowan-KV has higher performance under skewed workloads for two reasons. First, in our cluster of 6 servers, due to hash-based sharding, there is no observable load imbalancing across servers under skewed workloads. Second, threads enjoy much better cache locality under skewed workloads.

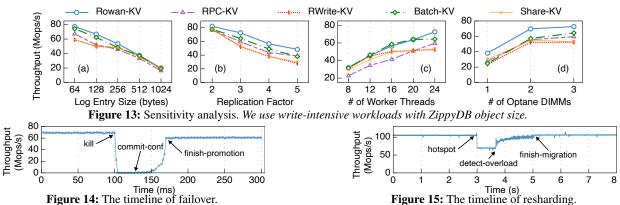
Performance with UP2X/UDB object size. Due to space limitations, here we only report the throughput under write-intensive workloads, as shown in Table 2. Rowan-KV delivers the highest throughput via powerful Rowan abstraction.

6.4 Sensitivity Analysis

We conduct experiments on sensitivity analysis using writeintensive workloads and ZippyDB object size.

Impact of object size. We change object size to generate varying log entry size. As shown in Figure 13(a), when log entry size is an integer multiple of XPLine size (e.g., 256B), all KVSs do not induce severe DLWA; thus, RWrite-KV and KVSs using WRITE have the similar throughput. RPC-KV consumes CPU cycles for replication RPCs, so it has 21% lower throughput against Rowan-KV with 1024B log entries. **Impact of replication factor.** Figure 13(b) presents throughput with varying replication factor. As replication factor increases, performance improvement between Rowan-KV and other KVSs increases. This is because, with higher replication factors, RPC-KV needs to consume more CPU cycles to handle a PUT request, and WRITE-enabled KVSs issue more WRITE and thus induce more DLWA. In contrast, Rowan-KV replicates objects in a one-sided manner and merges all remote writes into a single b-log in a sequential manner.

Impact of worker thread count. Figure 13(c) presents throughput with different worker thread counts. We make two observations. First, when the number of threads is small (i.e., ≤ 16), RPC-KV has the lowest throughput, since the CPU



is the bottleneck. Second, RWrite-KV and its variants yield poor scalability. This is because 1) for RWrite-KV and Batch-KV, the number of b-logs is proportional to thread count, and 2) for Share-KV, RNICs are more likely to receive and perform WRITE to b-logs in an out-of-order manner in case of high thread count; thus, they suffer more severe DLWA with higher thread count. DLWA harms PM performance (recall Figure 11), thus stalling throughput. In contrast, Rowan-KV exhibits superior throughput with different thread counts.

Impact of PM bandwidth. Figure 13(c) presents throughput with different number of Optane DIMMs per server. In case of one Optane DIMM, the PM bandwidth is bottleneck. Thus, RWrite-KV (which has the most severe DLWA) is outperformed by Rowan-KV, RPC-KV, Batch-KV, and Share-KV by 1.61×, 1.18×, 1.05×, and 1.28×, respectively. In case of three Optane DIMMs, CPU becomes the bottleneck and limits throughput, and PM bandwidth is not saturated (see Figure 10). Rowan-KV squeezes out CPU resources in two aspects: 1) it reduces CPU involvement via Rowan's one-sided semantic; 2) it largely eliminates DLWA, streamlining Optane DIMMs' internal operations and thus improving persistence efficiency of worker threads.

6.5 Failover and Cold Start

Failover. We kill a server to test Rowan-KV's failover mechanism. We use write-intensive workloads with ZippyDB objects and Rowan-KV runs for 50 seconds before the test. Figure 14 shows the timeline, where throughput is recorded per 2ms. The server is killed at time 100ms (i.e., kill). Rowan-KV uses 26ms to commit the new configuration (i.e., commit-config), which mainly includes detecting failure (8ms), writing new configuration to Zookeeper (4.3ms), and waiting for the failed server's lease to expire (10ms). Then, Rowan-KV consumes about 44ms to promote backups to primaries (i.e., finish-promotion). At this point, Rowan-KV can serve all requests from clients.

Cold start. We test cold start of a Rowan-KV instance, which contains 10 billion ZippyDB objects and thus occupies about 3TB PM space (6 servers). The time of cold start is 49.6s. Although cold start is slow, it is not common in datacenters. Periodically checkpointing DRAM-resident indexes can accelerate cold start, and we leave it for future work.

6.6 Dynamic Resharding

In this experiment, we evaluate Rowan-KV's dynamic resharding mechanism. We use read-intensive workloads with ZippyDB objects. Figure 15 presents the total throughput (6) servers) over time. At first, clients generate a uniform workload and each server has a similar CPU utilization (i.e., 90.2% $\sim 90.9\%$). At time 3s (i.e., hotspot), clients shift 80% requests for server A to a shard residing on server B, to make server B have a hotspot shard and overloaded. The throughput drops by 33% due to load imbalancing. At this time, server A and server B have a CPU utilization of 60.7% and 91% respectively. The average CPU utilization of the other 4 servers drops to 72.8%, since requests to overloaded server B suffer from long queueing and thus the limited number of clients cannot generate enough requests to other servers. CM detects the overload after 660ms (i.e., detect-overload) and produces a migration task that migrates the hotspot shard from server B to server A. The migration takes 1346ms and moves about 1.1 million objects. The throughput increases as the migration proceeds, since more GET requests to the hotspot shard can be served by server A. Finally, the system achieves a load-balanced state with steady throughput.

6.7 Comparison with Other Systems

We compare Rowan-KV with two state-of-the-art replicated KVSs designed for RDMA networks:

- Clover [63]. Clover runs on disaggregated PM, where PM servers do not have compute resources. Clients perform GET operations via RDMA READ verbs, and perform PUT operations (including replication) using a combination of RDMA WRITE and ATOMIC.
- HermesKV [45]. It is a DRAM-resident KVS built on Hermes [45], a broadcast-based replication protocol. HermesKV uses RPC for all inter-server communication (including replication). We modify the code to support PM: we store objects in PM and issue ntstore instructions for durability; indexes are in DRAM. In addition, we let clients generate KV requests to HermesKV servers.

We use ZippyDB objects and 4KB objects to test KVSs under small writes and large writes, respectively. The key distribution follows Zipfian with parameter 0.99. The replication factor is 3 and HermesKV runs with enabled DDIO.

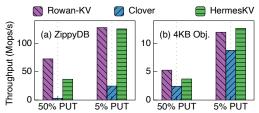


Figure 16: Comparison with Clover and HermesKV. (a) Throughput with ZippyDB objects. (b) Throughput with 4KB objects.

Figure 16(a) shows the results of small writes (ZippyDB objects). Under write-intensive workloads (i.e., 50% PUT), Rowan-KV outperforms Clover and HermesKV by 24.5× and 1.98×, respectively. Two reasons contribute to Clover's low throughput. First, due to the disaggregated architecture, every operation in Clover needs multiple network communications. Second, Clover uses RDMA ATOMIC to resolve conflicts between client threads, which leads to significant performance degradation when contention appears [35]. Using RDMA ATOMIC on PM is also considered slow due to its read-modify-write behavior [70]. HermesKV uses RPC for replication which consumes CPU cycles at backups, so it is outperformed by Rowan which uses one-sided Rowan for replication. We measure DLWA of these KVSs. Clover has 1.86× DLWA and HermesKV has 2.95× DLWA, since both of them generate a large number of random small writes on PM: for a PUT operation, Clover performs copy-on-write using WRITE and HermesKV performs in-place updates. In contrast, Rowan-KV adopts the log-structured approach to manage objects and exploits Rowan abstraction to minimize the number of write streams; thus, DLWA of Rowan-KV is less than 1.032×. Under read-intensive workloads (i.e., 5% PUT), Rowan-KV and HermesKV have similar throughput, which far exceeds that of Clover (about $5\times$).

Figure 16(b) reports the results of large writes (4KB objects). Under write-intensive workloads, Rowan-KV outperforms HermesKV by 1.42× and is bottlenecked by PM write bandwidth. HermesKV can not approach the limitation of PM write bandwidth, since its backups waste lots of CPU cycles to copy/persist large objects from RPC buffers to PM. Under read-intensive workloads, Rowan-KV and HermesKV are bottlenecked by the network bandwidth (11GB/s per server), which is much lower than PM read bandwidth (18GB/s).

7 Discussion

Although Intel killed Optane memory business for commercial reasons in summer 2022, we believe that Rowan is still applicable to future byte-addressable storage devices. For example, CXL storage devices (e.g., Samsung's Memory-Semantic SSD [13]), which are considered promising alternatives to Optane DIMMs, share similarities with Optane DIMMs: 1) limited write bandwidth; 2) byte interfaces with a block-level internal access granularity (e.g., flash page). Thus, when many remote threads concurrently access CXL storage devices with small IO size, Rowan can still effectively mitigate DLWA and thus boost system performance.

8 Related Work

PM KVSs. There are a host of works on PM KVSs, but most of them are single-machine (except Clover [63]). HiKV [71] and Bullet [39] are designed before the availability of real PM devices; both of them store objects into fine-grained PM hash tables. However, real PM devices have block-level internal access granularity (e.g., 256B in Optane DIMMs). To reduce DLWA, recent PM KVSs, including FlatStore [25], Viper [22], and Pacman [66], adopt log-structured approaches to manage objects. Rowan-KV also uses log-structured approach for the same reason, but focuses on distributed environments where objects are sharded and replicated.

RDMA replication. RDMA replication can be categorized into two groups, namely *backup-active* and *backup-passive*, depending on whether backups consume CPUs on the critical path of replication. Lots of systems [20, 21, 41, 45, 65, 75] belong to backup-active group, where backups' CPUs need to process messages during replication. For backup-passive group [17, 31, 46, 47, 57, 62, 69], primaries only need to wait for ACKs from the RNIC hardware of backups. For example, Hyperloop [47] uses RDMA WAIT and WRITE verbs to realize chain replication. Rowan-KV belongs to the backup-passive group, so it features low latency and high CPU efficiency. Yet, traditional backup-passive approaches can lead to DLWA on PM KVSs, driving us to design the Rowan abstraction.

RDMA abstraction. Due to limited expressivity of RDMA verbs, several works propose new RDMA abstractions [18, 19, 23, 27, 60, 72]. StRoM [60] and RMC [19] allow applications to define functions on NICs. Aguilera et al. [18] and PRISM [23] propose several new RDMA verbs to support far memory data structures and distributed systems. RedN [58] makes RDMA Turing complete using self-modifying chains. All above works (except RedN) require RNIC modification or specialized hardware (e.g., SmartNICs). In contrast, Rowan can be realized with off-the-shelf RNICs, leveraging RNIC features such as SRQ and MP RQ. Besides, Rowan targets handling high fan-in small PM writes.

9 Conclusion

This paper explored how to efficiently replicate PM KVSs using RDMA. We showed that existing approaches using RDMA WRITE cause severe device-level write amplification (DLWA) on PM. To this end, we proposed Rowan, a one-sided RDMA abstraction that can merge numerous remote writes into a single stream. Based on Rowan, we built Rowan-KV, a log-structured PM KVS; it outperforms RPC and RDMA WRITE alternatives in throughput and latency under write-intensive workloads, while achieving low DLWA.

References

- [1] Apache Cassandra Data Partitioning. https://www.instaclustr.com/blog/cassandra-data-partitioning/, 2022.
- [2] Choosing the Right DynamoDB Partition Key.

- https://aws.amazon.com/en/blogs/database/choosing-the-right-dynamodb-partition-key/, 2022.
- [3] Data Distribution and Movement in FoundationDB. https://github.com/apple/foundationdb/wiki/Data-Distribution-and-Movement, 2022.
- [4] Device Memory of RNICs. https://man7.org/linux/man-pages/man3/ibv_alloc_dm.3.html, 2022.
- [5] Intel Data Direct I/O Technology. https://www.intel.com/content/www/us/en/io/data-direct-i-o-technology-brief.html, 2022.
- [6] Intel Optane DC Persistent Memory Module (PMM). https://www.storagereview.com/news/inteloptane-dc-persistent-memory-module-pmm/, 2022.
- [7] Introduce Verbs API for Multi-packet Work Request. https://marc.info/?l=linux-rdma&m=151311334131294&w=2,2022.
- [8] ipmctl. https://github.com/intel/ipmctl, 2022.
- [9] Multi-Packet RQ. https://docs.mellanox.com/ display/rdmacore50/Multi-Packet+RQ, 2022.
- [10] Partitioning and horizontal scaling in Azure Cosmos DB. https://docs.microsoft.com/en-us/azure/cosmos-db/partitioning-overview, 2022.
- [11] RDMA Aware Networks Programming User Manual. https://www.mellanox.com/related-docs/prod_software/RDMA_Aware_Programming_user_manual.pdf, 2022.
- [12] RDMA Verbs Extensions for Persistency and Consistency. https://www.snia.org/sites/default/files/SDC/2016/presentations/persistent_memory/IdanBurstein_RDMA_VERBs_Extensions.pdf, 2022.
- [13] Samsung Electronics Unveils Far-Reaching, Next-Generation Memory Solutions at Flash Memory Summit 2022. https://news.samsung.com/global/samsung-electronics-unveils-far-reaching-next-generation-memory-solutions-at-flash-memory-summit-2022, 2022.
- [14] Twitter's doubling of character count from 140 to 280 had little impact on length of tweets. https://techcrunch.com/2018/10/30/twitters-doubling-of-character-count-from-140-to-280-had-little-impact-on-length-of-tweets/, 2022.

- [15] ZippyDB: the Architecture of Facebook's Strongly Consistent Key-Value Store. https://www.infoq.com/news/2021/09/facebook-zippydb/, 2022.
- [16] Atul Adya, Daniel Myers, Jon Howell, Jeremy Elson, Colin Meek, Vishesh Khemani, Stefan Fulger, Pan Gu, Lakshminath Bhuvanagiri, Jason Hunter, Roberto Peon, Larry Kai, Alexander Shraer, Arif Merchant, and Kfir Lev-Ari. Slicer: Auto-Sharding for Datacenter Applications. In 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16), pages 739–753, Savannah, GA, November 2016. USENIX Association.
- [17] Marcos K. Aguilera, Naama Ben-David, Rachid Guerraoui, Virendra J. Marathe, Athanasios Xygkis, and Igor Zablotchi. Microsecond Consensus for Microsecond Applications. In 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20), pages 599–616. USENIX Association, November 2020.
- [18] Marcos K. Aguilera, Kimberly Keeton, Stanko Novakovic, and Sharad Singhal. Designing Far Memory Data Structures: Think Outside the Box. In *Proceedings of the Workshop on Hot Topics in Operating Systems*, HotOS '19, page 120–126, New York, NY, USA, 2019. Association for Computing Machinery.
- [19] Emmanuel Amaro, Zhihong Luo, Amy Ousterhout, Arvind Krishnamurthy, Aurojit Panda, Sylvia Ratnasamy, and Scott Shenker. Remote Memory Calls. In Proceedings of the 19th ACM Workshop on Hot Topics in Networks, HotNets '20, page 38–44, New York, NY, USA, 2020. Association for Computing Machinery.
- [20] Thomas E. Anderson, Marco Canini, Jongyul Kim, Dejan Kostić, Youngjin Kwon, Simon Peter, Waleed Reda, Henry N. Schuh, and Emmett Witchel. Assise: Performance and Availability via Client-local NVM in a Distributed File System. In 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20), pages 1011–1027. USENIX Association, November 2020.
- [21] Jonathan Behrens, Sagar Jha, Ken Birman, and Edward Tremel. RDMC: A Reliable RDMA Multicast for Large Objects. In 2018 48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), pages 71–82, 2018.
- [22] Lawrence Benson, Hendrik Makait, and Tilmann Rabl. Viper: An Efficient Hybrid PMem-DRAM Key-Value Store. *Proc. VLDB Endow.*, 14(9):1544–1556, May 2021.
- [23] Matthew Burke, Sowmya Dharanipragada, Shannon Joyner, Adriana Szekeres, Jacob Nelson, Irene Zhang,

- and Dan R. K. Ports. PRISM: Rethinking the RDMA Interface for Distributed Systems. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, SOSP '21, page 228–242, New York, NY, USA, 2021. Association for Computing Machinery.
- [24] Zhichao Cao, Siying Dong, Sagar Vemuri, and David H.C. Du. Characterizing, Modeling, and Benchmarking RocksDB Key-Value Workloads at Facebook. In 18th USENIX Conference on File and Storage Technologies (FAST 20), pages 209–223, Santa Clara, CA, February 2020. USENIX Association.
- [25] Youmin Chen, Youyou Lu, Fan Yang, Qing Wang, Yang Wang, and Jiwu Shu. FlatStore: An Efficient Log-Structured Key-Value Storage Engine for Persistent Memory. In Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '20, page 1077–1091, New York, NY, USA, 2020. Association for Computing Machinery.
- [26] Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. Benchmarking Cloud Serving Systems with YCSB. In *Proceedings of the 1st* ACM Symposium on Cloud Computing, SoCC '10, page 143–154, New York, NY, USA, 2010. Association for Computing Machinery.
- [27] Alexandres Daglis, Dmitrii Ustiugov, Stanko Novaković, Edouard Bugnion, Babak Falsafi, and Boris Grot. SABRes: Atomic Object Reads for in-Memory Rack-Scale Computing. In *The 49th Annual IEEE/ACM Inter*national Symposium on Microarchitecture, MICRO-49. IEEE Press, 2016.
- [28] Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall, and Werner Vogels. Dynamo: Amazon's Highly Available Key-Value Store. In Proceedings of Twenty-First ACM SIGOPS Symposium on Operating Systems Principles, SOSP '07, page 205–220, New York, NY, USA, 2007. Association for Computing Machinery.
- [29] Siying Dong, Andrew Kryczka, Yanqin Jin, and Michael Stumm. Evolution of Development Priorities in Keyvalue Stores Serving Large-scale Applications: The RocksDB Experience. In 19th USENIX Conference on File and Storage Technologies (FAST 21), pages 33–49. USENIX Association, February 2021.
- [30] Aleksandar Dragojević, Dushyanth Narayanan, Miguel Castro, and Orion Hodson. FaRM: Fast Remote Memory. In 11th USENIX Symposium on Networked Systems Design and Implementation (NSDI 14), pages 401–414, Seattle, WA, April 2014. USENIX Association.

- [31] Aleksandar Dragojević, Dushyanth Narayanan, Edmund B. Nightingale, Matthew Renzelmann, Alex Shamis, Anirudh Badam, and Miguel Castro. No Compromises: Distributed Transactions with Consistency, Availability, and Performance. In *Proceedings of the 25th Symposium on Operating Systems Principles*, SOSP '15, page 54–70, New York, NY, USA, 2015. Association for Computing Machinery.
- [32] Alireza Farshin, Amir Roozbeh, Gerald Q. Maguire Jr., and Dejan Kostić. Reexamining Direct Cache Access to Optimize I/O Intensive Applications for Multihundred-gigabit Networks. In 2020 USENIX Annual Technical Conference (USENIX ATC 20), pages 673– 689. USENIX Association, July 2020.
- [33] Joshua Fried, Zhenyuan Ruan, Amy Ousterhout, and Adam Belay. Caladan: Mitigating Interference at Microsecond Timescales. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 281–297. USENIX Association, November 2020
- [34] Yixiao Gao, Qiang Li, Lingbo Tang, Yongqing Xi, Pengcheng Zhang, Wenwen Peng, Bo Li, Yaohui Wu, Shaozong Liu, Lei Yan, Fei Feng, Yan Zhuang, Fan Liu, Pan Liu, Xingkui Liu, Zhongjie Wu, Junping Wu, Zheng Cao, Chen Tian, Jinbo Wu, Jiaji Zhu, Haiyong Wang, Dennis Cai, and Jiesheng Wu. When Cloud Storage Meets RDMA. In 18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21), pages 519–533. USENIX Association, April 2021.
- [35] Stewart Grant and Alex C Snoeren. In-network Contention Resolution for Disaggregated Memory.
- [36] Shashank Gugnani, Arjun Kashyap, and Xiaoyi Lu. Understanding the Idiosyncrasies of Real Persistent Memory. *Proc. VLDB Endow.*, 14(4):626–639, December 2020.
- [37] Chuanxiong Guo, Haitao Wu, Zhong Deng, Gaurav Soni, Jianxi Ye, Jitu Padhye, and Marina Lipshteyn. RDMA over Commodity Ethernet at Scale. In *Proceedings of* the 2016 ACM SIGCOMM Conference, SIGCOMM '16, page 202–215, New York, NY, USA, 2016. Association for Computing Machinery.
- [38] Dongxu Huang, Qi Liu, Qiu Cui, Zhuhe Fang, Xiaoyu Ma, Fei Xu, Li Shen, Liu Tang, Yuxing Zhou, Menglong Huang, Wan Wei, Cong Liu, Jian Zhang, Jianjun Li, Xuelian Wu, Lingyu Song, Ruoxi Sun, Shuaipeng Yu, Lei Zhao, Nicholas Cameron, Liquan Pei, and Xin Tang. TiDB: A Raft-Based HTAP Database. *Proc. VLDB Endow.*, 13(12):3072–3084, August 2020.

- [39] Yihe Huang, Matej Pavlovic, Virendra Marathe, Margo Seltzer, Tim Harris, and Steve Byan. Closing the Performance Gap Between Volatile and Persistent Key-Value Stores Using Cross-Referencing Logs. In 2018 USENIX Annual Technical Conference (USENIX ATC 18), pages 967–979, Boston, MA, July 2018. USENIX Association.
- [40] Patrick Hunt, Mahadev Konar, Flavio P. Junqueira, and Benjamin Reed. ZooKeeper: Wait-Free Coordination for Internet-Scale Systems. In *Proceedings of the* 2010 USENIX Conference on USENIX Annual Technical Conference, USENIXATC'10, page 11, USA, 2010. USENIX Association.
- [41] Sagar Jha, Jonathan Behrens, Theo Gkountouvas, Matthew Milano, Weijia Song, Edward Tremel, Robbert Van Renesse, Sydney Zink, and Kenneth P. Birman. Derecho: Fast State Machine Replication for Cloud Services. ACM Trans. Comput. Syst., 36(2), apr 2019.
- [42] Anuj Kalia, David Andersen, and Michael Kaminsky. Challenges and Solutions for Fast Remote Persistent Memory Access. In *Proceedings of the 11th ACM Symposium on Cloud Computing*, SoCC '20, page 105–119, New York, NY, USA, 2020. Association for Computing Machinery.
- [43] Anuj Kalia, Michael Kaminsky, and David Andersen. Datacenter RPCs can be General and Fast. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*, pages 1–16, Boston, MA, February 2019. USENIX Association.
- [44] Anuj Kalia, Michael Kaminsky, and David G. Andersen. FaSST: Fast, Scalable and Simple Distributed Transactions with Two-Sided (RDMA) Datagram RPCs. In Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation, OSDI'16, page 185–201, USA, 2016. USENIX Association.
- [45] Antonios Katsarakis, Vasilis Gavrielatos, M.R. Siavash Katebzadeh, Arpit Joshi, Aleksandar Dragojevic, Boris Grot, and Vijay Nagarajan. Hermes: A Fast, Fault-Tolerant and Linearizable Replication Protocol. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '20, page 201–217, New York, NY, USA, 2020. Association for Computing Machinery.
- [46] Mikhail Kazhamiaka, Babar Memon, Chathura Kankanamge, Siddhartha Sahu, Sajjad Rizvi, Bernard Wong, and Khuzaima Daudjee. Sift: Resource-Efficient Consensus with RDMA. In Proceedings of the 15th International Conference on Emerging Networking Experiments And Technologies, CoNEXT '19, page

- 260–271, New York, NY, USA, 2019. Association for Computing Machinery.
- [47] Daehyeok Kim, Amirsaman Memaripour, Anirudh Badam, Yibo Zhu, Hongqiang Harry Liu, Jitu Padhye, Shachar Raindel, Steven Swanson, Vyas Sekar, and Srinivasan Seshan. Hyperloop: Group-Based NIC-Offloading to Accelerate Replicated Transactions in Multi-Tenant Storage Systems. In Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication, SIGCOMM '18, page 297–312, New York, NY, USA, 2018. Association for Computing Machinery.
- [48] Wonbae Kim, Chanyeol Park, Dongui Kim, Hyeongjun Park, Young ri Choi, Alan Sussman, and Beomseok Nam. ListDB: Union of Write-Ahead Logs and Persistent SkipLists for Incremental Checkpointing on Persistent Memory. In 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22), pages 161–177, Carlsbad, CA, July 2022. USENIX Association.
- [49] R. Madhava Krishnan, Wook-Hee Kim, Xinwei Fu, Sumit Kumar Monga, Hee Won Lee, Minsung Jang, Ajit Mathew, and Changwoo Min. TIPS: Making Volatile Index Structures Persistent with DRAM-NVMM Tiering. In 2021 USENIX Annual Technical Conference (USENIX ATC 21), pages 773–787. USENIX Association, July 2021.
- [50] Sangmin Lee, Zhenhua Guo, Omer Sunercan, Jun Ying, Thawan Kooburat, Suryadeep Biswal, Jun Chen, Kun Huang, Yatpang Cheung, Yiding Zhou, Kaushik Veeraraghavan, Biren Damani, Pol Mauri Ruiz, Vikas Mehta, and Chunqiang Tang. Shard Manager: A Generic Shard Management Framework for Geo-Distributed Applications. In *Proceedings of the ACM SIGOPS 28th Sympo*sium on Operating Systems Principles, SOSP '21, page 553–569, New York, NY, USA, 2021. Association for Computing Machinery.
- [51] Hyeontaek Lim, Dongsu Han, David G. Andersen, and Michael Kaminsky. MICA: A holistic approach to fast In-Memory Key-Value storage. In 11th USENIX Symposium on Networked Systems Design and Implementation (NSDI 14), pages 429–444, Seattle, WA, April 2014. USENIX Association.
- [52] Sara McAllister, Benjamin Berg, Julian Tutuncu-Macias, Juncheng Yang, Sathya Gunasekar, Jimmy Lu, Daniel S. Berger, Nathan Beckmann, and Gregory R. Ganger. Kangaroo: Caching Billions of Tiny Objects on Flash. In Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles, SOSP '21, page 243–262, New York, NY, USA, 2021. Association for Computing Machinery.

- [53] Maxim Naumov, John Kim, Dheevatsa Mudigere, Srinivas Sridharan, Xiaodong Wang, Whitney Zhao, Serhat Yilmaz, Changkyu Kim, Hector Yuen, Mustafa Ozdal, et al. Deep learning training in facebook data centers: Design of scale-up and scale-out systems. arXiv preprint arXiv:2003.09518, 2020.
- [54] Jinyoung Oh and Youngjin Kwon. Persistent Memory Aware Performance Isolation with Dicio, page 97–105. Association for Computing Machinery, New York, NY, USA, 2021.
- [55] Amy Ousterhout, Joshua Fried, Jonathan Behrens, Adam Belay, and Hari Balakrishnan. Shenango: Achieving High CPU Efficiency for Latency-Sensitive Datacenter Workloads. In *Proceedings of the 16th USENIX Conference on Networked Systems Design and Implementation*, NSDI'19, page 361–377, USA, 2019. USENIX Association.
- [56] John Ousterhout, Arjun Gopalan, Ashish Gupta, Ankita Kejriwal, Collin Lee, Behnam Montazeri, Diego Ongaro, Seo Jin Park, Henry Qin, Mendel Rosenblum, Stephen Rumble, Ryan Stutsman, and Stephen Yang. The RAM-Cloud Storage System. ACM Trans. Comput. Syst., 33(3), aug 2015.
- [57] Marius Poke and Torsten Hoefler. DARE: High-Performance State Machine Replication on RDMA Networks. In *Proceedings of the 24th International Sym*posium on High-Performance Parallel and Distributed Computing, HPDC '15, page 107–118, New York, NY, USA, 2015. Association for Computing Machinery.
- [58] Waleed Reda, Marco Canini, Dejan Kostić, and Simon Peter. RDMA is Turing complete, we just did not know it yet! In *Proceedings of NSDI* '22, Apr 2022.
- [59] Jiwu Shu, Youmin Chen, Qing Wang, Bohong Zhu, Junru Li, and Youyou Lu. TH-DPMS: Design and Implementation of an RDMA-Enabled Distributed Persistent Memory Storage System. ACM Trans. Storage, 16(4), oct 2020.
- [60] David Sidler, Zeke Wang, Monica Chiosa, Amit Kulkarni, and Gustavo Alonso. StRoM: Smart Remote Memory. In *Proceedings of the Fifteenth European Conference on Computer Systems*, EuroSys '20, New York, NY, USA, 2020. Association for Computing Machinery.
- [61] Adriana Szekeres, Michael Whittaker, Jialin Li, Naveen Kr. Sharma, Arvind Krishnamurthy, Dan R. K. Ports, and Irene Zhang. Meerkat: Multicore-Scalable Replicated Transactions Following the Zero-Coordination Principle. In *Proceedings of the* Fifteenth European Conference on Computer Systems, EuroSys '20, New York, NY, USA, 2020. Association for Computing Machinery.

- [62] Yacine Taleb, Ryan Stutsman, Gabriel Antoniu, and Toni Cortes. Tailwind: Fast and Atomic RDMA-based Replication. In 2018 USENIX Annual Technical Conference (USENIX ATC 18), pages 851–863, Boston, MA, July 2018. USENIX Association.
- [63] Shin-Yeh Tsai, Yizhou Shan, and Yiying Zhang. Disaggregating Persistent Memory and Controlling Them Remotely: An Exploration of Passive Disaggregated Key-Value Stores. In 2020 USENIX Annual Technical Conference (USENIX ATC 20), pages 33–48. USENIX Association, July 2020.
- [64] Shin-Yeh Tsai and Yiying Zhang. LITE Kernel RDMA Support for Datacenter Applications. In *Proceedings of* the 26th Symposium on Operating Systems Principles, SOSP '17, page 306–324, New York, NY, USA, 2017. Association for Computing Machinery.
- [65] Cheng Wang, Jianyu Jiang, Xusheng Chen, Ning Yi, and Heming Cui. APUS: Fast and Scalable Paxos on RDMA. In *Proceedings of the 2017 Symposium on Cloud Computing*, SoCC '17, page 94–107, New York, NY, USA, 2017. Association for Computing Machinery.
- [66] Jing Wang, Youyou Lu, Qing Wang, Minhui Xie, Keji Huang, and Jiwu Shu. Pacman: An Efficient Compaction Approach for Log-Structured Key-Value Store on Persistent Memory. In 2022 USENIX Annual Technical Conference (USENIX ATC 22), pages 773–788, Carlsbad, CA, July 2022. USENIX Association.
- [67] Qing Wang, Youyou Lu, Junru Li, and Jiwu Shu. Nap: A Black-Box Approach to NUMA-Aware Persistent Memory Indexes. In 15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21), pages 93–111. USENIX Association, July 2021.
- [68] Qing Wang, Youyou Lu, and Jiwu Shu. Sherman: A Write-Optimized Distributed B+Tree Index on Disaggregated Memory. In Proceedings of the 2022 International Conference on Management of Data, SIGMOD '22, page 1033–1048, New York, NY, USA, 2022. Association for Computing Machinery.
- [69] Xingda Wei, Zhiyuan Dong, Rong Chen, and Haibo Chen. Deconstructing RDMA-Enabled Distributed Transactions: Hybrid is Better. In *Proceedings of the* 13th USENIX Conference on Operating Systems Design and Implementation, OSDI'18, page 233–251, USA, 2018. USENIX Association.
- [70] Xingda Wei, Xiating Xie, Rong Chen, Haibo Chen, and Binyu Zang. Characterizing and Optimizing Remote Persistent Memory with RDMA and NVM. In 2021 USENIX Annual Technical Conference (USENIX ATC 21), pages 523–536. USENIX Association, July 2021.

- [71] Fei Xia, Dejun Jiang, Jin Xiong, and Ninghui Sun. HiKV: A Hybrid Index Key-Value Store for DRAM-NVM Memory Systems. In *Proceedings of the 2017 USENIX Conference on Usenix Annual Technical Conference*, USENIX ATC '17, page 349–362, USA, 2017. USENIX Association.
- [72] Jian Yang, Joseph Izraelevitz, and Steven Swanson. FileMR: Rethinking RDMA Networking for Scalable Persistent Memory. In 17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20), pages 111–125, Santa Clara, CA, February 2020. USENIX Association.
- [73] Jian Yang, Juno Kim, Morteza Hoseinzadeh, Joseph Izraelevitz, and Steve Swanson. An Empirical Guide to the Behavior and Use of Scalable Persistent Memory. In 18th USENIX Conference on File and Storage Technologies (FAST 20), pages 169–182, Santa Clara, CA, February 2020. USENIX Association.
- [74] Wenhui Zhang, Xingsheng Zhao, Song Jiang, and Hong Jiang. ChameleonDB: A Key-Value Store for Optane Persistent Memory. In *Proceedings of the Sixteenth Eu-*

- ropean Conference on Computer Systems, EuroSys '21, page 194–209, New York, NY, USA, 2021. Association for Computing Machinery.
- [75] Yiying Zhang, Jian Yang, Amirsaman Memaripour, and Steven Swanson. Mojim: A Reliable and Highly-Available Non-Volatile Memory System. In Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '15, page 3–18, New York, NY, USA, 2015. Association for Computing Machinery.
- [76] Jingyu Zhou, Meng Xu, Alexander Shraer, Bala Namasivayam, Alex Miller, Evan Tschannen, Steve Atherton, Andrew J. Beamon, Rusty Sears, John Leach, Dave Rosenthal, Xin Dong, Will Wilson, Ben Collins, David Scherer, Alec Grieser, Young Liu, Alvin Moore, Bhaskar Muppana, Xiaoge Su, and Vishesh Yadav. FoundationDB: A Distributed Unbundled Transactional Key Value Store. In *Proceedings of the 2021 International Conference on Management of Data*, SIGMOD/PODS '21, page 2653–2666, New York, NY, USA, 2021. Association for Computing Machinery.