







Maximizing Interconnect Bandwidth and Efficiency in NVMe-Based Key-Value SSDs with Fine-Grained Value Transfer

Junhyeok Park¹ , Chang-Gyu Lee¹ , Soon Hwang¹ , Seung-Jun Cha² , Woosuk Chung³ , Youngjae Kim¹ ,
¹Sogang University, Seoul, South Korea, ²ETRI, Daejeon, South Korea, ³Memory Systems Research, SK hynix.

Abstract—The Key-Value SSD (KV-SSD) redefines storage interfaces by integrating a key-value store directly within the device, offering native support for non-page-aligned key-value pairs. This architectural innovation enables KV-SSDs to offload storage management from the host system, positioning them as ideal candidates for resource disaggregation. However, KV-SSDs face significant challenges, notably I/O amplification caused by conflicts with traditional storage protocols like NVMe, which are designed around memory page units. Specifically, this results in a substantial increase in data traffic over interconnect between the host and SSD. This paper introduces BandSlim, a novel solution addressing these challenges in data transfer by leveraging (i) NVMe command piggybacking for universally compatible, fine-grained transfers without requiring interconnect-level support and (ii) the remote memory access capabilities of emerging Compute Express Link (CXL) interconnects for higher-performance fine-grained transfers. Through extensive evaluations, BandSlim achieves up to a 97.9% reduction in PCIe traffic compared to conventional NVMe KV-SSDs.

Introduction

Designing an efficient storage system necessitates reducing data movement costs from the host's memory to the storage media. However, traditional Key-Value Stores (KVSs) like RocksDB [1] function as middleware on top of file systems. Consequently, user I/O requests must navigate through the kernel's file system and block layers to execute data reads and writes to storage. This multi-layer traversal incurs significant memory copying and kernel context switch overheads during I/O operations. In contrast, Key-Value SSDs (KV-SSDs) offer a substantial reduction in these overheads. KV-SSDs implement KVSs at the device controller level and provide native support for key-value operations. This allows users to bypass the file system and block layer when processing I/O requests, enabling lower latency and higher throughput compared to traditional SSD-based host-side KVSs. Such KV-SSDs

are well-suited for modern resource disaggregation architectures as they can effectively decouple storage management from the host system [2].

KV-SSDs hold significant potential as a next-generation storage device but still face critical challenges, with the biggest issue being I/O amplification caused by limitations in existing protocols. Commercially and academically released KV-SSDs [3]–[6] utilize the Non-Volatile Memory Express (NVMe) protocol [7], which is specifically engineered for block-based storage devices. These KV-SSDs implement the Physical Region Page (PRP) list for conveying payload, essentially inducing I/O amplification originating from the size difference between the block and key-value. This results in a substantial increase in data traffic over the interconnect between the host and SSD.

One common solution to this problem has been host-side batching, where enough key-value pairs are grouped to align with the memory page size, as implemented in KV-SSDs like KV-CSD [4] through bulk PUT operations. However, buffering entries on the host side risks data loss during power failure, making this

approach unsuitable for mission-critical applications where data integrity is essential. For those scenarios where data persistence is critical and I/O transactions occur at the key-value pair level as defined by NVMe standard, a more fundamental solution is required.

To address these issues, we introduce *BandSlim*, which minimizes data traffic over host-SSD interconnect when transferring small key-value pairs. *BandSlim* employs two approaches: the first utilizes NVMe commands to achieve an inline value transfer at the storage protocol level, while the second leverages the emerging Compute eXpress Link (CXL) interconnect to address the transfer at the interconnect protocol level. The first approach piggybacks values smaller than a memory page to NVMe commands using reserved fields. We observed that this method significantly reduces data traffic over the host-SSD interconnect. The second approach, which uses CXL's remote memory access capability, aims to improve this further by offering a more fine-grained, higher-performance value transfer solution, bypassing the limitations of the piggybacking. It exposes the device-side NAND page buffer to the host's CPU and transfers values using CXL.mem operations. The first method can be implemented directly in existing NVMe protocols without requiring new interconnect support, while the second approach is optimized for systems that support CXL, allowing for more advanced and efficient fine-grained transfers.

In both methods, however, as the value size grows, the performance of fine-grained transfer becomes less favorable compared to the traditional PRP-based transfer. Therefore, *BandSlim* employs an adaptive transfer strategy, dynamically switching between fine-grained and PRP-based transfer to balance traffic reduction with optimal response times for varying value sizes.

For evaluation, we implemented *BandSlim* on a state-of-the-art FPGA-based NVMe KV-SSD [5] using the Cosmos+ OpenSSD platform [8], and estimated CXL.mem-based value transfer's benefits using real CXL.mem measurements. We demonstrated that *BandSlim* achieves up to 97.9% traffic reduction compared to NVMe KV-SSD without *BandSlim*.

The contributions of this paper are as follows:

- Identified traffic amplification in the host-SSD interconnect, specifically in KV-SSDs, and proposed NVMe and CXL-based solutions to minimize waste.
- Proposed a design of a CXL-based storage interface and demonstrated use cases where the interface shows clear advantages, particularly in KV-SSDs.
- Demonstrated the trade-off between data traffic and response time in fine-grained value transfer and effectively resolved it using an adaptive approach.

Background and Motivation

Storage Stack of KV-SSD

The storage stack of Key-Value SSDs (KV-SSDs) consists of user-level key-value APIs, a key-value device driver and controller based on protocols like NVMe, and an in-storage Key-Value Store (KVS). The key-value API offers point and range queries (e.g., PUT, SEEK). The size of the key and value in these APIs is handled as arbitrary length, not in block units (*key-value interface*). In case of LSM-based KV-SSDs with a key-value separation [3]–[6], a pair of key and value address is stored in the LSM-tree, and a value is stored in the Value Log (vLog). The vLog is a linear logical NAND flash address space which can be divided into multiple logical NAND pages. Each value is appended to the vLog sequentially, filling logical NAND pages which are mapped to physical NAND pages by the Flash Translation Layer (FTL). The entries of the LSM-tree point to corresponding values inside the vLog.

NVMe-Based Key-Value Pair Transfer

Within the NVMe key-value interface, when writing key-value pairs, the NVMe driver stores a key and metadata in the reserved fields of the NVMe command. The payload, which is the value in this context, is transferred via the Physical Region Page (PRP) as the block interface part of NVMe specification does [7]. The PRP is a linked list whose entry describes the addresses of physical memory pages of the host memory. One or more memory pages where the value is stored are specified to be transferred. Subsequently, the driver inserts the NVMe command into the submission queue and rings the doorbell to notify the device of the write request. The NVMe controller fetches the command from the queue, interprets it, and identifies the pages for copying from the received PRP list.

To initiate the value transfer, the controller triggers a Direct Memory Access (DMA) transaction, which copies pages from the host memory to the device memory. The controller later inserts the received key to the memory component of LSM-tree, MemTable, and writes the value to the vLog (see Figure 1). The reverse operation, involving the transfer of values from the KV-SSD to the host, follows a similar process.

PCIe Traffic Amplification in NVMe KV-SSD

As in typical KVSs, the key and value size are variable and not necessarily aligned to a memory page. According to Meta, RocksDB in a production environment experiences the size of values nearly not reaching a hundred bytes on average [9], which is far less than

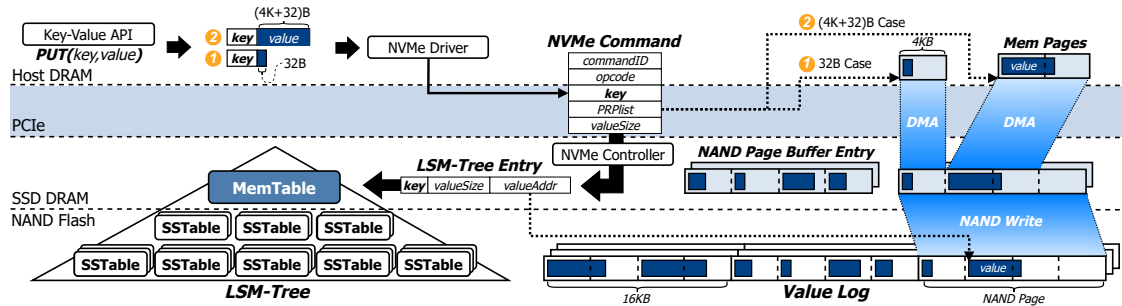
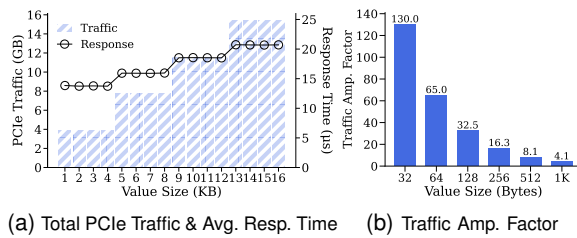


FIGURE 1: Data flow of two cases of key-value transfers with sub-4 KB payload (32 bytes) and over-4 KB payload ((4K+32) bytes) regarding PCIe interconnect traffic amplification in LSM-based NVMe KV-SSDs.



(a) Total PCIe Traffic & Avg. Resp. Time (b) Traffic Amp. Factor

FIGURE 2: Measurements of total PCIe traffic with average response time and PCIe traffic amplification factor for varying value sizes, with NAND I/O disabled.

the 4 KB memory page size. Consequently, a KV-SSD must be capable of effectively handling requests for such variable-sized, small values.

However, the current NVMe key-value interface causes significant inefficiencies because it adheres to the same procedure as the block-based NVMe protocol when transferring values to or from the device. *Specifically, its data transfer method, PRP, restricts DMA transfers to occur in units of 4 KB, a size of memory page.* As shown in Figure 1, consider a scenario where the value size is 32 bytes (1). In this case, one 4 KB memory page that temporarily holds the value is specified by the PRP, and a DMA copy of 4 KB occurs. On the other hand, if the value size slightly exceeds the memory page size, such as (4K+32) bytes (2), two memory pages are required to accommodate the value. Consequently, the DMA facilitated by the PRP transfers 8 KB of data. This amplified data traffic significantly raises energy consumption of the system, possibly increasing the Total Cost of Ownership (TCO).

Experimental Analysis of Traffic Amplification. We measured PCIe traffic from a host to a KV-SSD [5] by issuing 1 million key-value writes with variable-sized values, by using the Intel PCM [10] to track PCIe traffic.

Figure 2(a) shows total PCIe traffic during experiments, with traffic increasing stepwise at 4 KB value size boundaries. For example, total PCIe traffic for 1 KB and 4 KB values is around 4 GB, indicating constant transfer volumes up to 4 KB. This pattern repeats

for value size ranges like 5-8 KB. Average transfer response times display similar cascading patterns.

The issue becomes more severe with smaller values. The Traffic Amplification Factor, defined as the ratio of PCIe traffic to data size, surges for values under 1 KB. As shown in Figure 2(b), transferring a value of 32 bytes generates approximately 4 KB of PCIe data traffic, about 130 times the requested data size.

Limitations of Existing Methods. The NVMe protocol currently offers another transfer method besides PRP, called the Scatter-Gather List (SGL) [7]. Unlike PRP, SGL supports variable-sized DMAs across scattered memory segments. However, it has been reported that the cost of enabling the SGL outweighs the benefit for I/O smaller than 32 KB [11], indicating that using SGL for small value transfers is not advisable and realistic.

Proposed Design: *BandSlim*

Fine-Grained Value Transfer over NVMe

Due to space constraints, a detailed discussion of *BandSlim*'s first proposed technique, NVMe command piggybacking, is provided in the conference version of this paper [12]. A summary is as follows: the method leverages NVMe command, which is 64 bytes in size, to enable fine-grained transfer of small values. Considering most values in real-world are under 64 bytes [9], it repurposes up to 35 bytes of unused fields in the NVMe key-value write command for value transfer. For values exceeding 35 bytes, the piggybacking method introduces a `transfer` command, which sends the remaining bytes in 56 bytes per command after the initial write command. This approach significantly reduces PCIe traffic under real-world KVS workloads.

Fine-Grained Value Transfer using CXL

While the piggybacking method resolves the traffic amplification issue, it is still constrained by several key limitations. First, the overhead associated with creating, submitting, processing, and releasing NVMe

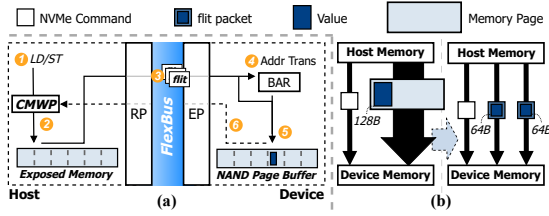


FIGURE 3: (a) CXL.mem memory write mechanism, and (b) Fine-grained value transfer via CXL.mem.

commands is significant, particularly when dealing with large value that demands issuing multiple commands. It makes this method effective only for transfers of tiny values, typically in the range of tens of bytes. Second, it cannot completely eliminate traffic bloating since it still depends on NVMe commands which requires extra traffic for signaling doorbells and completions.

To overcome the limitations of repetitive NVMe command overheads in the piggybacking method, we also propose a design that fundamentally supports fine-grained transfer at the interconnect level using CXL. For this, *BandSlim* defines the NAND page buffer as Host-managed Device Memory (HDM), making it directly accessible to the host through CXL.mem operations. During a CXL device enumeration, the driver queries the Base Address Register (BAR) and the size of the HDM to map the BAR and HDM within the host's system memory. The host CPU's system bus includes CXL Root Ports (RPs), establishing connections with KV-SSDs as End-Points (EPs). After that, the driver can transmit values via CXL.mem write operations.

Figure 3 illustrates the value transmission process in this method. The driver manages a CXL.mem Write Pointer (CMWP) to track the current position in the NAND page buffer. When a user initiates a key-value write, the NVMe command includes only the key, unlike the previous piggybacking method. ① The driver initiates a memory copy request to the CMWP. ② The host CPU then sends the value to the NAND page buffer at the CMWP location via CXL.mem. ③ The CXL RP receives this request, converting it into *flits*, CXL transaction units (*BandSlim* operates based on the default 68 bytes *flit* mode). ④ The CXL controller decompresses the *flits*, adjusts target addresses by subtracting the HDM base address, and forwards the request to the DRAM controller, ⑤ which sends the payload to the DRAM location indicated by the CMWP. ⑥ The controller notifies the driver to update the CMWP via NVMe command completion.

Since the FTL operates independently, the host cannot determine when each buffer entry is flushed. Thus, *BandSlim* always calls *clflush()* for CXL.mem writes to ensure the payload is sent to device memory.

Adaptive Value Transfer Method

As the value size increases, the performance of both fine-grained value transfer methods starts to degrade. The first method, in particular, suffers due to the accumulation of overheads in generating NVMe commands and synchronously handling them within the device, resulting in way longer transfer times compared to conventional PRP-based transfer. To tackle these issues, *BandSlim* utilizes a threshold-based reactive method that selects the most suitable transfer method from fine-grained transfer and PRP-based transfer based on the size of the value. The threshold is identified by exploratory runs using benchmarks depicted in the Evaluation. During the benchmark runs, various value sizes ranging from 4 bytes to 8 KB are tested through millions of PUT commands to compare transfer times.

The threshold is denoted as $\tau_\alpha = \alpha \cdot \tau$, where τ represents the value size at which fine-grained transfer becomes less efficient than PRP-based transfers. *BandSlim* employs the following strategy:

$$\text{Transfer Method} = \begin{cases} \text{Fine-Grained Transfer,} & \text{if Value Size} < \tau_\alpha, \\ \text{PRP-Based Transfer,} & \text{if Value Size} \geq \tau_\alpha. \end{cases}$$

The scaling factor α allows users to adjust the threshold: increasing α raises the threshold, while setting α to 1 retains the default threshold. For users prioritizing response time, α can be set to 1. For those focusing on traffic reduction, α can be increased to delay the transition point where PRP-based transfers become more efficient. This approach ensures efficient handling of value sizes ranging from sub-page to large.

In-Device Value Packing Mechanism. Our conference paper proposed the *Selective Packing with Back-filling Policy* to address the NAND page write amplification problem in NVMe KV-SSDs [12]. This policy aims to reduce the overhead of copying large values during the packing process within the device under adaptive value transfer by selectively packing only small values transmitted via fine-grained transfer, while large values are placed to page-aligned address via PRP-based DMA without packing. To minimize internal fragmentation, the policy allows small values to fill gaps between large values at page-aligned addresses and also avoid them with a DMA Log Table (DLT).

Figure 4 illustrates the process of packing under CXL.mem-based transfer. When a user requests a large value write that exceeds the threshold τ_α , the value is transferred via PRP DMA over NVMe (CXL.io) ①. The controller completes the DMA and creates a DLT entry based on the DMA destination address and value size ②. This entry is sent to the host in the NVMe Completion Queue Entry (CQE) ③, where the *BandSlim* driver caches it ④. For a small value below

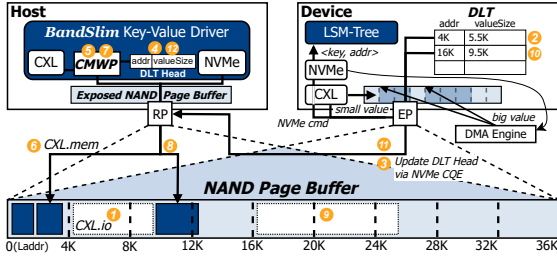


FIGURE 4: A process of *Selective Packing with Back-filling* over the NAND page buffer with CMWP and DLT.

the threshold τ_α , the driver checks if adding the value size to the CMWP exceeds the specified address in the cached DLT head entry. If it does not, the value is transferred to the CMWP address over CXL.mem with *clflush()* enabled, while only the key is sent in an NVMe command (5-6). If the next small value size exceeds the address in the cached DLT head entry, the CMWP updates to the DLT entry's address plus the value size, and the value is transferred to the updated CMWP address via CXL.mem (7-8). For another large value, the transfer occurs via PRP DMA, with the destination address aligned to the nearest page boundary, tracked by the controller (9). The new address and value size are added to the DLT and cached on the host (10-12).

Since only the NAND page buffer is exposed as HDM, while CMWP and DLT prevent conflicts between CXL.mem and PRP DMA, *BandSlim* preserves data consistency and metadata integrity. Plus, since *BandSlim* does not modify the core components of KV-SSDs, including key-value indexing and FTL, the existing read and query performance remain unaffected.

Evaluation

Evaluation Setup

We used the state-of-the-art NVMe KV-SSD [5] on the Cosmos+ OpenSSD platform [8] as the baseline to verify our solutions. The platform consists of an Arm-based Xilinx Zynq-7000 SoC, 1 TB of NAND (4 channels, 8 ways), and a PCIe Gen2 $\times 8$ interconnect, paired with a host node featuring 64 cores of Intel Xeon Gold 6226R CPU, 384 GB of DDR4 memory, and Ubuntu 22.04 OS. The SoC operates the *BandSlim* key-value controller, PCIe interface controller, DRAM controller, and NAND flash controller, while the host node runs the *BandSlim* key-value driver and benchmarks.

Since publicly accessible SSDs and FPGA boards that support CXL remains extremely scarce, including ours, we validated our designs by implementing the piggybacking mechanism on the OpenSSD-based setup, while the CXL.mem-based transfer was evalu-

ated through simulations using measured CXL.mem response times. The traditional PRP-based transfer and the NVMe piggybacking method underwent thorough system-level validation on the OpenSSD setup, providing reliable performance measurements. We then scaled the measured payload transfer times of baseline and piggybacking to PCIe Gen5, and incorporated the NVMe overheads into the CXL.mem response times to simulate the CXL.mem-based transfers. This approach allowed us to directly compare the performance of PRP-based transfers, NVMe piggybacking and CXL.mem-based transfers, demonstrating the potential benefits of CXL.mem optimizations.

For performance evaluations, we used *db_bench*, a benchmarking tool used in RocksDB [1]. We enabled *db_bench* to send NVMe key-value commands to the OpenSSD platform through the NVMe passthrough.

We conducted various workload patterns to verify our proposed design, described as follows.

- **Workload A ($W(A)$):** A *db_bench*'s *fillseq* pattern where keys are sequential and value sizes are fixed. It serves as a baseline for evaluating performance under uniform and predictable write patterns.
- **Workload B ($W(B)$):** This writes 1 million random key-value pairs with value sizes of 8 bytes or 2 KB at a 9:1 ratio. This tests the transfer method's efficiency in optimizing for frequent small value writes.
- **Workload C ($W(C)$):** Similar to $W(B)$ but reverses the value size ratio to 1:9, emphasizing scenarios with large value dominance. The goal is to reveal trade-offs in handling small versus large values.
- **Workload D ($W(D)$):** This workload writes values of sizes (8, 16, 32, 64, 128, 256, 512 bytes, 1 KB, and 2 KB) in random order, totaling 1 million, with each size having an equal ratio. This evaluates the transfer method's adaptability to diverse data sizes.
- **Workload M ($W(M)$):** This is a *db_bench*'s *mixgraph All_random* [9] workload which reflects real-world characteristics with a maximum value size of 1 KB and almost 70% of values being under 35 bytes. We have modified *mixgraph* to issue only 1 million PUTs.

In all experiments, we used unique keys of 4 bytes generated by a hash function with a random seed.

We conducted evaluations for the following designs.

- **Baseline:** the state-of-the-art NVMe KV-SSD [5]. It employs the PRP-based page-unit value transfer.
- **Piggyback:** transfers values only via piggybacking.
- **CXL-FGT:** transfers values only via CXL.mem.
- **Adaptive:** transfers values via the adaptive method.

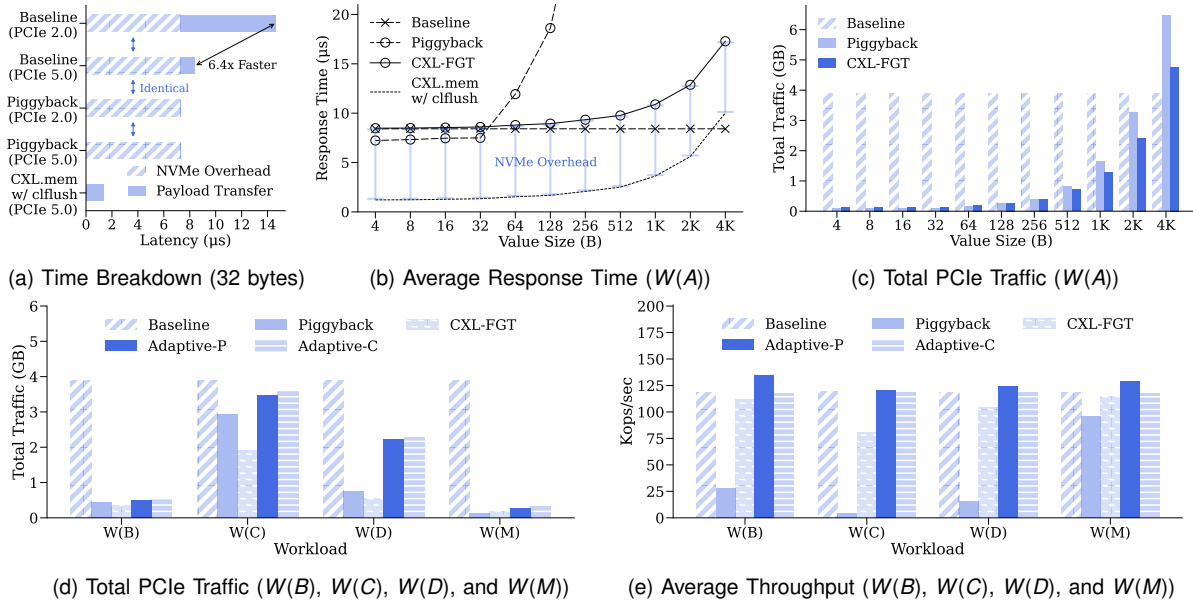


FIGURE 5: (a) Breakdown of estimated response times for values of 32 bytes on PCIe Gen5 using *Baseline* and *Piggyback*. (b) Estimated response times and (c) total PCIe traffic for *Piggyback* and *CXL-FGT* with 1 million key-value pairs (*W(A)*) of varying value sizes transferred from host to device memory. (d) Estimated PCIe traffic and (e) average throughput for various workloads (*W(B)*, *W(C)*, *W(D)*, and *W(M)*) on PCIe Gen5.

Effects of Fine-Grained Value Transfer

Piggybacking-Based Value Transfer. The evaluation of PCIe traffic and response times for *Piggyback* on a real platform is detailed in our conference paper [12], as space constraints prevent an in-depth discussion here. Briefly, *Piggyback* achieves up to 97.9% traffic reduction for values of 4 to 32 bytes but becomes less efficient as value size increases due to the overheads of trailing commands, approaching *Baseline* at around 2 KB and surpassing it for larger sizes. Response times are halved for values up to 32 bytes but degrade for larger sizes due to the trailing commands overheads.

CXL.mem-Based Value Transfer. To evaluate *CXL-FGT*, we first measured write response times for various value sizes using a 256 GB CXL memory device connected via PCIe Gen5, consistently invoking *clflush()*. These results, obtained on a dedicated server with an CXL-supported AMD EPYC 9754 configured as a CPU-less NUMA node, served as the foundation for a performance simulation model of *CXL-FGT*.

Next, we analyzed the response times [12] for *Baseline* and *Piggyback* measured from our OpenSSD (PCIe Gen2) setup. We extracted the payload transfer time for various payload sizes from the total response times by separating the time spent on generating, submitting, processing, and completing NVMe commands. We then scaled the extracted payload transfer times by a factor of 6.4, reflecting the higher gigatransfer rate of Gen5 (32 GT/s) compared to Gen2 (5 GT/s).

While actual performance gains may be smaller due to system overheads, this adjustment assumes an ideal scenario focused purely on bandwidth improvements. We added the NVMe overhead back to the scaled payload transfer times of *Baseline* and *Piggyback* to estimate their response times on PCIe Gen5. Finally, we created a performance simulation model for *CXL-FGT* by combining the measured CXL.mem write response times with the same NVMe overheads used for *Baseline* (*CXL-FGT* always require a single NVMe command). This approach enables a direct comparison of *Baseline*, *Piggyback*, and *CXL-FGT* designs. Figure 5(a) illustrates the time breakdown and adjustment used to estimate the response times of *Baseline* and *Piggyback* on PCIe Gen5 for values of 32 bytes.

Using the *CXL-FGT* simulation model and the estimated response times of *Baseline* and *Piggyback*, we evaluated the performance of each transfer mode across different value sizes using *Workload A*. The results, presented in Figure 5(b), reveal key insights. First, *CXL-FGT* performs worse than *Piggyback* for value sizes of 4 bytes to 32 bytes due to the mandatory NVMe command overhead separated from CXL transfers. This limitation could be mitigated by integrating piggybacking with *CXL-FGT*. Second, *CXL-FGT* does not outperform *Baseline* on PCIe Gen5, again, highlighting a limitation due to the mandatory NVMe command overhead. However, as expected, *CXL-FGT* significantly outperforms *Piggyback* starting

at the 64 bytes case, demonstrating how CXL.mem mitigates the performance degradation of *Piggyback* for larger data sizes. For instance, in the 256 bytes case, *CXL-FGT* achieves nearly 4× faster response times compared to *Piggyback*, with an even greater performance gap as value sizes increase.

Figure 5(c) shows that *CXL-FGT* generates less traffic compared to *Piggyback*. Using CXL *flit* units of 68 bytes, we estimated *CXL-FGT*'s traffic, while *Baseline* and *Piggyback* reflect measured values [12]. CXL enables finer granularity since *Piggyback* incurs additional traffic with each `transfer` command submission, such as doorbell rings and tail pointer reads. CXL avoids this overhead, allowing more efficient data movement. However, note that piggybacking's primary strength lies in resolving KV-SSD bottlenecks without requiring new interconnect technologies.

Effects of Adaptive Value Transfer

The evaluation results of the adaptive value transfer using *Workloads B, C, D, and M* on a real platform is detailed in our conference paper [12]. Briefly, *Adaptive* consistently outperforms *Piggyback* and *Baseline* by transitioning from piggybacking to PRP at 128 bytes (identified as the threshold based on experiments with *W(A)*), achieving the best performance in all workloads. It reduces traffic by 93.3% for *W(M)* while improving throughput by 12% over *Piggyback*, and in *W(C)*, it increases throughput nearly 13-fold despite generating 18% more traffic than *Piggyback*.

To assess the impact of CXL.mem-based transfer under various workloads, we simulated *Workloads B, C, D, and M* using results from Figure 5(b) and (c). Figure 5(d) shows the traffic results, where *CXL-FGT* achieved lower overall traffic than *Piggyback* in all workloads except *W(M)*, despite the additional traffic overhead compared to *Piggyback* for transferring values of 4 to 32 bytes. This reduction stems from avoiding frequent NVMe command overhead by using CXL.mem, with the benefit most evident in *W(C)* due to its larger values. For the *Adaptive* method, we evaluated scenarios employing piggybacking (*Adaptive-P*) and CXL.mem (*Adaptive-C*). The threshold for switching the transfer mode in *Adaptive-P* was set to 64 bytes, based on the results in Figure 5(b). The same threshold was applied to *Adaptive-C* to ensure a fair comparison between piggybacking and CXL.mem. As expected, both methods exhibited higher traffic than fine-grained transfers. Moreover, since *Piggyback* inherently resulted in lower traffic compared to *CXL-FGT* for value sizes between 4 and 32 bytes, *Adaptive-C* exhibited slightly higher traffic than *Adaptive-P*.

Figure 5(e) presents the average throughput derived by the simulated response times. Across all workloads, *CXL-FGT* consistently outperformed *Piggyback*, further demonstrating the effectiveness of CXL.mem-based transfers. Meanwhile, *Adaptive-P* achieved the best performance in all workloads, while *Adaptive-C* generally showed performance comparable to the *Baseline*. Again, this result is attributed to the inherent limitations of *CXL-FGT*, which retains the mandatory NVMe command overhead. In particular, workloads with a high proportion of value sizes between 4 and 32 bytes (e.g., *W(B)*, *W(M)*) highlighted the limitations of *CXL-FGT*, which adheres to NVMe routines, resulting in *Adaptive-C* performing worse than *Adaptive-P*.

In summary, *CXL-FGT* achieves finer-grained traffic and significantly better performance than *Piggyback* as value sizes increase. However, as the method maintains NVMe transactions, the performance advantage of fine-grained transfer over the *Baseline* diminishes with the increased bandwidth of PCIe Gen5. Furthermore, *CXL-FGT* does not fully utilize cache coherency of CXL as it actively invokes `clflush()` for each write. We plan to address these limitations by developing new storage protocols leveraging the full potential of CXL.

Conclusion

BandSlim tackles the PCIe traffic amplification challenges in KV-SSDs, which arise from conflicts between non-page-aligned key-value pairs and NVMe protocols designed for memory page units. *BandSlim* leverages NVMe command piggybacking and CXL.mem to transfer values in a fine-grained and efficient manner. Extensive evaluations showcases that *BandSlim* significantly reduces traffic and optimizes data transfer efficiency.

ACKNOWLEDGMENTS

This work was funded in part by the National Research Foundation of Korea (NRF) grant funded by the Korean Government (MSIT) (RS-2024-00453929), in part by the Institute of Information & Communications Technology Planning Evaluation (IITP) grant funded by the Korea Government (MSIT) under Grant RS-2021-11210528, and in part by an SK hynix research grant. Y. Kim is the corresponding author.

REFERENCES

1. Facebook, "RocksDB," <http://rocksdb.org>, 2014.
2. A. Tomlin, "Why KV SSD will replace ZNS," <https://www.snia.org/educational-library/why-kv-ssd-will-replace-zns-2022>, 2022.

3. C. G. Lee, H. Kang, D. Park, S. Park, Y. Kim, J. Noh, W. Chung, and K. Park, "iLSM-SSD: An Intelligent LSM-Tree Based Key-Value SSD for Data Analytics," in *Proceedings of International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*, 2019.
4. I. Park, Q. Zheng, D. Manno, S. Yang, J. Lee, D. Bonnie, B. Settlemeyer, Y. Kim, W. Chung, and G. Grider, "KV-CSD: A Hardware-Accelerated Key-Value Store for Data-Intensive Applications," in *Proceedings of the IEEE International Conference on Cluster Computing (CLUSTER)*, 2023, pp. 132–144.
5. S. Lee, C. G. Lee, D. Min, I. Park, W. Chung, A. Sivasubramaniam, and Y. Kim, "Iterator Interface Extended LSM-tree-based KVSSD for Range Queries," in *Proceedings of the 16th ACM International Systems and Storage Conference (SYSTOR)*, 2023.
6. D. Min, K. Kim, C. Moon, A. Khan, S. Lee, C. Yun, W. Chung, and Y. Kim, "A Multi-tenant Key-value SSD with Secondary Index for Search Query Processing and Analysis," *ACM Transactions on Embedded Computing Systems*, vol. 22, no. 4, Jul. 2023.
7. NVM Express Inc., "Key Value Command Set Specification," <https://nvmexpress.org/specifications/>, 2024.
8. J. Kwak, S. Lee, K. Park, J. Jeong, and Y. H. Song, "Cosmos+ OpenSSD: Rapid prototype for flash storage systems," *ACM Trans. Storage*, Jul. 2020.
9. H. Cao, S. Dong, S. Vemuri, and D. H. C. Du, "Characterizing, modeling, and benchmarking RocksDB key-value workloads at facebook," in *Proceedings of the 18th USENIX Conference on File and Storage Technologies (FAST)*, 2020, pp. 209–224.
10. Intel, "Intel PCM," <https://github.com/intel/pcm>, 2010.
11. "nvme : add Scatter-Gather List (SGL) support in NVMe driver," <https://merlin.infradead.org/pipermail/linux-nvme/2017-July/011956.html>, 2017.
12. J. Park, C.-G. Lee, S. Hwang, S. Yang, J. Noh, W. Chung, J. Lee, and Y. Kim, "BandSlim: A Novel Bandwidth and Space-Efficient KV-SSD with an Escape-from-Block Approach," in *Proceedings of the 53rd International Conference on Parallel Processing (ICPP)*, ser. ICPP '24, 2024, p. 1187–1196.

Junhyeok Park received the B.S. degree in Computer Science and Engineering from Sogang University, Seoul, South Korea, in 2024, where he is currently pursuing his M.S. degree. His research interests include next-generation NAND flash drives and file and storage systems. Email: junttang@sogang.ac.kr.

Chang-Gyu Lee received the B.S. degree in Software and Computer Engineering from Ajou University, Suwon, South Korea, in 2017. He is currently pursuing a

Ph.D. degree in Computer Science and Engineering at Sogang University, Seoul, South Korea. His research focuses on operating systems, and file and storage systems. Email: changgyu@u.sogang.ac.kr.

Soon Hwang received his B.S. and M.S. degrees in Computer Science and Engineering from Sogang University, Seoul, South Korea, in 2021 and 2023, respectively. After working as a software engineer at Samsung Electronics from 2023 to 2024, he began his Ph.D. studies in the same department at Sogang University. His research interests include file and storage systems for HPC and AI. Email: soonhw@sogang.ac.kr.

Seung-Jun Cha received his B.E., M.E., and Ph.D. degrees in Computer Engineering from Chungnam National University, Daejeon, South Korea, in 2006, 2008, and 2013, respectively. He is currently a senior researcher at the Electronics and Telecommunications Research Institute (ETRI) in Daejeon, South Korea. His research interests include enhancing performance in system software, such as operating systems and database management systems, and he is currently focused on advancing memory interconnect technologies. Email: seungjunn@etri.re.kr.

Woosuk Chung is a Vice President at SK hynix, leading the software solution group in the Memory Systems Research organization. With over 10 years of experience in the field, he has been at the forefront of research and development in AI and Big Data Analytics optimized memory and storage systems. His expertise spans the design and optimization of memory and storage architectures to meet the demands of emerging applications. Throughout his career, he has made significant contributions to the advancement of memory and storage technologies, driving innovation and excellence in the industry. Email: woosuk.chung@sk.com.

Youngjae Kim received the B.S. degree in Computer Science from Sogang University, Seoul, South Korea, in 2001, the M.S. degree in Computer Science from KAIST, in 2003, and the Ph.D. degree in Computer Science and Engineering from Pennsylvania State University, University Park, PA, USA in 2009. He is currently a professor with the Department of Computer Science and Engineering, Sogang University. Before joining Sogang University, he was a R&D Staff Scientist with the US Department of Energy's Oak Ridge National Laboratory, from 2009 to 2015. His research interests include operating systems, file and storage systems, parallel and distributed systems, and AI systems. Email: youkim@sogang.ac.kr.