

BPF for Storage: An Exokernel-Inspired Approach

Yuhong Zhong*
Columbia University

Hongyi Wang*
Columbia University

Yu Jian Wu*
Columbia University

Asaf Cidon
Columbia University

Ryan Stutsman
University of Utah

Amy Tai
VMware Research

Junfeng Yang
Columbia University

ABSTRACT

The overhead of the kernel storage path accounts for half of the access latency for new NVMe storage devices. We explore using BPF to reduce this overhead, by injecting user-defined functions deep in the kernel’s I/O processing stack. When issuing a series of dependent I/O requests, this approach can increase IOPS by over 2.5× and cut latency by half, by bypassing kernel layers and avoiding user-kernel boundary crossings. However, we must avoid losing important properties when bypassing the file system and block layer such as the safety guarantees of the file system and translation between physical blocks addresses and file offsets. We sketch potential solutions to these problems, inspired by exokernel file systems from the late 90s, whose time, we believe, has finally come!

“ As a dog returns to his vomit, so a fool repeats his folly. ”

Attributed to King Solomon

CCS CONCEPTS

• **Software and its engineering** → **Operating systems**;
• **Computer systems organization** → *Secondary storage organization*; • **Information systems** → *Data management systems*.

KEYWORDS

BPF, eBPF, storage, NVM, non-volatile memory

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

HotOS '21, May 31–June 2, 2021, Ann Arbor, MI, USA

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8438-4/21/05...\$15.00

<https://doi.org/10.1145/3458336.3465290>

ACM Reference Format:

Yuhong Zhong*, Hongyi Wang*, Yu Jian Wu*, Asaf Cidon, Ryan Stutsman, Amy Tai, and Junfeng Yang. 2021. BPF for Storage: An Exokernel-Inspired Approach. In *Workshop on Hot Topics in Operating Systems (HotOS '21), May 31–June 2, 2021, Ann Arbor, MI, USA*. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3458336.3465290>

1 INTRODUCTION

Storage devices have lagged behind networking devices in achieving high bandwidth and low latency. While 100 Gb/s bandwidth is now common for network interface cards (NICs) and the physical layer, storage devices only now support 2–7 GB/s bandwidth and 4–5 μs latencies [10, 16, 22, 23]. With such devices, the software stack is now a substantial overhead on every storage request. In our experiments this can account for about half of I/O operation latency, and the impact on throughput can be even more significant.

Kernel-bypass frameworks, e.g. SPDK [48], and near-storage processing reduce kernel overheads. However, there are clear drawbacks to both such as significant, often bespoke, application-level changes [44, 45], lack of fine-grained isolation, wasteful busy waiting, and the need for specialized hardware in the case of computational storage [13, 19, 46]. Therefore, we want a standard OS-supported mechanism that can reduce the software overhead for fast storage devices.

To address this, we turn to the networking community for inspiration. Linux’s eBPF [6] provides an interface for applications to embed simple functions directly in the kernel. When used to intercept I/O, these functions can perform processing that is traditionally done in the application while avoiding copying data and context switches between the kernel and user space. Linux eBPF is widely used for packet processing and filtering [5, 33], security [12] and tracing [32].

BPF’s ubiquity in the kernel and its wide acceptance make it a natural scheme for application-specific kernel-side extensions in layers outside of the networking stack. For example, it could be used to chain dependent I/Os such as traversing a

*These authors contributed equally to this work.

disk-based data structure like a B-tree where one block references another. Embedding such BPF functions deep enough in the kernel has the potential to eliminate nearly all of the software overhead of issuing I/Os such as in kernel-bypass, while avoiding the pitfalls of kernel-bypass such as polling and wasted CPU time.

To realize the performance gains of BPF we identify four substantial research challenges which are unique to the storage use case. First, for ease of adoption, our architecture must support Linux with standard file systems and applications with minimal modifications to either. It should be efficient and bypass as many software layers as feasible. Second, storage pushes BPF beyond simple packet processing uses. Packets are self-describing, so BPF can operate on them mostly in isolation. Traversing a structure on-disk is stateful and frequently requires consulting outside state. Storage BPF functions will need to understand applications’ on-disk formats and access outside state in the application or kernel, for example, to synchronize concurrent accesses or to access in-memory metadata. Third, BPF functions must not violate file system security guarantees while still allowing shared disk usage among applications. Storage blocks themselves typically do not record ownership or access control attributes, in contrast to network packets whose headers specify the flows to which the packets belong. Fourth, the approach must enable concurrency. Applications support concurrent accesses via fine-grained synchronization (e.g. lock coupling [31]) to avoid read-write interference with high throughput; synchronization from BPF functions may be needed.

Our approach is inspired by the work on exokernel file system designs. User-defined kernel extensions were the cornerstone of XN for the Xok exokernel. It supported mutually distrusting “libfs”es via code downloaded from user processes to the kernel [35]. In XN, these *untrusted deterministic functions* were interpreted to give the kernel a user-defined understanding of file system metadata. Then, the application and kernel design was clean slate, and absolute flexibility in file system layout was the goal. While we require a similar mechanism to allow users to enable the kernel to understand their data layout, our goals are different: we want a design that allows applications to define custom BPF-compatible data structures and functions for traversing, filtering and aggregating on-disk data, works with Linux’s existing interfaces and file systems, and substantially prunes the amount of kernel code executed per-I/O to drive millions of IOPS.

Finally, we make the observation, and empirically validate, that many popular storage data structures, such as B-trees or log-structured merge trees, are implemented as large, immutable files [11, 17, 24, 29], and that their extents get changed infrequently. This leads us to focus our initial

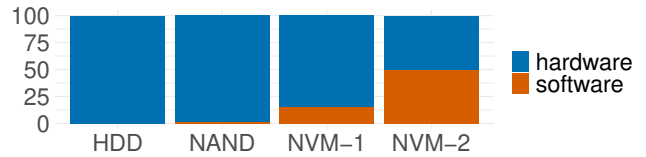


Figure 1: Kernel’s latency overhead with 512 B random reads. HDD is Seagate Exos X16, NAND is Intel Optane 750 TLC NAND, NVM-1 is first generation Intel Optane SSD (900P), and NVM-2 is second generation Intel Optane SSD (P5800X)

kernel crossing	351 ns	5.6%
read syscall	199 ns	3.2%
ext4	2006 ns	32.0%
bio	379 ns	6.0%
NVMe driver	113 ns	1.8%
storage device	3224 ns	51.4%
total	6.27 μ s	100.0%

Table 1: Average latency breakdown of a 512 B random read() syscall using Intel Optane SSD gen 2.

design on implementing BPF functions on immutable data structures with largely stable extents.

2 SOFTWARE IS THE STORAGE BOTTLENECK

The past few years have seen new memory technologies emerge in SSDs attached to high bandwidth PCIe using NVMe. This has led to storage devices that now rival the performance of fast network devices [1] with a few microseconds of access latency and gigabytes per second of bandwidth [10, 16, 22, 23]. Hence, just as the kernel networking stack emerged as a CPU bottleneck for fast network cards [26, 37, 41, 47], the kernel storage stack is now becoming a bottleneck for these new devices.

Figure 1 shows this phenomenon by breaking down the fraction of read I/O latency that can be attributed to the device hardware and the system software for increasingly fast storage devices. The results show the kernel’s added software overhead was already measurable (10-15% latency overhead) with the first generation of fast NVMe devices (e.g. first generation Optane SSD or Samsung’s Z-NAND). *In the new generation of devices software accounts for about half of the latency of each read I/O!* As devices get faster, the relative overhead of the kernel stack will only continue getting worse.

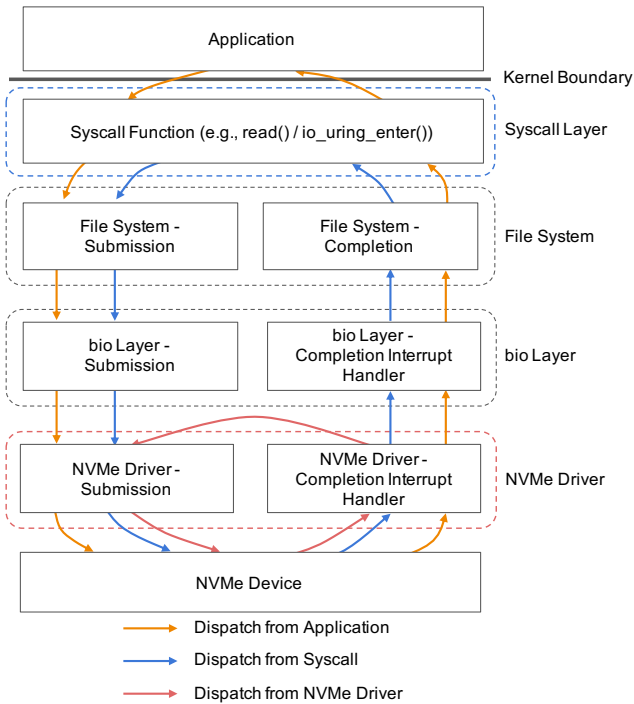


Figure 2: Dispatch paths for the application and the two kernel hooks.

The Overhead Source. To breakdown this software overhead, we measured the average latency of the different software layers when issuing a random 512 B read() system call with O_DIRECT on an Intel Optane SSD Gen 2 prototype (P5800X) on a 6-core i5-8500 3 GHz server with 16 GB of memory, Ubuntu 20.04, and Linux 5.8.0. We disable processor C-states and turbo boost and use the maximum performance governor. We use this setup throughout the paper. Table 1 shows that the layers that add the most latency are the file system (in this case ext4), followed by the transition from user space into kernel space.

Kernel-bypass allows applications to submit requests directly to devices, effectively eliminating all of these costs except the costs to post NVMe requests (“NVMe driver”) and the device latency itself [26, 36, 48, 49]. However, eliminating all of these layers comes at a high cost. The kernel can only delegate whole devices to processes, and applications must implement their own file systems [39]. Although NVMe supports SR-IOV, which allows partitioning one physical device into multiple namespaces and exposing each namespace as a separate virtual device, the maximum number of virtual devices is usually limited [8]. The same is true for building file system support for near-storage compute devices [44, 45]. Even when they re-implement file systems, they still cannot safely share files or capacity between distrusting processes.

Finally, lack of efficient application-level interrupt dispatching means these applications must resort to polling to be efficient at high load; this makes it impossible for them to efficiently share cores with other processes, resulting in wasteful busy polling when I/O utilization isn’t high enough.

There have been recent efforts to streamline kernel I/O submission costs via a new system call called *io_uring* [7]. It provides batched and asynchronous I/O submission/completion path that amortizes the costs of kernel boundary crossings and can avoid expensive scheduler interactions to block kernel threads. However, each submitted I/O still must pass through all of the kernel layers shown in Table 1. So, each I/O still incurs significant software overhead when accessing fast storage devices (we quantify this in §3).

3 BPF TO THE RESCUE?

With the rise of fast networks, Berkeley Packet Filter (BPF) has gained in popularity for efficient packet processing since it eliminates the need to copy each packet into userspace; instead, applications can safely operate on packets in the kernel. Common networking use cases include filtering packets [4, 5, 25, 33], network tracing [2, 3, 32], load balancing [4, 15], packet steering [30] and network security checks [12]. It has also been used as a way to avoid multiple network crossings when accessing disaggregated storage [38]. Linux supports BPF via the eBPF extension since Linux 3.15 [21]. The user-defined functions can be executed either using an interpreter or a just-in-time (JIT) compiler.

BPF for Storage. We envision similar use of BPF for storage by removing the need to traverse the kernel’s storage stack and move data back and forth between the kernel and user space when issuing dependent storage requests. Storage applications often generate many “auxiliary” I/O requests, such as index lookups. A key characteristic of these requests is that they occupy I/O bandwidth and CPU time to fetch data that is ultimately *not* returned to the user. For example, a search on a B-tree index is a series of pointer lookups that lead to the final I/O request for the user’s data page. Each of these lookups makes a roundtrip from the application through the kernel’s storage stack, only for the application to throw the data away after simple processing. Other similar use cases include database iterators that scan tables sequentially until an attribute satisfies a condition [18], aggregate queries that only need to return a final result (e.g. the maximum or average value), or graph databases that execute depth-first searches [9].

The Benefits. We design a benchmark that executes lookups on an on-disk B⁺-tree (which we call a B-tree for simplicity), a common data structure used to index databases [20, 34]. For simplicity, our experiments assume that the leaves of

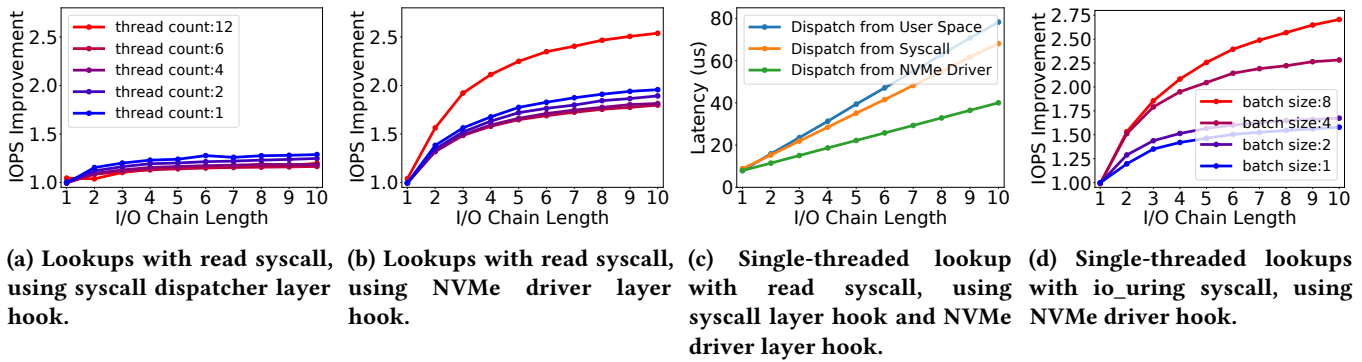


Figure 3: Search throughput improvement on B-tree with varying depth when reissuing lookups from different kernel layers.

the index contain user data rather than pointers [24, 40]. We assume each node of the B-tree occupies its own disk page, which means for a B-tree of depth d , a lookup requires reading d pages from disk. The core operation of a B-tree lookup is parsing the current page to find the offset of the next disk page, and issuing a read for the next page, a “pointer lookup”. Traditionally, a B-tree lookup requires d successive pointer lookups from userspace. To improve on the baseline, we reissue successive pointer lookups from one of two hooks in the kernel stack: the syscall dispatch layer (which mainly eliminates kernel boundary crossings) or the NVMe driver’s interrupt handler on the completion path (which eliminates nearly all of the software layers from the resubmission). Figure 2 shows the dispatch paths for the two hooks along with the normal user space dispatch path.

These two hooks call an eBPF function that simply issues a memcpy, simulating parsing the page to find the next pointer. Afterwards, the kernel hook dispatches a random read I/O request for the pointer extracted by the BPF function. This experiment represents an upper bound on the expected performance improvement from using eBPF to accelerate a single dependent sequence of storage I/Os.

Figures 3a and 3b show the throughput speedup of both BPF hooks relative to the baseline application traversal, and Figure 3c shows the latency of both hooks while varying the depth of the B-tree. When lookups are reissued from the syscall dispatch layer, the maximum speedup is 1.25 \times . The improvement is modest because each lookup still incurs the file system and block layer overhead; the speedup comes exclusively from eliminating kernel boundary crossings. As storage devices approach 1 μ s latencies, we expect greater speedups from this dispatch hook. On the other hand, reissuing from the NVMe driver makes subsequent I/O requests significantly less computationally expensive by bypassing nearly the entire software stack. Doing so achieves throughput improvements up to 2.5 \times and reduces latency by up to

49%. Relative throughput improvement actually goes down when adding more threads, because the baseline application also benefits from more threads until it reaches CPU saturation at 6 threads. Once the baseline hits CPU saturation, the computational savings due to reissuing at the driver becomes much more apparent. The throughput improvement from reissuing in the driver continues to scale with deeper trees, because each level of the tree compounds the number of requests that are issued cheaply.

What about io_uring? The previous experiments use Linux’s standard, synchronous read system call. Here, we repeat these experiments using the more efficient and batched io_uring submission path to drive B-tree lookups from a single thread. Like before, we reissue lookups within the NVMe driver and plot the throughput improvement against an application that simply batches I/Os using unmodified io_uring calls. Figure 3d shows the throughput speedup due to reissuing from within the driver relative to the application baseline.

As expected, increasing the batch size (number of system calls batched in each io_uring call), increases the speedup, since a higher batch size increases the number of requests that can be reissued at the driver. For example, for a batch size of 1 only 1 request (per B-tree level) can be reissued inexpensively, whereas for a batch size of 8, each B-tree level saves on 8 concurrent requests. Therefore, placing the hooks close to the device benefits both standard, synchronous read calls and more efficient io_uring calls. With deep trees, BPF coupled with io_uring delivers > 2.5 \times higher throughput; even three dependent lookups give 1.3–1.9 \times speedups.

4 A DESIGN FOR STORAGE BPF

Our experiments have given us reason to be optimistic about BPF’s potential to accelerate operations with fast storage devices; however, to realize these gains, I/O resubmissions must happen as early as possible, ideally within the kernel NVMe

interrupt handler itself. This creates significant challenges in using BPF to accelerate storage lookups for a practical system such as a key-value store.

We envision building a library that provides a higher level-interface than BPF and new BPF hooks in the Linux kernel as early in the storage I/O completion path as possible, similar to XDP [25]. This library would contain eBPF and helper functions that accelerate access and operations on popular data structures, such as B-trees and log-structured merge trees (LSM).

Within the kernel, these BPF functions that would be triggered in the NVMe driver interrupt handler on each block I/O completion. By giving these functions access to raw buffers holding block data, they could extract file offsets from blocks fetched from storage and immediately reissue an I/O to those offsets; they could also filter, project, and aggregate block data by building up buffers that they later return to the application. By pushing application-defined structures into the kernel these functions can traverse persistent data structures with limited application involvement. Unlike XN, where functions were tasked with implementing full systems, these storage BPF functions would mainly parse the storage block layout of application data structures and run simple, application-defined computations.

We outline some of the key design considerations and challenges for our preliminary design, which we believe we can realize without substantial re-architecture of the Linux kernel.

Installation & Execution. To accelerate dependent accesses, our library installs a BPF function using a special `ioctl`. Once installed, the application I/Os issued using that file descriptor are “tagged”; submissions to the NVMe layer propagate this tag. The kernel I/O completion path, which is triggered in the NVMe device interrupt handler, checks for this tag. For each tagged submission/completion, our NVMe interrupt handler hook passes the read block buffer into the corresponding BPF function.

When triggered, the function can perform a few actions. For example, it can extract a file offset from the block; then, it can “recycle” the NVMe submission descriptor and I/O buffer by calling a helper function that retargets the descriptor to the new offset and reissues it to the NVMe device submission queue. Hence, one I/O completion can determine the next I/O that should be submitted with minimal CPU cost or delay. This lets functions perform rapid traversals of structures without application-level involvement.

The function can also copy or aggregate data from the block buffer into its own buffers. This lets the function perform selection, projection, or aggregation to build results to return to the application. These buffers are hidden from and not returned to the application during intermediate calls

to the function, for example if the right block hasn’t been found yet. When the function completes it can indicate which buffer should be returned to the application.

Translation & Security. In Linux the NVMe driver doesn’t have access to file system metadata. If an I/O completes for a block at offset o_1 in a file, a BPF function might extract file offset o_2 as the next I/O to issue. However, o_2 is meaningless to the NVMe context, since it cannot tell which physical block this corresponds to without access to the file’s metadata and extents. Blocks could embed physical block addresses to avoid the need to consult the extents, but without imposing limits on these addresses, BPF functions could access any block on the device. Hence, a key challenge is imbuing the NVMe layer with enough information to efficiently and safely map file offsets to the file’s corresponding physical block offsets without restricting the file system’s ability to remap blocks as it chooses.

For simplicity and security in our design, each function only uses the file offsets in the file to which the `ioctl` attached the function. This ensures functions cannot access data that does not belong to the file. To do this without slow file system layer calls and without constraining the file system’s block allocation policy, we plan to only focus on cases *when the extents for a file do not change*.

We make the observation that many data center applications do not modify in-place persistent storage structures. For example, once an LSM-tree writes SSTable files to disk, they are immutable and their extents are stable [29]. Similarly, the index file extents remain nearly stable in on-disk B-tree implementations. In a 24 hour YCSB [28] (40% reads, 40% updates, 20% inserts, Zipfian 0.7) experiment on MariaDB running TokuDB [17], we found the index file’s extents only changed every 159 seconds on average with only 5 extent changes in 24 hours unmapping any blocks. Also note that in many index implementations, the index is stored on a single file, which helps further simplify our design.

We exploit the relative stability of file extents via a soft state cache of the extents at the NVMe layer. When the `ioctl` first installs the function on the file storing the data structure, its extents are propagated to the NVMe layer. If any block is unmapped from any of the file’s extents, a new hook in the file system triggers an invalidation call to the NVMe layer. Ongoing recycled I/Os are then discarded, and an error is returned to the application layer, which must rerun the `ioctl` to reset the NVMe layer extents before it can reissue tagged I/Os. This is a heavy-handed but simple approach. It leaves the file system almost entirely decoupled from the NVMe layer, and it places no restrictions on the file system block allocation policies. Of course, these invalidations need to be rare for the cache to be effective, but as the YCSB results

confirm, we expect this is true in most of the applications we target.

I/O Granularity Mismatches. When the BIO layer “splits” an I/O, e.g. across two discontinuous extents, it will generate multiple NVMe operations that complete at different times. We expect these cases to be rare enough that we can perform that I/O as a normal BIO and return the buffer and completion to the application. There, it can run the BPF function itself and restart the I/O chain with the kernel starting at the next “hop”. This avoids extra complexity in the kernel. Similarly, if application needs to generate more than one I/O in response to a single I/O completion, we propagate the completion up to the BIO layer which allocates and submits the multiple I/Os to the NVMe layer. This avoids returning to userspace.

Caching. As the caching of indices is often managed by the application [17, 27, 29], we assume the BPF traversal will not interact with the buffer cache directly and that applications manage caching and synchronizing with traversals. Cache eviction and management is increasingly done at the granularity of application-meaningful objects (e.g. individual data records) instead of whole pages. Our scheme fits well into this model, where BPF functions can return specific objects to the application rather than pages, after which the application can apply its own caching policies.

Concurrency. A write issued through the file system might only be reflected in the buffer cache and would not be visible to the BPF traversal. This could be addressed by locking, but managing application-level locks from within the NVMe driver, and especially from within the interrupt handler, can be expensive. Therefore, data structures that require fine-grained locking (e.g. lock coupling in B+trees [31]) require careful design to support BPF traversal.

To avoid read/write conflicts, we initially plan to target data structures that remain immutable (at least for a long period of time). Fortunately, many data structures have this property, including LSM SSTable files that remain immutable [29, 43], and on-disk B-trees that are not updated dynamically in-place, but rather in a batch process [17]. In addition, due to the difficulty of acquiring locks, we plan initially to only support read-only BPF traversals.

Fairness. BPF issued requests do not go through the file system or block layer, so there is no easy place to enforce fairness or QoS guarantees among processes. However, the default block layer scheduler in Linux is the noop scheduler for NVMe devices, and the NVMe specification supports command arbitration at hardware queues if fairness is a requirement [14]. Another challenge is that the NVMe layer may reissue an infinite number of I/O requests. The eBPF verifier prevents loops with unknown bounds [42], but

we would also need to prevent unbounded I/O loops at our NVMe hook.

For fairness purposes and to prevent unbounded traversals, we plan to implement a counter per process in the NVMe layer that will track the number of chained submissions, and set a bound on this counter. The counter’s values can periodically be passed to the BIO layer to account the number of requests.

5 CONCLUSIONS

BPF has the potential to significantly speed up dependent lookups to fast storage devices. However, it creates several key challenges, arising due to the loss of context when operating deep in the kernel’s storage stack. In this paper, we focused primarily on enabling the initial use case of index traversals. Notwithstanding, even for current fast NVMe devices (and more so for future ones), chaining a small number of requests using BPF provides significant gains. We envision a BPF for storage library could help developers offload many other standard storage operations to the kernel, such as aggregations, scans, compaction, compression and deduplication. We also believe that the interactions of BPF with the cache and scheduler policies create exciting future research opportunities.

6 ACKNOWLEDGMENTS

We would like to thank the anonymous reviewers for their comments. We thank Frank Hady and Andrew Ruffin for their generous equipment support, and to Adam Manzanares, Cyril Guyot, Qing Li and Filip Blagojevic for their feedback throughout the project. This work was in part supported by ARO grant W911NF-21-1-0078, ONR grants N00014-16-1-2263 and N00014-17-1-2788; NSF grants CNS-2104292, CNS-1750558 and CNS-15-64055; an Intel equipment donation; a JP Morgan Faculty Research Award; and a DiDi Faculty Research Award.

REFERENCES

- [1] 200Gb/s ConnectX-6 Ethernet Single/Dual-Port Adapter IC | NVIDIA. <https://www.mellanox.com/products/ethernet-adapter-ic/connectx-6-en-ic>.
- [2] bcc. <https://github.com/iovisor/bcc>.
- [3] bpftrace. <https://github.com/iovisor/bpftrace>.
- [4] Cilium. <https://github.com/cilium/cilium>.
- [5] Cloudflare architecture and how BPF eats the world. <https://blog.cloudflare.com/cloudflare-architecture-and-how-bpf-eats-the-world/>.
- [6] eBPF. <https://ebpf.io/>.
- [7] Efficient io with io_uring. https://kernel.dk/io_uring.pdf.
- [8] Flash storage for servers is getting faster and smarter. <https://www.samsung.com/semiconductor/newsroom/tech-leadership/flash-storage-for-servers-is-getting-faster-and-smarter/>.
- [9] Graph databases for beginners: Graph search algorithm basics. <https://neo4j.com/blog/graph-search-algorithm-basics/>.
- [10] Intel® Optane™ SSD DC P5800X Series. <https://ark.intel.com/content/www/us/en/ark/products/201859/intel-optane-ssd-dc-p5800x-series-1-6tb-2-5in-pcie-x4-3d-xpoint.html>.
- [11] LevelDB. <https://github.com/google/leveldb>.
- [12] MAC and Audit policy using eBPF. <https://lkml.org/lkml/2020/3/28/479>.
- [13] NGD systems newport platform. <https://www.ngdsystems.com/technology/computational-storage>.
- [14] NVMe base specification. https://nvmexpress.org/wp-content/uploads/NVM-Express-1_4b-2020.09.21-Ratified.pdf.
- [15] Open-sourcing Katran, a scalable network load balancer. <https://engineering.fb.com/2018/05/22/open-source/open-sourcing-katran-a-scalable-network-load-balancer/>.
- [16] Optimizing Software for the Next Gen Intel Optane SSD P5800X. <https://www.intel.com/content/www/us/en/events/memory-and-storage.html?videoId=6215534787001>.
- [17] Percona TokuDB. <https://www.percona.com/software/mysql-database/percona-tokudb>.
- [18] RocksDB iterator. <https://github.com/facebook/rocksdb/wiki/Iterator>.
- [19] SmartSSD computational storage drive. <https://www.xilinx.com/applications/data-center/computational-storage/smartssd.html>.
- [20] SQL server index architecture and design guide. <https://docs.microsoft.com/en-us/sql/relational-databases/sql-server-index-design-guide>.
- [21] A thorough introduction to eBPF. <https://lwn.net/Articles/740157/>.
- [22] Toshiba memory introduces XL-FLASH storage class memory solution. <https://business.kioxia.com/en-us/news/2019/memory-20190805-1.html>.
- [23] Ultra-Low Latency with Samsung Z-NAND SSD. https://www.samsung.com/semiconductor/global.semi.static/Ultra-Low_Latency_with_Samsung_Z-NAND_SSD-0.pdf.
- [24] WiredTiger storage engine. <https://docs.mongodb.com/manual/core/wiredtiger/>.
- [25] XDP. <https://www.iovisor.org/technology/xdp>.
- [26] Adam Belay, George Prekas, Ana Klimovic, Samuel Grossman, Christos Kozyrakis, and Edouard Bugnion. IX: A protected dataplane operating system for high throughput and low latency. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*, pages 49–65, Broomfield, CO, October 2014. USENIX Association.
- [27] Badrish Chandramouli, Guna Prasaad, Donald Kossmann, Justin Levandoski, James Hunter, and Mike Barnett. FASTER: A concurrent key-value store with in-place updates. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD '18*, page 275–290, New York, NY, USA, 2018. Association for Computing Machinery.
- [28] 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.
- [29] Siying Dong, Mark Callaghan, Leonidas Galanis, Dhruva Borthakur, Tony Savor, and Michael Strum. Optimizing space amplification in RocksDB. In *CIDR*, volume 3, page 3, 2017.
- [30] Pekka Enberg, Ashwin Rao, and Sasu Tarkoma. Partition-aware packet steering using XDP and eBPF for improving application-level parallelism. In *Proceedings of the 1st ACM CoNEXT Workshop on Emerging In-Network Computing Paradigms, ENCP '19*, page 27–33, New York, NY, USA, 2019. Association for Computing Machinery.
- [31] Goetz Graefe. A Survey of B-Tree Locking Techniques. *ACM Transactions on Database Systems*, 35(3), July 2010.
- [32] Brendan Gregg. *BPF Performance Tools*. Addison-Wesley Professional, 2019.
- [33] Toke Høiland-Jørgensen, Jesper Dangaard Brouer, Daniel Borkmann, John Fastabend, Tom Herbert, David Ahern, and David Miller. The express data path: Fast programmable packet processing in the operating system kernel. In *Proceedings of the 14th international conference on emerging networking experiments and technologies*, pages 54–66, 2018.
- [34] Stratos Idreos, Niv Dayan, Wilson Qin, Mali Akmanalp, Sophie Hilgard, Andrew Ross, James Lennon, Varun Jain, Harshita Gupta, David Li, et al. Design continuums and the path toward self-designing key-value stores that know and learn. In *CIDR*, 2019.
- [35] M. Frans Kaashoek, Dawson R. Engler, Gregory R. Ganger, Hector M. Briceño, Russell Hunt, David Mazières, Thomas Pinckney, Robert Grimm, John Jannotti, and Kenneth Mackenzie. Application performance and flexibility on exokernel systems. In *Proceedings of the Sixteenth ACM Symposium on Operating Systems Principles, SOSP '97*, page 52–65, New York, NY, USA, 1997. Association for Computing Machinery.
- [36] Kostis Kaffes, Timothy Chong, Jack Tigar Humphries, Adam Belay, David Mazières, and Christos Kozyrakis. Shinjuku: Preemptive scheduling for μ second-scale tail latency. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*, pages 345–360, 2019.
- [37] 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.
- [38] Kornilios Kourtis, Animesh Trivedi, and Nikolas Ioannou. Safe and Efficient Remote Application Code Execution on Disaggregated NVM Storage with eBPF. *arXiv preprint arXiv:2002.11528*, 2020.
- [39] Jing Liu, Andrea C Arpaci-Dusseau, Remzi H Arpaci-Dusseau, and Sudarsun Kannan. File systems as processes. In *11th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 19)*, 2019.
- [40] Lanyue Lu, Thanumalayan Sankaranarayanan Pillai, Hariharan Gopalakrishnan, Andrea C Arpaci-Dusseau, and Remzi H Arpaci-Dusseau. Wisckey: Separating keys from values in SSD-conscious storage. *ACM Transactions on Storage (TOS)*, 13(1):1–28, 2017.
- [41] Michael Marty, Marc de Kruijf, Jacob Adriaens, Christopher Alfeld, Sean Bauer, Carlo Contavalli, Michael Dalton, Nandita Dukkupati, William C. Evans, Steve Gribble, Nicholas Kidd, Roman Kononov, Gautam Kumar, Carl Mauer, Emily Musick, Lena Olson, Erik Rubow, Michael Ryan, Kevin Springborn, Paul Turner, Valas Valancius, Xi Wang, and Amin Vahdat. Snap: A microkernel approach to host networking. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP '19*, page 399–413, New York, NY, USA, 2019. Association for Computing Machinery.

- [42] Luke Nelson, Jacob Van Geffen, Emina Torlak, and Xi Wang. Specification and verification in the field: Applying formal methods to BPF just-in-time compilers in the Linux kernel. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 41–61. USENIX Association, November 2020.
- [43] Patrick O’Neil, Edward Cheng, Dieter Gawlick, and Elizabeth O’Neil. The log-structured merge-tree (LSM-tree). *Acta Informatica*, 33(4):351–385, 1996.
- [44] Weikang Qiao, Jieqiong Du, Zhenman Fang, Michael Lo, Mau-Chung Frank Chang, and Jason Cong. High-throughput lossless compression on tightly coupled CPU-FPGA platforms. In *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pages 37–44. IEEE, 2018.
- [45] Zhenyuan Ruan, Tong He, and Jason Cong. INSIDER: Designing in-storage computing system for emerging high-performance drive. In *2019 USENIX Annual Technical Conference (USENIX ATC 19)*, pages 379–394, 2019.
- [46] Sudharsan Seshadri, Mark Gahagan, Sundaram Bhaskaran, Trevor Bunker, Arup De, Yanqin Jin, Yang Liu, and Steven Swanson. Willow: A user-programmable SSD. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*, pages 67–80, Broomfield, CO, October 2014. USENIX Association.
- [47] Shin-Yeh Tsai and Yiyang Zhang. Lite kernel rdma support for data-center 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.
- [48] Ziyi Yang, James R Harris, Benjamin Walker, Daniel Verkamp, Changpeng Liu, Cunyin Chang, Gang Cao, Jonathan Stern, Vishal Verma, and Luse E Paul. SPDK: A development kit to build high performance storage applications. In *2017 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*, pages 154–161. IEEE, 2017.
- [49] Irene Zhang, Jing Liu, Amanda Austin, Michael Lowell Roberts, and Anirudh Badam. I’m not dead yet! the role of the operating system in a kernel-bypass era. In *Proceedings of the Workshop on Hot Topics in Operating Systems, HotOS '19*, page 73–80, New York, NY, USA, 2019. Association for Computing Machinery.