

# Multi-version Indexing in Flash-based Key-Value Stores

Pulkit A. Misra  
Duke University

Jeffrey S. Chase  
Duke University

Johannes Gehrke  
Microsoft Corporation

Alvin R. Lebeck  
Duke University

## Abstract

Maintaining multiple versions of data is popular in key-value stores since it increases concurrency and improves performance. However, designing a multi-version key-value store entails several challenges, such as additional capacity for storing extra versions and an indexing mechanism for mapping versions of a key to their values. We present SKIMPYFTL, a FTL-integrated multi-version key-value store that exploits the remap-on-write property of flash-based SSDs for multi-versioning and provides a tradeoff between memory capacity and lookup latency for indexing.

## 1 Introduction

Transactional key-value stores use a multi-version storage to increase concurrency [7] and reduce abort rate [3] as reads can be satisfied from a consistent snapshot in the past (old versions), while writes create new versions. Figure 1 illustrates the impact of single vs multi-versioning on transaction abort rate for a social network application; workload and methodology are described in §4. As seen from the figure, multi-versioning reduces abort rate by 2× vs single-version storage, and its benefit increases with offered load.

However, there are several challenges in designing a multi-version storage, including: 1) additional capacity, 2) index for mapping versions to values, and 3) version management. First, the extra versions require additional capacity, which can be prohibitively expensive with in-memory (DRAM) storage systems. Second, these systems need an index to map versions of a key to their value so reads can be serviced from a consistent snapshot. A naïve approach is to store the entire index in DRAM; this approach provides the lowest lookup latencies but has a high space overhead. An efficient indexing technique needs to find a tradeoff between lookup latencies and DRAM requirement for indexing. Third, multi-versioning necessitates a version management (garbage collection) scheme for effective capacity utilization. The scheme needs to strike a balance between capacity reclamation by discarding old versions and servicing reads from a consistent snapshot.

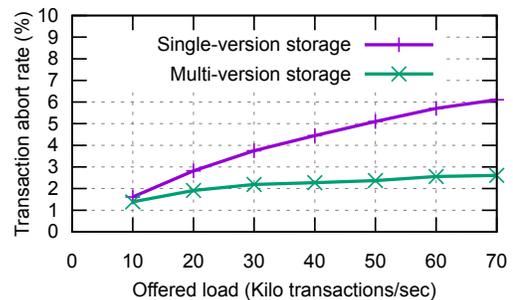


Figure 1: Transaction abort rate with single vs multi-version storage

Our previous work SEMEL [17] addresses the capacity and version management challenge with designing a multi-version storage system. It uses Solid State Drives (SSDs) for storing multiple versions of data and a watermark-based version management for effective capacity utilization as watermarks provide a bound on the oldest version that can be read by the application. This paper addresses the challenge with indexing.

SSDs are a more attractive proposition for multi-versioning than DRAM or newer non-volatile memory (NVM) technologies (e.g., Intel Optane) because they provide TB capacity per drive for less than \$1.00/GB. In addition, new standards for Software-Defined Flash (SDF) enable flash-based storage systems customized for application requirements [11, 12, 16, 17, 23, 26]. These advances in flash storage enable storing more data per server, while delivering better performance compared to spinning disks and at a lower cost compared to DRAM and NVMs.

Furthermore, SEMEL exploits an intrinsic property of flash-based SSDs — remap-on-write — to implement multi-versioning in an SSDs Flash Translation Layer (FTL) using SDF. SEMEL’s approach removes abstractions and provides better performance compared to a naïve approach of stacking a multi-version software layer over a standard FTL. However, SEMEL incurs a high space overhead since it maintains the

entire index for mapping a version of a key to its location on flash in host DRAM.

This paper presents SKIMPYFTL, a system that addresses all the 3 challenges with designing a multi-version key-value storage system. It builds on top of SEMEL by leveraging SSDs for FTL-integrated multi-versioning along with a watermark-based version management scheme and addresses the challenge with indexing by providing a tradeoff between DRAM capacity and lookup latency. SKIMPYFTL uses a hash table for mapping key versions to their values. It follows SkimpyStash [9] in offloading hash collision list on flash to reduce DRAM requirement for indexing and adds support for multi-versioning using version pointers, which are also stored on flash along with the collision list.

A SKIMPYFTL prototype utilizing LightNVM Open-Channel SSD emulation framework [4] reveals SKIMPYFTL provides 72-91% throughput of SEMEL for read-dominant key-value workloads (75-100% reads), while reducing the memory requirement for indexing by a factor of  $0.95\times$ . For a transactional YCSB [6] workload, SKIMPYFTL provides 85% peak throughput of SEMEL. Finally, SKIMPYFTL outperforms a naïve multi-version key-value store implemented over a standard FTL on both workloads.

This paper is a first step in exploring the design space of memory-efficient indexing in multi-version flash-based key-value stores. Our overarching goal for this work is designing a distributed flash management framework. In the future, we plan to explore other data structures, such as LSM trees [19], for implementing multi-versioning inside the FTL. We also plan to further explore garbage collection strategies while maintaining multiple versions of data.

## 2 Background

This section summarizes the internals and the remap-on-write property of a NAND-flash SSD (§2.1), provides background on Software-Defined Flash (SDF, §2.2) and shows how SEMEL uses SDF in a multi-version storage system (§2.3).

### 2.1 Internals of a NAND Flash-based SSD

NAND flash memory in an SSD is organized as an array of blocks where each block contains some number of pages. Typically pages are 2-16KB in size and each block contains 128-256 pages. The page size is the smallest unit for reads and writes. A page on flash must first be erased before it can be overwritten. However, erase operations on flash occur only at a block granularity.

To accommodate these properties, flash-based SSDs use a Flash Translation Layer (FTL) to map logical addresses dynamically to physical locations. This level of indirection allows the FTL to remap a logical block to a clean physical page on each write (*remap-on-write*), leaving the old value in place

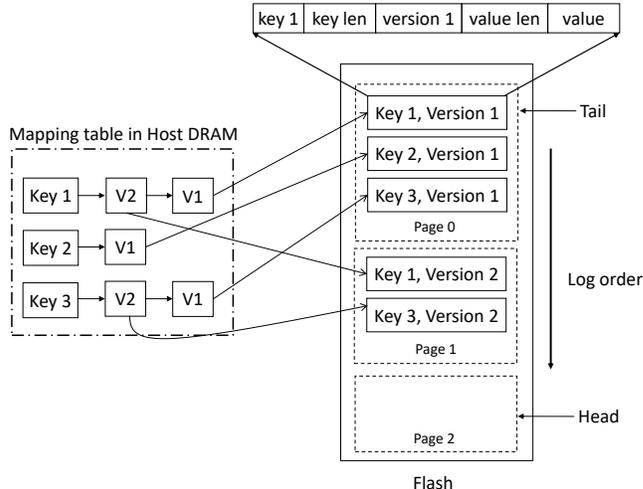


Figure 2: Mapping table and data layout on flash in SEMEL

pending garbage collection. The FTL’s garbage collector also remaps values as needed to clean blocks to erase.

This remap-on-write behavior naturally leaves previous versions of data in place for reading [24]. SEMEL and SKIMPYFTL exploit this property as a foundation for an FTL-integrated multi-version key-value store.

### 2.2 Software-Defined Flash

The recently proposed Software-Defined Flash (SDF) structure moves some FTL functionality to host software [5, 12, 20, 21]. Several vendors provide some form of SDF (e.g., CNEX Labs, SanDisk/FusionIO, Radian Memory).

SDF enables several optimizations across traditional system boundaries. First, flash-based key-value stores can map keys directly to pages on flash rather than addressing the SSD as a block device, eliminating one level of indirection [12, 16, 23, 26]. Next, customized mapping techniques may exploit system-specific information to improve performance and/or provide new functionality, such as snapshot capability [24]. Another work leverages SDF to exploit the inherent parallelism of SSDs by mapping log-structured merge (LSM) tree [19] operations to different SSD channels [20].

### 2.3 SEMEL Flash Translation Layer

SEMEL is a lightweight multi-version key-value store based on SDF. It writes new values in a log-structured fashion [22], densely packed in pages on flash. Figure 2 shows the mapping table and data layout in SEMEL. As seen from the figure, SEMEL maintains a linked list in DRAM with an entry for each version of a key. Each version is assigned a 64-bit create timestamp and maps directly to a page on flash and the version’s offset within the page, removing a level of indirection.

SEMEL’s DRAM-based mapping table is prohibitively expensive. Each version entry is 20B in size (4B page address, 8B version timestamp and 8B pointer to prior version): for a 1 TB SSD and 512B key-value pairs, SEMEL consumes 40 GB of DRAM to map the entire SSD. The goal of SKIMPYFTL is to provide DRAM-efficient version mapping with performance as close as possible to SEMEL.

### 3 SkimpyFTL

This section describes how SKIMPYFTL addresses the need for memory-efficient dynamic indexing of a multi-version store in SEMEL. Our approach combines a flash-based mapping table (§3.1) with a DRAM-based mapping translation cache for hot keys (§3.2). The two structures operate together to handle reads and writes efficiently in the common case (§3.3). The scheme also extends the garbage collection protocol (§3.4).

#### 3.1 Mapping Table

SKIMPYFTL indexes the key-value pairs on flash with a hash-based mapping table whose buckets are rooted in host DRAM. Multiple key-value pairs may hash to the same bucket; SKIMPYFTL handles such collisions using *linear chaining*, where key-value pairs in the same bucket are chained using a linked list. Rather than maintaining this collision list in DRAM, SKIMPYFTL offloads it to flash; a hash table bucket in DRAM points to the head of the linear list on flash, and each entry in the list on flash points to the next entry. This approach is inspired by a prior work [9]. SKIMPYFTL allows configuring the number of hash buckets to achieve a desired balance of DRAM cost and lookup time.

Figure 3 shows the mapping table and the data layout on flash. SKIMPYFTL writes data versions to flash in a log-structured fashion, as in SEMEL. In addition to the usual fields for each key-value pair (version), SKIMPYFTL stores two pointers: a *hash next* pointer to the next entry in the bucket’s collision list (hash chain) and a *prior version* pointer to a prior (older) version of the key. The prior version pointer is an optimization: SKIMPYFTL can traverse a hash chain to find a requested version, but it is often faster to traverse the prior version pointer from a later version of the desired key, as described next.

#### 3.2 Mapping Translation Cache

Prior characterization studies suggest that datacenter workloads tend to be read-dominated [2, 18]. Furthermore, the popularity of data items in real-world workloads often follows a power law distribution, where a small subset of the keys receive a large portion of the accesses [2, 6]. Such workloads would cause significant read amplification to traverse the collision lists. In particular, a key that is updated infrequently

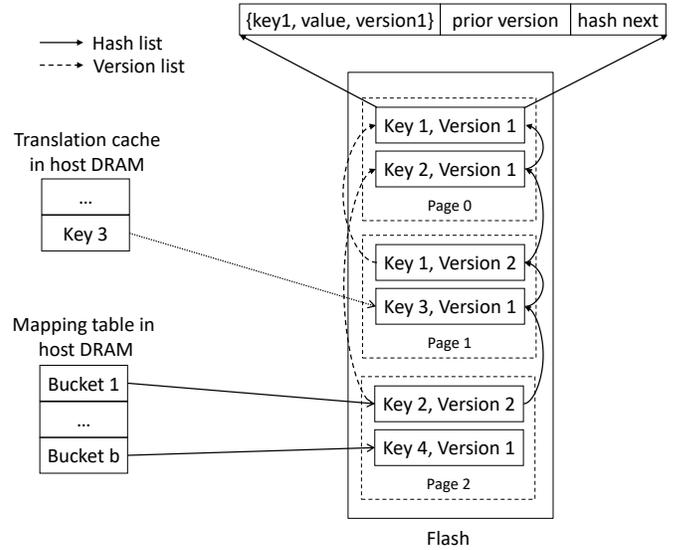


Figure 3: Mapping table, translation cache and data layout on flash in SKIMPYFTL

tends to migrate to the end of the list, since new versions are written to the front of the list.

To mitigate this cost, SKIMPYFTL caches key translations in a *mapping translation cache* in DRAM. A translation cache entry for a key stores the location of its latest version on flash.

#### 3.3 Request Life Cycle

Figure 3 illustrates how the DRAM-based translation cache operates in conjunction with the mapping table to handle GET (read) and PUT (write) requests efficiently.

A GET request first does a lookup to locate the latest version of the requested key. Then, if the request is a snapshot read for a previous version, it traverses the prior version pointer(s) to locate the requested version. A lookup for a hot or recent key hits in the translation cache, which returns the pointer to the latest version. On a cache miss, SKIMPYFTL hashes the key to a bucket and then follows the hash chain until it finds an entry for the key, which is the latest version. It then caches the mapping in the translation cache.

For example, for a GET of {key 1, version 1} in Figure 3, SKIMPYFTL first hashes key 1 to bucket 1 and traverses the hash chain: {key 2, version 2} → {key 3, version 1} → {key 1, version 2}, until it encounters key 1. It then updates the translation cache to point to the latest version of key 1.

A PUT request hashes the key to a bucket and links the new value to the front of the hash chain: it sets its hash next field to the current index pointed to by the bucket and updates the current index. To populate the prior version field, SKIMPYFTL probes the translation cache for the most recent version of the key. This lookup typically results in a hit in the common case of a read-modify-write operation on the key. On a miss,

SKIMPYFTL sets the prior version field to a special value that indicates to a GET request that prior versions may exist and must be retrieved by searching further in the chain. It fills in the missing prior version pointers during remapping (§3.4).

For example, for a PUT request for key 2 in Figure 3, SKIMPYFTL writes the new version to the end of the log on flash, hashes it to bucket 1, sets the hash next of key 2 to point to key 3 (on page 1) and updates bucket 1 to point to the new version. It sets the prior version field of key 2 to point to a prior version.

### 3.4 Garbage Collection

The garbage collection process starts from the tail of the log and remaps valid versions (key-value pairs) to the head of the log. For each key, SKIMPYFTL retains the youngest version with a timestamp less than the current *watermark*, and discards older versions. The watermark is a timestamp that advances continuously. In the SEMEL transactional key-value store, the watermark is the minimum timestamp that could appear in any future request from a client, following Centiman [10]. Each client periodically passes its timestamp for its last acknowledged operation. The minimum of these timestamps is the watermark.

SKIMPYFTL extends SEMEL’s garbage collection to use the bucket collision lists and version lists, and to maintain their integrity. To determine whether to remap or discard a version, SKIMPYFTL probes the translation cache and version list for a more recent version that is younger than the watermark. On a miss, it must traverse a bucket hash chain. For each remapped version, it must also update the hash next pointer of its predecessor in the chain to point at the new location. For this reason, SKIMPYFTL garbage collects an entire hash bucket at a time by sweeping its hash chain and remapping retained versions in reverse temporal order, updating their pointers in the usual way as it goes.

To maintain fidelity of the hash chains and version lists, SKIMPYFTL blocks any new writes to a bucket during the process of traversing its hash chain and remapping valid data. The remapping process packs the retained versions of a bucket’s keys densely into flash pages, reducing lookup time for later GET requests.

## 4 Evaluation

Here we present the preliminary results for our prototype implementation of SKIMPYFTL and compare it with SEMEL. To elucidate the advantages of implementing multi-versioning within the FTL, we also compare against a separate multi-version KV store on top of a standard FTL; we refer to this system as VFTL. We use SEMEL with multi-versioning disabled for single-version store.

**Implementation.** We use a modified Open-Channel SSD framework [4] from our prior work [17] for all FTL implemen-

tations. In software-only mode, the emulator supports storing data values in DRAM, and provides IOCTLS for get, put and erase functionality for flash blocks. It also allows specifying latencies for read page, write page and block erase operations.

All our FTL implementations use 32 bits (4 bytes) for storing the location of a key on flash. To allow an FTL to pack multiple key-value pair within a page; we divide each page into fixed number of chunks and use 3 out of the 32 bits for storing the start chunk for a key in a page. The remaining 29 bits are used for addressing a page.

**Workload.** We use two types of workloads for evaluation: 1) key-value operations workload, and 2) transactional YCSB [6] workload. We implement a micro-benchmark for the key-value operations workload; the micro-benchmark issues non-transactional get and put requests to single keys, for a varying get request percentage (%). Popularity of keys in the benchmark is controlled using a zipfian distribution; we set the zipfian coefficient  $\alpha$  to 0.99 in our micro-benchmark, a frequently used value in key-value workloads [2, 6]. Our transactional YCSB [6] workload models a social network application, where the data of popular users is read more often and users have different rate of posting updates. The workload models this behavior by using different values of  $\alpha$  for controlling popularity of keys in read-only ( $\alpha_r$ ) and read-write ( $\alpha_{rw}$ ) transactions. Each read-only transaction accesses from 1 to 10 keys, and a read-write transaction operates on 1 to 5.

**Experimental Setup.** All experiments are run on a server with a 16 core Intel Xeon E5-2640 processor clocked and 128 GB DRAM. The SSD emulator is backed by 32 GB DRAM, with a hardware queue depth of 128. The SSD has a page size of 4KB and there are 32 pages in a block. A flash page read, write time is 50 and 100  $\mu$ s respectively and it takes 1 ms to erase a flash block. To ease garbage collection, all FTL implementations reserve 10% capacity for remapping data.

For all experiments, we populate the system with 20M keys; the key size is 16B and packed data on flash is 512B, which includes key, value, version, and pointers to prior version and next entry in hash collision chain. As a flash page is 4KB in size, we employ a *packing logic* in the FTL that waits for up to 1 ms (tunable) to pack data of multiple keys into a page. There are up to 8 keys packed into a page, which increases the garbage collection overhead per page. Each run is of 15 mins and we pre-condition the SSD so garbage collection runs in the background in all experiments that perform writes.

Our SEMEL prototype stores the entire mapping table in DRAM. It needs  $\sim 67$ M mappings ( $=32\text{GB}/512\text{B}$ ) to address the entire SSD, and each mapping is 20B (4B page address, 8B version timestamp, and 8B prior version pointer) in size. Through sensitivity analysis, we set the translation cache size to 10% of the number of keys in the workload, which gives  $\sim 90\%$  hit rate for our current workload. Figure 4 shows the impact of average number of keys / bucket on throughput with our key-value micro-benchmark with 10% writes. Based on the results, we size the mapping table to have 5 keys / bucket

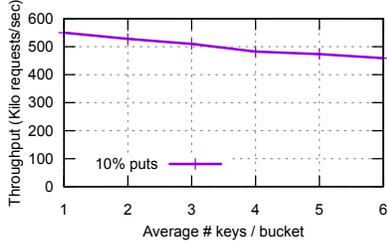


Figure 4: Impact of # keys / bucket

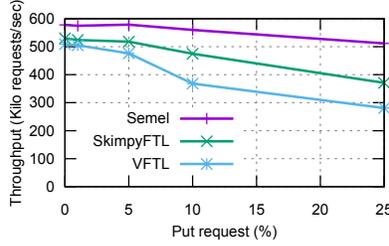


Figure 5: K-V operations throughput

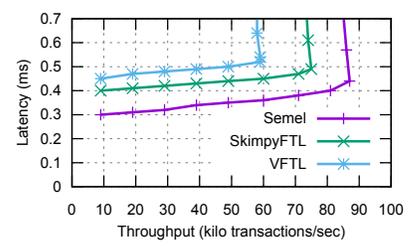


Figure 6: Txns: throughput vs latency

i.e., 4M buckets. Each entry in the mapping table is 4B (page address), and 20B (4B page address, 16B LRU pointers) in the translation cache. SKIMPYFTL uses  $\sim 5\%$  of the memory used by SEMEL.

#### 4.1 Key-Value Workload

We first evaluate the performance of all systems with a non-transactional key-value workload. Figure 5 shows the throughput for varying put request percentage (%) with our key-value micro-benchmark. The throughput of all 3 systems — SEMEL, VFTL and SKIMPYFTL— drops as put request percentage is increased because writes increase the garbage collection overhead. As expected, SEMEL provides the highest throughput as it implements multi-versioning in flash and maps all versions of all keys in DRAM. SKIMPYFTL provides from 72% - 91% of the throughput of SEMEL, while only using 5% of the memory. The throughput degradation is worse in SKIMPYFTL as put request percentage increases because garbage collection overhead — traversing entire hash collision chains for determining data to discard and remap — is more frequent. We plan to explore other approaches to garbage collection, with lower overheads, in future. Interestingly, the throughput of VFTL is even lower than SKIMPYFTL. VFTL provides the lowest throughput as it suffers from log stacking [25] — the log in multi-version layer and standard FTL operate independently, which leads to uncoordinated garbage collection and randomization of writes.

#### 4.2 Transactional Workload

To evaluate performance for a transactional workload, we stack a layer over the 3 systems that supports transactions using MVCC [3]. The memory overhead of the layer is the same for all systems. Our YCSB [6] workload issues 90% read-only transactions, with zipfian coefficient for read-only and read-write transactions set to  $\alpha_r = 0.99$  and  $\alpha_{rw} = 0.75$ , respectively.

Figure 6 shows the throughput and latency with the 3 storage systems for an increasing offered load. The trend in performance of the systems is similar to the key-value workload. SEMEL provides the best performance; SKIMPYFTL pro-

vides 85% of the peak throughput of SEMEL. VFTL provides the lowest peak throughput and highest latency. Thus, showing the disadvantage of the naïve approach of implementing multi-versioning, agnostic of the underlying storage medium. All approaches have similar transaction commit rates before saturation ( $\sim 96\%$ ).

### 5 Related Work

Prior works have proposed indexes for flash-based key-value stores [1, 8, 9, 13–16]. However, none of these works support multi-versioning. FD-tree [14] and WiscKey [13] use a tree-based approach for indexing key-value pairs on SSD, whereas FAWN [1], NVMKV [16] and FlashStore [8] use a hash table for indexing. SILT [15] stores data in multiple tiers; as key-value pairs age, they are compacted with other pairs and transitioned to other skimpy memory-optimized tiers. SKIMPYFTL is inspired from SkimpyStash [9], both approaches use configurable number of hash buckets to reduce the DRAM requirement for the mapping table, however SkimpyStash does not maintain multiple versions, nor does it have a translation cache.

### 6 Conclusion

We present SKIMPYFTL, a FTL-integrated multi-version key-value store that exploits remap-on-write property of flash-based SSDs for maintaining multiple versions of data and provides a tradeoff between DRAM capacity and lookup latency for indexing. Our evaluation reveals SKIMPYFTL provides 72-91% throughput of SEMEL, while reducing memory requirement by a factor of 0.95x. We also show the benefit of implementing multi-versioning in the FTL. In the future, we plan to explore other data structures for indexing, garbage collection strategies and storing data in multiple tiers, based on versions, frequency and type of access to keys.

### Acknowledgments

This work is supported in part by the National Science Foundation (CNS-1616947).

## References

- [1] David G. Andersen, Jason Franklin, Michael Kaminsky, Amar Phanishayee, Lawrence Tan, and Vijay Vasudevan. FAWN: A fast array of wimpy nodes. In *Proceedings of the ACM SIGOPS 22Nd Symposium on Operating Systems Principles, SOSP '09*, pages 1–14, New York, NY, USA, 2009. ACM.
- [2] Berk Atikoglu, Yuehai Xu, Eitan Frachtenberg, Song Jiang, and Mike Paleczny. Workload analysis of a large-scale key-value store. In *Proceedings of the 12th ACM SIGMETRICS/PERFORMANCE Joint International Conference on Measurement and Modeling of Computer Systems, SIGMETRICS '12*, pages 53–64, New York, NY, USA, 2012. ACM.
- [3] Philip A. Bernstein and Nathan Goodman. Multiversion concurrency control - theory and algorithms. *ACM Transactions on Database Systems*, 8(4):465–483, December 1983.
- [4] Matias Björling, Javier González, and Philippe Bonnet. LightNVM: The Linux Open-Channel SSD Subsystem. In *15th USENIX Conference on File and Storage Technologies (FAST)*. USENIX, 2017.
- [5] Adrian M. Caulfield, Todor I. Mollov, Louis Alex Eisner, Arup De, Joel Coburn, and Steven Swanson. Providing safe, user space access to fast, solid state disks. In *Proceedings of the Seventeenth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS XVII*, pages 387–400, New York, NY, USA, 2012. ACM.
- [6] 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*, pages 143–154, New York, NY, USA, 2010. ACM.
- [7] James C Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, Jeffrey John Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, et al. Spanner: Google’s Globally Distributed Database. In *OSDI*, 2012.
- [8] Biplob Debnath, Sudipta Sengupta, and Jin Li. Flashstore: High throughput persistent key-value store. *Proc. VLDB Endow.*, 3(1-2):1414–1425, September 2010.
- [9] Biplob Debnath, Sudipta Sengupta, and Jin Li. SkimpyS-tash: RAM space skimpy key-value store on flash-based storage. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data, SIGMOD '11*, pages 25–36, New York, NY, USA, 2011. ACM.
- [10] Bailu Ding, Lucja Kot, Alan Demers, and Johannes Gehrke. Centiman: Elastic, high performance optimistic concurrency control by watermarking. In *Proceedings of the Sixth ACM Symposium on Cloud Computing, SoCC '15*, pages 262–275. ACM, 2015.
- [11] Jian Huang, Anirudh Badam, Laura Caulfield, Suman Nath, Sudipta Sengupta, Bikash Sharma, and Moinuddin K. Qureshi. Flashblox: Achieving both performance isolation and uniform lifetime for virtualized ssds. In *FAST'17*. USENIX, February 2017.
- [12] William K. Josephson, Lars A. Bongo, Kai Li, and David Flynn. Dfs: A file system for virtualized flash storage. *Trans. Storage*, 6(3):14:1–14:25, September 2010.
- [13] Lanyue Lu and Thanumalayan Sankaranarayanan Pillai and Andrea C. Arpaci-Dusseau and Remzi H. Arpaci-Dusseau. WiscKey: Separating Keys from Values in SSD-conscious Storage. In *14th USENIX Conference on File and Storage Technologies (FAST 16)*, pages 133–148, Santa Clara, CA, February 2016. USENIX Association.
- [14] Yinan Li, Bingsheng He, Robin Jun Yang, Qiong Luo, and Ke Yi. Tree indexing on solid state drives. *Proc. VLDB Endow.*, 3(1-2):1195–1206, September 2010.
- [15] Hyeontaek Lim, Bin Fan, David G. Andersen, and Michael Kaminsky. Silt: A memory-efficient, high-performance key-value store. In *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles, SOSP '11*, pages 1–13, New York, NY, USA, 2011. ACM.
- [16] Leonardo Mármol, Swaminathan Sundararaman, Nisha Talagala, Raju Rangaswami, Sushma Devendrappa, Bharath Ramsundar, and Sriram Ganesan. NVMKV: a scalable and lightweight flash aware key-value store. In *6th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 14)*, 2014.
- [17] Pulkit A. Misra, Jeffrey S. Chase, Johannes Gehrke, and Alvin R. Lebeck. Enabling lightweight transactions with precision time. In *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '17*, pages 779–794, New York, NY, USA, 2017. ACM.
- [18] Rajesh Nishtala, Hans Fugal, Steven Grimm, Marc Kwiatkowski, Herman Lee, Harry C. Li, Ryan McElroy, Mike Paleczny, Daniel Peek, Paul Saab, David Stafford, Tony Tung, and Venkateshwaran Venkataramani. Scaling memcache at facebook. In *Presented as part of the 10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*, pages 385–398, Lombard, IL, 2013. USENIX.

- [19] Patrick O’Neil, Edward Cheng, Dieter Gawlick, and Elizabeth O’Neil. The log-structured merge-tree (lsm-tree). *Acta Inf.*, 33(4):351–385, June 1996.
- [20] Jian Ouyang, Shiding Lin, Song Jiang, Zhenyu Hou, Yong Wang, and Yuanzheng Wang. SDF: Software-defined Flash for Web-scale Internet Storage Systems. In *Proceedings of the 19th International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS ’14*, pages 471–484, New York, NY, USA, 2014. ACM.
- [21] Xiangyong Ouyang, David Nellans, Robert Wipfel, David Flynn, and Dhabaleswar K. Panda. Beyond block I/O: Rethinking traditional storage primitives. In *Proceedings of the 2011 IEEE 17th International Symposium on High Performance Computer Architecture, HPCA ’11*, pages 301–311, Washington, DC, USA, 2011. IEEE Computer Society.
- [22] Mendel Rosenblum and John K. Ousterhout. The design and implementation of a log-structured file system. In *Proceedings of the Thirteenth ACM Symposium on Operating Systems Principles, SOSP ’91*, pages 1–15. ACM, 1991.
- [23] Mohit Saxena, Michael M. Swift, and Yiying Zhang. Flashtier: A lightweight, consistent and durable storage cache. In *Proceedings of the 7th ACM European Conference on Computer Systems, EuroSys ’12*, pages 267–280. ACM, 2012.
- [24] Sriram Subramanian, Swaminathan Sundararaman, Nisha Talagala, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Snapshots in a flash with ioSnap. In *Proceedings of the Ninth European Conference on Computer Systems, EuroSys ’14*, pages 23:1–23:14. ACM, 2014.
- [25] Jingpei Yang, Ned Plasson, Greg Gillis, Nisha Talagala, and Swaminathan Sundararaman. Don’t stack your log on my log. In *2nd Workshop on Interactions of NVM/Flash with Operating Systems and Workloads (IN-FLOW 14)*, Broomfield, CO, 2014. USENIX Association.
- [26] Yiying Zhang, Leo Prasath Arulraj, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. De-indirection for flash-based SSDs with nameless writes. In *Proceedings of the 10th USENIX Conference on File and Storage Technologies, FAST’12*, pages 1–1. USENIX, 2012.