

ElasticBF: Fine-grained and Elastic Bloom Filter Towards Efficient Read for LSM-tree-based KV Stores

Yueming Zhang¹, Yongkun Li^{1,2}, Fan Guo¹, Cheng Li^{1,2}, Yinlong Xu^{1,2}

¹University of Science and Technology of China

²Anhui Province Key Laboratory of High Performance Computing, USTC

Abstract

LSM-tree based KV stores suffer from severe read amplification, especially for large KV stores. Even worse, many applications may issue a large amount of lookup operations to search for nonexistent keys, which wastes a lot of extra I/Os. Even though Bloom filters can be used to speedup the read performance, existing designs usually adopt a uniform setting for all Bloom filters and fail to support dynamic adjustment, thus results in a high false positive rate or large memory consumption. To address this issue, we propose ElasticBF, which constructs more small filters for each SSTable and dynamically load into memory as needed based on access frequency, so it realizes a fine-grained and elastic adjustment in running time with the same memory usage. Experiment shows that ElasticBF can achieve $1.94\times-2.24\times$ read throughput compared to LevelDB under different workloads, and preserves the same write performance. More importantly, ElasticBF is orthogonal to existing works optimizing the structure of KV stores, so it can be used as an accelerator to further speedup their read performance.

1 Introduction

Key-value (KV) store has become an important storage engine for many applications [12, 4, 16, 23]. The most common design of KV stores is based on Log-Structured Merge-tree (LSM-tree), which groups KV pairs into fixed-size files, e.g., the SSTables in LevelDB [7]. Files are further organized into multiple levels as shown in Figure 1. KV pairs are sorted in each level except level 0. With this level-based structure, data is first flushed from memory to the lowest level, and then merged into higher levels by using compaction when this level reaches its capacity limit. Compaction inevitably causes the write amplification problem, and many recent researches focus on addressing this issue [3, 22, 13].

On the other hand, LSM-tree based KV stores also suffer from severe read amplification problem[14, 21],

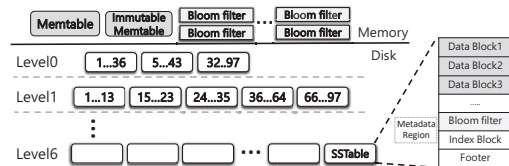


Figure 1: LSM-tree based structure.

especially for large KV stores. This is mainly because a KV store contains multiple levels, and reading a KV item needs to check from the lowest level to the highest level until the data is found or all levels are checked. This process inevitably incurs multiple I/Os and amplifies the read operation. In the worst case, 14 SSTables need to be inspected [13], and even worse, if the target KV item does not exist in the store [19, 20, 9], then all the I/O requests are totally wasted. We point out that only one file in each level need to be examined as data in each level are kept in a sorted order, and this file can also be easily located by checking the key ranges of each file.

To reduce extra I/O requests, Bloom filters are widely used in KV stores [19]. By first asking the filter to check if the requested data exists in the SSTable, extra I/Os can be reduced. However, Bloom filter has false positive, so it may return an “existence” answer even if the data does not really exist, still incurs extra I/Os to inspect the SSTable. Even though the false positive rate (FPR) can be reduced by increasing the length of filters [2, 7], it will increase the memory usage. As a result, some filters may be swapped out to disks as memory is usually a scarce resource in KV stores [11, 8]. If filters are not in memory, extra I/Os must be incurred to load the filter into memory before inspecting a SSTable, and this exacerbates the read amplification. A recent work proposes to adjust the length of filters for different levels [6], while it only uses a same setting for all filters in the same level and fails to dynamically adjust the setting according to data hotness. Therefore, there still remains a challenging problem for LSM-tree-like KV stores: How to reduce the false positive rate of Bloom filters with least memory usage?

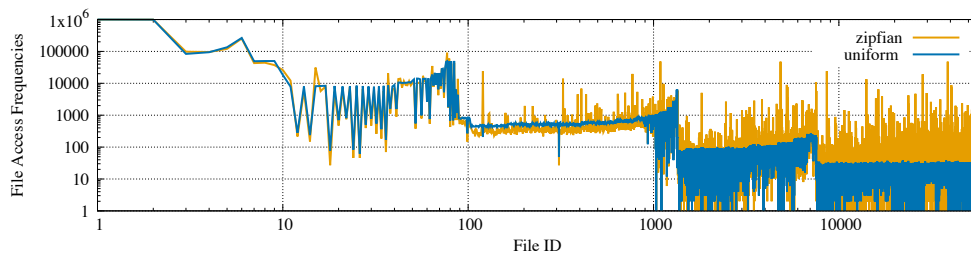


Figure 2: File access frequencies under two different workloads.

In this paper, we propose ElasticBF, realizes a fine-grained and elastic Bloom filter by constructing more small filters on disk and dynamically load some into memory as needed based on access frequency and data hotness. So ElasticBF can easily tune the size of filters for each SSTable in running time with the same memory usage. With the fine-grained and elastic feature, ElasticBF greatly reduces the false positive rate under limited memory space. Besides, ElasticBF is orthogonal to existing works focusing on optimizing KV store structure. Our prototype based on LevelDB demonstrates that ElasticBF can achieve up to $2.24\times$ read throughput compared to conventional Bloom filter design in LevelDB.

2 Background & Motivation

As introduced in §1, Bloom filters can be used to reduce the number of I/Os incurred in each read operation. For example, each SSTable in LevelDB is associated with a filter, which can be viewed as a bit array constructed from all the KV pairs within the SSTable by using several hash functions, and then the filter can be used to easily check the existence of a KV pair. However, Bloom filters always have false positive, and the false positive rate can be expressed as $(1 - e^{-k/b})^k$ [10], where b indicates the number of bits allocated for each key, denoted by *bits-per-key*, and k means the number of hash functions. We point out that b determines the memory usage of a Bloom filter. Since $(1 - e^{-k/b})^k$ is minimized when $k = \ln 2 \cdot b$, false positive rate can be represented as 0.6185^b . Based on this formula, we can easily observe the tradeoff between memory usage and false positive rate for a filter.

It is a common consensus that accesses to files usually possess locality, this is also true for KV stores. We further validate this by running an experiment with LevelDB [7]. Since there is no publicly available KV workload trace [15], and YCSB [5] is the standard benchmark for evaluating KV stores. We use YCSB to load a 100GB database and generate two representative workloads [12] containing one million Get requests with uniform and zipfian distributions, which are provided by YCSB to simulate real-world application scenarios [5]. Note that there are about 50K SSTables in a 100GB KV store, so issuing one million Get requests is enough for us to evaluate the stable behaviour of the KV store. Figure 2 shows the access frequency of all SSTables, which are numbered sequen-

tially from the lowest level to the highest level. As the level increases, the access frequency of SSTables tends to decrease, because the higher the level is, the more SSTables are accessed. Therefore, filters at lower levels need a smaller false positive rate than those at higher levels. A recent work Monkey [6] maximizes the read throughput under uniform distribution by allocating more *bits-per-key* to filters at lower levels. We emphasize that even for SSTables in the same level, unevenness of access frequency is still very common, but Monkey does not consider the skewness within the same level.

With the consideration of access locality, it is not cost-effective when using the same *bits-per-key* to configure all filters for different SSTables. Instead, allocating more *bits-per-key* for filters of hot SSTables, and allocating less bits for cold SSTables, can reduce the overall false positive rate for all requests and keep the same memory usage, because hot SSTables may receive most read requests but only account for a small portion. Therefore, it is of big significance to adjust the *bits-per-key* for every SSTable in a fine-grained and elastic manner so as to minimize the false positive rate with least memory usage. However, the configuration of Bloom filters is fixed, or it can only be adjusted by regenerating filters, and adjusting it dynamically through regenerating will incur a large overhead. Since SSTables are immutable, regenerating Bloom filters needs to re-read the data in SSTables.

3 ElasticBF

3.1 Main Idea

The access frequency varies significantly from different SSTables, and the number of hot SSTables (with high access frequency) is far less than that of cold SSTables. Thus, we could reduce the amount of extra I/Os caused due to false positive of filters by increasing the *bits-per-key* for filters of hot SSTables, meanwhile, we can still limit the memory overhead by decreasing the *bits-per-key* for filters of cold SSTables. However, the *bits-per-key* allocated to a filter is fixed and unable to change as long as the filter has been generated, so it is unable to dynamically adjust the allocation for different filters at running time. To address this issue, instead of directly adjusting the *bits-per-key*, we choose to build multiple filters when constructing each SSTable, and each filter

is usually allocated with smaller *bits-per-key* and called *filter unit*, we then dynamically adjust the number of filters in memory for each SSTable by enabling some filter units and loading them into memory or disabling some in-memory filter units and simply removing them from memory. Note that we do not need to write the filter units back to the disk, because copies of these filter units are also stored in the SSTables. Thus, we achieve the elastic feature to dynamically adjust the *bits-per-key* or the false positive rate, so we call the scheme ElasticBF.

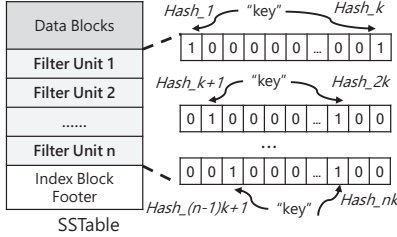


Figure 3: Construction of ElasticBF

In ElasticBF, each SSTable is initially associated with multiple filter units, as shown in Figure 3, which are generated by using different and independent hash functions based on the KV pairs within the SSTable and sequentially placed in the metadata region of the SSTable. We collectively call all the filter units assigned to a SSTable a *filter group*. Since multiple filter units within a filter group are independent, when querying for a key exists in the SSTable, the key must be nonexistent certainly as long as one filter unit gives a negative return. That is, only when all filter units indicate the existence of a key, we then really read out the data blocks to search the key.

One important feature of ElasticBF is that the overall false positive rate of a filter group is exactly the same as that of a single Bloom filter which uses the same number of *bits-per-key* allocated to all filter units within the filter group. We call this feature *separability*. This feature can also be easily proved. Assuming that each filter unit is a b/n *bits-per-key* filter, so its FPR can be expressed as $0.6185^{b/n}$. Since hash functions used by each filter unit are independent, the FPR of n filter units can be derived as $(0.6185^{b/n})^n$, which is exactly the same with that of a single filter with b *bits-per-key*. We also validate this feature via experiments by using the *bloom_test* provided by LevelDB, and the results conform with the analysis.

ElasticBF provides a way to dynamically adjust Bloom filters during running time, and based on the feature of separability, the memory usage can also be limited by enabling more filter units for hot SSTables and disabling some filter units for cold ones. However, to deploy ElasticBF in a KV store, there still remain two challenging issues: (1) *How to design an adjusting rule to determine the most appropriate number of filter units for each SSTable?* (2) *How to realize a dynamic adjustment with small overhead?*

3.2 Adjusting Rule

The goal of adjusting Bloom filters for each SSTable is to reduce the extra I/Os caused due to false positive, so we use a metric which is defined as the expected number of I/Os caused by false positive to guide the adjustment, and we denote this amount of extra I/Os as $E[\text{Extra_IO}]$. Specifically, $E[\text{Extra_IO}]$ can be expressed as

$$E[\text{Extra_IO}] = \sum_{i=1}^n f_i \cdot fp_i, \quad (1)$$

where n means the total number of SSTables in the KV store, f_i denotes the access frequency of SSTable i , fp_i denotes the false positive rate and it is determined by the number of filter units loaded in memory for SSTable i .

ElasticBF adjusts the number of filter units for each SSTable only when the metric $E[\text{Extra_IO}]$ could be reduced under the fixed memory usage. It limits memory usage by fixing the total length for all filter units in the memory, that is, we keep the average *bits-per-key* as a fixed value. Note that SSTables which have higher access frequency will contribute more to $E[\text{Extra_IO}]$ when using the same allocation of Bloom filters, so minimizing $E[\text{Extra_IO}]$ actually results in allocating more filter units for hot SSTables, which meets our goal exactly.

The adjustment of Bloom filter proceeds as follows. Each time when a SSTable is accessed, we first increase its access frequency by one and update $E[\text{Extra_IO}]$, then we check whether $E[\text{Extra_IO}]$ could be decreased if enabling one more filter unit for this SSTable and disabling some filter units in other SSTables so as to guarantee the same memory usage. However, one critical issue is how to quickly find which filters should be disabled but not incurring a large overhead. To address this problem, we extend Multi-Queue (MQ) [24, 17], and we describe its design details in next subsection.

3.3 Dynamic Adjustment with MQ

We maintain multiple Least-Recently-Used (LRU) queues to manage the metadata of each SSTable as shown in Figure 4, and we denote these queues as Q_0, \dots, Q_m , where m is equal to the maximum number of filter units allocated to each SSTable. Recall that all the filter units within a filter group are kept on disks along with the SSTable, but not all filter units are enabled, and only the enabled filter units are used for key-existence check. Each element of a queue corresponds to one SSTable and keeps the metadata of the SSTable, including the enabled filter units residing in memory. Q_i manages the SSTables which already enabled exactly i filter units, e.g., each SSTable in Q_2 enabled two filter units.

To determine which filter units should be disabled and then removed from memory, we use an *expiring* policy which associates a variable named *expiredTime* with each item in MQ, and it denotes the time point at which the

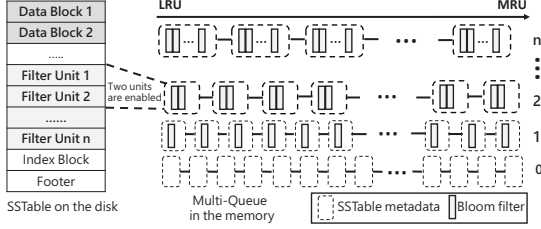


Figure 4: Multi-Queue Structure in ElasticBF

corresponding SSTable should be expired and selected as a candidate to adjust its Bloom filters. Precisely, *expiredTime* is defined as $currentTime + lifeTime$, where *currentTime* denotes the number of Get requests issued to the KV store so far, and *lifeTime* is a fixed constant. In practice, we let *lifeTime* have the same order of magnitude as the total number of SSTables, and we point out that ElasticBF is not sensitive to the value of *lifeTime* under this setting. Based on the above definitions, when a SSTable is inserted into a queue, its *expiredTime* will be initialized as $currentTime + lifeTime$, and at each time when the SSTable is read, its *expiredTime* will be updated accordingly based on the new value of *currentTime*. We define a SSTable as “expired” if the total number of Get requests issued to the KV store becomes larger than its *expiredTime*. The physical meaning is that the SSTable has not been read in the past *lifeTime* requests if it becomes “expired”, so this SSTable can be regarded as *cold* and some of its filter units could be disabled.

For each access, we follow the original MQ algorithm to find the “expired” SSTable [17]. To avoid “expired” SSTable enabling too many filter units, we search “expired” SSTables from Q_m to Q_1 , and for each queue, we search from the LRU side to the MRU side, since an “expired” SSTable must be the least recently used one. When we find an “expired” SSTable, we downgrade it to the next lower-level queue to release one filter unit, and if we encounter a not “expired” SSTable, we jump to lower-level to search “expired” SSTables. This process stops until we find enough memory space to load new filter units for the SSTable which tends to increase its filter units. If we can not find enough “expired” filters to make room for new filters, then we simply skip the Bloom filter adjustment this time.

3.4 Overhead Analysis

To support runtime adjustment, ElasticBF keeps multiple filter units in each SSTable. But it requires extra storage and generating these filters may also add latency to writes. First, assuming that the size of KV pairs is 1KB and ElasticBF uses four *bits-per-key*, so one filter unit only costs around 1KB storage, and it is just 0.05% of a SSTable which is usually 2MB. Second, our experiment shows that the time of building a filter takes only around 1% of the time for constructing a SSTable on the

disk. ElasticBF also leverages multi-threading to generate multiple filter units simultaneously so as to further reduce the computation time of generating filters. Therefore, the storage overhead is small and computation time is short. Besides, ElasticBF maintains multiple queues to quickly identify candidate SSTables for decreasing filter units. Since the queues we used are LRU queues, it is easy to find “expired” SSTables. The number of CPU cycles for searching is small, and it can be ignored when compared with I/O time. Extra memory overhead of MQ comes mainly from keeping *expiredTime*, since LSM-tree based KV stores originally use a linked list to manage the metadata of SSTables. Thus, the memory overhead is also small. For example, for a 100GB KV store, there are around 50K SSTables, assuming that 4B is used to record the *expiredTime* of each SSTable, then the total memory overhead is only around 200KB.

4 Performance Evaluation

We run experiments on a machine consisting of two Intel(R) Xeon(R) E5-2650 v4 CPUs with 48-cores running at 2.2GHz, 64GB memory and 480GB SSD. The operating system is Linux 3.10.0-514. We compare ElasticBF with LevelDB which is widely used KV store, and use the benchmark YCSB-C [14, 18], which is the C++ version of YCSB [5] for evaluation. We set the KV pair size as 1KB, and set the *lifeTime* of ElasticBF as 20K as there are around 50K SSTables in total in our evaluation. Since lookups of non-existent items are common in practical systems [20, 9, 19], we assume half of the Get requests are to lookup non-existent items, and we use direct I/O to minimize the cache influence.

4.1 Read Performance

We first show the read performance of ElasticBF and LevelDB. We focus on a 100GB KV store, and let the system start from an empty cache in which no filter is buffered so as to include the overhead of loading Bloom filters. The average Bloom filter space for each key (*bits-per-key*) is four bits. We issue one million Get requests under different workloads which are generated by adjusting the distribution of accessed keys.

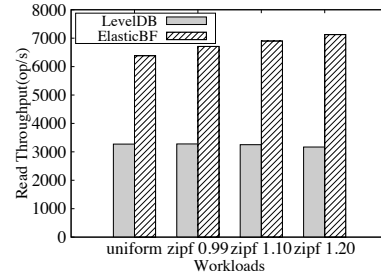


Figure 5: Read Throughput

Figure 5 shows the results. Because LevelDB uses

the same configuration for all Bloom filters, even if the workload changes, it has little effect on read performance. We can see that ElasticBF outperforms LevelDB under all workloads, e.g., the read throughput is increased to $1.94\times$ - $2.24\times$ under different workloads. Besides, when the workload is more skewed, the improvement becomes larger. This also validates the efficiency of taking into account access locality. To further explain the improvement of ElasticBF, we also count the total number of I/Os generated to serve the one million Get operations. We can see that the number of I/Os issued by LevelDB is around $2.42\times$ - $3.05\times$ larger than that of ElasticBF. This is the main reason why ElasticBF achieves higher read throughput than LevelDB.

	uniform	zipf 0.99	zipf 1.10	zipf 1.20
LevelDB	1525595	1585605	1634752	1667947
ElasticBF	628225	578553	550658	545345

Table 1: Number of I/Os for data access

4.2 Read Performance vs Memory Usage

In this experiment, we show the read throughput by varying the *bits-per-key*. Note that allocating more bits for each key can reduce the false positive rate and thus improve the read performance, but it also implies a larger memory usage, so this experiment actually shows the tradeoff between read performance and memory usage. We also enable a 8MB *block_cache* [1] to cache the recently accessed blocks in this experiment, and warm it up before evaluation. Figure 6 shows the experiment results. The x-axis represents the number of bits allocated for each Bloom filter for LevelDB, while for ElasticBF, we adjust the setting of filter units to guarantee the same memory usage for fair comparison.

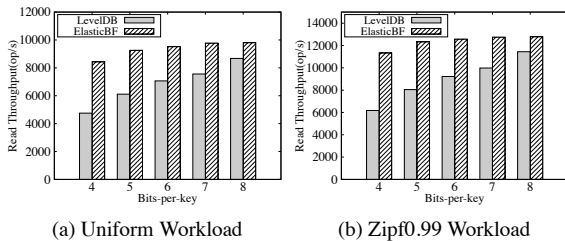


Figure 6: Read throughput w/ different memory usage

Results show that ElasticBF always has higher read throughput than LevelDB with the same memory usage. In other words, ElasticBF can achieve the same read throughput with much smaller memory usage. For example, the read throughput of ElasticBF under the setting of 4 *bits-per-key* is similar to that of LevelDB under the setting of 8 *bits-per-key*. This implies that ElasticBF can achieve a similar read performance with LevelDB with only a half memory usage. However, we can find that the benefit of ElasticBF decreases if more memory can be allowed to use. In particular, when the number of bits

allocated for each key is 8, the increase of read throughput reduces to 10%. However, we like to point out that allocating a large number of bits for each key may not be practical in real systems as KV stores are usually large. For example, if we allocate 8 *bits-per-key* for a 100TB KV store, then only the Bloom filter costs about 100GB of memory, which is too large in practical use as other metadata also needs to consume memory.

4.3 Write Performance

Now we evaluate the impact of ElasticBF on write performance. We compare the time to load a 100 GB KV store by using ElasticBF and LevelDB. We also consider different Bloom filter settings, while the memory usage is guaranteed to be the same. We run the experiments three times for each setting so as to obtain a stable result. As shown in Figure 7, ElasticBF has almost the same write throughput with LevelDB (4.818MB/s and 4.823MB/s separately when *bits-per-key* is four), and the performance difference is only around 0.1%. The main reason is that Bloom filters are organized into blocks in SSTables, which are further written to devices sequentially, and ElasticBF also uses multi-threading to speedup the generation of Bloom filters. In short, ElasticBF has negligible impact on writes.

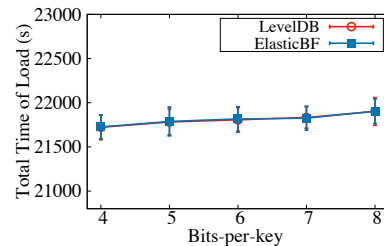


Figure 7: Time to load a KV store

5 Conclusion

In this paper, we developed ElasticBF, a fine-grained and elastic Bloom filter to minimize extra I/Os in KV stores. ElasticBF supports dynamic adjustment during running time by effectively determining the most appropriate number of filter units for each SSTable according to access patterns. Experimental results show that ElasticBF can greatly reduce the number of I/Os during key lookups and improve the read throughput without sacrificing write performance. More importantly, ElasticBF is orthogonal to the works optimizing the structure of KV stores, so it can be widely used to further speedup the read performance for these KV stores.

6 Acknowledgements

This work was supported in part by National Nature Science Foundation of China (61772484 and 61772486). Yongkun Li is the corresponding author.

References

- [1] Block Cache. <https://github.com/facebook/rocksdb/wiki/Block-Cache>, 2017.
- [2] Tuning Bloom filters. <https://docs.datastax.com/en/cassandra/3.0/cassandra/operations/opsTuningBloomFilters.html>, 2017.
- [3] BALMAU, O. M., DIDONA, D., GUERRAOU, R., ZWAENEPOEL, W., YUAN, H., ARORA, A., GUPTA, K., AND KONKA, P. Triad: creating synergies between memory, disk and log in log structured key-value stores. In *USENIX ATC'17* (2017), no. EPFL-CONF-228863.
- [4] BALMAU, O. M., GUERRAOU, R., TRIGONAKIS, V., AND ZABLOTCHI, M. I. Flodb: Unlocking memory in persistent key-value stores. In *EuroSys 2017* (2017), no. EPFL-CONF-227333.
- [5] COOPER, B. F., SILBERSTEIN, A., TAM, E., RAMAKRISHNAN, R., AND SEARS, R. Benchmarking cloud serving systems with ycsb. In *Proceedings of the 1st ACM symposium on Cloud computing* (2010), ACM, pp. 143–154.
- [6] DAYAN, N., ATHANASSOULIS, M., AND IDREOS, S. Monkey: Optimal navigable key-value store. In *Proceedings of the 2017 ACM International Conference on Management of Data* (2017), ACM, pp. 79–94.
- [7] GHEMAWAT, S., AND DEAN, J. Leveldb. <https://github.com/google/leveldb>, 2011.
- [8] HARTER, T., BORTHAKUR, D., DONG, S., AIYER, A. S., TANG, L., ARPACI-DUSSEAU, A. C., AND ARPACI-DUSSEAU, R. H. Analysis of hdfs under hbase: a facebook messages case study. In *FAST* (2014), vol. 14, p. 12th.
- [9] KANG, Y., PITCHUMANI, R., MARLETTE, T., AND MILLER, E. L. Muninn: A versioning flash key-value store using an object-based storage model. In *Proceedings of International Conference on Systems and Storage* (2014), ACM, pp. 1–11.
- [10] KIRSCH, A., AND MITZENMACHER, M. Less hashing, same performance: building a better bloom filter. In *ESA* (2006), vol. 6, Springer, pp. 456–467.
- [11] LAI, C., JIANG, S., YANG, L., LIN, S., SUN, G., HOU, Z., CUI, C., AND CONG, J. Atlas: Baidu’s key-value storage system for cloud data. In *Mass Storage Systems and Technologies (MSST), 2015 31st Symposium on* (2015), IEEE, pp. 1–14.
- [12] LIM, H., ANDERSEN, D. G., AND KAMINSKY, M. Towards accurate and fast evaluation of multi-stage log-structured designs. In *FAST* (2016), pp. 149–166.
- [13] LU, L., PILLAI, T. S., ARPACI-DUSSEAU, A. C., AND ARPACI-DUSSEAU, R. H. Wisckey: Separating keys from values in ssd-conscious storage. In *FAST* (2016), pp. 133–148.
- [14] PAPAGIANNIS, A., SALOUSTROS, G., GONZÁLEZ-FÉREZ, P., AND BILAS, A. Tucana: Design and implementation of a fast and efficient scale-up key-value store. In *USENIX Annual Technical Conference* (2016), pp. 537–550.
- [15] PITCHUMANI, R., FRANK, S., AND MILLER, E. L. Realistic request arrival generation in storage benchmarks. In *Mass Storage Systems and Technologies (MSST), 2015 31st Symposium on* (2015), IEEE, pp. 1–10.
- [16] RAJU, P., KADEKODI, R., CHIDAMBARAM, V., AND ABRAHAM, I. Pebblesdb: Building key-value stores using fragmented log-structured merge trees. In *Proceedings of the 26th Symposium on Operating Systems Principles* (2017), ACM, pp. 497–514.
- [17] RAMOS, L. E., GORBATOV, E., AND BIANCHINI, R. Page placement in hybrid memory systems. In *Proceedings of the international conference on Supercomputing* (2011), ACM, pp. 85–95.
- [18] REN, J. Ycsb-c. <https://github.com/basicthinker/YCSB-C>, 2015.
- [19] SEARS, R., AND RAMAKRISHNAN, R. blsm: a general purpose log structured merge tree. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data* (2012), ACM, pp. 217–228.
- [20] URDANETA, G., PIERRE, G., AND VAN STEEN, M. Wikipedia workload analysis for decentralized hosting. *Computer Networks* 53, 11 (2009), 1830–1845.
- [21] VO, H. T., WANG, S., AGRAWAL, D., CHEN, G., AND OOI, B. C. Logbase: a scalable log-structured database system in the cloud. *Proceedings of the VLDB Endowment* 5, 10 (2012), 1004–1015.
- [22] WU, X., XU, Y., SHAO, Z., AND JIANG, S. Lsm-trie: an lsm-tree-based ultra-large key-value store for small data. In *Proceedings of the 2015 USENIX Conference on Usenix Annual Technical Conference* (2015), USENIX Association, pp. 71–82.
- [23] ZHANG, H., DONG, M., AND CHEN, H. Efficient and available in-memory kv-store with hybrid erasure coding and replication. In *FAST* (2016), pp. 167–180.
- [24] ZHOU, Y., PHILBIN, J., AND LI, K. The multi-queue replacement algorithm for second level buffer caches. In *USENIX Annual Technical Conference, General Track* (2001), pp. 91–104.