

Reducing Bloom Filter CPU Overhead in LSM-Trees on Modern Storage Devices

Zichen Zhu, Ju Hyoung Mun, Aneesh Raman, Manos Athanassoulis

zczhu@bu.edu, jmun@bu.edu, aneeshr@bu.edu, mathan@bu.edu

Boston University

ABSTRACT

Bloom filters (BFs) accelerate point lookups in Log-Structured Merge (LSM) trees by reducing unnecessary storage accesses to levels that do not contain the desired key. BFs are particularly beneficial when there is a significant performance difference between probing a BF (hashing and accessing memory) and accessing data (on secondary storage). However, this gap is decreasing as modern storage devices (SSDs and NVMs) have increasingly lower latency, to the point that the cost of *accessing data* can be comparable to that of filter probing and *hashing*, especially for large key sizes that exhibit high hashing cost. In an LSM-tree, BFs are employed when querying each of the levels of the tree, thus, exacerbating the CPU cost as the data size grows (and, thus, the tree height). To address the increasing CPU cost of BFs in LSM-trees, we propose to *re-use hash calculations* aggressively within and across BFs, as well as between different levels, and we show both analytically and experimentally that we can maintain close-to-ideal false positive rate while significantly reducing the runtime. The reduced CPU cost of queries using the proposed *hash sharing* leads to 20% higher lookup performance in an LSM-tree with 22GB of data (5 levels) stored in a state-of-the-art PCIe SSD. The benefit further increases for faster underlying storage. Specifically, we show that when NVM devices will be available the improvement can increase up to 65%.

1 INTRODUCTION

LSM-trees are Everywhere. Log-Structured Merge-trees (LSM-trees) [27] are the core data structure of several state-of-the-art key-value engines like RocksDB [13] at Facebook, LevelDB [14] and BigTable [6] at Google, HBase [16] and Cassandra [3] at Apache, WiredTiger [37] at MongoDB, X-Engine [17] at Alibaba and DynamoDB [11] at Amazon. LSM-trees are widely adopted because they offer high ingestion rate and support fast reads. In addition to the systems developed in industry that are mentioned above, in the past few years, various LSM-tree optimizations on compaction, membership filtering, and memory management have been proposed [1, 2, 5, 7, 9, 10, 18, 22, 23, 25, 26, 32, 34, 38, 39, 41, 42].

The Structure of LSM-trees. LSM-trees maintain sorted runs across multiple levels with exponentially increasing capacity, which

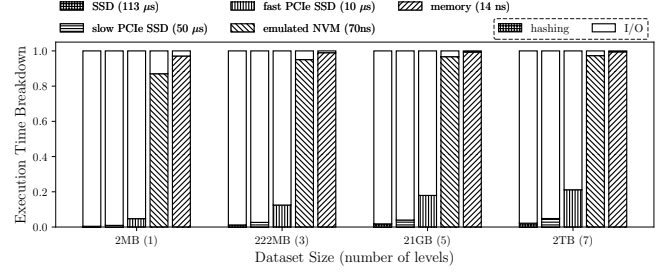


Figure 1: The percentage of time spent in hashing to probe BF in an LSM-tree increases as either the data size increases and/or as the underlying storage becomes faster. Thus, hashing becomes the main performance bottleneck for LSM-trees that hold large datasets on fast storage devices.

have potentially overlapping key ranges. The performance of LSM-trees is determined by many tuning knobs, including compaction policy, size ratio between levels, and metadata used to accelerate read queries. In particular, LSM-trees employ fence pointers and Bloom filters (BFs) to reduce unnecessary storage accesses [24].

Bloom Filters in LSM-trees. Since key-value pairs are spread across multiple levels, a point query might need to probe every level of a tree, thereby, requiring multiple I/Os for a single lookup. To avoid unnecessary accesses, LSM-trees typically employ BFs [4] to identify whether the target key belongs to a given level. A BF is associated with each level on the secondary storage (or a file that belongs to a given level, in case of partitioned LSM-trees [12]), and is often pre-fetched in main memory, to be readily available during a point query before accessing slow storage. The cost of querying a BF is two-fold: (a) the hash calculation and (b) the probing of the filter's bits. Often, the BFs fit in main memory, and thus, the probing cost is negligible. On the other hand, accessing data on secondary storage, e.g., hard disk drives (HDD) or solid-state drives (SSD), is several orders of magnitude more expensive than probing the filter in memory. This performance gap always renders it worthwhile to consult BFs before accessing data. BFs reduce the number of data accesses and the overall query latency at the price of additional memory footprint and hashing.

What About Faster Storage? Contrary to common perception, however, BFs are not always beneficial [33]. The rationale behind the ubiquitous use of BFs in LSM-trees is that there is a *considerable cost difference between accessing a BF (in memory) and accessing data (on disk)*. As the gap of access latency between BFs and data narrows, the advantage of using a BF weakens. If the data is already cached in main memory, BFs are detrimental. Further, as new storage devices like SSDs and non-volatile memories (NVMs) [31] emerge, the latency gap between memory and storage narrows. Typically, a BF probe requires an expensive hash calculation and one or more

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.

Under Submission, 2021

© Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM.

<https://doi.org/10.1145/1122445.1122456>

memory accesses for a total cost on the order of $1\mu\text{s}$ for 1KB keys. Meanwhile, a potentially unnecessary disk access is on the order of 2ms for HDD and $100\mu\text{s}$ for SSDs. As a result, for each query within a level, an additional $1\mu\text{s}$ can help to avoid an I/O and, thus, significantly reduces the query latency, especially for zero-result point lookups. However, faster PCIe SSD devices offer access latency as low as $10\mu\text{s}$ access latency or even lower (e.g., we measured $7\mu\text{s}$ when using Intel’s SPDK [35] to bypass the file system on our PCIe SSD device), making the access latency comparable to the hash calculation cost. In addition, future NVM devices are expected to have much faster latency, being only $5\times$ slower than DRAM, in the order of 70ns [28]. Overall, as storage devices become faster, they challenge the across-the-board benefit of using BFs.

Bloom Filters Incur High CPU Overhead. Ideally, a BF requires multiple independent hash functions in order to achieve the theoretical false positive rate (FPR). In practice, however, the different BF indexes are often computed using a single hash function to calculate a hash digest, followed by much cheaper bitwise operations (rotations and modulo) to generate the remaining index probes¹. Taking into account that current SSD devices have several orders of magnitude lower access latency than disk and that future SSDs and NVMs are bound to be faster, *hashing latency is on its way to become comparable with data access latency*. For example, accessing a 4KB data page on our off-the-shelf SSD needs $113\mu\text{s}$, while the cost of hashing a 1KB-key using MurmurHash64, which is used in production systems [13], is 235ns , making storage about $480\times$ more expensive than hashing. However, accessing a data page of our PCIe SSD device takes $10\mu\text{s}$ ($7\mu\text{s}$ when bypassing the file system), reducing this gap to $42\times$ ($30\times$ without file system). In addition, future NVM devices are expected to offer access latency in the order of 70ns , making accessing storage $3\times$ faster than hashing. Note that when data is cached in main memory, a single hash function calculation is about $17\times$ more expensive than accessing a memory page making the use of the BF detrimental. The LSM hashing overhead is further exacerbated as multiple BFs are queried per lookup (one per level), and repeated hash calculations turn querying over fast storage (or cached data) into a CPU-intensive operation. Figure 1 shows the execution time breakdown of point queries (half of which are empty queries) in a state-of-the-art LSM-tree. We focus on the time spent hashing and time spent waiting for I/O completion. We compare different devices for storing the data (an SSD, a PCIe SSD, a fast PCIe SSD, an emulated NVM, and main memory) and different dataset sizes, ranging from 22MB (2 level) to 2TB (7 levels), with size ratio 10. For a fast PCIe SSD that has access latency of $10\mu\text{s}$, about 80% of the lookup time is consumed by hashing for 2GB of data. As the access latency of NVM is expected to be reduced, hashing overhead will inevitably increase. The fraction of time spent on hashing accumulates over all levels, and eventually becomes the main bottleneck. Further, recent designs deliver higher and more flexible accuracy for point queries [20] and short range queries [26], using *more smaller* BFs. Both approaches introduce four or more additional BFs per SST file, thus increasing the hashing overhead.

Hash Sharing. To reduce the CPU overhead, we propose to aggressively re-use hash computations within a BF and across different

BFs residing in different levels. We reproduce state-of-the-art results for BFs that use only one hash calculation and perform cheaper computations (termed *pseudo-hashing*) for the remaining positions of the BF’s bitvector. We take this a step further by (i) showing how to share hash computations across multiple LSM levels, and (ii) across different independent modules of a single logical BF, in the case ElasticBF [20].

The aggregate cost of hashing in state-of-the-art LSM-Trees depends on the height of the tree, which depends on the data size, since for each point query a BF per level is typically accessed. However, *hash sharing across levels* decouples the aggregated hashing cost from data size, since, regardless of the number of LSM-tree levels, the amount of hashing remains constant. Similarly, hash sharing for ElasticBF allows us to decouple its CPU cost from its design, and internally use pseudo-hashing.

Contributions. Our work offers the following contributions.

- We identify that BFs dominate LSM query latency for *fast storage* and *high hashing cost*.
- We decouple the amount of hashing from the data size (height of a LSM-tree) by hash sharing across different levels.
- We decouple the amount of hashing from the design complexity of Elastic BF, an approach that can be used by other BF variants.
- We show through analytical and experimental results that hash sharing improves LSM query performance by 20% on PCIe SSD and 65% on emulated NVM devices.

2 BACKGROUND

LSM-tree Basics. Many modern key-value stores adopt LSM-trees as their storage layer in order to handle write-intensive workloads, because LSM-trees are designed for fast ingestion [3, 6, 11, 13, 14, 16, 17, 27, 37]. To support fast writes, LSM-trees buffer all inserts (including the ones updating or deleting existing entries) in a memory buffer, typically referred to as Level 0. When the buffer reaches a predetermined capacity, it is flushed to secondary storage in the form of a sorted *run*, consisting of multiple files stored as immutable Sorted-String Tables (SST files). All runs in the secondary storage are organized in a tree-like structure where each level has exponentially larger capacity according to a user-defined size ratio T . The number of LSM-tree levels L depends on the total data size, the size of the memory buffer, and the size ratio [26]. Shallower levels store more recent updates and have smaller capacity. Similar to buffer flushing, whenever a level fills up, a sort-merge operation is triggered between this newly-saturated level and the next one, and obsolete entries are removed during this process.

Point Queries in LSM-trees. As LSM-updates are out-of-place, multiple entries with the same key may exist. However, searching can terminate safely after finding the first matching entry, because matching keys in the older levels are guaranteed to be obsolete. Therefore, a point query first consults the in-memory buffer and then, traverses the tree from the shallowest to the deepest level until it finds the first match. In case of tiering, which has multiple overlapping runs per level, searching within a level goes from the youngest to the oldest run and terminates if there is a match.

Auxiliary In-Memory Data Structures. To boost query performance, LSM-trees maintain two in-memory auxiliary data structures for each SST file: fence pointers and Bloom filters.

¹For example, see the implementation of LegacyNoLocalityBloomImpl at https://github.com/rockset/rocksdb-cloud/blob/master/util/bloom_impl.h.

Fence Pointers: Since entries within a disk-resident run are sorted by key, the min-max range of each page does not overlap with any other page. Fence pointers are the min-max ranges for each disk page, along with their aggregation at the level of each SST file and each level. They ensure that at most one I/O occurs when searching for a target key within a single run.

Bloom Filters: Each SST file is also equipped with a BF to avoid unnecessary I/Os. A BF is a membership test data structure that uses an m -bit vector and originally k independent hash functions to store and query the membership of n elements [4, 36]. All negative responses to membership queries are always correct, however, positive responses might either be *true positives*, or *false positives* with a small probability which is a function of k , m , and n . The expected false positive rate (f_p) and the optimal number of hash functions to use are shown in Eq. (1).

$$f_p \approx \left(1 - e^{-kn/m}\right)^k \quad \text{where} \quad k_{opt} = \left\lceil \frac{m}{n} \ln 2 \right\rceil \quad (1)$$

Overall, the impact of false positives in LSM-trees can be calculated by considering the disk accesses due to false positives across all levels [8]. All LSM-based key-value stores employ BFs [26] or other variations like ElasticBF [20], SuRF [39, 40] and Rosetta [26].

Storage Access vs. Hashing. Next, we put into context, the comparison between storage access and hashing latency. Table 1 shows the access latency for a 4KB page in various devices (HDD, SSD, PCIe SSD, NVMe, and memory) as well as the hashing latency of a 1KB key using six representative hash functions: 64-bit MurmurHash64 (MM64), XXHash (XX), MD5, SHA-256, CRC and CITY64 (CITY). We use the RocksDB implementation of MM64, XX, and CRC, and Google’s implementation of CITY [15]. As MD5 and SHA-256 are more than one order of magnitude more expensive than other hash functions they are rarely used for practical implementations for BFs. Overall, we observe that even the most efficient hash functions are comparable with accessing data on PCIe SSDs, where hashing accounts for 10% of the access latency, and it is expected to worsen when NVMe devices with DRAM-like latency become available.

3 THE CPU COST OF BLOOM FILTERS

While BFs are employed virtually in all LSM-based data systems to address the storage bottleneck, their use comes at a cost. In this section, we analyze the cost and the benefit of BFs.

With respect to the BFs of an LSM-tree, point queries can be classified as empty queries or non-empty queries targeting existing keys [8]. We first concentrate on a single LSM-tree level, hence we refer to queries targeting the SST files at a specific level i . We represent lookup cost using the fraction of non-empty queries over all point queries as α_i , the BF access cost (CPU cost of hashing and memory cost of probing the BF indexes) as T_{BF} and the data access cost of a page as T_D . The average cost of a point read workload $\mathcal{T}(i)$ for level i with α_i non-empty queries is calculated using the cost of accessing both the BF and the data for the non-empty queries ($\alpha \cdot (T_{BF} + T_D)$), and accessing the BF and data access due to false positive for the empty queries ($(1 - \alpha) \cdot f_p \cdot T_D$ where f_p is false positive ratio), as shown in Eq. (2).

$$\begin{aligned} \mathcal{T}(i) &= \alpha_i \cdot (T_{BF} + T_D) + (1 - \alpha_i) \cdot (T_{BF} + f_p \cdot T_D) \\ &= T_{BF} + \alpha_i \cdot T_D + (1 - \alpha_i) \cdot f_p \cdot T_D \end{aligned} \quad (2)$$

Operation	Latency	Normalized
4KB I/O on HDD	4.6 ms	328571×
4KB I/O on SSD	113 μ s	8071×
4KB I/O on PCIe SSD	10 μ s	714×
4KB I/O on PCIe SSD (using SPDK)	7 μ s	500×
4KB I/O on NVMe	70 ns	5×
4KB access on Memory	14 ns	1×
CITY of 1KB-key	176 ns	13×
Murmur Hash 64 (MM64) of 1KB-key	235 ns	17×
CRC of 1KB-key	323 ns	23×
XXHash (XX) of 1KB-key	874 ns	62×
MD5 of 1KB-key	2.85 μ s	203×
SHA-256 of 1KB-key	5.17 μ s	378×

Table 1: The decreasing access latency of new storage devices makes the hashing cost of a 1KB-key comparable with accessing a page in NVMe (within one order of magnitude).

The BF cost, T_{BF} , consists of two components: (a) hash calculation T_H , and (b) BF probes T_p . The cost for probing a BF depends on where it is stored. Since BFs are usually cached in main memory, the cost for BF probing is often negligible compared to the I/Os. In fact, the hot BFs of an LSM-tree will reside in cache memory, making T_p negligible even compared to T_H . On the contrary, the hash calculation cost purely depends on the CPU power and the key size. Putting everything together, the cost for a read workload $\mathcal{T}(i)$ on level i with α_i fraction of non-empty queries is as shown in Eq. (3).

$$\mathcal{T}(i) = T_H + T_p + \alpha_i \cdot T_D + (1 - \alpha_i) \cdot f_p \cdot T_D \quad (3)$$

Using Eq. (3), we can now understand what is the main bottleneck for point lookups. When α_i is non-negligible, the time spent to retrieve data dominates the overall LSM-tree lookup cost, because T_D corresponds to expensive accesses on slow storage. Even when α_i is very small f_p contributes to a number of slow storage accesses, making them the bottleneck when T_D is very high. However, novel storage devices have dramatically reduced access latency. In this setting, even with high α_i , the bottleneck shifts to hashing.

To better understand the cost breakdown of querying the whole LSM-tree, we synthesize the overall cost using the cost per level. First, we present the overall existing ratio of the total lookups as α , while α_i is the existing ratio of each level i . In an LSM-tree, a point lookup proceeds to the next level only when it fails to find the matching key in that level. To compute α_i , the total number of queries to level i and the rate of existing results of that level are required. In order to calculate α_i , we introduce the probability that a lookup finished in level i (β_i), i.e., the probability that a level i has matching elements. Therefore, the sum of β_i is identical to the overall existing ratio of the total lookups as α . If we assume that keys in an LSM-tree with size ratio T is perfectly uniform, then, β_i depends on α and the size of level i , i.e., $\beta_i = T^{i-1} \cdot \beta_1$, where $\beta_1 = \frac{\alpha}{\sum_{j=1}^L T^{j-1}}$. The probability that a lookup finished in level i is β_i , while the number of lookups that can reach to level i is $1 - \sum_{j=1}^{i-1} \beta_j$. Thus, α_i becomes $\frac{\beta_i}{1 - \sum_{j=1}^{i-1} \beta_j}$, and α_1 is equal to β_1 since a point query to the first lookup is always essential. To sum up, the cost of lookup in the whole LSM-tree is shown in Eq. (4). The detailed

simplifying process is shown in Appendix.

$$\begin{aligned} \text{cost} &= \mathcal{T}(1) + \sum_{i=2}^L (1 - \sum_{j=1}^{i-1} \beta_j) \cdot \mathcal{T}(i) \\ &= (L - \frac{\alpha}{T-1}) \cdot T_H + \alpha \cdot T_D + (L - \frac{\alpha}{T-1}) \cdot f_p \cdot T_D \end{aligned} \quad (4)$$

As shown in Eq. (4), the storage access due to true positives stay constant, while the BF related costs pile up. Hence, the amount of hashing required accumulates with the number of levels, pushing the fraction of time spent on hashing higher with growing data size (and, hence, the tree height) as shown earlier in Figure 1.

4 SHARING BLOOM FILTER HASHING

This section introduces the benefits from bit-rotation optimization to reduce the hashing cost in a single Bloom filter. Then, hash sharing is proposed as a novel optimization scheme that can reduce the hashing cost within multiple Bloom filters.

4.1 Hash Sharing in a Single BF

Textbook BFs [4] rely on k independent hash functions to generate k indexes, which results in high CPU overhead. Practical BF implementations share a single hash calculation for their k indexes [13]. For example, RocksDB uses a single XXHash hash function and multiple indexes by rotating the hash digest. Specifically, given hash function $h(x)$, we define $\delta = h(x) \ll 17 | h(x) \gg 15$, and the i^{th} ($0 \leq i \leq k-1$) hash function $g_i(x)$ is calculated using $g_i(x) = h(x) + i \cdot \delta$. Such an optimization reduces T_H by a factor of k , since it computes only a single hash digest and the bit rotation cost is negligible.

To showcase the impact of this optimization on performance and the false positive ratio (FPR), we conduct a micro-benchmark on a single BF. We vary the underlying hash function, and the key size (8B to 512B). We populate a BF using 10K keys, and we use with 10 bits per key, thus 7 hash indexes, which is optimal. We execute 100K empty point queries to measure the performance and the query performance.

The first experiment measures the impact of hash sharing via bit-rotation on FPR. The key size is fixed to 512B. Table 2 compares the FPR for a BF that uses 7 different hash functions (all k), with a BF that uses a single hash function and bit-rotation. We observe that the bit-rotation optimization does not affect the FPR. In fact, it achieves the FPR close to the theoretical expectation.

Next, we compare the lookup latency with and without bit rotation varying the key sizes. As expected, the hashing cost increases as the key size grows. The cost is further magnified when multiple hash functions are used as shown in Figure 2. Using a single hash

Hash Function	FPR (%)	Hash Function	FPR (%)
MurmurHash (MM)	0.850%	SHA-256	0.868%
MurmurHash64 (MM64)	0.853%	CRC	0.819%
XXHash (XX)	0.794%	CITY	0.850%
MD5	0.921%	All k	0.899%

Table 2: Using k different hash functions to produce k indexes does not necessarily yield to better false positive rate than bit-rotation trick with only 1 hash digest

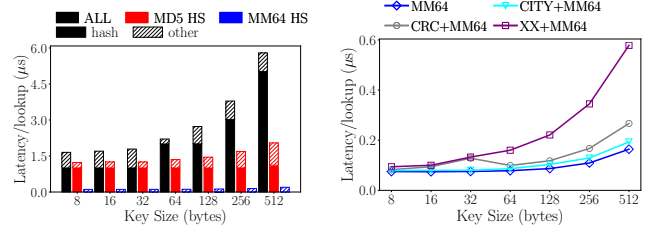


Figure 2: The benefits of hash sharing increase as the key size grows. Cheap hashing significantly reduces the average query latency.

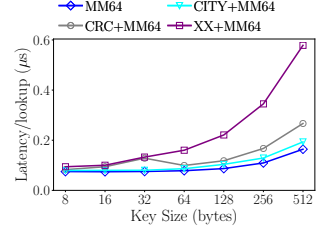


Figure 3: Double hashing and the bit rotation generally has much better efficiency than that all hash functions are computed

function is key to reduce the cost of calculating the k indexes. In addition, it is important to select a hash function that is computed fast, hence, for the remainder of the paper, we focus on MM64, XXHash, CRC, and CITY. Note that we use MM64 as our primary hash function because of its efficient execution and low FPR, as well as its wide usage in production, notably by RocksDB.

4.2 Hash Sharing in ElasticBF

In addition to the classical BF, several BF variants have been proposed for LSM-trees [20, 26] that increase the hashing overhead to address more complex workloads. We focus on ElasticBF which consists of multiple small filter units per BF to address access skewness. By design, each filter unit employs unique hash function to ensure that they are independent, thus increasing substantially the hashing overhead. The hashing cost of ElasticBF increases with the number of filter units and the number of levels in the LSM-tree.

We first address the increased hashing cost due to the number of units, and we then present a solution for the increase of the cost due to the number of levels, which affects all BF-variants in LSM-trees. The bit-rotation optimization is directly applicable to ElasticBF. In addition, we use the double hashing scheme [19] that ensures that each unit will get an provably independent hash function. The double hashing scheme bounds the expected FPR compared to the standard BF by $O(1/n)$ where n is the number of inserted elements. Formally, the double hashing scheme can be defined as follows: Given two independent hash functions $h_1(x)$ and $h_2(x)$, the i^{th} ($0 \leq i \leq k-1$) hash function $g_i(x)$ is defined as $g_i(x) = h_1(x) + i \cdot h_2(x)$.

Design	Hash Function	FPR (%)
One hash function and bit rotations	MM64	0.829%
	XX	0.897%
	CRC	0.841%
	CITY	0.834%
Double hashing	XX + MM64	0.761%
	CRC + MM64	0.808%
	CITY + MM64	0.842%

Table 3: ElasticBF normally achieves almost the same false positive rate as normal Bloom filter

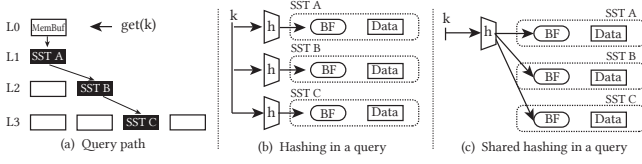


Figure 4: Hash sharing across BFs of different levels.

To investigate how much hash sharing can affect the FPR and the CPU overhead, we emulate ElasticBF and conduct a micro-benchmark that compares the bit-rotation and the double-hashing scheme. The benchmark is similar to the previous one. We populate the ElasticBF using 10K keys and then feed it with 100K empty point queries. Note that bits-per-key is 10, the number of filter unit is 7 and each filter unit uses a single index. For the double hashing scheme, we use MM64 as the primary and three hash functions (XX, CRC, and CITY) as secondary, respectively. As shown in Table 3, the FPR of double-scheme hashing is closer to the theoretical expectation than the single hashing with bit-rotation, while the bit-rotation is faster than the double hashing.

4.3 Hash sharing in an LSM Tree

Now we apply our hash sharing design to BFs across multiple levels of an LSM-tree. When issuing point lookups, we first search the memory buffer, and then, continue searching to the last level, until the target key is found. For each level, the BF is probed first, and if the result is positive (true or false), the data pages are accessed as well. In particular, the fraction of time spent on hashing cost increases for faster devices, which is further exacerbated as data size grows. The latter affects the height of the tree, and as a result, the number of BFs probed during each query.

As the access latency of modern storage devices is getting comparable with hashing – that dominates the Bloom filter cost – hashing becomes the main bottleneck, especially for empty queries that access data only when there is a false positive. In order to mitigate this overhead, we propose to *share the hashing* across BF residing multiple LSM-levels, as shown in Figure 4. When searching for a key in the tree, the hash digest is computed only once, and it is re-used across the entire tree hierarchy, to probe BFs of the SST files that contain the relevant range. As a result, the hashing cost stays constant regardless of the number of levels, shaving off a factor L from the hashing cost in Eq. (4). The new cost is shown in Eq. (5).

$$cost^{share} \approx T_H + \alpha \cdot T_D + (L - \frac{\alpha}{T-1}) \cdot f_p \cdot T_D \quad (5)$$

Performance Benefits. By sharing the hash calculation, the hashing cost is decoupled from the number of levels in LSM-trees, and as a result, from the data size. In our experiments, we have seen that there is no difference in the measured false positive across the different levels of the LSM-tree between the state-of-the-art design and hash sharing, while the hashing cost of an empty query drops by a factor of L . Similarly, for ElasticBF in an LSM-tree, the hash sharing reduces hashing cost further by a L and the number of filter units without harming the FPR. In the following section, we experimentally show the benefits of hash sharing.

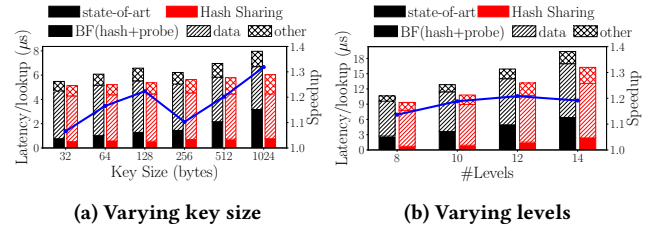


Figure 5: Hash sharing reduces hashing overhead, which is more pronounced for larger keys and higher trees.

5 EXPERIMENTAL EVALUATION

We now present the experimental evaluation of *hash sharing* in LSM-trees. We present its benefits as we vary the key size, the height of the LSM-Tree, the bits per key, and the workload characteristics.

Hardware Environment. We run our experiments in our in-house server, which is equipped with two sockets each with an Intel Xeon Gold 6230 2.1GHz processor with 20 hardware threads and 40 threads with virtualization enabled. The total capacity for main memory is 384GB, and for L3 cache is 27.5MB. The server is equipped with two 7200RPM hard drives, one off-the-shelf SSD, and two state-of-the-art PCIe SSD devices that can offer 600K-1M IOPS and access latency less than 15μs for 4KB page accesses. Unless otherwise specified, we use the PCIe SSD with direct I/O in our experimentation.

Experimental Platform, Workloads and Metrics. We build an in-house LSM-tree prototype system² based on the architecture of RocksDB [13], which uses the fast local Bloom filter (format_version = 5, only supports 64-bit hash digest) from RocksDB. We conduct our experiments with different workloads to stress-test all approaches. Since our algorithm only affects read performance, all experiments execute read-only workloads with 1M empty point lookups after bulk-loading a 22GB dataset unless otherwise noted. **Describe what is the key and the value pair, the ranges and the distribution in 1-2 sentences. The default entry size is 1KB with 512B key and we vary the key size (from 8B to 1KB) or dataset size (from 340MB to 22GB) according to the experiments.** In addition, we set the file size as 2MB, the size ratio of the LSM-tree as 10 and bits-per-key for BFs as 10. Our experiments focus on the lookup latency and its time breakdown, which can best illustrate the CPU bottleneck throughout the lookup process. For each experiment, we report the average across five executions.

Hash Sharing Scales Better with Key Size. The hashing cost inevitably increases as the key size grows. In order to highlight the impact of key size on hashing overhead, we conduct an experiment varying the key size from 8B to 1024B. For this experiment, we populate an LSM-tree using 22GB of key-value pairs with each key-value pair of size 2KB. Thus, the resulting LSM-Tree has 5 levels. Figure 5a shows the lookup latency (y-axis) of empty queries for variable key sizes (x-axis). Here, we compare the state-of-the-art with a system that employs hash sharing. The time breakdown shows the amount of time spent for BFs (including hashing and filter probing), on performing an I/O to retrieve data, and on other operations (e.g., binary search). Longer keys will cause a higher hashing overhead, therefore, hash sharing becomes more beneficial as the key size grows. Notably, the hash cost for key size 1024B

²Our codebase can be found in <https://github.com/BU-DISC/BF-Shared-Hashing>.

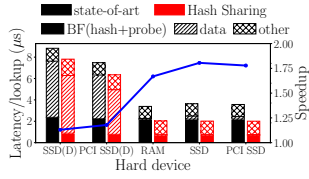


Figure 6: For faster storage, the benefit of hash sharing is pronounced.

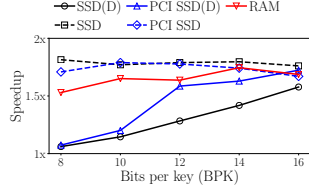


Figure 7: For fewer data accesses (lower FPR), the benefit of hash sharing is larger.

has been reduced by more than one half using hash sharing and the overall query performance is improved by 20%. **is this true? we only have 20% improvement. Is this only the hashing cost maybe?**

Hash Sharing Scales Better with Data Size. The data size growth results in LSM-trees with more levels. In this experiment, we vary the bulk loaded data size, resulting in LSM-Trees with sizes between 8 and 14 levels with size ratio 2. Figure 5b shows that hash sharing is increasingly helpful as the number of levels increases. The state-of-the-art system performs one hash calculation per level, thus the hashing cost accumulates with the number of levels. On the contrary, by sharing the hashing calculation across levels, every query performs only one hash calculation, hence decoupling the hashing cost from the height of the tree. In hash sharing, the increase of the BF cost in the breakdown is a result of the unavoidable filter probing, which is much cheaper than hashing. Overall, hash sharing scales better than the state of the art with the data size.

Hash Sharing Has Higher Impact for Faster Devices. As the storage access latency reduces, the fraction of the time spent hashing increases, to the point it dominates point query latency. We validate this through another experiment by varying the storage device (SSD, PCIe SSD and RAM-disk). We use the RAM-disk to mimic the behavior of a future non-volatile memory with performance close to DRAM. We also experiment with our SSD and PCIe SSD with the direct I/O turned off the cases when direct I/O is turned off for the SSD and the PCIe SSD devices. In Figure 6, SSD(D) and PCI SSD(D) indicates that direct I/O is enabled, while for SSD and PCI SSD it is disabled. This experiment shows that as the data set will reside in increasingly faster devices, the overhead of hashing is pronounced, and as a result, the benefit from hash sharing will lead to a speedup close to 1.8x.

Hash Sharing Has Higher Impact for Lower FPR. Using more bite per key leads to lower FPR, and reduces the number of time a query would read data due to a false positive. We experiment by varying the bits per key between 8 and 14, and in Figure 7 we show that the benefit of hash sharing is higher for decreasing FPR. The reason is that with lower FPR, and hence less false positive data accesses, the fraction of hashing cost dominates the query latency. For example, when BPK is 10 (a typical setting in state-of-the-art LSM-trees), hash sharing can improve the query performance by 20% on PCIe SSD, and for future NVM devices (emulated by RAM disk), such speedup can increase up to 65%.

Hash Sharing Hash Higher Impact for Empty Queries. In the previous experiments we focus on empty point queries, for which hash sharing is most effective. We now examine the benefit of hash sharing as we increase the ratio of queries that will return a

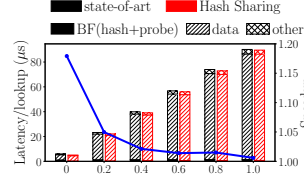


Figure 8: Hash sharing benefits reduce as the existing point lookup ratio grows.

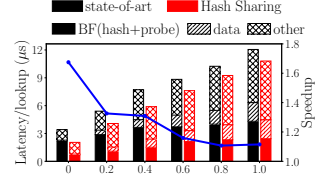


Figure 9: The speedup on RAM disk is pronounced and has similar pattern as SSD

positive result, α , and as a result, they will make at least on disk access. Figure 8 shows the latency breakdowns of the state-of-the-art and hash sharing on the SSD. We observe that that percentage of time spent hashing reduces as α increases. The cost of hashing is constant, however, the average data retrieval cost increases as the number of non-empty queries increases. Figure 9 shows the same experiment on RAM-disk that emulates a future NVM device. We now observe that while the benefit reduces as α increases, it starts from 65% for $\alpha = 0$ and it stays more than 15% for $\alpha = 1$.

6 RELATED WORK

The hashing cost of BFs has been identified as a key optimization, especially when the cost a false positive is low. The textbook implementation of a BF requires k independent hash functions, however, their cost is prohibitively high. To mitigate this cost Less Hashing Bloom Filters (LHBF) [19] and One-Hashing Bloom filters (OHBF) [21] aim to achieve the same false positive ratio while reducing the hashing cost. LHBF divides the filter into k partition with identical size, and calculates the index for partition i using a hash function, $g_i(x) = h_1(x) + i \cdot h_2(x)$, based on two hash functions, $h_1(x)$ and $h_2(x)$. Similarly, OHBF divides the filter into k partitions of uneven sizes and calculates the index for partition i using a single hash function, $g_i(x) = h(x) \% m_i$. The OHBF can be implemented by using only one hash function and a few modulo operations.

Orthogonally to the hashing cost, there have been efforts to reduce the probing cost focusing on the locality of bit-vector accesses [29, 30]. Blocked Bloom filters (BBF) [29] split the filter into a sequence of blocks to reduce memory probing for different locations generated by the k hash functions. BBFs partitions have a small fixed size of one (or a few) cache lines, and unlike classical BFs, the first hash calculation points to a specific block, and all subsequent probes are performed in the same cache line (or group of cache lines). Bloom-1 [30] filter maps k bits in a single word, instead of mapping to an entire filter, in order to reduce the probing cost. Thus, Bloom-1 can achieve membership identification with only one memory access. While BBFs and Bloom-1 are a great match for in-memory workloads, their locality does not benefit disk-resident workloads where the benefit from being cache-efficient is masked by the latency to retrieve data from the disk.

The aforementioned approaches aim to optimize a single BF, while our work aims to optimize a collection of multiple BFs, by sharing hashing not only *within a BF*, but also *across BFs*. Hence, our design can be combined with any techniques that reduces the hashing cost of a single BF.

7 CONCLUSIONS

In this paper we observe that that as we move to faster storage devices, hashing for BFs in LSM-Trees becomes the main bottleneck, and we address this by decoupling the hashing overhead from the number of distinct levels in the tree (and as a result the data size) by sharing a single hash digest across different levels. Our technique reduces the fraction of time spent on hashing during lookups.

REFERENCES

- [1] Ildar Absalyamov, Michael J Carey, and Vassilis J Tsotras. 2018. Lightweight Cardinality Estimation in LSM-based Systems. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. 841–855.
- [2] Wail Y Alkawaileet, Sattam Alsubaiee, and Michael J Carey. 2020. An LSM-based Tuple Compaction Framework for Apache AsterixDB. *Proceedings of the VLDB Endowment* 13, 9 (2020), 1388–1400.
- [3] Apache. [n. d.]. Cassandra. <http://cassandra.apache.org> ([n. d.]).
- [4] Burton H Bloom. 1970. Space/Time Trade-offs in Hash Coding with Allowable Errors. *Commun. ACM* 13, 7 (1970), 422–426.
- [5] Edward Bortnikov, Anastasia Braginsky, Eshcar Hillel, Idit Keidar, and Gali Sheffi. 2018. Accordion: Better Memory Organization for LSM Key-Value Stores. *Proceedings of the VLDB Endowment* 11, 12 (2018), 1863–1875.
- [6] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber. 2006. Bigtable: A Distributed Storage System for Structured Data. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI)*. 205–218.
- [7] Niv Dayan, Manos Athanassoulis, and Stratos Idreos. 2017. Monkey: Optimal Navigable Key-Value Store. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. 79–94.
- [8] Niv Dayan, Manos Athanassoulis, and Stratos Idreos. 2018. Optimal Bloom Filters and Adaptive Merging for LSM-Trees. *ACM Transactions on Database Systems (TODS)* 43, 4 (2018), 16:1–16:48.
- [9] Niv Dayan and Stratos Idreos. 2018. Dostoevsky: Better Space-Time Trade-Offs for LSM-Tree Based Key-Value Stores via Adaptive Removal of Superfluous Merging. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. 505–520.
- [10] Niv Dayan and Stratos Idreos. 2019. The Log-Structured Merge-Bush & the Wacky Continuum. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD)*. 449–466.
- [11] Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall, and Werner Vogels. 2007. Dynamo: Amazon's Highly Available Key-value Store. *ACM SIGOPS Operating Systems Review* 41, 6 (2007), 205–220.
- [12] Siying Dong, Mark Callaghan, Leonidas Galanis, Dhruva Borthakur, Tony Savor, and Michael Strum. 2017. Optimizing Space Amplification in RocksDB. In *Proceedings of the Biennial Conference on Innovative Data Systems Research (CIDR)*.
- [13] Facebook. [n. d.]. RocksDB. <https://github.com/facebook/rocksdb> ([n. d.]).
- [14] Google. [n. d.]. LevelDB. <https://github.com/google/leveldb/> ([n. d.]).
- [15] Google. 2021. CityHash. <https://github.com/google/cityhash> (2021).
- [16] HBase. 2013. Online reference. <http://hbase.apache.org/> (2013).
- [17] Gui Huang, Xuntao Cheng, Jianying Wang, Yujie Wang, Dengcheng He, Tieying Zhang, Feifei Li, Sheng Wang, Wei Cao, and Qiang Li. 2019. X-Engine: An Optimized Storage Engine for Large-scale E-commerce Transaction Processing. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. 651–665.
- [18] Taewoo Kim, Alexander Behm, Michael Blow, Vinayak Borkar, Yingyi Bu, Michael J. Carey, Murtadha Hubail, Shiva Jahangiri, Jianfeng Jia, Chen Li, Chen Luo, Ian Maxon, and Pouria Pirzadeh. 2020. Robust and efficient memory management in Apache AsterixDB. *Software - Practice and Experience* 50, 7 (2020), 1114–1151.
- [19] Adam Kirsch and Michael Mitzenmacher. 2008. Less hashing, same performance: Building a better Bloom filter. *Random Structures & Algorithms* 33, 2 (2008), 187–218.
- [20] Yongkun Li, Chengjin Tian, Fan Guo, Cheng Li, and Yinlong Xu. 2019. ElasticBF: Elastic Bloom Filter with Hotness Awareness for Boosting Read Performance in Large Key-Value Stores. In *Proceedings of the USENIX Annual Technical Conference (ATC)*. 739–752.
- [21] Jianyuan Lu, Tong Yang, Yi Wang, Huichen Dai, Xi Chen, Linxiao Jin, Haoyu Song, and Bin Liu. 2018. Low Computational Cost Bloom Filters. *IEEE/ACM Trans. Netw.* 26, 5 (2018), 2254–2267.
- [22] Chen Luo. 2020. Breaking Down Memory Walls in LSM-based Storage Systems. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. 2817–2819.
- [23] Chen Luo and Michael J Carey. 2019. On Performance Stability in LSM-based Storage Systems. *Proceedings of the VLDB Endowment* 13, 4 (2019), 449–462.
- [24] Chen Luo and Michael J. Carey. 2020. LSM-based Storage Techniques: A Survey. *The VLDB Journal* 29, 1 (2020), 393–418.
- [25] Chen Luo, Pinar Tözün, Yuanyuan Tian, Ronald Barber, Vijayshankar Raman, and Richard Sidle. 2019. Umzi: Unified Multi-Zone Indexing for Large-Scale HTAP. In *Proceedings of the International Conference on Extending Database Technology (EDBT)*. 1–12.
- [26] Siqiang Luo, Subarna Chatterjee, Rafael Ketsetsidis, Niv Dayan, Wilson Qin, and Stratos Idreos. 2020. Rosetta: A Robust Space-Time Optimized Range Filter for Key-Value Stores. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. 2071–2086.
- [27] Patrick E. O'Neil, Edward Cheng, Dieter Gawlick, and Elizabeth J. O'Neil. 1996. The log-structured merge-tree (LSM-tree). *Acta Informatica* 33, 4 (1996), 351–385.
- [28] Georgios Psaropoulos, Ismail Oukid, Thomas Legler, Norman May, and Anastasia Ailamaki. 2019. Bridging the Latency Gap between NVM and DRAM for Latency-bound Operations. In *Proceedings of the International Workshop on Data Management on New Hardware (DAMON)*. 13:1–13:8.
- [29] Felix Putze, Peter Sanders, and Johannes Singler. 2009. Cache-, hash-, and space-efficient bloom filters. *ACM Journal of Experimental Algorithmics* 14 (2009).
- [30] Yan Qiao, Tao Li, and Shigang Chen. 2011. One memory access bloom filters and their generalization. In *INFOCOM 2011. 30th IEEE International Conference on Computer Communications, Joint Conference of the IEEE Computer and Communications Societies, 10-15 April 2011, Shanghai, China*. 1745–1753.
- [31] Simone Raoux, Geoffrey W Burr, Matthew J Breitwisch, Charles T Rettner, Yi-Chou Chen, Robert M Shelby, Martin Salinga, Daniel Krebs, Shih-Hung Chen, Hsiang-Lan Lung, and Chung Hon Lam. 2008. Phase-change random access memory: A scalable technology. *IBM Journal of Research and Development* 52, 4-5 (2008), 465–480.
- [32] Kai Ren, Qing Zheng, Joy Arulraj, and Garth Gibson. 2017. SlimDB: A Space-Efficient Key-Value Storage Engine For Semi-Sorted Data. *Proceedings of the VLDB Endowment* 10, 13 (2017), 2037–2048.
- [33] Ori Rottenstreich and Isaac Keslassy. 2015. The Bloom Paradox: When Not to Use a Bloom Filter. *IEEE/ACM Trans. Netw.* 23, 3 (2015), 703–716.
- [34] Subhadeep Sarkar, Tarikul Islam Papon, Dimitris Staratzis, and Manos Athanassoulis. 2020. Letha: A Tunable Delete-Aware LSM Engine. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. 893–908.
- [35] Jonathan Stern. 2016. Introduction to the Storage Performance Development Kit (SPDK). <https://software.intel.com/content/www/us/en/develop/articles/introduction-to-the-storage-performance-development-kit-spdk.html> (2016).
- [36] Sasu Tarkoma, Christian Esteve Rothenberg, and Emil Lagerspetz. 2012. Theory and Practice of Bloom Filters for Distributed Systems. *IEEE Communications Surveys & Tutorials* 14, 1 (2012), 131–155.
- [37] WiredTiger. [n. d.]. Source Code. <https://github.com/wiredtiger/wiredtiger> ([n. d.]).
- [38] Lei Yang, Hong Wu, Tieying Zhang, Xuntao Cheng, Feifei Li, Lei Zou, Yujie Wang, Rongyao Chen, Jianying Wang, and Gui Huang. 2020. Leaper: A Learned Prefetcher for Cache Invalidation in LSM-tree based Storage Engines. *Proceedings of the VLDB Endowment* 13, 11 (2020), 1976–1989.
- [39] Huanchen Zhang, Hyeontaek Lim, Viktor Leis, David G Andersen, Michael Kaminsky, Kimberly Keeton, and Andrew Pavlo. 2018. SuRF: Practical Range Query Filtering with Fast Succinct Tries. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. 323–336.
- [40] Huanchen Zhang, Hyeontaek Lim, Viktor Leis, David G Andersen, Michael Kaminsky, Kimberly Keeton, and Andrew Pavlo. 2020. Succinct Range Filters. *ACM Transactions on Database Systems (TODS)* 45, 2 (2020), 5:1–5:31.
- [41] Teng Zhang, Jianying Wang, Xuntao Cheng, Hao Xu, Nanlong Yu, Gui Huang, Tieying Zhang, Dengcheng He, Feifei Li, Wei Cao, Zhongdong Huang, and Jianling Sun. 2020. FPGA-Accelerated Compactions for LSM-based Key-Value Store. In *18th USENIX Conference on File and Storage Technologies, FAST 2020, Santa Clara, CA, USA, February 24-27, 2020*. 225–237.
- [42] Yueming Zhang, Yongkun Li, Fan Guo, Cheng Li, and Yinlong Xu. 2018. ElasticBF: Fine-grained and Elastic Bloom Filter Towards Efficient Read for LSM-tree-based KV Stores. In *Proceedings of the USENIX Conference on Hot Topics in Storage and File Systems (HotStorage)*.

Appendices

The cost of lookup in the whole LSM-tree in Eq. (??) can be expanded to (6) using Eq. (3) assuming that memory probing (T_P) is negligible.

$$\begin{aligned} \text{cost} &\approx T_H + \alpha_1 \cdot T_D + (1 - \alpha_1) \cdot f_p \cdot T_D \\ &+ \sum_{i=2}^L \left(1 - \sum_{j=1}^{i-1} \beta_j \right) \cdot [T_H + \alpha_i \cdot T_D + (1 - \alpha_i) \cdot f_p \cdot T_D] \end{aligned} \quad (6)$$

Since $\alpha_i = \frac{\beta_i}{1 - \sum_{j=1}^{i-1} \beta_j}$ and $\alpha_1 = \beta_1$, Eq. (6) becomes Eq. (7).

$$\begin{aligned} \text{cost} &\approx T_H + \beta_1 \cdot T_D + (1 - \beta_1) \cdot f_p \cdot T_D \\ &+ \sum_{i=2}^L \left(1 - \sum_{j=1}^{i-1} \beta_j \right) \cdot \left[T_H + \frac{\beta_i \cdot T_D}{1 - \sum_{j=1}^{i-1} \beta_j} + \left(1 - \frac{\beta_i}{1 - \sum_{j=1}^{i-1} \beta_j} \right) \cdot f_p \cdot T_D \right] \\ &= \left(L - \sum_{i=2}^L \sum_{j=1}^{i-1} \beta_j \right) \cdot T_H + \sum_{i=1}^L \beta_i \cdot T_D + \left(L - \sum_{i=2}^L \sum_{j=1}^{i-1} \beta_j - \sum_{i=1}^L \beta_i \right) \cdot f_p \cdot T_D \end{aligned} \quad (7)$$

Since $\sum_{i=1}^L \beta_i = \alpha$, Eq. (7) becomes Eq. (8).

$$\text{cost} \approx \left(L - \sum_{i=2}^L \sum_{j=1}^{i-1} \beta_j \right) \cdot T_H + \alpha \cdot T_D + \left(L - \sum_{i=2}^L \sum_{j=1}^{i-1} \beta_j - \alpha \right) \cdot f_p \cdot T_D \quad (8)$$

If we assume that keys in the LSM-tree is perfectly uniform, then, β_i depends on α and the size of level i , i.e., $\beta_i = T^{i-1} \cdot \beta_1$, while $\beta_1 = \frac{\alpha}{\sum_{j=1}^L T^{j-1}}$. Thus, $\sum_{i=2}^L \sum_{j=1}^{i-1} \beta_j$ can be approximated as Eq. (9).

$$\begin{aligned} \sum_{i=2}^L \sum_{j=1}^{i-1} \beta_j &= \sum_{i=2}^L \sum_{j=1}^{i-1} T^{j-1} \cdot \beta_1 = \\ &= \beta_1 \cdot \sum_{i=2}^L \frac{T^{i-1} - 1}{T - 1} = \frac{\beta_1}{T - 1} \cdot \left(\sum_{i=2}^L (T^{i-1}) - (L - 1) \right) = \\ &= \frac{\beta_1}{T - 1} \cdot \left(\sum_{i=1}^L (T^{i-1}) - 1 - (L - 1) \right) = \frac{\beta_1}{T - 1} \cdot \left(\sum_{i=1}^L (T^{i-1}) - L \right) = \\ &= \frac{\alpha}{\sum_{k=1}^L T^{k-1}} \cdot \frac{1}{T - 1} \cdot \left(\sum_{i=1}^L (T^{i-1}) - L \right) = \frac{\alpha}{T - 1} \cdot \left(1 - \frac{L}{\sum_{k=1}^L T^{k-1}} \right) = \\ &= \frac{\alpha}{T - 1} \cdot \left(1 - \frac{L}{\frac{T^L - 1}{T - 1}} \right) = \frac{\alpha}{T - 1} \cdot \left(1 - \frac{L \cdot (T - 1)}{T^L - 1} \right) \\ &(\text{for } T > 3 \text{ or } L > 3) \approx \frac{\alpha}{s - 1} \end{aligned} \quad (9)$$

Therefore, the cost of lookup in the whole LSM-tree in Eq. (??) can be approximated as Eq. (4).