

Constructing and Analyzing the LSM Compaction Design Space [Experimental Analysis Paper]

Subhadeep Sarkar

Dimitris Staratzis

Zichen Zhu

Manos Athanassoulis

Boston University

ABSTRACT

Log-structured merge (LSM) trees offer efficient ingestion by appending incoming data. LSM-trees are widely used as the storage layer of production NoSQL data stores. To enable competitive read performance, LSM-trees periodically re-organize data to form a tree with levels of exponentially increasing capacity, through iterative *compactions*. Compactions fundamentally influence the performance of an LSM-engine in terms of write amplification, write throughput, point and range lookup performance, space amplification, and delete performance. Hence, choosing the *appropriate compaction strategy* is crucial and, at the same time, hard as the LSM-compaction design space is vast, largely unexplored, and has not been formally defined in the literature. As a result, most LSM-based engines use a fixed compaction strategy, typically hand-picked by an engineer, which decides *how* and *when* to compact data.

In this paper, we present the design space of LSM-compactions, and evaluate state-of-the-art compaction strategies with respect to their key performance metrics. Toward this goal, our first contribution is to introduce a set of design primitives that can formally define any compaction strategy. We identify four primitives, namely (i) compaction trigger, (ii) compaction eagerness, (iii) compaction granularity, and (iv) data movement policy, which can synthesize any existing and also completely new compaction strategies. Our second contribution is to experimentally analyze compaction strategies. We present 23 observations and 10 high-level takeaway messages, which fundamentally show how LSM systems can navigate the design space of compactions.

1 INTRODUCTION

LSM-based Key-Value Stores. Log-structured merge (LSM) trees are widely used today as the storage layer of modern NoSQL key-value stores [31, 35, 37]. LSM-trees employ the *out-of-place* paradigm to achieve fast ingestion. Incoming key-value pairs are buffered in main memory, and are periodically flushed to persistent storage as *sorted immutable runs*. As more runs accumulate on disk, they are sort-merged to construct fewer yet longer sorted runs. This process is known as *compaction* [25, 35]. To facilitate fast *point lookups*, LSM-trees use auxiliary in-memory data structures (Bloom filters and fence pointers) that help to reduce the average number of disk I/Os performed per lookup [18, 19]. Owing to these advantages, LSM-trees are adopted by several production key-value stores including

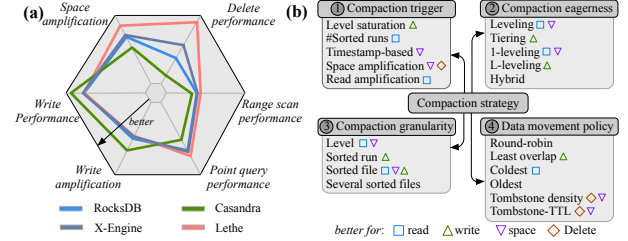


Fig. 1: (a) The different compaction strategies adopted in state-of-the-art LSM-engines lead to the diverse performances offered by the engines; (b) The taxonomy of LSM compactions in terms of the design primitives.

LevelDB [27] and BigTable [15] at Google, RocksDB [25] at Facebook, X-Engine [29] at Alibaba, WiredTiger at MongoDB [53], CockroachDB at Cockroach Labs [16], Voldemort [34] at LinkedIn, DynamoDB [22] at Amazon, AsterixDB [2], Cassandra [6], HBase [7], Accumulo [5] at Apache, and bLSM [49] and cLSM [26] at Yahoo. Academic systems based on LSM-trees include Monkey [18], SlimDB [40], Dostoevsky [19, 20], LSM-Bush [21], Lethe [44], Silk [10, 11], LSbM-tree [51], SifrDB [36], and Leaper [54].

Compactions in LSM-Trees. Compactions in LSM-trees are employed periodically to *reduce read and space amplification at the cost of write amplification* while ensuring data consistency and query correctness [8, 9]. A compaction merges two or more sorted runs, between one or multiple levels to ensure that the LSM-tree maintains levels with exponentially increasing sizes [37]. Compactions are typically invoked when a level reaches its capacity, at which point, the compaction routine moves data from the saturated level to the next one, that has an exponentially larger capacity. Any duplicate entries (resulting from *updates*) and invalidated entries (resulting from *deletes*) are removed during a compaction, retaining only the logically correct (latest valid) version [24, 38]. A compaction either (a) sort-merges two consecutive levels in their entirety – *full compaction* or (b) sort-merges a part of a level only with the overlapping part from the next level – *partial compaction*. Compactions dictate *how* and *when* disk-resident data is re-organized, and thereby, influence the physical data layout on the disk. Fig. 1(a) presents a qualitative comparison of the performance implications of the various compaction strategies adopted in state-of-the-art LSM engines.

The Challenge: Hand-Picking Compaction Strategies. Despite compactions being critical to the performance of LSM-engines, the process of *choosing an appropriate compaction strategy* requires a human in the loop. In practice, decisions on “*how to (re-)organize data on disk*”, and thereby, “*which compaction strategies to implement or use*” in a production LSM-based data store are often subject to the expertise of the engineers or the database administrators (DBAs). This is largely due to two reasons. First, the process of compaction in LSM-trees is often treated as a black-box and is rarely exposed as

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.
doi:XX.XX/XXX.XX

a tunable knob. While the LSM-compaction design space is vast, the lack of a formal template for compactions leads to heavily relying on individual expertise, leaving a large part of the design space unexplored. Second, there is a lack of analytical and experimental data on how compactions influence the performance of an LSM-engine subject to the underlying design of the storage engine and the workload characteristics. Hence, it is difficult, even for experts, to answer design questions such as:

- (i) My LSM-engine is offering lower write performance than expected: *Would a change in the compaction strategy help? If yes, which strategies should be used?*
- (ii) The workload we used to process has changed: *How does this affect the read throughput of my system? Is there a compaction strategy that can improve the read throughput?*
- (iii) We are due to design a new LSM-engine for processing a specific workload: *How should I compact my data for best overall performance? Is there a compaction strategy that I must avoid?*

Relying on human expertise to hand-pick the appropriate compaction strategies for each application does not scale, especially for large-scale system deployments.

Contributions. To this end, in this work, we formalize the design space of compactions in LSM-based storage engines. Further, we experimentally explore this space, and based on this, we present 10 high-level takeaway message, and 23 observations that serve as a comprehensive set of guidelines for LSM-compactions, and lay the groundwork for compaction tuning and automation.

Conceptual Contribution: Constructing the Compaction Design Space. We identify the defining characteristics of a compaction, or compaction *primitives*: (i) the **trigger** (i.e., *when* to compact), (ii) the **eagerness** (i.e., *how* to organize the data after compaction), (iii) the **granularity** (i.e., *how much* data to compact at a time), and (iv) the **data movement policy** (i.e., *which* data to compact). Together, the four primitives define when and how to compact data in an LSM-tree. Fig. 1(b) presents the taxonomy of LSM-compactions along with the various options for each of the design primitives.

Experimental Contribution 1: Unify the Experimental Infrastructure of Multiple Compaction Strategies. To establish a consistent experimental platform, we integrate several state-of-the-art compaction strategies into a unified codebase, based on the widely adopted open-source RocksDB [25] LSM-engine. This integration bridges wild variations of implementation and configuration knobs of different compaction strategies across different LSM-engines. Further, we implement each compaction strategy (1) through the prism of the aforementioned four primitives, and (2) on top of the same data store to ensure an apples-to-apples comparison. We implement *nine state-of-the-art compaction strategies* that are popular among production and academic systems and are key to the understanding of the LSM-compaction design space. We implement these strategies through significant modifications to the latest RocksDB codebase, and expose *more than a hundred design knobs* to enable custom configuration and to ensure a fair evaluation.

Experimental Contribution 2: Analyzing the Compaction Design Space. We provide a comprehensive experimental analysis of the LSM-compaction design space, which quantifies the impact of each of the design primitives on a number of performance metrics. This experimental analysis also serves as a roadmap for selecting

a compaction strategy subject to the workload characteristics and performance goals. We perform more than 2000 experiments with the nine compaction strategies to take a deep dive on the following.

- **Performance Implications of Compactions.** We quantify the impact of compactions on LSM performance in terms of ingestion throughput, query latency, space and write amplification, and delete efficacy in §5.1.
- **Workload Influence on Compactions.** While the composition and (ingestion and access) distribution of the workload influence the compaction performance, deciding which compaction strategy to employ is workload-agnostic in existing systems. To analyze the workload’s impact on compactions performance, we experiment with a number of representative workloads by varying (i) the size of ingested data, (ii) the proportion of ingestion and lookups, (iii) the proportion of empty and non-empty point lookups, (iv) the selectivity of range queries, (v) the fraction of updates and (vi) deletes, (vii) the key-value size, as well as (viii) the workload distribution (uniform, normal, and Zipfian) in §5.2.
- **Tuning Influence on Compactions.** LSM tuning typically focuses on knobs like memory buffer size, page size, and size ratio which are not believed to be connected with compaction performance. We experiment with these knobs to uncover when compactions are affected (and when not) by these knobs in §5.3.
- **Answering Design Questions.** Finally, throughout §5 we present various observations and key insights of our experimental evaluation, and in §6 we discuss a roadmap for designing and choosing compaction in LSM engines.

This work defines the LSM compaction design space and presents a thorough account of how the different primitives affect compactions and how each compaction strategy, in turn, affects the overall performance of a storage engine.

Key Takeaways. Finally, the high-level key takeaways from our analysis are the following.

A. There is no perfect compaction strategy. When it comes to selecting a compaction strategy for an LSM engine, there is *no single best*. Thus, a compaction strategy needs to be custom-tailored to specific combinations of workload, LSM tuning, and performance goals.

B. It is important to look into the compaction “black-box”. To understand the performance implications of LSM compactions, it is crucial to stop treating them as a black-box; rather, as an ensemble of fundamental design primitives. Following this approach we reason about the performance implications of each design primitive independently. We identify and avoid common pitfalls given a workload and a target performance.

C. Switching compaction strategies can significantly boost the performance of an LSM engine. Switching between compaction strategies as the workload changes and/or the performance goals shift can boost the performance of an LSM engine significantly. Understanding the behavior and performance implications of the compaction primitives allows for minor modifications to existing codebases to invoke the appropriate compaction strategy, as necessary.

2 BACKGROUND

We now present the necessary background of LSM-trees and its operations as well as the terminology we use in throughout the

paper. A more detailed discussion/survey on LSM-basics can be found in the literature [18, 35].

LSM-Basics. To support fast data ingestion, LSM-trees buffer incoming inserts, updates, and deletes (i.e., ingestion, in general) within main memory. Once the memory buffer becomes full, the entries contained are sorted on the key and the buffer is flushed as a *sorted run* to the disk-component of the tree. In practice, a sorted run is a collection of one or more *immutable files* that have typically the same size. For an LSM-tree with L levels, we assume that its first level (Level 0) is an in-memory buffer and the remaining levels (Level 1 to $L - 1$) are disk-resident [18, 35]. On disk, each Level i ($i > 1$) has a capacity that is larger than that of Level $i - 1$ by a factor of T , where T is the size ratio of the tree.

LSM-Compactions. To limit the number of sorted runs on disk (and thereby, to facilitate fast lookups and better space utilization), LSM-trees periodically sort-merge runs (or parts of a run) from a Level i with the overlapping runs (or parts of runs) from Level $i + 1$. This process of data re-organization and creating fewer longer sorted runs on disk is known as *compaction*. However, the process of sort-merging data from different levels requires the data to be moved back and forth between the disk and main memory. This results in write amplification, which can be as high as $40\times$ in state-of-the-art LSM-based data stores [39].

Partial compactions. To amortize data movement, and thus, avoid latency spikes, state-of-the-art LSM engines organize data into smaller files, and perform compactions at the granularity of files instead of levels [24]. If Level i grows beyond a threshold, a compaction is triggered and one file (or a subset of files) from Level i is chosen to be compacted with files from Level $i + 1$ that have an overlapping key-range. This process is known as *partial compaction*. The decision on which file(s) to compact depends on the design of the storage engine. Fig. 2 presents a comparative illustration of the full compaction and partial compaction routines in LSM-trees.

Querying LSM-Trees. Since LSM-trees realize updates and deletes in an out-of-place manner, multiple entries with the same key may exist in a tree with only the recent-most version being valid. Thus, a query must find and return the recent-most version of an entry.

Point lookups. A point lookup starts at the memory buffer and traverses the tree from the smallest level to the largest one, and from the youngest run to the oldest one within a level. The lookup terminates immediately after a match of the target key is found. To limit the number of runs a lookup probes, state-of-the-art LSM-based data stores use in-memory data structures, such as Bloom filters and fence pointers [20, 25].

Range scans. A range scan requires sort-merging the runs qualifying for a range query across all levels of the tree. The runs are sort-merged in memory and the latest version for each qualifying entry is returned while discarding all older, logically invalidated versions.

Deletes in LSM-Trees. Deletes in LSM-trees are performed logically without necessarily disturbing the target entries. A point delete operation is realized by inserting a special type of key-value entry, known as a *tombstone*, that logically invalidates the target entry. During compactions, a tombstone purges any older entries with a matching key. A delete is eventually considered as *persistent* once the corresponding tombstone reaches the last tree-level, at which point the tombstone can be safely dropped. The time taken

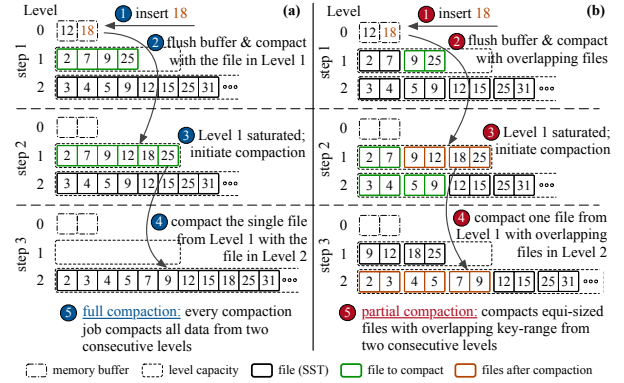


Fig. 2: (a) When invoked, the classical full compaction routine compacts whole levels at a time, while (b) in partial compaction, data is stored in equi-sized files which allows compactions to be performed at the granularity of files.

to persistently delete a data object from an LSM-based data store depends on process of data re-organization. Compactions, thus, also play a critical role in timely and persistent deletion of entries.

3 THE COMPACTION DESIGN SPACE

In this section, we identify the *design primitives* that provide a structured decomposition of arbitrary compaction strategies. This allows us to create the taxonomy of the universe of LSM compaction strategies including all the classical as well as new ones.

3.1 Compaction Primitives

We define a compaction strategy as an *ensemble of design primitives that represents the fundamental decisions about the physical data layout and the data (re-)organization policy*. Each primitive answers a fundamental design question.

- 1) **Compaction trigger:** **When** to re-organize the data layout?
- 2) **Compaction eagerness:** **How** to lay out the data physically on the persistent storage media?
- 3) **Compaction granularity:** **How much** data to move at-a-time during layout re-organization?
- 4) **Data movement policy:** **Which** block of data to be moved during re-organization?

Together, these design primitives define *when* and *how* an LSM-engine re-organizes the data layout on the persistent media. The proposed primitives capture any state-of-the-art LSM-compaction strategy and also enables synthesizing new or unexplored compaction strategies. Below, we define these four design primitives.

3.1.1 Compaction Trigger. Compaction triggers refer to the set of events that can initiate a compaction job. The most common compaction trigger is based on the *degree of saturation* of a level in an LSM-tree [2, 25–27, 30, 49, 50]. The degree of saturation for Level i ($1 \leq i \leq L - 1$) is typically measured as the ratio of the number of bytes of data stored in Level i to the theoretical capacity in bytes for Level i . Once the degree of saturation goes beyond a pre-defined threshold, one or more immutable files from Level i are marked for compaction. Some LSM-engines use the file count in a level to compute degree of saturation [27, 29, 30, 42, 48]. Note that, the file count-based degree of saturation works only when all immutable files are of equal size or in systems that have a tunable file size.

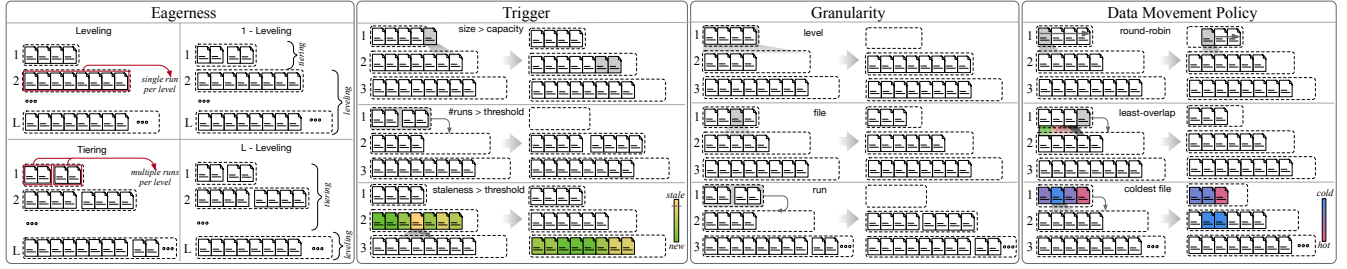


Fig. 3: The primitives that define LSM compactions: trigger, eagerness, granularity, and data movement policy.

In practice, a compaction job starts as soon as all the necessary resources, i.e., memory, CPU, and device bandwidth, are available.

Other compaction triggers include the *staleness of a file*, the *tombstone-based time-to-live*, *space amplification*, and even *read amplification*. For example, to ensure propagation of updates and deletes to the deeper levels of a tree, some LSM-engines assign a time-to-live (TTL) for each file during its creation. Each file can live in a level for a bounded time, and once the TTL expires, the file is marked for compaction [25]. Another delete-driven compaction trigger ensures bounded persistence latency of deletes in LSM-trees through a different timestamp-based scheme. Each file containing at least one tombstone is assigned a special time-to-live in each level, and up on expiration of this timer, the file is marked for compaction [44]. Compaction triggers based on space amplification [42] and read amplification [33] that can facilitate specific applications have also been implemented in production data stores. Below, we present a list of the most common **compaction triggers**:

- i) **Level saturation**: level size goes beyond a nominal threshold
- ii) **#Sorted runs**: sorted run count for a level reaches a threshold
- iii) **Timestamp-based**: a file lives in a level for too long
- iv) **Space amplification (SA)**: overall SA surpasses a threshold
- v) **Read amplification (RA)**: expected RA surpasses a threshold

3.1.2 Compaction Eagerness. The compaction eagerness determines the data layout on storage by controlling the number of sorted runs per level. Compactions move data back and forth between storage and memory, consuming a significant proportion of the device bandwidth. There is, thus, an inherent competition for the device bandwidth between ingestion (external) and compaction (internal) – a trade-off depending on the eagerness of compactions.

The compaction eagerness is commonly classified into two categories: leveling and tiering [18, 19]. With leveling, once a compaction is triggered in Level i , the file(s) marked for compaction are merged with the overlapping file(s) from Level $i + 1$, and the result is written back to Level $i + 1$. As a result, Level $i + 1$ ends up with a (single) longer sorted run of immutable files [25–27, 29, 30, 49]. For tiering, each level may contain more than one sorted runs with overlapping key domains. Once a compaction is triggered in Level i , all sorted runs in Level i are merged together, the result of which is written to Level $i + 1$ as a new sorted run without disturbing the existing runs in that level [2, 6, 7, 25, 48, 50].

A hybrid design is proposed in Dostoevsky [20] where the last level is implemented as leveling and all the remaining levels on disk are tiered. The number of runs in each of these levels is configurable, subject to the workload composition. A generalization of this idea is proposed in the as a continuum of designs [21, 32] that allows each level to separately decide between leveling and tiering. Among

production systems, RocksDB implements the first disk-level (Level 1) as tiering [42], and it is allowed to grow perpetually in order to avoid write-stalls [10, 11, 13] in ingestion-heavy workloads. Note that some production LSM-engines internally assign a higher priority to writes over compactions to avoid write stalls; however, this tuning leads to a tree structure that violates the LSM-tree properties, and thus, is out of the scope of our analysis. Below is a list of the most common options for **compaction eagerness**:

- i) **Leveling**: every time a sorted run arrives to a level
- ii) **Tiering**: when the #sorted runs per level reaches a threshold
- iii) **1-leveling**: tiering for Level 1; leveling otherwise
- iv) **L-leveling**: leveling for last level; tiering otherwise
- v) **Hybrid**: a level can be tiering or leveling independently

3.1.3 Compaction Granularity. Compaction granularity refers to the amount of data moved during a single compaction job, once the compaction trigger fires. In the classical LSM-design [37], once a level reaches its capacity, all data from that level are moved to the next level through compaction. To amortize the I/O costs due to compactions, state-of-the-art leveled LSM-based engines employ *partial compaction* [25, 27, 29]. In partial compaction, instead of moving a whole level, a smaller granularity of data participates in every compaction. The granularity of data can be a single file [24, 29, 44] or multiple files [1, 2, 6, 37] depending on the system design and the workload. Note that, partial compaction does not radically change the total amount of data movement due to compactions, but amortizes this data movement uniformly over time, thereby preventing undesired latency spikes. The exact data movement due to compaction, is determined by the size ratio of a tree and by the data movement policy we discuss next. Below, we present a list of the most common **compaction granularity** options:

- i) **Level**: all data in two consecutive levels
- ii) **Sorted run**: all sorted runs in a level
- iii) **Sorted file**: one sorted file at a time
- iv) **Several sorted files**: several sorted files at a time

3.1.4 Data Movement Policy. The data movement policy is specifically relevant for LSM-engines adopting *partial compaction* as it involves a decision-making on which file(s) to choose for compaction. While the literature commonly refers to this decision for partial compaction as *file picking policy* [23], we use the term *data movement* to generalize for any possible data movement granularity. Moving a file from a shallower level (Level i) to a deeper level (Level $i + 1$) affects point lookup performance, write and space amplification, as well as delete performance.

A naïve way to choose file(s) is at random or by using a round-robin policy [27, 30]. These data movement policies do not focus on optimizing for any particular performance metric, but help in

reducing space amplification. To optimize for read throughput, many production data stores [25, 29] select the “coldest” file(s) in a level once a compaction is triggered. Another common optimization goal is to minimize write amplification. In this policy, files with the least overlap with the target level are marked for compaction which minimizes the data movement [12, 23]. To reduce space amplification, some storage engines choose files with the highest number of tombstones and/or updates [25]. Another delete-aware approach introduces a tombstone-timestamp driven file picking policy that aims to timely persist logical deletes [44]. Note that, as LSM-trees store data in immutable files, the smallest granularity of data movement is typically a file. Below, we present the list of the most common **data movement policies**:

- i) **Round-robin**: chooses files in a round-robin manner
- ii) **Least overlapping parent**: file with least overlap with “parent”
- iii) **Least overlapping grandparent**: as above with “grandparent”
- iv) **Coldest**: the least recently accessed file
- v) **Oldest**: the oldest file in a level
- vi) **Tombstone density**: file with #tombstones above a threshold
- vii) **Tombstone-TTL**: file with expired tombstones-TTLs

3.2 Compaction as an Ensemble of Primitives

Synthesizing Diverse Compaction Strategies. By definition, every compaction strategy takes one or more values for each of the four primitives. The trigger, granularity, and data movement policy are multi-valued primitives, whereas eagerness is single-valued.

For example, the classical LSM-tree [37] is a **leveled** tree (*eagerness*) that compacts **whole levels** at a time (*granularity*) once a **level reaches a nominal size** (*trigger*). The classical LSM-design does not implement many subtle optimizations including partial compactions, and by definition, does not need a data movement policy. A more complex example is the compaction strategy for a **leveled** LSM-tree (*eagerness*) in which compactions are performed at the *granularity* of a **file**. A compaction is *triggered* if either (a) a **level reaches its capacity** or (b) a **file containing tombstones is retained in a level longer than a pre-set TTL** [44]. Once triggered, the *data movement policy* chooses (a) **the file with the highest density of tombstones**, if there is one or (b) **the file with the least overlap with the parent level**, otherwise.

The Compaction Design Space Cardinality. Two compaction strategies are considered different from each other if they differ in at least one of the four primitives. Two compaction strategies that are identical in terms of three of the four primitives, but only differ in one (say, data movement), can have vastly different performance when subject to the same workload while running on identical hardware. Plugging in some typical values for the cardinality of the primitives, we estimate the cardinality of the compaction universe as $>10^4$, a vast yet largely unexplored design space. Table 1 shows a representative part of this space, detailing the compaction strategies used for more than twenty academic and production systems.

Compactions Analyzed. For our analysis and experimentation, we select nine representative compaction strategies that are prevalent in production and academic LSM-based systems. These strategies capture the wide variety of the possible compaction designs. We codify and present these candidate compaction strategies in Table 2. Full1 represents the classical compaction strategy for leveled LSM-trees that compacts two consecutive levels upon invocation. LO+1

Database	Compaction Eagerness	Compaction Trigger				Compaction Granularity			Data Movement Policy									
		Level saturation	#Sorted runs	File staleness	Space amp.	Tombstone-TTL	Level	Sorted run	File (single)	File (multiple)	Round-robin	Least overlap (+1)	Least overlap (+2)	Coldest file	Oldest file	File w/ most TS	Expired TS-TTL	N/A
RocksDB [25], Monkey [19]	Leveling / L-leveling	✓	✓					✓	✓		✓		✓	✓	✓	✓		
	Tiering		✓	✓	✓		✓											✓
LevelDB [27], Monkey (J) [18]	Leveling	✓						✓		✓	✓	✓						
SlimDB [40]	Tiering	✓							✓	✓								✓
Dostoevsky [20]	L-leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T				✓ ^L						✓ ^T
LSM-Bush [21]	Hybrid leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T				✓ ^L						✓ ^T
Lethe [44]	Leveling	✓				✓		✓	✓	✓	✓							✓
Silk [10], Silk+ [11]	Leveling	✓						✓	✓	✓	✓							
HyperLevelDB [30]	Leveling	✓						✓	✓	✓	✓	✓						
PebblesDB [39]	Hybrid leveling	✓						✓	✓									✓
Cassandra [6]	Tiering		✓	✓		✓		✓				✓						✓
	Leveling	✓						✓	✓	✓	✓					✓	✓	
WiredTiger [53]	Leveling	✓						✓										✓
X-Engine [29], Leaper [54]	Hybrid leveling	✓						✓	✓		✓					✓		
HBase [7]	Tiering		✓					✓										✓
AsterixDB [2]	Leveling	✓						✓										✓
	Tiering		✓					✓										✓
Tarantool [50]	L-leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T										✓
ScyllaDB [48]	Tiering		✓	✓		✓		✓	✓		✓						✓	✓
	Leveling	✓						✓	✓	✓						✓		
blSM [49], cLSM [26]	Leveling	✓						✓			✓							
Accumulo [5]	Tiering	✓	✓					✓										✓
LSBM-tree [51, 52]	Leveling	✓						✓										✓
SiftDB [36]	Tiering	✓							✓									✓

Table 1: Compaction strategies in state-of-the-art systems. [✓^L: for levels with leveling; ✓^T: for levels with tiering.]

and LO+2 denote two partial compaction routines that choose a file for compaction with the smallest overlap with files in the parent ($i + 1$) and grandparent ($i + 2$) level, respectively. RR chooses files for compaction in a round-robin fashion from each level. ColD and Old are read-friendly compaction strategies that mark the coldest and oldest file(s) in a level for compaction, respectively. TSD and TSA are delete-driven compaction strategies with triggers and data movement policies that are determined by the density of tombstones and the age of the oldest tombstone contained in a file, respectively. Finally, Tier represents the variant of tiering compaction with a trigger of space amplification.

4 BENCHMARKING COMPACTIONS

We now discuss the benchmarking process, discussing our experimental platform, how we integrated new compactions policies, and our measurement methodology.

4.1 Standardization of Compaction Strategies

We choose RocksDB [25] as our experimental platform, as it (i) is open-source, (ii) is widely used across industry and academia, (iii) has a large active community. To ensure fair comparison we implement all compaction strategies under the same LSM-engine.

Implementation. We integrate our codebase into RocksDB v6.11.4. We assign to *compactions* a *higher priority than writes* to accurately measure their performance implications, while guaranteeing that the LSM structure is maintained during workload execution [47].

Compaction Trigger. The default compaction trigger for (hybrid) leveling in RocksDB is level saturation [41], and for the universal compaction is space amplification [42]. RocksDB also supports delete-driven compaction triggers, specifically whether the #tombstones in a file goes beyond a threshold. We further implement a trigger based on the tombstones age to facilitate timely deletes [44].

Primitives	Full [2, 37, 53]	LO+1 [19, 25, 44]	Cold [25]	Old [25]	TSO [25, 44]	RR [26, 27, 30, 49]	LO+2 [27, 30]	TSA [43]	Tier [6, 28, 40, 48]
Compaction trigger	level saturation	level saturation	level saturation	level saturation	1. TS-density 2. level saturation	level saturation	level saturation	1. TS age 2. level saturation	level saturation
Compaction eagerness	leveling	leveling	leveling	leveling	leveling	leveling	leveling	leveling	tiering
Compaction granularity	levels	files	files	files	files	files	files	files	sorted runs
Data movement policy	N/A	least overlapping parent	coldest file	oldest file	1. most tombstones 2. least overlapping parent	round-robin	least overlapping grandparent	1. expired TS-TTL 2. least overlapping parent	N/A

Table 2: Compaction strategies evaluated in this work.

Compaction Eagerness. By default, RocksDB supports only two degrees of eagerness: *hybrid leveling* (tiered first level, leveled otherwise) [41] and a variation of *tiering* (with a different trigger), termed universal compaction [42]. We also implement pure *leveling* by limiting the number of first-level runs to one, and triggering a compaction when the number of first-level files is more than one.

Compaction Granularity. The granularity for leveling is *file* and *sorted runs* for tiering. To implement classical leveling, we mark all files of a level for compaction. We ensure that ingestion may resume only after all the compaction-marked files are compacted thereby replicating the behavior of the full compaction routine.

Data Movement Policy. RocksDB (v6.11.4) provides four different data movement policies: a file (i) with least overlap with its parent level, (ii) least recently accessed, (iii) with the oldest data in a level, and (iv) that has more tombstones than a threshold. We also implement partial compaction strategies that choose a file (v) in a round-robin manner, (vi) with the least overlap with its grandparent level, and (vii) based on the age of the tombstones in a file.

Designing the Compaction API. We expose the compaction primitives through a new API as configurable knobs. An application can configure the desired compaction strategy and initiate workload execution. The API also allows the application to change the compaction strategy for an existing database. Overall, our experimental infrastructure allows us (i) to ensure an identical underlying structure while setting the compaction benchmark, and (ii) to tune and configure the design of the LSM-engine as necessary.

4.2 Performance Metrics

Compactions affect both LSM reads and writes. Below, we present the performance metrics used in our experimental evaluation.

Write Amplification (WA). The repeated reads and writes due to compaction cause high WA [39]. We formally define WA as *the number of times an entry is (re-)written without any modifications to disk during its lifetime. We use the RocksDB metric `compact.write.bytes` and the actual data size to compute WA.*

Write Latency. Write latency is driven by the device bandwidth utilization, which depends on (i) write stalls due to compactions and (ii) the sustained device bandwidth. *We use the `db.write.micros` histogram to measure the average and tail of the write latency.*

Read Amplification (RA). RA is the ratio between the number of disk pages read for a point lookup and the number of point lookups. *We use `rocksdb.bytes.read` to compute RA.*

Point Lookup Latency. Compactions determine the position of the files in an LSM-tree which affects point lookups on entries contained in those files. *Here, we use the `db.get.micros` histogram.*

Range Lookup Latency. The range lookup latency depends on the selectivity of the range query, but is affected by compaction eagerness. *We also use the `db.get.micros` histogram for range lookups.*

Space Amplification (SA). SA depends on compaction eagerness, compaction granularity, and the data movement policy. SA is defined as *the ratio between the size of logically invalidated entries and the size of the unique entries in the tree* [20]. *We compute SA using the size of the database and the size of the logically valid entries.*

Delete Performance. We measure the degree to which the tested approaches persistently delete entries within a time-limit [44]. *We use the RocksDB file metadata age and a delete persistence threshold.*

4.3 Benchmarking Methodology

Having discussed the benchmarking approaches and the metric measured, we now discuss the methodology for varying the key input parameters for our analysis: *workload* and the *LSM tuning*.

4.3.1 Workload. A typical key-value workload comprises of five primary operations: inserts, updates, point lookups, range lookups, and deletes. Point lookups target keys that may or may not exist in the database – we refer to these as *non-empty* and *empty point lookups*, respectively. Range lookups are characterized by their *selectivity*, and deletes target keys that are present in the database. To analyze the impact of each operation in our experiments we vary the *fraction* of each operation as well as their qualitative characteristics (i.e., empty vs. non-empty, selectivity, key size, and value size). We further vary the *data distribution* of ingestion and read queries focusing on (i) uniform, (ii) normal, and (iii) Zipfian distributions. Overall, our custom-built benchmarking suite is a superset of the influential YCSB benchmark [17] as well as the insert benchmark [14], and supports a number of important knobs that are missing from existing workload generators, including deletes. Our workload generator exposes over 64 degrees of freedom, and is available via GitHub [46] for dissemination, testing, and adoption.

4.3.2 LSM Tuning. We further study the interplay of LSM tuning and compaction strategies. We consider questions like *which compaction strategy is appropriate for a specific LSM design and a given workload?* To answer such questions we vary in our experimentation key LSM tuning parameters, like (i) the memory buffer size, (ii) the block cache size, and (iii) the size ratio of the tree.

5 EXPERIMENTAL EVALUATION

We now present the key experimental results using the nine compaction strategies listed in Table 2.

Goal of the Study. Our analysis aims to answer the following three fundamental questions:

- i) *Performance implications:* How do compactions affect the overall performance of LSM-engines?
- ii) *Workload influence:* How do workload distribution and composition influence compactions, and thereby, the performance of LSM-engine?

- iii) *Tuning influence*: What is the interplay between LSM compactions and tuning?

Ultimately, the goal of this study is to help practitioners and researchers to make informed decisions when deciding which compaction strategies to support and use in an LSM-based engine.

Experimental Setup. All experiments were run on AWS EC2 server with instances of type t2.xlarge (virtualization type: hardware virtual machine) [4]. The virtual machines (VMs) used were supported by 8 Intel Scalable Processors (vCPUs) clocked at 3.0GHz each, 32GB of DIMM RAM, 45MB of L3 cache, and were running Ubuntu 20.04 LTS. Each VM had an attached 40GB SSD volume with 4000 provisioned IOPS (volume type: io2) [3].

Default Setup. Unless otherwise mentioned all experiments are performed on a RocksDB setup with an LSM-tree of size ratio 10. The memory buffer is implemented as a skiplist. The size of the write buffer is set as 8MB which can hold up to 512 16KB disk pages. Fence pointers are maintained for each disk page and Bloom filters are constructed for every file with 10 bits memory allocated for every entry. Additionally, we have 8MB block cache (RocksDB default) assigned for data, filter, and index blocks. To capture the true raw performance of RocksDB as an LSM-engine, we (i) assign compactions a higher priority than writes, (ii) enable direct I/Os for both read and write operations, and (iii) set the number of background threads responsible for compactions to 1.

Workloads. We experiment with a YCSB variation. The average size of a key-value entry is 128B with the key size of 8B. We vary the number of inserts, going up to 2^{26} operations. As compaction performance proves to be agnostic to data size, and in the interest of experimenting with many configurations, we perform our base experiments with 10M inserts. Unless otherwise mentioned, ingestion and lookups are uniform and interleaved.

Presentation. For each experiment, we present the primary observations (O) along with key takeaway (TA) messages. In the interest of space, we limit our discussion to the most interesting results and observations. More results are available in our technical report [45]. Further, note that the delete-driven compaction strategies, TSD and TSA, fall back to LO+1 in absence of deletes, and thus, are omitted from the experiments without deletes.

5.1 Performance Implications

We first analyze the implications of compactions on the ingestion, lookup, and overall performance of an LSM-engine.

5.1.1 Data loading. In this experiment, we insert 10M key-value entries uniformly generated into an empty database to quantify the raw ingestion performance. The influence of the different compaction policies on ingestion is shown in Fig. 4(a)-(d).

O1: High Data Movement due to Compactions. Fig. 4(a) shows that the overall (read and write) data movement due to compactions is significantly larger than the actual size of the data ingested. We observe that among the leveled LSM- designs, the Full moves $63\times$ ($32\times$ for reads and $31\times$ for writes) the data originally ingested. The data movement is significantly smaller for Tier, however, it remains $12\times$ of the data size. These observations confirm the remarks on write amplification in LSMs presented in the literature [39], but also highlight the problem of *read amplification due to compactions* leading to poor device bandwidth utilization.

TA I: Compactions with higher eagerness move more data. Eager compactions (leveling) compact $2.5 - 5.5\times$ more data than lazy (tiering).

O2: Partial Compaction Reduces Data Movement. We now shift our attention to the differences between the different variations of leveling. Contrary to Full, all other approaches do not compact entire levels but only a small number of overlapping files. Fig. 4(a) shows that leveled partial compaction leads to 34% - 56% less data movement than Full. The reason is twofold: (1) A file with no overlap with its parent level, is only logically merged. Such *pseudo-compactions* require simple metadata (file pointer) manipulation in memory, and no I/Os. (2) Small compaction granularity allows for choosing a file with (i) the least overlap, (ii) the highest number of updates, or (iii) the highest # tombstones for compaction. Such data movement policies reduce overall data movement. Specifically, LO+1 (LO+2) are designed to pick for compaction files with least overlap with the parent $i+1$ level (grandparent $i+2$ level). They move 10%-23% less data than other partial compaction strategies.

TA II: Smaller compaction granularity reduces data movement. Partial compactions move $\sim 42\%$ less data than full-level compactions.

O3: The Compaction Count is Higher for Partial Compaction Routines. Fig. 4(b) shows that partial compactions initiate $4\times$ more compaction jobs than Full, as many as the number of levels in the tree. Note that for a steady-state LSM-Tree with partial compactions every memory buffer flush triggers cascading compactions to all L levels, while in a full-level compaction system this happens only when a level is full (every T compactions). Finally, since both Tier and Full are full-level compactions the compaction count is similar.

TA III: Larger compaction granularity leads to fewer but larger compactions. Full-level compactions perform about $1/L$ times fewer compactions than partial compaction routines, however, full-level leveling moves nearly $2L$ more data per compaction.

O4: Full Leveling has the Highest Mean Compaction Latency. As expected, Full compactions have the highest average latency ($1.2 - 1.9\times$ higher than partial leveling, and $2.1\times$ than tiering). Full can neither take advantage of pseudo-compactions nor optimize the data movement during compactions, hence, each compaction job moves a large amount of data. Fig. 4(c) shows the mean compaction latency for all strategies as well as the median (P50), the 90th percentile (P90), the 99th percentile (P99), and the maximum (P100). We observe that the tail latency (P90, P99, P100) is more predictable for full leveling while partial compactions and especially tiering have high variability due to the different data movement policies and the key-overlap depending on the data distribution.

The compaction latency presented in Fig. 4(c) can be broken to IO time and CPU time. We observe that the CPU effort is about 50% regardless of the compaction strategy. During a compaction, CPU cycles are spent in (1) obtaining locks and taking snapshots, (2) merging the entries, (3) updating file pointers and metadata, and (4) synchronizing output files post compaction. Among these, the time spent to sort-merge the data in memory dominates.

O5: The Tail Write Stalls is Highest for Tiering. Fig. 4(d) shows the tail write latency – the worst case latency for a write operation – is highest for tiering. Specifically, the tail write latency for tiering

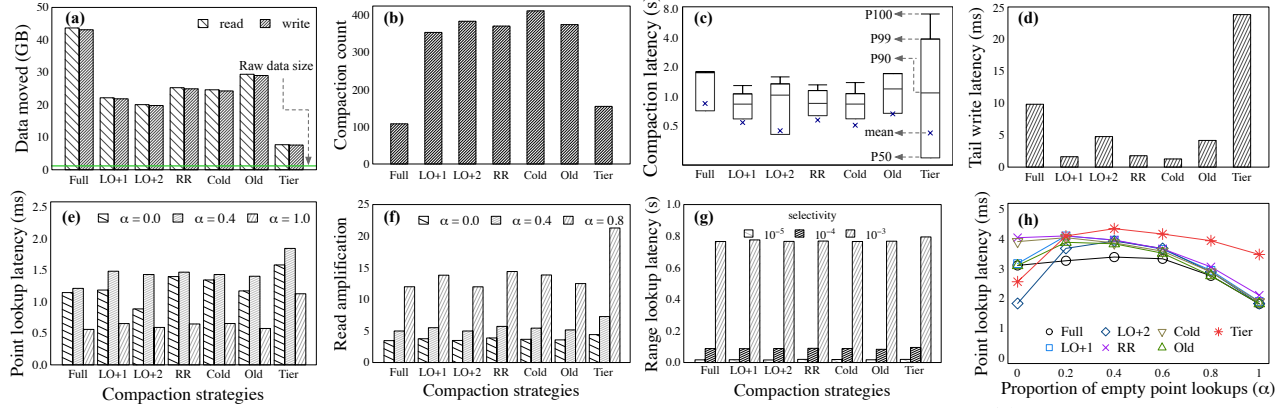


Fig. 4: Compactions influence the ingestion performance of LSM engines heavily in terms of (a) the overall data movement, (b) the compaction count, (c) the compaction latency, and (d) the tail latency for writes, as well as (e, f) the point lookup performance. (g) The range scan performance remains agnostic of compactions as the amount of data read by a scan remains the same. (h) shows the variation in the point lookup latency as the proportion of empty queries (α) increases in a workload.

is about $2.5\times$ greater than Full and $5 - 12\times$ greater than partial compaction strategies. Tiering in RocksDB [42] is optimized for writes and opportunistically seeks to compact all data to a large single level as a sorted run. This design achieves lower average write latency at the expense of prolonged write stalls in the worst case, which is when overlap between two consecutive levels is very high. In addition, Full has $2 - 5\times$ higher tail write stalls than partial compactions because when multiple consecutive levels are close to saturation, a buffer flush can result in a cascade of compactions.

TA IV: Tier and Full may cause prolonged write latency. Tail write stall for tiering is up to 25ms, while for partial leveling (Old) as low as 1.3ms.

5.1.2 Querying the Data. In this experiment, we perform 1M point lookups on the previously generated preloaded database (with 10M entries). The lookups are uniformly distributed in the domain and we vary the fraction of empty lookups α between 0 and 1. Specifically, $\alpha = 0$ indicates that we consider only non-empty lookups, while for $\alpha = 1$ we have lookups on non-existing keys. We also execute 1000 range queries, while varying their selectivity.

O6: The Point Lookup Latency is Highest for Tiering and Lowest for Full-Level Compaction. Fig. 4(e) shows that point lookups perform best for Full, and worst for tiering. The mean latency for point lookups with tiering is between $1.1 - 1.9\times$ higher than that with leveled compactions for lookups on existing keys, and $\sim 2.2\times$ higher for lookups on non-existing keys. Note that lookups on existing keys must always perform at least one I/O per lookup (unless they are cached), and therefore, takes longer than that on non-existing keys. When considering existing keys on a tree with a size ratio T , theoretically, the lookup cost for tiering should be $T\times$ higher than its leveling equivalent [18]. However, this worst-case is not always accurate as, in practice the cost depends on the (i) block cache size and the caching policy, (ii) temporality of the lookup keys, and (iii) the implementation of the compaction strategy. RocksDB tiering has overall fewer sorted runs than textbook tiering. Taking into account the block cache and temporality in the lookup workload, the observed tiering cost is less than $T\times$ the cost observed for Full. In addition, Full is 3%-15% lower than the partial compaction routines, because during normal operation of Full some levels might be entirely empty, while for partial compaction

all levels are always close to being full. Finally, we note that the various data movement policies used in the different partial leveling compaction strategies do not affect significantly point lookup latency, which always benefits from Bloom filters (10 bits-per-key) and the block cache, which is about 0.05% of the data size.

O7: The Point Lookup Latency Increases for Workloads with Comparable Number of Empty and Non-Empty Queries. A surprising result for point lookups that is also revealed in Fig. 4(e) is that they perform worse when the fraction of empty and non-empty lookups is balanced. Intuitively, one would expect that as we have more empty queries (that is, as α increases) the latency would decrease since the only data accesses needed by empty queries are the ones due to Bloom filter false positives [18]. To further investigate this result, we plot in Fig. 4(h) the 90th percentile (P90) latency which shows a similar curve for point lookup latency as we vary α . In our configuration each file uses 20 pages for its Bloom filters, 4 pages for its index blocks, and that the false positive is $FPR = 0.8\%$. A non-empty query ($\alpha = 0$) needs to load the Bloom filters of the levels it visits until it terminates. For all intermediate levels, it accesses the index and data blocks with FPR , and for the level that finally finds its result, it loads the Bloom filter, the index blocks, and the corresponding data. On the other hand, an empty query ($\alpha = 1$) accesses the Bloom filters of all levels before returning an empty result. Note that for each level it also accesses the index and data blocks with FPR . The counter-intuitive shape is a result of the non-empty lookups not needing to load the Bloom filters for all levels when $\alpha = 0$ and the empty lookups accessing index and data only when there is a false positive when $\alpha = 1$.

TA V: The Point Lookup Latency is Largely Unaffected by the Data Movement Policy. In presence of Bloom filters (with high enough memory) and small enough block cache, the point query latency remains largely unaffected by the data movement policy as long as the number of sorted runs in the tree remains the same. This is because block-wise caching of the filter and index blocks reduces the time spent performing disk I/Os significantly.

O8: Read Amplification is Influenced by the Block Cache Size and File Structure, and is Highest for Tiering. Fig. 4(f) shows that the read amplification across different compaction strategies for non-empty queries ($\alpha = 0$) is between 3.5 and 4.4. This is

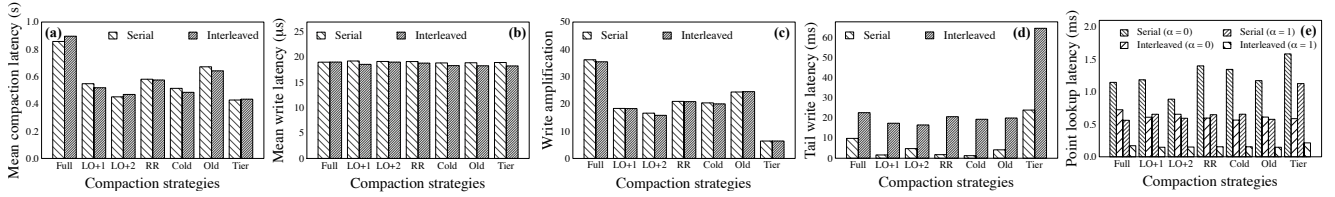


Fig. 5: (a-c) The average ingestion performance for workloads with interleaved inserts and queries is similar to that of an insert-only workload, but (d) with worse tail performance. However, (e) interleaved lookups are significantly faster.

attributed to the size of filter and index blocks which are $5\times$ and $1\times$ the size of a data block, respectively. Each non-empty point lookup fetches between 1 and L filter blocks depending on the position of the target key in the tree, and up to $L \cdot FPR$ index and data blocks. Further, the read amplification increases exponentially with α , reaching up to 14.4 for leveling and 21.3 for tiering (for $\alpha = 0.8$). Fig. 4(f) also shows that the estimated read amplification for point lookups is between $1.2\times$ and $1.8\times$ higher for Tier than for leveling strategies. This higher read amplification for Tier is owing to the larger number of sorted runs in the tree, and is in line with O7.

O9: The Effect of Compactions on Range Scans is Marginal.

To answer a range query the LSM-Tree instantiates multiple *run-iterators* that sequentially scan each sorted run that contains pages with matching data. Thus, the performance of a range query depends on (i) the scan time of the iterators (which is related to the query selectivity) and (ii) the time to merge the data streamed by each iterator. Regardless of the exact compaction strategy, the number of sorted runs in a leveled LSM-tree remains the same (one per level), which results in largely similar range query latency for all leveled compaction variations, especially for larger selectivity (Fig. 4(g)). Note that in the absence of updates and deletes, the overall amount of data qualifying for a range query is virtually identical for leveling and tiering despite the number of sorted runs being higher in tiering. The $\sim 5\%$ higher average range query latency for Tier is attributed to the additional I/Os needed to handle partially consumable disk pages from each run ($O(L \cdot T)$ in the worst case).

TA VI: In Absence of Updates/Deletes, the Range Query Latency is Unaffected by Compactions. For workloads with unique inserts, the total amount of data sort-merged by the iterators is largely similar regardless of the compaction strategy, leading to similar range query performance.

5.1.3 Executing mixed workloads. We now discuss performance implications when ingestion (thus compaction) and read queries are interleaved. In this experiment we interleave the ingestion of 10M unique key-value entries with 1M point lookups on existing keys. All read queries are performed after $L - 1$ levels are full.

O10: Tail Latency for Writes is Higher for Mixed Workloads.

Fig. 5(a) and (b) show that the mean latency of compactions that are interleaved with point queries is only marginally affected for all compaction strategies. This is also corroborated by the write amplification remaining unaffected by mixing reads and writes as shown in Fig. 5(c). On the other hand, Fig. 5(d) shows that the tail write latency is significantly increased between $2 - 15\times$ when read queries are interleaved with compactions. This increase is attributed to (1) the need of point queries to access filter and index blocks that requires disk I/Os that compete with writes and saturate the device, and (2) the delay of memory buffer flushing during lookups.

O11: Interleaving Compactions with Point Queries Keeps the Cache Warm.

Since in this experiment we start the point queries when $L - 1$ levels of the tree are full we expect that the interleaved read query execution will be faster than the serial one, by $1/L$ (25% in our configuration) which corresponds to the difference in the height of the trees. However, Fig. 5(e) shows this difference to be between 26% and 63% for non-empty queries and between 69% and 81% for empty queries. The reasons interleaved point query execution is faster than expected are that (1) about 10% of lookups terminate within the memory buffer, without requiring any disk I/Os, and (2) the block cache is pre-warmed with filter, index, and data blocks cached during compactions. Putting together the marginal impact of interleaved execution to compactions and the significant benefits for point queries, the total workload execution time is $1.6 - 2.1\times$ faster than serial execution as shown in Fig. 5(f).

TA VII: Compactions Help Lookups by Pre-Warming the Block Cache with Index and Filter Blocks for Mixed Workloads. As the file metadata including the index and filter blocks needs to be updated during compactions, the block cache is pre-warmed with the filter, index, and data blocks, which helps the subsequent point lookups.

5.2 Workload Influence

Next, we analyze the implications of the workload distribution and composition on compactions.

5.2.1 Varying the Ingestion Distribution. In this experiment we use a mixed workload that varies the ingestion distribution (uniform, Zipfian, normal) and has uniform lookup distribution.

O12: Ingestion Performance is Agnostic to Insert Distribution.

Fig. 4(a), 6(a), and 6(e) show that the total data movement during compactions remains virtually identical for (unique) insert-only workloads generated using uniform, Zipfian, and normal distributions, respectively. Further, we observe that the mean and tail compaction latencies are agnostic of the ingestion distribution (Fig. 4(c), 6(b), and 6(f) are almost identical as well). As long as the data distribution does not change over time, the entries inserted in each level follow the same distribution and the expected overlap between files of different levels remain the same. Therefore, for an ingestion-only workload the data distribution does not influence the choice of compaction strategy.

O13: Insert Distribution Influences Point Queries. Fig. 6(c) shows that while tiering has a slightly higher latency for point lookups, the relative performance of the compaction strategies is close to each other for any fraction of non-empty queries in the workload (all values of α). As $\alpha \rightarrow 1$, the filter and index block misses in the block cache diminishes sharply and approaches zero. Note that in this experiment the inserts are generated using Zipfian

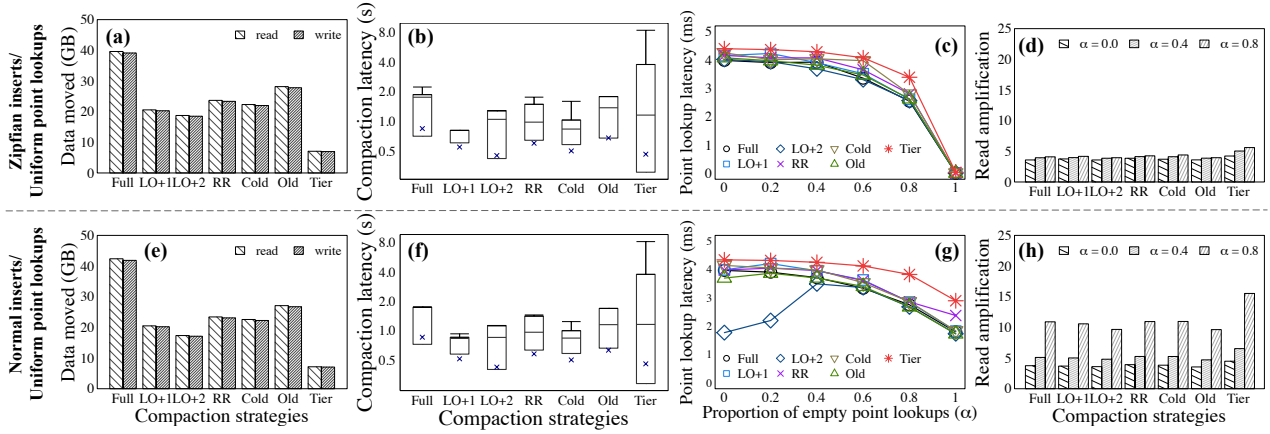


Fig. 6: As the ingestion distribution changes to (a-d) Zipfian and (e-h) normal, the ingestion performance of the database remains nearly identical with improvement in the lookup performance.

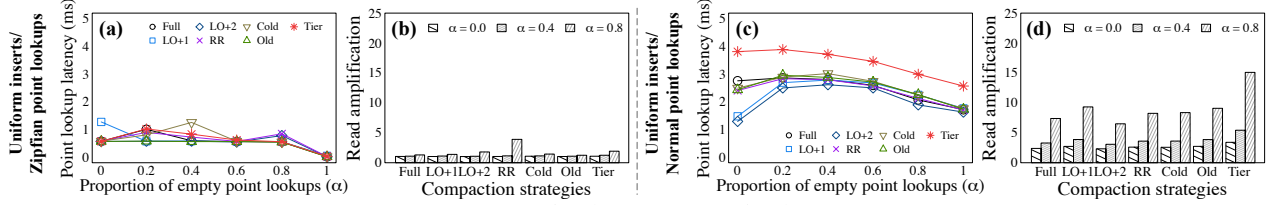


Fig. 7: A skewed lookup distribution, such as Zipfian (a, b) and normal (c, d) improves the lookup performance dramatically in presence of block cache and with assistance from the Bloom filters.

ingestion and target a small part of the domain. This allows a large fraction of the non-existing point lookups to terminate in memory by simply utilizing file metadata (i.e., the min/max keys contained in each file) without requiring to fetch any filter (and index) blocks. In Fig. 6(d), we observe that the read amplification remains comparable to that in Fig. 4(f) (uniform ingestion) for $\alpha = 0$ and even $\alpha = 0.4$. However, for $\alpha = 0.8$, the read amplification in Fig. 6(d) becomes 65%-75% smaller than in the case of uniform inserts. The I/Os performed to fetch the filter blocks is close to zero. *This shows that all compaction strategies perform equally well while executing an empty query-heavy workload on a database pre-populated with Zipfian inserts.* In contrast, when performing lookups on a database pre-loaded with normal ingestion, the point lookup performance (Fig. 6(g)) largely resembles its uniform equivalent (Fig. 4(h)), as the ingestion-skewness is comparable. The filter and index block hits are $\sim 10\%$ higher for the normal distribution compared to uniform for larger values of α , which explains the comparatively lower read amplification shown in Fig. 6(h). This plot also shows the first case of unpredictable behavior of LO+2 for $\alpha = 0$ and $\alpha = 0.2$. We observe more instances of such unpredictable behavior for LO+2, which probably explains why it is rarely used in new LSM stores.

5.2.2 Varying the Point Lookup Distribution. In this experiment, we change the point lookup distribution to **Zipfian** and **normal**, while keeping the ingestion distribution as uniform.

O14: Point Lookup Distribution Significantly Affects Performance. Zipfian point lookups on uniformly populated data leads to low latency point queries for all compaction strategies, as shown in Fig. 7(a) because the block cache is enough for the popular blocks in all cases, as also shown by the low read amplification in Fig. 7(b). On the other hand, when queries follow the normal distribution

partial compaction strategies LO+1 and LO+2 dominate all other approaches, while Tier is found to be perform significantly slower than all other approaches, as shown in Fig. 7(c) and 7(d).

TA VIII: For Skewed Ingestion/Lookups, all Compaction Strategies Behave Similarly in terms of Lookup Performance. While the ingestion distribution does not influence its performance, heavily skewed ingestion or lookups impacts query performance due to block cache and file metadata.

5.2.3 Varying the Proportion of Updates. We now vary the update-to-insert ratio, while interleaving queries with ingestion. An update-to-insert ratio 0 means that all inserts are unique, while a ratio 8 means that each unique insert receives 8 updates on average.

O15: For Higher Update Ratio Compaction Latency for Tiering Drops; LO+2 Dominates the Leveling Strategies. As the fraction of updates increases the compaction latency decreases significantly for tiering because we discard multiple updated entries in every compaction (Fig. 8(a)). We observe similar but less pronounced trends for Full and LO+2, while the remaining leveling strategies remain largely unchanged. *Overall, larger compaction granularity helps to exploit the presence of updates by invalidating more entries at a time.* Among the leveling strategies, LO+2 performs best as it moves $\sim 20\%$ less data during compactions, which also affects write amplification as shown in Fig. 8(b).

O16: As the Fraction of Updates Increase, Compactions Show More Stable Performance. We previously discussed that Tier has the highest tail compaction latency. As the fraction of updates increase all compaction strategies including Tier have lower tail compaction latency Fig. 8(c) shows that Tier's tail compaction latency drops from $6\times$ higher than Full to $1.2\times$ for update-to-insert ratio 8, which demonstrates that Tier is most suitable for

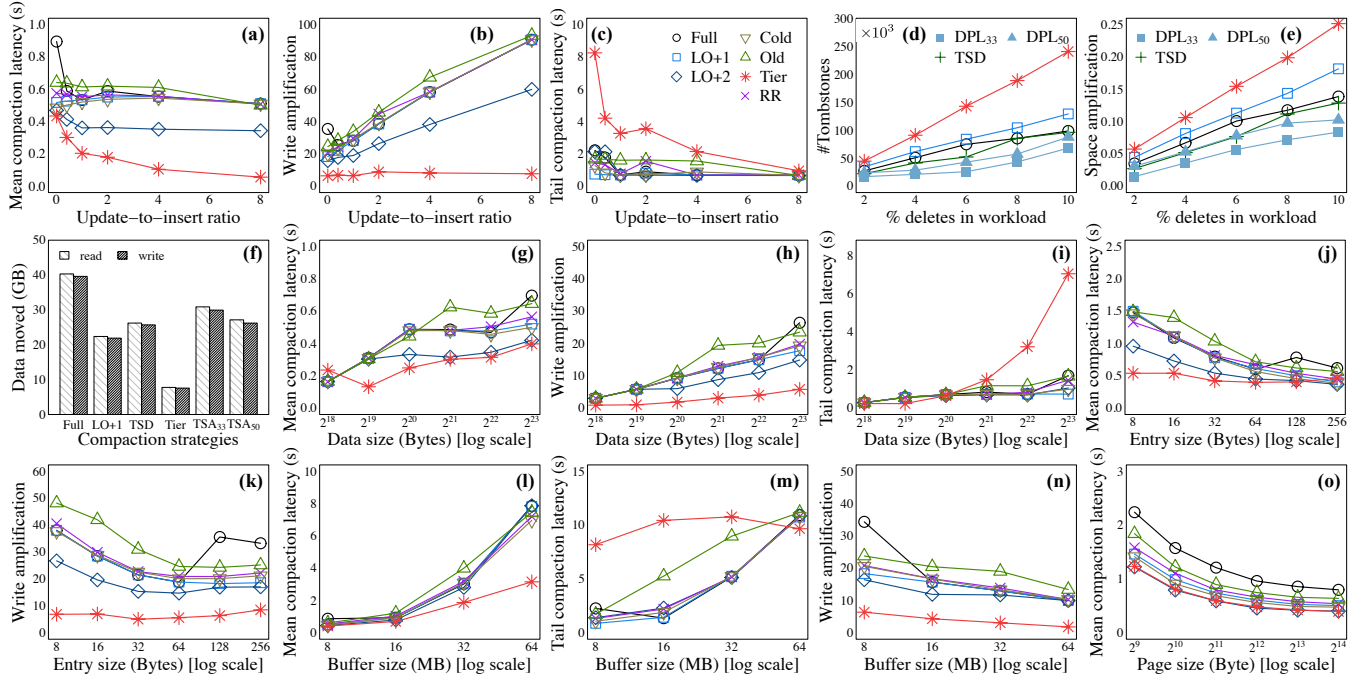


Fig. 8: Varying in the composition of the workload (a-k) and the LSM tuning (l-o) demonstrates that there is no perfect compaction strategy – choosing the appropriate compaction strategies is subject to the workload and the performance goal.

update-heavy workloads. We also observe that lookup latency and read amplification also decrease for update-heavy workloads.

TA IX: Tiering Dominates the Performance Space for Update-Intensive Workloads. When subject to update-intensive workloads, Tier exhibits superior compaction performance on average along with comparable lookup performance (as leveled LSMs), which allows Tier to dominate the overall performance space.

5.2.4 Varying Delete Proportion. We now analyze the impact of deletes, which manifest as out-of-place invalidations with special entries called tombstones [44]. We keep the same data size and vary the proportion of point deletes in the workload. All deletes are issued on existing keys and are interleaved with the inserts.

O17: TSD and TSA Offer Superior Delete Performance. We quantify the efficacy of deletion using the number of tombstones at the end of the workload execution. The lower this number the faster deleted data have been purged from the database, which in turn reduces space, write, and read amplification. Fig. 8(d) shows that TSD and TSA, the partial compaction strategies that are optimized for deletes, maintain the fewer tombstones at the end of the experiment. For a workload with 10% deletes, TSD purges 16% more tombstones than Tier and 5% more tombstones by picking the files that have a tombstone density above a pre-set threshold for compaction. For TSA, we experimented with two different thresholds for delete persistence: TSA₃₃ and TSA₅₀ is set to 25% and 50% of the experiment run-time, respectively. As TSA guarantees persistent deletes within the thresholds set, it compacts more data aggressively, and ends up with 7-10% fewer tombstones as compared to TSD. Full manages to purge more tombstones than any partial compaction routine, as it periodically compacts entire levels. Tier retains the highest number of tombstones as it maintains the highest number of sorted

runs overall. As the proportion of deletes in the workload increases, the number of tombstones remaining the LSM-tree (after the experiment is over) increases. TSA and TSD along with Full scale better than the partial compaction routines and tiering. By compacting more tombstones, TSA and TSD also purge a larger amount of invalid data reducing space amplification, as shown in Fig. 8(e).

O18: Optimizing for Deletes comes at a (Write) Cost. The reduced space amplification offered by TSA and TSD is achieved by compacting the tombstones eagerly, which increases the overall amount of data moved due to compaction. Fig. 8(f) shows that TSD and TSA₅₀ compacts 18% more data than the write optimized LO+1 (for TSA₃₃ this becomes 35%). Thus, TSD and TSA are useful when the objective is to (i) persist deletes timely or (ii) reduce space amplification caused by deletes.

TA X: TSD and TSA are Tailored for Deletes. In presence of deletes, TSA and TSD purges significantly more tombstones eagerly reducing space amplification. TSA ensures timely persistent deletion; however, for smaller persistence threshold, more data needs to be compacted to meet the threshold.

5.2.5 Varying the Ingestion Count. In this set of experiments we vary the ingestion count to report scalability.

O19: Tiering Scales Better than Leveling in terms of Compaction Latency But has Very High Tail Latency. The mean compaction latency scales sub-linearly for all compaction strategies, as shown in Fig. 8(g). The relative advantages of the different compaction strategies, however, remain similar, which shows that data size is not a determining factor when selecting the appropriate compaction strategy. This is further backed up by Fig. 8(h) which shows how write amplification scales. However, for latency-sensitive applications, as data size grows, the worst-case overlap of files in consecutive levels increases significantly for Tier (Fig. 8(i)).

5.2.6 Varying Entry Size. Here, we keep the key size constant (4B) and vary the value field from 4B to 254B to vary the entry size.

O20: For Smaller Entry Size, Leveling Compactions Are More Expensive. Smaller entry size increases the number of entries per page, which, in turn, leads to (i) more keys to compare during merge, and (ii) bigger Bloom filters that require more space and more CPU time per file. Fig. 8(j) shows these trends. We also observe similar trends for write amplification in Fig. 8(k) and for query latency. They both decrease as the entry size increases.

5.3 LSM Tuning Influence

In the final part of our analysis, we discuss the interplay of compactions with the standard LSM tunings knobs. We vary knobs like memory buffer size, page size, and size ratio, and analyze their impact on the compaction strategies employed. In the interest of space, we only include the most important results, while the rest can be found in an extended version of this paper [45].

O21: Tiering Offers Enhanced Ingestion Performance as the Buffer Size Increases. Fig. 8(l) shows that as the buffer size increases, the mean compaction latency increases across all compaction strategies. The size of buffer dictates the size of the files on disk; larger file size leads to more data being moved per compaction. Also, for larger file size, the Bloom filter size per file also increases along with the time spent for hashing, which increases compaction latency. Further, as the buffer size increases, Tier mean compaction latency scales better than the other strategies. Fig. 8(m) shows that the high tail compaction latency for Tier plateaus quickly as the buffer size is increased, and eventually crossovers with that for the eagerer compaction strategies when the buffer size becomes 64MB.

We also observe in Fig. 8(n) that among the partial compaction routines Old experiences an increased write amplification through-out, while L0+1 and L0+2 consistently offer lower write amplification, while guaranteeing predictable ingestion performance.

O22: All Compaction Strategies React Similarly to Page Size Change. In this experiment, we only vary the logical page size, which in turn, changes the number of entries per page. Therefore, the smaller the page size, the larger is the number of pages per file – which means more I/Os are required to access a file on the disk. For example, when the page size shrinks from 2^{10} to 2^9 , the number of pages doubles. With smaller page size, the index block size per file increases as more pages should be indexed, which also contributes to the increasing I/Os. Thus, an increase in the logical page size, reduces the mean compaction latency, as shown in Fig. 8(o).

O23: Miscellaneous Observations. We also vary LSM tuning parameters such as the size ratio, the memory allocated to Bloom filters, and the size of the block cache. We observe that these changing the values of these knobs affect the different compaction strategies similarly, and hence, does not influence the choice of the appropriate compaction strategy for any particular set up. In the interest of space we do not include the detailed results here, but the interested reader can find them in the Appendix of our technical report [45].

6 DISCUSSION

The design space detailed in Section 3 and the experimental analysis presented in Section 5 aim to offer to database researchers and practitioners the necessary insights to make educated decisions when selecting compaction strategies for LSM-based data stores.

Know Your LSM Compaction. LSM-trees are considered “write-optimized”, however, in practice their performance strongly depends on *when and how compactions are performed*. We depart from the notion of *treating compactions as a black-box*, and instead, we formalize *LSM compactions* as an ensemble of four fundamental compaction primitives. This allows us to reason about each of these primitives and navigate the *LSM compaction design space* in search of the appropriate compaction strategy for a workload or for custom performance goals. Further, the proposed compaction design space provides the necessary intuition about how simple modifications to an existing engine (like data movement policy or compaction granularity) can be key to achieving significant performance improvement or cost benefits.

Avoiding the Worst Choices. We discuss how to avoid common pitfalls. For example, tiering is often considered as the write-optimized variant, however, we show that it comes with high tail latency, making it unsuitable for applications that need worst-case performance guarantees. On the other hand, partial compactions with leveling, and especially hybrid leveling (e.g., 1-leveling) offer the most stable performance.

Adapting with Workloads. In prior work tiering is used for write-intensive use-cases, while leveling offers better read performance. However, in practice, in mixed HTAP-style workloads, lookups have a strong temporal locality, and are essentially performed on recent hot data. In such cases, the block cache is frequently proved to be enough for holding the working set and eliminate the need for other costly optimizations for read queries.

Exploring New Compaction Strategies. Ultimately, this work lays the groundwork for exploring the vast design space of LSM compactions. A key intuition we developed during this analysis is that contrary to existing designs, LSM-based systems can benefit by employing different compaction primitives in different levels, depending on the exact workload and the performance goals. The compaction policies we experimented with already support a wide range of metrics they optimize for including system throughput, worst-case latency, read, space, and write amplification, and delete efficiency. Using the proposed design space, new compaction strategies can be designed with new or combined optimization goals. We also envision systems that automatically select compaction strategies on the fly depending on the current context and workload.

7 CONCLUSION

LSM-based engines offer efficient ingestion and competitive read performance, while being able to manage various optimization goals like write and space amplification. A key internal operation that is at the heart of how LSM-trees work is the process of *compaction* that periodically re-organizes the data on disk.

We present the *LSM compaction design space* that uses four primitives to define compactions: (i) compaction trigger, (ii) compaction eagerness, (iii) compaction granularity, and (iv) data movement policy. We map existing approaches in this design space and we select several representative policies to study and analyze their impact on performance and other metrics including write/space amplification and delete latency. We present an extensive collection of observations, and we lay the groundwork for LSM systems that can more flexibly navigate the design space for compactions.

REFERENCES

- [1] W. Y. Alkowiileet, S. Alsubaiee, and M. J. Carey. An LSM-based Tuple Compaction Framework for Apache AsterixDB. *Proceedings of the VLDB Endowment*, 13(9):1388–1400, 2020.
- [2] S. Alsubaiee, Y. Altowim, H. Altwaijry, A. Behm, V. R. Borkar, Y. Bu, M. J. Carey, I. Cetindil, M. Cheelangi, K. Faraaz, E. Gabriellova, R. Grover, Z. Heilbron, Y.-S. Kim, C. Li, G. Li, J. M. Ok, N. Onose, P. Pirzadeh, V. J. Tsotras, R. Vernica, J. Wen, and T. Westmann. AsterixDB: A Scalable, Open Source BDMS. *Proceedings of the VLDB Endowment*, 7(14):1905–1916, 2014.
- [3] Amazon. EBS volume types. <https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ebs-volume-types.html>.
- [4] Amazon. EC2 Instance Types. <https://aws.amazon.com/ec2/instance-types/>.
- [5] Apache. Accumulo. <https://accumulo.apache.org/>.
- [6] Apache. Cassandra. <http://cassandra.apache.org>.
- [7] Apache. HBase. <http://hbase.apache.org/>.
- [8] M. Athanassoulis and S. Idreos. Design Tradeoffs of Data Access Methods. In *Proceedings of the ACM SIGMOD International Conference on Management of Data, Tutorial*, 2016.
- [9] M. Athanassoulis, M. S. Kester, L. M. Maas, R. Stoica, S. Idreos, A. Ailamaki, and M. Callaghan. Designing Access Methods: The RUM Conjecture. In *Proceedings of the International Conference on Extending Database Technology (EDBT)*, pages 461–466, 2016.
- [10] O. Balmau, F. Dinu, W. Zwaenepoel, K. Gupta, R. Chandhiramoorthi, and D. Di-dona. SILK: Preventing Latency Spikes in Log-Structured Merge Key-Value Stores. In *Proceedings of the USENIX Annual Technical Conference (ATC)*, pages 753–766, 2019.
- [11] O. Balmau, F. Dinu, W. Zwaenepoel, K. Gupta, R. Chandhiramoorthi, and D. Di-dona. SILK+ Preventing Latency Spikes in Log-Structured Merge Key-Value Stores Running Heterogeneous Workloads. *ACM Trans. Comput. Syst.*, 36(4):12:1–12:27, 2020.
- [12] M. Callaghan. Compaction priority in RocksDB. <http://smallldatum.blogspot.com/2016/02/compaction-priority-in-rocksdb.html>, 2016.
- [13] M. Callaghan. Compaction stalls: something to make better in RocksDB. <http://smallldatum.blogspot.com/2017/01/compaction-stalls-something-to-make.html>, 2017.
- [14] M. Callaghan. The Insert Benchmark. <http://smallldatum.blogspot.com/2017/06/the-insert-benchmark.html>, 2017.
- [15] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber. Bigtable: A Distributed Storage System for Structured Data. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 205–218, 2006.
- [16] CockroachDB. Open Issue: Storage: Performance degradation caused by kv tombstones. <https://github.com/cockroachdb/cockroach/issues/17229>, 2017.
- [17] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, and R. Sears. Benchmarking cloud serving systems with YCSB. In *Proceedings of the ACM Symposium on Cloud Computing (SoCC)*, pages 143–154, 2010.
- [18] N. Dayan, M. Athanassoulis, and S. Idreos. Monkey: Optimal Navigable Key-Value Store. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 79–94, 2017.
- [19] N. Dayan, M. Athanassoulis, and S. Idreos. Optimal Bloom Filters and Adaptive Merging for LSM-Trees. *ACM Transactions on Database Systems (TODS)*, 43(4):16:1–16:48, 2018.
- [20] N. Dayan and S. Idreos. 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*, pages 505–520, 2018.
- [21] N. Dayan and S. Idreos. The Log-Structured Merge-Bush & the Wacky Continuum. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 449–466, 2019.
- [22] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Vossell, and W. Vogels. Dynamo: Amazon’s Highly Available Key-value Store. *ACM SIGOPS Operating Systems Review*, 41(6):205–220, 2007.
- [23] S. Dong. Option of Compaction Priority. https://rocksdb.org/blog/2016/01/29/compaction_pri.html, 2016.
- [24] S. Dong, M. Callaghan, L. Galanis, D. Borthakur, L. Savor, and M. Strum. Optimizing Space Amplification in RocksDB. In *Proceedings of the Biennial Conference on Innovative Data Systems Research (CIDR)*, 2017.
- [25] Facebook. RocksDB. <https://github.com/facebook/rocksdb>.
- [26] G. Golan-Gueta, E. Bortnikov, E. Hillel, and I. Keidar. Scaling Concurrent Log-Structured Data Stores. In *Proceedings of the ACM European Conference on Computer Systems (EuroSys)*, pages 32:1–32:14, 2015.
- [27] Google. LevelDB. <https://github.com/google/leveldb/>.
- [28] HBase. Online reference. <http://hbase.apache.org/>, 2013.
- [29] G. Huang, X. Cheng, J. Wang, Y. Wang, D. He, T. Zhang, F. Li, S. Wang, W. Cao, and Q. Li. 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*, pages 651–665, 2019.
- [30] HyperLevelDB. Online reference. <https://github.com/rescrv/HyperLevelDB>.
- [31] S. Idreos and M. Callaghan. Key-Value Storage Engines. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 2667–2672, 2020.
- [32] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuums and the Path Toward Self-Designing Key-Value Stores that Know and Learn. In *Proceedings of the Biennial Conference on Innovative Data Systems Research (CIDR)*, 2019.
- [33] A. Kryczka. Compaction Styles. <https://github.com/facebook/rocksdb/blob/gh-pages-old/talks/2020-07-17-Brownbag-Compactions.pdf>, 2020.
- [34] LinkedIn. Voldemort. <http://www.project-voldemort.com>.
- [35] C. Luo and M. J. Carey. LSM-based Storage Techniques: A Survey. *The VLDB Journal*, 29(1):393–418, 2020.
- [36] F. Mei, Q. Cao, H. Jiang, and J. Li. SifrDB: A Unified Solution for Write-Optimized Key-Value Stores in Large Datacenter. In *Proceedings of the ACM Symposium on Cloud Computing, SoCC 2018, Carlsbad, CA, USA, October 11-13, 2018*, pages 477–489, 2018.
- [37] P. E. O’Neil, E. Cheng, D. Gawlick, and E. J. O’Neil. The log-structured merge-tree (LSM-tree). *Acta Informatica*, 33(4):351–385, 1996.
- [38] A. Pantelopoulos and N. G. Bourbakis. Prognosis: a wearable health-monitoring system for people at risk: methodology and modeling. *IEEE Trans. Information Technology in Biomedicine*, 14(3):613–621, 2010.
- [39] P. Raju, R. Kadekodi, V. Chidambaram, and I. Abraham. PebblesDB: Building Key-Value Stores using Fragmented Log-Structured Merge Trees. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP)*, pages 497–514, 2017.
- [40] K. Ren, Q. Zheng, J. Arulraj, and G. Gibson. SlimDB: A Space-Efficient Key-Value Storage Engine For Semi-Sorted Data. *Proceedings of the VLDB Endowment*, 10(13):2037–2048, 2017.
- [41] RocksDB. Leveled Compaction. <https://github.com/facebook/rocksdb/wiki/Leveled-Compaction>, 2020.
- [42] RocksDB. Universal Compaction. <https://github.com/facebook/rocksdb/wiki/Universal-Compaction>, 2020.
- [43] S. Sarkar, J.-P. Banâtre, L. Rilling, and C. Morin. Towards Enforcement of the EU GDPR: Enabling Data Erasure. In *Proceedings of the IEEE International Conference of Internet of Things (iThings)*, pages 1–8, 2018.
- [44] S. Sarkar, T. I. Papon, D. Staratzis, and M. Athanassoulis. Lethe: A Tunable Delete-Aware LSM Engine. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 893–908, 2020.
- [45] S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. Constructing and Analyzing the LSM Compaction Design Space (Technical Report). *BU DiSC Technical Report*, 2021.
- [46] S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. K-V Workload Generator. <https://github.com/BU-DiSC/K-V-Workload-Generator>, 2021.
- [47] S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. LSM Compaction Analysis. <https://github.com/BU-DiSC/LSM-Compaction-Analysis>, 2021.
- [48] ScyllaDB. Online reference. <https://www.scylladb.com/>.
- [49] R. Sears and R. Ramakrishnan. bLSM: A General Purpose Log Structured Merge Tree. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 217–228, 2012.
- [50] Tarantool. Online reference. <https://www.tarantool.io/>.
- [51] D. Teng, L. Guo, R. Lee, F. Chen, S. Ma, Y. Zhang, and X. Zhang. LSbM-tree: Re-Enabling Buffer Caching in Data Management for Mixed Reads and Writes. In *37th IEEE International Conference on Distributed Computing Systems, ICDCS 2017, Atlanta, GA, USA, June 5-8, 2017*, pages 68–79, 2017.
- [52] D. Teng, L. Guo, R. Lee, F. Chen, Y. Zhang, S. Ma, and X. Zhang. A Low-cost Disk Solution Enabling LSM-tree to Achieve High Performance for Mixed Read/Write Workloads. *ACM Trans. Storage*, 14(2):15:1–15:26, 2018.
- [53] WiredTiger. Source Code. <https://github.com/wiredtiger/wiredtiger>.
- [54] L. Yang, H. Wu, T. Zhang, X. Cheng, F. Li, L. Zou, Y. Wang, R. Chen, J. Wang, and G. Huang. Leaper: A Learned Prefetcher for Cache Invalidation in LSM-tree based Storage Engines. *Proceedings of the VLDB Endowment*, 13(11):1976–1989, 2020.

A SUPPLEMENTARY EXPERIMENTS

A.1 Varying Block Cache

We vary the block cache for insert-only and mixed workloads (10% existing point lookups interleaved with insertions). For mixed workload, the mean compaction latency remains stable when block cache varies from 8 to 256 MB. However, for insert-only workload, the mean compaction latency increases sharply when block cache is more than 32 MB (Fig. 9a and 10a). We observe that for insert-only

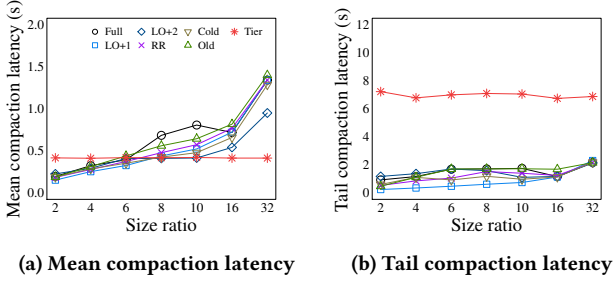


Fig. 12: Varying Size Ratio (insert-only)

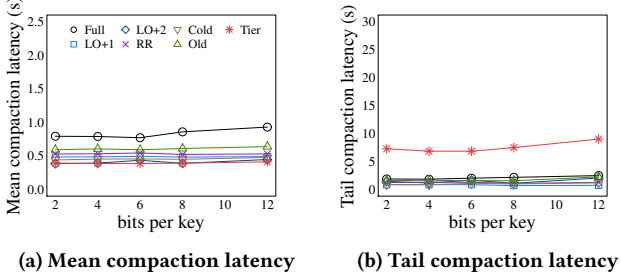


Fig. 13: Varying Size Ratio (insert-only)

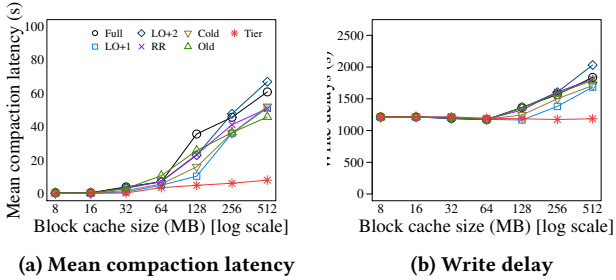


Fig. 9: Varying Block Cache (insert-only)

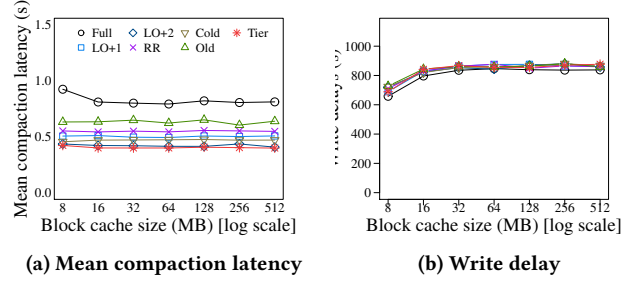


Fig. 10: Varying Block Cache (interleaving with 10% point lookups)

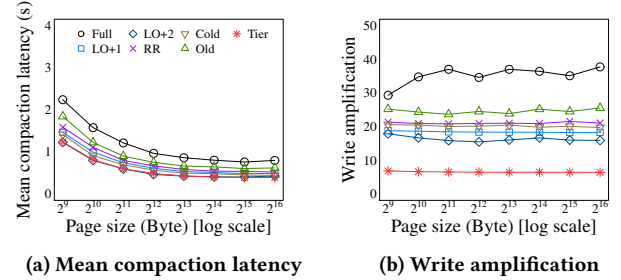


Fig. 11: Varying Page Size (insert-only)

workload, the write delay (also termed as write stall) is more than twice that of mixed workload (Fig. 9b and 10b). We leave this interesting phenomenon for future discussion. Compared to full and partial compaction, tiering is more stable with respect to different block cache size.

A.2 Varying Page Size

When we vary the page size, we observe almost consist patterns across different compaction strategies for all metrics (Fig. 11a and 11b). It turns out that compaction strategy does not play a big role for different page sizes.

A.3 Varying Size Ratio

We also compare the performance for different size ratio. According to Fig. 12a, tiering has higher mean compaction latency compared to other strategies when the size ratio is no more than 6 and after 6, full compaction and oldest compaction become the top-2 time-consuming strategies. In terms of tail compaction strategy in Fig. 12b, tiering is still the worst one compared to other strategies.

A.4 Varying Bits Per Key (BPK)

We also conduct the experiment to investigate the BPK's influence over compaction. From Fig. 12a and 13b, the mean and tail compaction latency may increase a little bit with increasing bits per key since larger filter blocks should be written but this increasing is very tiny since the increasing filter block is quite smaller than all data blocks. At the same, we also observe that the query latency even increases with increasing BPK (see Fig. 13c and 13d). This might come from higher filter block misses (Fig. YY) and this pattern becomes more obvious for existing queries in which case, accessing filter blocks is completely an extra burden.