



저작자표시-비영리-변경금지 2.0 대한민국

이용자는 아래의 조건을 따르는 경우에 한하여 자유롭게

- 이 저작물을 복제, 배포, 전송, 전시, 공연 및 방송할 수 있습니다.

다음과 같은 조건을 따라야 합니다:



저작자표시. 귀하는 원저작자를 표시하여야 합니다.



비영리. 귀하는 이 저작물을 영리 목적으로 이용할 수 없습니다.



변경금지. 귀하는 이 저작물을 개작, 변형 또는 가공할 수 없습니다.

- 귀하는, 이 저작물의 재이용이나 배포의 경우, 이 저작물에 적용된 이용허락조건을 명확하게 나타내어야 합니다.
- 저작권자로부터 별도의 허가를 받으면 이러한 조건들은 적용되지 않습니다.

저작권법에 따른 이용자의 권리는 위의 내용에 의하여 영향을 받지 않습니다.

이것은 [이용허락규약\(Legal Code\)](#)을 이해하기 쉽게 요약한 것입니다.

[Disclaimer](#)

M.S. THESIS

Validity Tracking Based Log Management for
In-Memory Databases

유효성 추적을 통한 인 메모리 데이터 베이스 로그 관리

August 2021

DEPARTMENT OF COMPUTER SCIENCE AND
ENGINEERING
COLLEGE OF ENGINEERING
SEOUL NATIONAL UNIVERSITY

Kwangjin Lee

Validity Tracking Based Log Management for
In-Memory Databases

유효성 추적을 통한 인 메모리 데이터 베이스 로그 관리

지도교수 염 현 영

이 논문을 공학석사학위논문으로 제출함

2021 년 6 월

서울대학교 대학원

컴퓨터공학부

이 광 진

이광진의 공학석사 학위논문을 인준함

2021 년 8 월

위 원 장	_____	전 병 곤
부위원장	_____	염 현 영
위 원	_____	이 영 기

Abstract

With in-memory databases (IMDBs), where all data sets reside in main memory for fast processing speed, logging and checkpointing are essential for achieving persistence in data. Logging of IMDBs has evolved to reduce run-time overhead to suit the systems, but this causes an increase in recovery time. Checkpointing technique compensates for these problems with logging, but existing schemes often incur high costs due to reduced system throughput, increased latency, and increased memory usage.

In this paper, we propose a checkpointing scheme using validity tracking-based compaction (VTC), the technique that tracks the validity of logs in a file and removes unnecessary logs. The proposed scheme shows extremely low memory usage compared to existing checkpointing schemes, which use consistent snapshots. Our experiments demonstrate that checkpoints using consistent snapshot increase memory footprint by up to two times in update-intensive workloads. In contrast, our proposed VTC only requires 2% additional memory for checkpointing. That means the system can use most of its memory to store data and process transactions.

Keywords: In-Memory Database, Persistence, Logging, Checkpointing, Snapshot

Student Number: 2019-23670

Contents

Abstract	1
1 Introduction	7
2 BACKGROUND AND MOTIVATION	12
2.1 Persistence in In-Memory Databases	12
2.2 Fork-Based Checkpointing	14
3 Design and Implementation	16
3.1 Design Overview	16
3.2 Distributed Logging and Log Data Format	18
3.3 Log File Compaction	19
3.4 Lazy Invalidation	24
3.5 Recovery	25
3.6 Correctness	27
3.7 Implementation	28
4 Evaluation	30
4.1 Experimental Setup	30
4.2 Performance	32

4.2.1	Throughput	32
4.2.2	Memory Footprint	33
4.2.3	Checkpointing Time	35
4.2.4	File Size	36
4.2.5	Restoring Time	37
5	Related Work	39
6	Conclusion	42
	초록	48

List of Figures

1.1	Redis memory footprint and throughput during checkpointing for the Yahoo! Cloud Serving Benchmark (YCSB) (50% update).	8
3.1	Overall procedure.	17
3.2	The log format.	18
3.3	Log file compaction.	20
3.4	Example of logging for insert and update.	21
3.5	The recovery order.	25
4.1	Throughput for varying record count (50% update proportion).	32
4.2	Throughput for varying update proportion (10M records).	32
4.3	Memory footprint comparison with existing schemes (50% update proportion).	34
4.4	Increased memory footprint (8M records).	34
4.5	Checkpointing time.	35
4.6	File size.	37

4.7 Restoring time	38
------------------------------	----

List of Tables

4.1	Parameters of YCSB workloads	31
-----	--	----

Chapter 1

Introduction

In-memory databases (IMDBs) are designed to achieve fast response time by processing data using the main memory, without accessing the disk. For this reason, IMDBs are widely adopted for various applications [1], such as e-commerce online transaction processing (OLTP) services, online games [2], finance [3], and more. The entire data residing in memory guarantees fast processing, but there is a risk of data loss due to system crashes, hardware failures, and power outages. To improve fault tolerance in long-running applications, IMDBs provide persistence through a variety of strategies. Checkpointing and logging are widely used techniques for the durability of IMDBs. Disk-based databases prefer ARIES-style [4] logging and checkpointing protocols, while most IMDBs record only redo logs excluding undo logs to reduce logging overhead and help performance. In addition, IMDBs need to checkpoint much more data than disk-based databases, so it is common to use an algorithm suitable for this, such as consistent snapshots [5], [6].

Many systems provide persistence by combining logging and checkpointing.

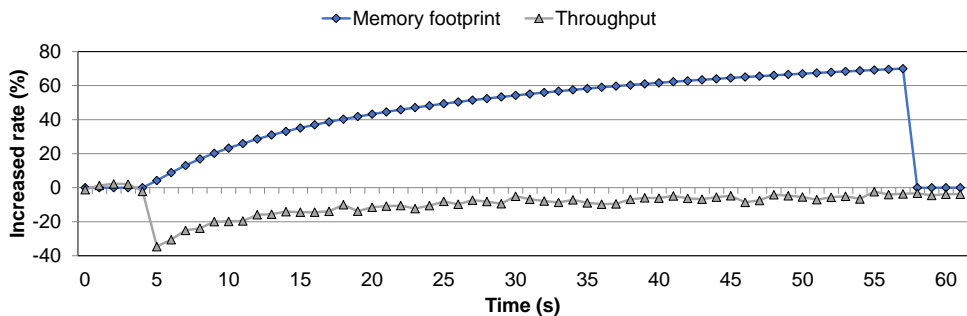


Figure 1.1 Redis memory footprint and throughput during checkpointing for the Yahoo! Cloud Serving Benchmark (YCSB) (50% update).

Systems that guarantee data durability only with periodic checkpointing can reduce run-time overhead, but the trade-off is that the system can lose a large amount of data due to system failure. Systems that use logging can lower the risk of data loss. Although logging increases the recovery interval and requires more storage space due to logs that accumulate over time, these problems can be alleviated by using checkpointing together. The combination of logging and checkpointing reduces recovery time by loading the latest checkpoint file and rerunning only subsequent logs. Furthermore, it allows for the space used by the logs to be reused.

Checkpointing plays an important role in effectively providing persistence for IMDBs, but it incurs significant costs for system throughput, latency, and peak memory usage. Figure. 1.1 shows the memory footprint and throughput of Redis, one of the most popular commercial IMDBs, during checkpointing. When using Redis, the memory footprint continues to grow over time until checkpointing is complete. As a result, memory usage at the end of checkpointing has increased up to 67%. Moreover, throughput decreases by 38% at the beginning of checkpointing due to frequent copy-on-write (CoW) operations.

Redis performs checkpointing by periodically taking consistent snapshots and storing them to stable storage. Data in a consistent snapshot should not be overwritten during checkpointing, so changes made to a database during client update requests are handled by CoW semantics. Since physical pages copied by CoW are not reclaimed until the checkpoint is completed, the memory footprint may increase up to two times during an update-intensive workload.

In a broad sense, a checkpoint is a technique that aims to keep the persistent state of a database up to date with the goal of reducing recovery time and reusing log space. Redis [7] and Hyper [8] simply take a consistent snapshot and store the contents of the snapshot through the `fork()` system call [9] and CoW semantics supported by the OS. However, this method has problems such as latency spikes [10] and increased memory usage. DMC [11] allows pages to be returned to the OS sooner before checkpointing is complete. This lowers the peak memory usage during checkpointing, but does not completely solve the increase in memory according to the update rate. Hekaton [12] and CALC [13] attempt to reduce the checkpointing overhead by using a partial checkpoint algorithm. This algorithm reduces cost by taking partial checkpoints that contain only some of the latest records. However, the process of merging partial checkpoints to create a complete checkpoint incurs another overhead.

In this paper, we propose a validity tracking-based log management scheme to provide improved durability and efficient checkpointing by minimizing the use of additional memory. By distributing and storing logs across multiple storage devices, we can hide the latency caused by the log buffer flush so that logging does not affect the throughput of the system. It also improves durability by reducing the flush cycle of the log buffer. Instead of generating checkpoints from the data on the main memory, our validity tracking-based compaction (VTC) scheme creates checkpoints by identifying valid logs in log files. VTC

uses only a small amount of extra memory because it does not require physical page duplication or multiple versioning in main memory.

We provide the proposed scheme as a simple user-level API. To test this system, we applied our scheme to Redis 5.0.6 and evaluate it with the Yahoo! Cloud Serving Benchmark (YCSB) [14] to compare its performance against the persistence schemes of Redis. The experimental results show that VTC consumes little memory to perform checkpoints and does not adversely affect system throughput and checkpoint time. With update-intensive workloads, VTC uses less than 2% of the size of the data set, while the memory footprint of the checkpointing scheme using consistent snapshots increases memory usage by up to 200%. This means that VTC permits most of an IMDB’s primary resource, system memory, to be used for data storage.

Our contributions can be summarized as follows.

- We analyze the persistence scheme for an existing in-memory database.
- We propose a persistence scheme that distributes and stores logs on multiple storage devices and removes unnecessary logs in the file by tracking log validity.
- We provide high-level APIs that can be easily applied to the existing IMDB with little modification.
- The experimental results show that our scheme offers slightly better throughput than the existing Redis logging scheme and maintains a more stable memory footprint than the existing Redis checkpointing scheme.

The rest of the paper is organized as follows. Section 2 describes the background and motivation. Section 3 introduces the design and implementation of

our proposed scheme. Section 4 evaluates VTC and other persistence schemes. Section 5 discusses the related works. Finally, Section 6 concludes this paper.

Chapter 2

BACKGROUND AND MOTIVATION

2.1 Persistence in In-Memory Databases

Because all of the data for IMDBs is in DRAM, which is a volatile memory, it is crucial to guarantee data durability with a fault-tolerance mechanism that will help prevent data loss in case of system failure [15]. IMDBs prefer data replication for fast failover and generally maintain replicas across multiple nodes to achieve high availability [12], [16], [17], [18]. However, catastrophic failures such as cluster-wide power outages can cause data loss if the data are not in stable storage. To avoid this issue, data must be kept in stable storage to ensure durability. The traditional techniques used for database durability are logging and checkpointing.

Most disk-based databases guarantee transaction durability with ARIES-style [4] logging. The ARIES protocol uses a write-ahead logging (WAL) scheme that sequentially records changes before modified pages are written to disk, and log records include redo and undo. Early IMDBs used similar techniques [19], [20].

However, logging to IMDBs is gradually optimized for high throughput and low latency, and the traditional logging scheme, which is relatively expensive compared to light transaction processing without disk access, is simplified for in-memory systems. In general, IMDBs reduce log volume by recording only the redo log and minimizing the log record’s information to mitigate the effects from log creation and log I/O overhead. However, replaying the log for recovery increases the recovery time. Additionally, WAL can recycle log space after applying all logs to the data file, whereas logging in IMDBs consumes more space over time. Therefore, periodic checkpointing is required to reduce the recovery time and recycle log space.

Checkpointing is also different in IMDBs. Because the entire data sets for IMDBs are kept in the main memory, it is common for their checkpoints to be larger than those of disk-based databases. Many IMDBs use a checkpointing algorithm that takes a consistent snapshot and stores it in stable storage. The wide application of consistent snapshots has led to extensive research in academia, and various algorithms [2], [21], [22] have been proposed. In fact, the commercial systems Redis and Hyper use the `fork()` system call as a consistent snapshot algorithm.

Redis, the most popular key-value IMDB [23], provides persistence via Redis data backup (RDB) and an append-only file (AOF). RDB is a feature that backs up the entire database in memory. It obtains a consistent snapshot using `fork()` and stores the contents of the snapshot in stable storage in the background through the child process created by `fork()`. AOF, a logging feature supported by Redis, appends a log of events that have changed the database to the log file. When a log file exceeds a specific size, the system acquires a snapshot, converts its contents to log format, and saves them as a file. As a result, it creates a new log file consisting of only the logs needed for recovery. This process mitigates

the increase in log space and recovery time.

2.2 Fork-Based Checkpointing

Fork-based checkpointing is a simple but efficient scheme that creates point-in-time consistent snapshots and store them in stable storage with OS supports. It has been demonstrated that the fork-based consistent snapshot algorithm outperforms other algorithms [2], [21], [22] for update intensive workloads [6]. In fact, many industrial IMDBs like Hyper [8] and Redis [7] employ the algorithm for checkpointing.

The `fork()` system call is used to create a child process by duplicating a process. Physical pages are not actually copied by `fork()`, and both processes refer to the same physical pages through virtual memory pages. If a page shared by both processes needs to be updated, the CoW technique copies the physical page to a new memory space and modifies it. Thanks to this CoW technique, point-in-time data on the child process (checkpointer) is not affected even as the parent process (worker) handles the client's update request. After `fork()`, the checkpointer traverses the snapshot and saves point-to-time data to a file.

However, there are well-known problems with fork-based checkpointing. The first problem is the latency spike due to the blocking operation, `fork()`, which occurs because IMDBs cannot process or respond to client requests while creating a process by `fork()`. Moreover, latency due to `fork()` increases as the size of the data set increases. The second problem is increased latency due to CoW. During the checkpointing period, the overhead due to CoW for processing update requests increases latency and affects throughput. If the update rate of the workload is high, CoW will frequently have a greater impact on the throughput. Finally, the increase in memory footprint by CoW is the most crucial problem.

The parent process handles client requests while the child process created by `fork()` writes the snapshot data to the file. If the request is an update, the physical page is copied by the CoW, which causes an increase in the memory footprint. Moreover, the memory footprint increases proportionally to the update rate of the workload. In the worst case of all pages being updated, the required memory size is twice that of the data set. Furthermore, the increased memory cannot be reclaimed until the child process is terminated. If there is no more available memory, either the transaction processing and checkpoint speed will be significantly slowed during the swap, or the out-of-memory killer will kill the processes. For this reason, many IMDB vendors [24] recommend that users take into consideration the memory increase due to `fork()` and set the swap to prevent out-of-memory problems.

We focused on a persistence scheme that combines logging and checkpointing, along with a checkpointing algorithm to minimize the memory use increase and provide stable throughput.

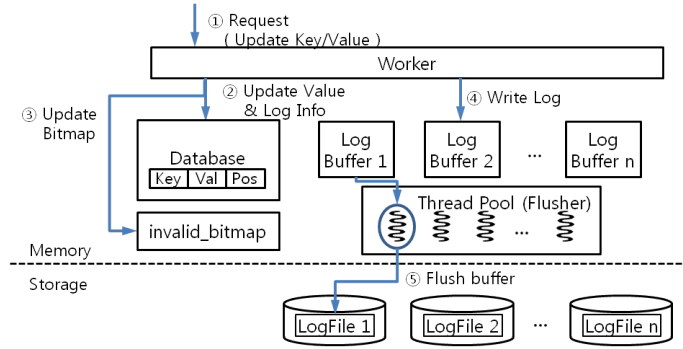
Chapter 3

Design and Implementation

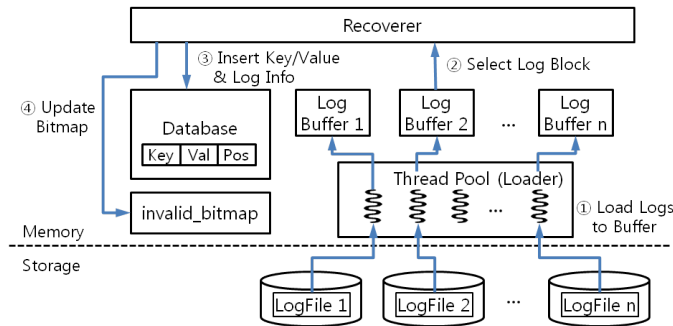
3.1 Design Overview

Our key idea is to distribute the logs into multiple files and reorganize them by identifying valid logs in the files. Figure. 3.1 shows the overall procedure of the proposed scheme. The worker creates logs of events that change the database and writes them to the log buffers. The flushers then flush the log buffer to the stable storage in the background. We hid the latency caused by flush by dividing the log into several SSDs and improved durability by reducing the flush interval.

The VTC scheme performs checkpointing based on logs stored in storage rather than on data in the main memory. VTC leaves only the logs needed for recovery through file-to-file copying, so to achieve this, we maintained an up-to-date log location for each database entry. An `invalid_bitmap` is allocated per log file, and each bit indicates the validity of each log in the file. The `invalid_bitmap` makes it simple to examine the validity of the log during check-



(a) Logging with multiple SSDs



(b) Recovery with multiple buffers

Figure 3.1 Overall procedure.

pointing. When the number of invalid logs in each log file reaches the threshold, checkpointing is triggered. Only one file can be checkpointed at a time, and logs are not stored in the file until checkpointing is complete. The separation of I/O between logging and checkpointing reduces checkpointing time and avoids log flush delays due to latency spikes caused by checkpointing.

Recovery works by sequentially replaying logs read from log files. To maximize the I/O bandwidth during recovery, several loaders simultaneously read logs from files and fill the buffers. Since the logs in each log block are guaranteed to be serialized, the recoverer compares the timestamps of the log blocks and processes them in order, starting with the log block having the smallest value.

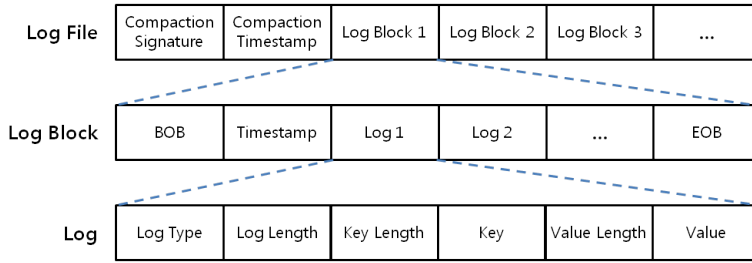


Figure 3.2 The log format.

We will explain each scheme in more detail in the following sections.

3.2 Distributed Logging and Log Data Format

For strong durability, logs should be immediately stored in stable storage. Unfortunately, synchronous durability leads to performance degradation. Therefore, many systems adopt asynchronous durability that buffers logs and flushes them to stable storage periodically. If the interval is too long, a large amount of data may be lost in case of system failure. If the interval is too short, the system throughput may be degraded due to write stalls. Write stalls occur when the storage device operation for the previous buffer flush is not completed when attempting to flush the buffer.

We use multiple storage devices to store logs in order to overcome the limitation on durability due to the storage device's performance. In addition, the use of multiple storage devices makes it possible to separate storage I/O for logging and checkpointing. This strategy ensures a stable log buffer flush cycle and system throughput by avoiding the effect of latency spikes on logging that can occur due to large data storage during checkpointing.

Figure. 3.1(a) shows the processing and logging for client requests. When a request such as insert or update is received from a client (①), the worker

reflects the processing result to the corresponding entry in the database and also stores the location information where the log will be stored in a variable called `log_pos` (②). The variable exists for each entry in the database and is used for checkpointing. If the entry already exists, the worker updates the `invalid_bitmap` to invalidate its old log (③). After that, the worker creates a log and writes it to the active log buffer (④). When the log buffer is filled to more than the minimal number of logs, the worker requests that the log buffer be flushed and then activates the next log buffer. Finally, durability is guaranteed when the buffer is flushed to stable storage by the flusher (⑤). For recovery, it is necessary to identify the order of logs distributed across multiple SSDs, but because the logs stored in each buffer are serialized, we only need to clarify the order between log blocks, which is a unit flushed from buffer to storage. To do this, we add a 4-byte timestamp to the header of the log block.

Figure. 3.2 shows the format of the log. The log file is composed of log blocks, and a compaction signature and compaction timestamp to indicate recent compaction histories are placed at the beginning of the file. The log block is a collection of logs that are flushed at a time. 1-byte BOB (Beginning of Block) and 1-byte EOB (End of Block) indicate the log block’s start and end. Between them are 4 bytes of timestamp and logs that include type, length, key, and value. The log type indicates the command or data type, and the length information is used to parse the log or determine the size to be copied during compaction.

3.3 Log File Compaction

We propose the VTC, a checkpointing scheme to reduce recovery time and recycle log space. Figure. 3.3 presents the key idea of VTC. The VTC identifies

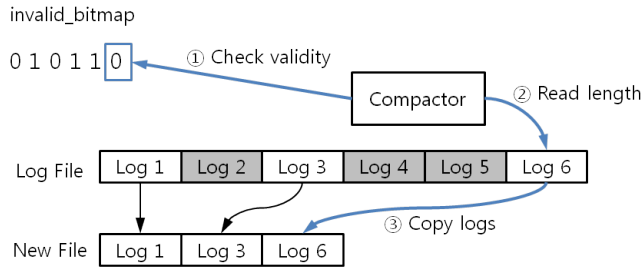


Figure 3.3 Log file compaction.

the log needed for recovery by referring to `invalid_bitmap` (①). It then checks the length of the log (②) and copies it to a new file if it is valid (③).

To manage the validity of logs in files, it is necessary to know the location of the latest log for each entry in the database. For this reason, the variable `log_pos` was defined to indicate the location of the log within the structure of the existing database entries. Logs in the file are accessed sequentially from the first log when restoring data or checkpointing. Because there is no need to search for a specific log, `log_pos` represents the order of logs in the file rather than a byte address. The upper few bits of `log_pos` are used to indicate which file the log is stored in. The `invalid_bitmap` allows the VTC to immediately recognize whether the logs are valid. One `invalid_bitmap` exists for each log file, and its bits indicate the validity of each log in the file. If the bit is 1, it means that the log is no longer needed for recovery.

Figure. 3.4 shows how the worker processes inserts and updates from clients. First, when a client requests an insert for a new key, the worker creates an entry with the key and value and adds it to the database. In addition, the location where the log will be stored is recorded in the variable `log_pos` in the entry (①). Then, the worker creates a log for insert and appends it the active log file (②). Later, when an update request for the same key is received, the worker

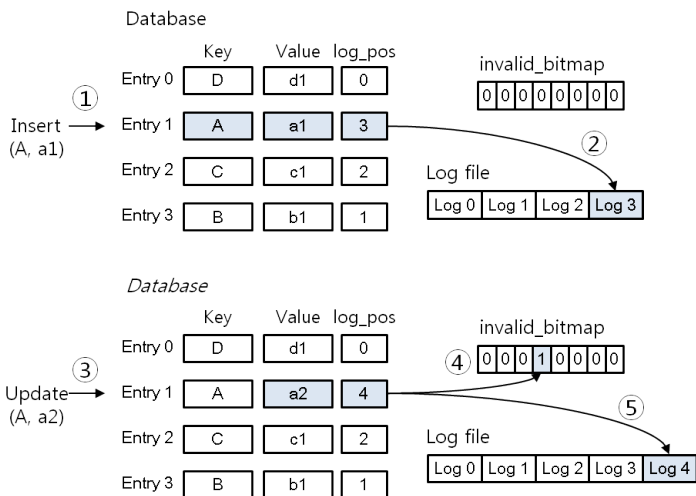


Figure 3.4 Example of logging for insert and update.

updates the value and `log_pos` in entry (③). Because the previous insert log for the key is no longer needed for recovery, the worker sets the bit in the `invalid_bitmap` corresponding to the old value of `log_pos` (④). Finally, a log for the update request is created and stored in the active log file (⑤). The insert log and update log for the same entry may be saved in different files. The procedure for delete requests is similar. The worker deletes the entry and sets the `invalid_bitmap` in the same way as for update. Subsequently, the worker appends the delete log to the active log file for recovery.

The gradually accumulated logs not only require more disk space but also take more time to recover. VTC prevents this problem by removing unnecessary logs from log files. Algorithm 1 and Algorithm 2 describe the VTC procedure, which requires four steps: preparing, copying, remapping, and completing. In the preparing step, VTC creates a new file, `temp.log`, to copy the valid log and writes the current timestamp as the compaction timestamp. In addition, VTC allocates memory for `delta`, a temporary array for calculating the locations of

Algorithm 1 Overview of the VTC procedure (step 1–2)

```
1: cp_time = get_curr_time() ▷ step 1
2: min_cp_time = get_min_cp_time()
3: dst_fp = create_file(temp.log)
4: write_cp_header(dst_fp, cp_time)
5: delta = allocate(log_count)
6: for each log blocks stored in file do ▷ step 2
7:     blk_ts = get_timestamp(block)
8:     for each logs stored in logblock do
9:         type = get_type(log)
10:        len = get_len(log)
11:        if type = delete then
12:            if blk_ts > min_cp_time then
13:                copy_log(dst_fp, log, len)
14:            else
15:                removal = removal + 1
16:        else
17:            if is_valid(bitmap, i) then
18:                copy_log(dst_fp, log, len)
19:            else
20:                removal = removal + 1
21:            delta[log_num] = removal
22:            i = i + 1
23: flush(dst_fp)
24: clear_invalid_bitmap(bitmap)
```

logs moved by copying (Algorithm 1, step 1, lines 1–5). The size of the array is determined by the number of logs in the file. Then, in the second step, copying, VTC sequentially reads logs from the target log file and copies only valid logs to `temp.log` file. VTC identifies the validity of each log by referring to the `invalid_bitmap` and copies the logs by referring to the log length in the header, without complicated parsing (Algorithm 1, step 2, lines 6–20). This step also fills the `delta` array to be referred to in the next step, remapping (Algorithm 1, step 2, line 21). Each element of the `delta` array corresponds to the logs in the file at the start of compaction. The values of the elements indicating the number of previously removed logs are referred to update the location of logs moved by compaction. After copying all valid logs, the VTC completes the second step

Algorithm 2 Overview of the VTC procedure (step 3–4)

```
1: for each entries stored in DB do ▷ step 3
2:   old_pos = entry→log_pos
3:   if is_on_compaction(old_pos) then
4:     entry→log_pos = old_pos - delta[old_pos]
5:   if entry→lazy = 1 then
6:     set_invalid_bitmap(bitmap, entry→log_pos)
7:     entry→lazy = 2
8: release(delta) ▷ step 4
9: rename_file(dst_fp)
10: delete_file(src_fp)
```

by clearing all bits of the `invalid_bitmap` for reuse (Algorithm 1, step 2, line 24). The next step is remapping to update the log location of each entry in the database (Algorithm 2, step 3, lines 1–4). VTC traverses all the entries in the database and adjusts the value of `log_pos`, which has the location of the latest log. VTC determines how much to change the `log_pos` value of each entry by referring to the `delta` array filled in step 2. For example, if `log_pos` is n , VTC gets the value of `delta[n]` and then decreases `log_pos` by the value of `delta[n]`. After updating the log location of entry, if the entry is lazily invalidated, the VTC reflects it in the `invalid_bitmap` (Algorithm 2, step 3, lines 5–7). Lazy invalidation will be discussed in more detail in the next section. When the above steps are completed, the VTC completes compaction by releasing the temporary array, deleting old log file, and renaming the new log file (Algorithm 2, step 4, lines 8–10).

The VTC’s handling of delete logs is different from write logs. Because the delete log is not invalidated by other logs, VTC does not refer to the `invalid_bitmap` when removing the delete log. Instead, the delete log can be removed when all other logs for the same entry are removed through checkpointing. The VTC determines whether to remove the delete log by comparing the delete log’s timestamp with the compaction timestamp of each file. If the

timestamp value of the delete log is smaller than the compaction timestamps, the VTC removes it. Otherwise, the VTC should keep the delete log in the log file to ensure correct recovery (Algorithm 2, step 2, lines 11–15). We will explain the removal of the delete log in more detail in section 3.6.

3.4 Lazy Invalidation

Even if the worker does not append logs to the log file where compaction is in progress, the logs in the target file may be invalidated by an update or deletion. This may cause problems in log management. For example, if the worker updates an `invalid_bitmap` to invalidate a log that has already been copied during the copying step, the information is lost due to the `invalid_bitmap` initialization at the end of the step. As a result, the log is not removed even by subsequent compaction. This inconsistency continues until the `invalid_bitmap` is rebuilt by replaying the logs on recovery.

We solve this problem by applying a lazy invalidation strategy. If the log in the file where compaction is in progress needs to be invalidated by an update or delete request, we delay it until compaction is complete. Before updating the entry, the worker checks to see if the old log of the entry belongs to a file that is undergoing compaction. If it is true, the worker sets the entry’s lazy variable to 1 and does not change the `invalid_bitmap`. Instead, the VTC creates a new entry and adds it to the database. Lazy invalidation immediately releases the memory for the entry’s key and value, thereby mitigating the increase in memory usage caused by entries that are delayed for deletion.

Lazily invalidated entries are dealt with in the VTC’s remapping step. In the remapping step, when after the `log_pos` of the lazily invalidated entry is adjusted, the compactor notes the new location of the log with the

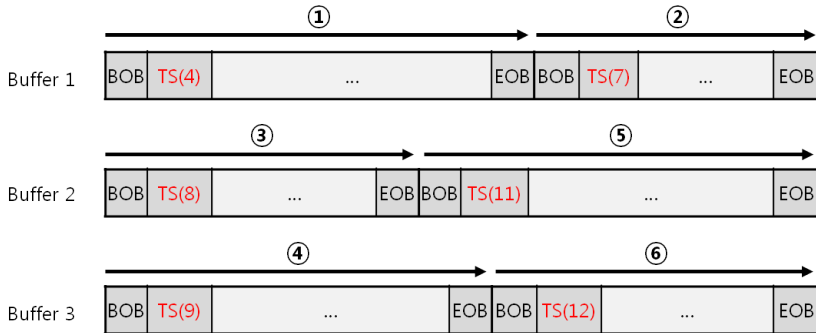


Figure 3.5 The recovery order.

`invalid_bitmap`, and the log is removed in the next compaction. The compactor completes the processing of lazy invalidation entries by setting the lazy variable of those entries to 2, indicating that they should be deleted by the worker later. If the compactor were to delete those entries, there could be a conflict with the worker’s add entry or delete entry, which is why we give the worker the role of deleting those entries.

3.5 Recovery

As shown in Figure. 3.1(b), the recoverer restores data by sequentially executing logs. In the VTC, we allocate loader and buffer per storage to maximize I/O bandwidth.

The loaders read log blocks from their respective files and fill the designated buffers (①). The recoverer selects log blocks in timestamp order and recovers data by replaying the logs in the log blocks (②). As the recoverer inserts the key and value into the database, it also stores the location of the log read for compaction (③). If overwrite or delete occurs during recovery, the recoverer records

the event in the `invalid_bitmap` (④). For example, as shown in Figure. 3.5, the recoverer compares the timestamp of each buffer's first log block and selects the log block with $TS(4)$ as the first recovery block. After completing the recovery of log block with $TS(4)$, the recoverer determines the next log block to recover by comparing the timestamp of log block with $TS(7)$ with others. In general, the use of timestamps causes an increase in recovery time. However, our recovery scheme uses a timestamp to determine the recovery order of log blocks. Since timestamp comparison is unnecessary while logs in a log block are sequentially executed, overhead due to timestamp use is insignificant.

When the recoverer encounters the write log, it checks whether an entry containing the same key exists in the database, and replays the log by selecting between insert and update. Therefore, a lookup is required for every log execution. To reduce the overhead caused by lookup, we divide the log into two groups based on whether a lookup is needed. The smallest value among compaction timestamps in each log file is the criterion for separating the two groups. The front group logs are guaranteed to have no duplicate keys, so keys and values are inserted immediately, without lookup.

Furthermore, we apply optimizations for the recovery of an sorted sets [25]. They spend a lot of time on recovery because they need to be sorted in ascending order. To optimize this, an additional log buffer is allocated, and the recoverer collects the logs that need to be sorted in the log buffer and batches them. This method helps reduce the recovery time by increasing the cache hit rate during sorting.

3.6 Correctness

Having explained how the VTC performs checkpointing by reconstructing log files using only the latest log for each entry, we will now provide proof of how the VTC guarantees correctness in all scenarios. The VTC maintains logs by entry, which is the smallest unit that the system can modify. This ensures that the system can restore data by replaying only the latest log of all entries. For instance, consider an entry A that is initially inserted with the value of a1, then updated to a2, and finally updated with the value of a3. These three write logs are stored across multiple log files. In this state, if a recovery proceeds due to system failure, the system executes three consecutive writes by referring to the timestamp of the log block. However, as a result, the final state of entry A is determined by the a3 update log, and the other two logs do not affect it. The VTC performs checkpointing individually for each file, and thus one or both of the a1 and a2 logs can be removed. Nevertheless, entry A can be restored correctly using the a3 update log remaining in the log file.

Next, we look at the process of the VTC removing the delete log. As we discussed, to properly remove the delete log, the VTC needs to confirm that all other logs of the deleted entry have been removed by checkpointing. If the VTC deletes the delete log without going through this process, at the time of recovery, entry A may be erroneously revived by a write log that may have remained in another log file. For example, consider the process in which VTC removes its logs after entry A is deleted. To remove the delete log of entry A, the VTC must first confirm that logs for three writes do not exist in other log files, which necessitates tracing all of the logs for each entry, incurring significant overhead. To avoid this, we use a compaction timestamp that represents the last checkpointing time for each file.

The VTC compares the delete log's timestamp with all compaction timestamps to determine removing the delete log of entry A. The timestamp of the delete log can be found by referring to the timestamp of the log block to which it belongs. If the delete log timestamp is less than the compaction timestamp from all of the log files, the VTC can safely remove the delete log because it is guaranteed that all three writes of entry A have been removed. Thus, in this case, no log for entry A remains, so no processing for entry A will occur during recovery. Conversely, if the compaction timestamp of any file has a value smaller than the delete log's timestamp, correct recovery can be guaranteed by maintaining the delete log for the corresponding entry. To do this, the VTC retains the delete log by copying it to a new file. In this case, the write log for entry A may be executed during recovery, but the delete log also remains, so entry A can be deleted and restored to the correct state.

3.7 Implementation

We implemented the proposed scheme on Redis 5.6.0. In Redis, write operations either create entries for new key-value pairs or update existing ones. Then write operations generate logs and write them to the log buffer. For these operations, we used the code path of Redis. To handle overwriting when the old log of an entry is stored in a file in which compaction is in progress, we insert codes for lazy invalidation, which sets a lazy variable and inserts a new entry instead of updating the entry. In addition, the entries contain the `log_pos` variable, which is 8 bytes in size, to keep track of their latest log. The upper 2 bits of the variable are reserved for the lazy invalidation of the entry. For the read operation we follow the Redis code path and add only the code to handle lazily invalidated entries. We also add functions for new algorithms and change the call path in

order to replace the existing Redis algorithms such as logging, checkpointing, and recovery.

We allocate as many threads as the number of files (storage devices) to flush the log buffer. These threads are responsible for reading log blocks from a file during recovery; we also add one thread for checkpointing. We only allow workers to add or remove entries in the database. This restriction prevents performance degradation due to contention between the worker and the compactor. Additionally, we use atomic operations to ensure atomicity for some variables that are shared between threads. We count the number of logs and the number of invalidation processes for each file for checkpointing with an appropriate frequency, and checkpointing is triggered when the number of invalidated logs and their ratio to total logs reach the thresholds. The user can determine thresholds in consideration of checkpointing frequency and execution time.

Chapter 4

Evaluation

4.1 Experimental Setup

In this section, we describe the experiments conducted to measure the performance of the proposed scheme under various conditions. We conducted all experiments using two machines as a client and server, each of which is equipped with an Intel Xeon W-2245 CPU running at 3.9 GHz; the CPU had 8 physical cores and 16 logical cores with hyper-threading and 32 GB of DRAM memory. The machines were connected through a 10 Gbps network. We used Samsung 860 PRO [26] SATA SSDs to store logs and checkpoints. Our scheme distributes logs to three SSDs, and Redis-AOF, which uses the existing logging scheme, stores logs in a single SSD or a RAID-0 array with three SSDs. One additional SSD was used as a swap device to prevent the process from being killed by an out-of-memory killer. The machines ran Ubuntu 18.04.4 LTS distribution with the Linux kernel 4.15.0.

To demonstrate the efficiency of our scheme, we applied the VTC to Redis-

Table 4.1 Parameters of YCSB workloads

Parameters	Setting
Record Count	2M, 4M, 6M, 8M, 10M, 12M
Update proportion	10%, 30%, 50%, 70%, 90%
Record Size	default
Distribution	zipfian
Number of threads	128

5.0.6 and compared the VTC performance to the following Redis protocols:

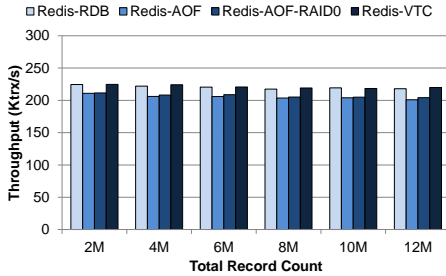
Redis-RDB: No logging, and all of the data were periodically backed up to a file through checkpointing.

Redis-AOF: This records the log of all events that change the database and manages the size of the log file through periodic checkpointing.

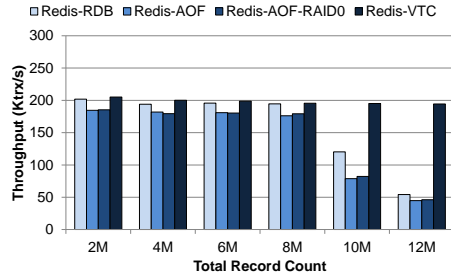
Redis-AOF with RAID-0 Setup: This has the same configuration as Redis-AOF except that it uses a RAID-0 array across the three SSDs handled through a software RAID driver in Linux.

To evaluate the performance of each scheme, we used Yahoo! Cloud Serving Benchmark (YCSB) [14] as the target workload. Table 4.1 summarizes the parameters of YCSB used for throughput evaluation. In order to evaluate the performance of each scheme with various configurations, we changed the number of records and the update proportions as shown in the table. After loading YCSB data into Redis, we measured the performance while the YCSB workload was running. For fair comparison we forced checkpointing at the same time. If the checkpointing time increased rapidly due to swap, we measured the performance for up to 600 seconds.

We disable the RDB compression option for a fair comparison because our prototype does not currently support data compression.

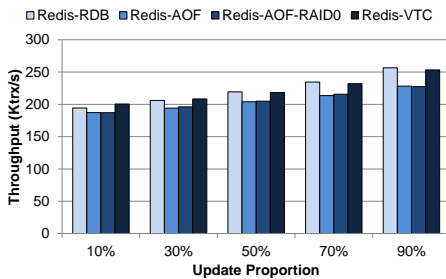


(a) Normal Throughput

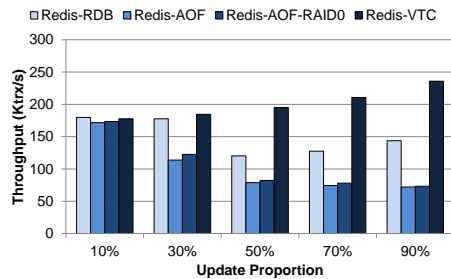


(b) Throughput during Checkpointing

Figure 4.1 Throughput for varying record count(50% update proportion)



(a) Normal Throughput



(b) Throughput during Checkpointing

Figure 4.2 Throughput for varying update proportion(10M records)

4.2 Performance

4.2.1 Throughput

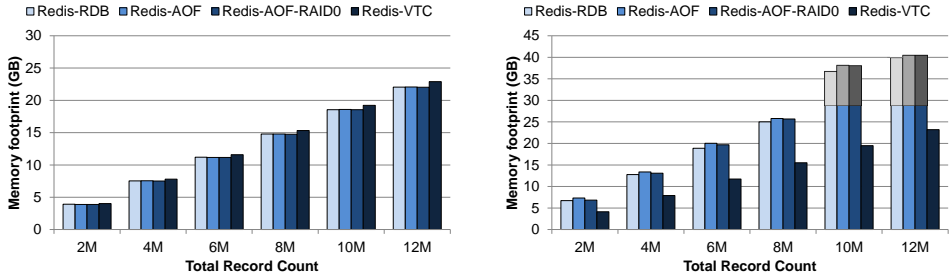
Figure. 4.1 shows the throughput of each system for the YCSB workload with various record counts. As shown in Figure. 4.1(a), Redis-VTC offers an average of 8% higher throughput than Redis-AOF with logging and similar throughput as Redis-RDB using only checkpoints. This means that Redis-VTC has little overhead for logging processing. Also, as the throughput of AOF-RAID0 is similar to that of Redis-AOF, we can see that the performance of the storage device does not affect the system throughput. Overall, for all of the schemes, the

throughput was not appreciably affected by the size of the data set. However, throughput during checkpointing tends to be different for each system depending on the size of the data set. Figure. 4.1(b) shows the throughput of the YCSB workload for a 50% update proportion while checkpointing was performed in the background. When the system memory was sufficient (a record count of 8M or less), each scheme exhibited a similar trend to normal throughput. Conversely, when the size of the data set increased, there was a difference in results that corresponded to the checkpointing scheme. In contrast to Redis-VTC, which showed stable throughput regardless of data set size, other schemes had severely degraded throughput when the size of the data set was larger than about 60% of system memory. Specifically, systems that use fork-based snapshots for checkpointing increase memory usage by CoW when processing an update request from a client. If the record count is 10M and the update proportion is 50%, more than 10 GB of memory is used by CoW during checkpointing. Eventually, swap due to insufficient system memory causes throughput degradation. In contrast, Redis-VTC performs checkpointing based on the validity of the logs and thus requires only a small amount of additional memory.

Figure. 4.2(b) shows the throughput for varying update proportions with a record count of 10M. As shown in Figure. 4.2(a), a high update proportion usually has a positive effect on performance. However, when the system memory is marginal, a high update proportion causes frequent CoW during checkpointing, which can cause swap. This means that using fork-based snapshots requires more extra system memory for update-intensive workloads.

4.2.2 Memory Footprint

Figure. 4.3(a) shows the memory footprint of the data set with varying record counts. Redis-VTC had a 3.5% larger memory footprint than other schemes be-



(a) Memory Footprint for Data Set

(b) Peak Memory Footprint (Gray bars denote swap memory)

Figure 4.3 Memory footprint comparison with existing schemes (50% update proportion)

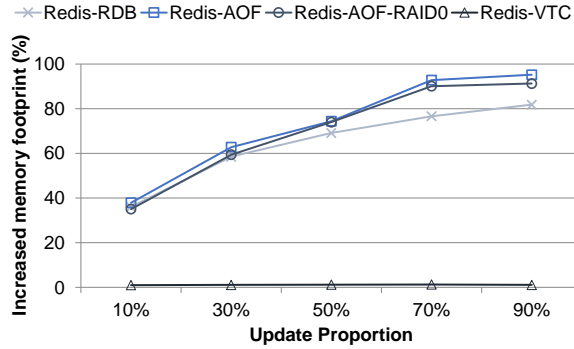


Figure 4.4 Increased memory footprint(8M records)

cause variables for managing logs were added for each entry. However, during checkpointing, Redis-VTC required less additional memory than other protocols. As shown in Figure. 4.3(b), during checkpointing, the memory footprint of Redis-RDB and Redis-AOF increased by 69.6% and 79.7%, respectively, while the memory increase of Redis-VTC is less than 2%. Figure. 4.4 shows that this gap can be wider as the update proportion is increased. This is because the higher the update proportion, the more CoW that occurs during checkpointing, which consumes more memory.

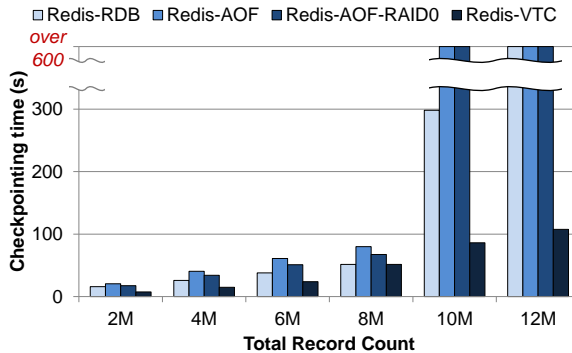


Figure 4.5 Checkpointing time

Redis-VTC creates a temporary array to update the location of the logs during checkpointing, but this is insignificant compared to the size of the entire data set because it only requires 8 bytes per entry. However, with Redis-RDB and Redis-AOF, memory usage continuously increases from CoW during checkpointing. Moreover, the increased memory cannot be reclaimed until the process in charge of storing the snapshot is terminated.

4.2.3 Checkpointing Time

The VTC performs checkpointing independently for each log file, but in order to have a fair comparison, we measured time by sequentially performing checkpointing for all log files. Redis-RDB was also configured not to use incremental-`fsync` option as Redis-VTC uses `fsync()` only once, after the last write when saving checkpoints to a file for optimization. However, this option was enabled in Redis-AOF because logging can be affected by latency spikes, and enabling this option increases the checkpoint time by about 20% because `fsync()` is called for every 32 MB write.

Figure. 4.5 shows the time taken for checkpointing with varying record

counts. In all schemes, under conditions of sufficient memory (record counts of less than 6M), the processing time increased as the size of the data set rose. We found that Redis-VTC required less than half the time of the other schemes under this condition. Because Redis-VTC performs checkpointing through file-to-file copy, it requires more I/O than other schemes to read the logs from storage device. However, if the system memory is sufficient, Redis-VTC can read logs to be copied from the buffer cache. Furthermore, a simple way to determine which log to copy by bitmap reference reduces checkpointing time.

We can see that the processing time for all three schemes increases as memory becomes insufficient. In particular, for schemes other than Redis-VTC, checkpointing time increases rapidly by swap. In this case, we only plotted up to 600 seconds, but checkpointing would normally take ten or more minutes. Redis-VTC also takes more time for checkpointing if the record count is 8M or more, because some logs are read from disk due to insufficient memory used as the buffer cache. However, Redis-VTC alleviates the increase in processing time because it reads the contents of a file sequentially, without random access.

4.2.4 File Size

Redis allows insert or update requests with set data [25] that include several sub-key and value pairs in the key. Because Redis-VTC has to manage the validity of logs by entry, it records the set data as separate logs for each sub-key/value pair. In this case, the file size increased because each separate log must include the parent key's information. For this reason, the YCSB workload, which uses bundled requests for data set loading before performance measurement, is not good for Redis-VTC in terms of file size.

Figure. 4.6 is the result of measuring the file size with a varying record count for the three schemes. We measured the file size when the YCSB workload ran

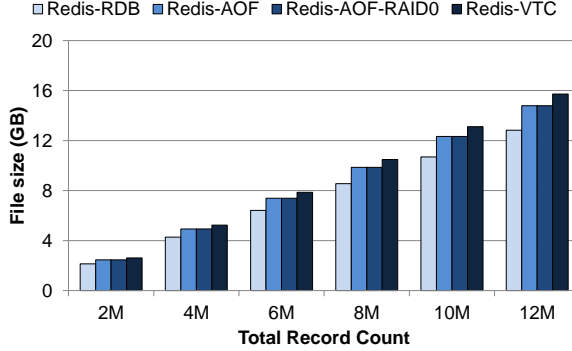


Figure 4.6 File size

the workload at a 50% update rate for 60 seconds after loading the data. The file sizes of Redis-VTC and Redis-AOF were measured before checkpointing. The file size of all schemes increases in proportion to the record count. For the reasons described above, the file size of Redis-VTC was, on average, 22% larger than Redis-RDB and 6% larger than Redis-AOF. For all schemes other than Redis-RDB, which does not use logging, the file size may be larger depending on the checkpointing interval and running time.

4.2.5 Restoring Time

Figure. 4.7 shows the restoring time by record count. We measured the data recovery time under the same conditions as measuring the file size. Redis-VTC took longer by an average of 53% to restore as compared to Redis-RDB. As described in the result of the file size, because the set data are logged separately, Redis-VTC executes each log independently for recovery. Conversely, Redis-RDB reduces the recovery time by handling the data set all at once. Moreover, as a property of checkpointing using snapshot, all data with the same parent key are stored together. This enables fast sorting with a high cache hit rate

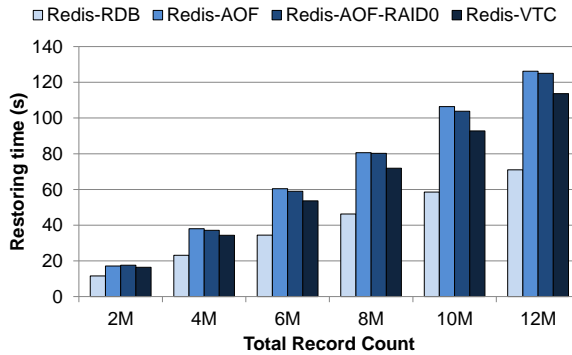


Figure 4.7 Restoring time

when recovering the sorted data set, which is one of the data types used by the YCSB workload. According to the results, Redis-VTC had 10% faster recovery time than Redis-AOF. However, if recovery was attempted after checkpointing, Redis-AOF recovered as quickly as Redis-RDB. Although our scheme takes more time to restore than other schemes with consistent snapshots, we believe that reducing memory usage and maintaining stable throughput are preferable for IMDBs. Moreover, data restoration may not occur frequently.

Chapter 5

Related Work

In this section, we epitomize the techniques required to provide persistence in IMDBs. There have been several studies to provide logging suitable for IMDBs. The fast processing speed of IMDBs makes the run-time overhead of traditional logging relatively large. To avoid the throughput degradation caused by logging, many IMDBs [8], [27] attempt lightweight logging based on logical logging. Command logging [28] used by H-Store is a logical logging variant that records a single log record including only a procedure ID and input parameters. While redo-only logging such as command logging reduces run-time overhead, it takes more time to recover because logs are replayed for recovery.

Adaptive logging [29] is a logging method that focuses on the balance between run-time and recovery performance. It extends command logging to distributed systems and allows all nodes to perform recovery in parallel. Moreover, It reduces the recovery time by re-executing only transactions related to the failed node through dependency analysis between transactions. Taurus [30] and SiloR [31] also use distributed logging to solve the performance bottleneck prob-

lem caused by logging. Taurus performs logging in parallel using a log sequence numbers vector to manage transaction dependency. Taurus performs logging in parallel and manages the dependency of logs scattered in several files using a log sequence numbers vector. SiloR allocates one thread for each storage to write and flush logs in parallel. The system performs group commit in the epoch using optimistic concurrency control.

Our logging system is in line with these approaches [29], [30], [31], in terms of storing logs across multiple storage devices. However, in general, single-stream logging records are logged in one storage device, but our scheme improves durability by distributing logs across multiple storage devices. Furthermore, our scheme avoids logging and checkpointing I/O from affecting each other by separating them into different storage devices.

An efficient checkpointing method has been continuously proposed by several studies. Systems for some applications are sufficient to guarantee durability with only periodic checkpoints. However, most systems use a combination of logging and checkpointing to minimize data loss due to system failure. The checkpointing method employed by many IMDBs is to use consistent snapshots. It takes a consistent snapshot of the in-memory and stores its contents in stable storage. Representative consistent snapshot algorithms include naive snapshot [32], copy-on-update(COU) [2], [21], Zigzag [22], and PingPong [22]. Naive snapshot stops the system for a consistent snapshot. It is the simplest algorithm, but it is not suitable for in-memory database systems that need to process transactions even during checkpointing. COU is the most widely used algorithm for non-blocking checkpointing, and a number of variants have been studied. In contrast to the general COU algorithm using physical page shadowing, SIREN [21] proposes a COU algorithm based on tuple units smaller than pages. Small duplication granularity has the effect of reducing memory usage,

but the average latency may increase due to tuple-level locking. Algorithms using `fork()` are applied to many IMDBs [7], [8] because it can easily implement COU with OS support, but there are problems with latency spike and memory increase. DMC [11] uses a memory dump to overcome the increase in memory usage, returning the page to the OS before the checkpoint is complete. However, at high update proportions, this scheme also requires a significant memory footprint.

Incremental or partial checkpointing reduces cost by limiting the amount of data processed at one time. Hekaton [12] creates a checkpoint file from the transaction logs not covered by a previous checkpointing and manages updates or deletions by recording them in delta files. The incremental checkpoints used by Hekaton can lower the cost by creating a checkpoint file only for new transactions. In contrast, a large number of files and a high ratio of deleted contents in the checkpoint file degrades recovery performance. To alleviate this, an additional process such as merging between checkpoints files is required. CALC [13] supports partial checkpointing. It performs checkpointing, including only records that have changed since the most recent checkpoint. It is effective in workloads where updates are not frequent, whereas in the opposite case, it may be inefficient in comparison to complete checkpointing due to overhead for merging files.

These systems that checkpoint data in main memory require page copying or version control, which causes increased memory usage. However, our scheme consumes less memory than the existing checkpointing scheme because it performs checkpointing using the log in the file. This enables efficient use of memory, the most important resource for IMDBs.

Chapter 6

Conclusion

This paper proposed an effective transaction log-based persistence scheme for IMDBs. Our key idea is to distribute logs across multiple storage devices and use log file compaction for checkpointing. The use of multiple storage devices improves the durability of the log without sacrificing system throughput. The log file compaction by managing log validity can keep the peak memory usage lower than the scheme using consistent snapshots. We implemented and evaluated our scheme in a famous IMDB, Redis, and the experimental results show that our proposed scheme can more stably manage memory during checkpointing compared to the scheme currently applied to commercial IMDBs. This is very desirable for IMDBs, because it allows the memory reserved for peak usage to be used for data storage.

Bibliography

- [1] H. Zhang, G. Chen, B. C. Ooi, K.-L. Tan, and M. Zhang, “In-memory big data management and processing: A survey,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 7, pp. 1920–1948, 2015.
- [2] M. Vaz Salles, T. Cao, B. Sowell, A. Demers, J. Gehrke, C. Koch, and W. White, “An evaluation of checkpoint recovery for massively multi-player online games,” *Proceedings of the VLDB Endowment*, vol. 2, no. 1, pp. 1258–1269, 2009.
- [3] B. K. Park, W.-W. Jung, and J. Jang, “Integrated financial trading system based on distributed in-memory database,” in *Proceedings of the 2014 Conference on Research in Adaptive and Convergent Systems*, pp. 86–87, 2014.
- [4] C. Mohan, D. Haderle, B. Lindsay, H. Pirahesh, and P. Schwarz, “Aries: A transaction recovery method supporting fine-granularity locking and partial rollbacks using write-ahead logging,” *ACM Transactions on Database Systems (TODS)*, vol. 17, no. 1, pp. 94–162, 1992.
- [5] L. Li, G. Wang, G. Wu, and Y. Yuan, “Consistent snapshot algorithms for in-memory database systems: experiments and analysis,” in *2018 IEEE*

- 34th International Conference on Data Engineering (ICDE)*, pp. 1284–1287, IEEE, 2018.
- [6] L. Li, G. Wang, G. Wu, Y. Yuan, L. Chen, and X. Lian, “A comparative study of consistent snapshot algorithms for main-memory database systems,” *IEEE Transactions on Knowledge and Data Engineering*, 2019.
- [7] Redis, “Redis.” <http://Redis.io/>. (2021).
- [8] A. Kemper and T. Neumann, “Hyper: A hybrid oltp&olap main memory database system based on virtual memory snapshots,” in *2011 IEEE 27th International Conference on Data Engineering*, pp. 195–206, IEEE, 2011.
- [9] Wikipedia, “Fork (system call).” [https://en.wikipedia.org/wiki/Fork_\(systemcall\)](https://en.wikipedia.org/wiki/Fork_(systemcall)). (2021).
- [10] Redis, “latency problems troubleshooting.” <http://redis.io/topics/latency>. (2021).
- [11] J. Park, Y. Lee, H. Y. Yeom, and Y. Son, “Memory efficient fork-based checkpointing mechanism for in-memory database systems,” in *Proceedings of the 35th Annual ACM Symposium on Applied Computing*, pp. 420–427, 2020.
- [12] C. Diaconu, C. Freedman, E. Ismert, P.-A. Larson, P. Mittal, R. Stonecipher, N. Verma, and M. Zwillig, “Hekaton: Sql server’s memory-optimized oltp engine,” in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pp. 1243–1254, 2013.
- [13] K. Ren, T. Diamond, D. J. Abadi, and A. Thomson, “Low-overhead asynchronous checkpointing in main-memory database systems,” in *Proceedings*

- of the 2016 International Conference on Management of Data, pp. 1539–1551, 2016.
- [14] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, and R. Sears, “Benchmarking cloud serving systems with ycsb,” in *Proceedings of the 1st ACM symposium on Cloud computing*, pp. 143–154, 2010.
- [15] F. Faerber, A. Kemper, P.-Å. Larson, J. Levandoski, T. Neumann, A. Pavlo, *et al.*, *Main memory database systems*. Now Publishers, 2017.
- [16] T. Wang, R. Johnson, and I. Pandis, “Query fresh: Log shipping on steroids,” *Proceedings of the VLDB Endowment*, vol. 11, no. 4, pp. 406–419, 2017.
- [17] T. Mühlbauer, W. Rödiger, A. Reiser, A. Kemper, and T. Neumann, “Scyper: Elastic olap throughput on transactional data,” in *Proceedings of the Second Workshop on Data Analytics in the Cloud*, pp. 11–15, 2013.
- [18] M. Stonebraker and A. Weisberg, “The voltdb main memory dbms.,” *IEEE Data Eng. Bull.*, vol. 36, no. 2, pp. 21–27, 2013.
- [19] H. V. Jagadish, D. Lieuwen, R. Rastogi, A. Silberschatz, and S. Sudarshan, “Dali: A high performance main memory storage manager,” in *VLDB*, vol. 94, pp. 48–59, 1994.
- [20] H. Jagadish, A. Silberschatz, and S. Sudarshan, “Recovering from main-memory lapses.,” in *VLDB*, vol. 93, pp. 391–404, Citeseer, 1993.
- [21] A.-P. Liedes and A. Wolski, “Siren: A memory-conserving, snapshot-consistent checkpoint algorithm for in-memory databases,” in *22nd International Conference on Data Engineering (ICDE’06)*, pp. 99–99, IEEE, 2006.

- [22] T. Cao, M. Vaz Salles, B. Sowell, Y. Yue, A. Demers, J. Gehrke, and W. White, “Fast checkpoint recovery algorithms for frequently consistent applications,” in *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data*, pp. 265–276, 2011.
- [23] D.-E. Ranking, “Db-engines ranking.” <http://db-engines.com/en/ranking>. (2021).
- [24] Redis, “Administration.” <https://Redis.io/topics/admin>. (2021).
- [25] Redis, “Data types.” <https://Redis.io/topics/data-types>. (2021).
- [26] Samsung, “Samsung ssd 860 pro.” <https://www.samsung.com/semiconductor/minisite/ssd/product/consumer/860pro>. (2021).
- [27] R. Kallman, H. Kimura, J. Natkins, A. Pavlo, A. Rasin, S. Zdonik, E. P. Jones, S. Madden, M. Stonebraker, Y. Zhang, *et al.*, “H-store: a high-performance, distributed main memory transaction processing system,” *Proceedings of the VLDB Endowment*, vol. 1, no. 2, pp. 1496–1499, 2008.
- [28] N. Malviya, A. Weisberg, S. Madden, and M. Stonebraker, “Rethinking main memory oltp recovery,” in *2014 IEEE 30th International Conference on Data Engineering*, pp. 604–615, IEEE, 2014.
- [29] C. Yao, D. Agrawal, G. Chen, B. C. Ooi, and S. Wu, “Adaptive logging: Optimizing logging and recovery costs in distributed in-memory databases,” in *Proceedings of the 2016 International Conference on Management of Data*, pp. 1119–1134, 2016.
- [30] Y. Xia, X. Yu, A. Pavlo, and S. Devadas, “Taurus: lightweight parallel logging for in-memory database management systems,” *Proceedings of the VLDB Endowment*, vol. 14, no. 2, pp. 189–201, 2020.

- [31] W. Zheng, S. Tu, E. Kohler, and B. Liskov, “Fast databases with fast durability and recovery through multicore parallelism,” in *11th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 14)*, pp. 465–477, 2014.
- [32] G. Bronevetsky, R. Fernandes, D. Marques, K. Pingali, and P. Stodghill, “Recent advances in checkpoint/recovery systems,” in *Proceedings 20th IEEE International Parallel & Distributed Processing Symposium*, pp. 8–pp, IEEE, 2006.

초록

인-메모리 데이터베이스는 메인 메모리에 상주해 있는 데이터셋에서 트랜잭션을 처리하기 때문에 클라이언트 요청에 대한 빠른 응답시간을 달성할 수 있다. 처리 속도의 향상은 인해 트랜잭션의 내구성을 보장하기 위한 기존의 로깅 기법과 체크포인팅 기법의 비용을 상대적으로 크게 만든다. 많은 인-메모리 데이터베이스가 로그의 부피를 줄이는 것으로 통해 로그 생성과 로그 저장 IO에 의한 오버헤드를 감소시키지만 그것은 복구 시간의 증가를 가져온다. 주기적인 체크포인팅은 복구 시간을 감소시키고 로그의 저장 공간을 재사용할 수 있도록 한다. 하지만 기존의 체크포인트 방법은 종종 시스템의 작업량 저하, 지연 증가, 메모리 사용량 증가로 인해 상당한 비용이 발생한다.

이 논문에서는 파일 내 로그의 유효성을 추적하고 불필요한 로그를 제거하는 기술인 validity tracking-based compaction (VTC)를 사용한 체크포인팅을 제안한다. 우리가 제안 하는 방식은 스냅샷을 사용하는 기존 체크 포인트 방식에 비해 메모리 사용량이 매우 낮출 수 있다. 우리의 실험에 따르면 기존의 체크포인팅 방법은 업데이트가 집중되는 워크로드에서 메모리 사용량이 최대 2배까지 증가 하는 것이 비하여 VTC는 2% 미만의 증가를 보인다. 그것은 시스템이 메모리의 대부분을 데이터를 보관하고 트랜잭션을 처리하기 위해서 사용할 수 있다는 것을 의미한다.

주요어: In-Memory Database, Persistence, Logging, Checkpointing, Snapshot

학번: 2019-23670