

# DSM: A Low-Overhead, High-Performance, Dynamic Stream Mapping Approach for MongoDB

TRONG-DAT NGUYEN AND SANG-WON LEE

*College of Information and Communication Engineering*

*Sungkyunkwan University*

*Suwon 16419, Korea*

*E-mail: {datnguyen; swlee}@skku.edu*

For write-intensive workloads, reclaiming free blocks in flash SSDs is expensive due to data fragmentation problem that leads to performance degradation. This paper addresses that problem in MongoDB, a popular document store in the current market, by introducing a novel stream mapping scheme that exploits unique characteristics of MongoDB and multi-streamed technology. It dynamically assigns streams for corresponding writes according to their hotness values and distinguishes writes on primary index files from writes on secondary index files. The proposed method is high-performance, low-overhead, and independent with data models or workloads. Empirical results in Linkbench benchmark show that compared to the original WiredTiger our approach improves the throughput and reduces the 99th-percentile latency by up to 65% and 46.2% respectively. Compared to the best-performance in the prior research, our approach improves the throughput and reduces the 99th-percentile latency by up to 23% and 28.5% respectively. Distinguishing writes on primary index files from writes on secondary index files enhances the throughput and the 99th-percentile latency by up to 11.7% and 15.7% respectively. Moreover, by tuning the leaf page size in B+Tree of MongoDB, we can significantly improve the throughput by 1.6x – 2.1x in Linkbench.

**Keywords:** data fragmentation, hot/cold data identification, multi-streamed SSD, NoSQL database, MongoDB

## 1. INTRODUCTION

NAND flash Solid state drives (SSDs) perform "erase-before-update" i.e., one data block should be erased before writing on new data pages [1, 2]. Because erase operation itself takes orders of magnitude slower than read and write operations [3], flash SSDs write updated data in empty blocks and mark the old data pages as invalid instead of erasing the current data blocks. When the number of empty blocks is lower than a threshold, *Garbage Collection (GC)*, a component of Flash Translation Layer (FTL) inside flash SSD that responsible for maintaining free blocks, is background processed to reclaim free blocks. During this process, if a non-empty data block is selected as a victim for reclaiming, valid pages from that block need to be copied back to another empty block before the actual erase operation is done. That lead to increasing the overhead of the reclaiming process. Moreover, NAND flash blocks have limited number of erase cycles; Flash FTLs use *wear-leveling*, a technique that ensures writes are distributed evenly among flash blocks, to enhance the flash SSD lifespan.

The locality of data access has a significant impact on the performance of flash memory and its lifetime due to the high-overhead of reclaiming free blocks and wear-leveling. In practice, IO workloads from clients exist skewness, i.e., a small proportion of data has

frequently accessed [4, 5, 6]. That forms the hot logical block addresses (LBAs) (LBAs have frequently accessed) and the cold LBAs (the remain LBAs have less frequently accessed) which are called in short as *hot data* and *cold data* respectively. *Data fragmentation* in flash SSD happens when one physical block includes hot data and cold data which in turn, increases the overhead of reclaiming blocks significantly.

Data fragmentation problem can be solved by identifying hot/cold data either based on update frequency [7, 8, 9] or based on history address information [10]. However, those approaches need to keep track of metadata in DRAM and high-cost of CPU for identifying hot/cold blocks thus increase the overhead of the system. One prior research designed a Flash-oriented file system that groups hot and cold segments according to write frequencies [11]. In another research, a write buffer cache module is used in SSD controller to separate sequential writes from random writes, achieved a higher level of parallelism for both patterns, thus improved the performance of the system [12].

With the high volumes and varieties of generated data, not-only SQL (NoSQL) solutions have become the alternatives for traditional RDBMSs [13, 14, 15, 16, 17]. Data fragmentation in DBMSs is a common problem not only in RDBMS but also in NoSQL solutions. For instance, Cassandra and RocksDB take the log-structured merge (LSM) tree [18] approach that has different update lifetime for files in each level of the LSM tree. *Multi-streamed SSD (MSSD)* technique [4, 5] assign different *streams* to different file types in Cassandra, RocksDB, so that writes on the similar update lifetimes are occurred on the same physical data blocks. In the virtualization environment, data fragmentation exists in flash SSDs when virtual machines have different IO workloads but share the same physical storage devices. One solution to that problem is that writes from one virtual machine are mapped to the same stream so that the scheduler can provide the corresponding resource [19].

To the best of our knowledge, there are a few studies have investigated on data fragmentation problem in MongoDB [20] - one of the common document stores with WiredTiger [21] as the default storage engine. Recently, most of the researchers compared between RDBMSs and NoSQLs. [22, 23, 24, 25], addressed data modeling transformation [15, 26, 27, 28, 29] or archived load-balanced sharding [30, 31]. In another research, authors argued different logging techniques in MongoDB and proposed a simple log management model that is useful for profiling the system [32]. Based on unique characteristics of space management in WiredTiger, TRIM command technique is used to reduce the overhead of MongoDB [33]. However, TRIM command does not entirely solve the data fragmentation [4]. In another research, file-based approach and boundary-based approach are proposed to address the data fragmentation in MongoDB [34]. However, those methods still have a level of data fragmentation such that writes from regions with different lifetimes are mapped to the same stream. Also, writes on primary indexes have different patterns and lifetimes compare to writes on secondary indexes. The boundary-based method is inadequate to distinguish writes on those two index types.

To solve those problems, we propose a novel dynamic stream mapping (DSM) that is an online high effective stream mapping based on hot/cold values of each data block. We summarize our contributions as below:

- First, by investigating WiredTiger's block management in detail and revisiting boundary-based approach with a complex data model, i.e., Linkbench benchmark [35], we pointed out two flaws existing in the approach: (1) writes from regions with different lifetime update are mapped to the same stream, and (2) writes from primary indexes mixed with writes from secondary indexes.
- Based on this observation, we proposed a novel mapping scheme named Dynamic

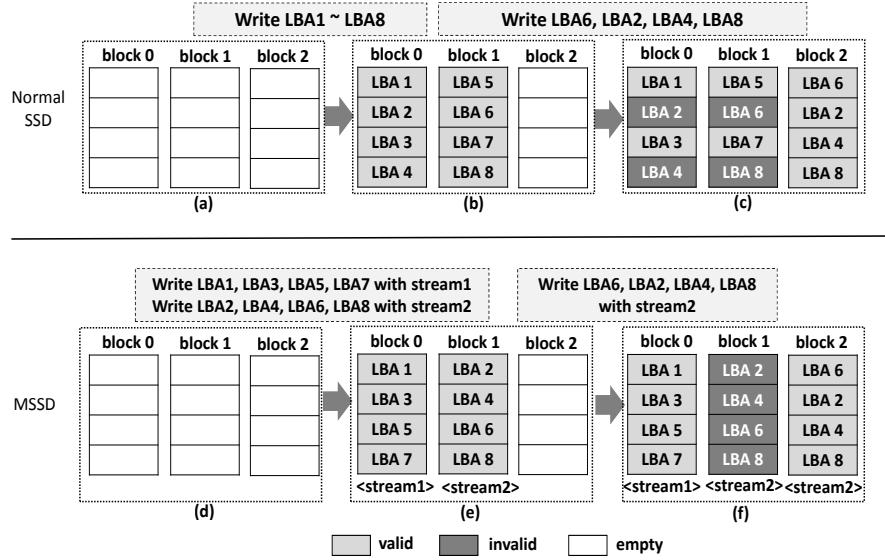


Fig. 1. Comparison between normal SSD and multi-streamed SSD

Stream Mapping (DSM) that grouped writes on corresponding streams based on their hotness values and separates writes on primary index files in distinguishing stream with writes on secondary index files. In Linkbench, compare to the original WiredTiger, DSM improves the throughput and the 99th-percentile latency by up to 65% and 46.2% respectively. Compare to the best-performance method in the prior research, our proposed method enhances the throughput and the 99th-percentile latency by up to 23% and 28.5% respectively. Moreover, index filtering gains additional 11.7% and 15.7% improvement in term of throughput and 99th-percentile latency respectively.

- Lastly, we combine the leaf page size tuning with DSM; the final result improved the throughput by 1.6x – 2.1x for Linkbench benchmark.

The rest of this paper is organized as follow. Section 2 explains the background of multi-streamed SSD and MongoDB in detail. We revisit the prior works in Section 3. Proposed methods and related algorithms are described in Section 4. Section 5 discusses evaluation results and analysis. Lastly, the conclusion is given in Section 6.

## 2. BACKGROUND

This section provides a background of Multi-streamed SSD technique that originally used in Cassandra and exploited in our proposed method. Also, we briefly introduce the relationship between data model in RDBMS and MongoDB and discuss the block management mechanism of WiredTiger in detail.

### 2.1 Multi-streamed SSD

Multi-streamed SSD is a technique allow applications from *user space* or *kernel space* explicitly assigned a *stream* along with a *write* or *pwrite* system call. The FTL inside flash SSD use that stream information to group writes mapped with the same stream

into one physical SSD block [5]. Multi-streamed technique requires the modification from both OS's kernel and the firmware in physical SSD devices that available in commercial products [36]. Fig. 1 compares the differences between the regular SSD and the multi-streamed SSD. Suppose that LBA2, LBA4, LBA6, and LBA8 are hot data; LBA1, LBA3, LBA5, and LBA7 are cold data. There are two write sequences occur on both SSDs. The first writes continuously from LBA1 to LBA8 and the second writes only hot data followed by order LBA6, LBA2, LBA4, LBA8. In the normal SSD, after the first sequence, the writes append in an empty block in order regardless of hot/cold data as shown in Fig. 1(b). After the second write sequence, new writes append in block 2 according to their order. The old LBAs are marked as invalid in block 0 and block 1 as illustrated in Fig. 1(c). In this case, if block 1 is a candidate for discarding in GC processing, there are overheads in the FTL that searching for a new empty block then copying back two valid LBAs, i.e., LBA5 and LBA7 to that new empty block before erasing the block 1.

In Multi-streamed SSD, however, writes are assigned to corresponding streams based on their frequency, i.e., hotness value. After the first write sequence, all hot LBAs are located in block 1, and all cold LBAs are located in block 0 as shown in Fig. 1(e). The second write sequence append data in block 2 as in regular SSD; however, the block 1 has all invalid pages as illustrated in Fig. 1(f). Erasing such kind of block is fast because the FTL update the mapping table and there is no cost for searching new empty block and copying valid data pages.

## 2.2 MongoDB and WiredTiger

MongoDB is a popular document store in NoSQL solutions that shares many characteristics with RDBMS i.e., secondary index supporting, transaction processing and concurrency controlling. The data model in MongoDB can be mapped to the one in RDBMS. *Database* concept is same for both models; *collections*, *documents*, and *key-value pairs* in MongoDB are mapped *tables*, *rows*, and *columns* in RDBMS respectively.

WiredTiger reinforces MongoDB as the default storage engine since the version 3.0. We research the internal block management of WiredTiger in detail to identify causes of data fragmentation. WiredTiger keeps the metadata of free pages, allocated pages, and invalid pages in a special page called checkpoint page. Only one checkpoint page is maintained in DRAM (i.e., live checkpoint), the other checkpoint pages are kept in non-volatile devices such as disks. The dirty pages are flushed from DRAM to disks by either page eviction from the buffer pool or checkpoint processing. WiredTiger neither adopts *update-in-place* approach as in traditional RDBMS nor *append-only* mode as in LSM-based DBMS. The space management first searching for a free block to write on, after the data page is successfully written, that block is marked as valid, the old version page is marked as invalid. During the checkpoint time, invalid pages are reclaimed and reused in the next checkpoint.

Fig. 2 illustrates how invalid data are reused at the checkpoint time. Suppose page B is flushed to disk through buffer pool's eviction process, the old version of page B is marked invalid accordingly. Other dirty data pages are kept in DRAM along with the live checkpoint, before flushing those pages to disk, the previous checkpoint pages are fetched from disk to DRAM (step 1) as shown in Fig. 2(a). Then it is merged (step 2) with the current live checkpoint, so invalid data of the same pages exist in the current version are reclaimed and can be reused again after the checkpoint process finished. After the merging phase, all dirty pages along with the live checkpoint are flushed to disk (step 3) as illustrated in Fig. 2(b).

This approach has two advantages: (1) Avoiding expensive compaction operations

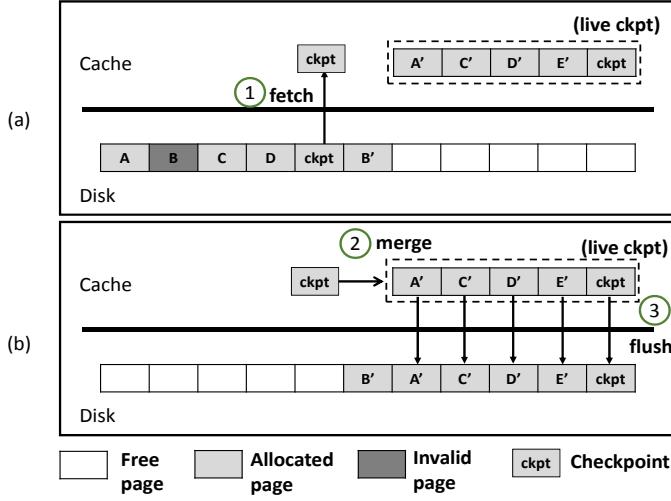


Fig. 2. Reusing data at checkpoint time in WiredTiger

that are popular in LSM-based approaches, and (2) old versions can serve as the backups used in the recovery process. However, that approach also leads to one drawback such that the valid data and invalid data of the same page are switched after each checkpoint, which forms an *internal fragmentation* in the SSD.

### 3. File-based Stream Mapping and Boundary-based Stream Mapping

We revisit the prior MSSD-based techniques in this section and argue the flaws in those methods by researching the IO patterns achieved from blktrace. Moreover, we define the requirements that an optimized MSSD method should satisfy.

#### 3.1 File-based Stream Mapping

**Table 1. The proportions of data written to file types with various of workloads [34]**

Benchmark	Operation ratio C:R:U:D	Colls	Pri. Indexes	2nd Indexes	Journal
Y-Update-Heavy	0:50:50:0	93.6	n/a	n/a	6.4
Y-Update-Only	0:0:100:0	89.6	n/a	n/a	10.4
LB-Original	12:69:15:4	58.6	3.1	37.22	1.08
LB-Mixed	12:0:84:4	66.1	0.5	31.13	2.27
LB-Update-Only	0:0:100:0	67.6	0.02	30.2	2.18

Typically, different file types in DBMS have different data accessed frequencies and write patterns. As shown in Table 1, we use workloads with different operations such as create, read, update, delete (CRUD) from YCSB [37] and Linkbench [35] as in the previous research [34]. For simple data model in YCSB, workload A (Y-Update-Heavy),

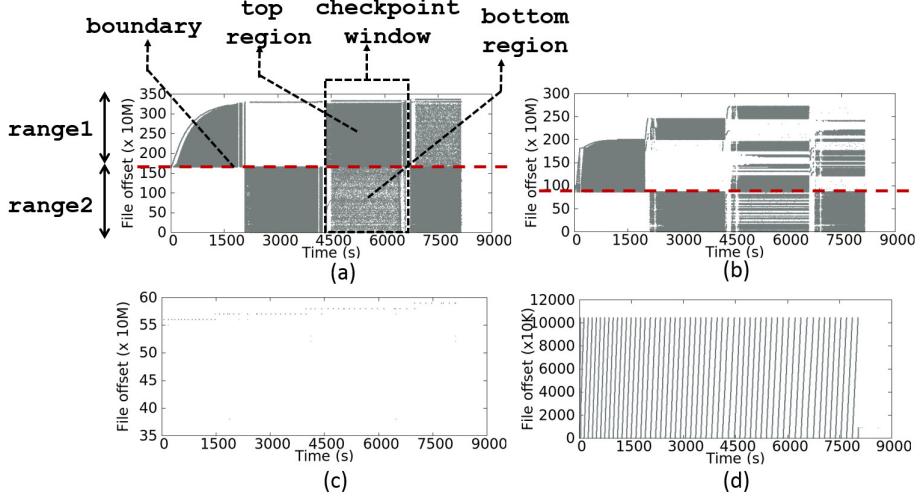


Fig. 3. Write patterns of various file types in WiredTiger with Linkbench benchmark, (a) Collection file,(b) secondary index , (c) primary index file, and (d) journal file

and only update workload (Y-Update-Only) are carried out. For complex data model in Linkbench, we use original Linkbench workload (LB-Original), mixed operations workload (LB-Mixed), and only update Linkbench workload (LB-Update-Only). Note that the table only includes information of most updated file types, the other minor updated file types, e.g., metadata files, system files, lock files are excluded.

Almost updates from workloads are carried on the collection files and secondary index files. The data model in YCSB has only one collection file and one primary index file, for workloads have only read and update operations, there is no write occurs on the primary index files. For MongoDB systems used SSD as the storage device, this asymmetric writes in file types forms the data fragmentation such that frequently written file types, i.e., collection and secondary index are considered of hot data, and the others file types are considered as cold data. *File-based* MSSD approach solves this problem by mapping each file type to a distinguish stream. This simple optimization gains moderate improvement of performance; however, there is another kind of data fragmentation called *internal data fragmentation* that is discussed detail in the next subsection.

### 3.2 Boundary-based Stream Mapping

To illustrate the internal data fragmentation, we use *blktrace* while running the LB-Update-Only workload in two hours to keep track of information of each write command on the underlying storage device. Fig. 3 shows written patterns of different file types in the system. The x-axis is the elapsed time in seconds, and the y-axis is the file offset. To avoid the effect of Operating System's cache, we enable *DirectIO*. There are two types of written patterns: (1) heavy randomly writes that occur on the collection file and secondary index file as shown in Fig. 3(a) and Fig. 3(b) respectively, and (2) sequential writes that occur in the primary index file and journal file as illustrated in Fig. 3(c) and Fig. 3(d) respectively. *Checkpoint window* is the period between a checkpoint and its very next checkpoint. *Write region (region in short)* is the area limited by two logical file offsets (on the y-axis) and two time-points (on the x-axis), usually is a checkpoint window.

At a given checkpoint window, the amount of data written to the *top* region and the

```

1: Require: boundary of each collection file and index file has computed
2: Input: file, and offset to write on
3: Output: sid - stream assign for this write
4: boundary  $\leftarrow$  getboundary(file)
5: if file is collection then
6:   if offset  $<$  boundary then
7:     sid  $\leftarrow$  COLL_SID1
8:   else
9:     sid  $\leftarrow$  COLL_SID2
10:  else if file is index then
11:    if offset  $<$  boundary then
12:      sid  $\leftarrow$  IDX_SID1
13:    else
14:      sid  $\leftarrow$  IDX_SID2
15:  else ▷ Other files i.e., metadata
16:    sid  $\leftarrow$  OTHER_SID

```

Fig. 4. Boundary-based stream mapping algorithm

*bottom* region are asymmetric and switches in next checkpoint window. For instance, in Fig. 3(a), after the first checkpoint, the bottom region has high frequency written than the top region. Then after the second checkpoint, the trend is reversed such that top region receives more writes compare with the bottom.

We revisit the *boundary-based* stream mapping approach in Fig. 4 [34]. The input is a file and file offset to write on, and the output is the mapped stream *sid*. The first step is retrieving the boundary of each collection file or index file which is the last file offset at the end of loading phase (line 4). Then from the rest of the algorithm, the corresponding stream is assigned to *sid* based on the type of input file, the input file offset is compared with the boundary to determine whether the write is to the top region or the bottom region. After stream id is mapped, the write command to the underlying file is given as *posix.fadvise(fid, offset, sid, advice)*, where *fid* is file identify, *offset* is the offset to write on, *sid* is stream id mapped and *advice* is passed as a predefined constant.

Unfortunately, empirical results show that boundary-based approach remains some flaws in a complicated data model, i.e., Linkbench rather than the simple one as in YCSB. Fig. 5 illustrates IO patterns of different file types using LB-Update-Only workload with y-axis is the file offset and x-axis is the elapsed time by second. We map writing occur on collection and index files to four streams which decoded with four different colors as shown in Fig. 5(a). Fig. 5(b) and Fig. 5(c) show a single-stream view of *stream1* and *stream2* that extracted from Fig. 5(a) which are bottom regions and top regions of collection files respectively. Inside one stream, however, regions of particular files have different written-lifetimes that lead to another unexpected overlapped written as shown in Fig. 5(d) and Fig. 5(e). We named that phenomenon as *cross-region fragmentation*.

In addition, as shown in Fig. 3 and Table 1, collection files and secondary index files have heavily randomly write patterns while primary indexes have scattered sequential write patterns. Therefore, the sequential writes should map to different stream from randomly writes. In summary, an optimal stream mapping scheme should satisfies below requirements:

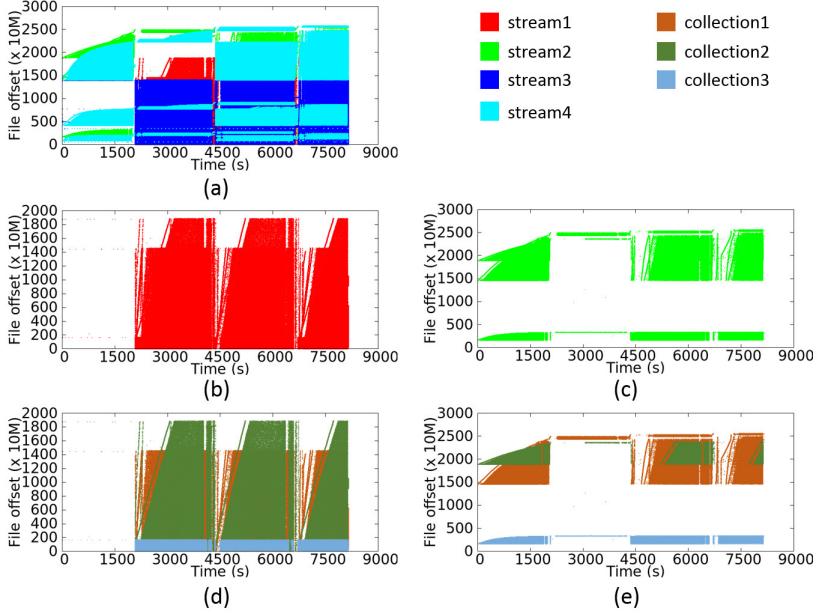


Fig. 5. Write patterns of different file types in WiredTiger with Linkbench: (a) Combined-all-stream write pattern; (b) and (c) are single-view write pattern and separately files write pattern respectively of bottom region of collection files; (d) and (e) are single-view write pattern and separately files write pattern respectively of bottom region of collection files

1. Solves *internal fragmentation* of collection files and index files.
2. Maps data written on primary index files to different stream from secondary index files, and collection files.
3. Solves *cross-region fragmentation*.
4. Works independently with various data models (number of collection files, primary files, and secondary index files) and the limitation of the number of streams supported by SSD.

*Boundary-based* approach satisfies the first and the last requirement, however, suffers from the requirement (2) and (3). The *file-based* approach that maps each file types to different streams solve only the second requirement. There is an extended solution from the file-based method that map each physical collection file or index file to a stream. This approach solves requirement (2) and (3) but becomes impossible when the number of files is larger than the maximum number of streams that the underlying SSD supported.

Due to those limitations of current approaches, we introduce a novel stream mapping scheme that satisfies all the requirements that discussed in detail in next section.

#### 4. Dynamic Stream Mapping

To solve the *cross-region fragmentation* problem in boundary-based approach, we propose a novel online stream mapping scheme that classifies all regions of collection files into  $K$  groups based on their hotness values, then assign each group to a distinguish

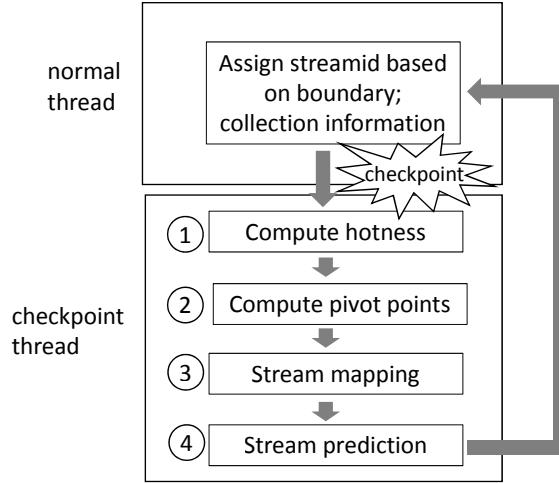


Fig. 6. Overall process of DSM approach

stream. We adopt similar stream mapping scheme for all regions of index files. Moreover, we separate writes on primary index files from writes on secondary index files by using statistical information. Note that it requires the underlying SSDs support at least 2K streams. The hotness value of one region changes after each checkpoint based on many aspects, e.g., current workload, data model, and cache size. Regardless of those aspects, our proposed approach dynamically classifies writes on regions into groups according to their hotness values and maps groups to corresponding streams. Therefore we named our proposed method as *Dynamic Stream Mapping (DSM)*. Note that, the boundary-based approach is a particular case of DSM where K is equal to two.

Fig. 6 illustrates the proposed idea in sequential steps. While the system is working at a normal thread, before writing data buffer to the storage system, based on the boundary of each file, we assign predicted streams for according regions, collect statistical information for each collection or index files. When issuing checkpoints, the system blocks normal threads until the checkpoint process finished. Our proposed method follows below steps:

- Step 1: Computes hotness values for two regions of each file based on statistical information collected previously in the normal thread.
- Step 2: For each collection file or index file, computes pivot points that aid to grouping regions according to their hotness values computed in step 1.
- Step 3: Maps streams to regions using hotness values and pivot points computed in step 1 and step 2.
- Step 4: For each region, predicts corresponding stream for the next checkpoint based on current mapped stream in step 3 and predicted-stream in the previous checkpoint.

After the remain works of checkpoint process finished, the system resumes the normal threads. We discuss the detail of each step in the below subsections.

#### 4.1 Stream assignment at the normal thread

```

1: Require:
    + boundary of each collection file and index file has computed
    + Predicted-streams for two regions of each file i.e., file.sid1 and file.sid2 are
        computed
2: Input: file, and offset to write on
3: Output: sid - stream assign for this write
4: boundary  $\leftarrow$  getboundary(file)
5: if file is collection OR file is index then
6:   if offset  $<$  boundary then
7:     sid  $\leftarrow$  file.sid1
8:     file.numw1  $\leftarrow$  file.numw1 + 1
9:     file.off_min1  $\leftarrow$  min(file.off_min1, offset)
10:    file.off_max1  $\leftarrow$  max(file.off_max1, offset)
11:   else
12:     sid  $\leftarrow$  file.sid2
13:     file.numw2  $\leftarrow$  file.numw2 + 1
14:     file.off_min2  $\leftarrow$  min(file.off_min2, offset)
15:     file.off_max2  $\leftarrow$  max(file.off_max2, offset)
16: else                                 $\triangleright$  Other files i.e., metadata
17:   sid  $\leftarrow$  OTHER_SID

```

Fig. 7. Dynamic stream mapping algorithm (at normal thread)

We assign streams during the system working time in a normal thread as in boundary-based method. The algorithm in Fig. 4 is rewritten in Fig. 7. The main differences are in line 7 and line 12 where the streams (*file.sid1* and *file.sid2*) assigned for the top region and the bottom region of each file is no longer constant but recomputed after each checkpoint. For the first checkpoint, the assigned streams are set by initial values. From the second checkpoint, assignment for *file.sid1* and *file.sid2* are done at the checkpoint time that discussed in subsection 4.4. Statistical information per each region (the number of writes, lower-bound offset, upper-bound offset of written files) are saved as illustrated at lines 7-9 and 13-15 in DSM algorithm in Fig. 7.

#### 4.2 Hotness Computing

In DSM approach, hotness value of a region is defined as the average number of writes occur per 4-KB data page on that region. Notice that hotness value of writing data in one region changed after each checkpoint, so the stream assignment should base on relative values rather than absolute values. To illustrate the idea, we use Fig. 3(a) which is the IO pattern of a collection file as discussed from the previous experiment. For a given collection file or index file, there is a *boundary* partition the file into two regions i.e., top and bottom that located by file offsets named *range1* and *range2* respectively. For a given region, the writing density is computed as in Eq.(1). *numwi* is the number of writes on that region within a checkpoint window, *PAGESIZE* is set as 4KB.

$$\text{density}_i = (\text{numwi}/\text{range}_i) \times \text{PAGESIZE}, \quad (1)$$

$$\text{hotness}_i = \lg(\text{density}_i / (t_2 - t_1)) \quad (2)$$

where  $i = 1, 2$  if writing on top region or bottom region respectively

Then, we compute the hotness value as in Eq.(2). Note that the checkpoint window is depended on: (1) type of workload, (2) ratio between cache size and database size, and (3) data model. Therefore, to eliminate those constraints, we divide the density by the interval time to get the writing density in a unit of time. Moreover, we adopt logarithmic scaling to solve the common sub-optimize problem that inspired from proportional selection phase in Genetic Algorithm [38]. Fig. A.2 in the Appendix describes the detail implementation of hotness computing.

### 4.3 Compute pivot points and stream mapping

After hotness values of all regions are computed, the next step is classifying all regions of one file type, e.g., collection into  $K$  groups. We partition a min–max range into  $K - 1$  evenly disjoint sections using  $K - 1$  pivot points  $p_i$  such that:  $\min \leq p_i \leq \max$ , where  $1 \leq i \leq K - 1$  and  $\min$  and  $\max$  are minimum hotness value and maximum hotness value of all regions respectively.

The pivot points are computed using Eq.(3). To increase the flexibility of the approach for various workload, we use positive-integer weight parameter  $\alpha$  by treating the first group and the last group (the "coldest" group) and the "hottest" group separately with the remains groups (warm groups). Note that in a particular case when  $\alpha$  equals to  $K$ , all groups have equal weight.

$$\begin{cases} p_1 = \min + (\max - \min) \times (1/\alpha) \\ p_{K-1} = \min + (\max - \min) \times ((\alpha - 1)/\alpha) \\ p_j = p_1 + ((p_{K-1} - p_1)/(k - 2)) \times (j - 1), \quad \text{where } 2 \leq j \leq K \end{cases} \quad (3)$$

Next, for a given region, we assign it to a corresponding group according to its hotness value as described in Eq.(4). Each group in turned mapped to a stream id

$$\text{group} = \begin{cases} \text{group}_1, & \min \leq \text{hotness} \leq p_1 \\ \text{group}_K, & p_{K-1} \leq \text{hotness} \leq \max \\ \text{group}_{g+1}, & p_g \leq \text{hotness} \leq p_{g+1}, 1 \leq g \leq K - 2 \end{cases} \quad (4)$$

The detail implementation for pivot points computing and stream mapping are described in Fig. A.3 and Fig. A.4 in the Appendix section respectively.

### 4.4 Stream prediction

Before writing to a region, we assign a stream to that writing using *posix\_fadvise(fd, offset, sid, advice)* command. In order words, we must predict the hot-cold trends for each region before the writes come. One region is hot in the current checkpoint may keep hot trend or reversely become cold in next checkpoint. For example, in Fig. 8 that plot write patterns with LB-Update-Only workload during four hours for *link* collection and *count* collection, the boundaries between two regions are marked as dash lines. In Fig. 8 (a), the hot-cold trend switches between regions in some first checkpoints then keep on the same trend after around two hours of executing time. On the other hand, in Fig. 8 (b), the hotness levels between two regions are almost similar during the execution time except for the first checkpoint. We propose a low overhead hot-cold data prediction as described in Eq.(5), where *cursid1*, *cursid2* are current streams assigned to corresponding groups in the previous step; *sid<sub>i</sub>* are stream finally map to each region that uses for the

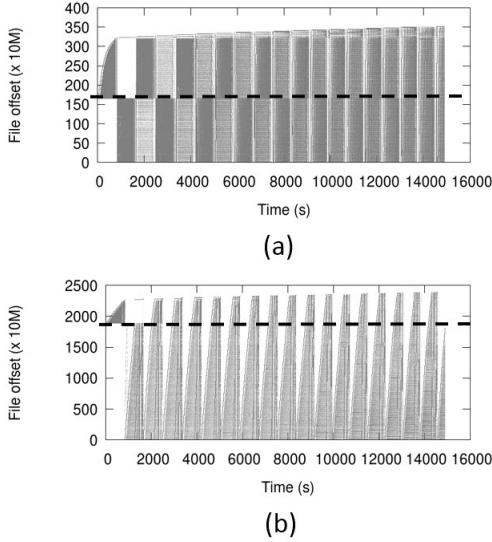


Fig. 8. The difference between hot-cold trends: (a) *link* collection file and (b) *count* collection file

next checkpoint. The detail of stream prediction is described in Fig. A.5 in the Appendix section.

$$sid_i = \begin{cases} cursid1, & sid_i == cursid1 \\ cursid2 & \text{otherwise} \end{cases} \quad (5)$$

where  $i = 1, 2$

Finally, we summaries all steps discussed until now in DSM algorithm in Fig. 9. Firstly, we use function *ProcessStat(F)* (see Fig. A.1 in the Appendix for detail) to extract supported information from the number of writes on each region and other statistical data (computed in the normal thread as described in the previous section). Then we compute hotness values and pivot points in line 8, 9, and 10 respectively. The *for* loop from line 11 to line 20 classifies two regions in each collection file or index file into corresponding groups. Due to primary indexes have different write patterns compare with the secondary indexes, we handle them separately by using the *globalpct1* and *globalpct2* of each file that are the percentage of writes on the bottom regions and the top regions respectively over total writes on all files during a checkpoint. Function *ProcessStat(F)* describes the detail of how to compute *globalpct1* and *globalpct2*. For a given index file, if there is a low-frequency write region (percentage of writes less than the THRESHOLD2) or low-frequency write on the whole file (total number of writes on two regions less than the THRESHOLD1) then that index file is considered as a primary index file. The last step predicts stream for each region of underlying file at line 21.

Our proposed DSM approach is low overhead and dynamically adapts to any workload, any data model. Firstly, for each collection file or index file, we use a data structure named *mssd-object* to capture the statistical information (i.e., number of writes, minimum and maximum write offset) for each region. Because the data model in NoSQL DBMS is flexibly changing at any time rather than fixed as in RDBMS, we allocate those *mssd-objects* dynamically based on the current number of files (i.e., number of collections or number of indexes). Moreover, during the normal thread, instead of updating statistical information for each data block as in previous studies, we update on region-based

```

1: Require: number of writes on bottom region  $numw1$  and bottom region  $numw2$   

   are computed for each file within the checkpoint  

2: Input:  $F$  - List of collection files and index files,  $K$  - number of groups  

3: Output: Stream ids of each collection or index file are mapped  

4:  $CP \leftarrow \emptyset$                                       $\triangleright$  Collection pivot points  

5:  $IP \leftarrow \emptyset$                                       $\triangleright$  Index pivot points  

6:  $cursid1 \leftarrow cursid2 \leftarrow nul$                    $\triangleright$  streams in current checkpoint  

7: ProcessStat(F)  

8: ComputeHotness(F, coll_min, coll_max, idx_min, idx_max)  

9: ComputePivots(CP, K - 1, coll_min, coll_max)  

10: ComputePivots(IP, K - 1, idx_min, idx_max)  

11: for each file  $f$  in  $F$  do            $\triangleright$  phase1: Compute the total writes in a stream id  

12:   if  $f$  is collection then  

13:      $cursid1 \leftarrow MapSIDByHotness(hotness1_f, CP, COLL_INIT_SID)$   

14:      $cursid2 \leftarrow MapSIDByHotness(hotness2_f, CP, COLL_INIT_SID)$   

15:   else                                          $\triangleright$  Index files  

16:     if  $(globalpct1_f + globalpct2_f \leq THRESHOLD1) \text{ OR}$   

         $(globalpct1_f \leq THRESHOLD2) \text{ OR}$   

         $(globalpct1_f \leq THRESHOLD2)$  then           $\triangleright$  Primary index  

17:        $cursid2 \leftarrow cursid1 \leftarrow PRIMARY_IDX_SID$   

18:     else  

19:        $cursid1 \leftarrow MapSIDByHotness(hotness1_f, IP, IDX_INIT_SID)$   

20:        $cursid2 \leftarrow MapSIDByHotness(hotness2_f, IP, IDX_INIT_SID)$   

21:   PredictStream(f, cursid1, cursid2)

```

Fig. 9. Dynamic Stream Mapping algorithm(at checkpoint thread)

*mssd-objects* that have small memory footprint (i.e., lower than 100B). Moreover, hotness values and stream mapping are based on relative value (i.e., proportion) rather than absolute value. Thus our proposed method work dependently with workloads.

## 5. Evaluation and Analysis

This section describes experiment settings, evaluation results of our proposed method compared to the original WiredTiger (as the base line) and prior methods. We also analyze the effectiveness of distinguishing writes on primary index files from writes on secondary index files as well as the effectiveness of reducing the maximum leaf page size of collection file as done in the prior research.

### 5.1 Experimental Settings

To fairly compare our proposed method with file-based approach and boundary-based approach, we adopt the same experimental setup with the previous research [34]. To enable multi-streamed SSD technique, we use both modified Linux kernel 3.13.11 and customized-firmware Samsung 840 Pro SSD as in [4]. For eliminating network latency, we set up both the client and the server in the same commodity computer with 48 cores

Intel Xeon 2.2GHz processor, 32GB DRAM. In the client layer, we use YCSB 0.5.0<sup>1</sup> and LinbenchX 0.1<sup>2</sup> (an extended version of Linkbench that support MongoDB) with diversity workloads as shown in Table 1. The number of documents in YCSB is set to 23 million, and *maxid1* in Linkbench is set to 80 million. All benchmarks are executed during two hours with 40 client threads. In the server layer, we use a stand-alone MongoDB 3.2.1<sup>3</sup> server with DirectIO and various cache sizes from 5GB to 30GB. WiredTiger is used as the storage engine with all default settings.

## 5.2 Multi-streamed SSD optimization evaluation

**Table 2. Stream mapping schemes in multi-streamed approaches included the original WiredTiger**

Method	Kernel	Metadata	Journal	Primary Index	Collection	2nd Index
<b>Original</b>	0	0	0	0	0	0
<b>File-based</b>	0	0	1	3	2	3
<b>Boundary-based</b>	0	0	1	4,5	2,3	4,5
<b>DSM</b>	0	0	1	1	2,3,4	5,6,7

Table 2 summarizes mapping schemes of all methods. There is no stream mapping in the original WiredTiger, so all writes are mapped to the default stream 0 (reserved for files in the Linux kernel). In file-based approach, each file type is mapped to a distinguish stream. In the case of boundary-based approach, all collection files are mapped to two streams: one for top regions and another for bottom regions. The same mapping scheme is adopted for index files. In DSM approach, we use there-group mapping DSM i.e., set K equal to 3 with three streams for collection files and three streams for secondary index files. Writes on primary index files are mapped to a distinguish stream to writes on secondary index files. There are some important notes:

- Except for the DSM method, the remains map writes on primary index files and secondary index files to the same stream.
- The stream mapping for metadata files and journal files are same for all methods.
- The Samsung 840 Pro SSD support maximum only eight streams from 0 to 7. In DSM approach six streams have already mapped for collection files and secondary index files. Therefore, writes on journal files and on primary index files share the same stream (i.e., stream1). It is adequate because writes on journal files and primary files follow lightly sequential patterns, hence can be considered as cold data and can be mapped in the same stream.

Fig. 10 and Fig. 11 show the performance results in term of throughput and the 99th-percentile latency improvement respectively. In particular, for the LB-Original workload, because the ratio of read operations are high (i.e., 69%) and the data size is approximate 32GB. Pages tend to fetched in (flush out) to (from) the buffer pool with high frequency, we keep the reasonable large buffer sizes (20GB, 25GB, and 30GB) to keep secondary indexes in DRAM as long as possible and avoid the performance degradation due to lacking

<sup>1</sup><https://github.com/brianfrankcooper/YCSB/releases/tag/0.5.0>

<sup>2</sup><https://github.com/Percona-Lab/linkbenchX>

<sup>3</sup><https://github.com/mongodb/mongo/archive/r3.2.1.tar.gz>

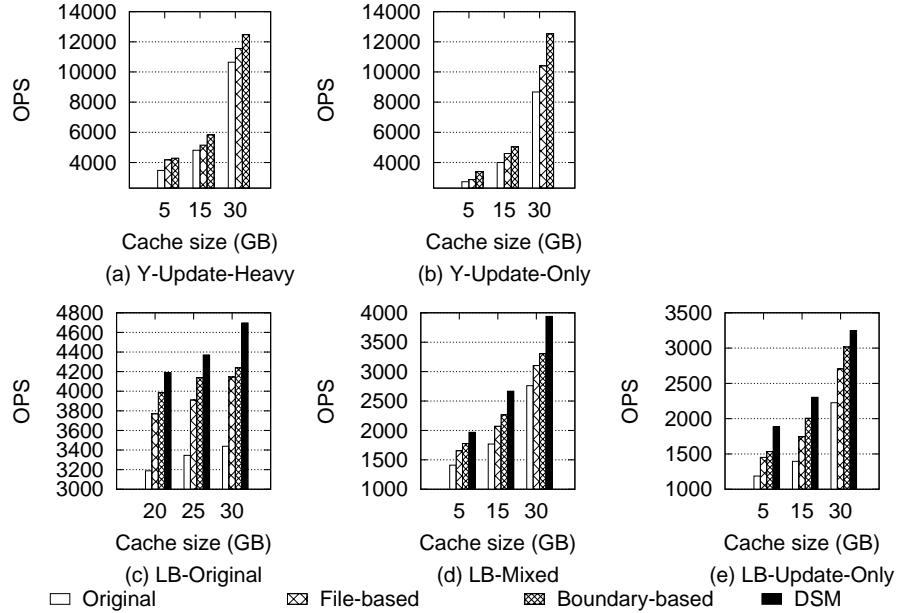


Fig. 10. Throughput of optimized methods compare with the original

space in the buffer pool. Note that DSM approach is not carried out in YCSB benchmark due to YCSB's simple data model such that there are a collection file and a primary index file, DSM in this benchmark is almost similar to boundary-based method.

Empirical results have shown boundary-based approach lost its effectiveness in solving internal fragmentation problem when the data model becomes complex. In YCSB, the boundary-based has additional 2.8%–13.1% and 9.73%–20.4% of throughput improved rate compare to the file-based method. However, those gaps become smaller in Linkbench that just 2.2%–5.8%, 6.5%–9.4% and 5.6%–15% for LB-Original, LB-Mixed, and LB-Update-Only respectively. There is the same trend with 99th-percentile latency. Compare to the file-based approach, the boundary-based has additional 3.6%–4.6% and 12.7%–20.9% of improvement rate for Y-Update-Heavy and Y-Update-Only respectively. With heavily-read workload (i.e., LB-Original), the additional improved rate is similar or even better than the Y-Update-Heavy with 3.5%–9.8%. With heavier write workloads (i.e., LB-Mixed and LB-Update-Only), the additional rates are 2.1%–3.9%, and 2.5%–10% respectively, that are significantly lower than the Y-Update-Only. The reason is with complex data model in Linkbench, cross-region phenomenon is dominant thus the boundary-based approach becomes less effective.

Conversely, compared to the best-performing method in prior work (i.e., boundary-based), DSM approach effectively solves the internal fragmentation and cross-region fragmentation. DSM improves the throughput by up to 10.8%, 19.2%, and 23%; enhances the 99th-percentile latency by up to 8.9%, 21.2%, and 28.5% for LB-Original, LB-Mixed, and LB-Update-Only respectively. Overall, DSM is the best method, compared to the original WiredTiger, DSM has throughput improve by up to 36.5%, 50.7%, and 65%; 99th-percentile latency grow by up to 27.6%, 31.8%, and 46.2% for LB-Original, LB-Mixed, and LB-Update-Only respectively.

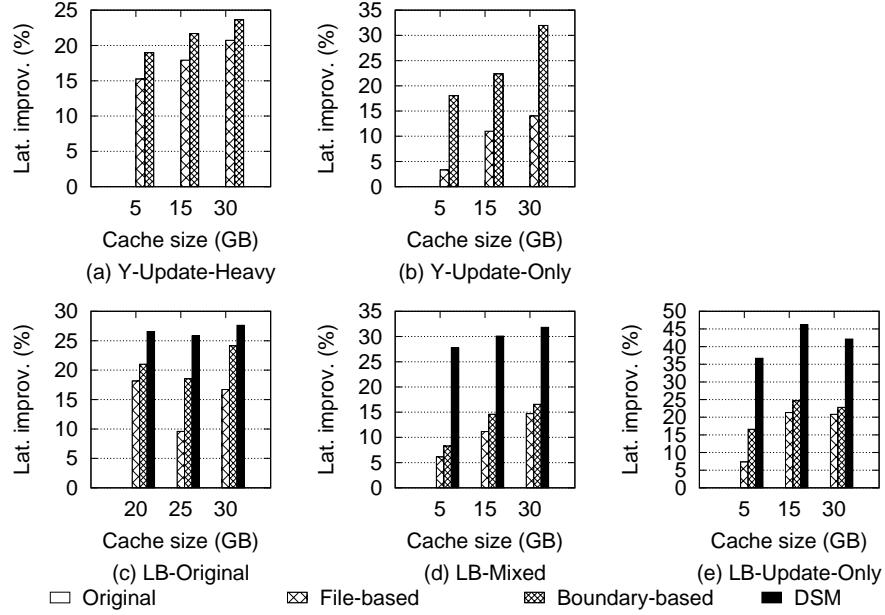


Fig. 11. Latency of optimized methods compare with the original

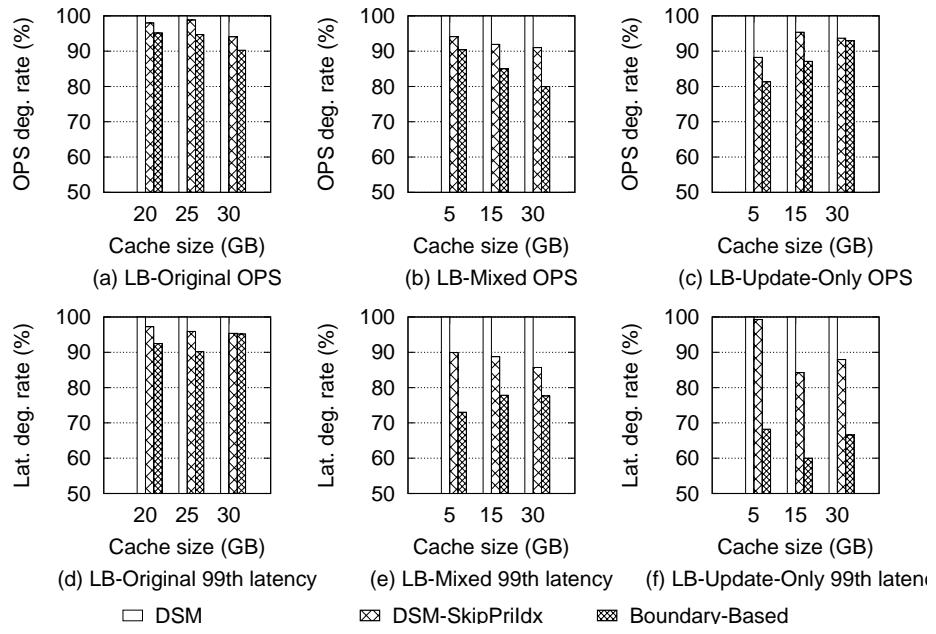


Fig. 12. Effects of primary index filtering in term of throughput and 99th latency

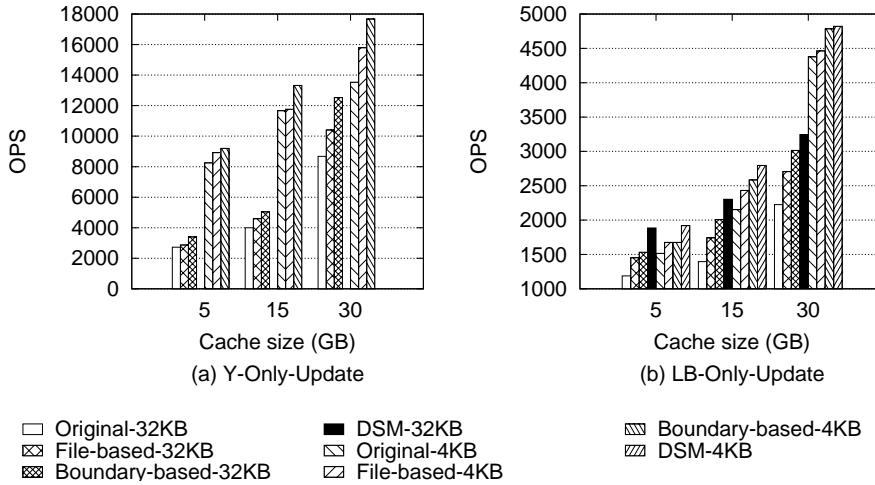


Fig. 13. Leaf page size optimization evaluation of methods with LB-Update-Only

### 5.3 Effects of primary index filtering

DSM not only solves cross-region fragmentation but also distinguish writes on primary index files from secondary index files. Fig. 12 shows the effects of primary index filtering by using typical DSM and the DSM without primary index filtering, i.e., *DSM-SkipPriIdx*. The y-axis shows the degradation rate of throughput and latency of a method compare with the based line method (i.e., DSM), the boundary-based results are included as the reference. Typically, heavier write workloads have more effective in primary index filtering. Without primary indexing filtering, the throughput benefit decreases by up to 5.9%, 9%, 11.7%, and the 99th latency benefit lost by up to 4.6%, 14.2%, 15.7% for LB-Original, LB-Mixed, and LB-Update-Only respectively. Even though without primary index filtering, DSM still shows better performance than boundary-based both on throughput and 99th latency.

### 5.4 Leaf page size optimization evaluation

In this subsection, we further optimize MongoDB by reducing the maximum leaf page size of collection files (the leaf page sizes of index files are unchanged). To evaluate the impact of leaf page size to the performance of our proposed method, we adopt the same experiment with the prior research. We setup YCSB benchmark and Linkbench benchmark using various workloads and cache sizes for 32KB leaf page size (default) and 4KB leaf page size. For the space limitation in the paper, we only show the throughput results of the heaviest write workloads, i.e., Y-Update-Only and LB-Update-Only as in Fig. 13. Overall, DSM-4KB still is the best method that improves throughput 1.6x–2.1x, and 1.8%–48.6% compared to the Original-32KB, and DSM-32KB respectively. Small leaf page size optimization is effective with simple data model in YCSB; however, lost its effectiveness in the complex data model in Linkbench. For instance, in Y-Only-Update, the boundary-based method shows significantly improving the throughput performance

by up to 2x–3.3x, but only improves 1.4x–2.1x in LB-Only-Update.

### 5.5 Discussion

The proposed method is based on the statistic information i.e., write regions, write frequencies collected from files during a checkpoint period. Besides, distinguishing data types (primary index and secondary indexes) from one large data file also aids to improve the performance. Thus, any storage engine which has various data objects contain inside one large data file and those data objects have different access lifetimes can gain benefits from DSM. For instance, InnoDB storage engine in MySQL uses one large per-table user tablespace for both primary index and secondary indexes and one large system tablespace file for many data types e.g., metadata, double write buffer, rollback segments.

There is a trade-off between the performance gain and buffer pool usage when reducing the leaf page size of collection files. In the original WiredTiger, the system must fetch the whole 32KB-page from the storage device to the buffer pool in order to read a some-bytes record. The buffer pool size is usually fixed, with the smaller page size, there is higher data pages contain the queried records in DRAM, thus reduce the number of IO accesses. On the other hand, with the same benchmark parameter (number of documents in YCSB or maxid1 in Linkbench), reducing the leaf page size from 32KB to 4KB lead to increasing number of leaf pages and number of internal pages in the B+Tree. Therefore, collection files and index files become larger and require additional spaces from the buffer pool to keep the hot index files in DRAM. For instance, reducing the leaf page size of collection files from 32KB to 4KB lead to the database size increases from 51.7GB to 58.4GB (+12.8%) in YCSB, but it grows from 55GB to 108.5GB (approximate double) in Linkbench.

## 6. Conclusion

In this paper, we have discussed data fragmentation in MongoDB as well as the proposed methods in detail. The file-based method is the simplest one that solves the data fragmentation due to the different lifetime of writes on file types but remains internal fragmentation caused by asymmetric regions writing. For simple data model in YCSB, the boundary-based approach is adequate to solve the internal fragmentation that shows good performance improvement in YCSB but struggling with complex data model in Linkbench due to cross-region fragmentation. To address that challenge, we extended boundary-based method by introducing DSM, a novel low-overhead stream mapping scheme that dynamically grouping writes on corresponding streams based on hotness values in each checkpoint period. DSM works independently with data models and workloads regardless of the limitation of the number of streams that the physical SSD supported. The number of groups and other parameters are configurable to gain the best performance. Stream mapping using primary index filtering in DSM has considerable performance improvement. Moreover, simple data model in YCSB gains more benefits from decreasing B+tree leaf page size than complex data model in Linkbench. In practical applications, the data models are complex with many collection files, index files rather than simply as in YCSB. Our proposed method is adequate for such applications thus it works effectively regardless of data models or workloads.

In the next research, we plan to evaluate the performance of our proposed method with emerging multi-streamed SSD devices i.e., NVMe SSD. We also optimize the algorithms for distributed environment regard to replica sets and shards.

## ACKNOWLEDGMENT

This research was supported by the MSIP (Ministry of Science, ICT and Future Planning), Korea, under the "SW Starlab" (IITP-2015-0-00314) supervised by the IITP (Institute for Information & communications Technology Promotion).

## REFERENCES

1. Y. Deng and J. Zhou, "Architectures and optimization methods of flash memory based storage systems," *Journal of Systems Architecture*, Vol. 57, no. 2, 2011, pp. 214–227.
2. J.-U. Kang, J.-S. Kim, C. Park, H. Park, and J. Lee, "A multi-channel architecture for high-performance nand flash-based storage system," *Journal of Systems Architecture*, Vol. 53, no. 9, 2007, pp. 644–658.
3. S.-W. Lee, B. Moon, C. Park, J.-M. Kim, and S.-W. Kim, "A case for flash memory ssd in enterprise database applications," in *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*. ACM, 2008, pp. 1075–1086.
4. F. Yang, K. Dou, S. Chen, M. Hou, J.-U. Kang, and S. Cho, "Optimizing nosql db on flash: A case study of rocksdb," in *2015 IEEE 12th Intl Conf on Ubiquitous Intelligence and Computing and 2015 IEEE 12th Intl Conf on Autonomic and Trusted Computing and 2015 IEEE 15th Intl Conf on Scalable Computing and Communications and Its Associated Workshops (UIC-ATC-ScalCom)*. IEEE, 2015, pp. 1062–1069.
5. J.-U. Kang, J. Hyun, H. Maeng, and S. Cho, "The multi-streamed solid-state drive," in *6th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 14)*, 2014.
6. S.-W. Lee and B. Moon, "Design of flash-based dbms: an in-page logging approach," in *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*. ACM, 2007, pp. 55–66.
7. J. Kim, D. H. Kang, B. Ha, H. Cho, and Y. I. Eom, "Mast: Multi-level associated sector translation for nand flash memory-based storage system," *Computer Science and its Applications*. Springer, 2015, pp. 817–822.
8. T. Jung, Y. Lee, J. Woo, and I. Shin, "Double hot/cold clustering for solid state drives," *Advances in Computer Science and its Applications*. Springer, 2014, pp. 141–146.
9. D. Park and D. H. Du, "Hot and cold data identification for flash memory using multiple bloom filters," in *USENIX Conference on File and Storage Technologies (FAST)*, 2011.
10. J.-W. Hsieh, T.-W. Kuo, and L.-P. Chang, "Efficient identification of hot data for flash memory storage systems," *ACM Transactions on Storage (TOS)*, Vol. 2, no. 1, 2006, pp. 22–40.
11. C. Min, K. Kim, H. Cho, S.-W. Lee, and Y. I. Eom, "Sfs: random write considered harmful in solid state drives." in *FAST*, 2012, p. 12.
12. S.-H. Park, J.-W. Park, S.-D. Kim, and C. C. Weems, "A pattern adaptive nand flash memory storage structure," *IEEE transactions on computers*, Vol. 61, no. 1, 2012, pp. 134–138.
13. T. I. Damaiyanti, A. Imawan, F. I. Indikawati, Y.-H. Choi, and J. Kwon, "A similarity query system for road traffic data based on a nosql document store," *Journal of Systems and Software*, Vol. 127, 2017, pp. 28–51.

14. H. Shim, “Phash: A memory-efficient, high-performance key-value store for large-scale data-intensive applications,” *Journal of Systems and Software*, Vol. 123, 2017, pp. 33–44.
15. Y.-T. Liao, J. Zhou, C.-H. Lu, S.-C. Chen, C.-H. Hsu, W. Chen, M.-F. Jiang, and Y.-C. Chung, “Data adapter for querying and transformation between sql and nosql database,” *Future Generation Computer Systems*, Vol. 65, 2016, pp. 111–121.
16. D. G. Chandra, “Base analysis of nosql database,” *Future Generation Computer Systems*, Vol. 52, 2015, pp. 13–21.
17. R. Dharavath and C. Kumar, “A scalable generic transaction model scenario for distributed nosql databases,” *Journal of Systems and Software*, Vol. 101, 2015, pp. 43–58.
18. P. O’Neil, E. Cheng, D. Gawlick, and E. O’Neil, “The log-structured merge-tree (lsm-tree),” *Acta Informatica*, Vol. 33, no. 4, 1996, pp. 351–385.
19. B. Jun and D. Shin, “Workload-aware budget compensation scheduling for nvme solid state drives,” in *Non-Volatile Memory System and Applications Symposium (NVMSA), 2015 IEEE*. IEEE, 2015, pp. 1–6.
20. MongoDB, “Mongodb architecture,” <https://www.mongodb.com/mongodb-architecture>, Aug. 2017.
21. WiredTiger, “Wiredtiger, making big data roar,” <http://www.wiredtiger.com/>, Aug. 2017.
22. C. Győrödi, R. Győrödi, G. Pecherle, and A. Olah, “A comparative study: Mongodb vs. mysql,” in *Engineering of Modern Electric Systems (EMES), 2015 13th International Conference on*. IEEE, 2015, pp. 1–6.
23. S. H. Aboutorabi, M. Rezapour, M. Moradi, and N. Ghadiri, “Performance evaluation of sql and mongodb databases for big e-commerce data,” in *Computer Science and Software Engineering (CSSE), 2015 International Symposium on*. IEEE, 2015, pp. 1–7.
24. C. Băzăr, C. S. Iosif *et al.*, “The transition from rdbms to nosql. a comparative analysis of three popular non-relational solutions: Cassandra, mongodb and couchbase,” *Database Syst J*, Vol. 5, no. 2, 2014, pp. 49–59.
25. B. Alexandru, R. Florin, and I. A. Laura, “Mongodb vs. oracle-database comparison,” in *Third International Conference on Emerging Intelligent Data and Web Technologies, IEEE*, 2012, pp. 330–335.
26. C.-H. Lee and Y.-L. Zheng, “Sql-to-nosql schema denormalization and migration: A study on content management systems,” in *Systems, Man, and Cybernetics (SMC), 2015 IEEE International Conference on*. IEEE, 2015, pp. 2022–2026.
27. G. Zhao, Q. Lin, L. Li, and Z. Li, “Schema conversion model of sql database to nosql,” in *P2P, Parallel, Grid, Cloud and Internet Computing (3PGCIC), 2014 Ninth International Conference on*. IEEE, 2014, pp. 355–362.
28. A. Kanade, A. Gopal, and S. Kanade, “A study of normalization and embedding in mongodb,” in *Advance Computing Conference (IACC), 2014 IEEE International*. IEEE, 2014, pp. 416–421.
29. G. Zhao, W. Huang, S. Liang, and Y. Tang, “Modeling mongodb with relational model,” in *Emerging Intelligent Data and Web Technologies (EIDWT), 2013 Fourth International Conference on*. IEEE, 2013, pp. 115–121.
30. X. Wang, H. Chen, and Z. Wang, “Research on improvement of dynamic load balancing in mongodb,” in *Dependable, Autonomic and Secure Computing (DASC), 2013 IEEE 11th International Conference on*. IEEE, 2013, pp. 124–130.

31. Y. Liu, Y. Wang, and Y. Jin, “Research on the improvement of mongodb auto-sharding in cloud environment,” in *Computer Science & Education (ICCSE), 2012 7th International Conference on*. IEEE, 2012, pp. 851–854.
32. P. Murugesan and I. Ray, “Audit log management in mongodb,” in *Services (SERVICES), 2014 IEEE World Congress on*. IEEE, 2014, pp. 53–57.
33. T.-D. Nguyen and S.-W. Lee, “I/o characteristics of mongodb and trim-based optimization in flash ssds,” in *Proceedings of the Sixth International Conference on Emerging Databases: Technologies, Applications, and Theory*. ACM, 2016, pp. 139–144.
34. T.-D. Nguyen and S.-W. Lee, “Optimizing mongodb using multi-streamed ssd,” in *To appear in Proceedings of the Seventh International Conference on Emerging Databases: Technologies, Applications, and Theory*. Springer, 2017.
35. T. G. Armstrong, V. Ponnekanti, D. Borthakur, and M. Callaghan, “Linkbench: a database benchmark based on the facebook social graph,” in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*. ACM, 2013, pp. 1185–1196.
36. Samsung, “Multi-stream technology,” <http://www.samsung.com/semiconductor/insights/article/25465/multistream>, Aug. 2016.
37. 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*. ACM, 2010, pp. 143–154.
38. D. E. Goldberg and K. Deb, “A comparative analysis of selection schemes used in genetic algorithms,” *Foundations of genetic algorithms*, Vol. 1, 1991, pp. 69–93.

**Trong-Dat Nguyen** received the MS degree from the School of Computer Science and Engineering, Kyungpook National University, Korea, in 2014. He is currently working toward the Ph.D degree at Sungkyunkwan University, Suwon, Korea. His research interests include NoSQL DBMSs, flash-based database technology.





**Sang-Won Lee** received the Ph.D degree from the Computer Science Department, Seoul National University, Korea, in 1999. He is an associate professor with the College of Information & Communication Engineering, Sungkyunkwan University, Suwon, Korea. He was a research professor at Ewha Womans University and a technical staff at Oracle, Korea. His research interest includes flash-based database technology.

## APPENDIX

Fig. A.1 describes the algorithm of processing the statistical information. The total number of writes for each region is computed as in the first *for* loop in line 4 – 10. For a given file,  $numw1_f$  and  $numw2_f$  are numbers of writes in *region1* and *region2* of that file respectively. In the second *for* loop, the global percentage of each region ( $globalpct1_f$  or  $globalpct2_f$ ) is computed which is used for distinguishing primary index from secondary index.

```

1: function PROCESSSTAT( $F$ )            $\triangleright$  Processing statistical data for each file i.e.,  

   hotness value, global percentage of write  

2:    $coll\_count2 \leftarrow coll\_count1 \leftarrow 0$   

3:    $idx\_count2 \leftarrow idx\_count1 \leftarrow 0$   

4:   for each file  $f$  in  $F$  do            $\triangleright$  phase1: Compute the total writes  

5:     if  $f$  is collection then  

6:        $coll\_count1 \leftarrow coll\_count1 + numw1_f$   

7:        $coll\_count2 \leftarrow coll\_count2 + numw2_f$   

8:     else  

9:        $idx\_count1 \leftarrow idx\_count1 + numw1_f$   

10:       $idx\_count2 \leftarrow idx\_count2 + numw2_f$   

11:   for each file  $f$  in  $F$  do           $\triangleright$  phase2: Compute the hotness  

12:     if  $f$  is collection then  

13:        $globalpct1_f \leftarrow numw1_f / coll\_count1 * 100$   

14:        $globalpct2_f \leftarrow numw2_f / coll\_count2 * 100$   

15:     else  

16:        $globalpct1_f \leftarrow numw1_f / idx\_count1 * 100$   

17:        $globalpct2_f \leftarrow numw2_f / idx\_count2 * 100$ 
```

Fig. A.1. Algorithm of processing statistical information

Fig. A.2 is the detailed implementation of computing density value and hotness value as described in Eq.(1) and Eq.(2) respectively.

Fig. A.3 is the detailed implementation of Eq.(3). Remind that alpha is the positive integer. The first and the last pivot point are computed as in line 2 and line 3 respectively. Then other pivot points are computed in the *for* loop in line 5 – 6.

We describe the detailed implementation of Eq.(4) in Fig. A.4. For a given region, we find the pair of pivot points such that  $pivots_j \leq hotness \leq pivots_{j+1}$ , then assign the

```

1: function COMPUTEHOTNESS( $F, coll\_min, coll\_max, idx\_min, idx\_max$ )           ▷
   Processing statistical data for each file i.e., hotness value, global percentage of
   write
2:   for each file  $f$  in  $F$  do                                              ▷ phase2: Compute the hotness
3:     //Number of 4KB page writes on each range in a unit of time
4:      $density1 \leftarrow (numw1_f * 4096) / range1_f / (t2 - t1)$ 
5:      $density2 \leftarrow (numw2_f * 4096) / range2_f / (t2 - t1)$ 
6:      $hotness1_f \leftarrow \lg(density1/t)$ 
7:      $hotness2_f \leftarrow \lg(density2/t)$ 
8:     if  $f$  is collection then
9:        $coll\_min \leftarrow \min(coll\_min, hotness1_f, hotness2_f)$ 
10:       $coll\_max \leftarrow \max(coll\_max, hotness1_f, hotness2_f)$ 
11:    else
12:       $idx\_min \leftarrow \min(idx\_min, hotness1_f, hotness2_f)$ 
13:       $idx\_max \leftarrow \max(idx\_max, hotness1_f, hotness2_f)$ 

```

Fig. A.2. Algorithm of computing hotness

```

1: function COMPUTEPIVOTS( $pivots, n, min, max$ )
2:    $pivots[0] \leftarrow min + (max - min) / \alpha$ 
3:    $pivots[n - 1] \leftarrow min + (max - min) * (\alpha - 1) / \alpha$ 
4:    $step \leftarrow (pivots[n - 1] - pivots[0]) / (n - 1)$ 
5:   for  $i \leftarrow 1$  to  $(n - 2)$  do
6:      $pivots[i] \leftarrow pivots[0] + step * i$ 

```

Fig. A.3. Algorithm of computing pivot points

corresponding stream  $sid$  to that region in line 3, 5, and 9.

Fig. A.5 describes the detailed implementation of Eq.(5). For a given file, if the current computed stream e.g.,  $cursid1$  is same with the predicted stream in the previous checkpoint ( $sid1_{file}$ ), it means the hot-cold trends is unchanged then the predicted stream for next checkpoint is kept same as before. Otherwise, we map that region to the other stream ( $cursid2$ ).

```

1: function MAPSIDBYHOTNESS(hotness, pivots, n, initsid)
2:   if hotness  $\leq$  pivots[0] then                                 $\triangleright$  the most left
3:     sid  $\leftarrow$  initsid
4:   else if pivots[n - 1]  $\leq$  hotness then                 $\triangleright$  the most right
5:     sid  $\leftarrow$  initsid + n
6:   else
7:     Find the pivot point j in array pivots such that:
8:     pivots[j]  $\leq$  hotness  $\leq$  pivots[j + 1]
9:     sid  $\leftarrow$  initsid + (j + 1)
10:  return sid

```

Fig. A.4. Algorithm of mapping SID by hotness

```

1: function PREDICTSTREAM(file, cursid1, cursid2)
2:   if sid1file == cursid1 then       $\triangleright$  the hot-cold trend is same, do not swap
3:     sid1file  $\leftarrow$  cursid1
4:   else
5:     sid1file  $\leftarrow$  cursid2
6:   if sid2file == cursid2 then       $\triangleright$  the hot-cold trend is same, do not swap
7:     sid2file  $\leftarrow$  cursid2
8:   else
9:     sid2file  $\leftarrow$  cursid1

```

Fig. A.5. Algorithm of stream prediction