

Advocating for Key-Value Stores with Workload Pattern Aware Dynamic Compaction

<u>Heejin Yoon</u>, Jin Yang, Juyoung Bang, Sam H. Noh^{*}, Young-ri Choi

UNIST, *Virginia Tech





Contents

- Background and Motivation
- DOPA-DB
- Evaluation
- Discussion & Future works
- Conclusion



Background: LSM-based KVS

- Log-Structured Merge-tree based KV store
 - High write performance due to sequential write
 - Used for various data-intensive applications









Flush of LSM-based KVS

• Data is saved to disk in the form of file through Flush





Compaction of LSM-based KVS

• Disk-level data is sorted by Compaction





Dynamic workload in real world

• In real life, workload pattern regarding ratio of write and read operations changes over time



The QPS variation in a 14-day time span of UDB [1]



6

Dynamic workload in real world

• In real life, workload pattern regarding ratio of write and read operations changes over time





7

Write stall on write intensive workload

- In write intensive workloads, write stalls may occur
 - L_0 stall \rightarrow most crucial!
 - MemTable Stall
 - Pending stall





Visualize L₀ status with fill-to-threshold ratio

- Fill-to-threshold ratio = $\frac{L_0 \ size}{L_0 \ stall \ threshold}$
 - L₀ stall threshold : Size that triggers L₀ stall
 - Fill-to-threshold ratio > 1 \rightarrow L_0 stall triggered





L₀ stall on write intensive workload

- Large upper-level compaction causes L_0 stall
 - L₀ compaction cannot be triggered on time



< L₀ status on write intensive workload>



Compaction strategy for write intensive workload

Adjusting upper-level compaction (ULC) size could avoid L_0 stall •



• With ULC latency, next L_0 compaction latency must be considered



Multiple disk access on read intensive workload

• For read, LSM KVS has to check multiple components





No L₀ stall on read intensive workload

- On read intensive workload, L_0 stall is rare
 - No flush on read-only workload



 $< L_0$ status on read intensive workload>



Compaction strategy for read intensive workload

- Upper-level compaction size can be very large
 - Helpful in reducing the number of levels to access







Goal of dynamic compaction





Compaction size depending on workload pattern







Compaction Size Recommender :

Dynamic compaction size recommendation for dynamic workload

+

Tiered LSM-tree based KV store with multiple key ranges



Dynamic wOrkload Pattern Aware LSM-Based KV store



DOPA-DB based on tiered LSM KVS

- Each level, except for L_0 and last level, consists of p sub-levels
- Tiered LSM KVS avoids compaction at neighboring levels except for the last level



DOPA-DB with multiple key ranges

- To support dynamic compaction strategy, range is used for compaction unit
 - Not overlapped
 - Uniformly sized





Tiered LSM KVS with multiple key ranges for DOPA-DB

- Every level has multiple key ranges
 - Number of range grows based on range amplification factor



Tiered compaction with multiple key range

- Tiered compaction with multiple key ranges is performed based on sub-levels
 - Number of compaction units is decided by CSR



...



Placement after tiered compaction

• Placement of compaction output maintains non overlapped status of each sub level



...



DOPA-DB with CSR

• Compaction Size Recommender(CSR) suggests size of upper-level compaction





How CSR calculates ULC size

- To calculate ULC (upper-level compaction) size, values of parameters are estimated
 - Monitor flush & compaction speed during fixed time window (5 sec)
 - Remaining L₀ size
 - Estimated next L₀ compaction size



Total L_0 file size: 400MBRemaining L_0 size : 400MBEstimated next L_0 compaction size = 200MB

Flush speed: 80MB/s Compaction speed: 100MB/s



- Step 1. Calculate latency for ULC (upper-level compaction)
 - A = Latency until L_0 stall threshold reached = $\frac{Remaining L_0 \ size}{Flush \ speed}$





- Step 1. Calculate latency for ULC (upper-level compaction)
 - B = Next L_0 compaction latency = $\frac{Estimated next L_0 compaction size}{Compaction speed}$



- Step 1. Calculate latency for ULC (upper-level compaction)
 - Latency until L₀ stall threshold reached(A) Next L₀ compaction latency(B)
 - A B = 3 seconds allowed for upper-level compaction



- Step 2. Calculate compaction size within latency for ULC
 - Based on result of Step 1, we can obtain compaction size
 - Compaction size = latency of compaction × compaction speed



➔ Compaction size = 3 seconds × 100MB/s = 300MB



• Step 3. Find number of ranges corresponding to the size



Evaluation setup



Evaluation on mixed workload

- For write/read intensive workloads, throughput of DOPA-DB is higher
 - Dynamic compaction size recommendation helps to optimize compaction for each operation



Evaluation on mixed workload

- Multiple sub levels increases the number of disk access when workload pattern changes
- Large amount of upper-level compaction causes disk bandwidth competition





Discussion & Future work





Discussion & Future work

Nimble and accurate adaptation

- L_0 size •
- Flush speed ٠
- Compaction speed ٠
- Operation ratio? ٠
- Bandwidth competition? ٠





Discussion & Future work

Current DOPA-DB **only** supports single-threaded compaction Using multiple compaction threads Controlling compaction with multiple threads has to be considered • Parameters of current DOPA-DB are chosen **empirically** with 2 considerations Parameters of DOPA-DB Range must not make size of SST too small ٠ Size of compaction unit size must not be too large ٠ Key order distribution of real-world dataset is dynamic Skewed datasets Dynamic range readjustment is necessary for uniformly sized ranges



Conclusion

• DOPA-DB : Dynamic wOrkload Pattern Aware LSM-based KV store

- Advocate dynamic compaction size for dynamic workload
- Tiered based LSM KVS can control fine-grained compaction size

Evaluation on mixed workload

- Potential benefits of dynamic compaction
- Insight into research directions for dynamic compaction strategies

• Envisioned dynamic approach can further drive performance improvements







heejin5178@unist.ac.kr

