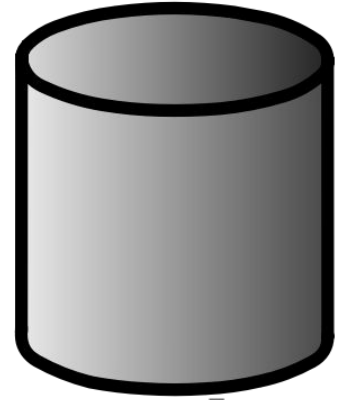# Range Deletes in LSM-tree

Benchmark and Analysis

Xiao Zhang, Wangyi Chen, Xingjian Zhang

# Introduction

- LSM-trees in modern databases

- Problem of range deletes

- Exploring the algorithm for efficient range deletions

# LSM-Trees: Origins and Challenges

- History of LSM-trees

- Role in database management systems

- Advantages and Challenges

# Problem Statement

- Need for efficient range deletion in LSM-trees

- Performance challenges in range deletion

- Goal: minimize I/O and CPU overhead while maintaining consistency
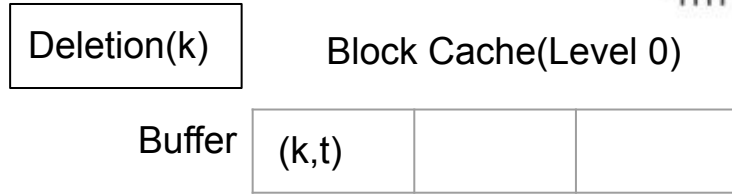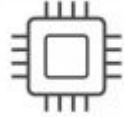
# Goals

Benchmarking on RocksDB

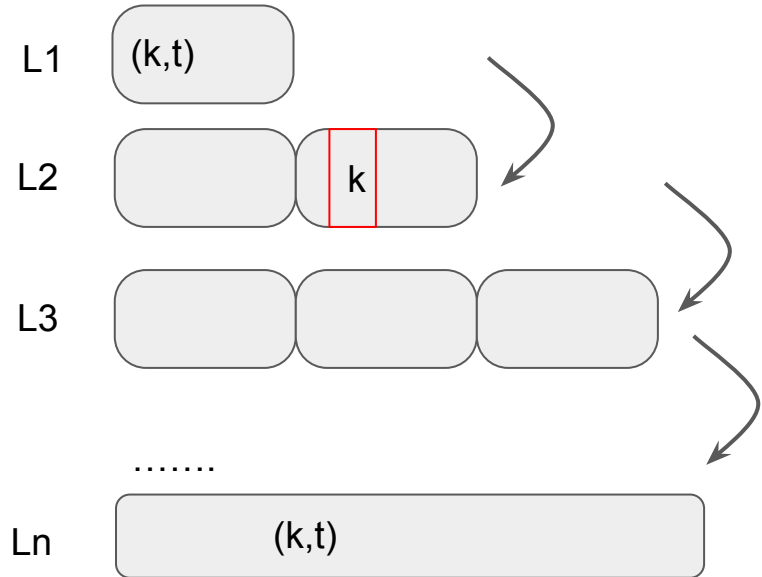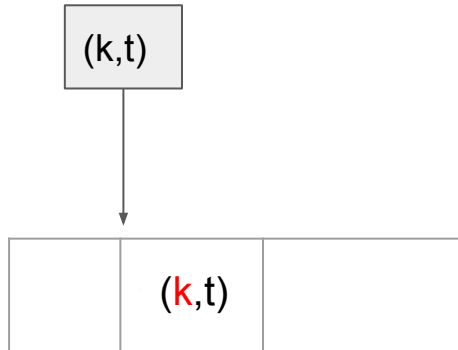Give intuitions on performance bound and potential improvements

Focus:

- Skyline Performance
- Memory Footprint
- Cost of Consistency

# LSM tree- Point Deletion

Deletion(k)

Block Cache(Level 0)

Buffer | (k,t) | | |

Compaction:

(k,t)

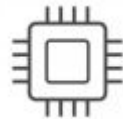| | (k,t) | |

L1 (k,t)

L2 | | k |

L3 | | | |

.......

Ln (k,t)

# Range Deletion

Are we going to add point_delete * N?   NO!

To avoid the memory buffer explode, we only need to insert the (lower, upper, tombstone)

To keep this range delete information, RocksDB provides an implementation of RangeDeleteBlock
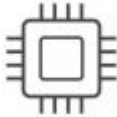
Block Cache(Level 0)

Buffer

| lower | upper | Tomb |
|-------|-------|------|

**X**   Not Match!

| Key | Value |
|-----|-------|

# Range Deletion

Are we going to add point_delete * N?   NO!

Block Cache(Level 0)

**Range Delete Block**

Buffer

| (k,v) | | | |
|---|---|---|---|

| Begin key | End key | Seqnum |
|---|---|---|
| "a" | "c" | 6 |
| "b" | "d" | 2 |
| "e" | "g" | 4 |
| … | … | |

Range Tombstone Block

| a | **6** | c | b | **2** | d | e | **4** | g | … |

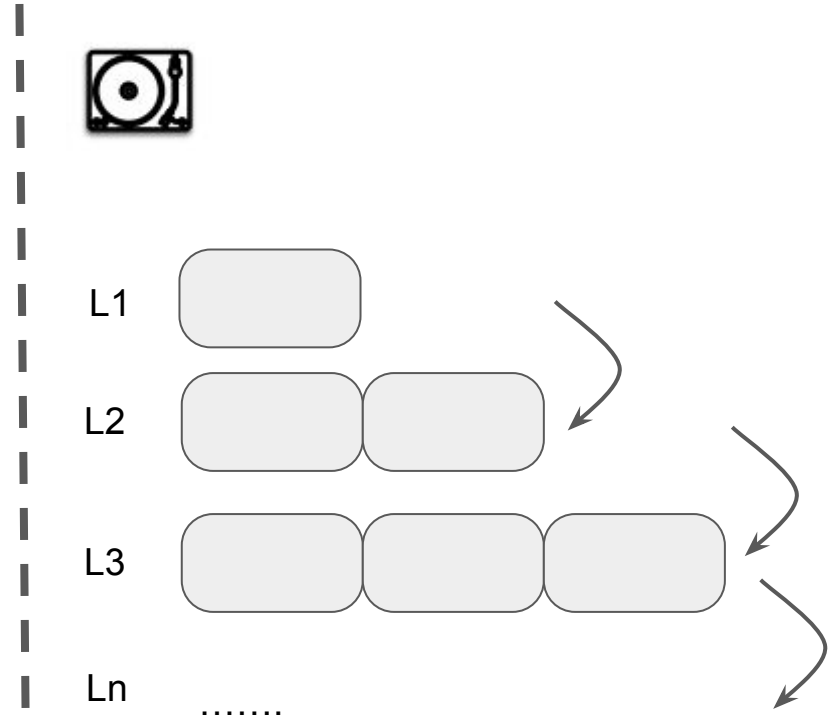| key bytes | seqnum | value bytes | …

# Range Deletion

Compaction

Pages in L1:
[1,8],[11,18],[20,30]

Pages in L2: [20,28],[30,48]

>>[10,30] from RDBlock
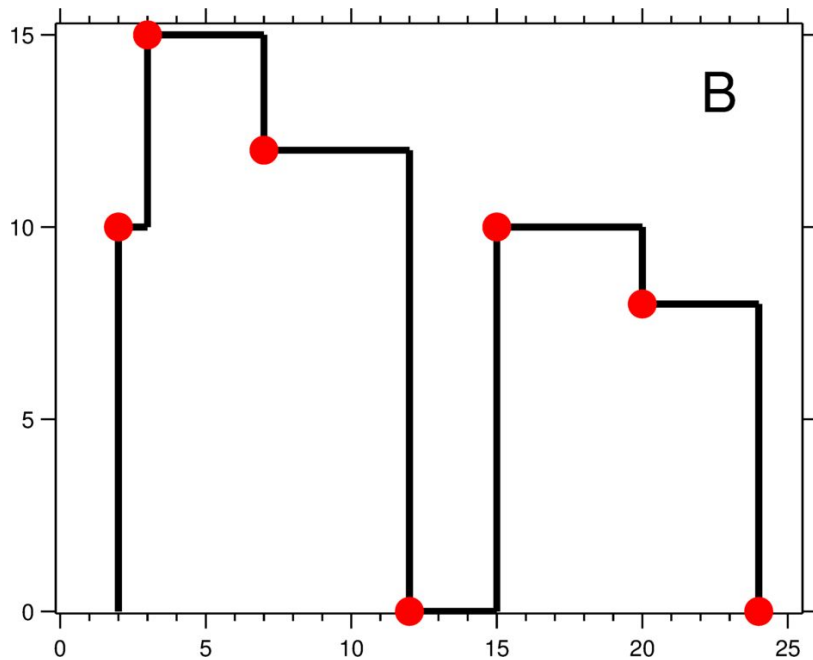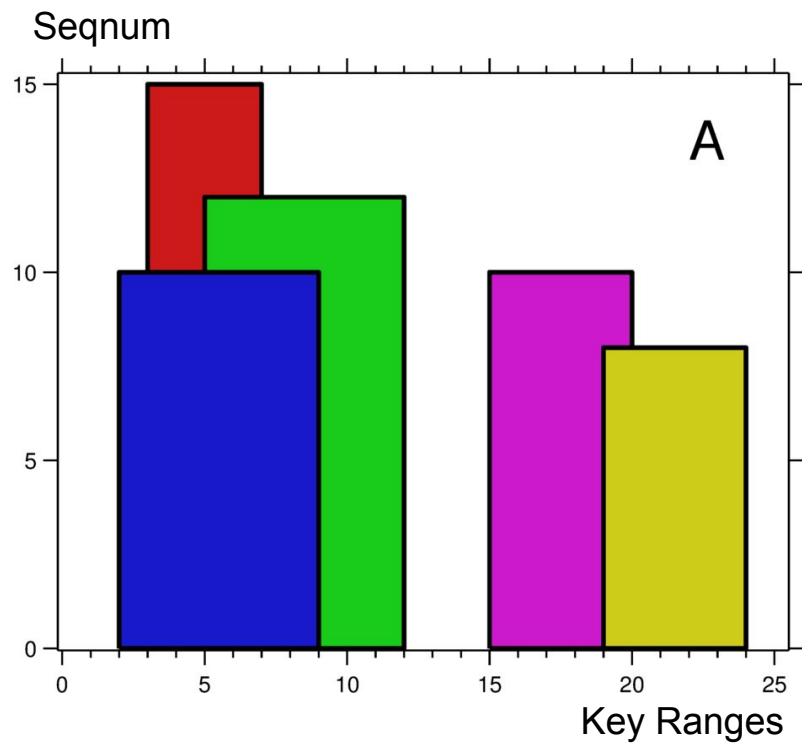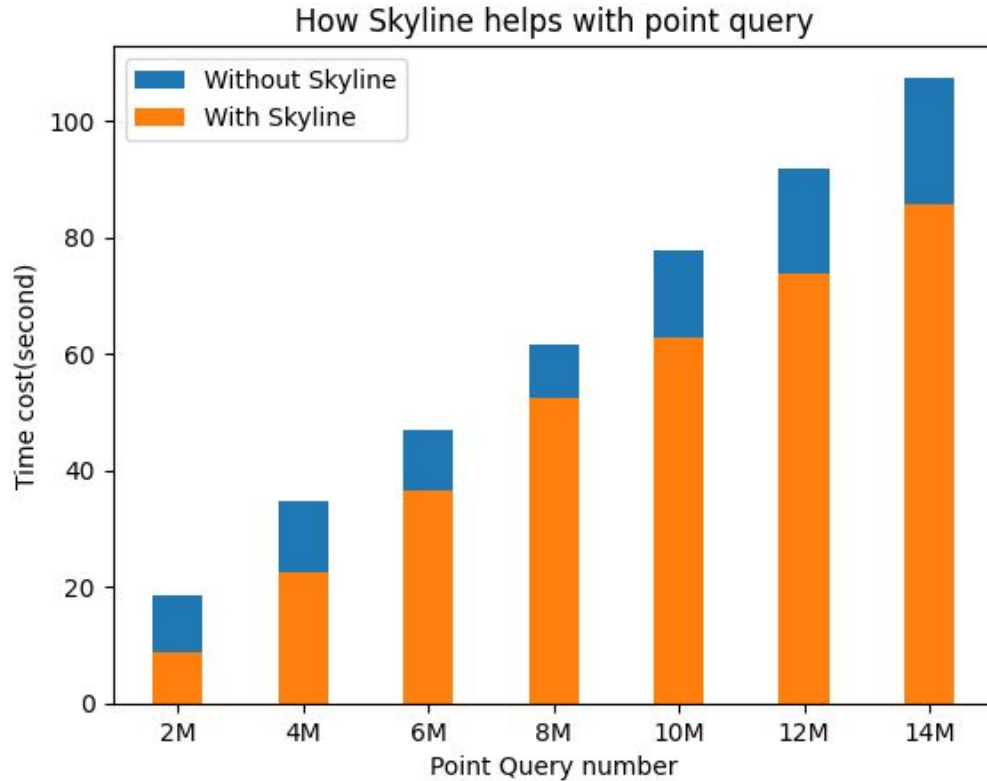
[11,18],[20,30],[20,28]

Pages in L2: [1,8],[30,48]

L1

L2

L3

Ln    .......

# Skyline Model

# Skyline Model

# Skyline Performance Tests



How Skyline helps with point query

20M values in Database.
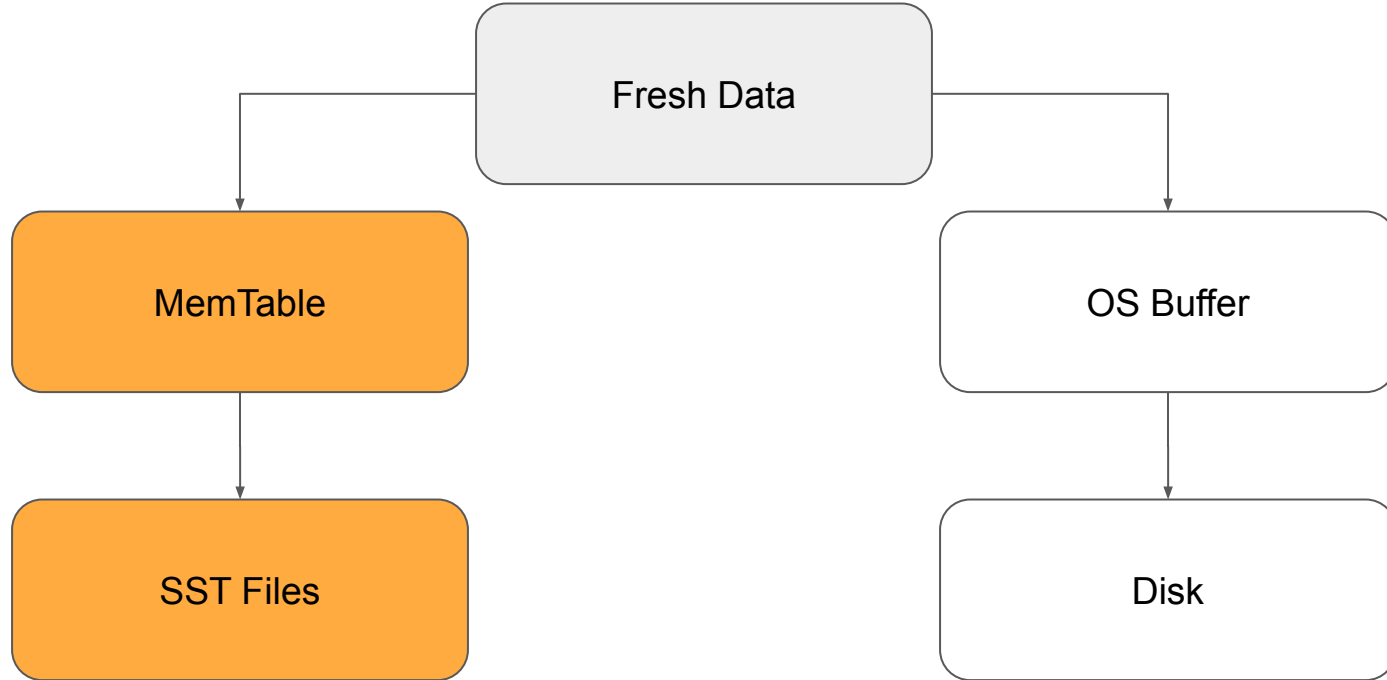2M inserts of Deletion

Reduces the time cost of
point query by 20%

# Cost of Consistency

# Consistency

- Write Append Log
- Checksum

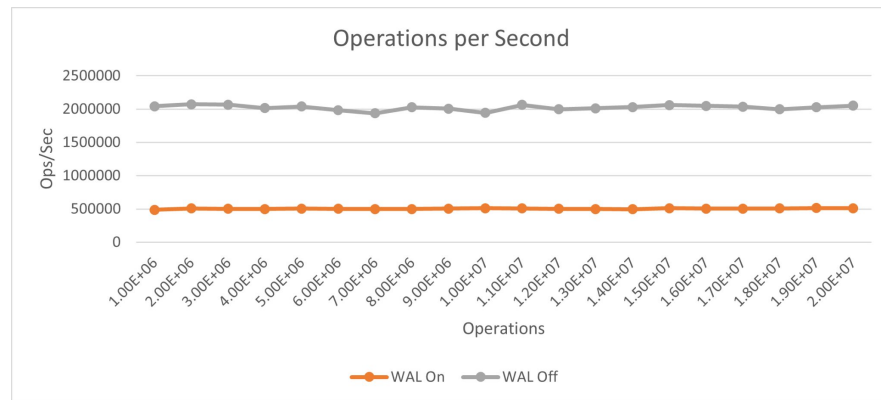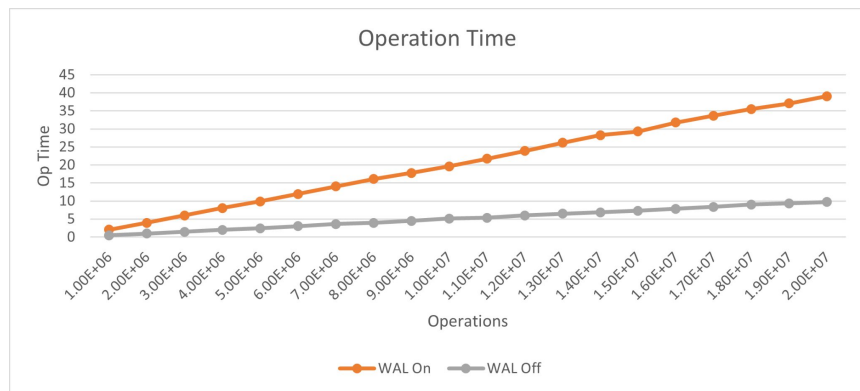# Write Ahead Log (WAL)

# Experiment Setup

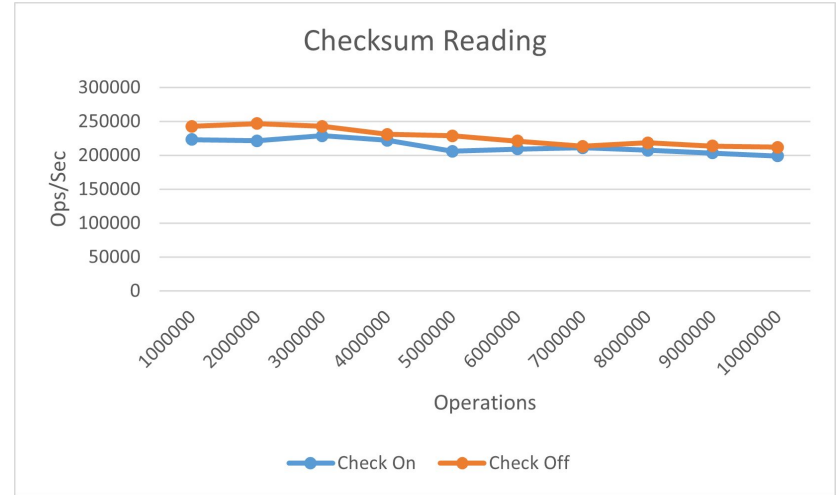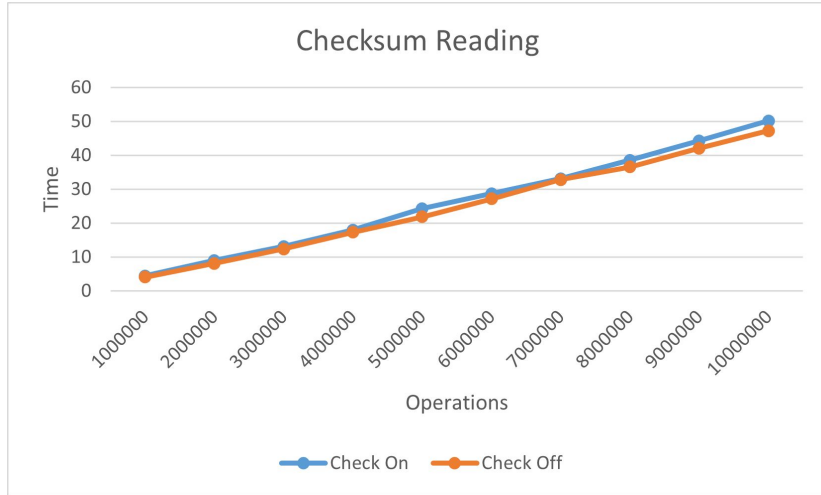Number of keys: [1:20] * 1,000,000

Number of levels: 7

Platform: db_bench

- fillseq: Fill the database with sequential data
- deleteseq: Delete a range of data
- readrandom: Read random data

# Results - WAL

# Results - Checksum

# Thank you

Q & A