# Benchmarking, Analyzing, and Optimizing Write Amplification of Partial Compaction in Rocksdb

Ran Wei*          Zichen Zhu*          Andrew Kryczka

Jay Zhuang          Manos Athanassoulis

# Log-Structured Merge-tree (LSM-tree)

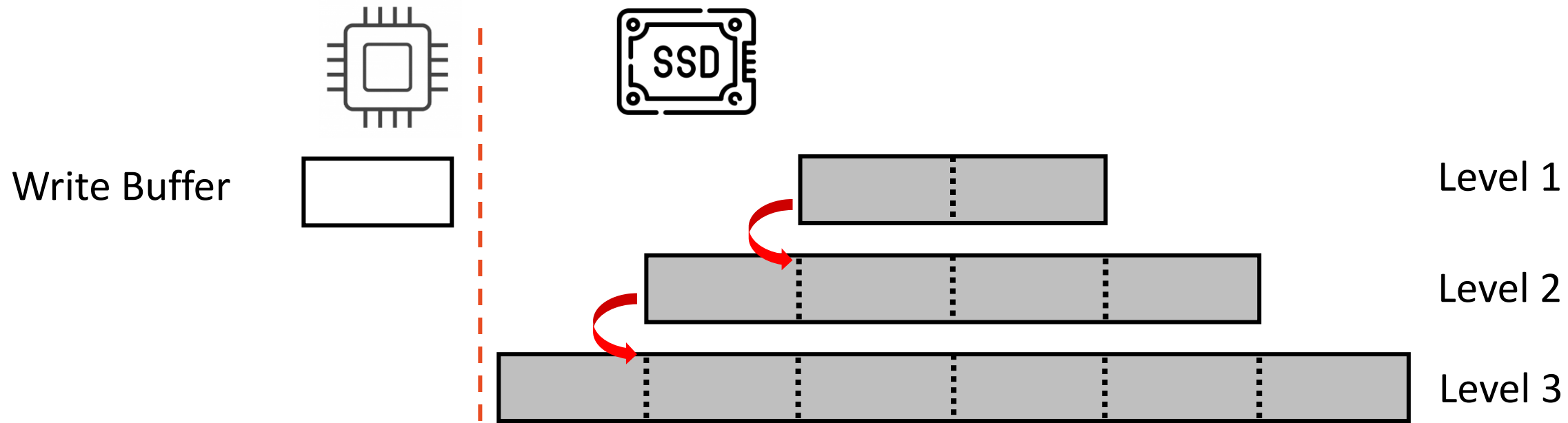Widely adopted because it offers fast ingestion rate and competitive read latency



**NoSQL**

**Relational**

**Time-series**

# Log-Structured Merge-Trees (LSM-Trees)
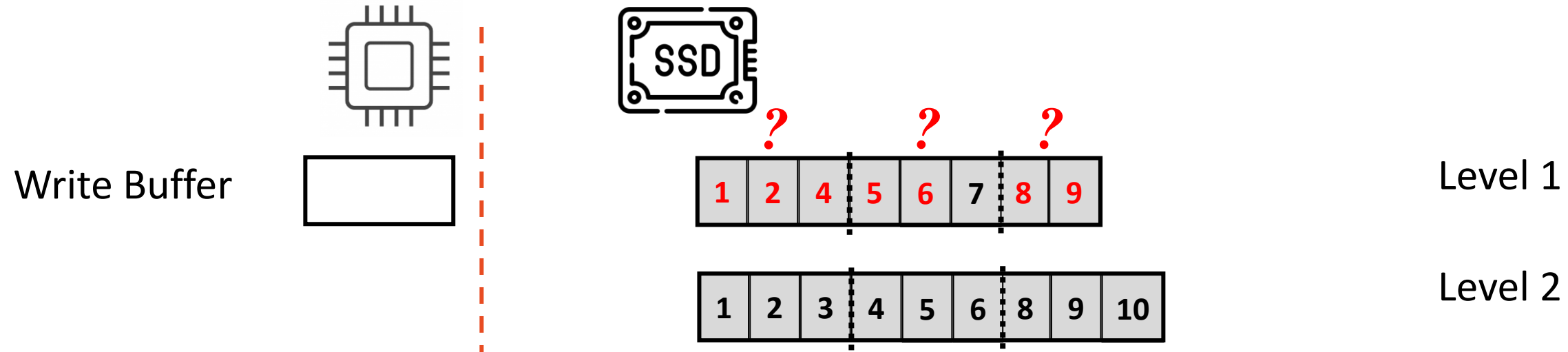
Write Buffer

Level 1

Level 2

Level 3

*Level capacity grows by size ratio T*

*Leveling (Classic): one sorted run per level*

Compaction introduces **Write Amplification**, short as

Write Amp, defined as $\dfrac{\#bytes\ written\ to\ storage}{\#bytes\ inserted\ by\ user}$

3

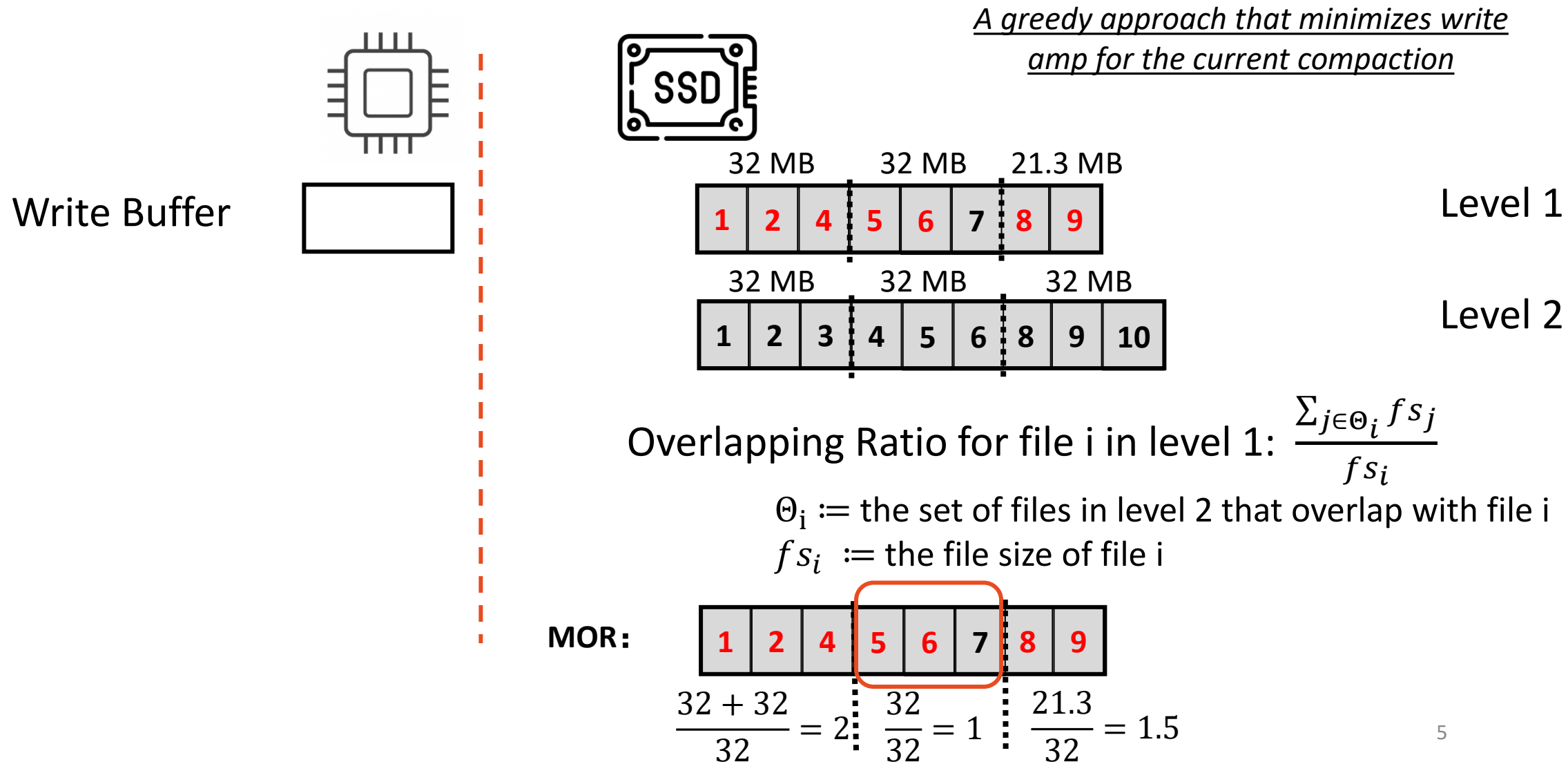# LSM Basics - Compaction

Write Buffer

Level 1

Level 2

Full compaction -> *High compaction latency*

Partial compaction -> Selecting one file to compact
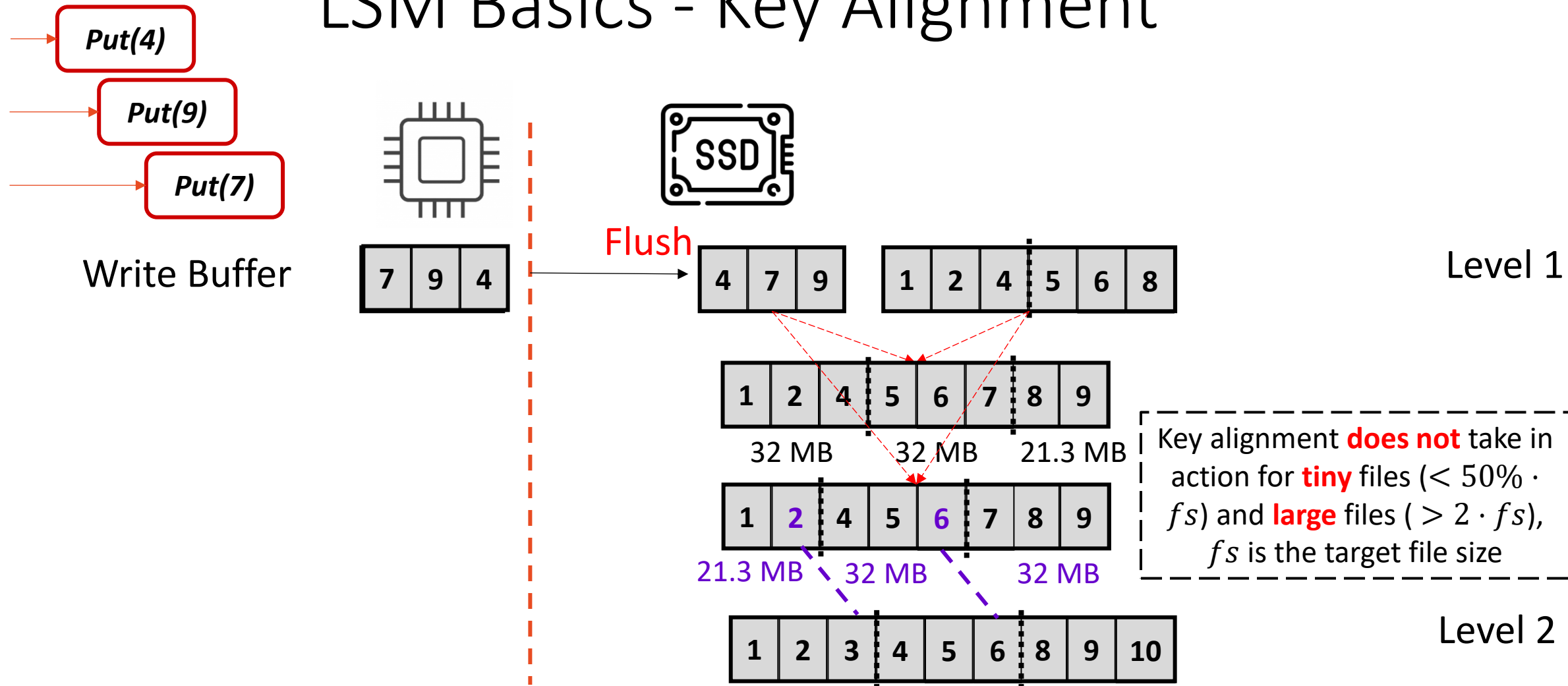
Which file in level 1 should we choose?

# LSM Basics - MinOverlappingRatio (MOR)

*A greedy approach that minimizes write amp for the current compaction*

Write Buffer

SSD

32 MB     32 MB     21.3 MB

| 1 | 2 | 4 | 5 | 6 | 7 | 8 | 9 |

Level 1

32 MB     32 MB     32 MB

| 1 | 2 | 3 | 4 | 5 | 6 | 8 | 9 | 10 |

Level 2

Overlapping Ratio for file i in level 1: $\dfrac{\sum_{j \in \Theta_i} f s_j}{f s_i}$

$\Theta_i :=$ the set of files in level 2 that overlap with file i

$f s_i :=$ the file size of file i

**MOR:**

| 1 | 2 | 4 | 5 | 6 | 7 | 8 | 9 |

$$\frac{32 + 32}{32} = 2 \qquad \frac{32}{32} = 1 \qquad \frac{21.3}{32} = 1.5$$

5

# LSM Basics - Key Alignment

Write Buffer

Level 1

Level 2

Some entries (e.g., 4) are unnecessarily written multiple times in the same level

# LSM Basics - Key Alignment

**Put(4)**

**Put(9)**

**Put(7)**

Write Buffer

SSD

Flush

**Level 1**

Write Buffer: 7 9 4

Flush to: 4 7 9 | 1 2 4 5 6 8

1 2 4 5 6 7 8 9

32 MB    32 MB    21.3 MB

1 **2** 4 5 **6** 7 8 9

21.3 MB    32 MB    32 MB

1 2 3 4 5 6 8 9 10

**Level 2**

Key alignment **does not** take in action for **tiny** files ($< 50\% \cdot fs$) and **large** files ($> 2 \cdot fs$), $fs$ is the target file size

Key alignment reduces unnecessary rewritings in the same level

# What is the Minimum Write-Amp?

Write Buffer

Level 1

Level 2

Level 3

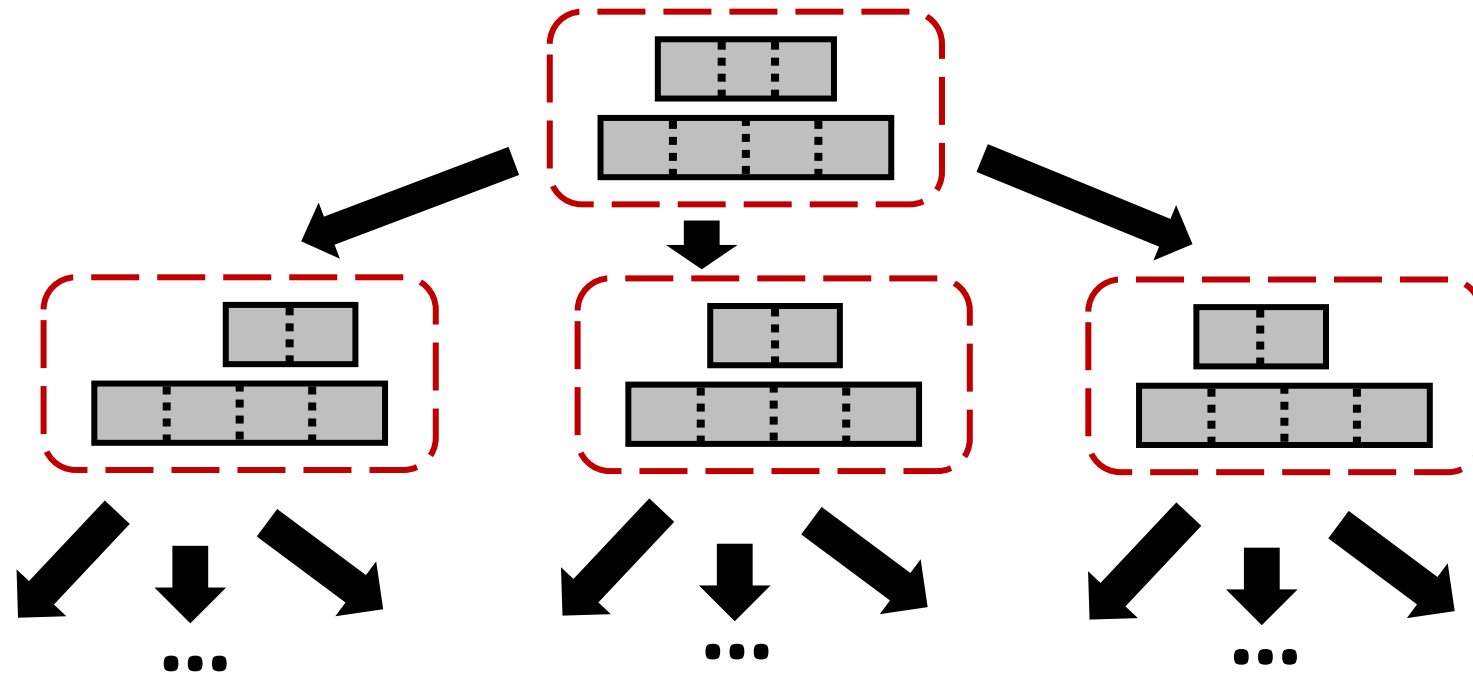...                                                                                    ...

What is the *minimum write amplification* if we can
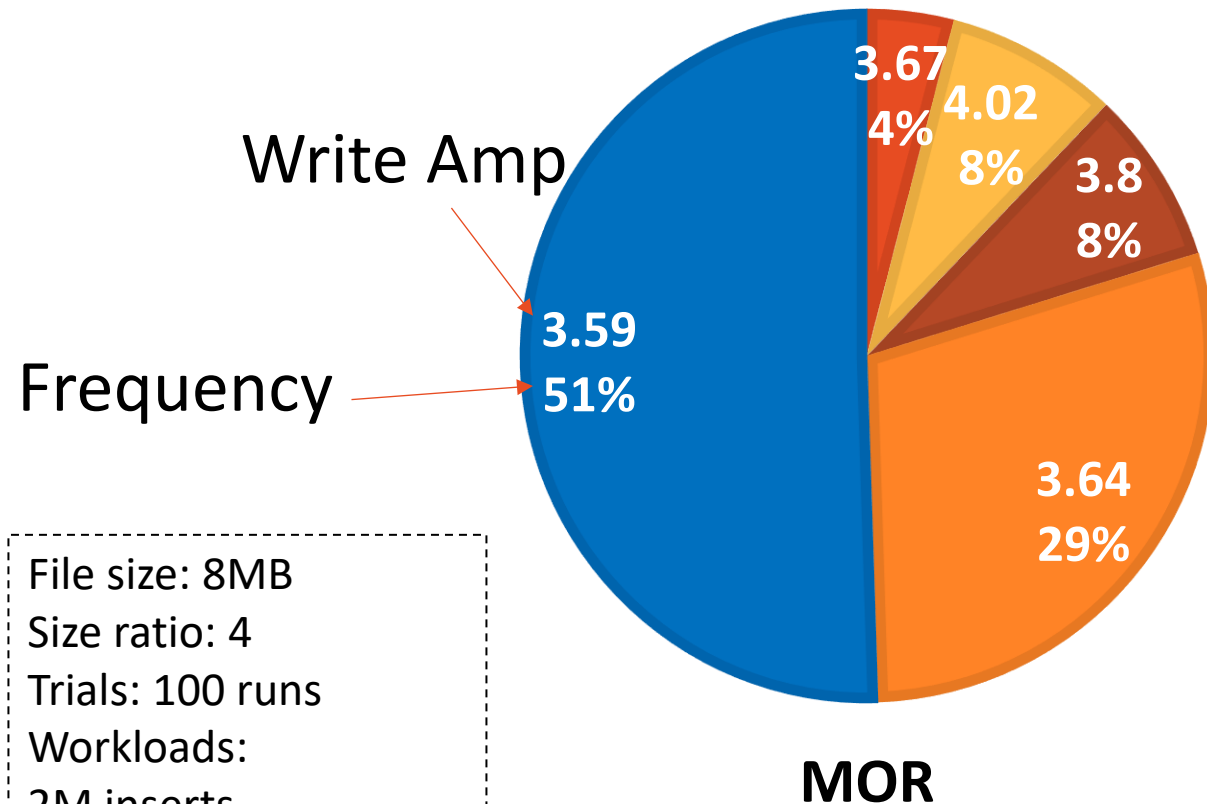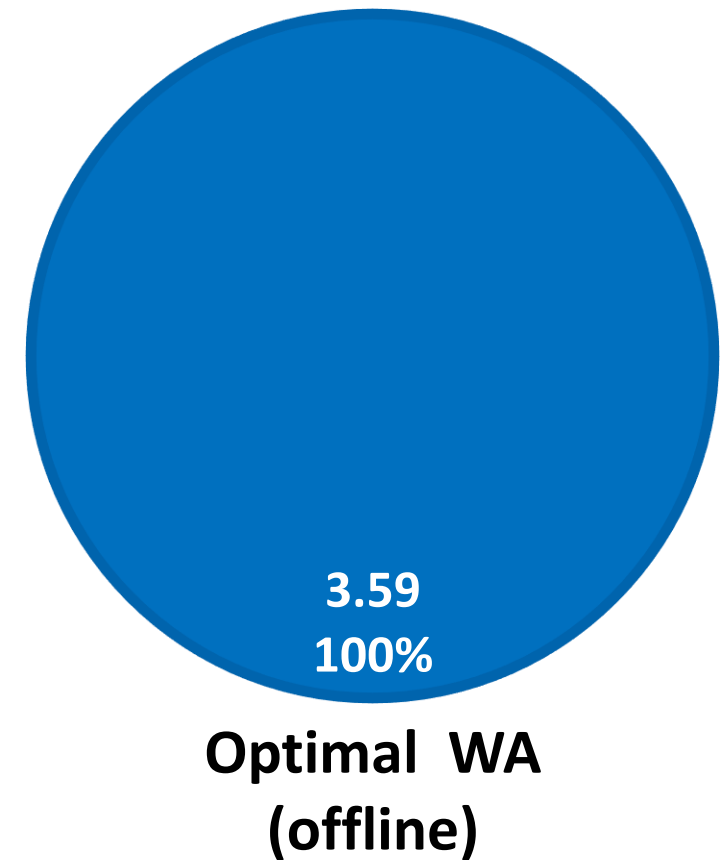smartly choose the file to compact?

# Write-Amp is Hard to Model



Picking a different file can have cascadingly different impact over the LSM-tree shape

# Write-Amp is Hard to Optimize

Brute Force: Depth-First Searching/Breadth-First Searching

⇓          ⇓          ⇓

*Exponential **searching space!***
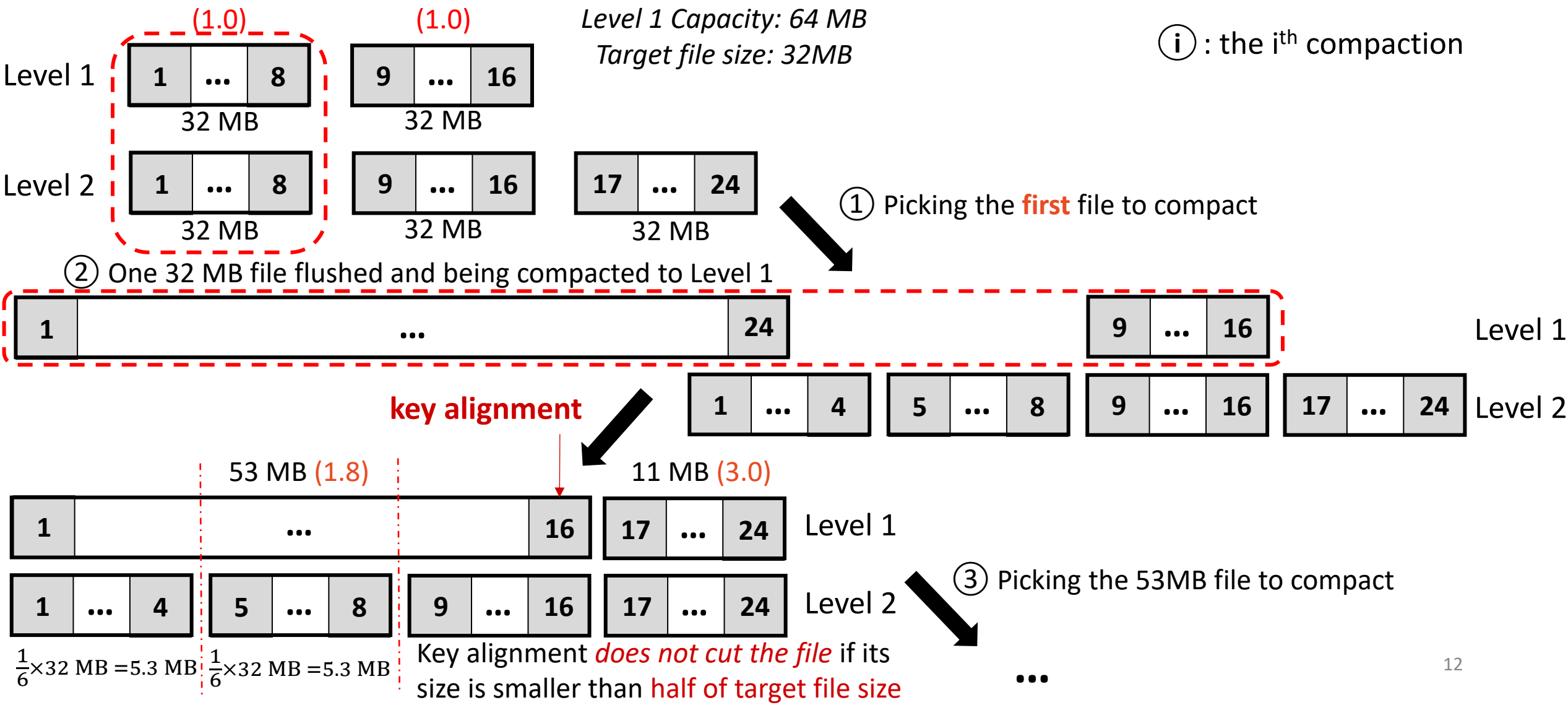
# MinOverlappingRatio (MOR)



Write Amp

Frequency

File size: 8MB
Size ratio: 4
Trials: 100 runs
Workloads:
2M inserts
8-byte key
56-byte value

**MOR**

**Optimal WA (offline)**

MinOverlappingRatio (MOR) file picking policy can **reach the minimum write amp** but it is **unstable**.
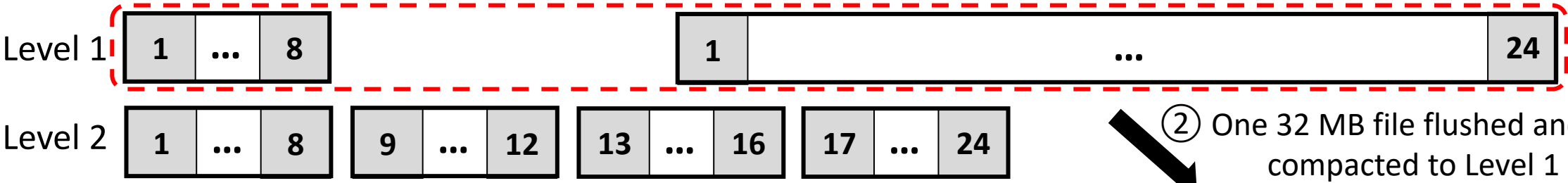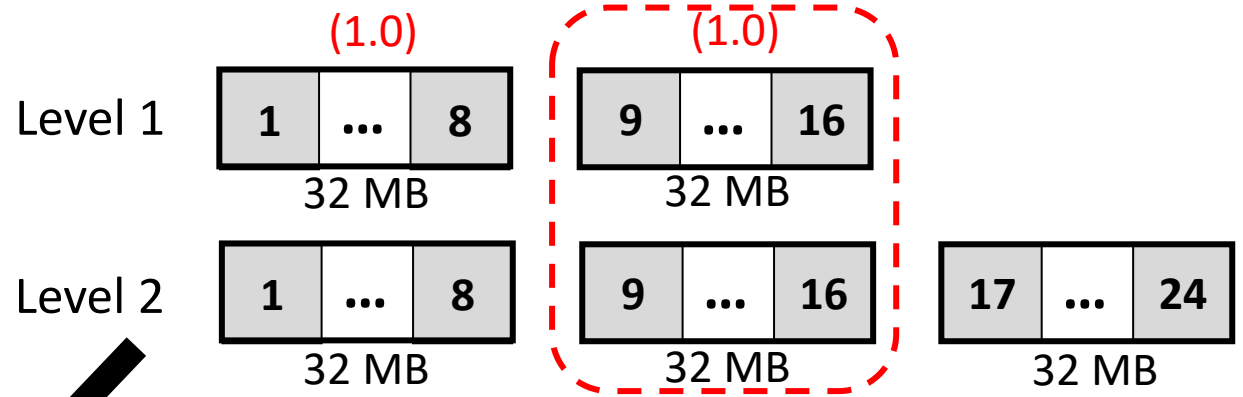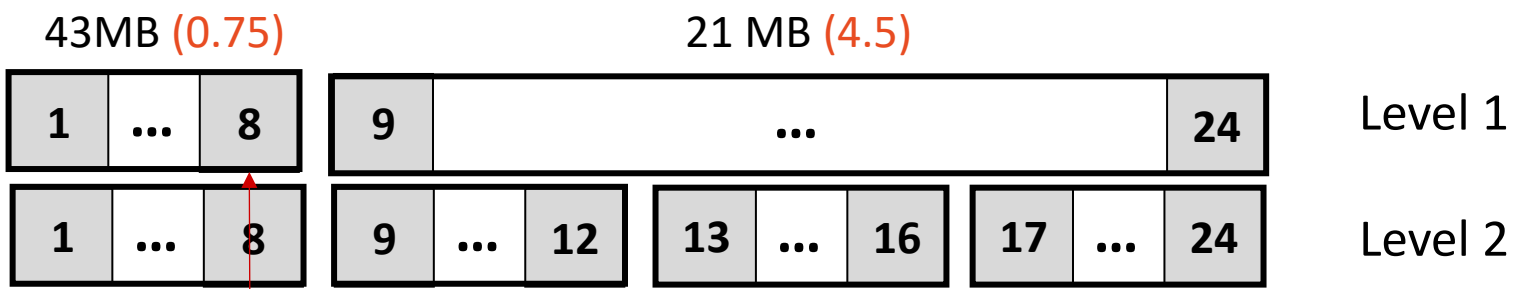
# Example of MinOverlappingRatio



Level 1 Capacity: 64 MB
Target file size: 32MB

(i) : the i^th compaction

(1.0)    (1.0)

Level 1    | 1 ... 8 |    | 9 ... 16 |
           32 MB          32 MB

Level 2    | 1 ... 8 |    | 9 ... 16 |    | 17 ... 24 |
           32 MB          32 MB           32 MB

① Picking the **first** file to compact

② One 32 MB file flushed and being compacted to Level 1

| 1 ... 24 |    | 9 ... 16 |    Level 1

| 1 ... 4 |  | 5 ... 8 |  | 9 ... 16 |  | 17 ... 24 |    Level 2

**key alignment**

53 MB (1.8)    11 MB (3.0)

| 1 ... 16 |  | 17 ... 24 |    Level 1

| 1 ... 4 | | 5 ... 8 | | 9 ... 16 | | 17 ... 24 |    Level 2

$\frac{1}{6} \times 32$ MB = 5.3 MB   $\frac{1}{6} \times 32$ MB = 5.3 MB

③ Picking the 53MB file to compact

Key alignment *does not cut the file* if its size is smaller than half of target file size

...

12

# Example of MinOverlappingRatio

*Level 1 Capacity: 64 MB*
*Target file size: 32MB*

(i) : the $i^{th}$ compaction

① Picking the **second** file to compact



Level 1

(1.0) | (1.0)

| 1 ... 8 | 9 ... 16 |
| 32 MB | 32 MB |

Level 2

| 1 ... 8 | 9 ... 16 | 17 ... 24 |
| 32 MB | 32 MB | 32 MB |

Level 1

| 1 ... 8 | 1 ... 24 |

Level 2

| 1 ... 8 | 9 ... 12 | 13 ... 16 | 17 ... 24 |

② One 32 MB file flushed and being compacted to Level 1

43MB (0.75)          21 MB (4.5)

| 1 ... 8 | 9 ... 24 |     Level 1

③ Picking the 43MB file to compact

| 1 ... 8 | 9 ... 12 | 13 ... 16 | 17 ... 24 |     Level 2
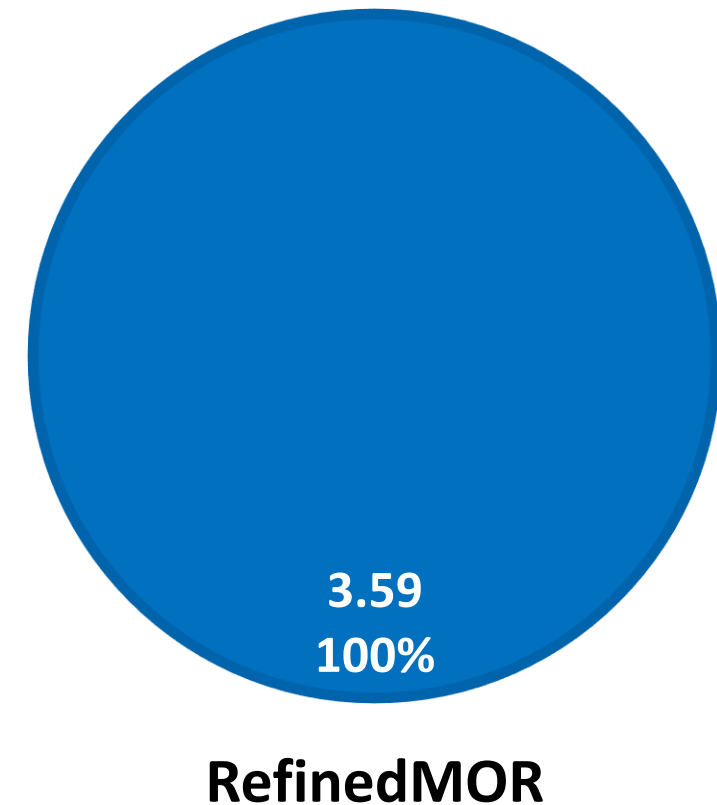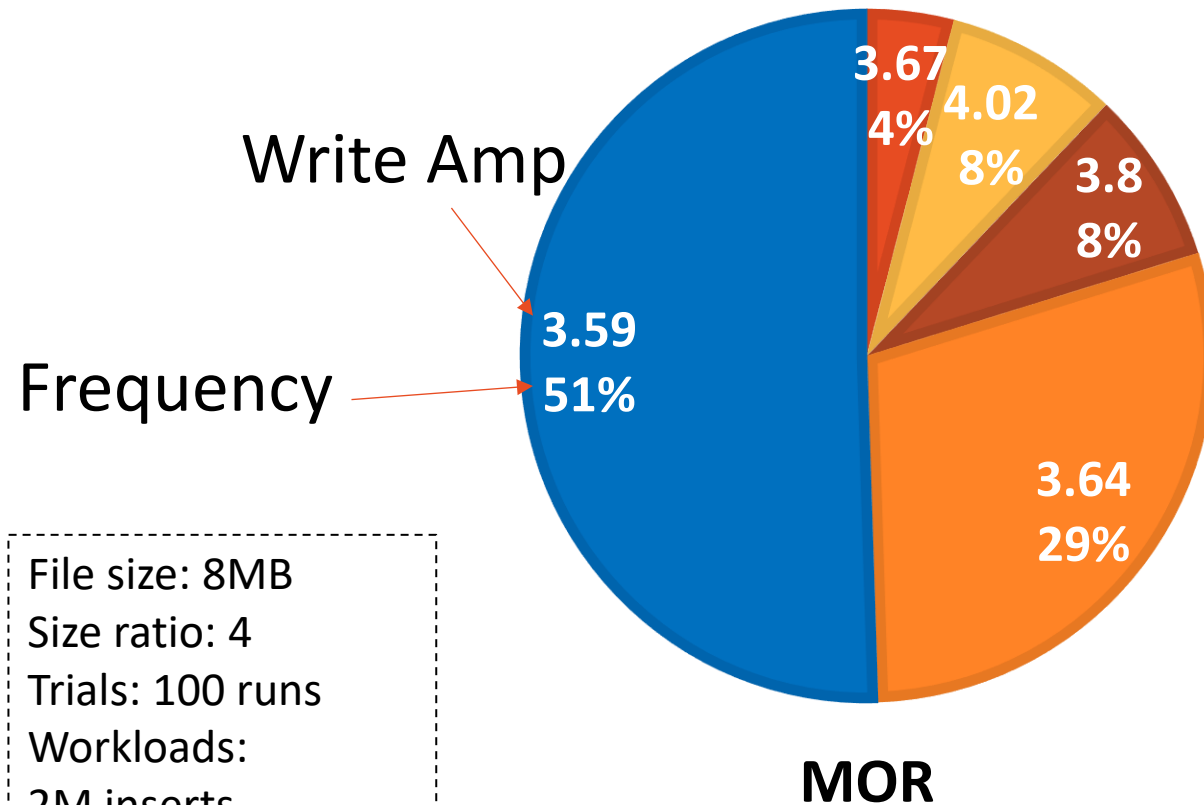
...

**key alignment**

13

# Example



To achieve key alignment, future compaction merges the subsequent file of the picked one.

⇓  ⇓  ⇓

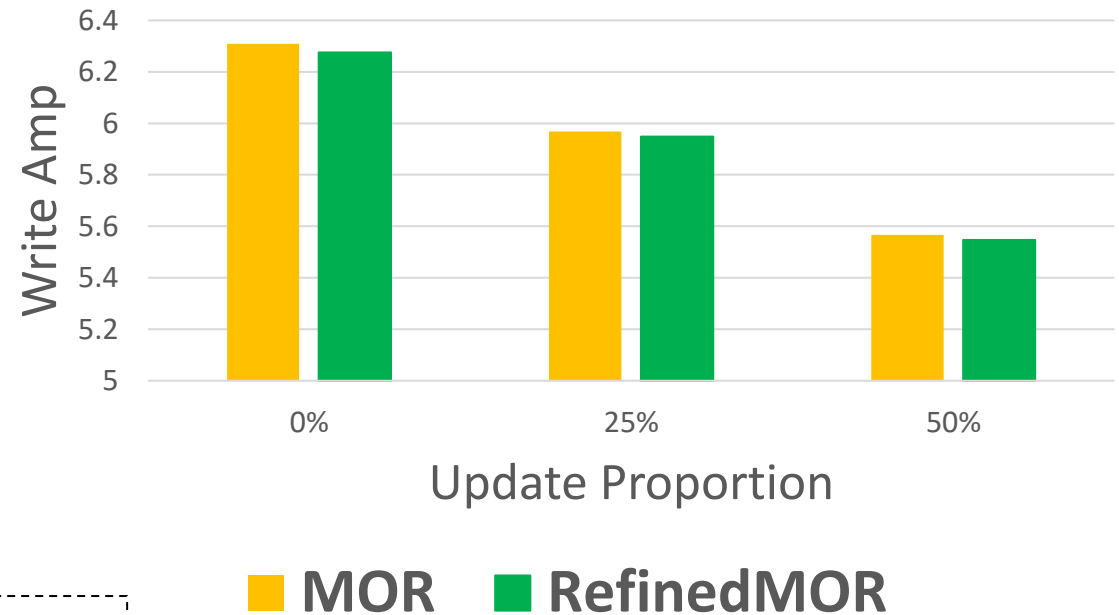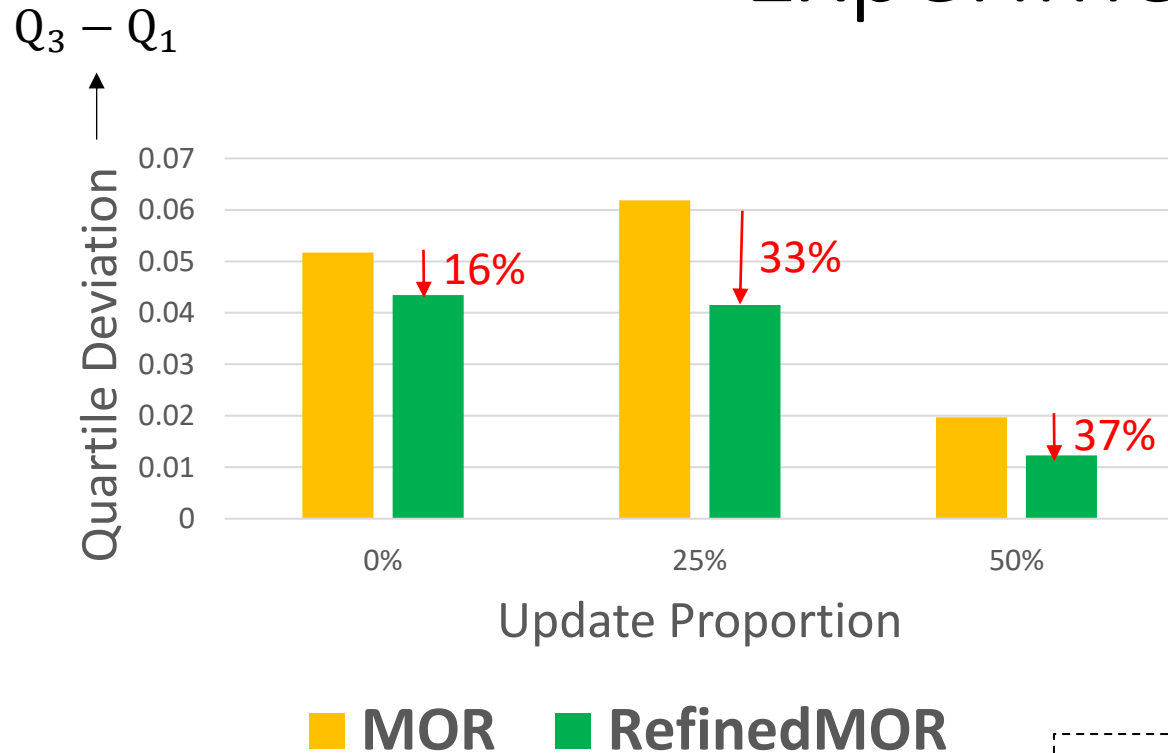RefinedMOR: Pick the one of which the subsequent file (if exists) has the largest overlapping ratio from minimum ones.

# Experimental Results



Write Amp

Frequency

File size: 8MB
Size ratio: 4
Trials: 100 runs
Workloads:
2M inserts
8-byte key
56-byte value

**MOR**

**RefinedMOR**

RefinedMOR achieves the minimum write amp in all 100 runs for small workloads.

# Experimental Results



$Q_3 - Q_1$

Quartile Deviation / Update Proportion

↓ 16%   ↓ 33%   ↓ 37%

0% 25% 50%

■ MOR   ■ RefinedMOR

Write Amp / Update Proportion

0% 25% 50%

■ MOR   ■ RefinedMOR

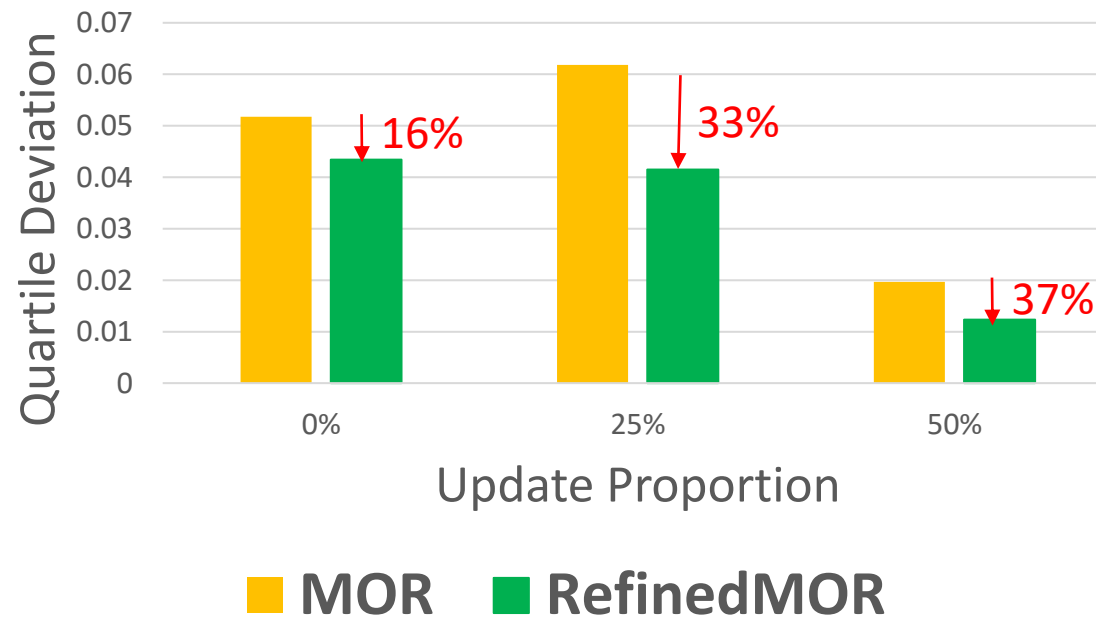File size: 8MB
Size ratio: 4
Trials: 20 runs

RefinedMOR has **similar average write amp** and **much lower quartile deviation** compared to MOR.

16

# Summary of RefinedMOR

RefinedMOR reduces the quartile deviation (i.e., $Q_3 - Q_1$) by up to 37% without worsening the average write amp.

# Other Observations

Round-Robin selection policy favors workloads with update skew.

Slower storage devices have lower WA but higher space amplification.

Picking policy has low impact over WA for update-intensive workloads.

Trivial move should be always prioritized to compaction.

More can be found in our full paper *Benchmarking, Analyzing, and Optimizing Write Amplification of Partial Compaction in RocksDB.*

# Q & A