# PebblesDB: Building Key-Value Stores using Fragmented Log Structured Merge Trees

Pandian Raju[1], Rohan Kadekodi[1], Vijay Chidambaram[1,2], Ittai Abraham[2]
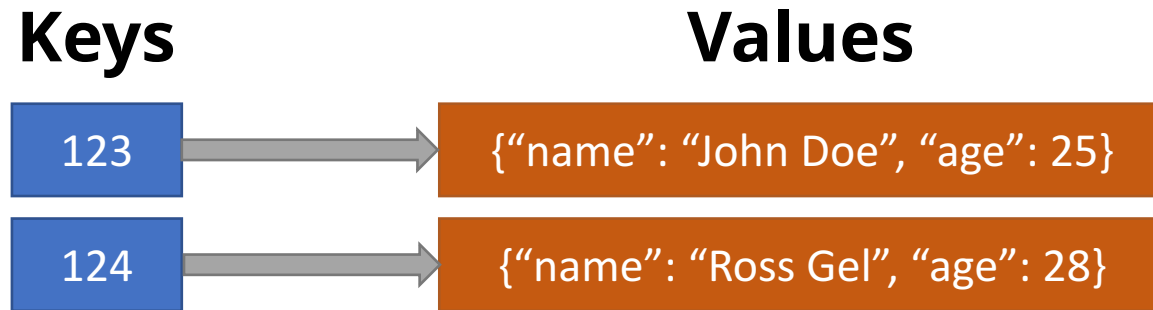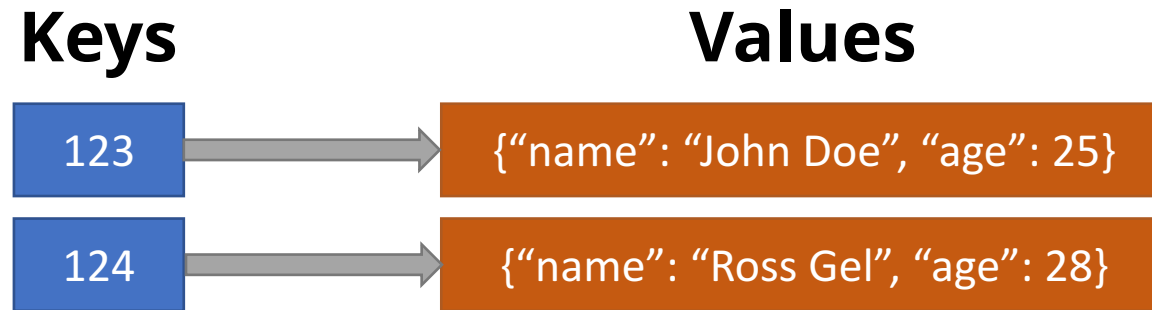
[1]The University of Texas at Austin

[2]VMware Research

**vmware**®

**TEXAS**
The University of Texas at Austin

# What is a key-value store?

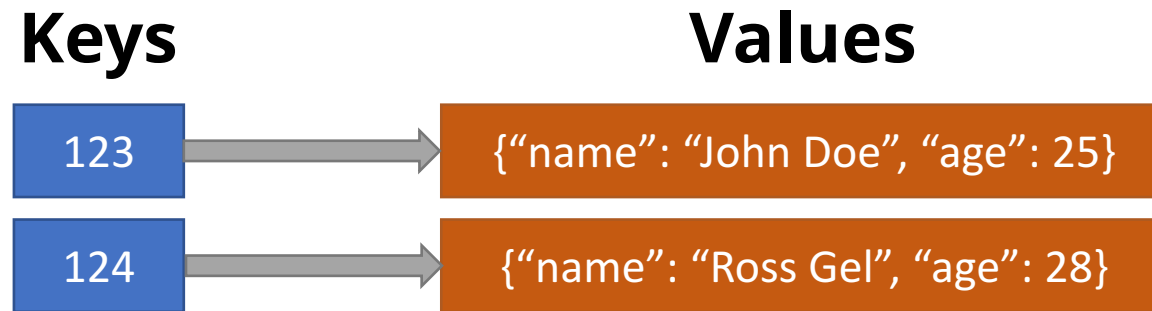- Store any arbitrary value for a given key

| Keys | Values |
|:---:|:---:|
| 123 | {"name": "John Doe", "age": 25} |
| 124 | {"name": "Ross Gel", "age": 28} |

# What is a key-value store?

- Store any arbitrary value for a given key

**Keys**                    **Values**

| 123 | → | {"name": "John Doe", "age": 25} |
| 124 | → | {"name": "Ross Gel", "age": 28} |

- **Insertions**:
- **Point lookups**:
- **Range Queries**:

# What is a key-value store?

- Store any arbitrary value for a given key

**Keys**                    **Values**

| 123 | → | {"name": "John Doe", "age": 25} |

| 124 | → | {"name": "Ross Gel", "age": 28} |

- **Insertions**: put(key, value)
- **Point lookups**:
- **Range Queries**:

# What is a key-value store?

- Store any arbitrary value for a given key

| **Keys** | **Values** |
|---|---|
| 123 | {"name": "John Doe", "age": 25} |
| 124 | {"name": "Ross Gel", "age": 28} |

- **Insertions**: put(key, value)
- **Point lookups**: get(key)
- **Range Queries**:

# What is a key-value store?

- Store any arbitrary value for a given key

**Keys** **Values**

| 123 | → | {"name": "John Doe", "age": 25} |
| 124 | → | {"name": "Ross Gel", "age": 28} |

- **Insertions**: put(key, value)
- **Point lookups**: get(key)
- **Range Queries**: get_range(key1, key2)

# Key-Value Stores - widely used

- Google's <span style="color:red">BigTable</span> powers Search, Analytics, Maps and Gmail
- Facebook's <span style="color:red">RocksDB</span> is used as storage engine in production systems of many companies
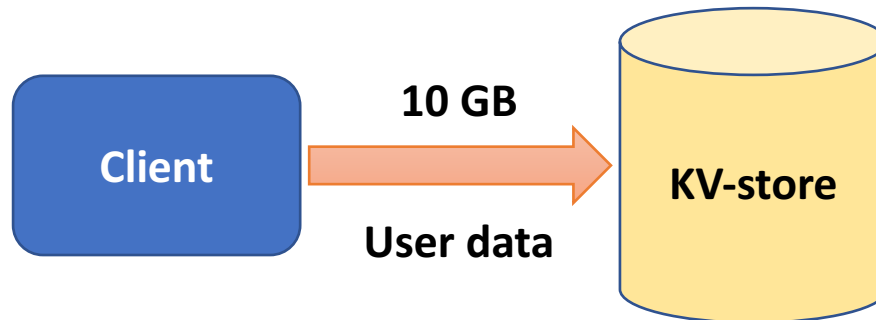
# Write-optimized data structures

- **Log Structured Merge Tree (LSM)** is a write-optimized data structure used in key-value stores
- Provides high write throughput with good read throughput, but suffers high write amplification

# Write-optimized data structures

- Log Structured Merge Tree (LSM) is a write-optimized data structure used in key-value stores

- Provides high write throughput with good read throughput, but suffers high write amplification

- Write amplification - Ratio of amount of write IO to amount of user data

```
            10 GB
Client  ─────────────►  KV-store
            User data
```
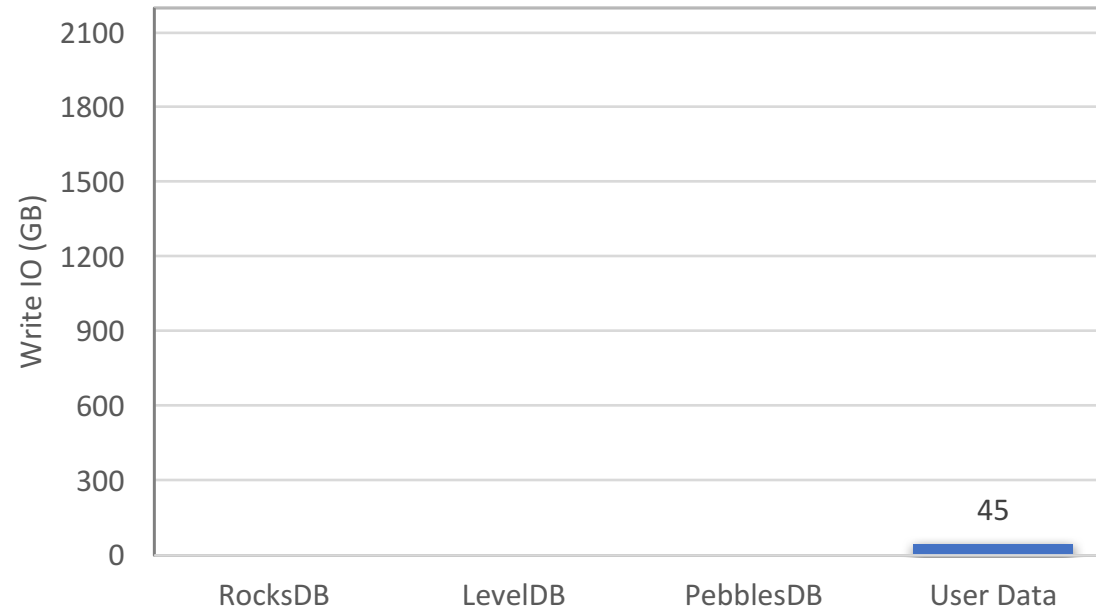
If total write I/O is 200 GB
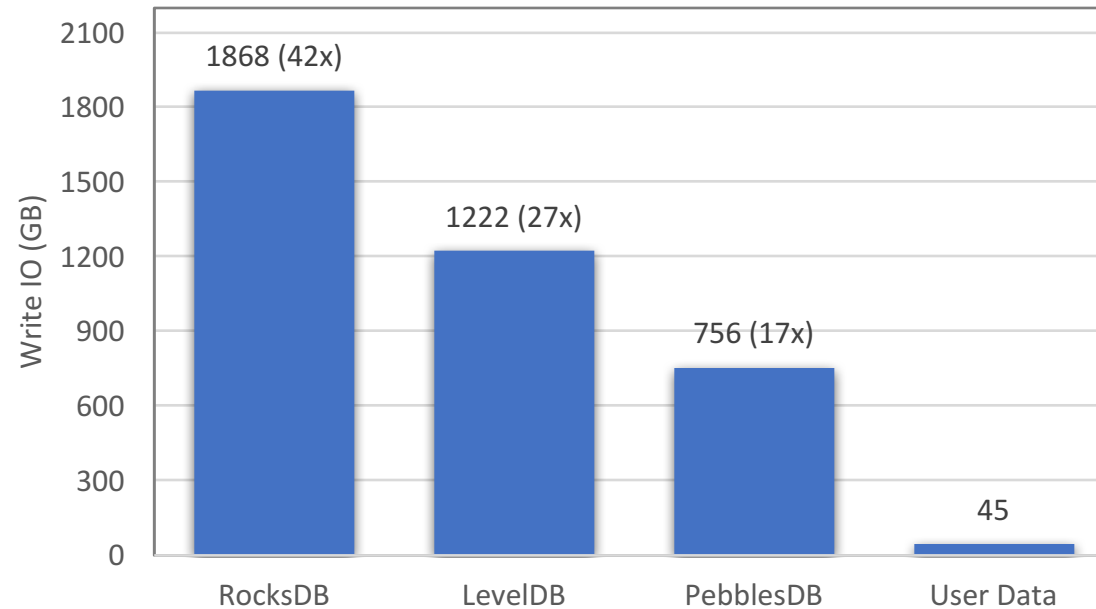
Write amplification = 20

# Write amplification in LSM based KV stores

- Inserted 500M key-value pairs
- Key: 16 bytes, Value: 128 bytes
- Total user data: ~45 GB

# Write amplification in LSM based KV stores

- Inserted 500M key-value pairs
- Key: 16 bytes, Value: 128 bytes
- Total user data: ~45 GB

# Why is write amplification bad?

- Reduces the write throughput
- Flash devices wear out after limited write cycles

(Intel SSD DC P4600 – can last ~5 years assuming ~5 TB write per day)

RocksDB can write ~500 GB of user data per day to a SSD to last 1.25 years

# **PebblesDB**

High performance write-optimized key-value store

Built using new data structure
Fragmented Log-Structured Merge Tree

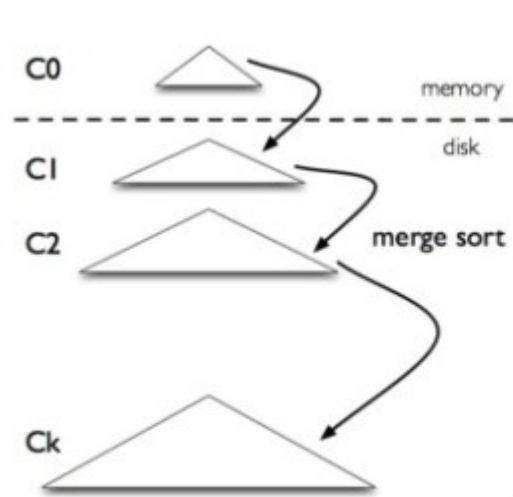Achieves 3-6.7x higher write throughput and 2.4-3x lesser write amplification compared to RocksDB

Gets the highest write throughput and least write amplification as a backend store to MongoDB
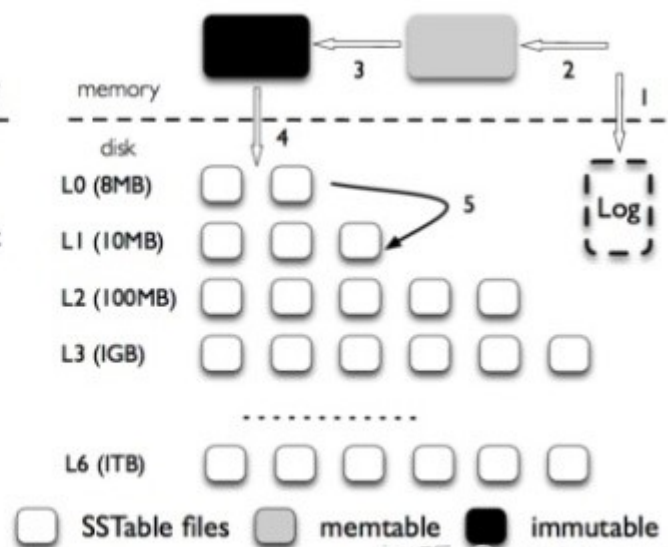
# Outline

- Log-Structured Merge Tree (LSM)
- Fragmented Log-Structured Merge Tree (FLSM)
- Building PebblesDB using FLSM
- Evaluation
- Conclusion

# Outline

- **Log-Structured Merge Tree (LSM)**
- Fragmented Log-Structured Merge Tree (FLSM)
- Building PebblesDB using FLSM
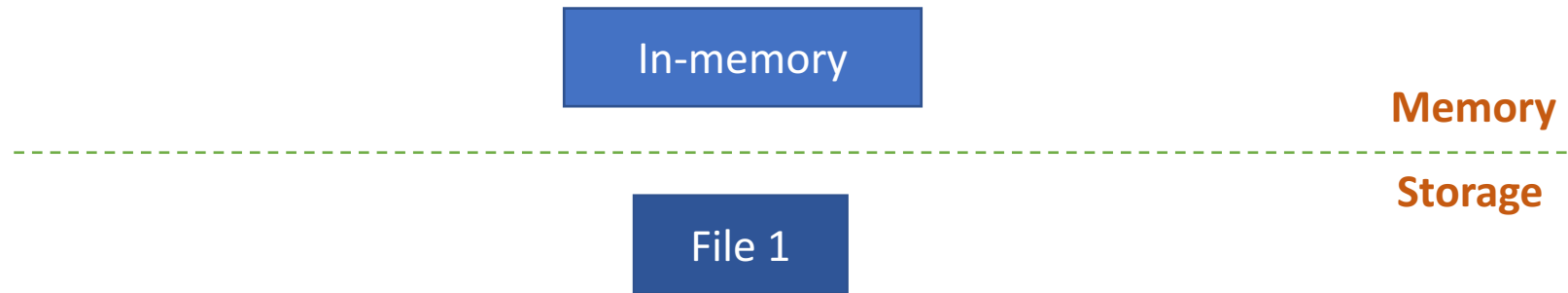- Evaluation
- Conclusion

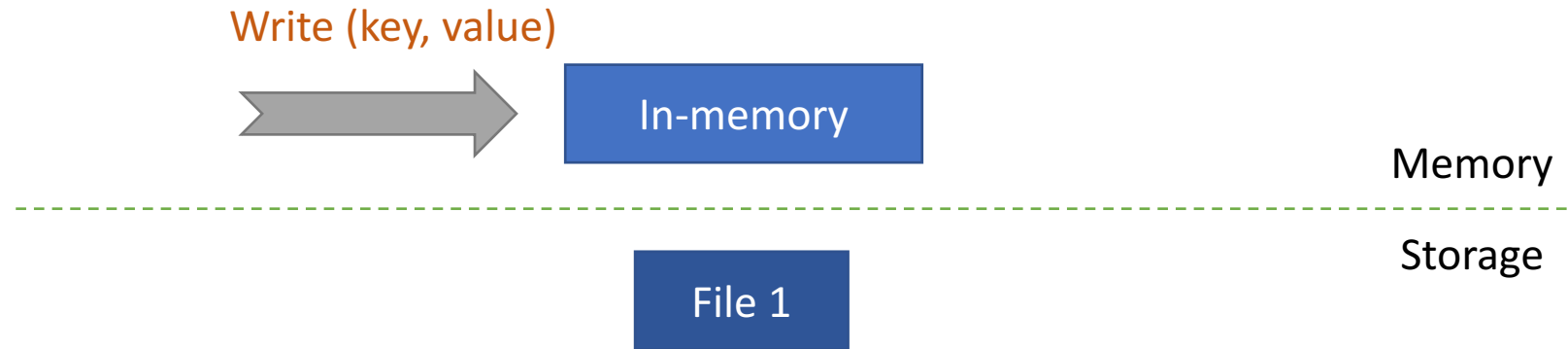| | | |
|---|---|---|
| C0 | | memory |
| C1 | | disk |
| C2 | merge sort | |
| Ck | | |

memory

disk

L0 (8MB)
L1 (10MB)
L2 (100MB)
L3 (1GB)

L6 (1TB)

Log

SSTable files    memtable    immutable

(a) LSM-tree          (b) LevelDB

知乎 @gushitong

# Log Structured Merge Tree (LSM)

In-memory

**Memory**

**Storage**

File 1

Data is stored both in memory and storage

# Log Structured Merge Tree (LSM)

Write (key, value)
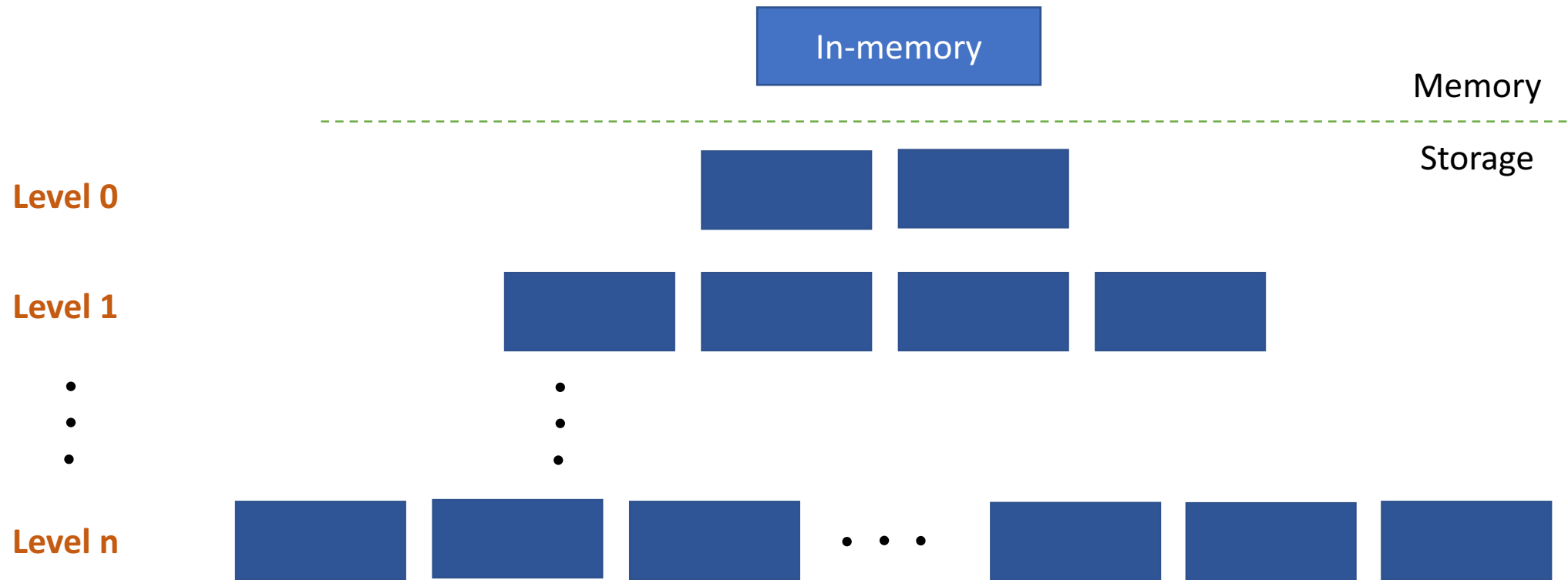
In-memory

Memory

Storage

File 1

Writes are directly put to memory
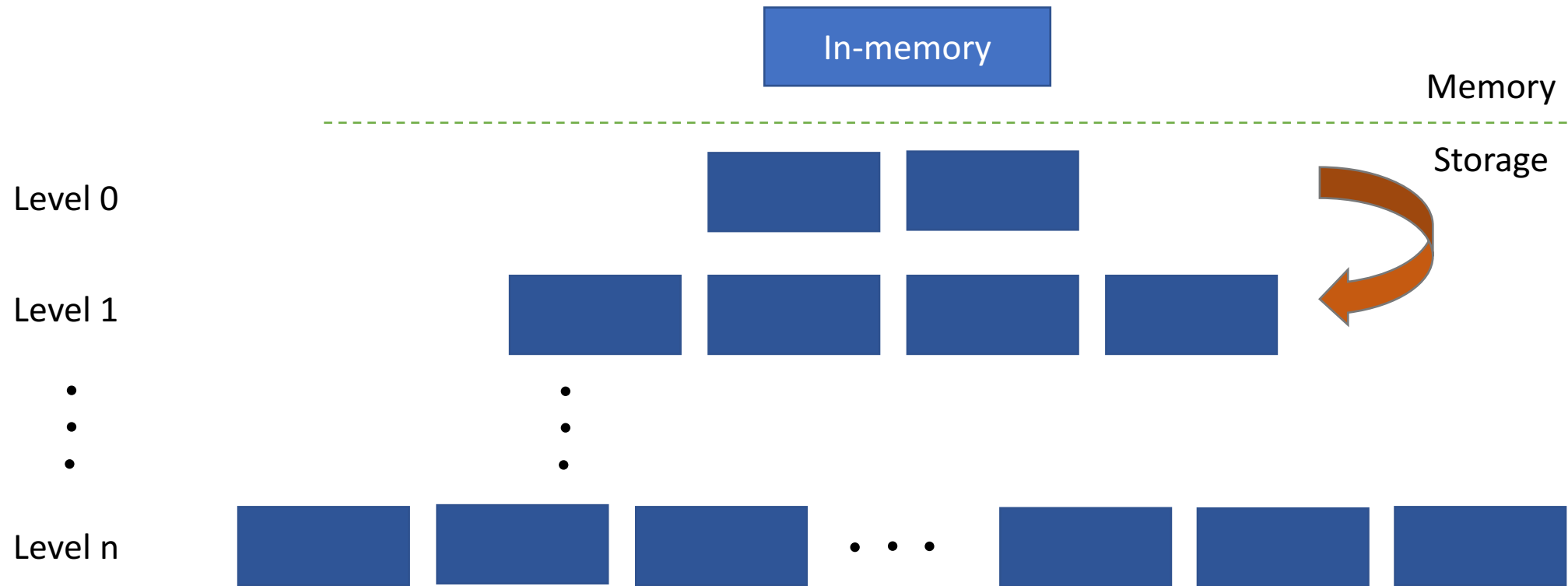
# Log Structured Merge Tree (LSM)



In-memory data is periodically written as files to storage (sequential I/O)

# Log Structured Merge Tree (LSM)
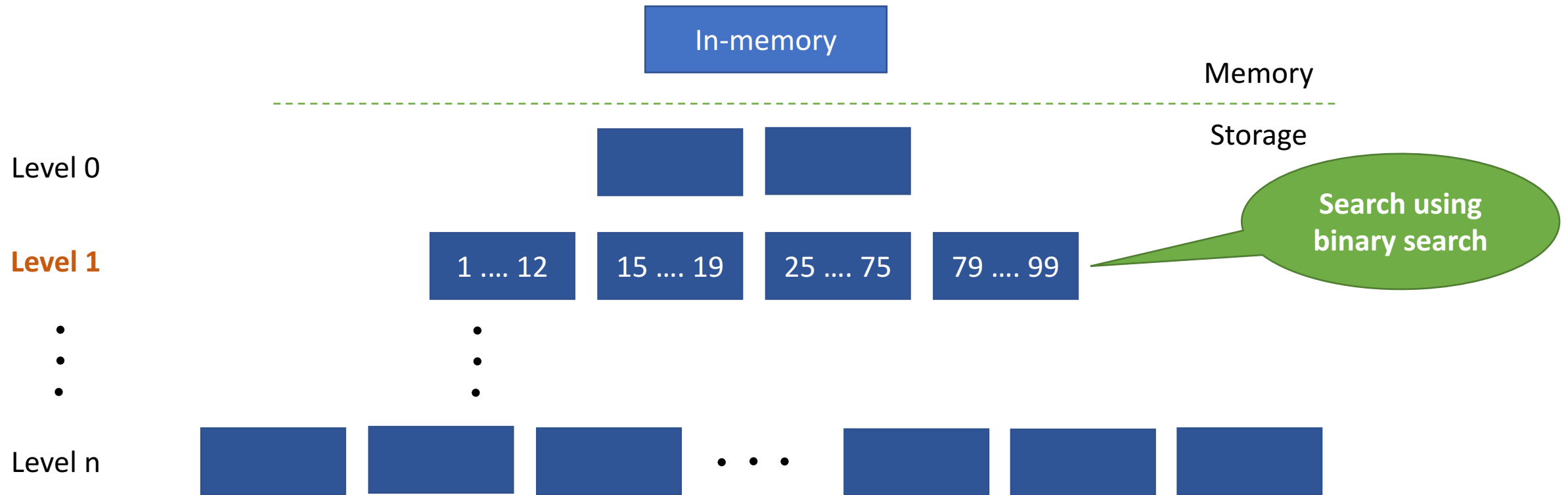


Files on storage are logically arranged in different levels

# Log Structured Merge Tree (LSM)



Compaction pushes data to higher numbered levels

# Log Structured Merge Tree (LSM)

In-memory

Memory

Storage

Level 0

Level 1

| 1 .... 12 | 15 .... 19 | 25 .... 75 | 79 .... 99 |

Search using binary search

Level n

...

Files are sorted and have non-overlapping key ranges

# Log Structured Merge Tree (LSM)

In-memory

Limit on number of level 0 files

Memory

Storage

**Level 0**

2 …. 57    23 …. 78

Level 1

Level n

Level 0 can have files with overlapping (but sorted) key ranges

# Write amplification: Illustration

| Time: $t_1$ New sstable in Level 0 | Level 0 | 10 210 | |
| | Level 1 | 1 100 | 200 400 |
| Time: $t_2$ After compacting Level 0 into Level 1 | Level 0 | | |
| | Level 1 | 1 **10** 100 | 200 **210** 400 |
| Time: $t_3$ New sstable in Level 0 | Level 0 | 20 220 | |
| | Level 1 | 1 10 100 | 200 210 400 |
| Time: $t_4$ After compacting Level 0 into Level 1 | Level 0 | | |
| | Level 1 | 1 10 **20** 100 | 200 210 **220** 400 |
| Time: $t_5$ New sstable in Level 0 | Level 0 | 30 330 | |
| | Level 1 | 1 10 20 100 | 200 210 220 400 |
| Time: $t_6$ After compacting Level 0 into Level 1 | Level 0 | | |
| | Level 1 | 1 10 20 **30** 100 | 200 210 220 **330** 400 |

**Figure 2: LSM Compaction. The figure shows sstables being inserted and compacted over time in a LSM.**

# Root cause of write amplification

Rewriting data to the same level multiple times

To maintain sorted non-overlapping files in each level

# Outline

- Log-Structured Merge Tree (LSM)
- **Fragmented Log-Structured Merge Tree (FLSM)**
- Building PebblesDB using FLSM
- Evaluation
- Conclusion

# Naïve approach to reduce write amplification

- Just append the file to the end of next level
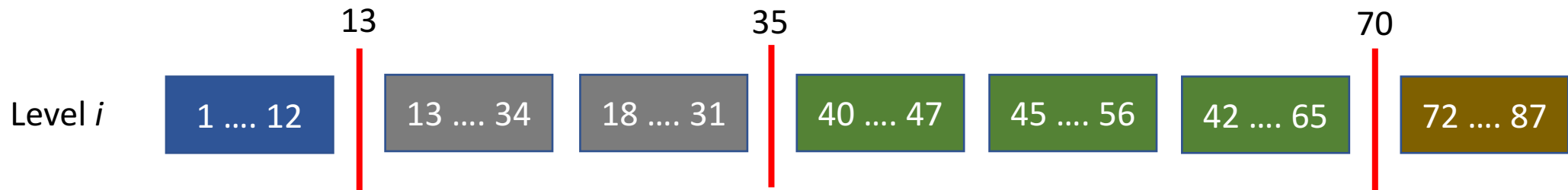- Many (possibly all) overlapping files within a level

Level *i* | 1 …. 89 | 5 …. 65 | 6 …. 91 | 8 …. 95 | 9 …. 99 | 1 …. 102 | 1 … 271

(all files have overlapping key ranges)

- Affects the read performance

# Partially sorted levels

- **Hybrid** between all non-overlapping files and all overlapping files
- Inspired from **Skip-List** data structure
- Concrete boundaries (**guards**) to group together overlapping files



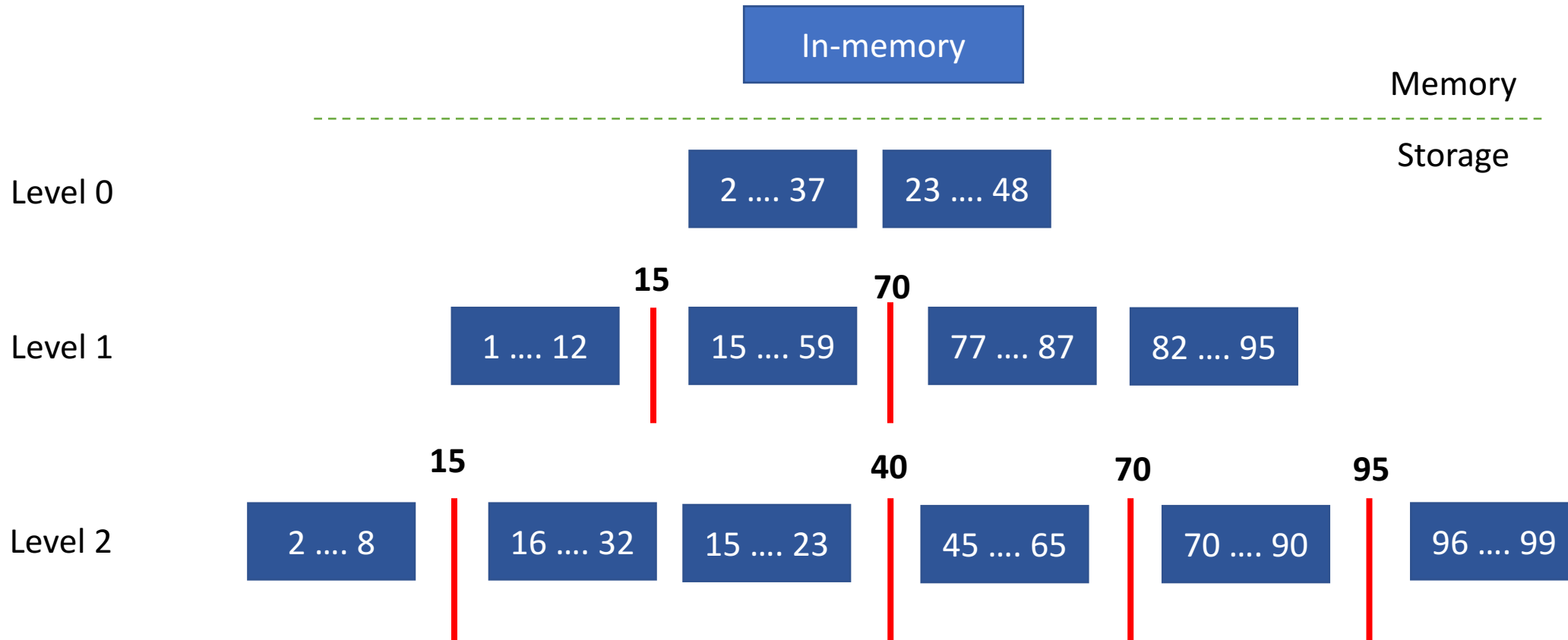(files of same color can have overlapping key ranges)

# Fragmented Log-Structured Merge Tree

Novel modification of LSM data structure

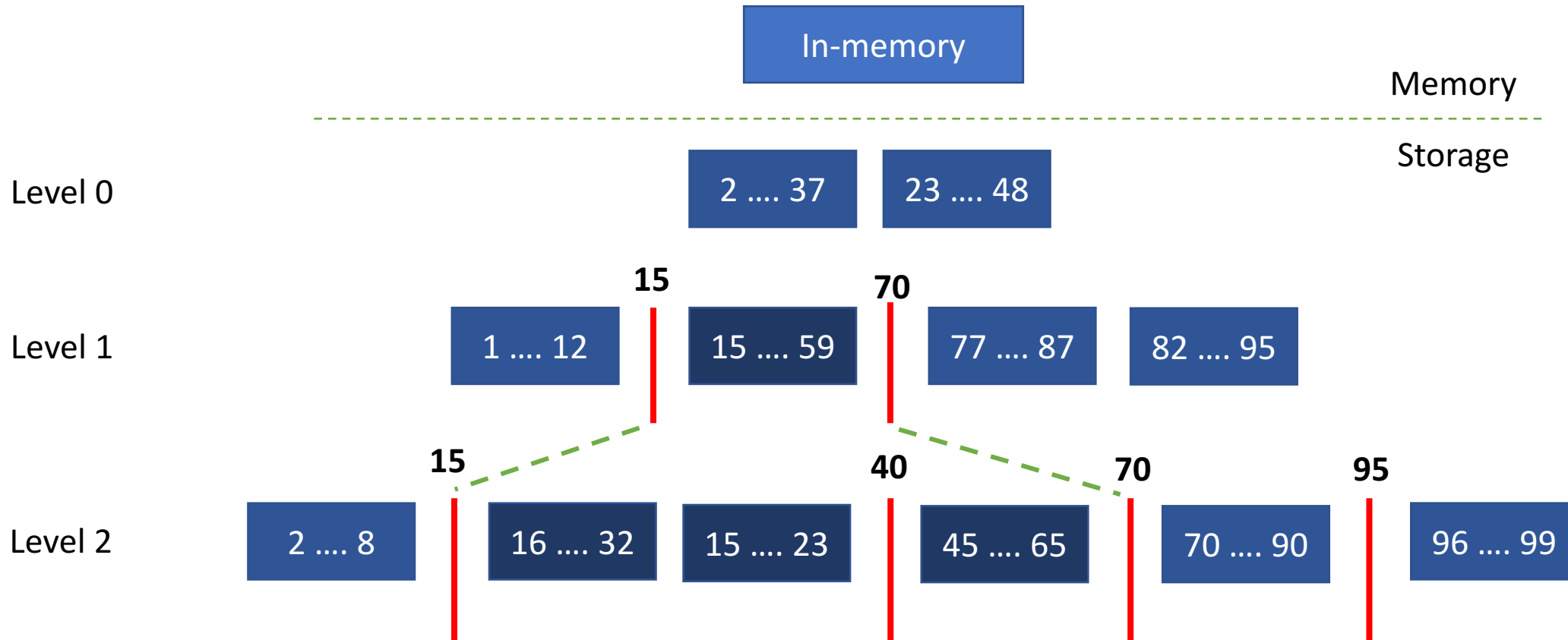Uses guards to maintain partially sorted levels

Writes data only once per level in most cases

# FLSM structure
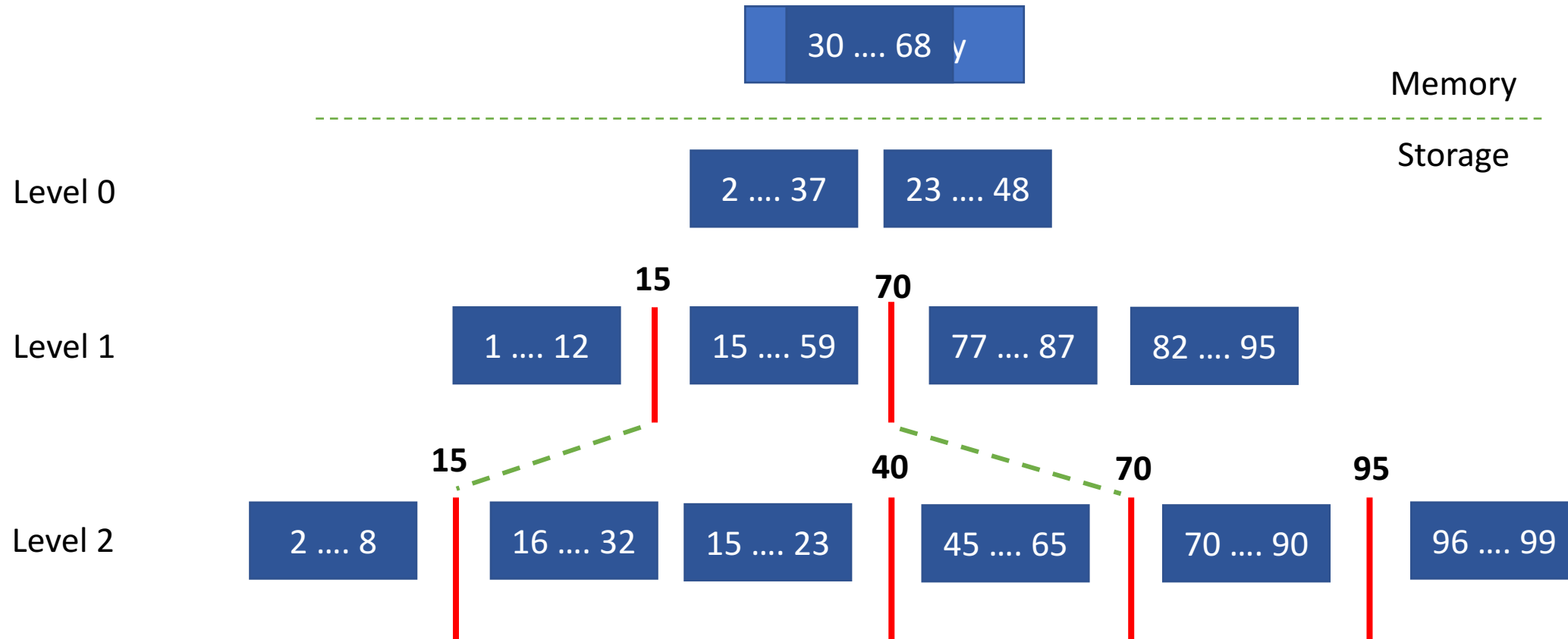


Note how files are logically grouped within guards

# FLSM structure



Guards get more fine grained deeper into the tree

# How does FLSM reduce write amplification?

# How does FLSM reduce write amplification?



Memory

Storage

Level 0

| 2 …. 37 | 23 …. 48 |

**15**    **70**

Level 1

| 1 …. 12 | 15 …. 59 | 77 …. 87 | 82 …. 95 |

**15**    **40**    **70**    **95**

Level 2

| 2 …. 8 | 16 …. 32 | 15 …. 23 | 45 …. 65 | 70 …. 90 | 96 …. 99 |

30 …. 68   y

Max files in level 0 is configured to be 2

# How does FLSM reduce write amplification?
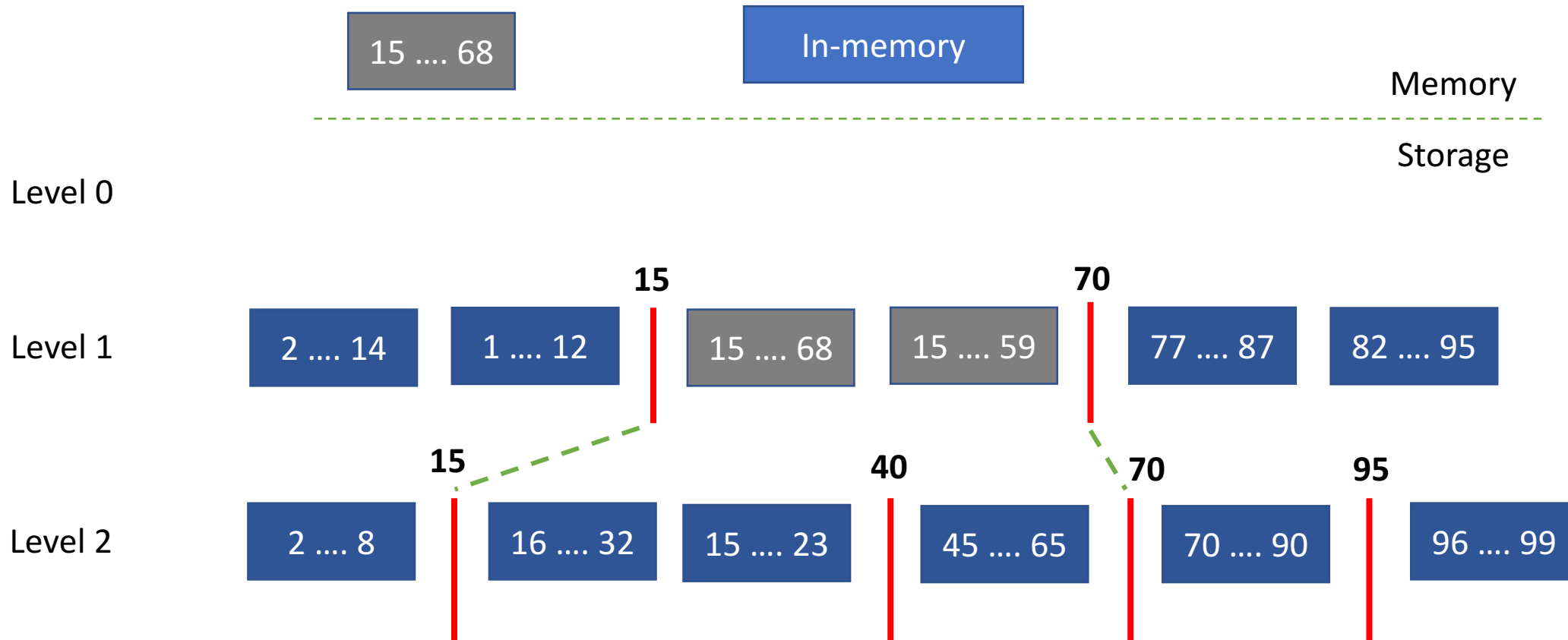


Compacting level 0
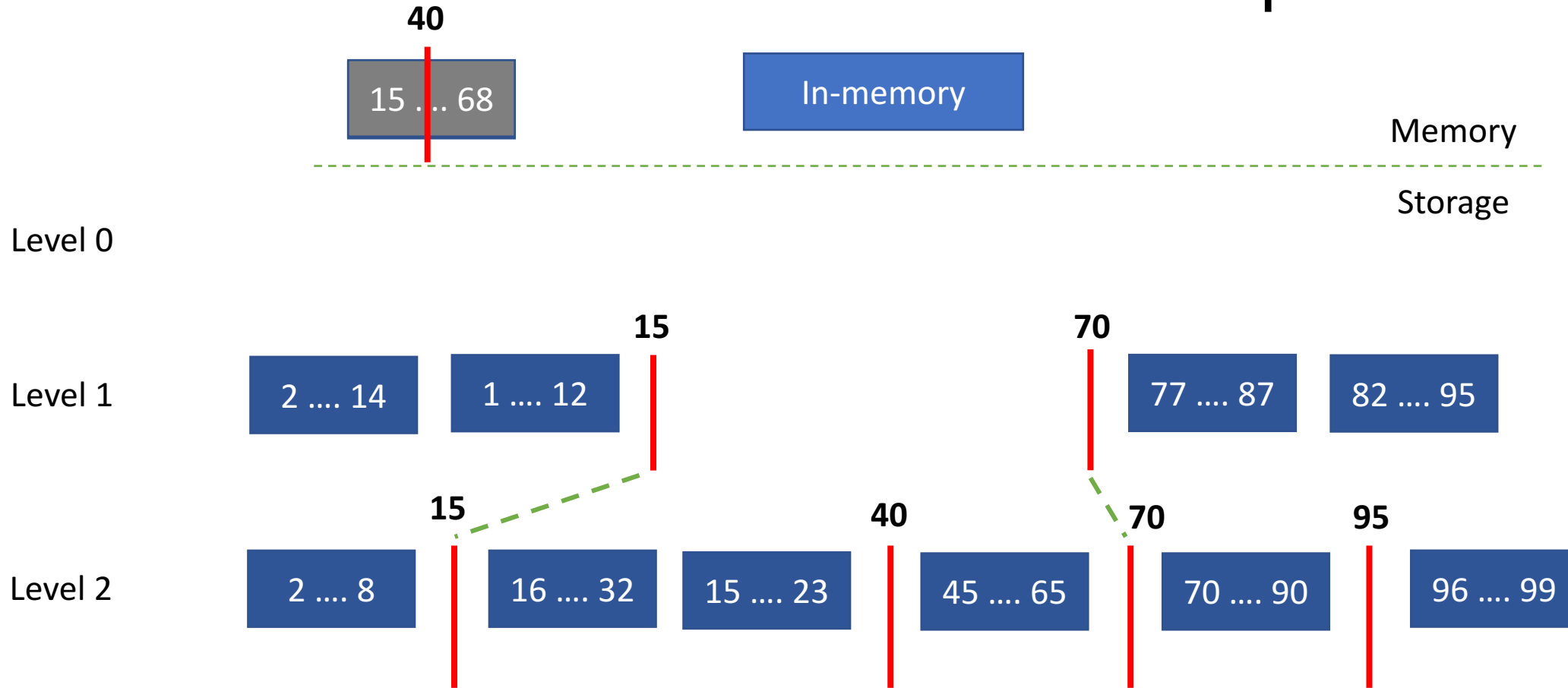
# How does FLSM reduce write amplification?



Fragmented files are just appended to next level

# How does FLSM reduce write amplification?



Guard 15 in Level 1 is to be compacted

# How does FLSM reduce write amplification?

**40**

| 15 .... 68 |

In-memory

Memory

- - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Storage

Level 0

**15**                              **70**

Level 1

| 2 .... 14 |  | 1 .... 12 |                              | 77 .... 87 |  | 82 .... 95 |

**15**                    **40**              **70**            **95**

Level 2

| 2 .... 8 |  | 16 .... 32 |  | 15 .... 23 |   | 45 .... 65 |  | 70 .... 90 |   | 96 .... 99 |

Files are combined, sorted and fragmented

# How does FLSM reduce write amplification?



Fragmented files are just appended to next level

# How does FLSM reduce write amplification?

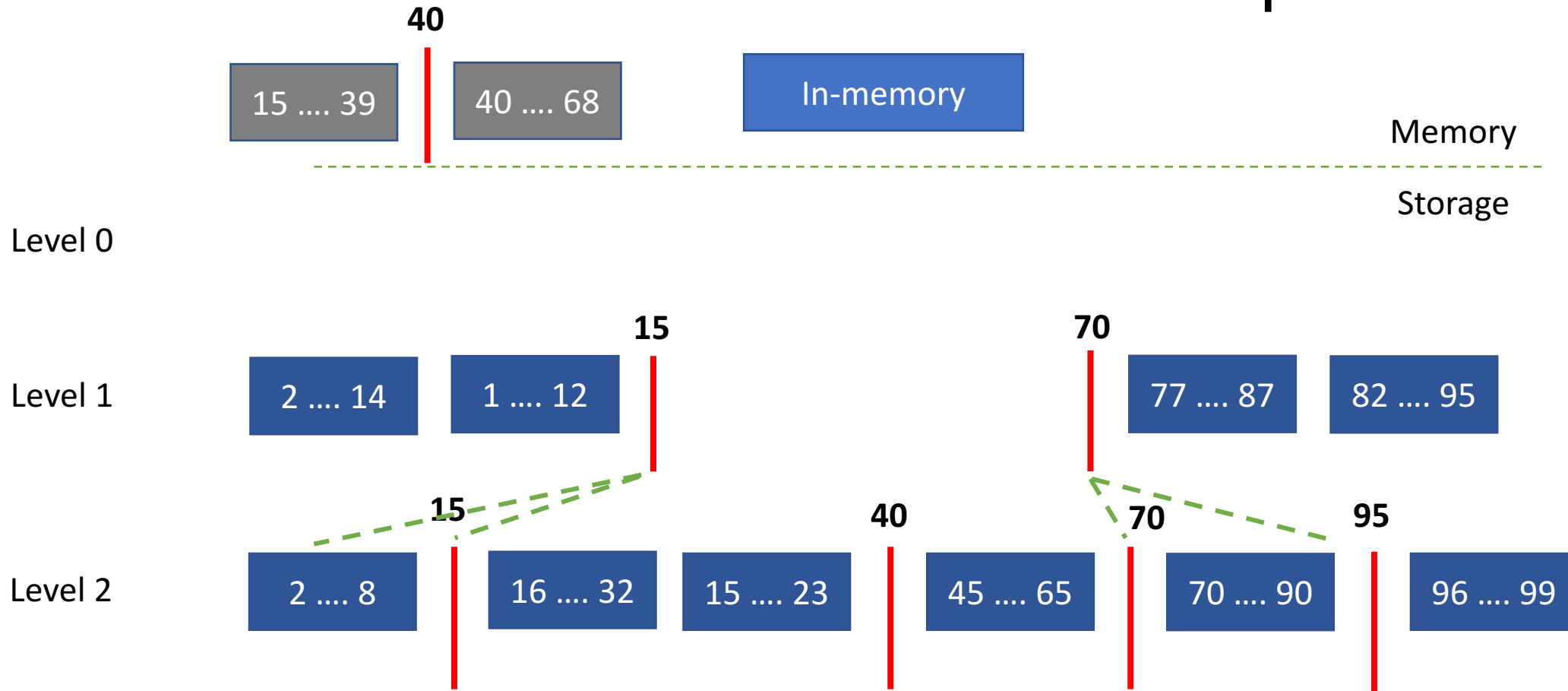FLSM <span style="color:red">doesn't re-write data</span> to the same level
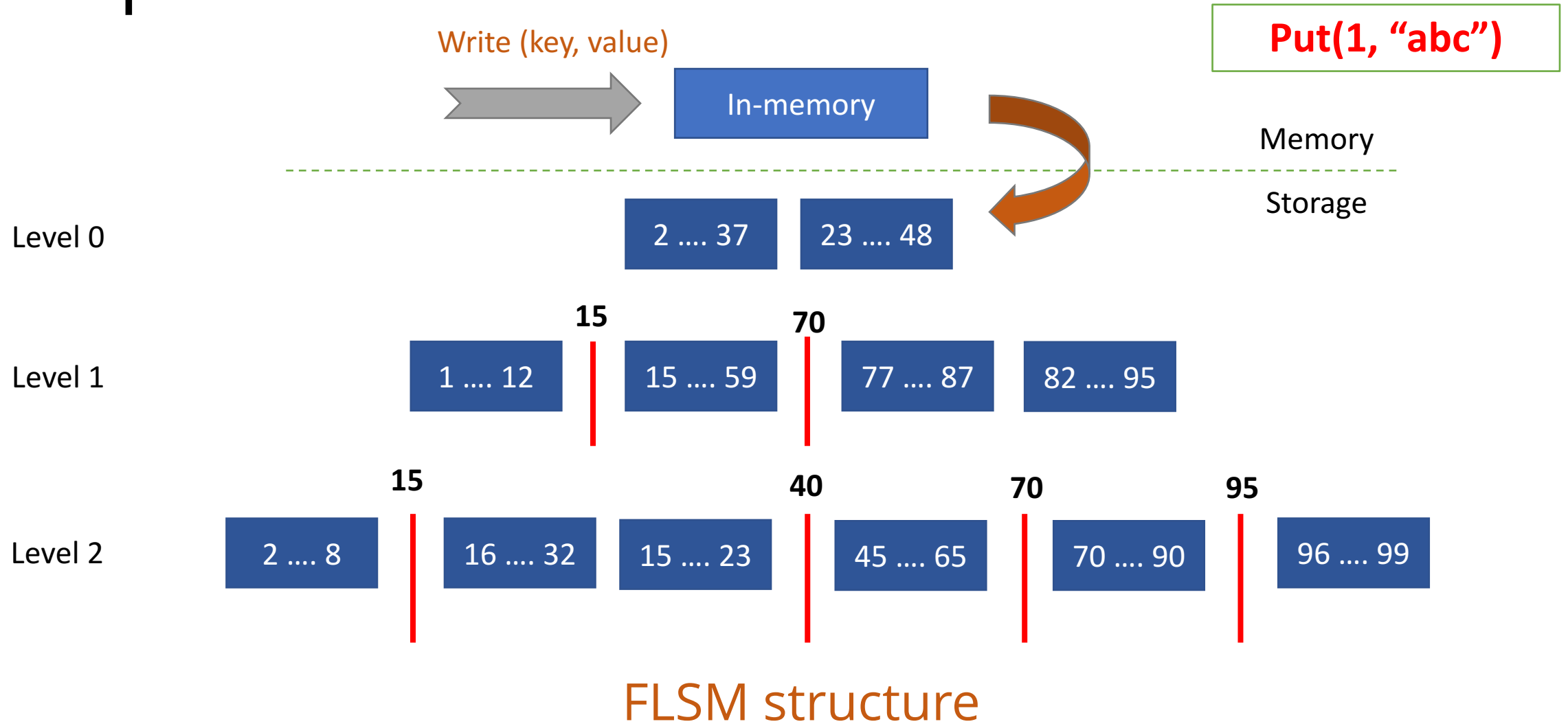in most cases

# How does FLSM maintain read performance?

FLSM maintains <span style="color:red">partially sorted levels</span> to efficiently
reduce the search space

# Selecting Guards

- Guards are chosen randomly and dynamically
- Dependent on the distribution of data

# Operations: Write

Write (key, value)

Put(1, "abc")

In-memory

Memory

Storage

Level 0     2 …. 37    23 …. 48

**15**     **70**

Level 1     1 …. 12    15 …. 59    77 …. 87    82 …. 95

**15**     **40**     **70**     **95**

Level 2     2 …. 8    16 …. 32    15 …. 23    45 …. 65    70 …. 90    96 …. 99

## FLSM structure

54

# Operations: Get



Get(23)

In-memory

Memory

Storage

Level 0

| 2 …. 37 | 23 …. 48 |

**15** **70**

Level 1

| 1 …. 12 | 15 …. 59 | 77 …. 87 | 82 …. 95 |

**15** **40** **70** **95**

Level 2

| 2 …. 8 | 16 …. 32 | 15 …. 23 | 45 …. 65 | 70 …. 90 | 96 …. 99 |

FLSM structure

# Operations: Get

# Operations: Get

Get(23)

In-memory

Memory
- - - - - - - - - - - - - - - - - - - - - - - - - - -
Storage

Level 0

| 2 …. 37 | 23 …. 48 |

**15**          **70**

Level 1

| 1 …. 12 | 15 …. 59 | 77 …. 87 | 82 …. 95 |

**15**          **40**      **70**      **95**

Level 2

| 2 …. 8 | 16 …. 32 | 15 …. 23 | 45 …. 65 | 70 …. 90 | 96 …. 99 |

All level 0 files need to be searched

57

# Operations: Get



Get(23)

In-memory

Memory

Storage

Level 0    2 …. 37    23 …. 48

**15**    **70**

Level 1    1 …. 12    15 …. 59    77 …. 87    82 …. 95

**15**    **40**    **70**    **95**

Level 2    2 …. 8    16 …. 32    15 …. 23    45 …. 65    70 …. 90    96 …. 99

Level 1: File under guard 15 is searched

# Operations: Get

In-memory

Memory

Storage

Level 0

| 2 .... 37 | 23 .... 48 |

**15**      **70**

Level 1

| 1 .... 12 | 15 .... 59 | 77 .... 87 | 82 .... 95 |

**15**      **40**      **70**      **95**

Level 2

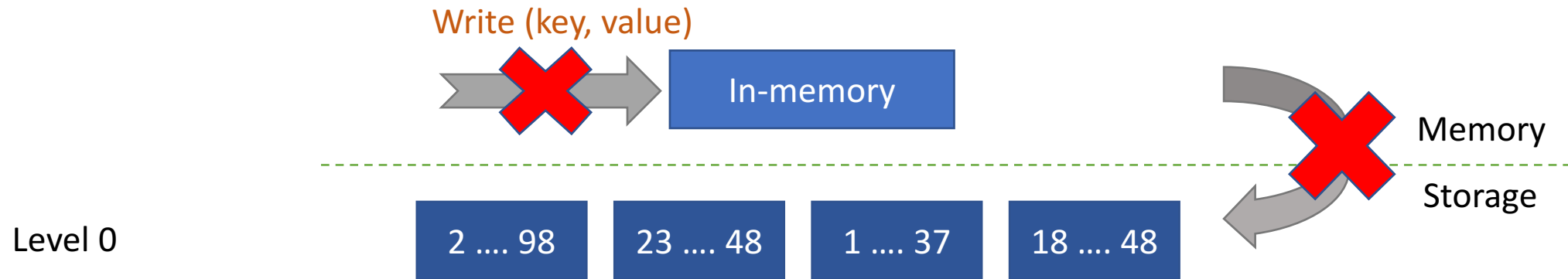| 2 .... 8 | 16 .... 32 | 15 .... 23 | 45 .... 65 | 70 .... 90 | 96 .... 99 |

Level 2: Both the files under guard 15 are searched

59

# High write throughput in FLSM

- Compaction from memory to level 0 is stalled
- Writes to memory is also stalled

Write (key, value)

In-memory

Memory

Storage

Level 0

| 2 …. 98 | 23 …. 48 | 1 …. 37 | 18 …. 48 |

If rate of insertion is higher than rate of compaction, write throughput depends on the rate of compaction

# High write throughput in FLSM

- Compaction from memory to level 0 is stalled
- Writes to memory is also stalled

**FLSM has faster compaction because of lesser I/O and hence higher write throughput**

If rate of insertion is higher than rate of compaction, write throughput depends on the rate of compaction

# Challenges in FLSM

- Every read/range query operation needs to examine multiple files per level

- For example, if every guard has 5 files, read latency is increased by 5x (assuming no cache hits)

Trade-off between write I/O and read performance

# Outline

- Log-Structured Merge Tree (LSM)
- Fragmented Log-Structured Merge Tree (FLSM)
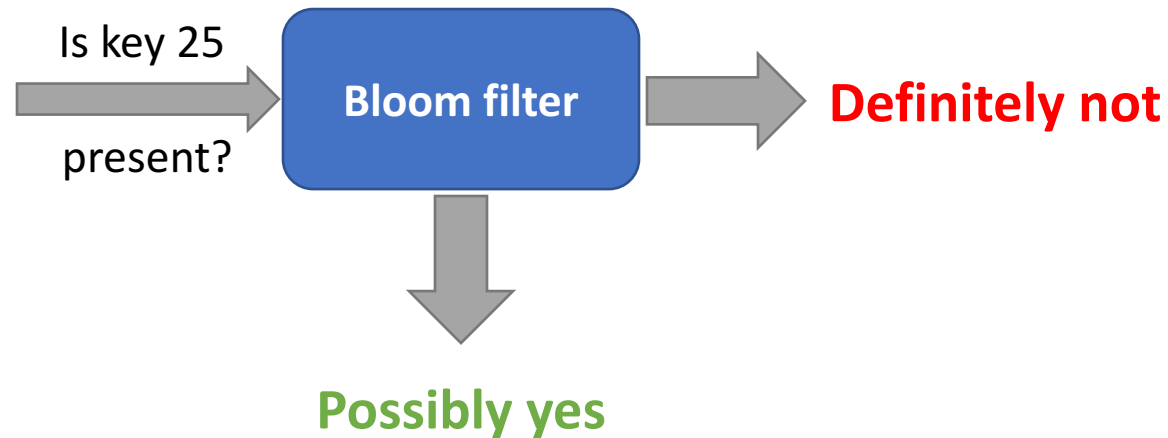- **Building PebblesDB using FLSM**
- Evaluation
- Conclusion

# PebblesDB

- Built by modifying **HyperLevelDB** (±9100 LOC) to use FLSM
- HyperLevelDB, built over LevelDB, to provide improved parallelism and compaction
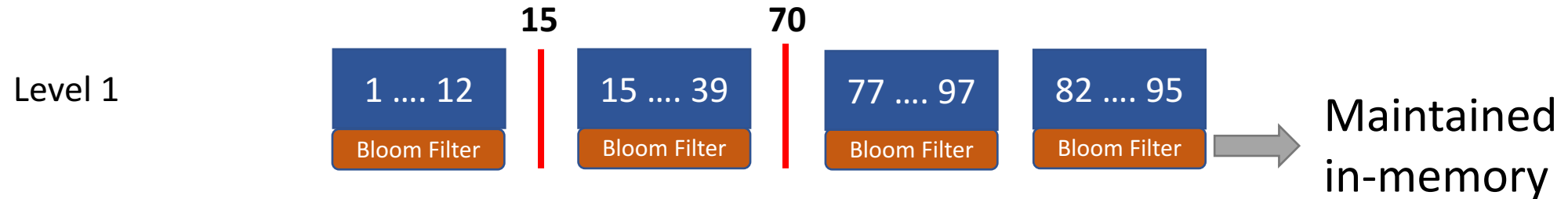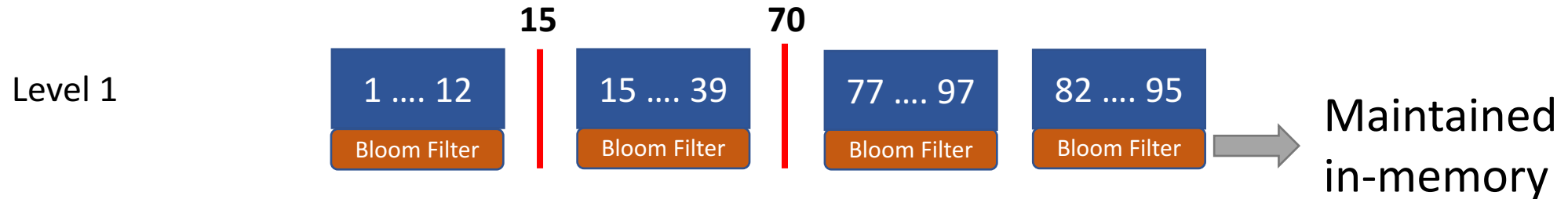- API compatible with LevelDB, but not with RocksDB

# Optimizations in PebblesDB

- **Challenge (get/range query):** Multiple files in a guard
- Get() performance is improved using <span style="color:red">file level bloom filter</span>

# Optimizations in PebblesDB

- **Challenge (get/range query):** Multiple files in a guard
- Get() performance is improved using file level bloom filter

Is key 25 present? → **Bloom filter** → **Definitely not**

**Possibly yes**

# Optimizations in PebblesDB

- **Challenge (get/range query):** Multiple files in a guard
- Get() performance is improved using file level bloom filter

# Optimizations in PebblesDB

- **Challenge (get/range query):** Multiple files in a guard
- Get() performance is improved using file level bloom filter



PebblesDB reads same number of files as any LSM based store

# Optimizations in PebblesDB

- **Challenge (get/range query):** Multiple files in a guard
- Get() performance is improved using <span style="color:red">file level bloom filter</span>
- Range query performance is improved using parallel threads and better compaction

# Outline

- Log-Structured Merge Tree (LSM)
- Fragmented Log-Structured Merge Tree (FLSM)
- Building PebblesDB using FLSM
- **Evaluation**
- Conclusion

# Evaluation

Micro-benchmarks

Real world workloads - YCSB
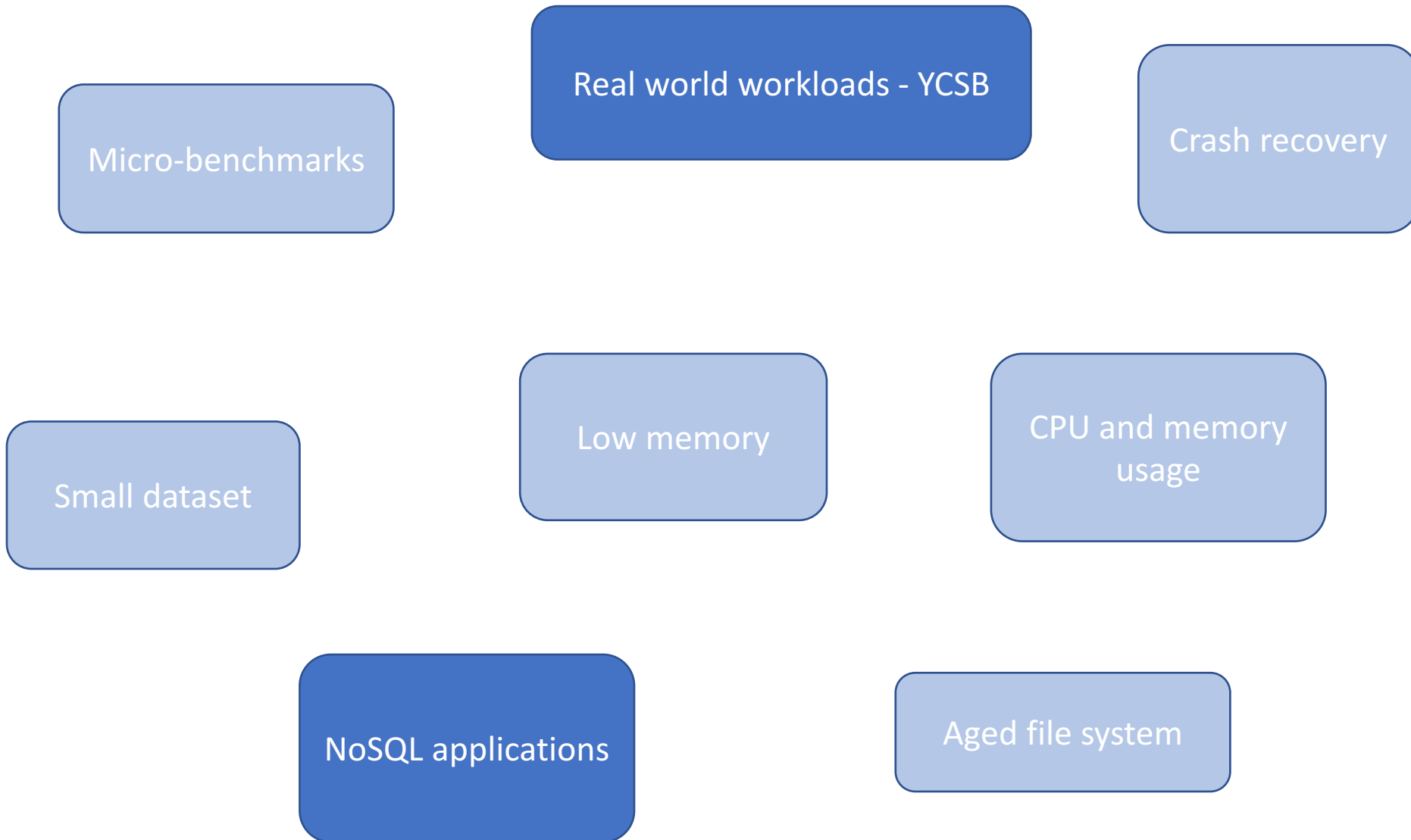
Crash recovery

Small dataset

Low memory

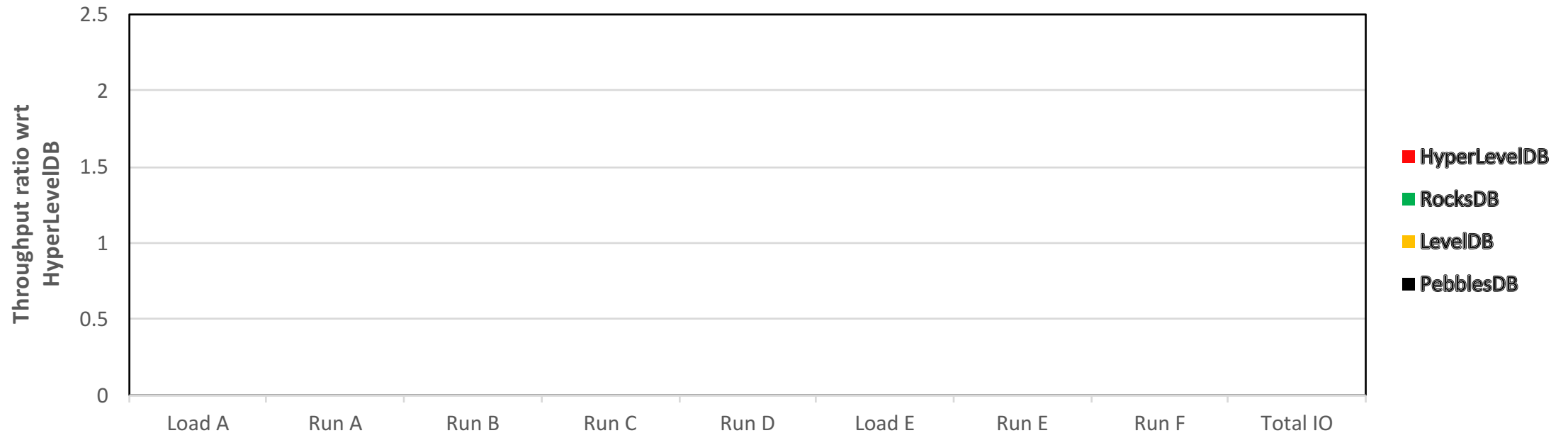CPU and memory usage

NoSQL applications

Aged file system

# Evaluation

Micro-benchmarks

Real world workloads - YCSB

Crash recovery

Small dataset

Low memory

CPU and memory usage

NoSQL applications

Aged file system

# Real world workloads - YCSB

- Yahoo! Cloud Serving Benchmark - Industry standard macro-benchmark
- Insertions: 50M, Operations: 10M, key size: 16 bytes and value size: 1 KB



Load A - 100 % writes      Run D    - 95% reads (latest), 5% writes
Run A    - 50% reads, 50% writes      Load E   - 100% writes
Run B    - 95% reads, 5% writes      Run E    - 95% range queries, 5% writes
Run C    - 100% reads      Run F    - 50% reads, 50% read-modify-writes

# Real world workloads - YCSB

- Yahoo! Cloud Serving Benchmark - Industry standard macro-benchmark
- Insertions: 50M, Operations: 10M, key size: 16 bytes and value size: 1 KB



Load A - 100 % writes
Run A   - 50% reads, 50% writes
Run B   - 95% reads, 5% writes
Run C   - 100% reads
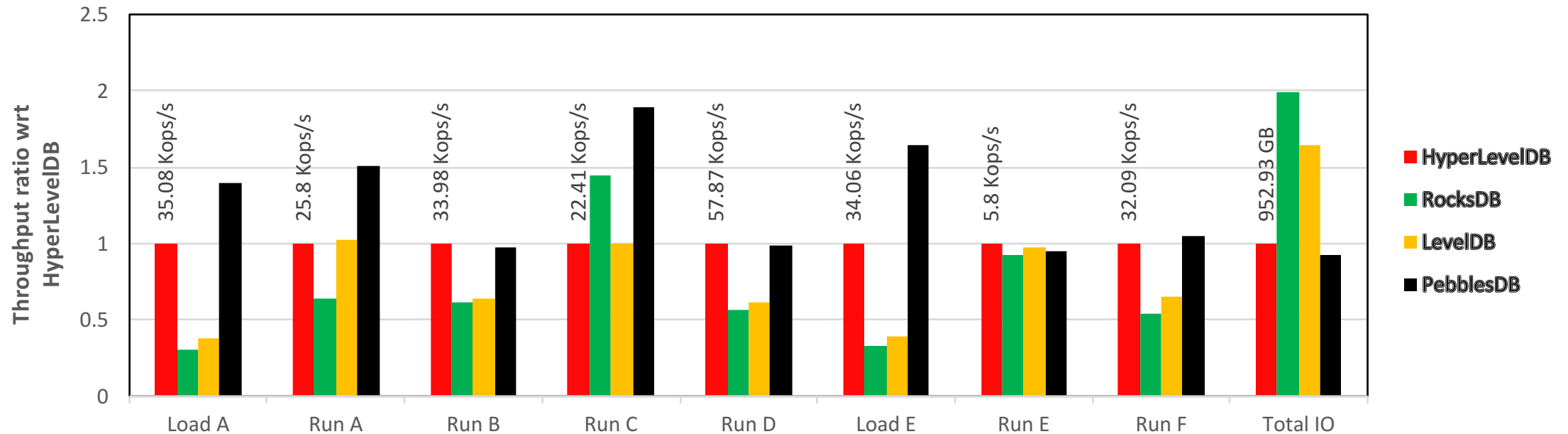
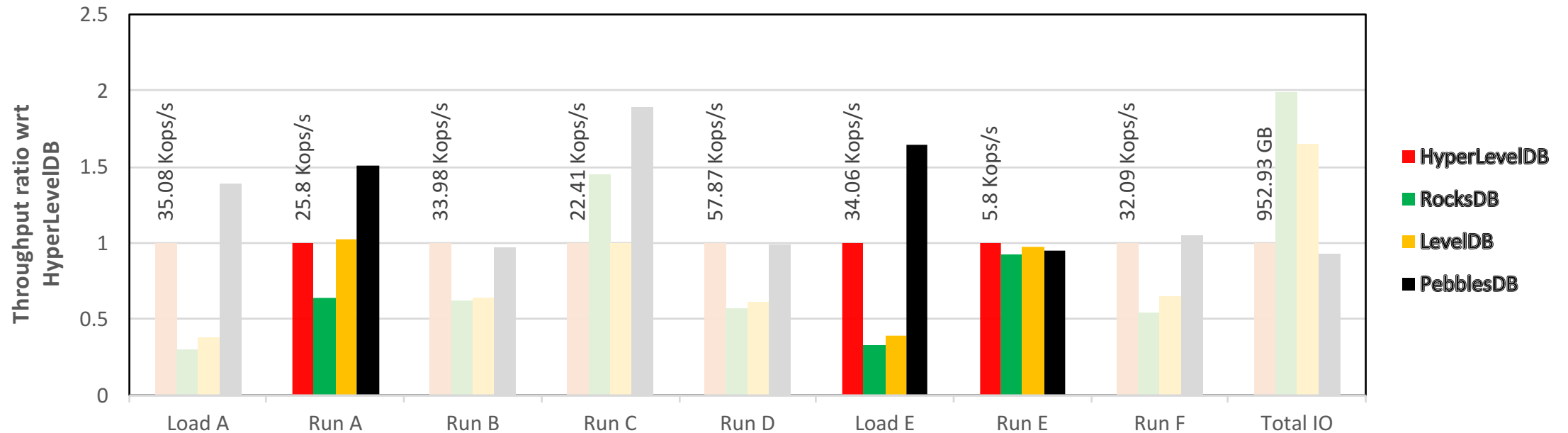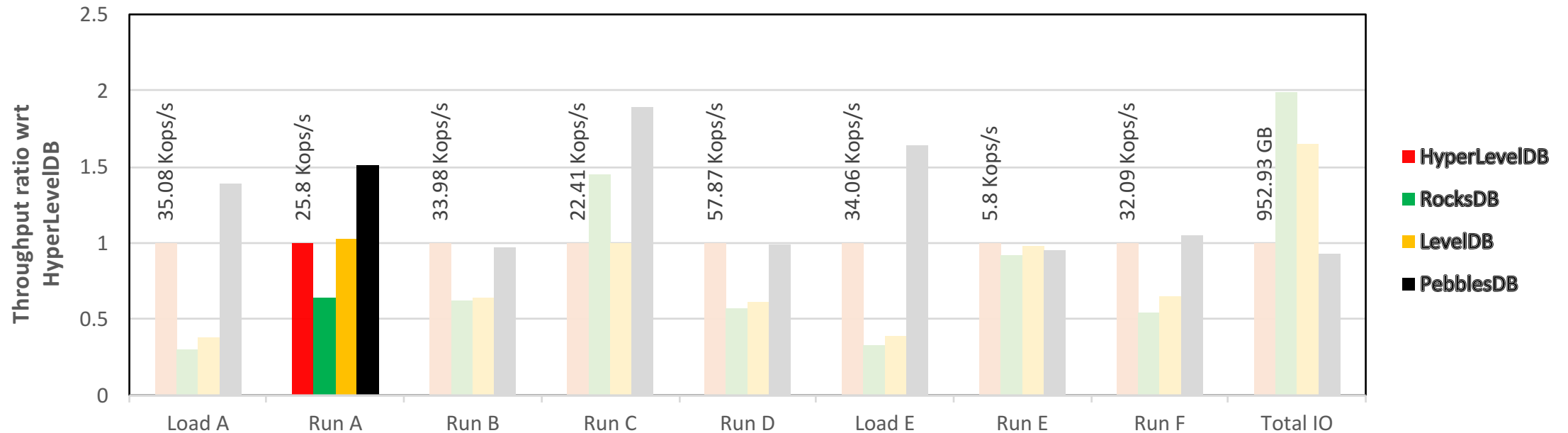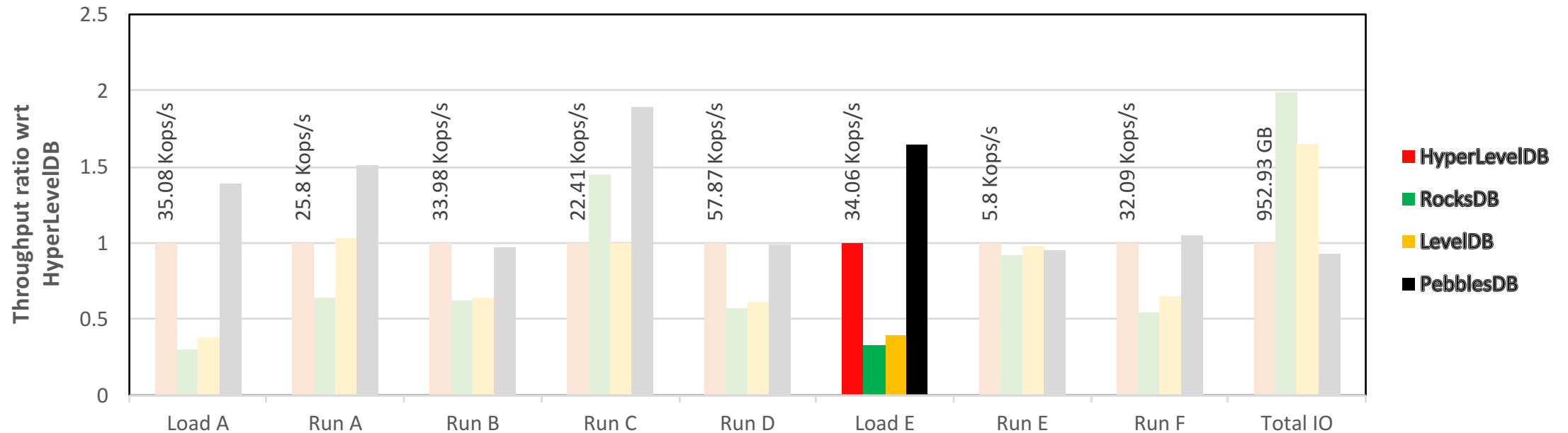Run D   - 95% reads (latest), 5% writes
Load E  - 100% writes
Run E   - 95% range queries, 5% writes
Run F   - 50% reads, 50% read-modify-writes

# Real world workloads - YCSB

- Yahoo! Cloud Serving Benchmark - Industry standard macro-benchmark
- Insertions: 50M, Operations: 10M, key size: 16 bytes and value size: 1 KB



Legend:
- HyperLevelDB (red)
- RocksDB (green)
- LevelDB (yellow)
- PebblesDB (black)

Y-axis: Throughput ratio wrt HyperLevelDB

Categories (with throughput labels):
- Load A — 35.08 Kops/s
- Run A — 25.8 Kops/s
- Run B — 33.98 Kops/s
- Run C — 22.41 Kops/s
- Run D — 57.87 Kops/s
- Load E — 34.06 Kops/s
- Run E — 5.8 Kops/s
- Run F — 32.09 Kops/s
- Total IO — 952.93 GB

Load A - 100 % writes
**Run A** - 50% reads, 50% writes
Run B - 95% reads, 5% writes
Run C - 100% reads

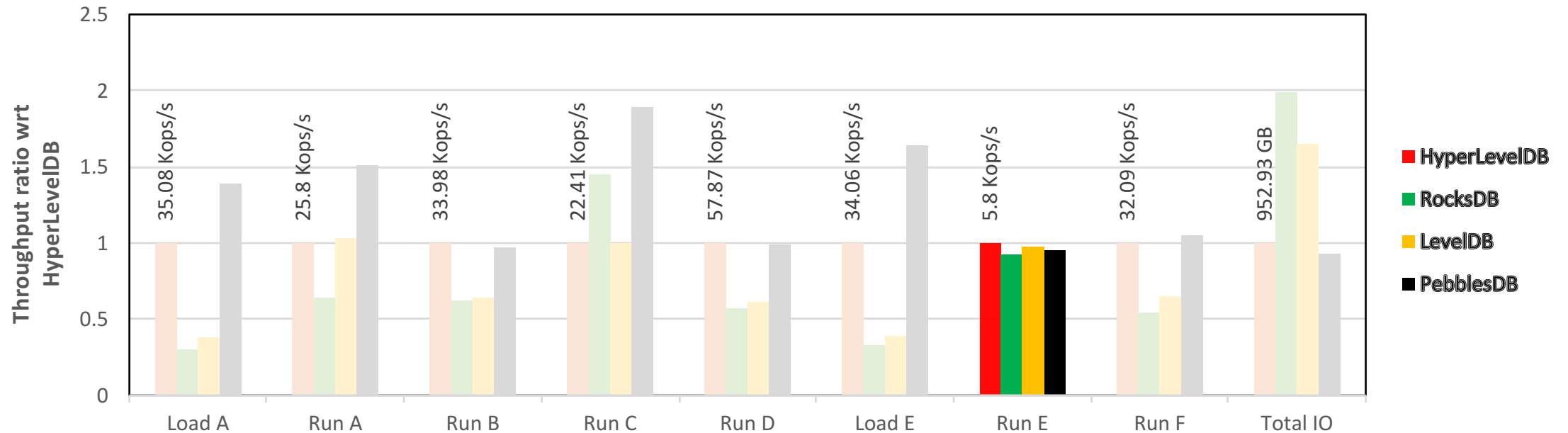Run D - 95% reads (latest), 5% writes
**Load E** - 100% writes
**Run E** - 95% range queries, 5% writes
Run F - 50% reads, 50% read-modify-writes

# Real world workloads - YCSB

- Yahoo! Cloud Serving Benchmark - Industry standard macro-benchmark
- Insertions: 50M, Operations: 10M, key size: 16 bytes and value size: 1 KB



| | |
|---|---|
| Load A - 100 % writes | Run D  - 95% reads (latest), 5% writes |
| **Run A**  - 50% reads, 50% writes | Load E  - 100% writes |
| Run B   - 95% reads, 5% writes | Run E   - 95% range queries, 5% writes |
| Run C   - 100% reads | Run F   - 50% reads, 50% read-modify-writes |

# Real world workloads - YCSB

- Yahoo! Cloud Serving Benchmark - Industry standard macro-benchmark
- Insertions: 50M, Operations: 10M, key size: 16 bytes and value size: 1 KB



Load A - 100 % writes      Run D   - 95% reads (latest), 5% writes

Run A   - 50% reads, 50% writes     **Load E**   - 100% writes

Run B   - 95% reads, 5% writes     Run E   - 95% range queries, 5% writes

Run C   - 100% reads      Run F   - 50% reads, 50% read-modify-writes

# Real world workloads - YCSB

- Yahoo! Cloud Serving Benchmark - Industry standard macro-benchmark
- Insertions: 50M, Operations: 10M, key size: 16 bytes and value size: 1 KB



Load A - 100 % writes

Run A  - 50% reads, 50% writes

Run B  - 95% reads, 5% writes

Run C  - 100% reads
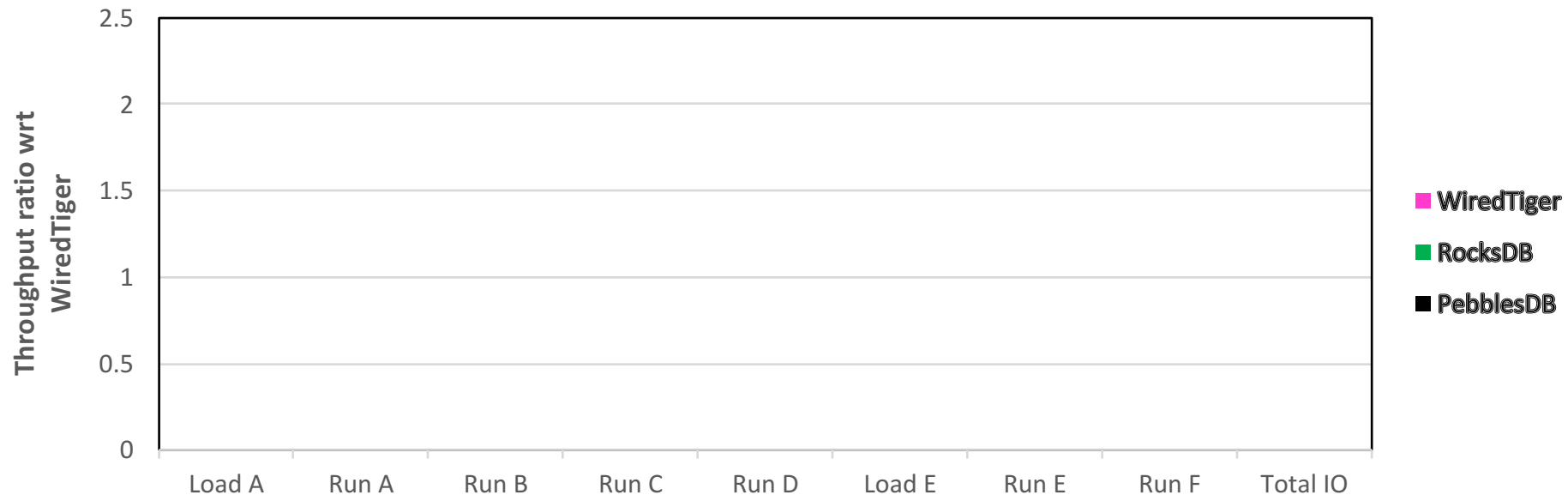
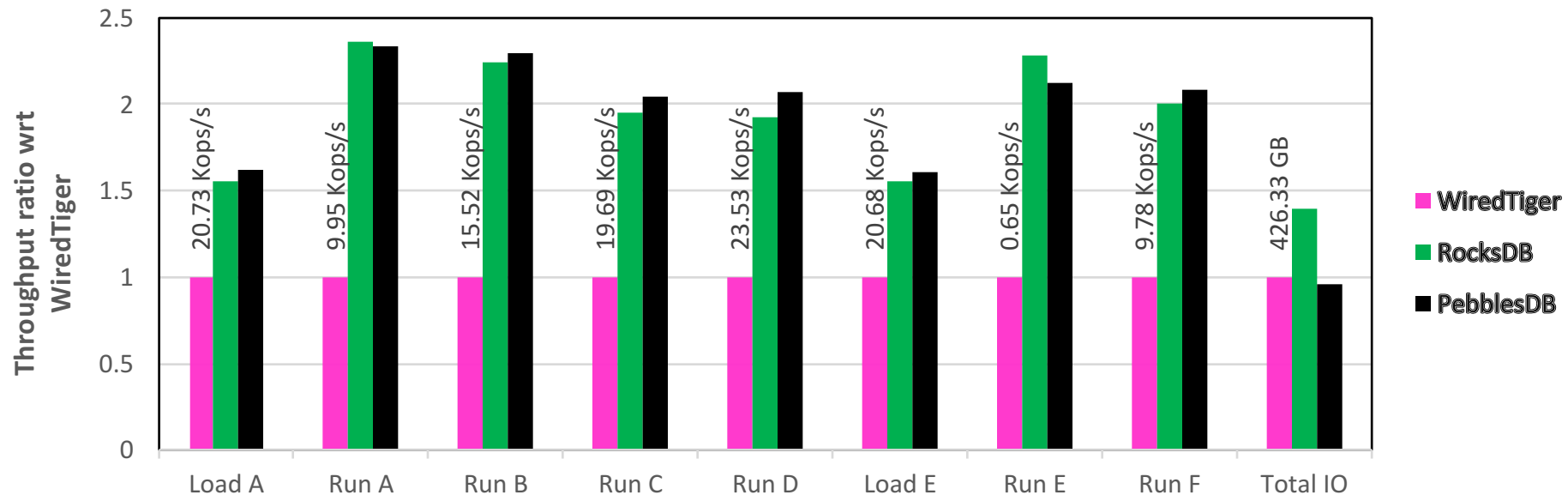Run D   - 95% reads (latest), 5% writes

Load E  - 100% writes

Run E   - 95% range queries, 5% writes

Run F   - 50% reads, 50% read-modify-writes

# NoSQL stores - MongoDB

- YCSB on MongoDB, a widely used key-value store
- Inserted 20M key-value pairs with 1 KB value size and 10M operations



Chart legend:
- WiredTiger
- RocksDB
- PebblesDB

Y-axis: Throughput ratio wrt WiredTiger (0, 0.5, 1, 1.5, 2, 2.5)

X-axis: Load A, Run A, Run B, Run C, Run D, Load E, Run E, Run F, Total IO

Load A - 100 % writes
Run A  - 50% reads, 50% writes
Run B  - 95% reads, 5% writes
Run C  - 100% reads
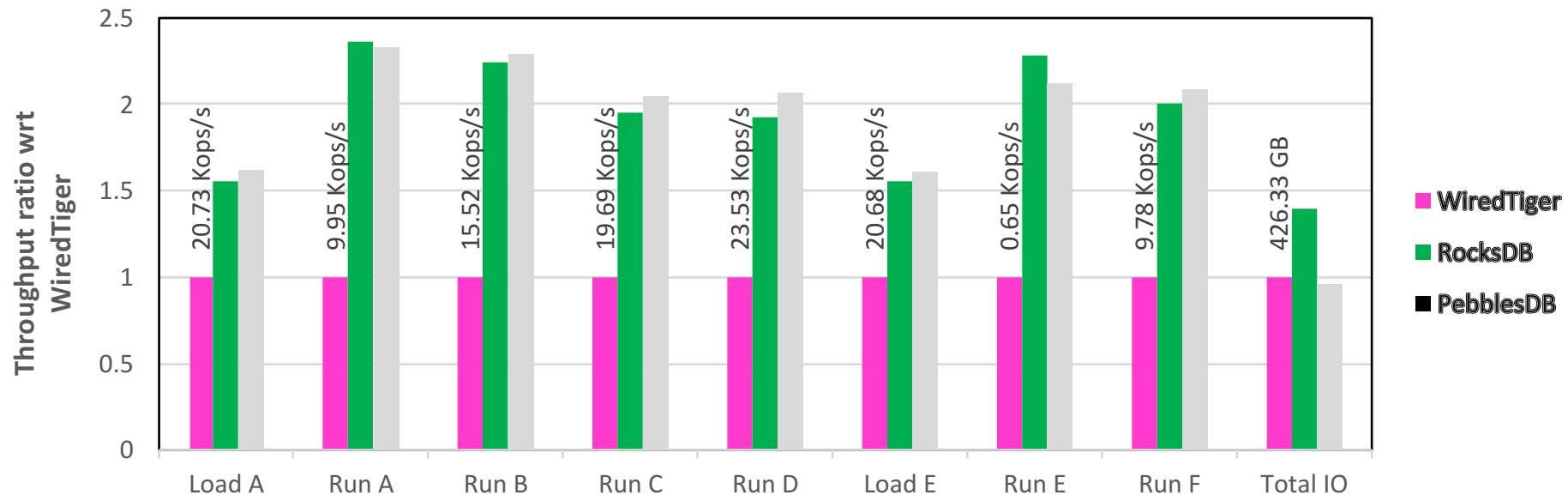
Run D  - 95% reads (latest), 5% writes
Load E  - 100% writes
Run E  - 95% range queries, 5% writes
Run F  - 50% reads, 50% read-modify-writes

# NoSQL stores - MongoDB

- YCSB on MongoDB, a widely used key-value store
- Inserted 20M key-value pairs with 1 KB value size and 10M operations



Throughput ratio wrt WiredTiger

Legend: WiredTiger, RocksDB, PebblesDB

Bar labels (left to right):
- Load A — 20.73 Kops/s
- Run A — 9.95 Kops/s
- Run B — 15.52 Kops/s
- Run C — 19.69 Kops/s
- Run D — 23.53 Kops/s
- Load E — 20.68 Kops/s
- Run E — 0.65 Kops/s
- Run F — 9.78 Kops/s
- Total IO — 426.33 GB

Load A - 100 % writes
Run A   - 50% reads, 50% writes
Run B   - 95% reads, 5% writes
Run C   - 100% reads
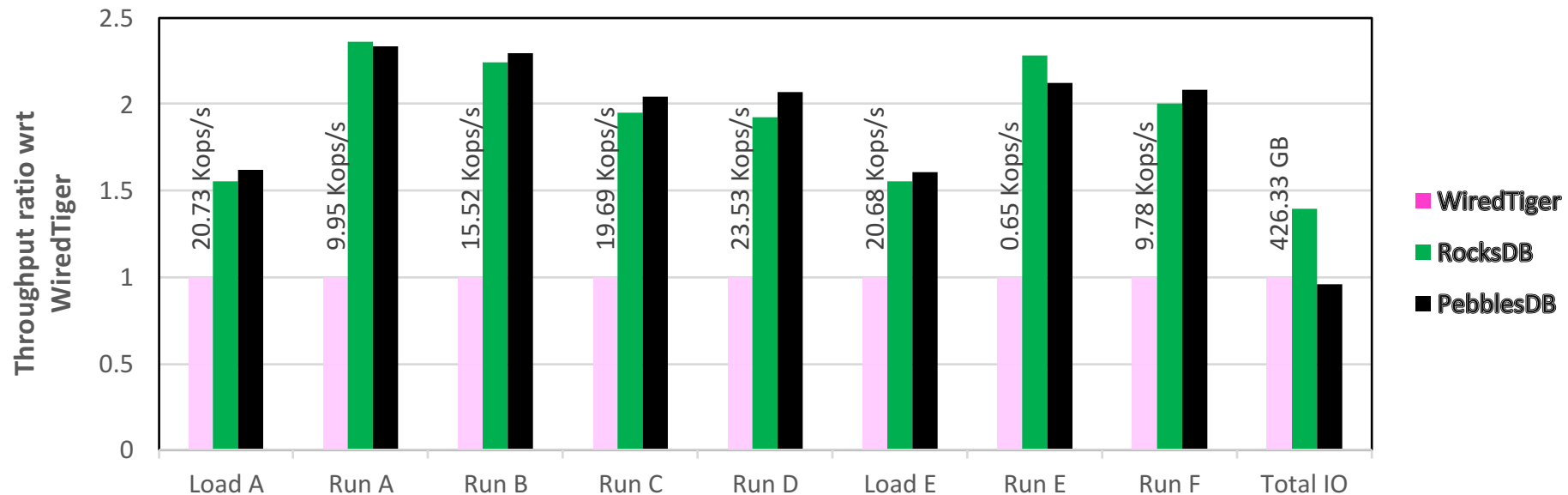
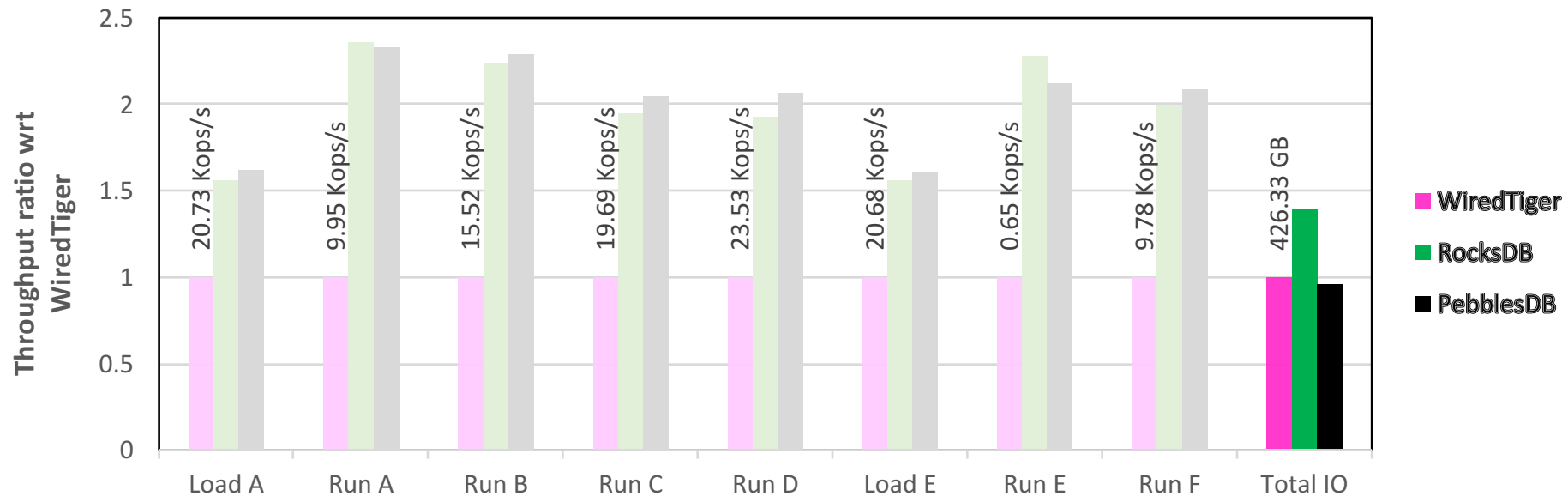Run D   - 95% reads (latest), 5% writes
Load E  - 100% writes
Run E   - 95% range queries, 5% writes
Run F   - 50% reads, 50% read-modify-writes

# NoSQL stores - MongoDB

- YCSB on MongoDB, a widely used key-value store
- Inserted 20M key-value pairs with 1 KB value size and 10M operations



Load A - 100 % writes          Run D   - 95% reads (latest), 5% writes
Run A   - 50% reads, 50% writes     Load E  - 100% writes
Run B   - 95% reads, 5% writes      Run E   - 95% range queries, 5% writes
Run C   - 100% reads           Run F   - 50% reads, 50% read-modify-writes

# NoSQL stores - MongoDB

- YCSB on MongoDB, a widely used key-value store
- Inserted 20M key-value pairs with 1 KB value size and 10M operations



Throughput ratio wrt WiredTiger

- 20.73 Kops/s (Load A)
- 9.95 Kops/s (Run A)
- 15.52 Kops/s (Run B)
- 19.69 Kops/s (Run C)
- 23.53 Kops/s (Run D)
- 20.68 Kops/s (Load E)
- 0.65 Kops/s (Run E)
- 9.78 Kops/s (Run F)
- 426.33 GB (Total IO)

Legend:
- WiredTiger
- RocksDB
- PebblesDB

Load A - 100 % writes
Run A   - 50% reads, 50% writes
Run B   - 95% reads, 5% writes
Run C   - 100% reads

Run D   - 95% reads (latest), 5% writes
Load E  - 100% writes
Run E   - 95% range queries, 5% writes
Run F   - 50% reads, 50% read-modify-writes

# NoSQL stores - MongoDB

- YCSB on MongoDB, a widely used key-value store
- Inserted 20M key-value pairs with 1 KB value size and 10M operations



Load A - 100 % writes
Run A   - 50% reads, 50% writes
Run B   - 95% reads, 5% writes
Run C   - 100% reads

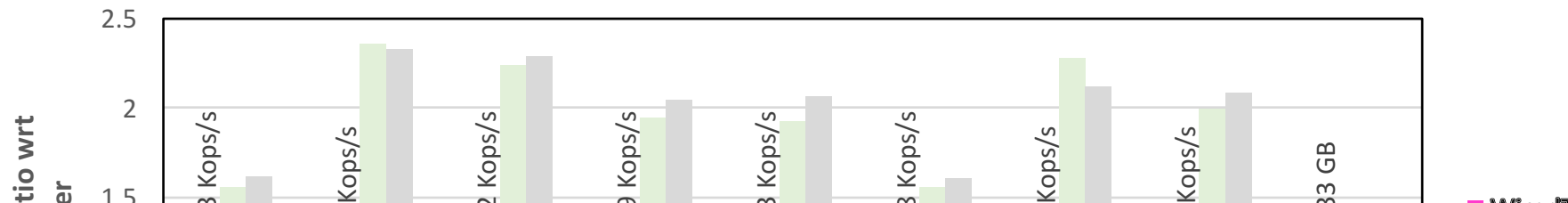Run D   - 95% reads (latest), 5% writes
Load E  - 100% writes
Run E   - 95% range queries, 5% writes
Run F   - 50% reads, 50% read-modify-writes

# NoSQL stores - MongoDB

- YCSB on MongoDB, a widely used key-value store
- Inserted 20M key-value pairs with 1 KB value size and 10M operations



**PebblesDB combines low write IO of WiredTiger with high performance of RocksDB**

Load A - 100 % writes
Run A   - 50% reads, 50% writes
Run B   - 95% reads, 5% writes
Run C   - 100% reads

Run D   - 95% reads (latest), 5% writes
Load E  - 100% writes
Run E   - 95% range queries, 5% writes
Run F   - 50% reads, 50% read-modify-writes

# Outline

- Log-Structured Merge Tree (LSM)
- Fragmented Log-Structured Merge Tree (FLSM)
- Building PebblesDB using FLSM
- Evaluation
- **Conclusion**

# Conclusion

- PebblesDB: key-value store built on Fragmented Log-Structured Merge Trees
  - Increases write throughput and reduces write IO at the same time
  - Obtains 6X the write throughput of RocksDB
- As key-value stores become more widely used, there have been several attempts to optimize them
- PebblesDB combines algorithmic innovation (the FLSM data structure) with careful systems building

# https://github.com/utsaslab/pebblesdb

https://github.com/utsaslab/pebblesdb



**Thank You!**

# Backup slides

# Operations: Seek

- **Seek(target):** Returns the smallest key in the database which is >= target
- Used for range queries (for example, return all entries between 5 and 18)

Level 0   –   1, 2, 100, 1000

Level 1   –   1, 5, 10, 2000

Level 2   –   5, 300, 500

Get(1)

# Operations: Seek

- **Seek(target):** Returns the smallest key in the database which is >= target
- Used for range queries (for example, return all entries between 5 and 18)
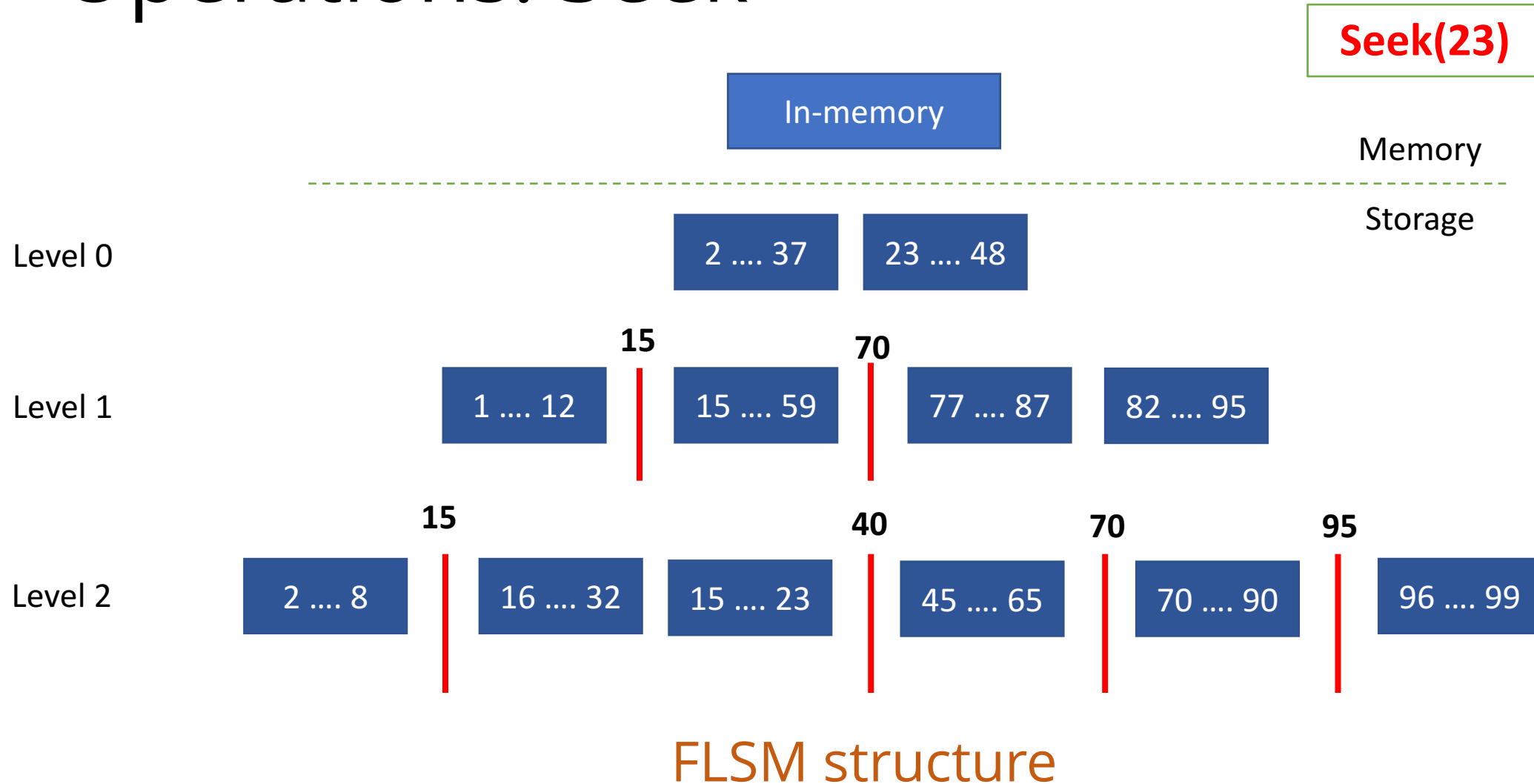
Level 0  –  1, 2, 100, 1000

Level 1  –  1, 5, 10, 2000
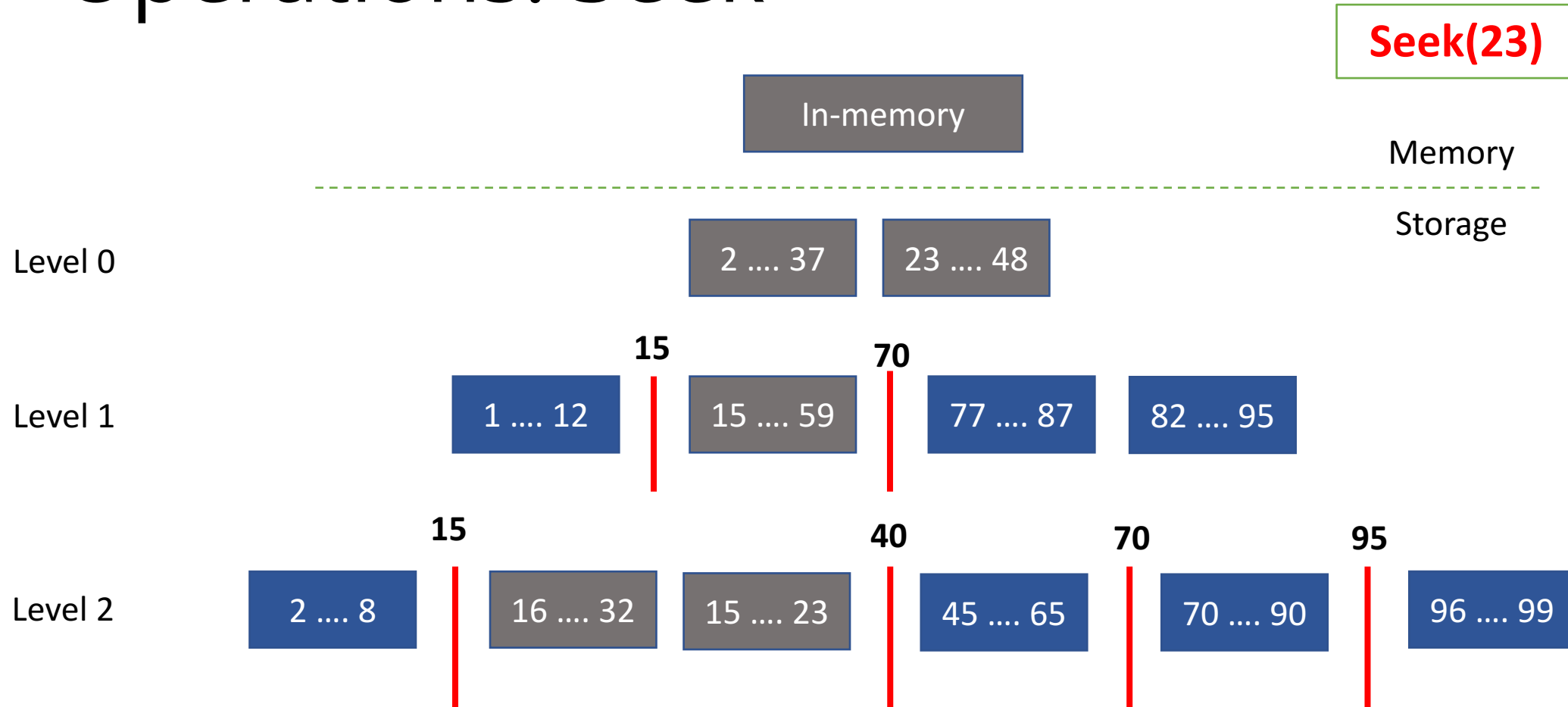
Level 2  –  5, 300, 500

Seek(200)

# Operations: Seek

- **Seek(target):** Returns the smallest key in the database which is >= target
- Used for range queries (for example, return all entries between 5 and 18)

# Operations: Seek



Seek(23)

In-memory

Memory

Storage

Level 0

2 .... 37    23 .... 48

**15**    **70**

Level 1

1 .... 12    15 .... 59    77 .... 87    82 .... 95

**15**    **40**    **70**    **95**

Level 2

2 .... 8    16 .... 32    15 .... 23    45 .... 65    70 .... 90    96 .... 99

FLSM structure

# Operations: Seek

Seek(23)

In-memory

Memory

Storage

Level 0

| 2 …. 37 | 23 …. 48 |

**15**          **70**

Level 1

| 1 …. 12 | 15 …. 59 | 77 …. 87 | 82 …. 95 |

**15**          **40**          **70**          **95**

Level 2
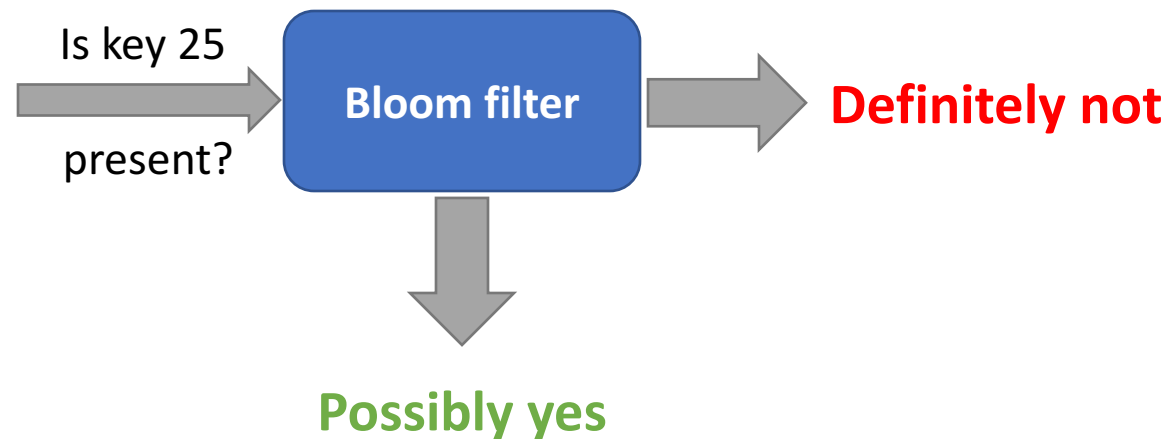
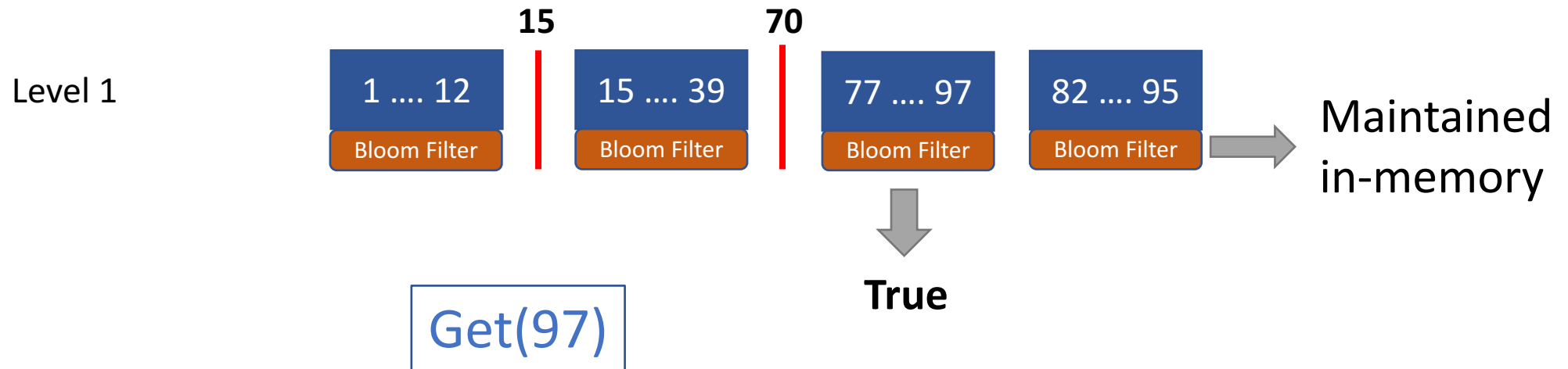| 2 …. 8 | 16 …. 32 | 15 …. 23 | 45 …. 65 | 70 …. 90 | 96 …. 99 |

All levels and memtable need to be searched

94

# Optimizations in PebblesDB

- **Challenge with reads:** Multiple sstable reads per level
- Optimized using sstable level bloom filters
- Bloom filter: determine if an element is in a set

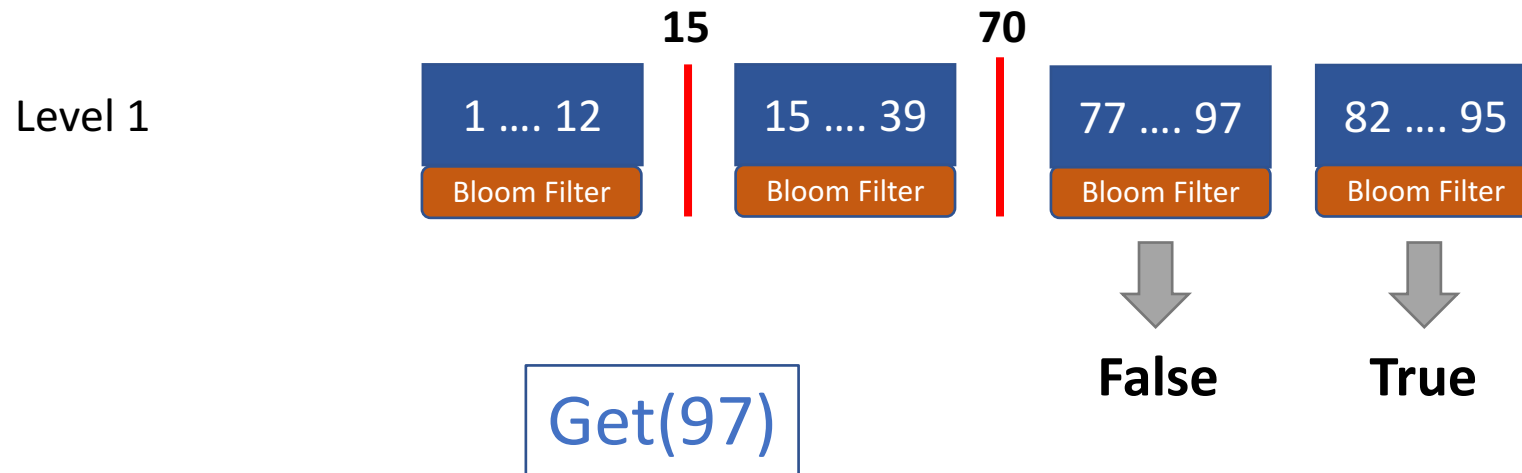Is key 25 present? → **Bloom filter** → **Definitely not**

**Possibly yes**

# Optimizations in PebblesDB

- **Challenge with reads:** Multiple sstable reads per level
- Optimized using sstable level bloom filters
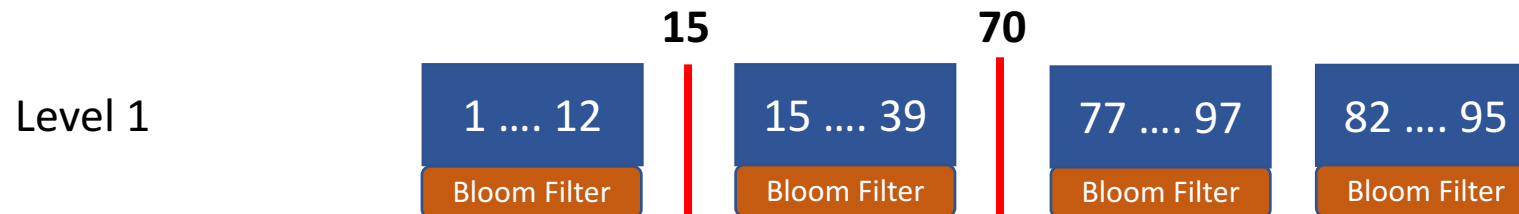- Bloom filter: determine if an element is in a set

# Optimizations in PebblesDB

- **Challenge with reads:** Multiple sstable reads per level
- Optimized using <span style="color:red">sstable level bloom filters</span>
- Bloom filter: determine if an element is in a set

**15**    **70**

Level 1

| 1 …. 12 | | 15 …. 39 | | 77 …. 97 | 82 …. 95 |
|---------|---|----------|---|----------|----------|
| Bloom Filter | | Bloom Filter | | Bloom Filter | Bloom Filter |

**False**       **True**

Get(97)
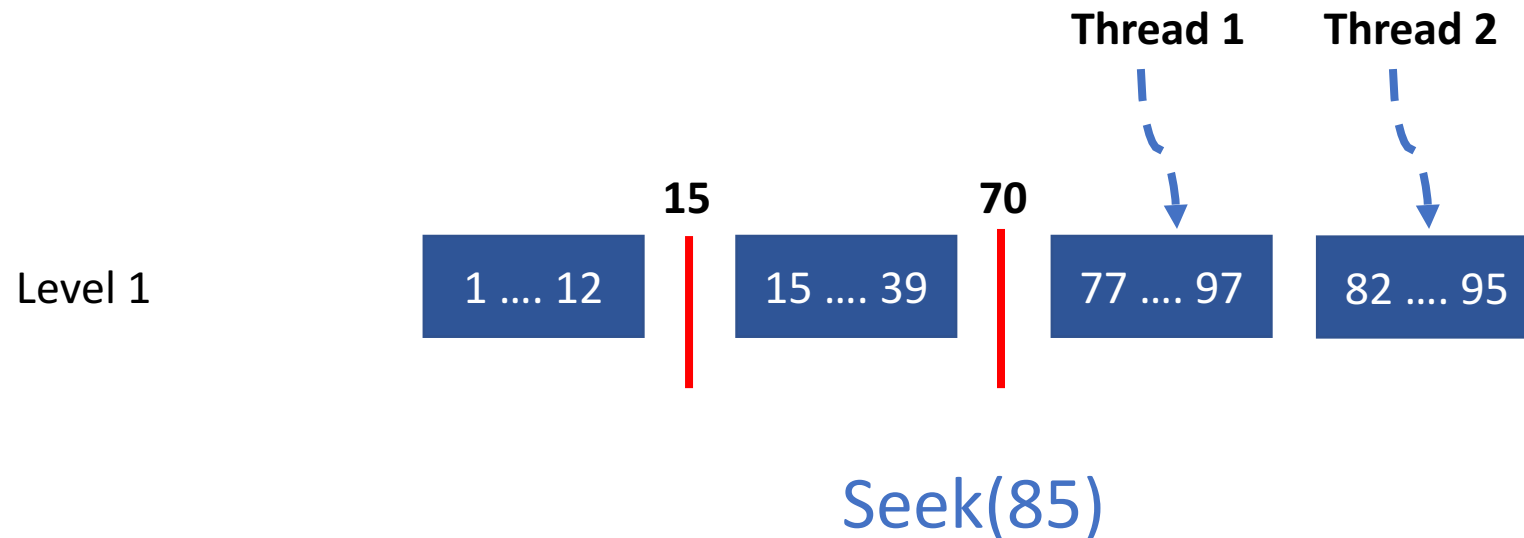
# Optimizations in PebblesDB

- **Challenge with reads:** Multiple sstable reads per level
- Optimized using <span style="color:red">sstable level bloom filters</span>
- Bloom filter: determine if an element is in a set

**15**        **70**

Level 1

| 1 …. 12 |
| --- |
| Bloom Filter |

| 15 …. 39 |
| --- |
| Bloom Filter |

| 77 …. 97 |
| --- |
| Bloom Filter |

| 82 …. 95 |
| --- |
| Bloom Filter |

PebblesDB reads <span style="color:red">at most one file</span> per guard with high probability

# Optimizations in PebblesDB

- **Challenge with seeks:** Multiple sstable reads per level
- **Parallel seeks:** Parallel threads to seek() on files in a guard

**Thread 1**    **Thread 2**

**15**    **70**

Level 1    | 1 .... 12 | | 15 .... 39 | | 77 .... 97 | | 82 .... 95 |

Seek(85)

# Optimizations in PebblesDB

- **Challenge with seeks:** Multiple sstable reads per level
- **Parallel seeks:** Parallel threads to seek() on files in a guard
- **Seek based compaction:** Triggers compaction for a level during a seek-heavy workload
  - Reduce the average number of sstables per guard
  - Reduce the number of active levels

Seek based compaction increases write I/O but as a trade-off to improve seek performance
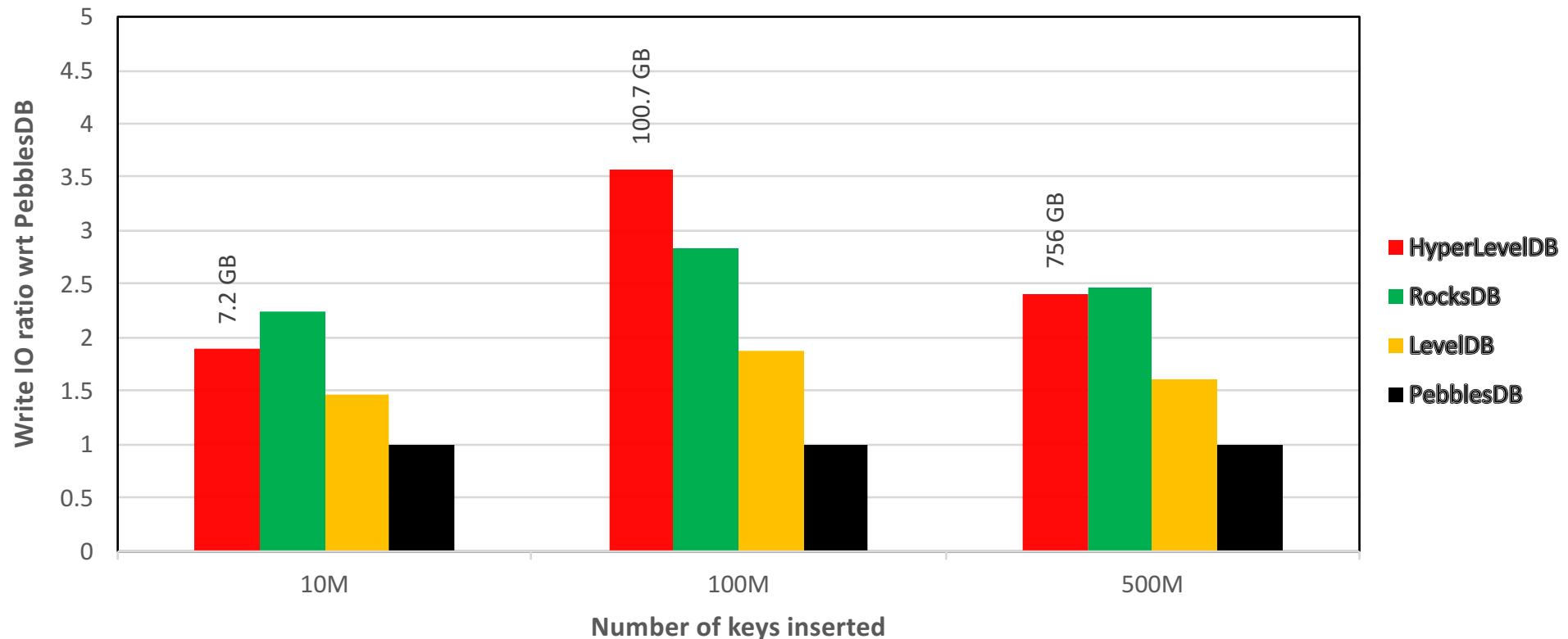
# Tuning PebblesDB

- PebblesDB characteristics like
  - Increase in write throughput,
  - decrease in write amplification and
  - overhead of read/seek operation

  all depend on one parameter, **maxFilesPerGuard** (default 2 in PebblesDB)
- Setting this to a very high value favors write throughput
- Setting this to a very low value favors read throughput

# Experimental setup

- Intel Xeon 2.8 GHz processor
- 16 GB RAM
- Running Ubuntu 16.04 LTS with the Linux 4.4 kernel
- Software RAID0 over 2 Intel 750 SSDs (1.2 TB each)
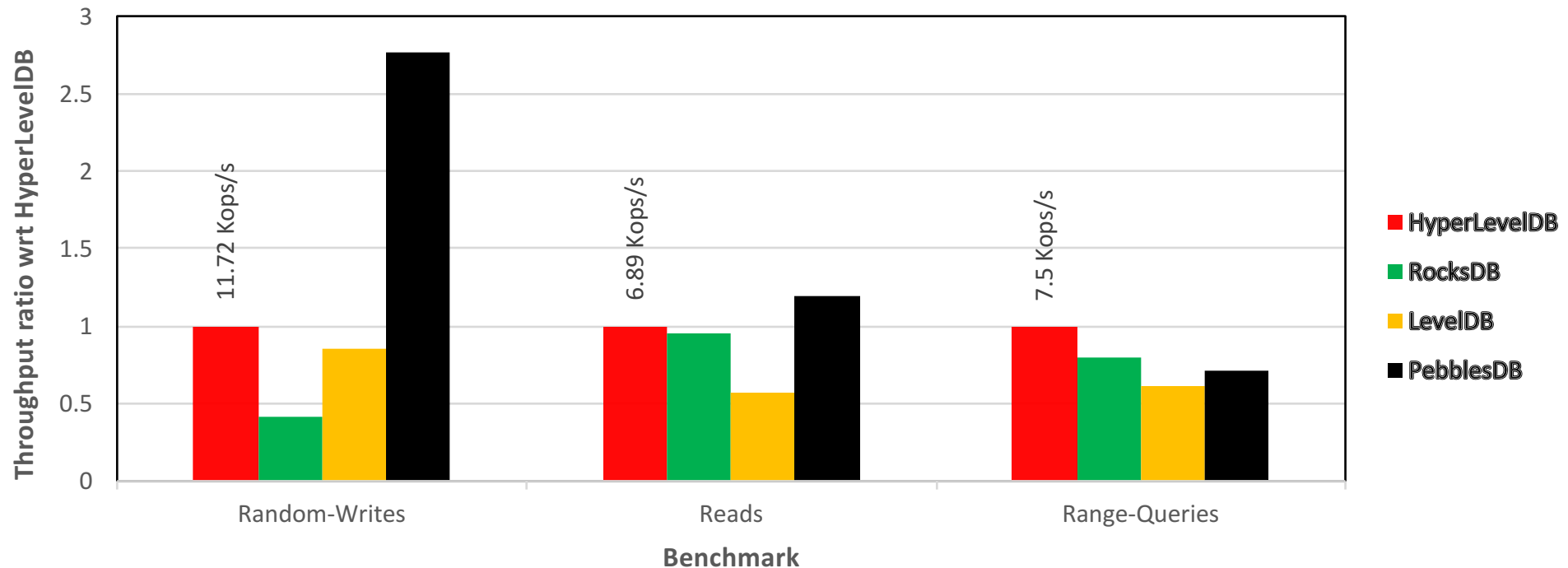- Datasets in experiments 3x bigger than DRAM size

# Write amplification

- Inserted different number of keys with key size 16 bytes and value size 128 bytes
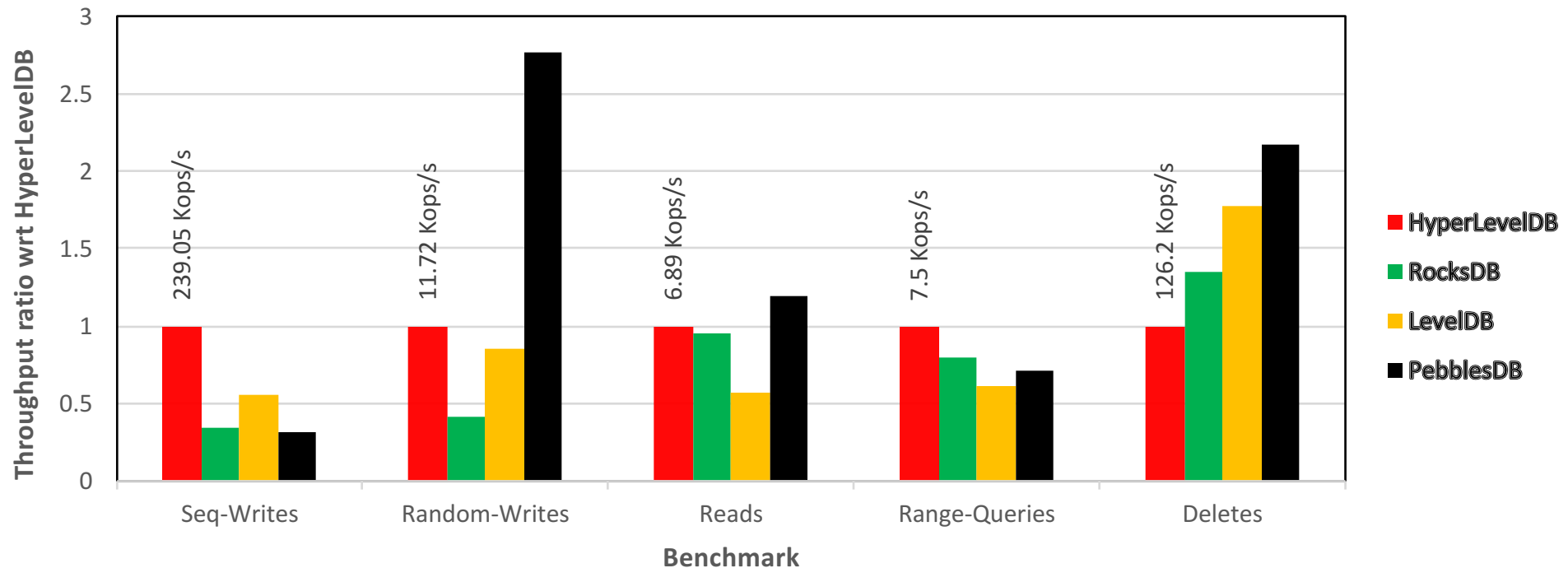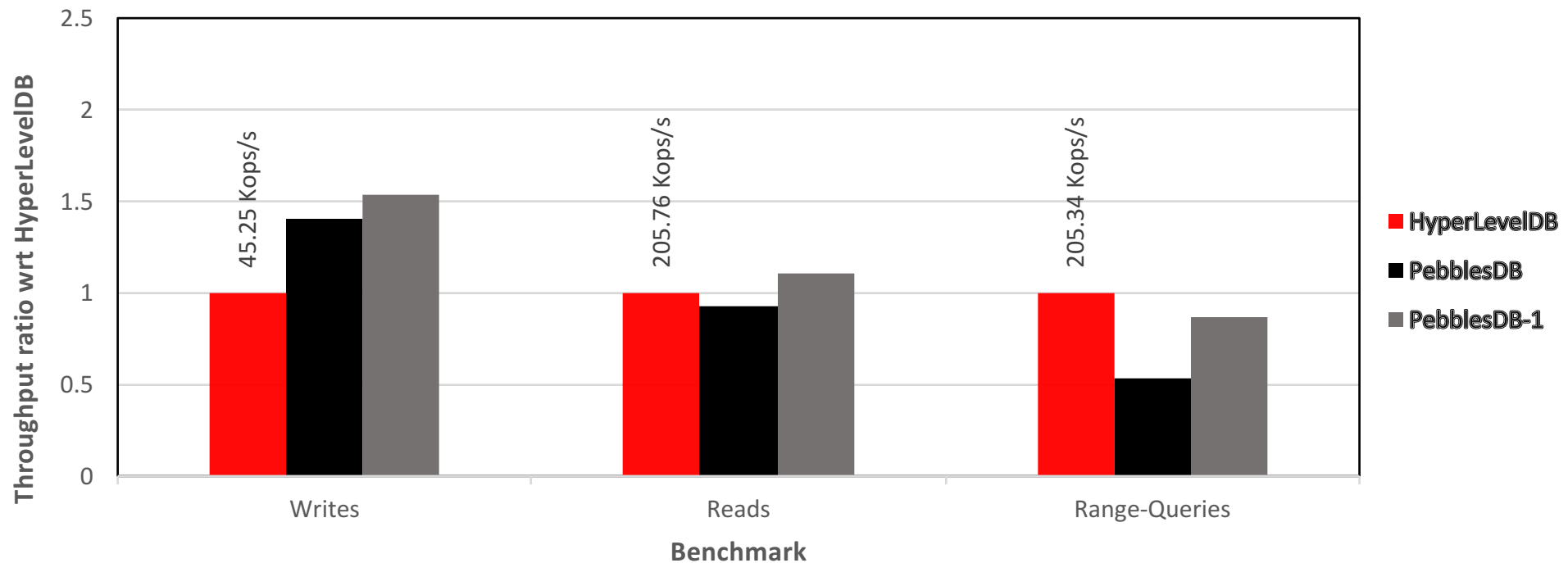
# Micro-benchmarks

- Used **db_bench** tool that ships with LevelDB
- Inserted 50M key-value pairs with key size 16 bytes and value size 1 KB
- Number of read/seek operations: 10M

# Micro-benchmarks

- Used **db_bench** tool that ships with LevelDB
- Inserted 50M key-value pairs with key size 16 bytes and value size 1 KB
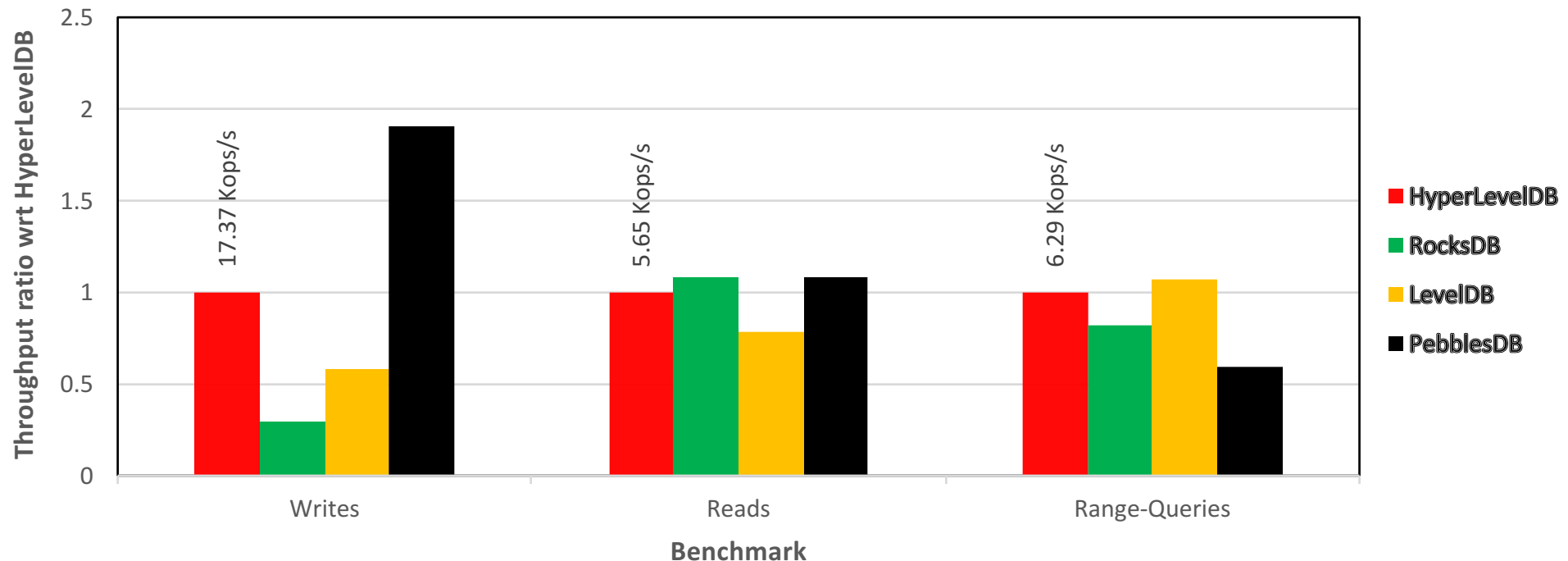- Number of read/seek operations: 10M

# Small cached dataset

- Insert 1M key-value pairs with 16 bytes key and 1 KB value
- Total data set (~1 GB) fits within memory
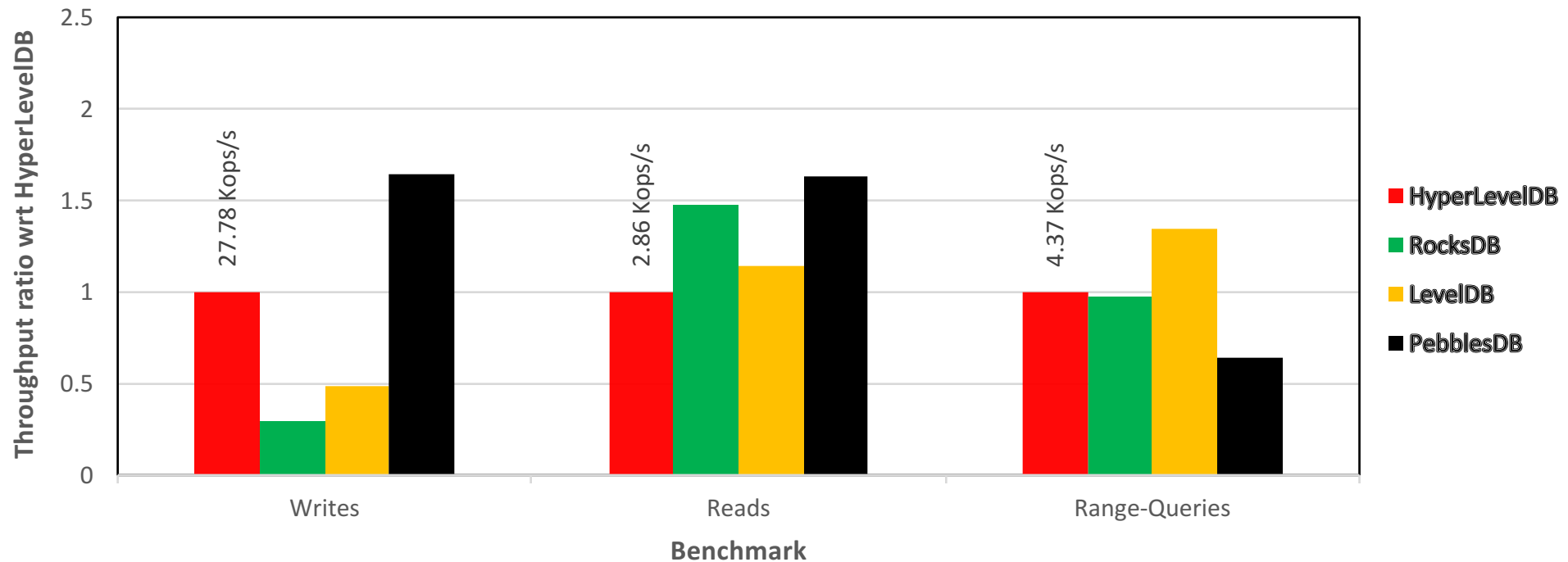- PebblesDB-1: with maximum one file per guard

# Aged FS and KV store

- File system aging: Fill up 89% of the file system
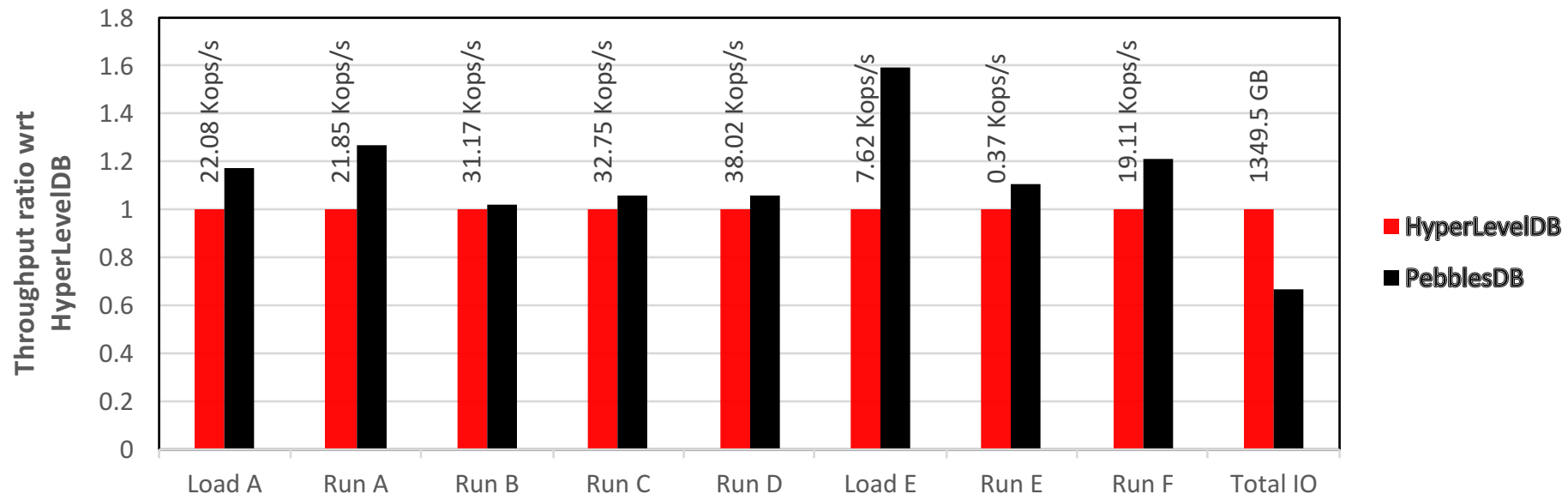- KV store aging: Insert 50M, delete 20M and update 20M key-value pairs in random order

# Low memory micro-benchmark

- 100M key-value pairs with 1KB (~65 GB data set)
- DRAM was limited to 4 GB

# NoSQL stores - HyperDex

- HyperDex – distributed key-value store from Cornell
- Inserted 20M key-value pairs with 1 KB value size and 10M operations



Load A - 100 % writes          Run D   - 95% reads (latest), 5% writes

Run A   - 50% reads, 50% writes          Load E  - 100% writes

Run B   - 95% reads, 5% writes          Run E   - 95% range queries, 5% writes

Run C   - 100% reads          Run F   - 50% reads, 50% read-modify-writes

# CPU usage

- Median CPU usage by inserting 30M keys and reading 10M keys
- PebblesDB: ~171%
- Other key-value stores: 98-110%
- Due to aggressive compaction, more CPU operations due to merging multiple files in a guard

# Memory usage

- 100M records (16 bytes key, 1 KB value) – 106 GB data set
  - 300 MB memory space
  - 0.3% of data set size

# Bloom filter calculation cost

- 1.2 sec per GB of sstable
- 3200 files – 52 GB – 62 seconds

# Impact of different optimizations

- Sstable level bloom filter improve read performance by 63%
- PebblesDB without optimizations for seek – 66%

# Thank you!

## Questions?