# RocksDB
## Key-Value Store Optimized For Flash

Siying Dong

Software Engineer, Database Engineering Team @ Facebook
April 20, 2016

# Agenda

RocksDB

# What is RocksDB?

- Key-Value persistent store
- Point / range lookup
- Optimized for flash
- Also work for pure-memory and spinning disks
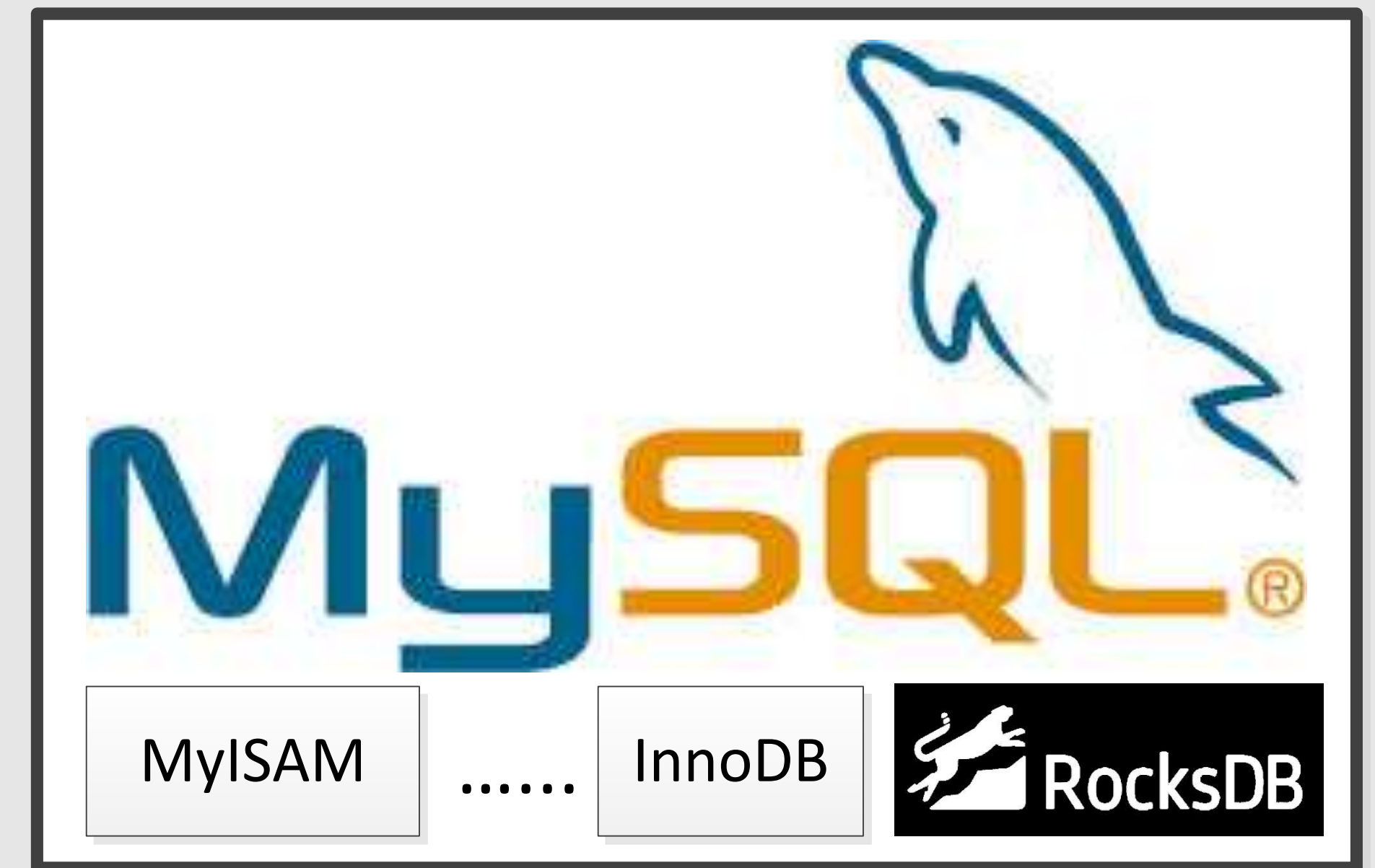- C++ library
- Other language bindings
- Fork of LevelDB

# RocksDB Interface

- Keys and Values are byte arrays
- Keys have total order.
- Update Operation: Put/Delete/Merge
- Queries: Get/Iterator
- Consistency: atomic multi-put, multi-get, iterator, snapshot read, transactions

# How Is RocksDB Used?

- As Storage Engine of Databases
- As Embedded Storage of Applications

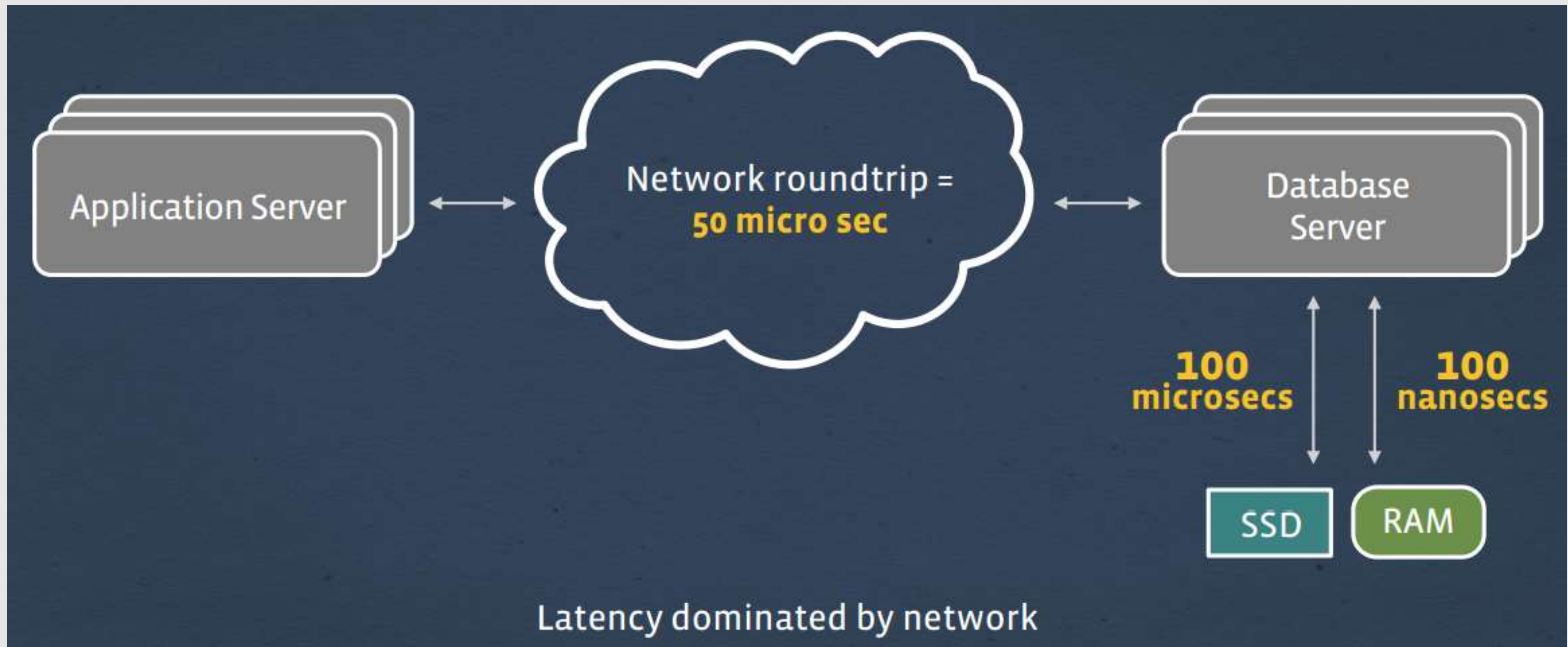# RocksDB As Storage Engine of Data Management Systems

mongoDB
RocksDB | Mmap | WiredTiger

MySQL
MyISAM ...... InnoDB RocksDB

Yahoo Sherpa
RocksDB

ZippyDB

*And many more ...*

# RocksDB As Embedded Storage

- Facebook: many backend services
- LinkedIn's FollowFeed
- Apache Samza
- Iron.io
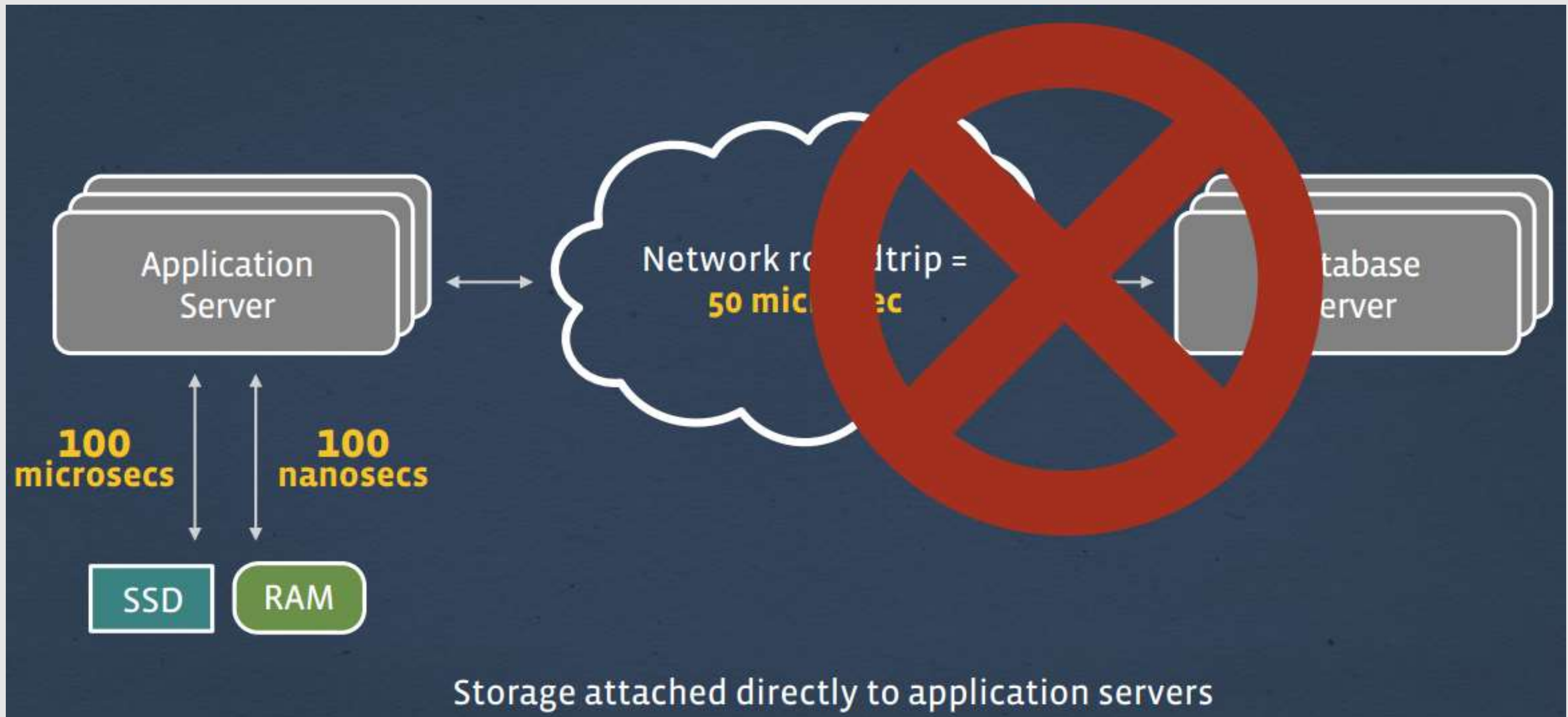- Tango Me
- Ceph
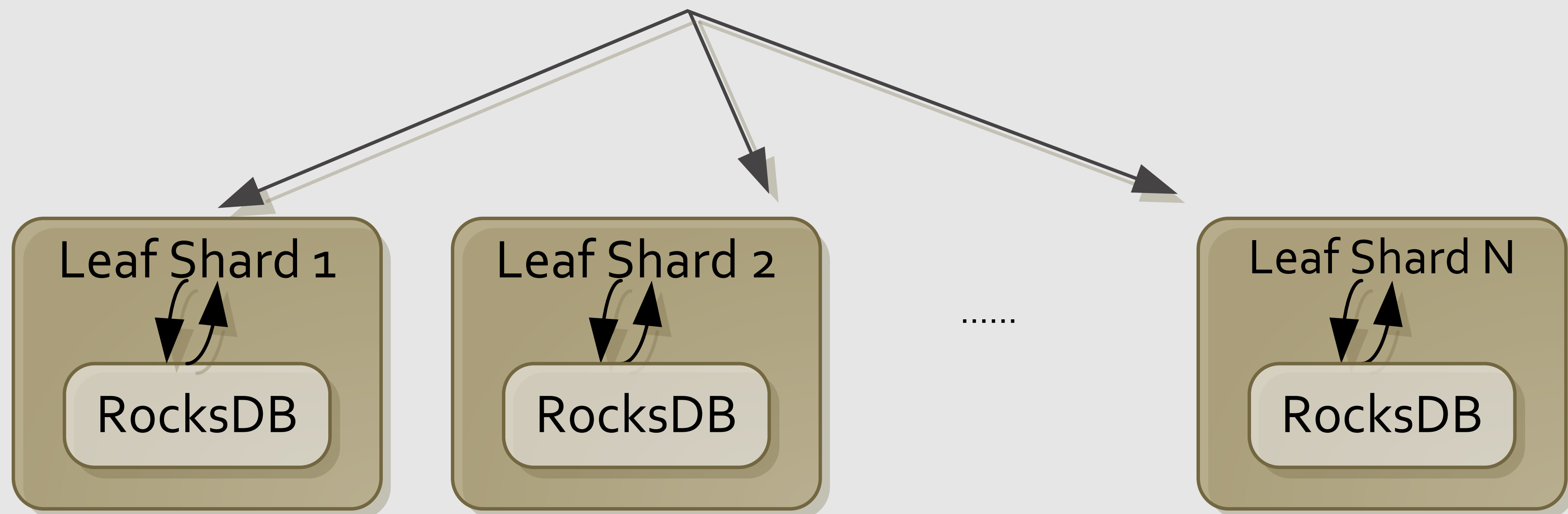- And more...

# Why Embedded Storage in Application?



Application Server ↔ Network roundtrip = **50 micro sec** ↔ Database Server

100 **microsecs** — 100 **nanosecs**

SSD — RAM

Latency dominated by network

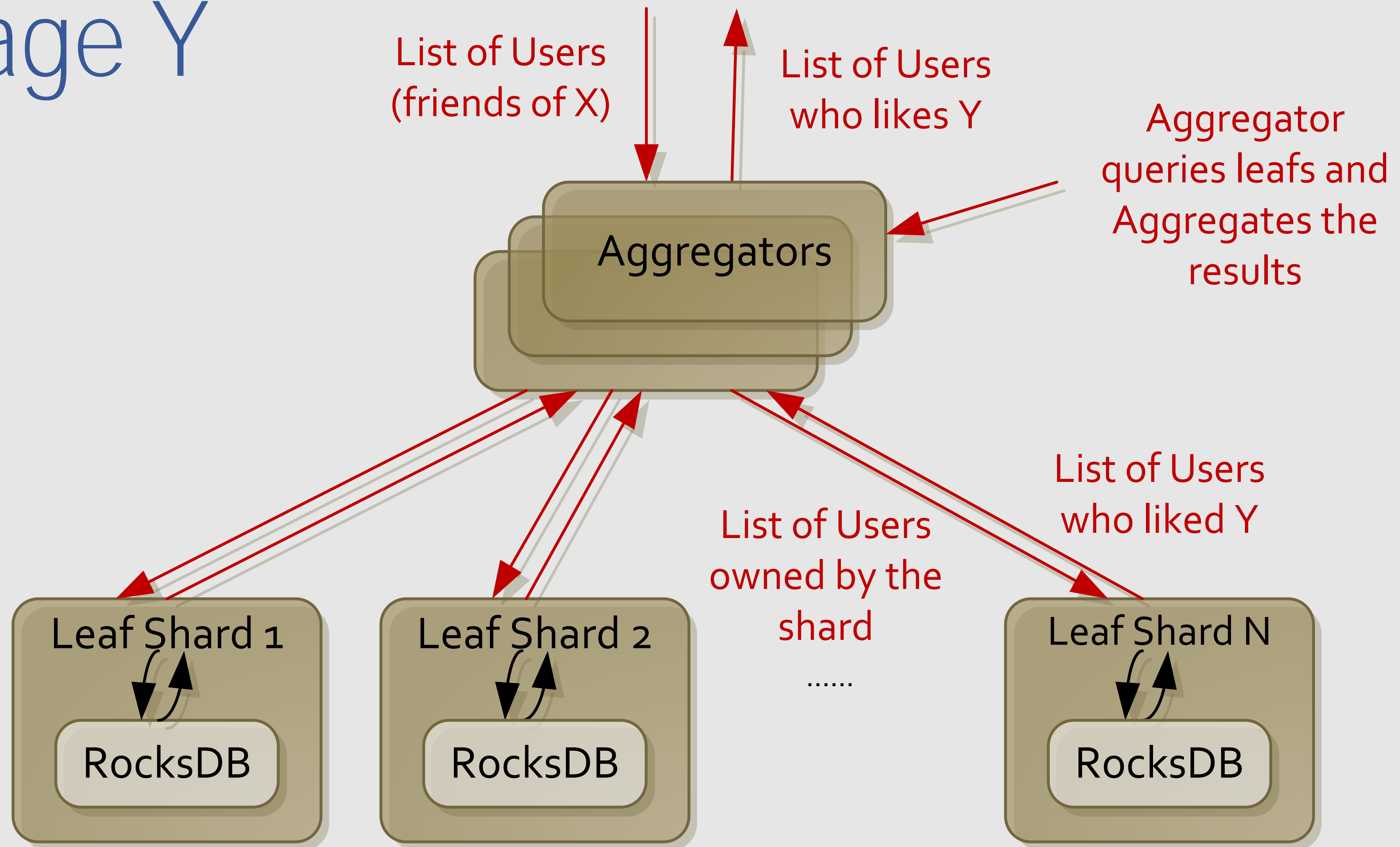# Why Embedded Storage in Application?



Storage attached directly to application servers

# Example: Find all friends of user X who like Page Y

Maintains liker -> page
mapping in RocksDB.
Sharded by likers.

| Leaf Shard 1 | Leaf Shard 2 | ...... | Leaf Shard N |
|:---:|:---:|:---:|:---:|
| RocksDB | RocksDB | | RocksDB |

# Example: Find all friends of user X who like Page Y

List of Users (friends of X)

List of Users who likes Y

Aggregator queries leafs and Aggregates the results

Aggregators

List of Users owned by the shard

......

List of Users who liked Y

Leaf Shard 1

RocksDB

Leaf Shard 2

RocksDB
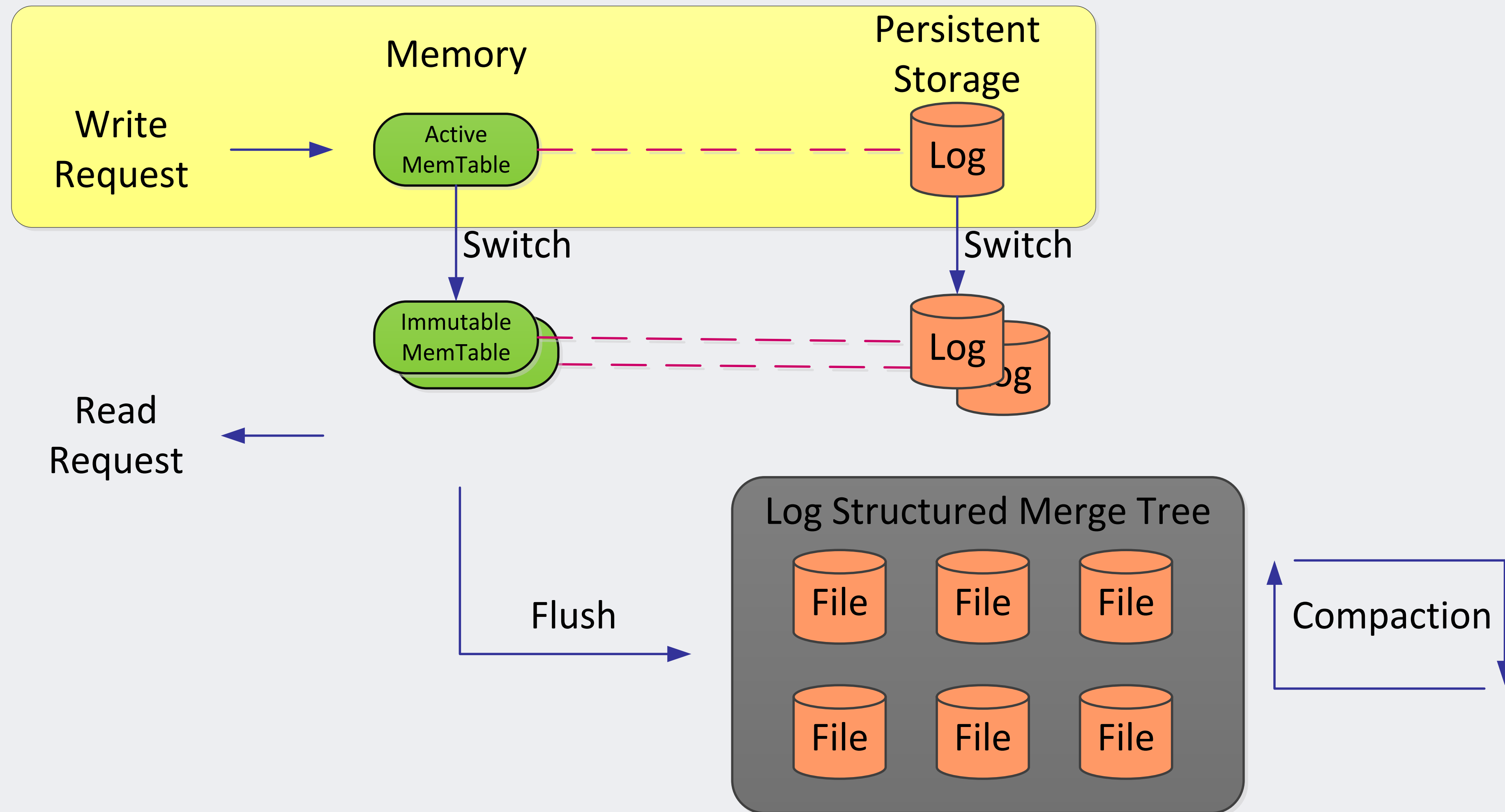
Leaf Shard N

RocksDB

# RocksDB Design

# RocksDB Design

# RocksDB Architecture
## Log-Structured Merge-Tree
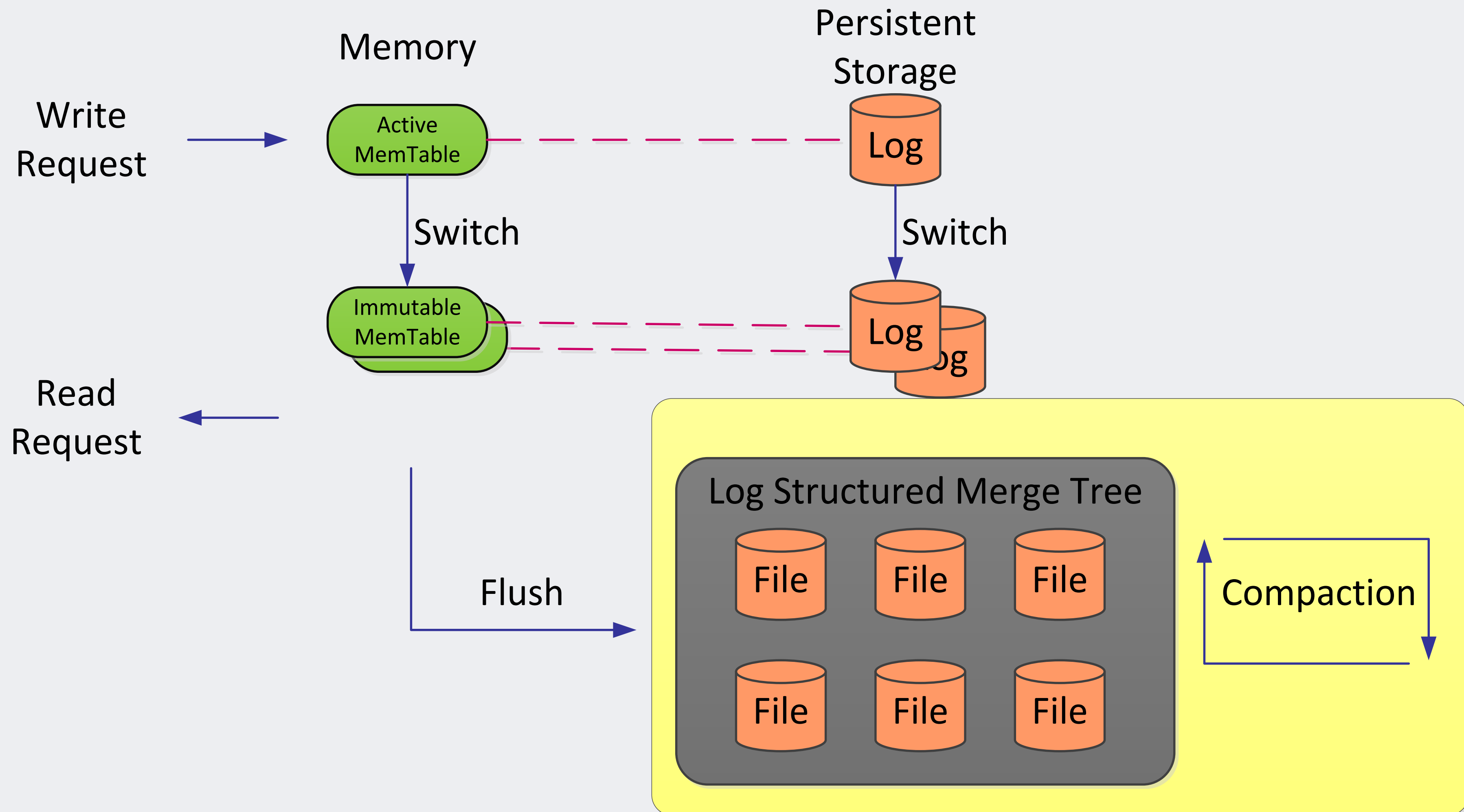
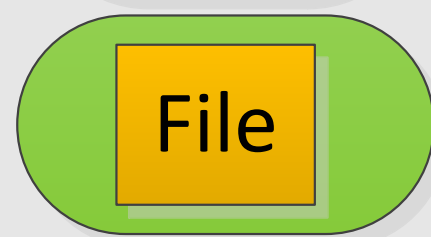# Write Path (1)

# Write Path (2)

# Write Path (3)

Memory

Persistent Storage

Write Request → Active MemTable

Active MemTable ---- Log

Switch

Immutable MemTable

Switch

Log

Log

Read Request ←

Flush →

**Log Structured Merge Tree**

| File | File | File |
| File | File | File |

Compaction

# Write Path (4)

Memory

Persistent Storage

Write Request → Active MemTable ----- Log

↓ Switch

↓ Switch

Immutable MemTable ----- Log Log

Read Request ←

Flush →

## Log Structured Merge Tree

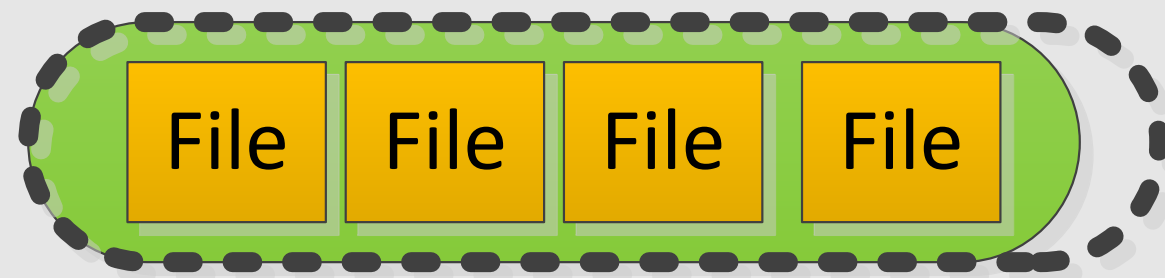File   File   File
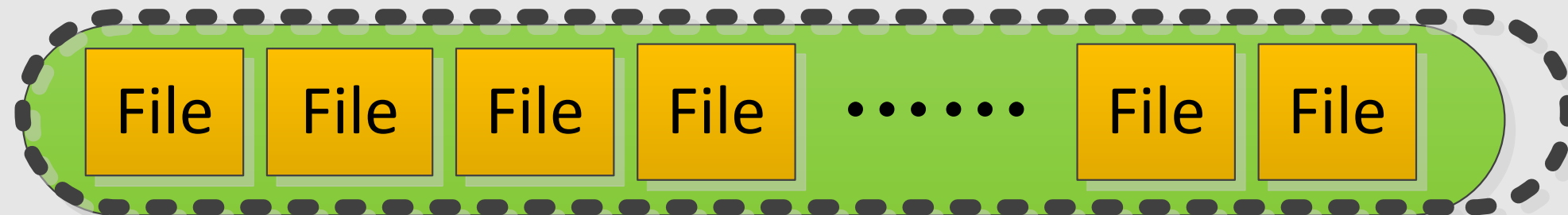
File   File   File

Compaction
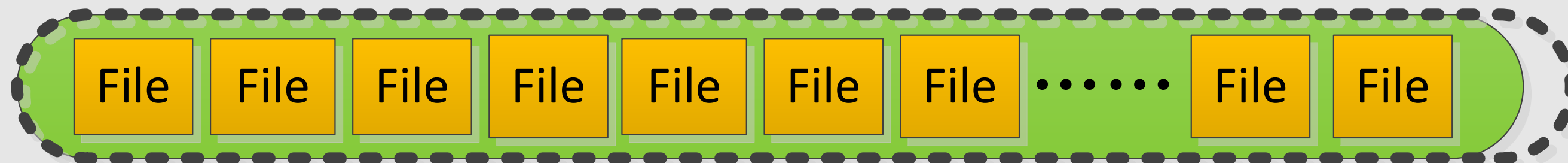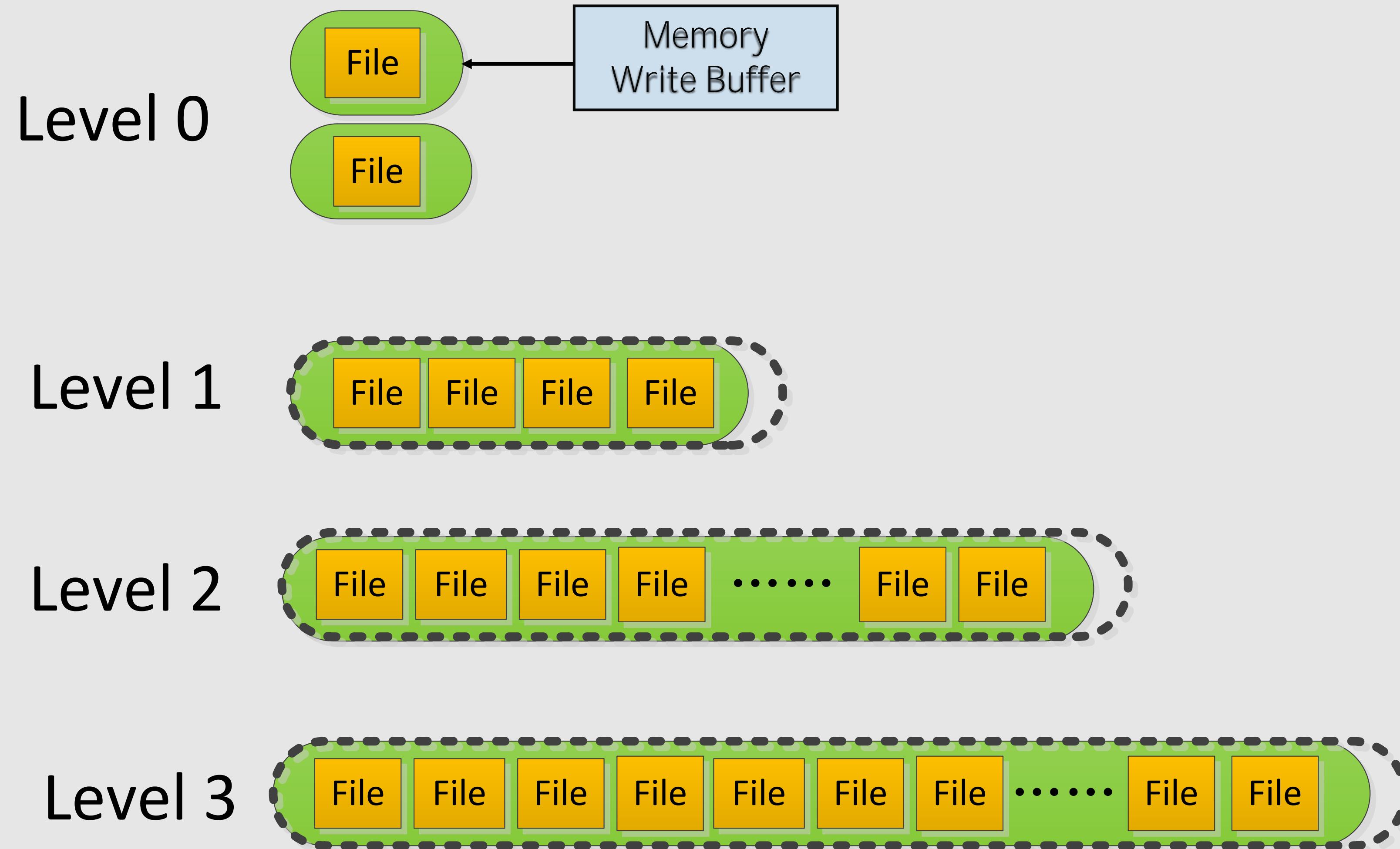
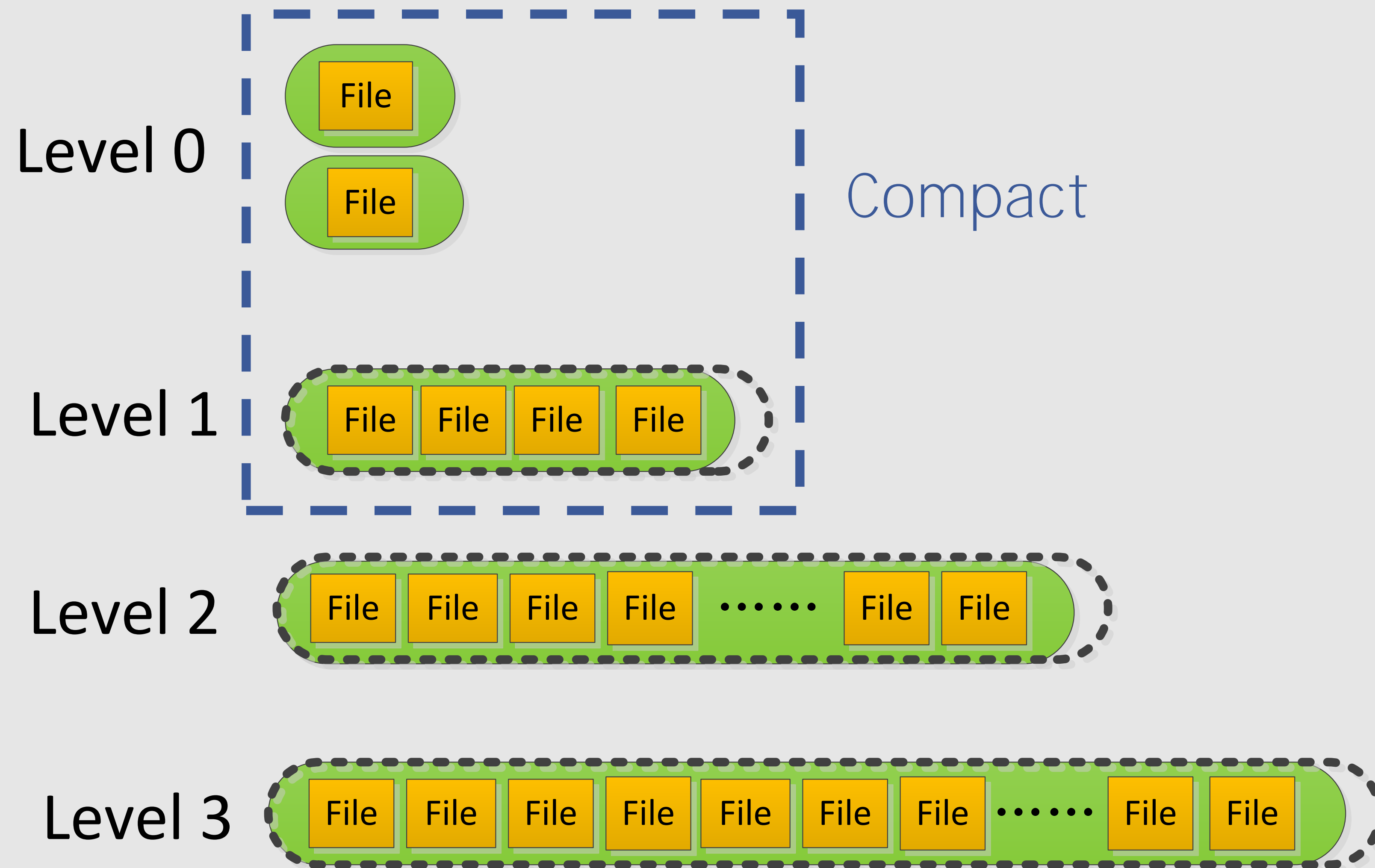# Level-Based Compaction

Level 0



Level 1



Level 2



Level 3

# Level-Based Compaction

# Level-Based Compaction

**Level 0**
- File
- File

Compact

**Level 1**
- File | File | File | File

**Level 2**
- File | File | File | File | ······ | File | File

**Level 3**
- File | File | File | File | File | File | File | ······ | File | File

# Level-Based Compaction

Level 1

| New File | New File | New File | New File | New File |

Level 2

| File | File | File | File | ⋯⋯ | File | File |

Level 3

| File | File | File | File | File | File | File | ⋯⋯ | File | File |

# Level-Based Compaction

Compact

Level 1

| File | File | File | File | File |

Level 2

| File | File | File | File | ...... | File | File |

Level 3

| File | File | File | File | File | File | File | ...... | File | File |

# Level-Based Compaction

**Level 1**

| File | File | | File |

**Level 2**

| File | New File | New File | New File | $\cdots\cdots$ | File | File | File |

**Level 3**

| File | File | File | File | File | File | File | $\cdots\cdots$ | File | File |

# Level-Based Compaction

Level 1

| File | File | File | File |

Compact

Level 2

| File | File | File | File | ...... | File | File | File |

Level 3

| File | File | File | File | File | File | File | ...... | File | File |

# Level-Based Compaction

Level 0

Level 1

| File | File | File | File |

Level 2

| File | | File | File | ...... | File |

Level 3

| New File | New File | New File | File | File | New File | New File | ...... | File | File |

# Example of Level Base Targets

Level 0

File

Memory
Write Buffer

File

Level 1    File File File File    Target: 1 GB

Level 2    File File File File ······ File File    Target: 10 GB

Level 3    File File File File File File File ······ File File    Target: 100GB

# Why is it flash-friendly?
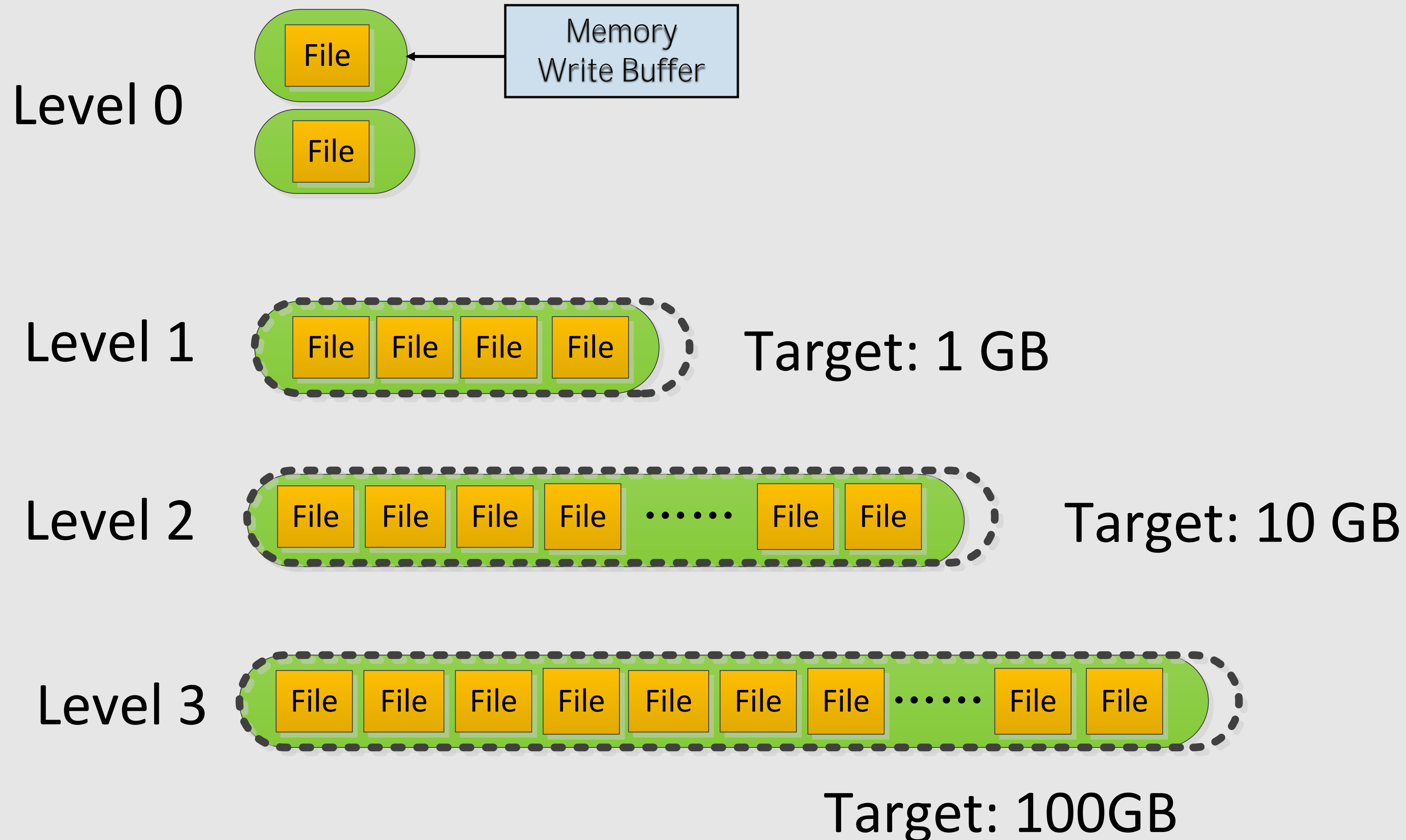
# Tuning Flexibility for Flash

Performance Metrics for applications on flash devices

- Write Amplification – wear out devices slower

- Space Amplification – store more data

- Read Amplification – better read IOPs

# Compactions' Impact on Amplifications

| | Space Amplification | Write Amplification | Memory Cache Required for ReadAmp = 1 |
|---|---|---|---|
| More Aggressive Compactions | 🙂 | 🙁 | 🙂 |
| Less Aggressive Compactions | 🙁 | 🙂 | 🙁 |

# Space Amplification is the bottleneck

- Example: our MySQL host on InnoDB:
  - *Read IOPS: < 10%*
  - *Write IOPS: < 35%*
  - *Peak Write Bandwidth: < 25%*
  - *CPU: < 40%*
  - *Write Endurance: last more than 3 years.*

Everything except space has room to go!

# Space Amplification of RocksDB

## Only 10% Extra Space

### How?

# Space efficiency in LSM?

Write Buffer in Memory

Persistent Store

File

File

File ......

File   File   ...... File

File   File   File   ...... File

10% Extra Space

Size Similar to User Data Size

# How Did We Guarantee 10%?

## A Space-Efficient Approach

Level 0    File

Level 1    File

⋮

Level N-2    File ......    Target: 8.76 GB < 1%

Level N-1    File File ...... File    Target: 87.6 GB    9%

Level N    File File File ...... File    876 GB    90% of total size

# Lower Write Amplification

## InnoDB

Row

Row

Row

Modify

Row

Row

Row

Row

Row

Row

Row

Row

Row

Read

Write

Write Amp = Page size / row size

## RocksDB

*Write amp 1*

flush Level 0

Merge Level 1    Target 1GB

*Write Amp 10*

Merge Level 2    Target 10 GB

*Write Amp 10*

Merge Level 3    Target 100 GB

*Write Amp 10*

Merge Level 4

Target 1000 GB

# How About Other Metrics?

- Read QPS
- Write Throughput

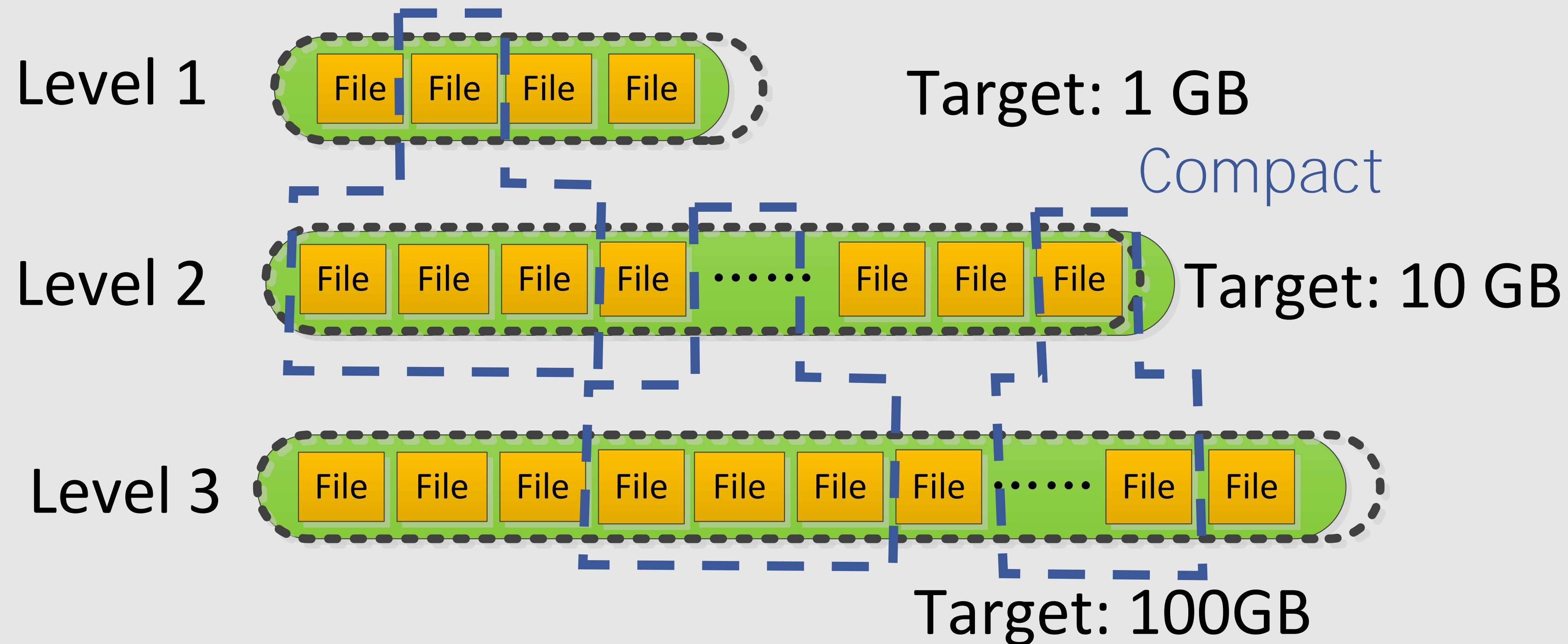# Make Read Throughput High: Reduced Locking in Reads

- Memtable: skip list

- Data Files: immutable

- LSM tree change: thread-local cache of the tree

- Synchronize opened files: allow to keep all files open

- Block cache mutex: sharded; more optimization coming.

# Write Throughput

- Throughput of Compactions
- Throughput of Memtable Inserts
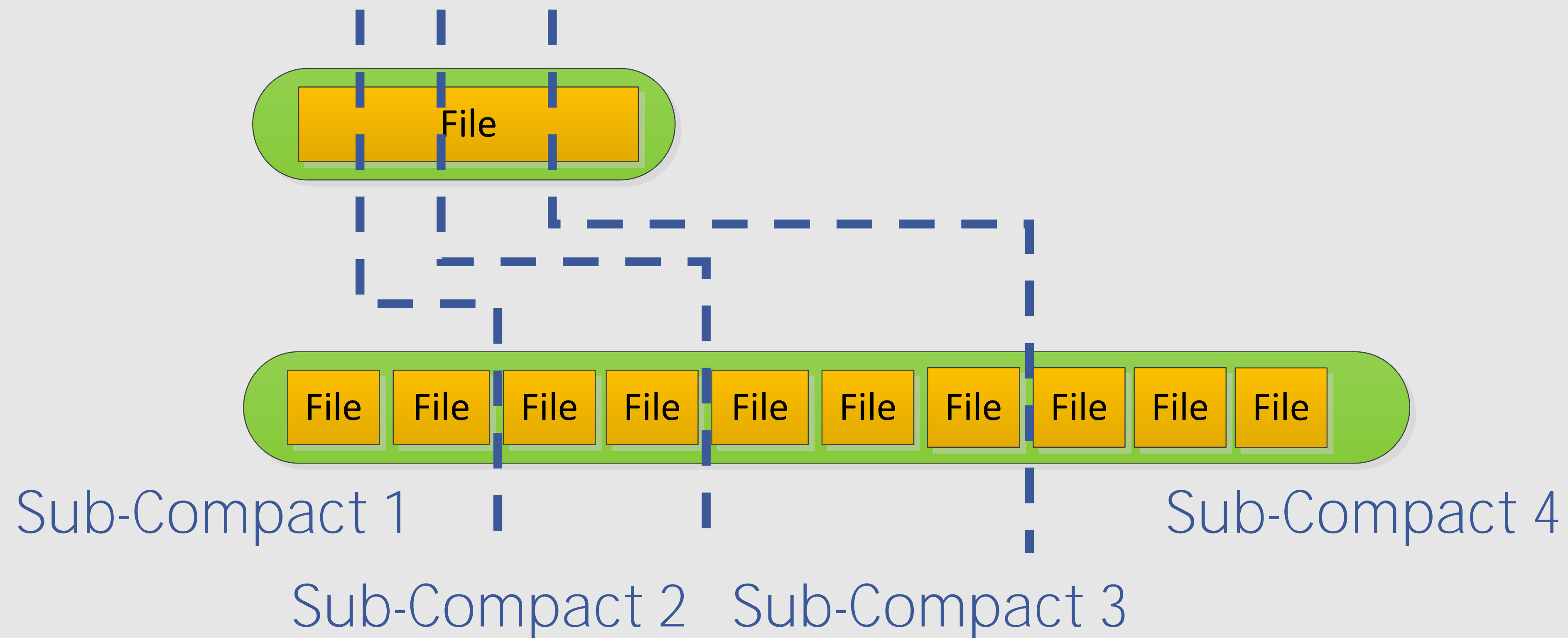
# Multi-thread compactions

Compact non-overlapping files

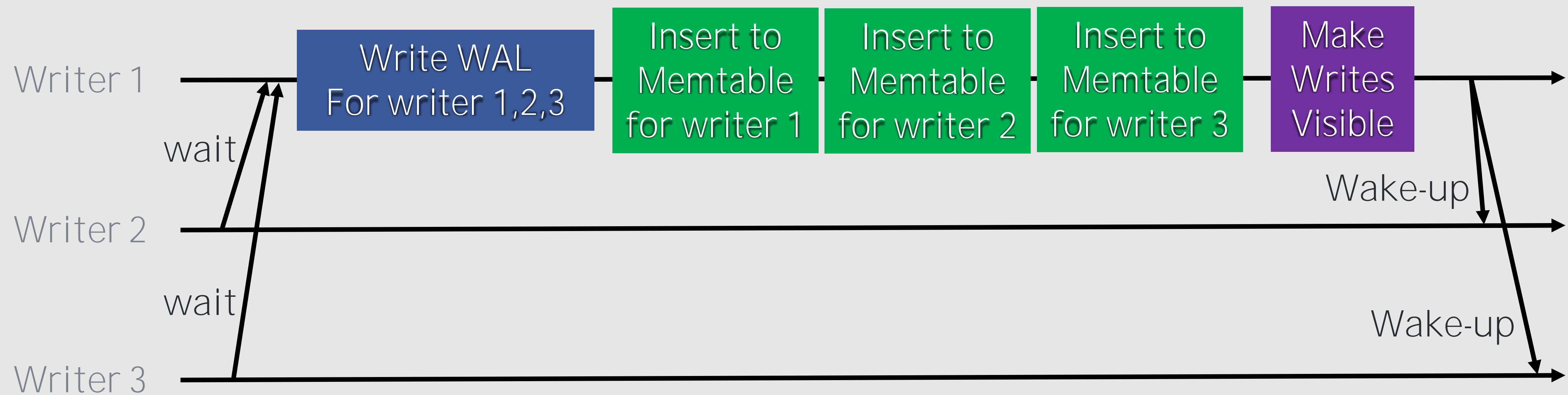Level 1    | File | File | File | File |    Target: 1 GB

Compact

Level 2    | File | File | File | File | ······ | File | File | File |    Target: 10 GB

Level 3    | File | File | File | File | File | File | File | ······ | File | File |    Target: 100GB

# Multi-thread compactions

Divide One Compaction to sub-compactions



Sub-Compact 1

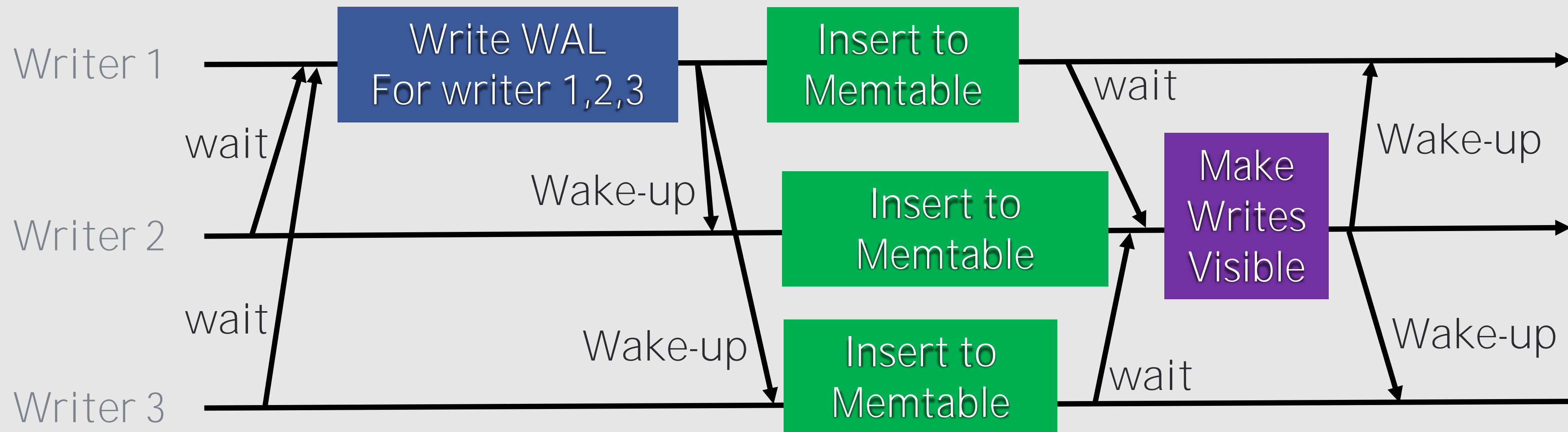Sub-Compact 2   Sub-Compact 3

Sub-Compact 4

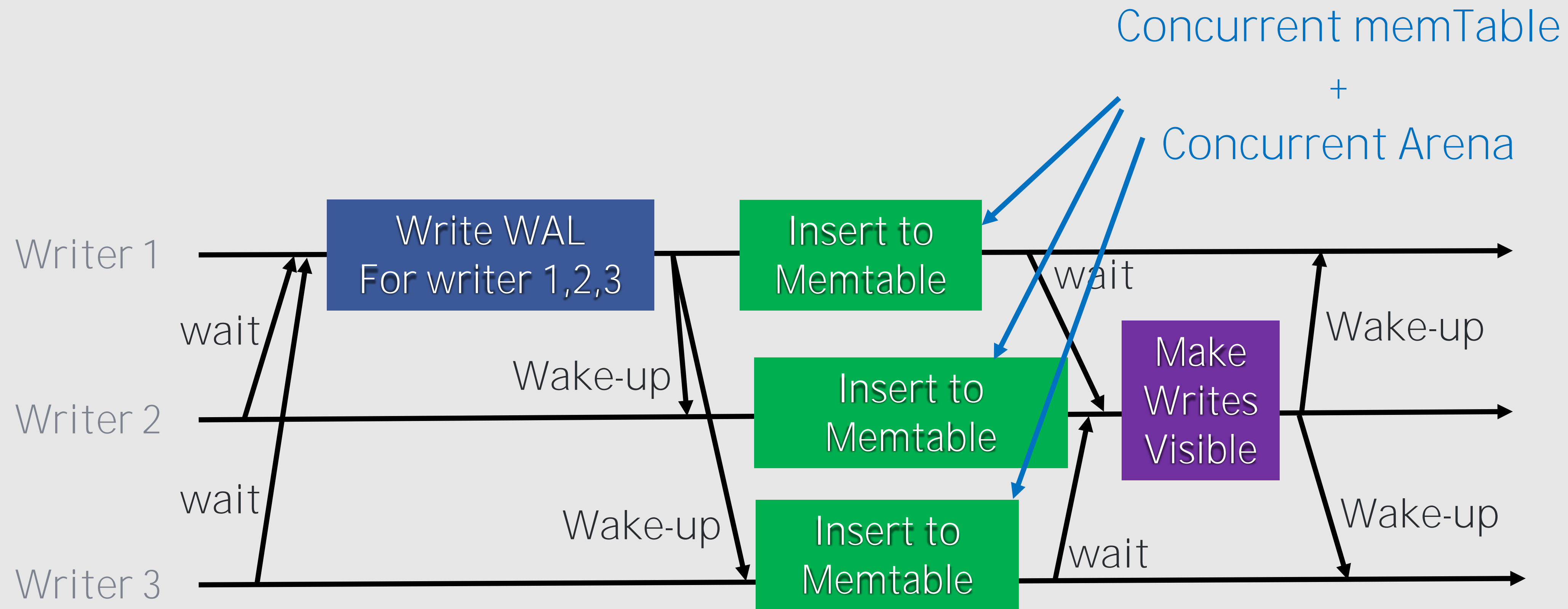# Parallel Memtable Insert

## Previous Workflow of Group Commit

# Parallel Memtable Insert

## Parallel Memtable Insert for Group Commit

# Parallel Memtable Insert

## Parallel Memtable Insert for Group Commit

Concurrent memTable
+
Concurrent Arena

Writer 1 —— Write WAL For writer 1,2,3 —— Insert to Memtable —— wait ——

wait

Wake-up

Make Writes Visible

Wake-up

Writer 2 —— Wake-up —— Insert to Memtable ——

wait

Writer 3 —— Wake-up —— Insert to Memtable —— wait —— Wake-up

# RocksDB Performance On Flash

- Space, Read And Write Amplificaton Trade-offs
- Low Space Amplification
- High Read QPS: Reduced Mutex Locking
- High Write Throughput: Parallel Compaction; Concurrent Memtable Insert

# Other Storage Media?

# RocksDB On Other Storage Media

- Memory-Only:
  - *Memory Efficiency*
  - *7 million reads/s in single host benchmark*
- Spinning Disk:
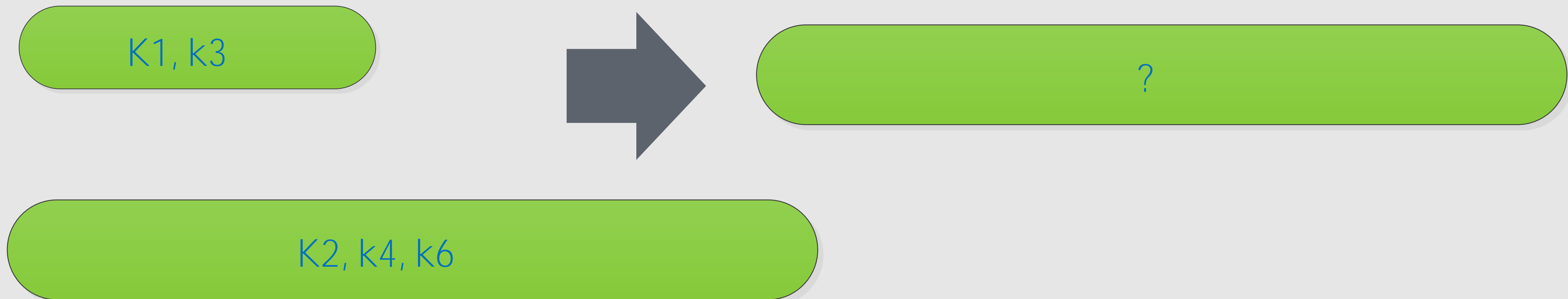  - *Write-Optimized*
  - *Reasonable Read Performance*

# Other Benefits?

# Features Enabled by LSM
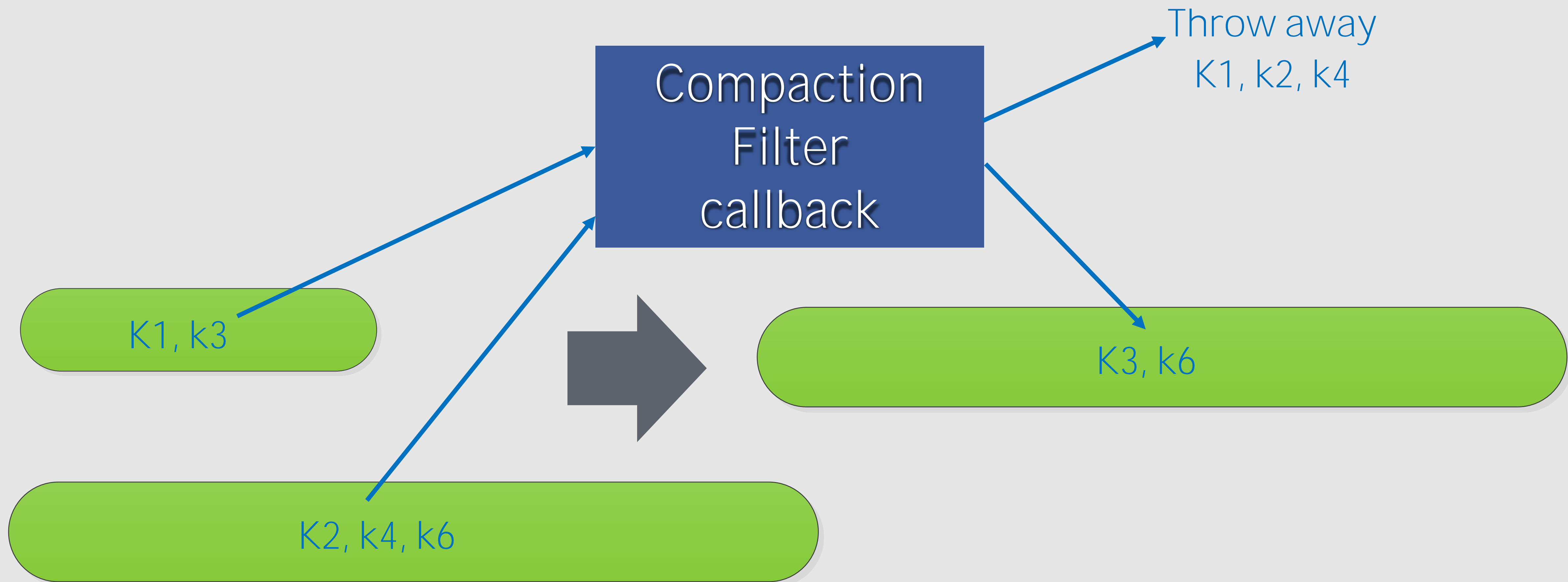
- Compaction Filter
- Merge Operator

# Compaction Filter

## Expire Stale Keys In Compactions

K1, k3 → ?

K2, k4, k6

# Compaction Filter

## Expire Stale Keys In Compactions

Compaction
Filter
callback

Throw away
K1, k2, k4

K1, k3

K2, k4, k6

K3, k6

# Merge Operator

## Avoid Read-Modify-Write

memtable

Merge k1 = +2

Merge k1 = -1

Put k1 = 5

# Merge Operator

Get(k1)

$5 + 2 - 1 = 6$

memtable
Merge k1 = +2
Merge k1 = -1

Merge Operator
callback

Put k1 = 5

# Merge Operator

memtable

Merge k1 = +2

Merge k1 = -1

Flush

Put k1 = 5

# Merge Operator

memtable
Merge k1 = +2
Merge k1 = -1

Merge Operator callback

Flush

+ 2 − 1 = +1

Merge k1 = +1

Put k1 = 5

# Merge Operator

Merge k1 = +1

Put k1 = 5

# Merge Operator

Merge k1 = +2

Merge k1 = +1

Put k1 = 5

# Merge Operator



Merge k1 = +2

Merge k1 = +1

Merge k1 = +3

Put k1 = 5

# Merge Operator

Merge k1 = +3

Put k1 = 8

Put k1 = 5

# Other Features

# Other Features

- Transactions
- Column Families
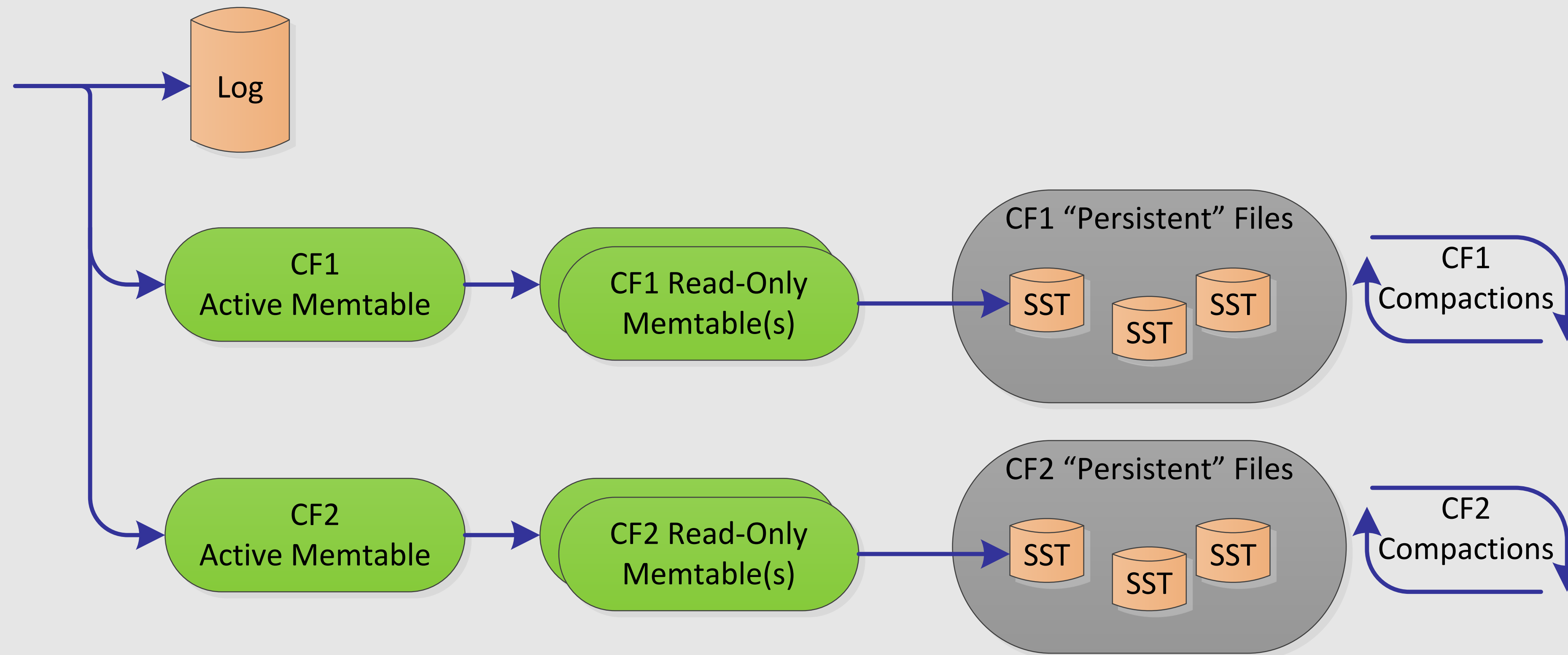- Monitoring
- And many more...

# Transactions

- Optimistic and Pessimistic Transactions
- Optimistic: verify conflicts with memtable values
- Pessimistic: lock keys being modified.

# Column Families

- Separate key spaces
- Allow atomic multiput, multiget, transactions
- Can apply different compaction setting, comparators, compaction filters, merge operators, etc.

# Column Families

# Monitoring

- Statistics
- DB Properties
- Per request profiling (perf context)
- User defined per SST file properties

# Conclusion

- RocksDB is widely used
- RocksDB uses LSM-tree
- RocksDB is highly tunable for flash
- RocksDB can be tuned to be space efficient
- RocksDB has good performance
- RocksDB has nice features

# Thank You!

- Portal: http://rocksdb.org/
- Github: https://github.com/facebook/rocksdb
- Discussion Group:

https://www.facebook.com/groups/rocksdb.dev/