

SLM-DB: Single Level Merge Key-Value Store with Persistent Memory

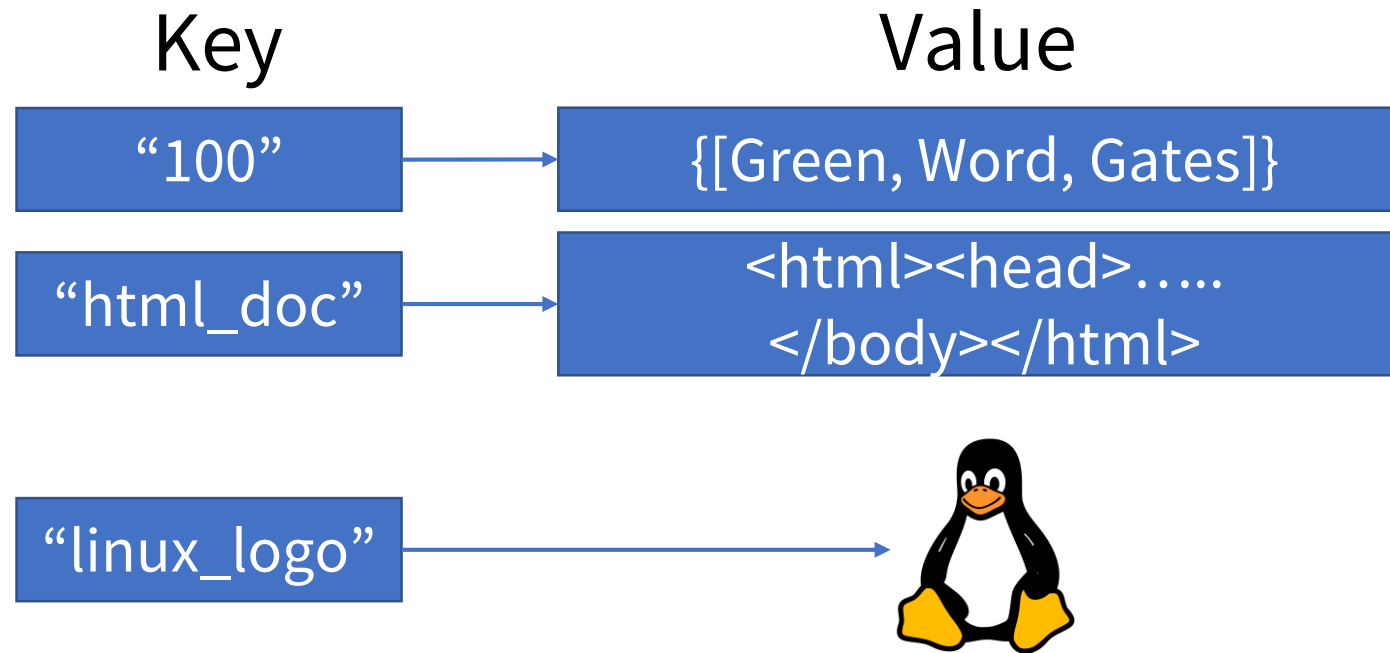
Olzhas Kaiyrakhmet, Songyi Lee, Beomseok Nam, Sam H. Noh, Young-ri Choi



Outline

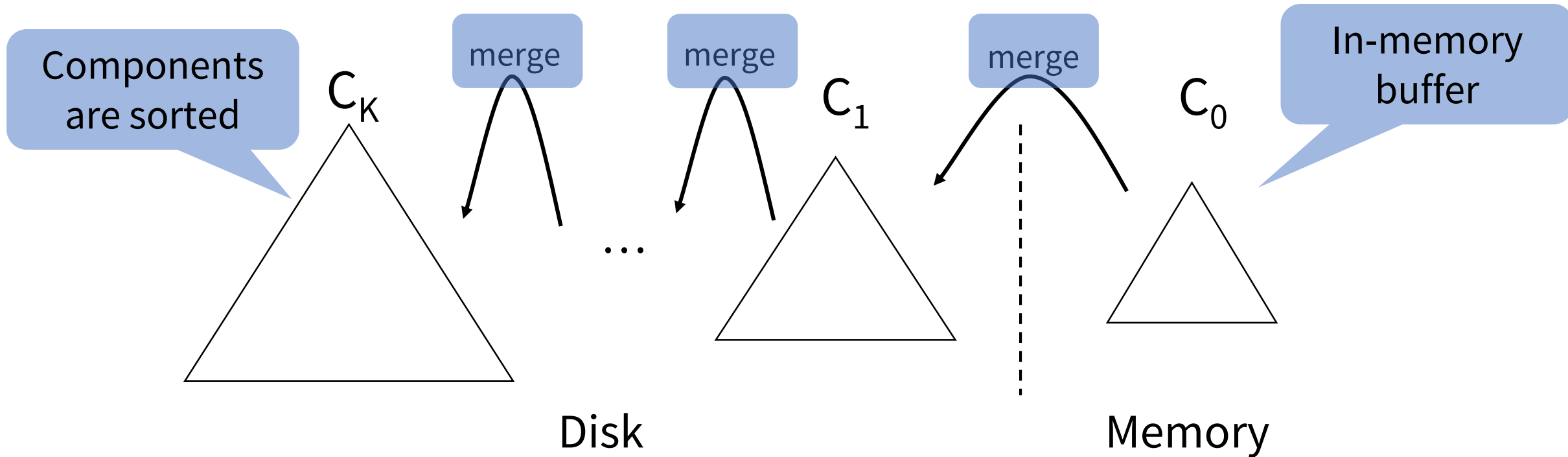
- Background
- Contributions
- Architecture
- Evaluation
- Conclusion

Key-Value Databases

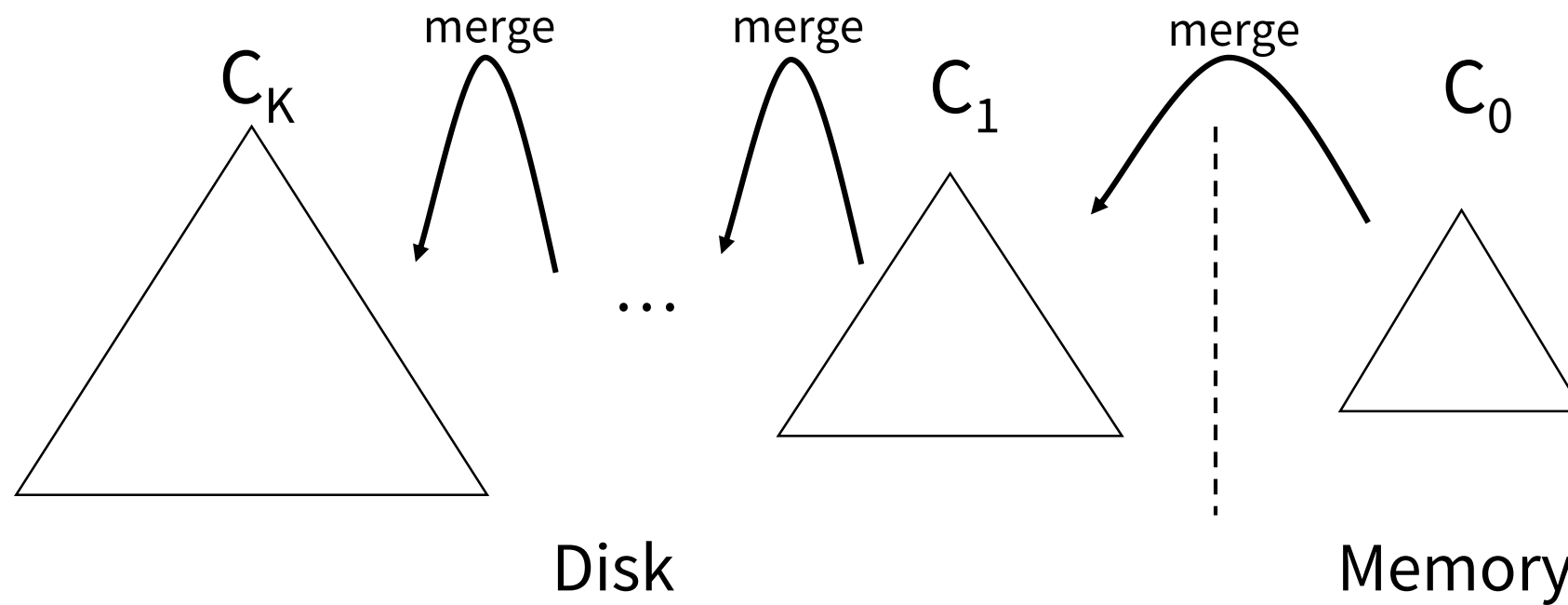


Log-Structured Merge (LSM) Tree

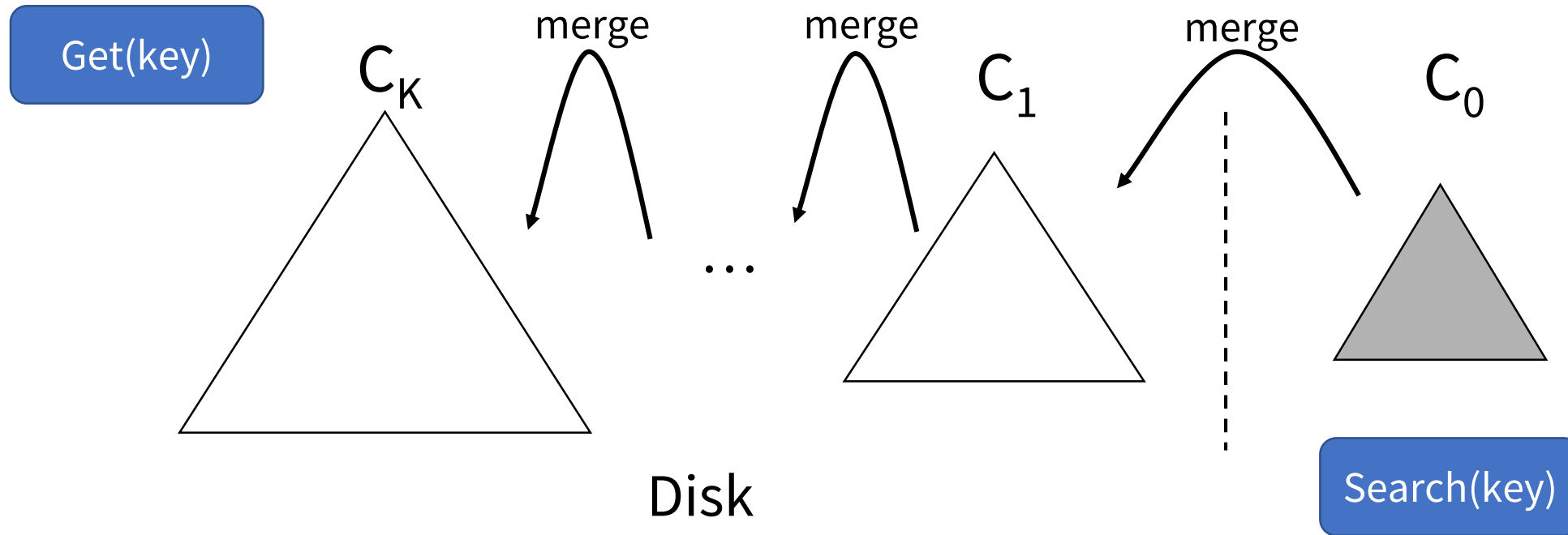
- Optimized for heavy write application usage
- Designed for slow hard drives



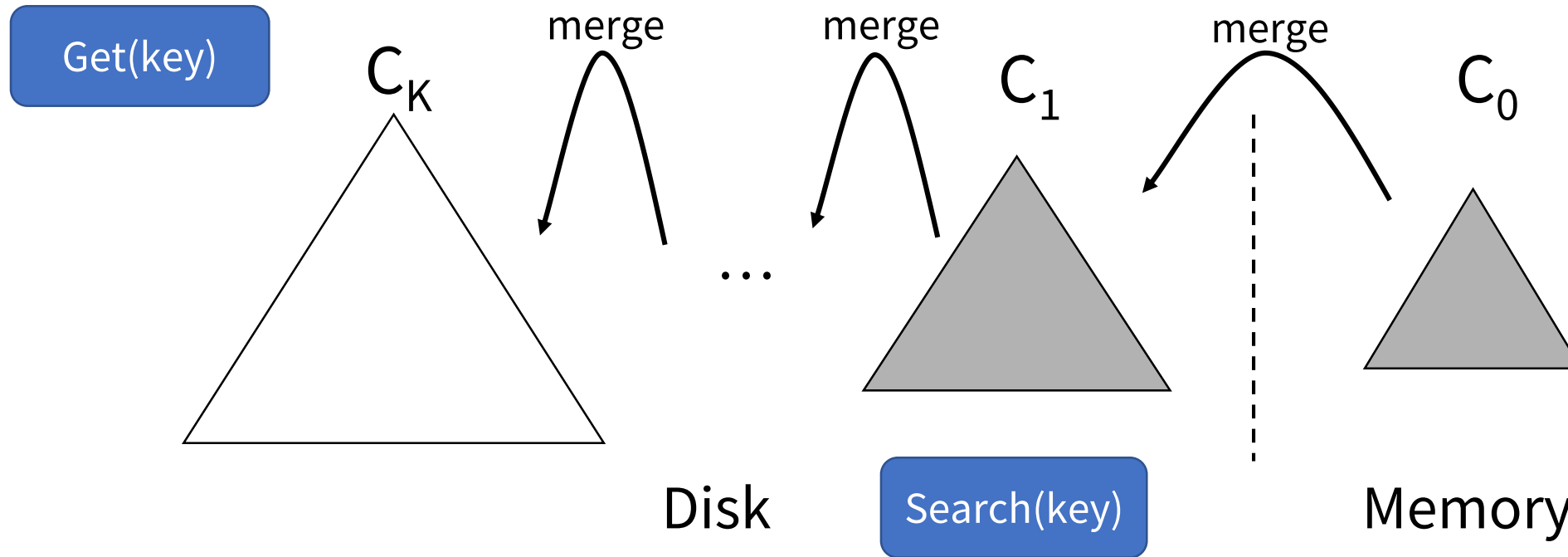
LSM-tree: disadvantages



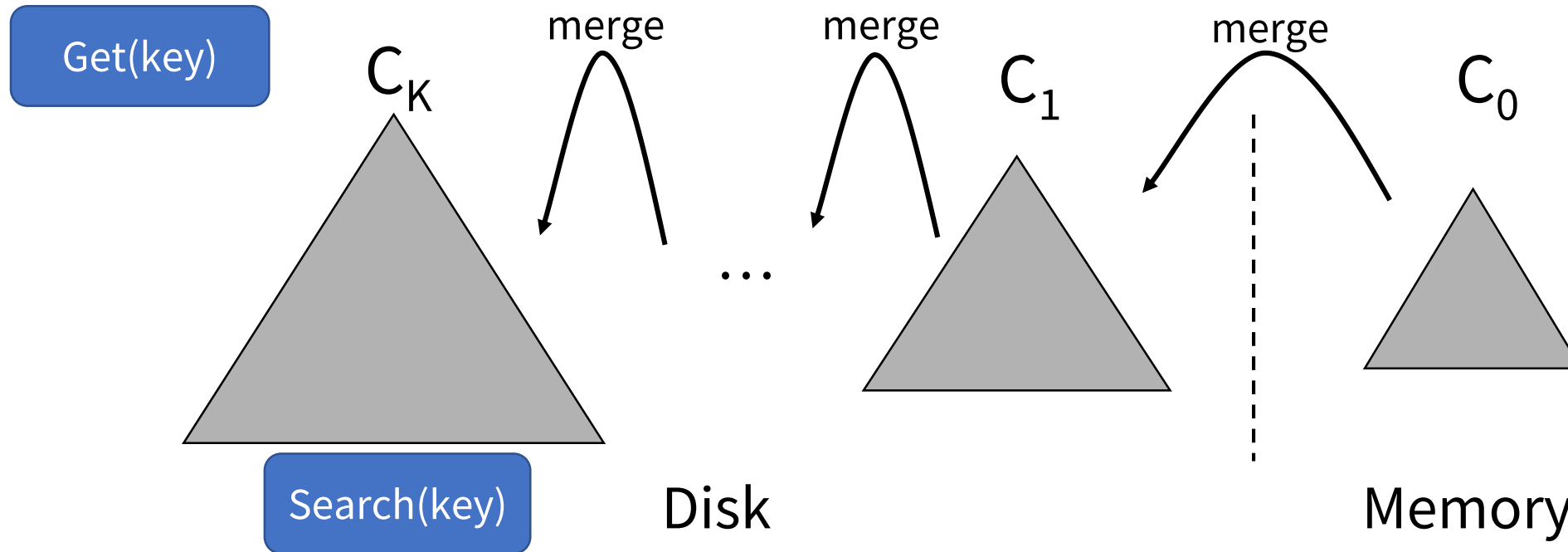
LSM-tree: disadvantages



LSM-tree: disadvantages

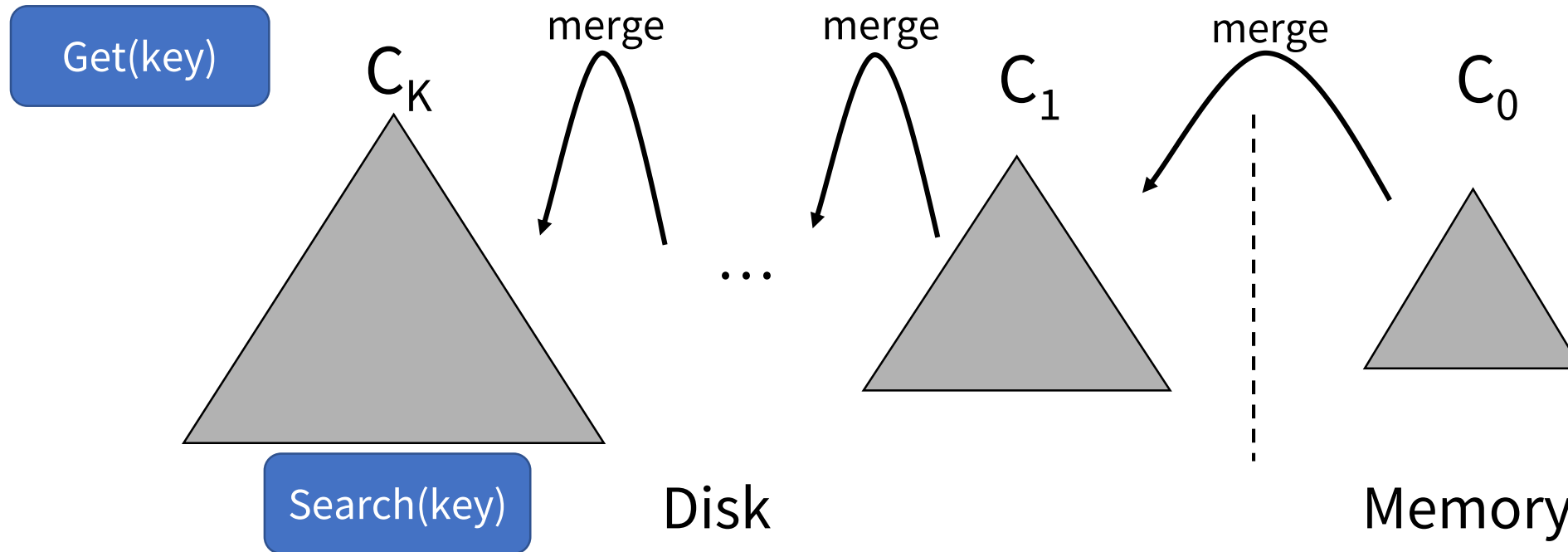


LSM-tree: disadvantages



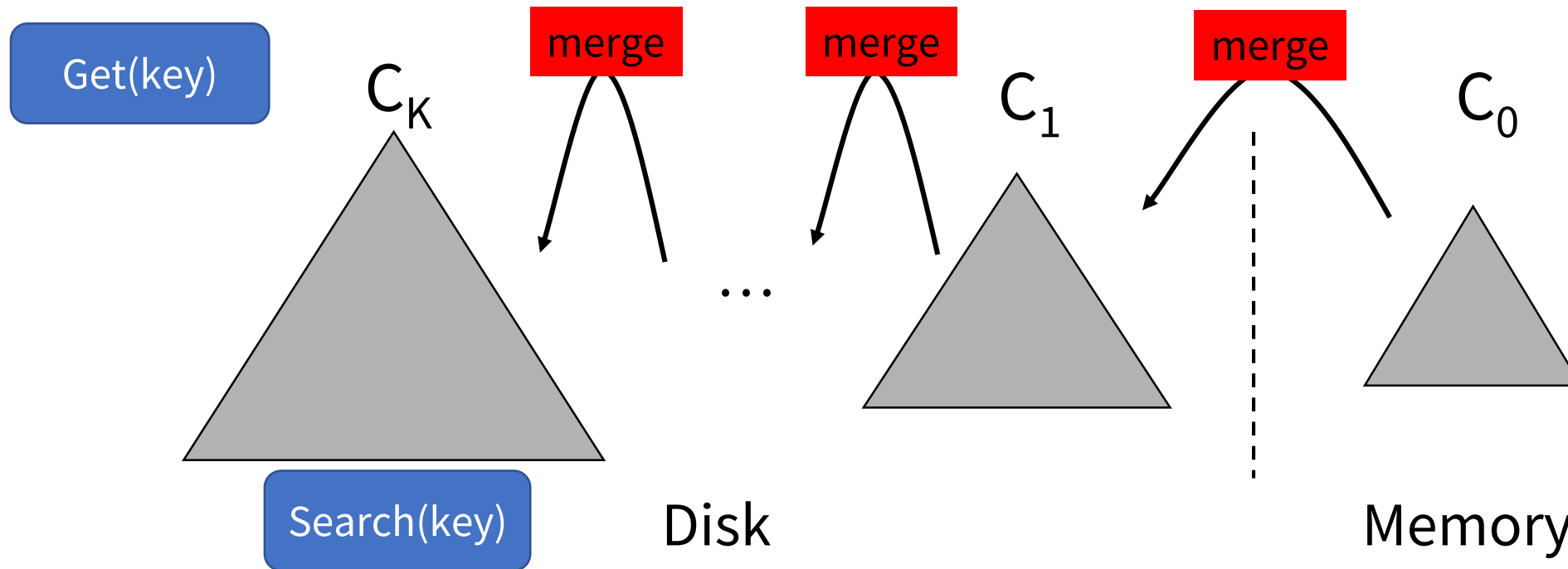
LSM-tree: disadvantages

- Large overhead to locate needed data



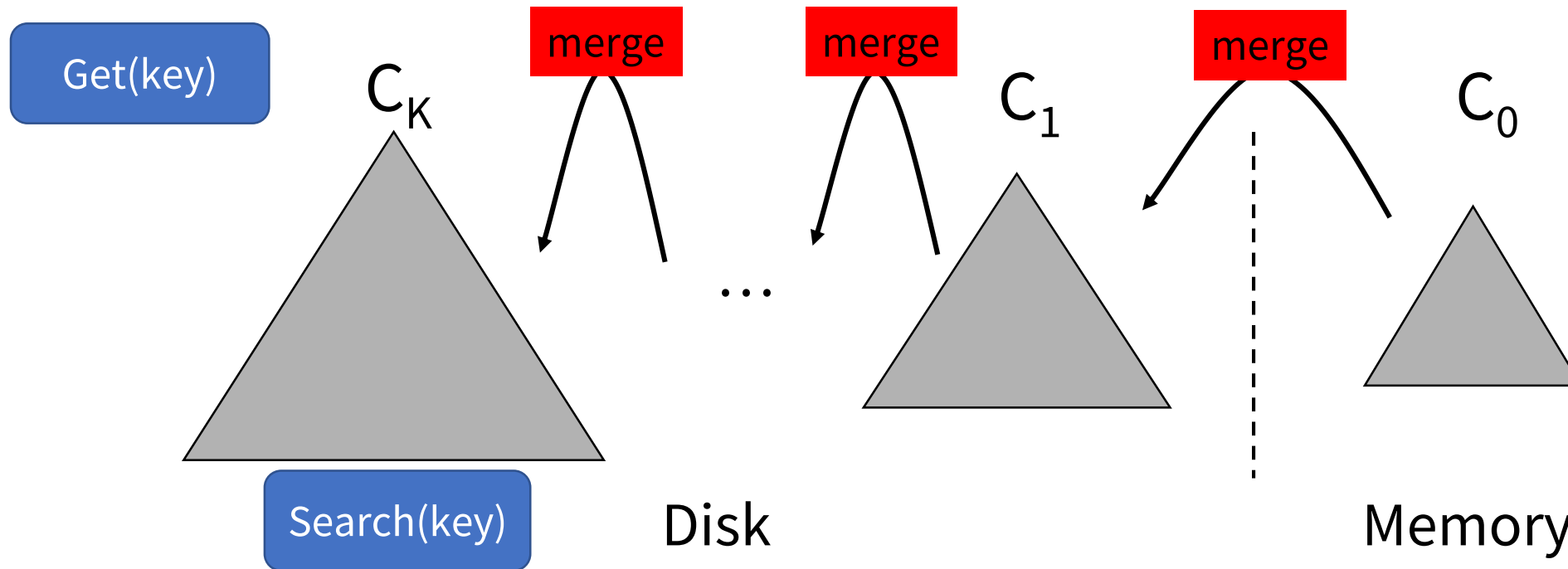
LSM-tree: disadvantages

- Large overhead to locate needed data

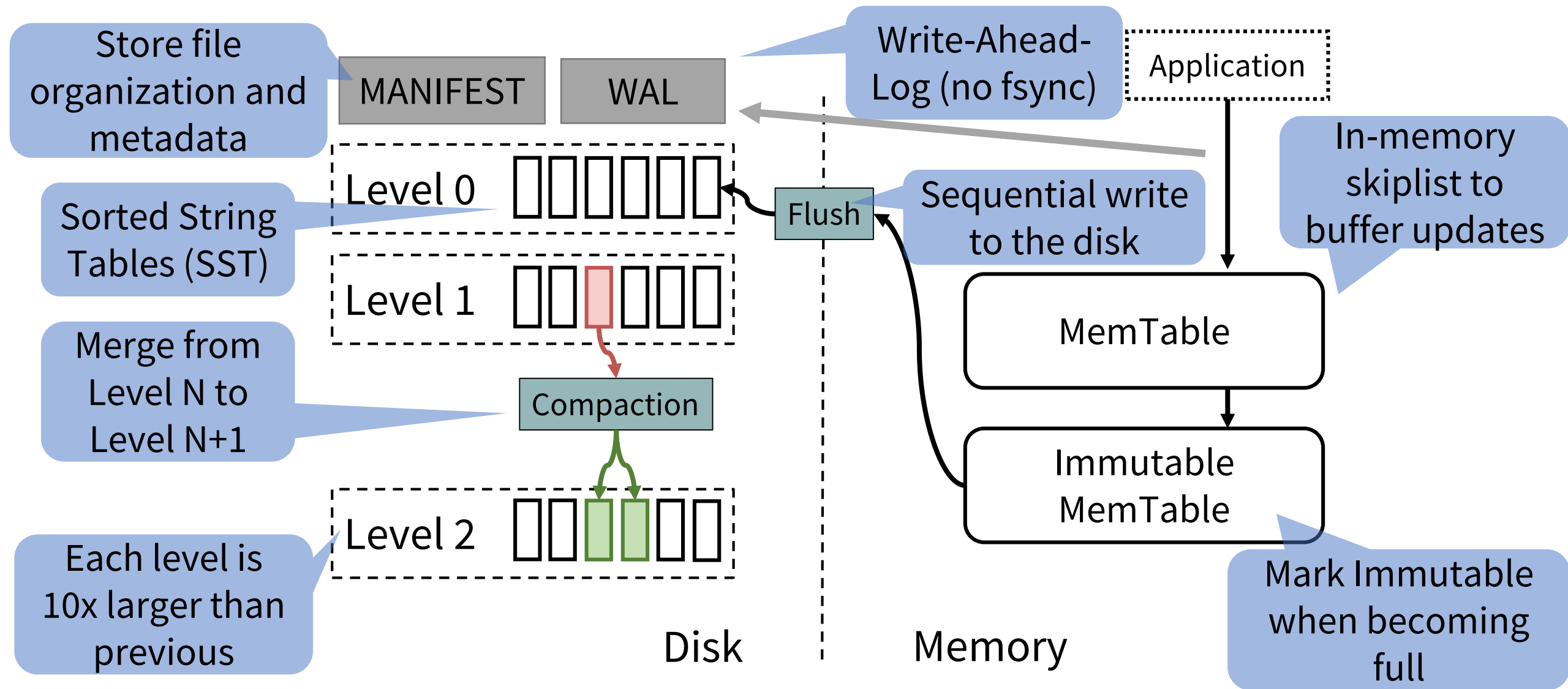


LSM-tree: disadvantages

- Large overhead to locate needed data
- High disk write amplification



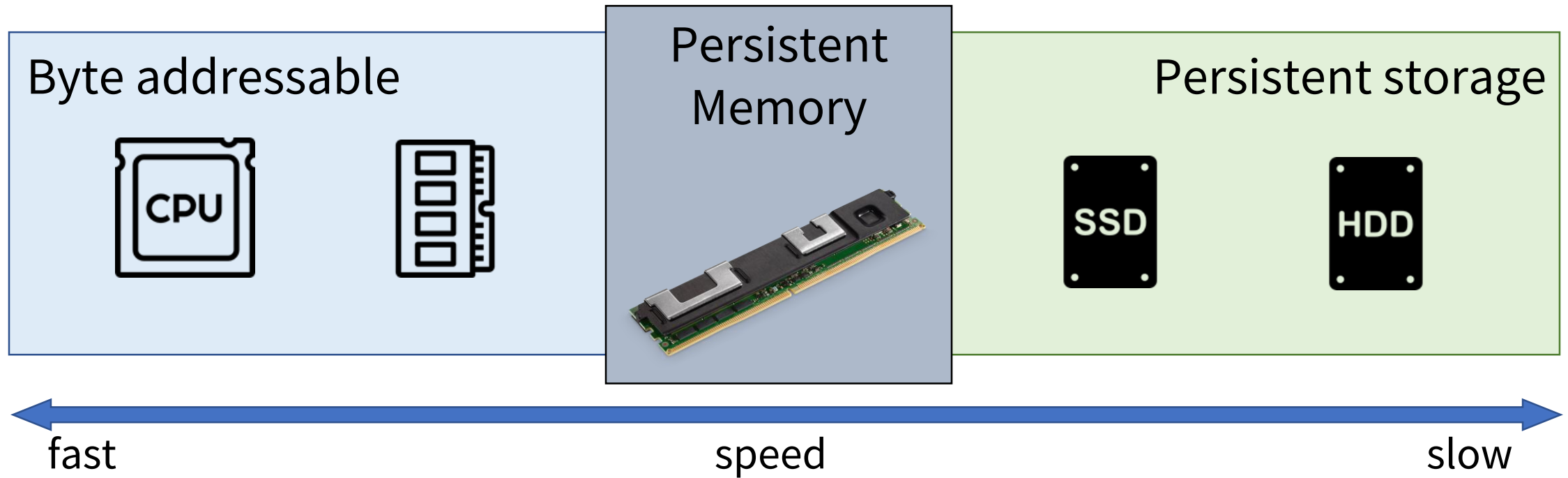
State-of-the-art LSM-tree: LevelDB



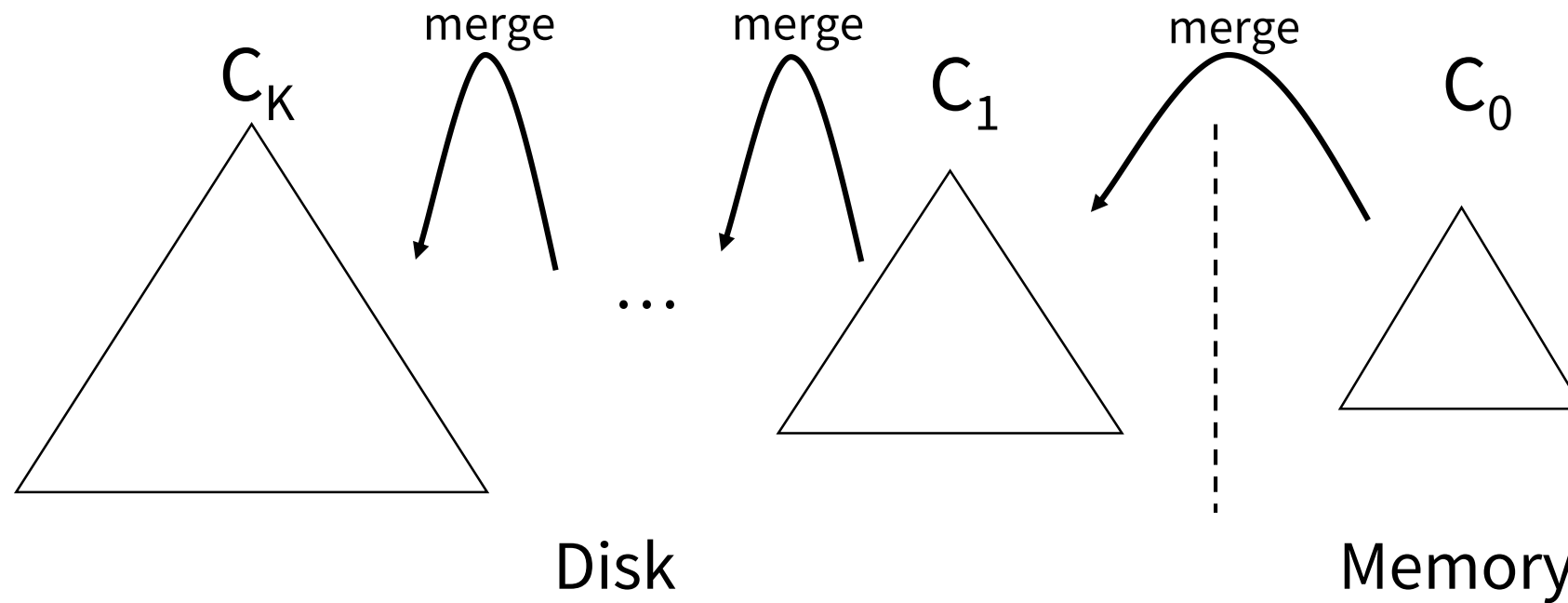
LSM-tree optimizations

- Improve parallelism:
 - RocksDB (Facebook)
 - HyperLevelDB
- Reduce write amplification:
 - PebblesDB (SOSP '17)
- Optimize for hardware(SSD):
 - VT-tree (FAST '13)
 - WiscKey (FAST '16)

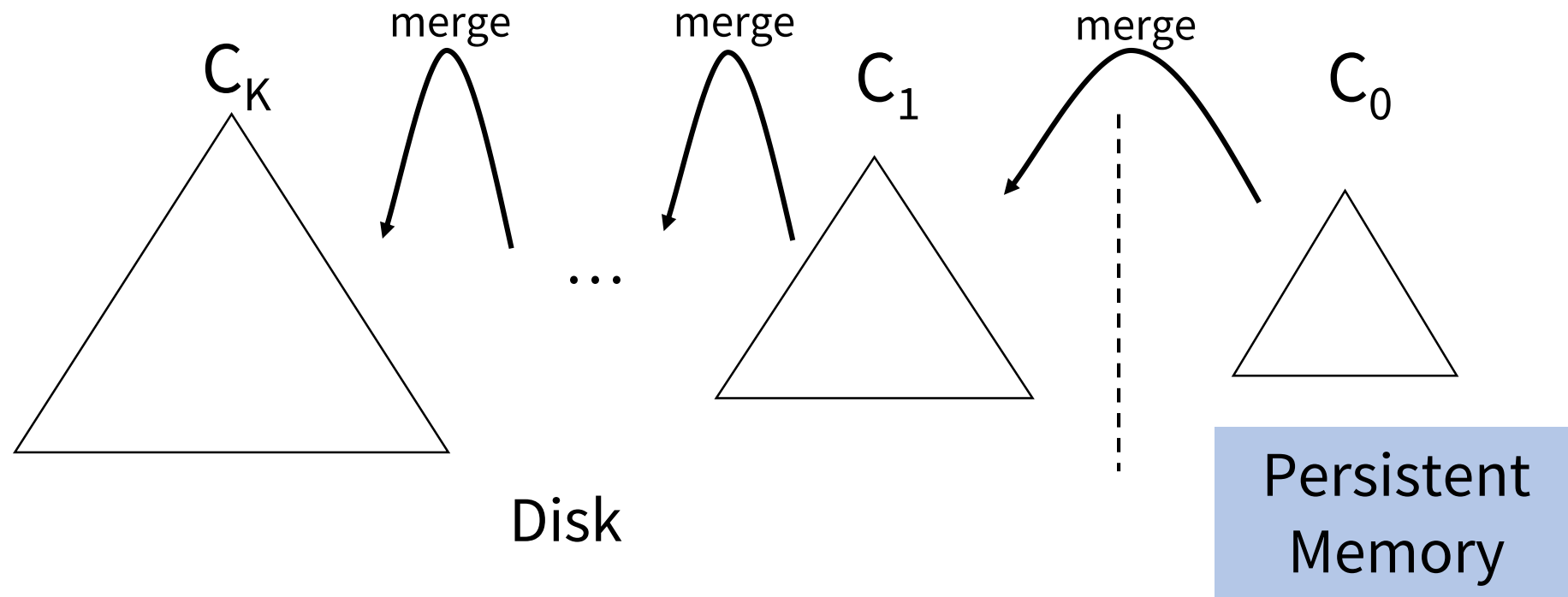
New era



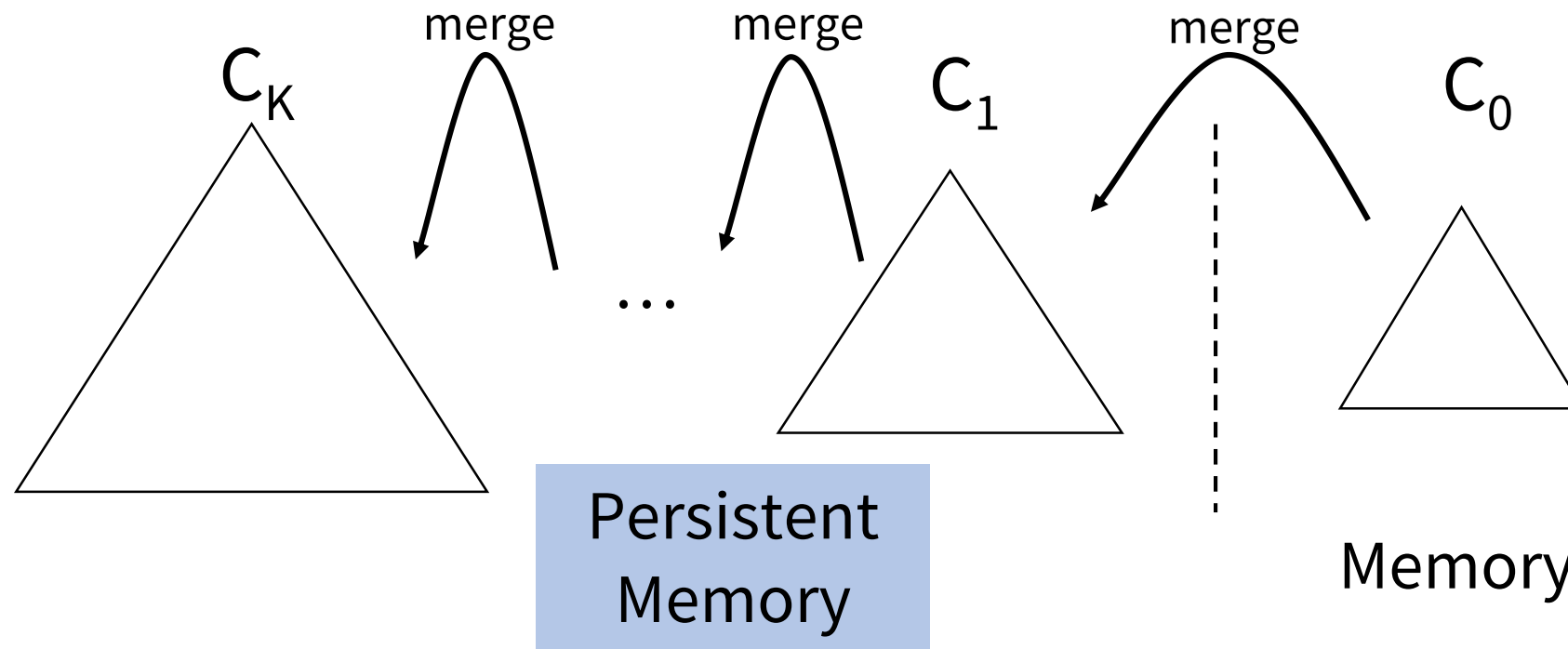
Simple approach



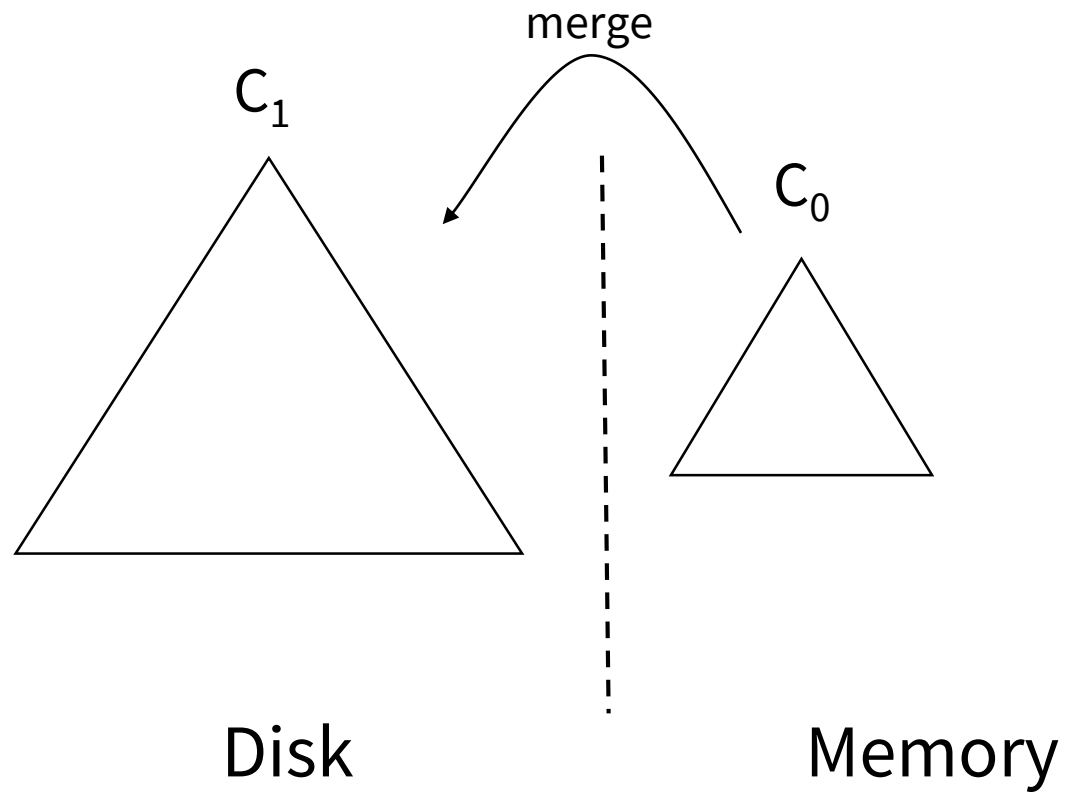
Simple approach



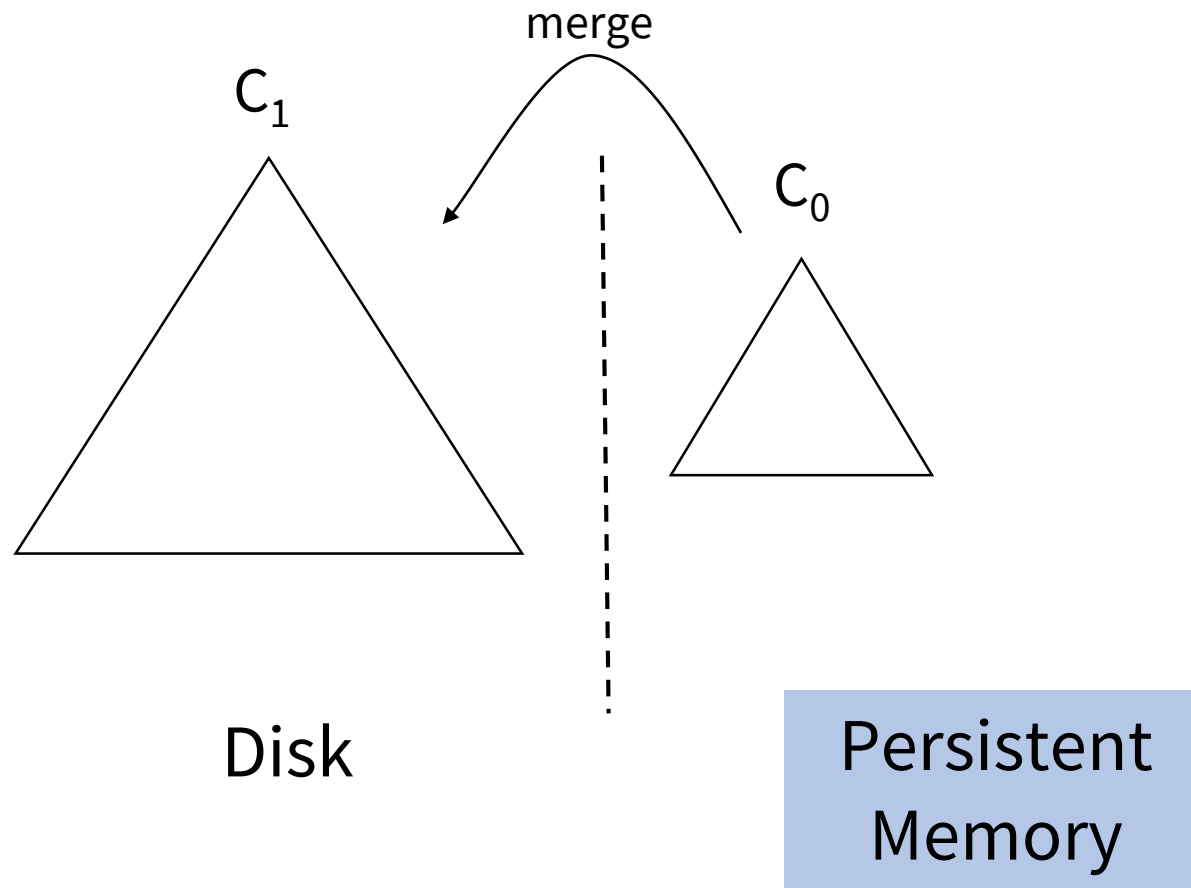
Simple approach



Our approach

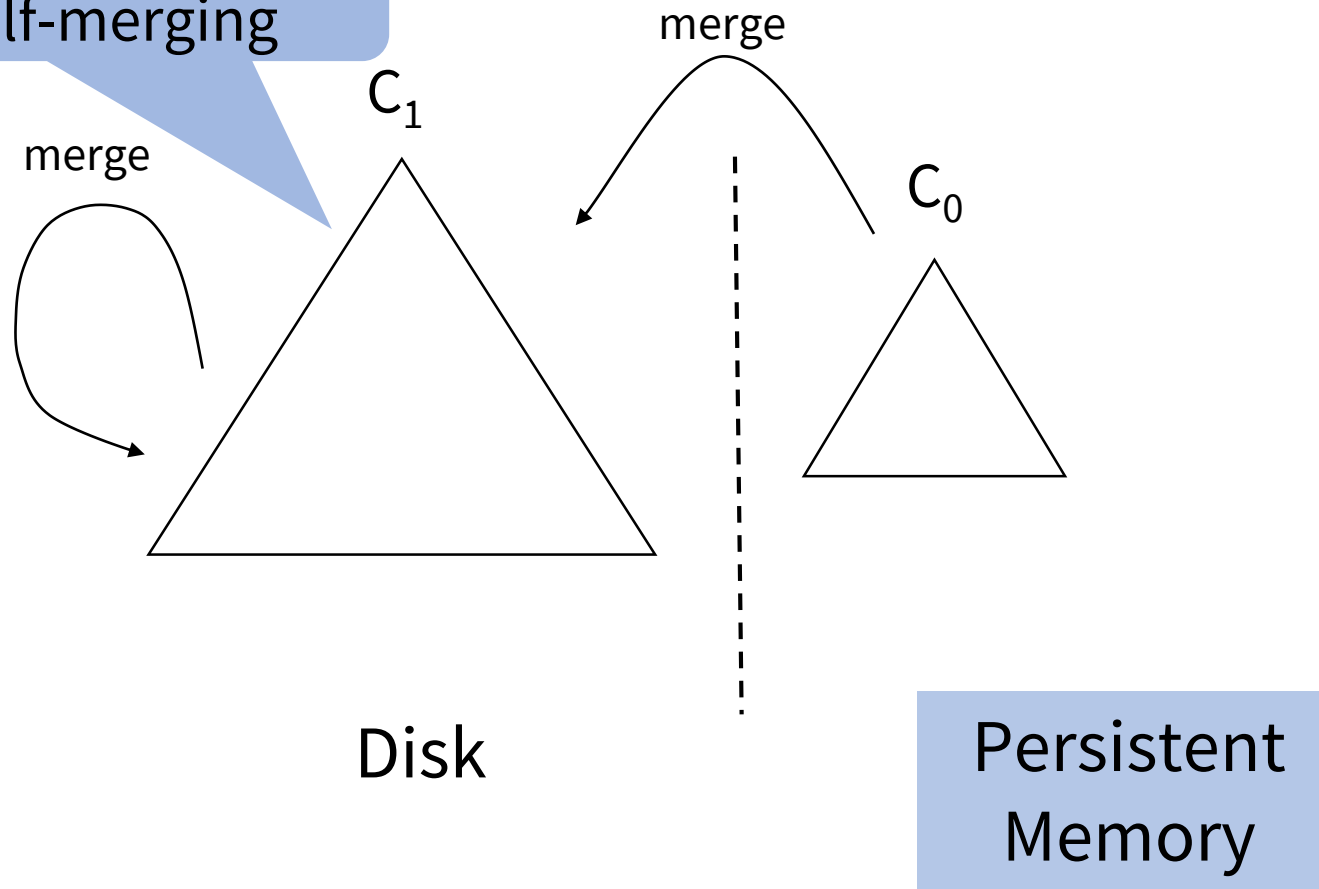


Our approach



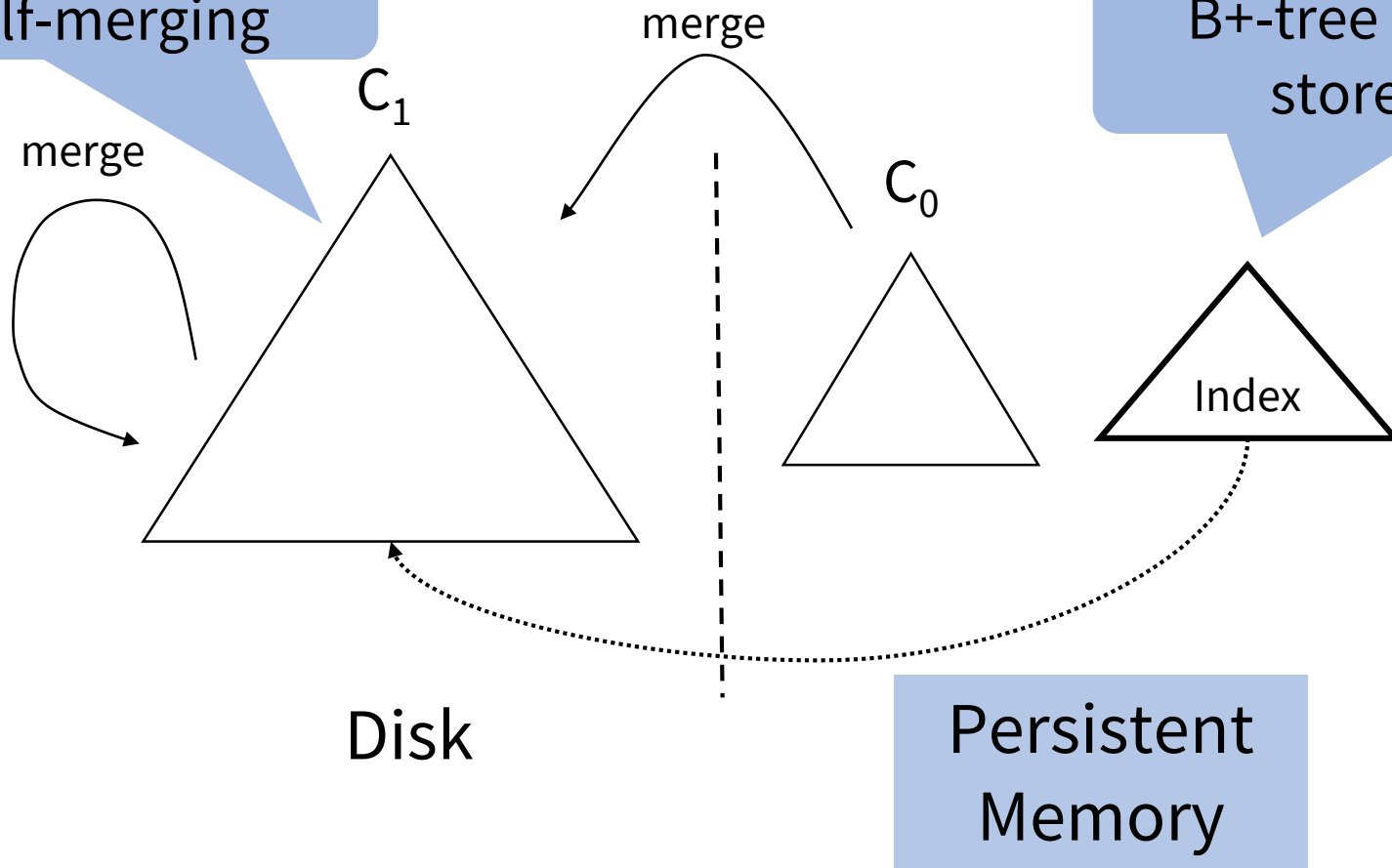
Our approach

Single disk component C_1
that does self-merging

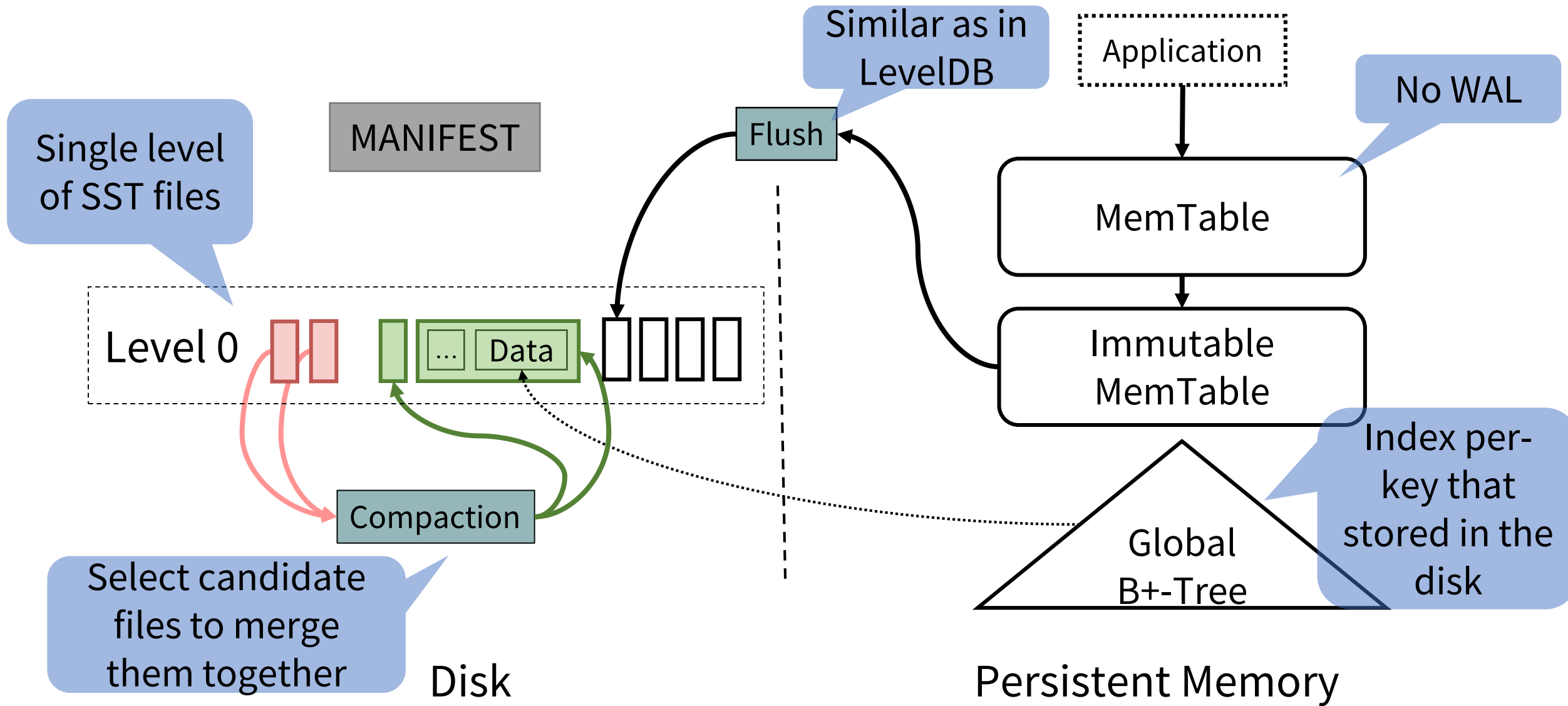


Our approach

Single disk component C_1
that does self-merging



Single-Level Merge DB (SLM-DB)



Contributions

Persistent MemTable

No Write-Ahead Logging (WAL)
Stronger consistency compared to LevelDB

Persistent B+-tree Index

Per-key index for fast search
No multi-leveled merge structure

Selective Compaction

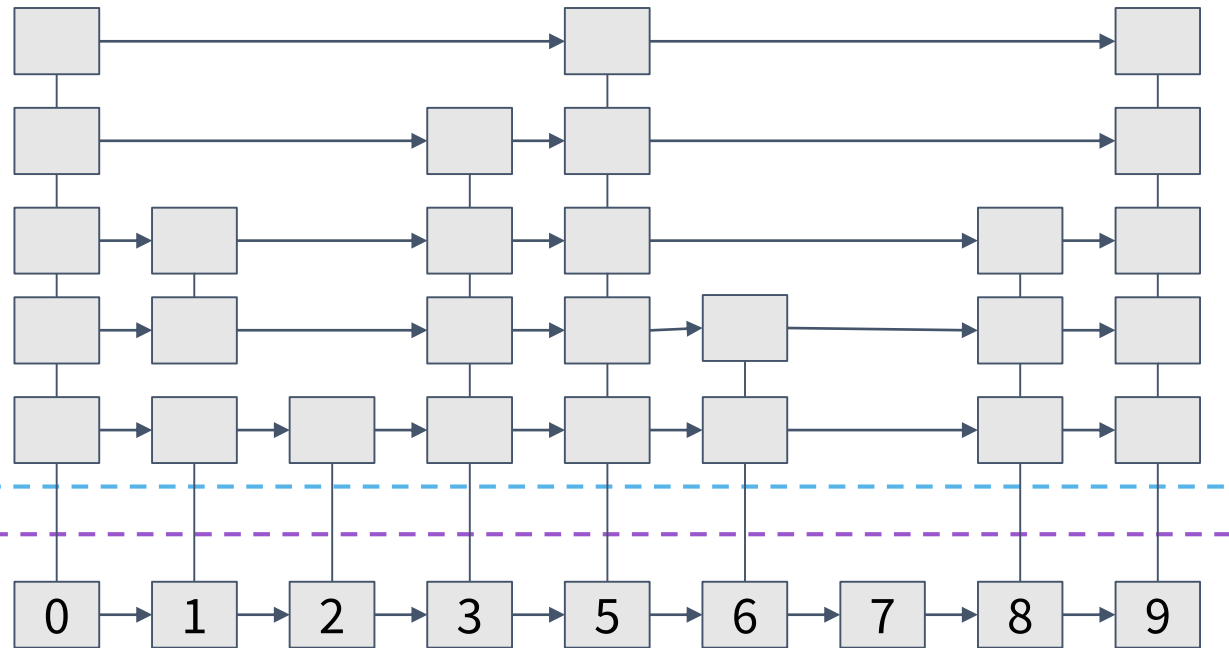
Live-key ratio of a Sorted-String Table
Leaf node scan in the B+-tree
Degree of sequentiality per range query

Persistent MemTable

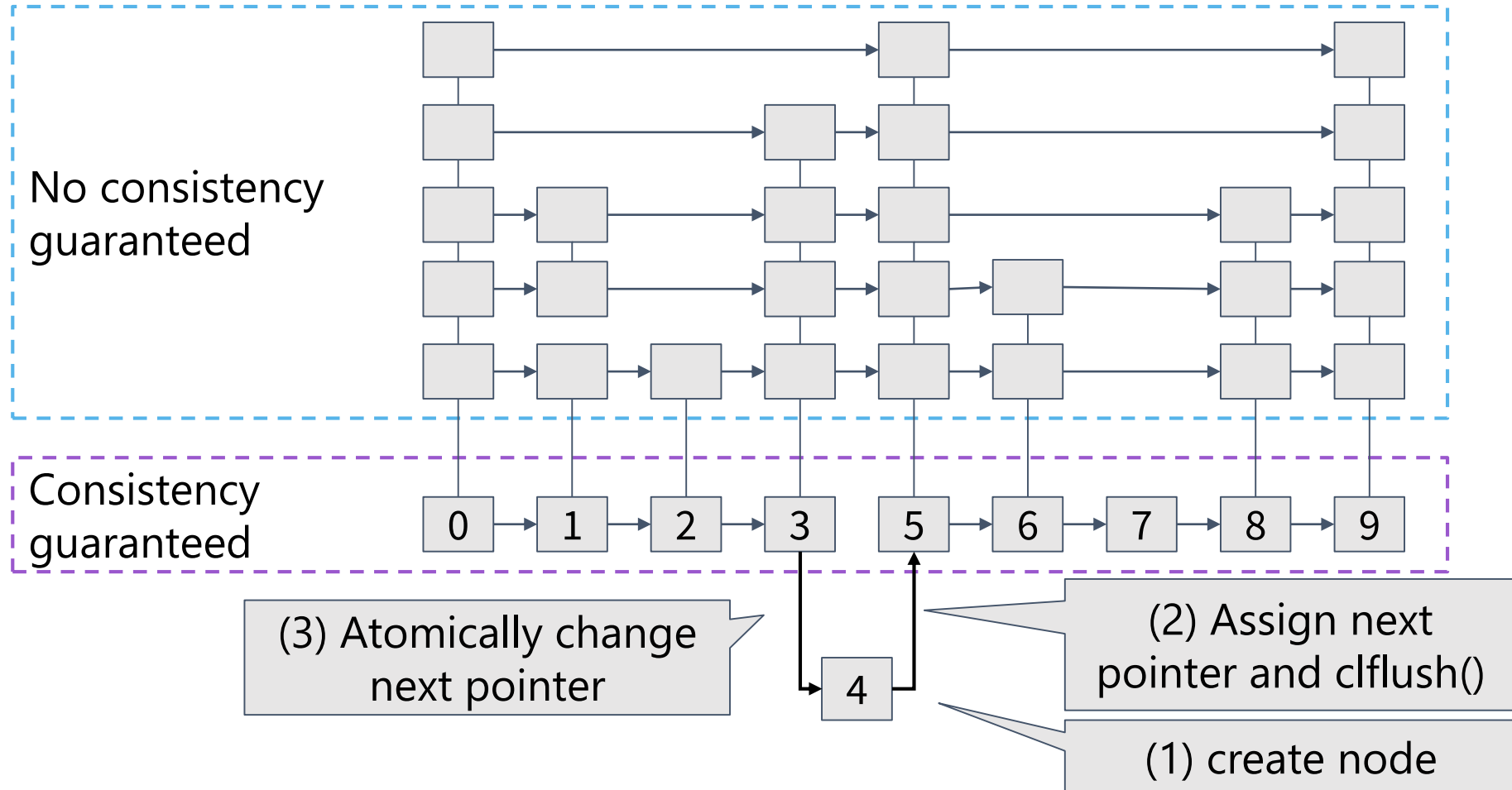
Recoverable
after failure

No consistency
guaranteed

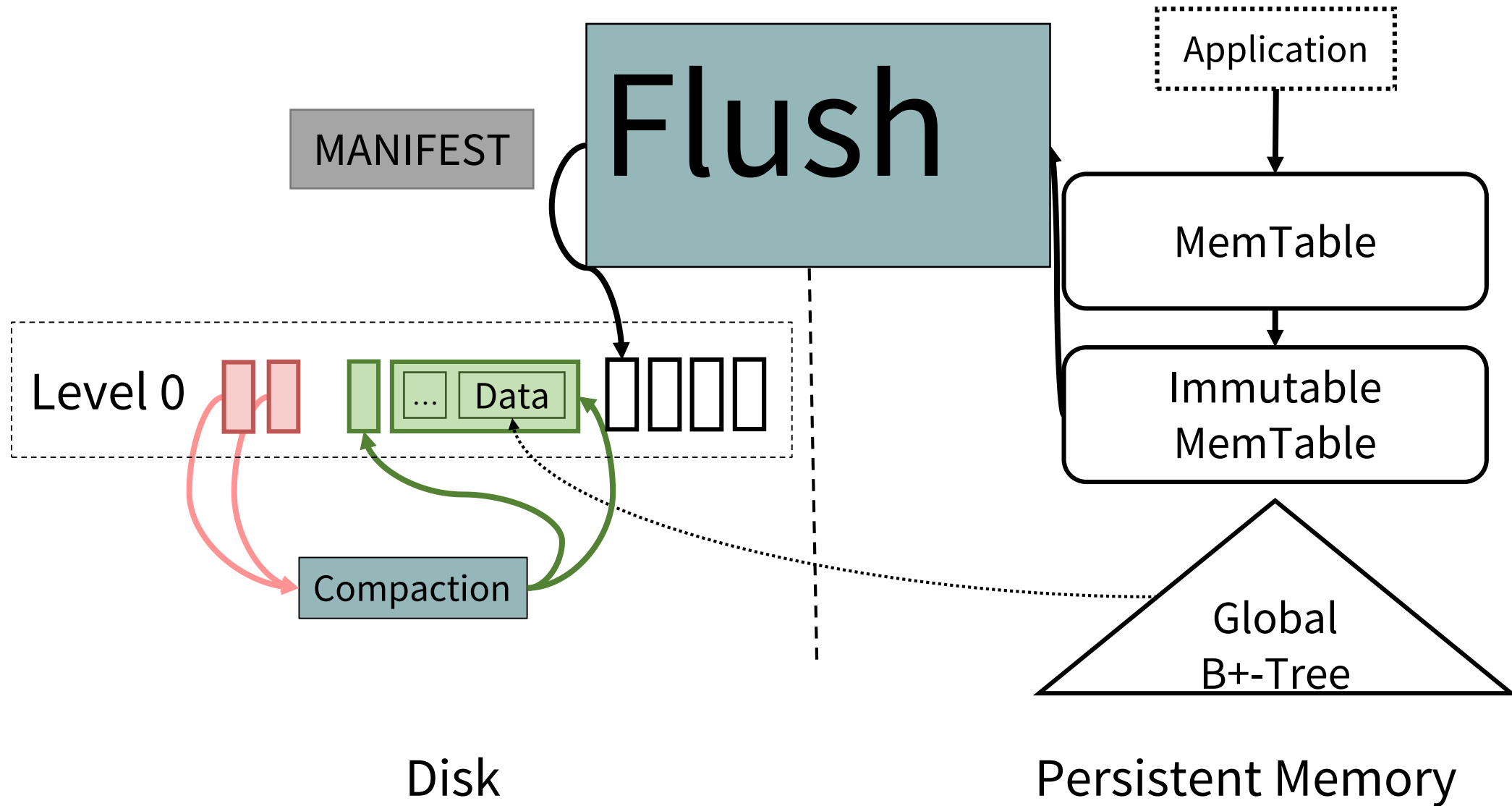
Consistency
guaranteed



Insert into Persistent MemTable

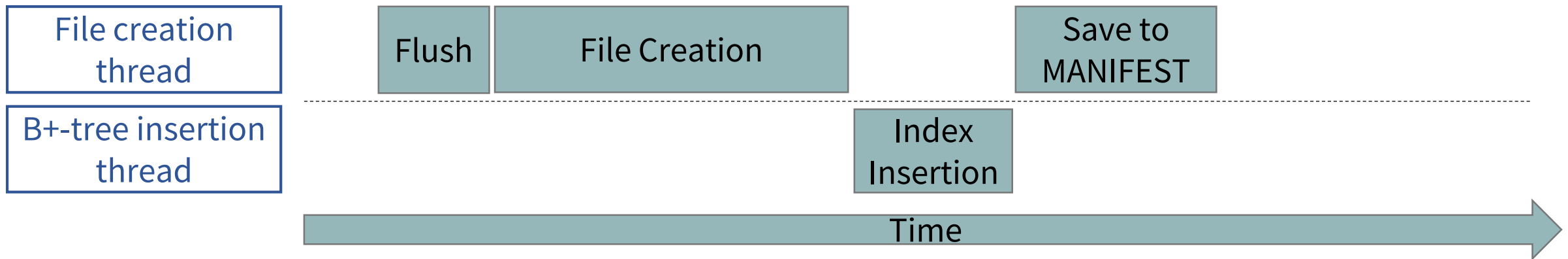


Single-Level Merge DB

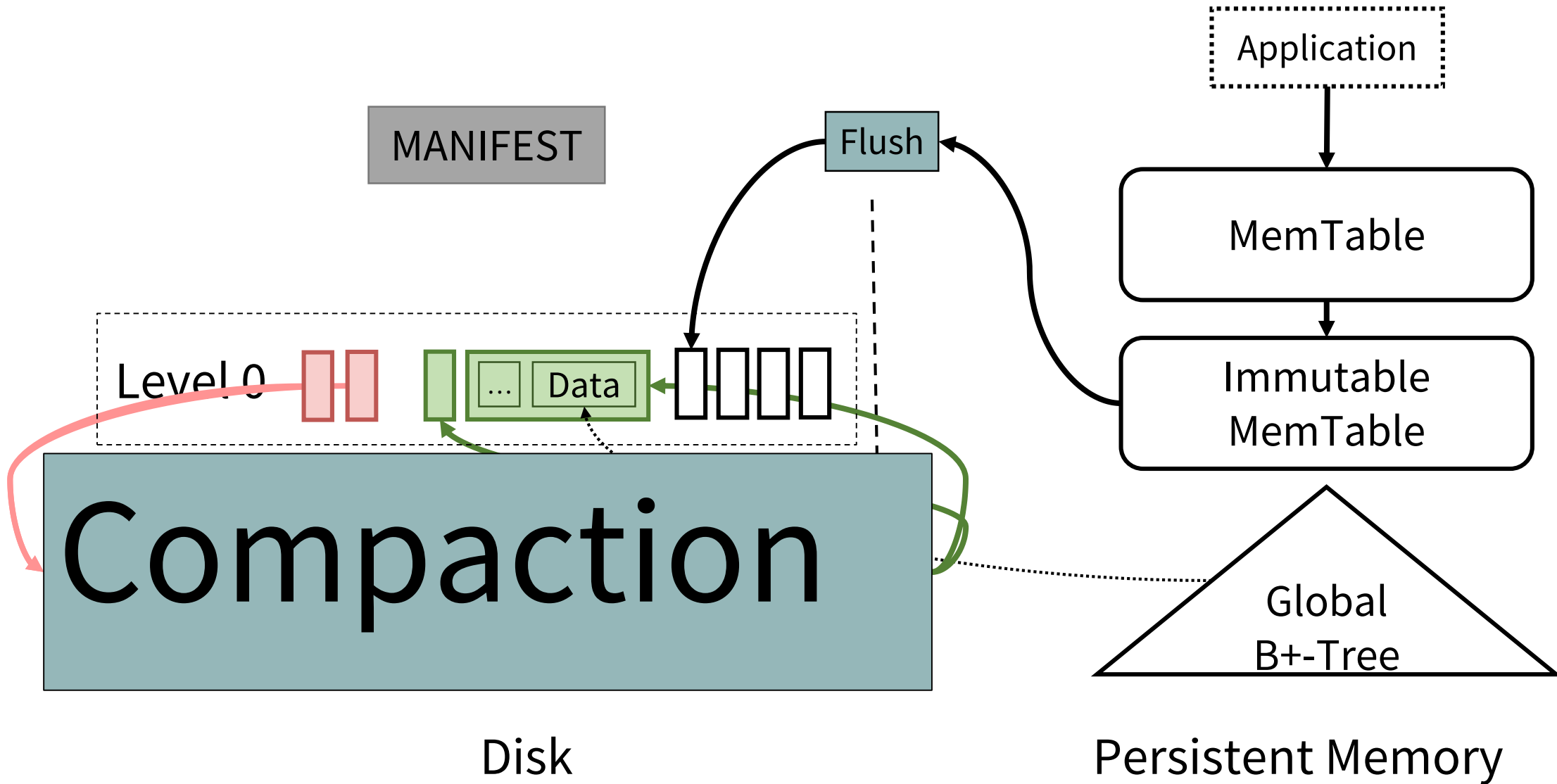


Flush



- Key-Index insertion into B+-tree happens during Immutable Memtable Flush to disk
- FAST-FAIR B+-tree (Hwang et al., FAST '18)

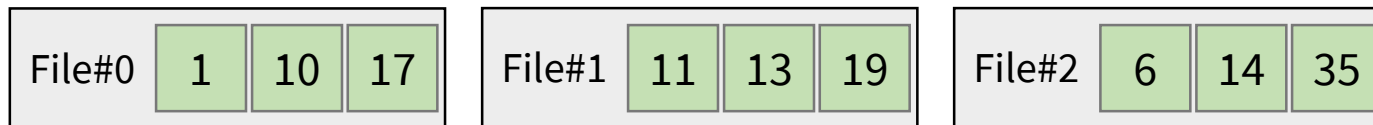


Single-Level Merge DB

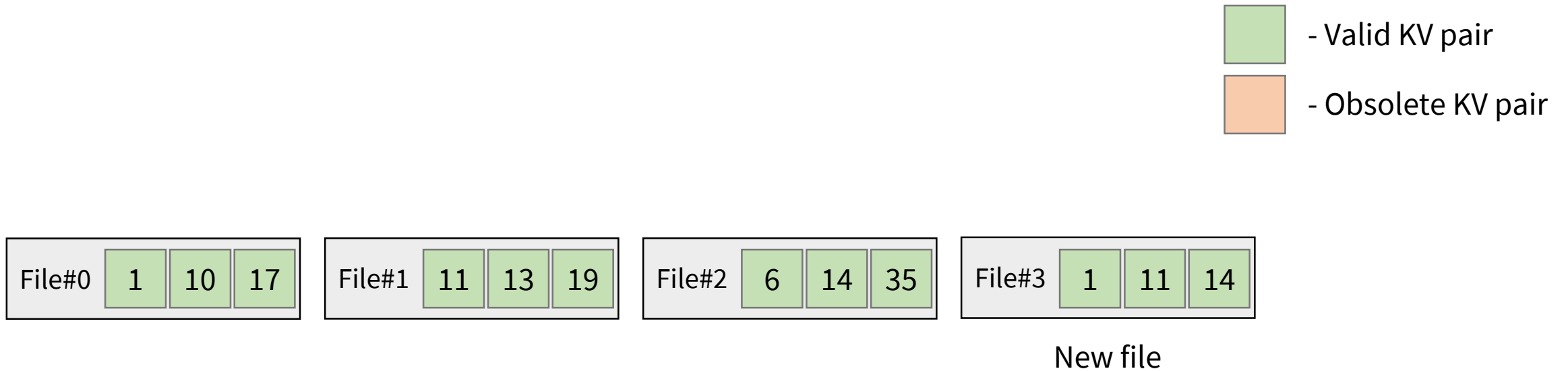


Why we need **Compaction**?

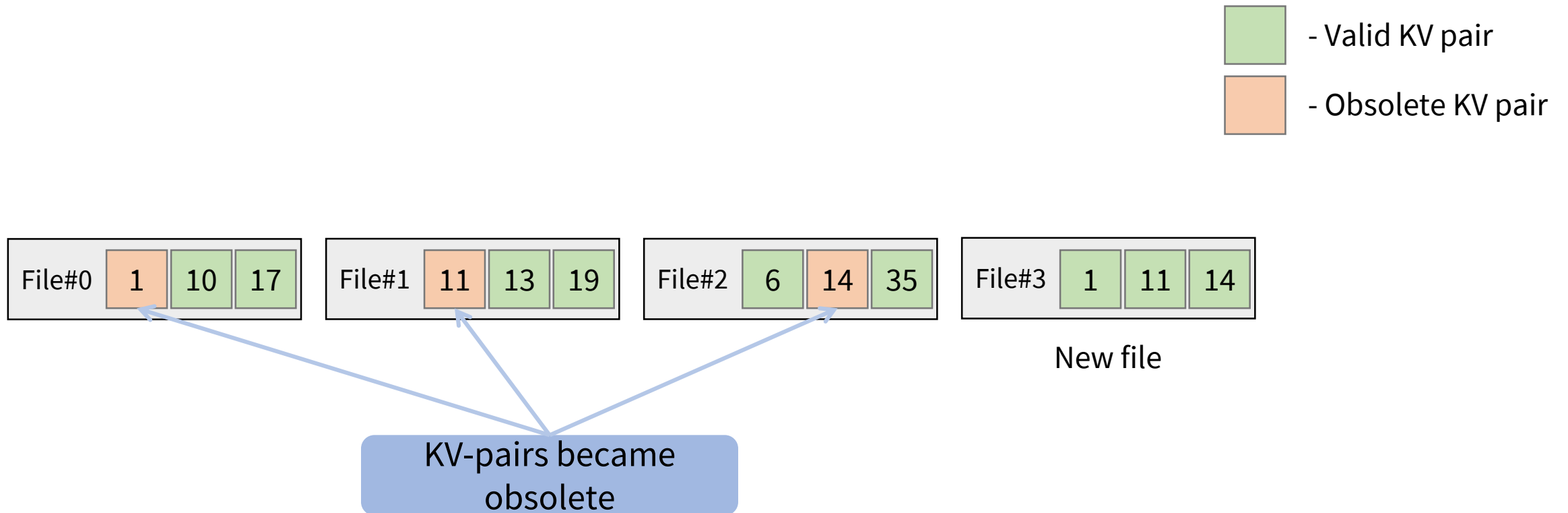
-  - Valid KV pair
-  - Obsolete KV pair



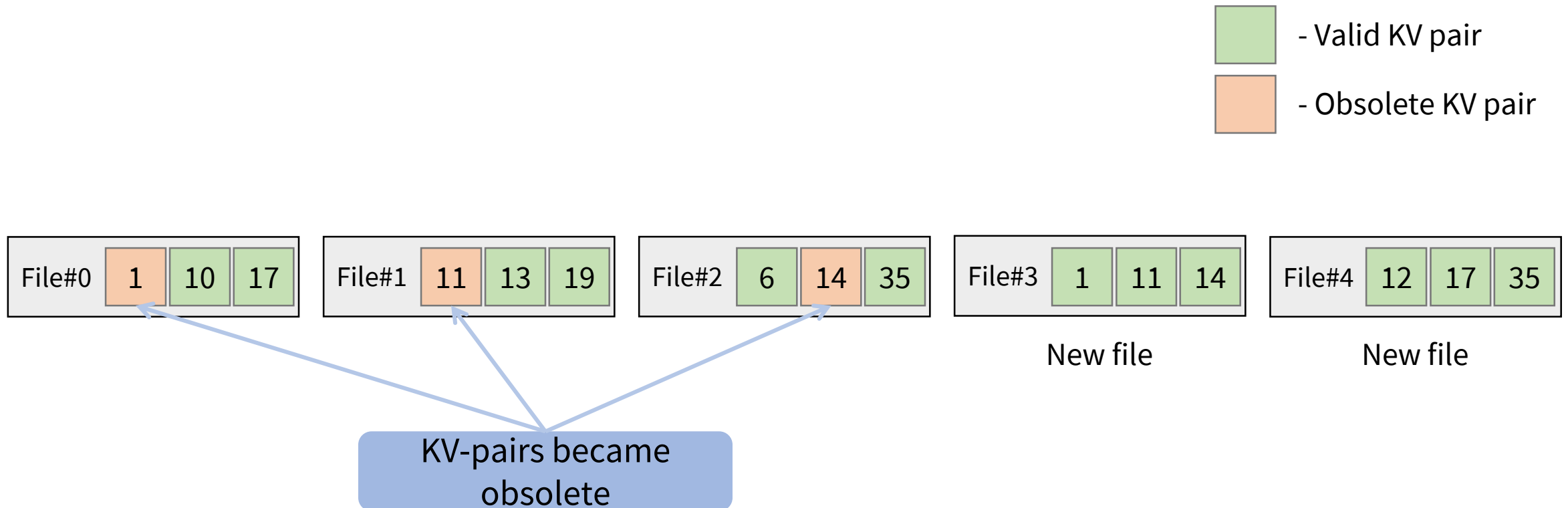
Why we need **Compaction**?



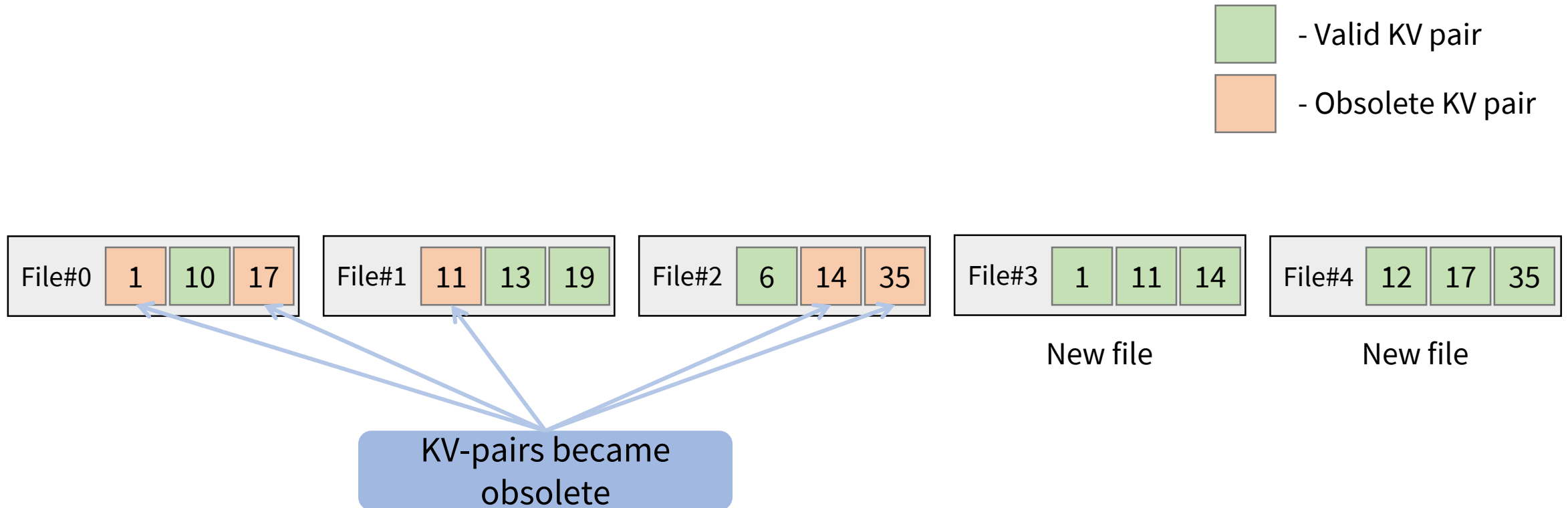
Why we need **Compaction**?



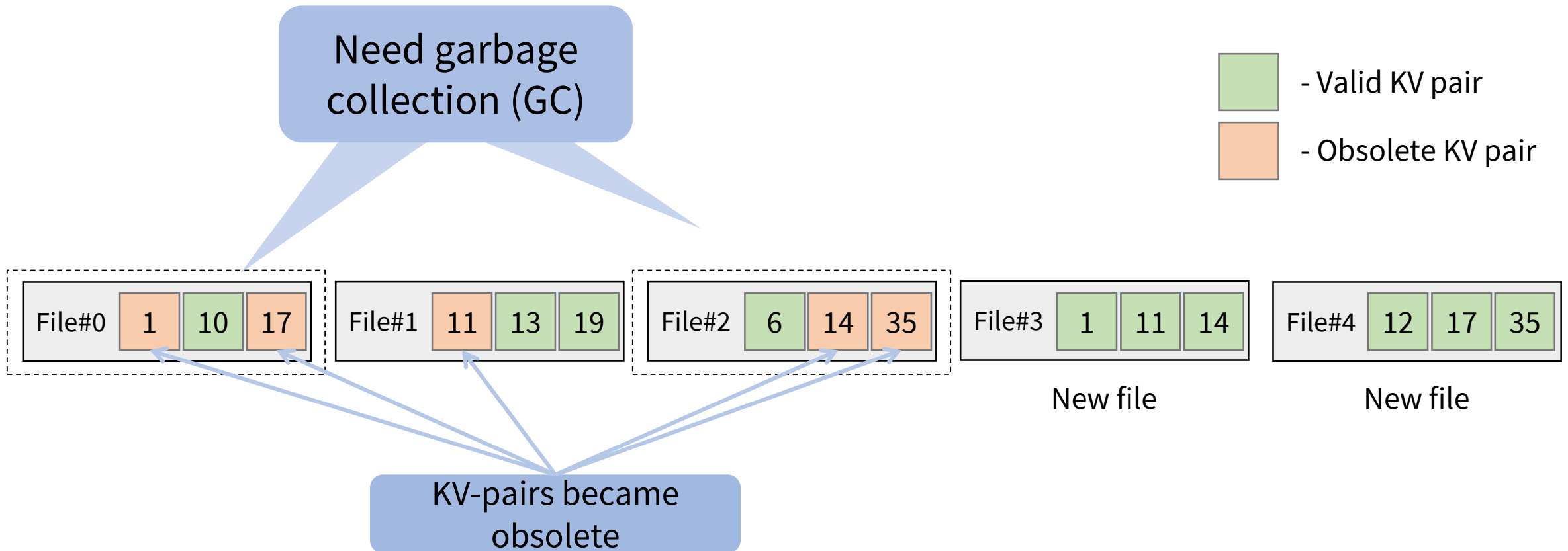
Why we need **Compaction**?



Why we need **Compaction**?



Why we need **Compaction**?



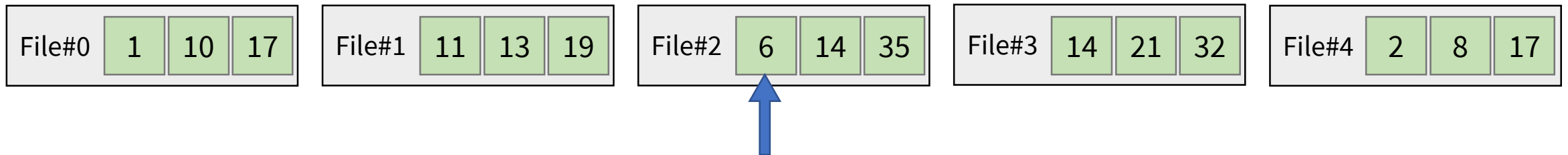
Why else?

RangeQuery(5, 12)

File#0	1	10	17
File#1	11	13	19
File#2	6	14	35
File#3	14	21	32
File#4	2	8	17

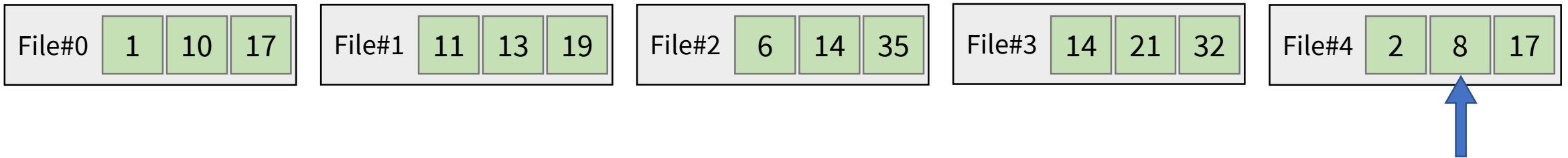
Why else?

RangeQuery(5, 12)



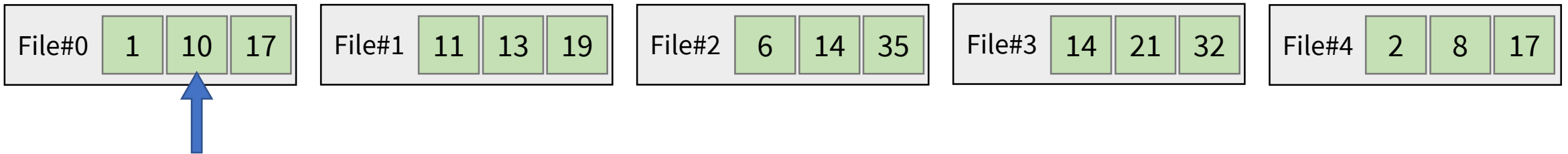
Why else?

RangeQuery(5, 12)



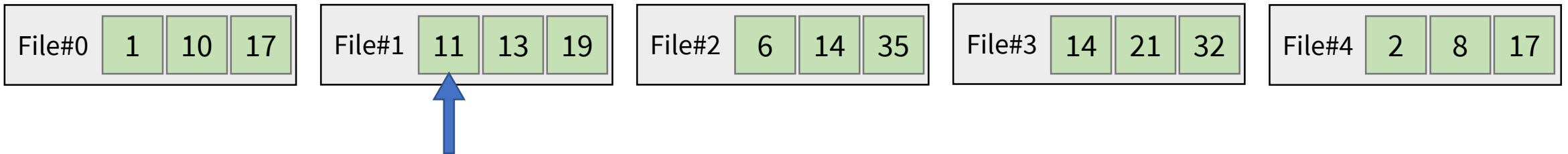
Why else?

RangeQuery(5, 12)



Why else?

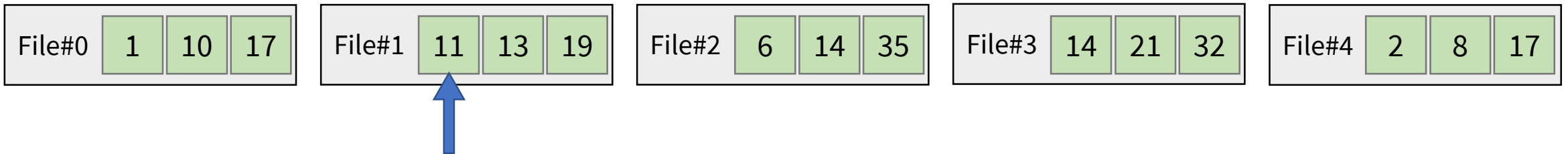
RangeQuery(5, 12)



Why else?

RangeQuery(5, 12)

Need to improve
sequentiality

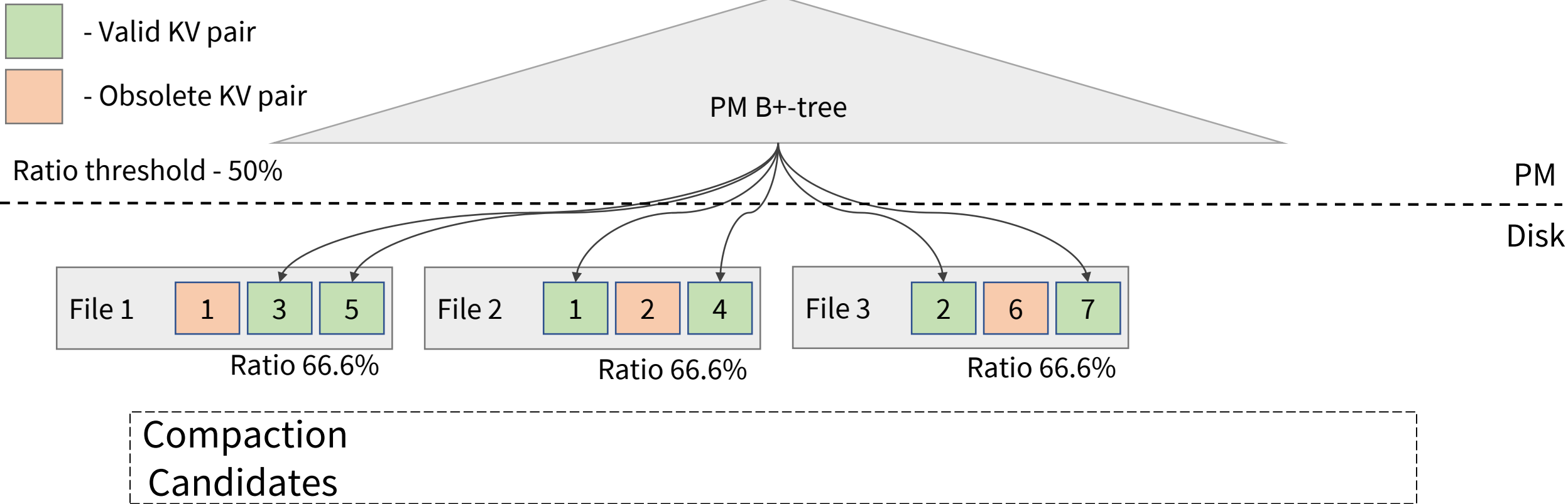


Selective compaction

- Selectively pick SSTable files
- Make those files as compaction candidates
- Merge together most overlapping compaction candidates
- Selection schemes for compaction candidates:
 - Live-key ratio selection of an SSTable (for GC)
 - Leaf node scans in the B+-tree (for sequentiality) [see paper]
 - Degree of sequentiality per range query (for sequentiality) [see paper]

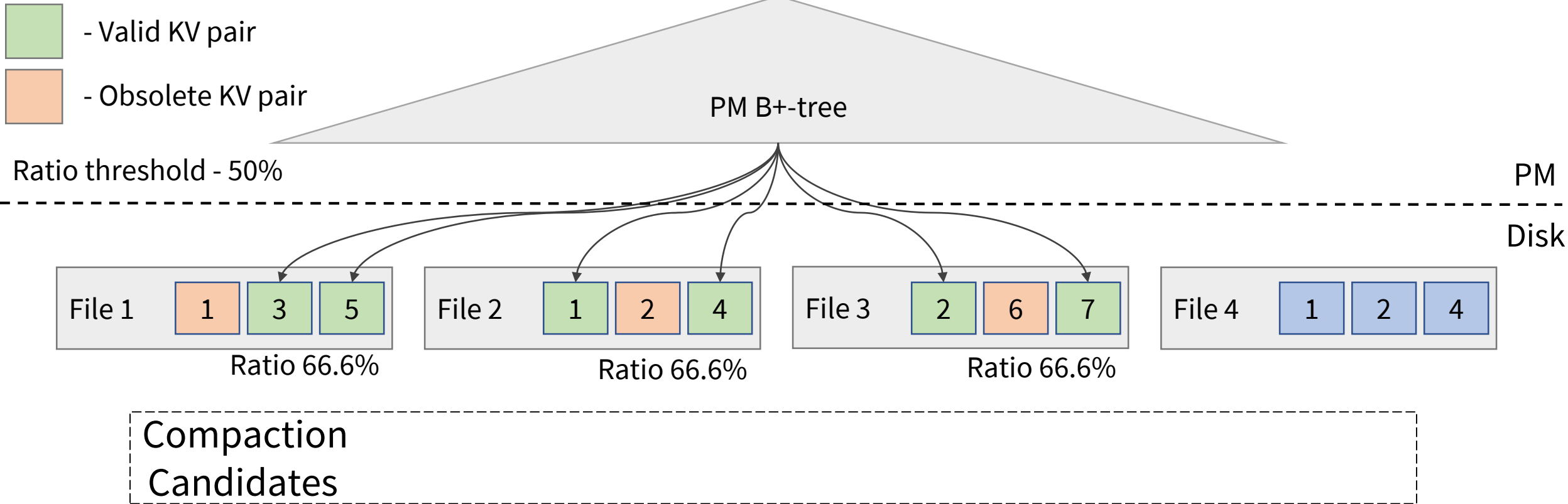
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



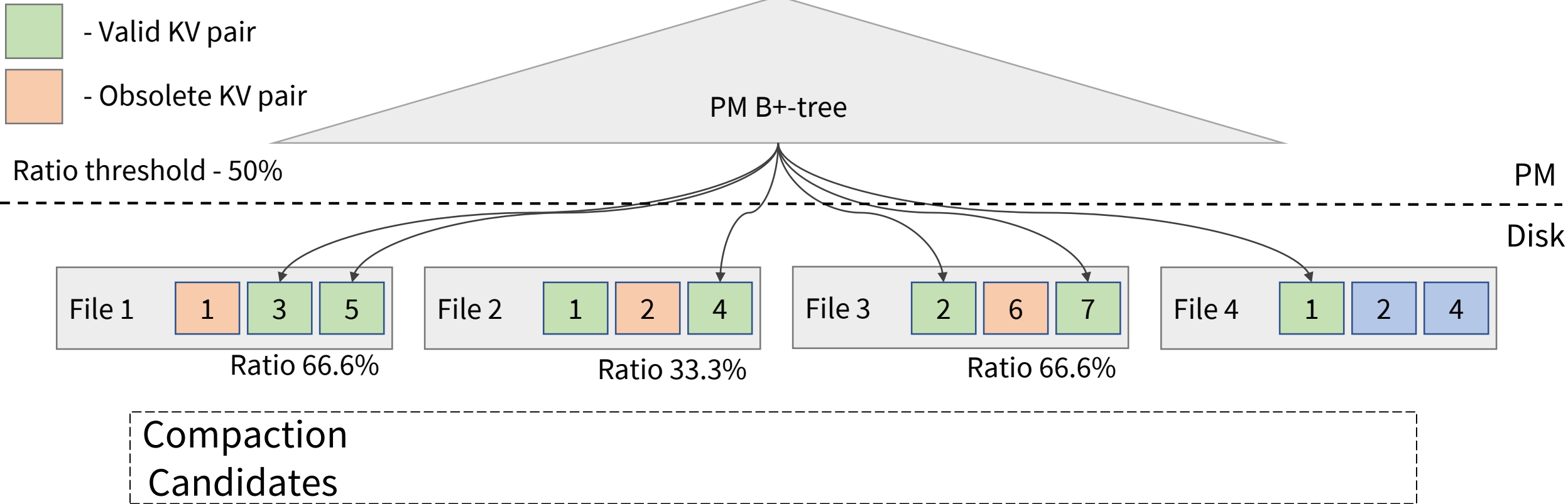
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



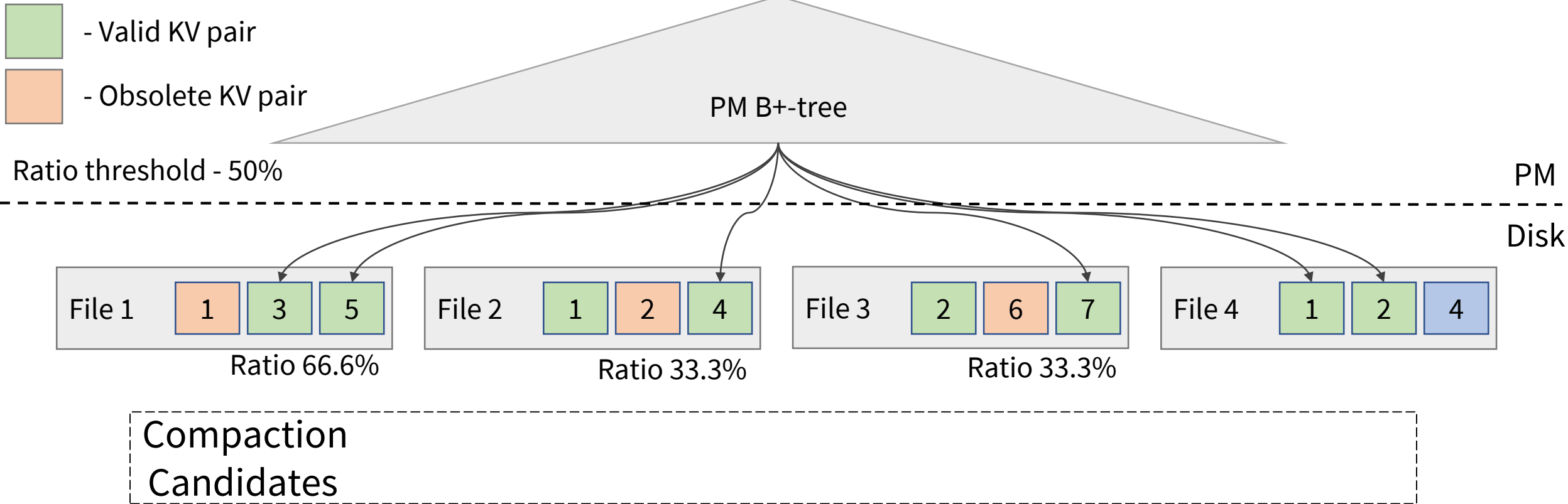
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



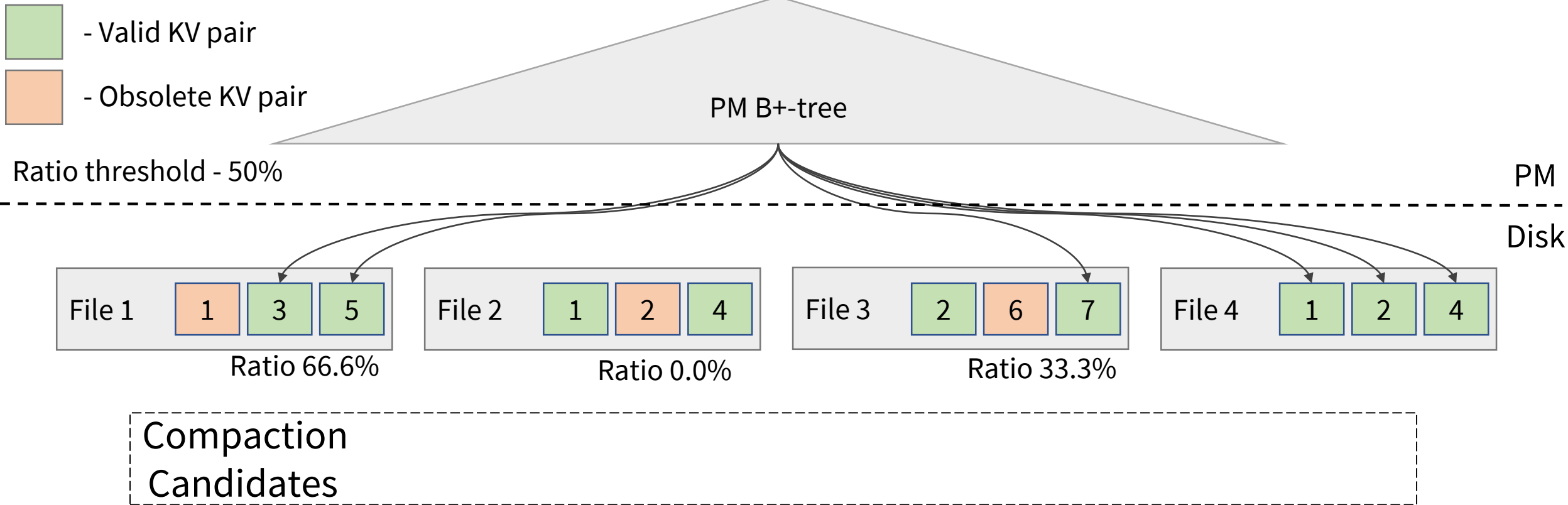
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



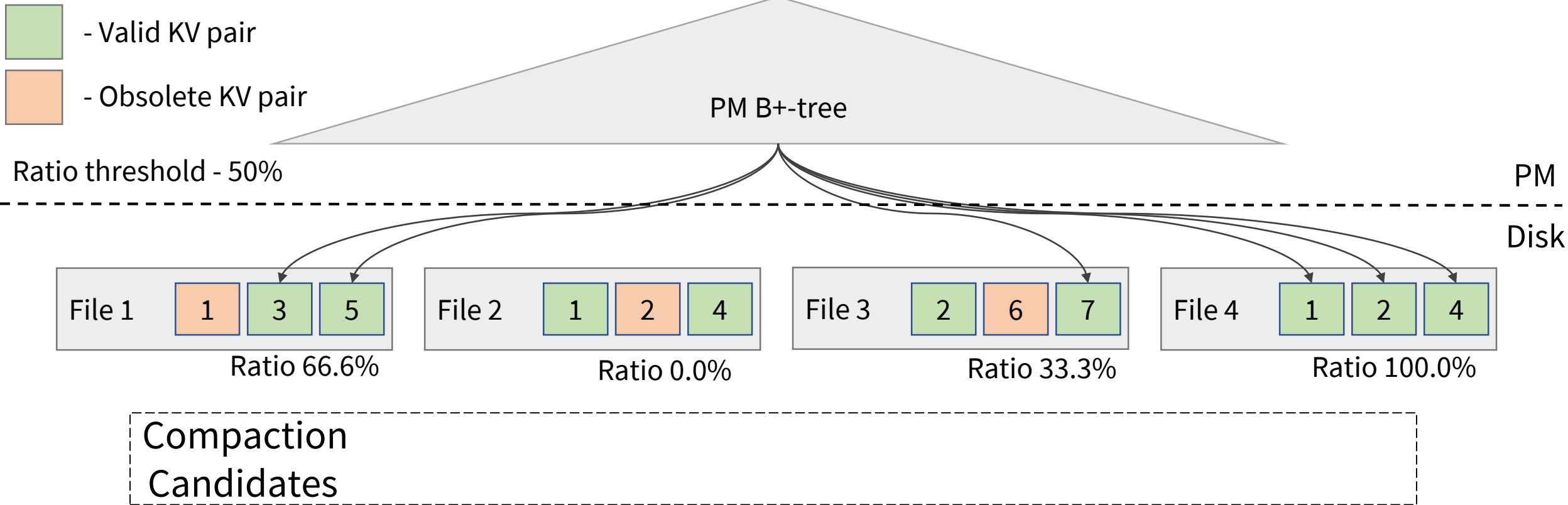
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



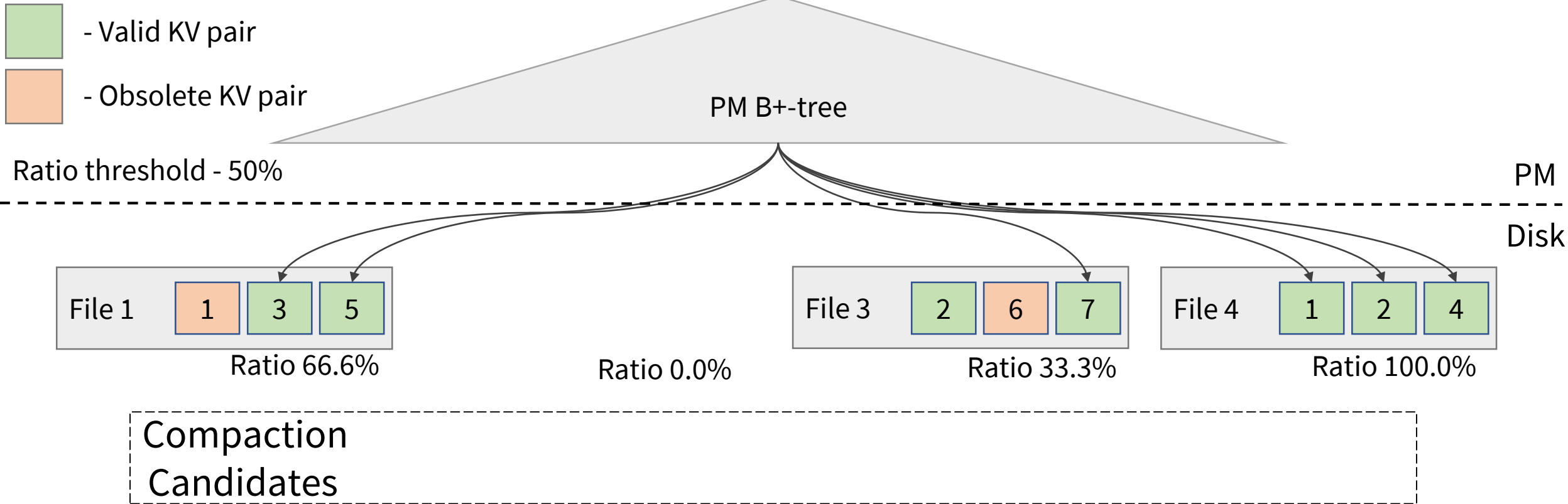
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



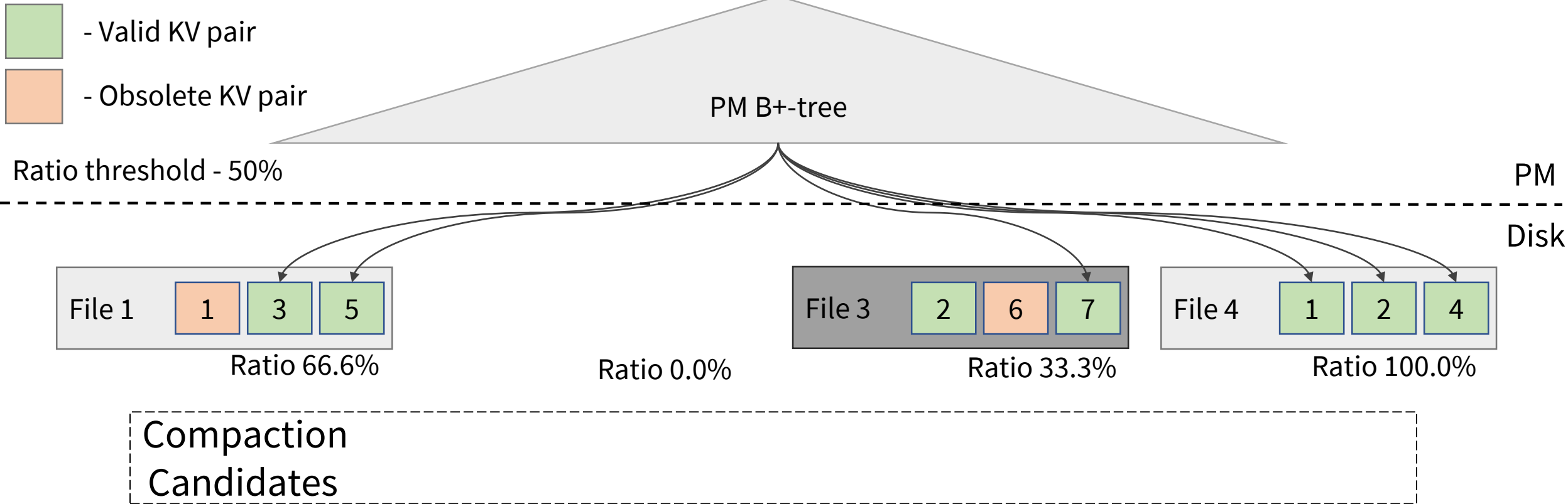
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



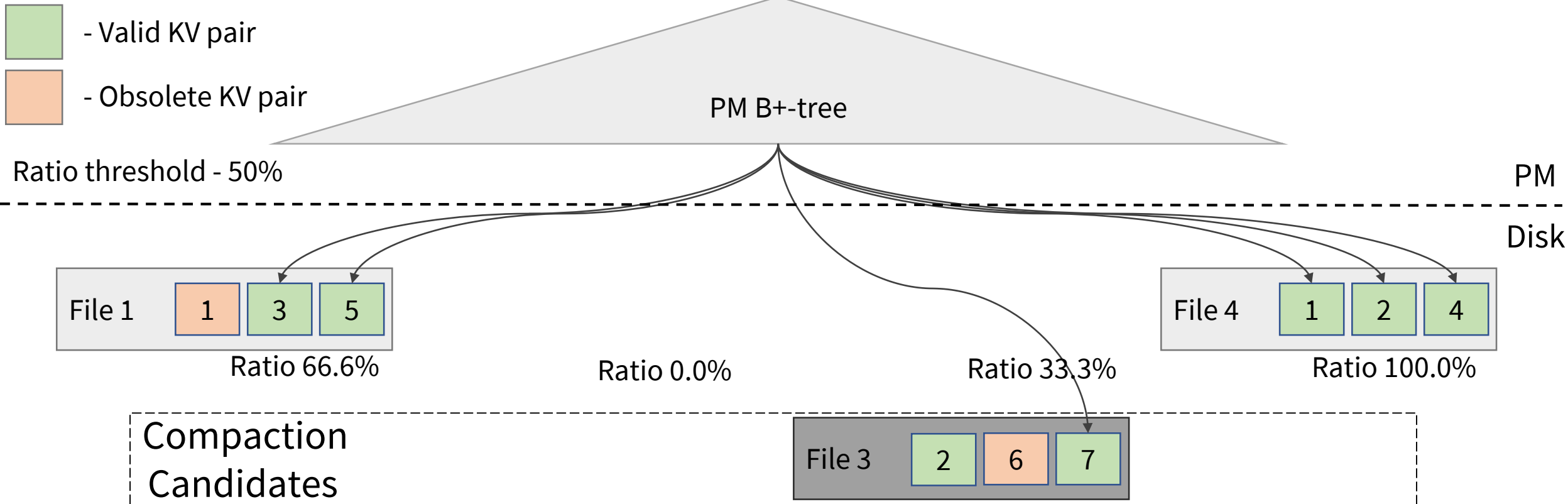
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



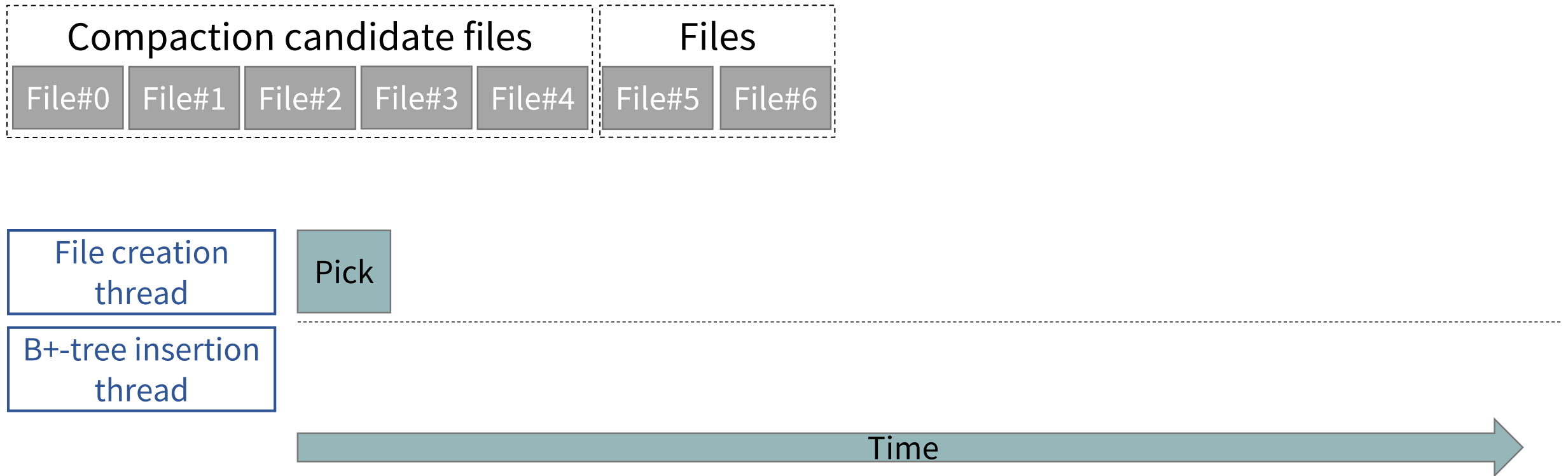
Live-key ratio selection

- To collect garbage
- If live (valid) to total key ratio is below threshold, then add to candidates



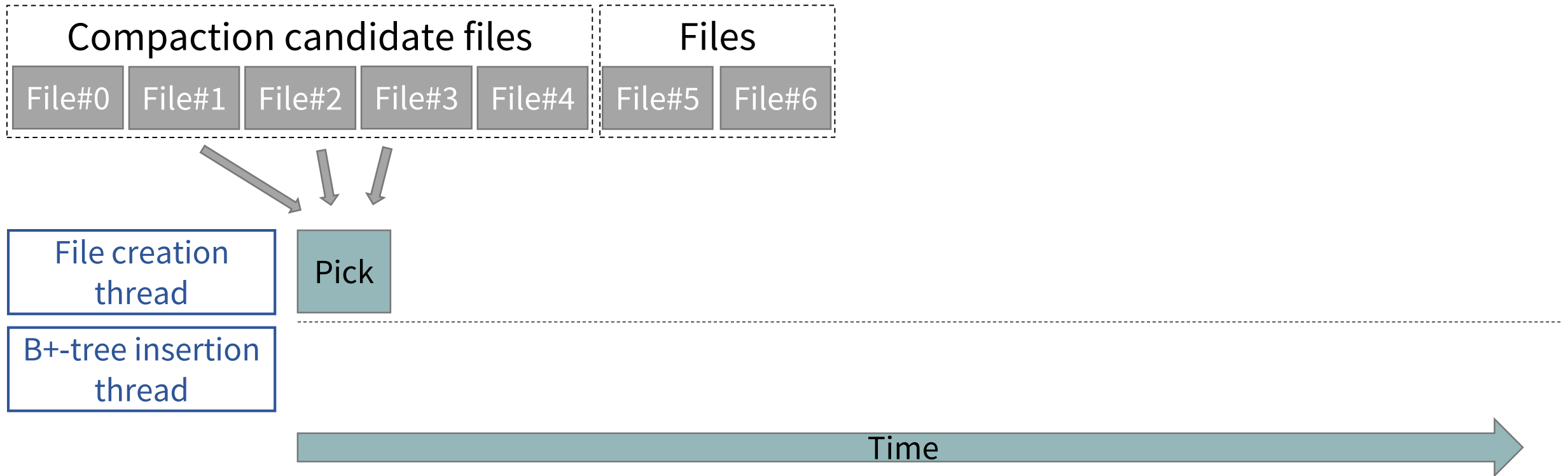
Compaction

- Compaction triggered when there are too many compaction candidate files



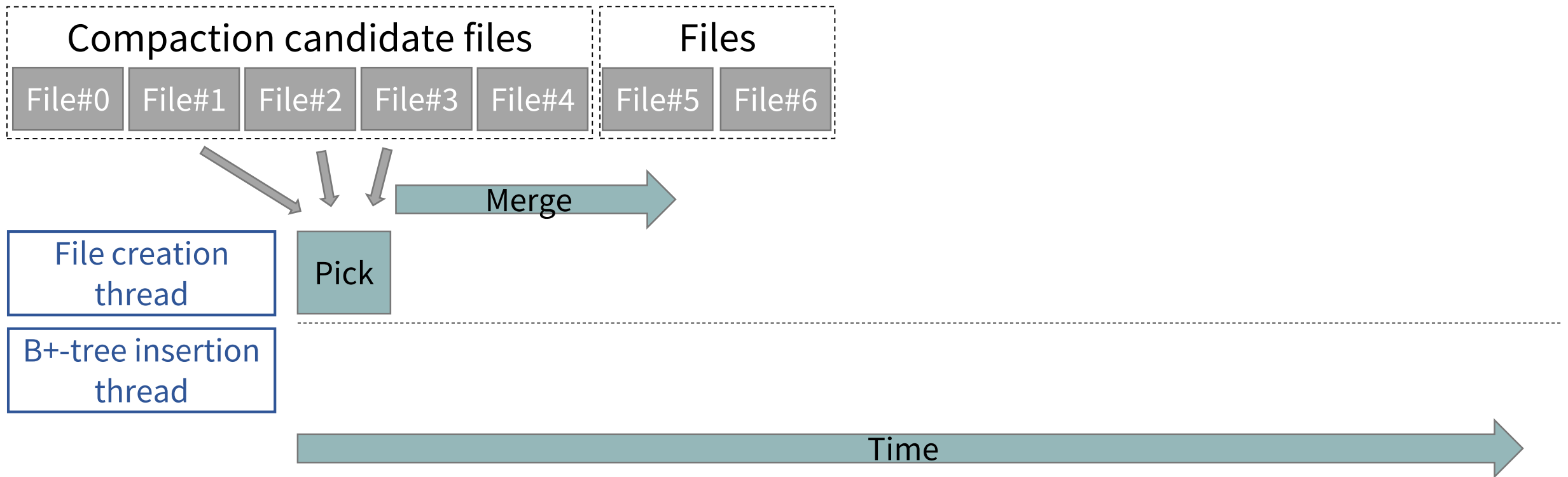
Compaction

- Compaction triggered when there are too many compaction candidate files



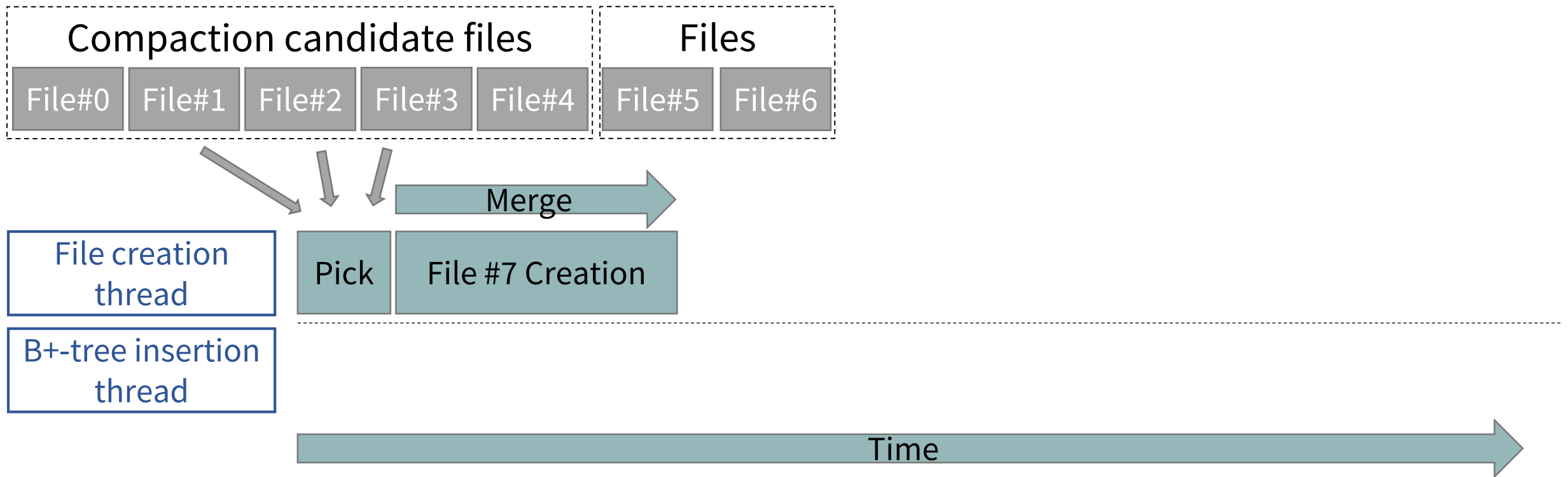
Compaction

- Compaction triggered when there are too many compaction candidate files



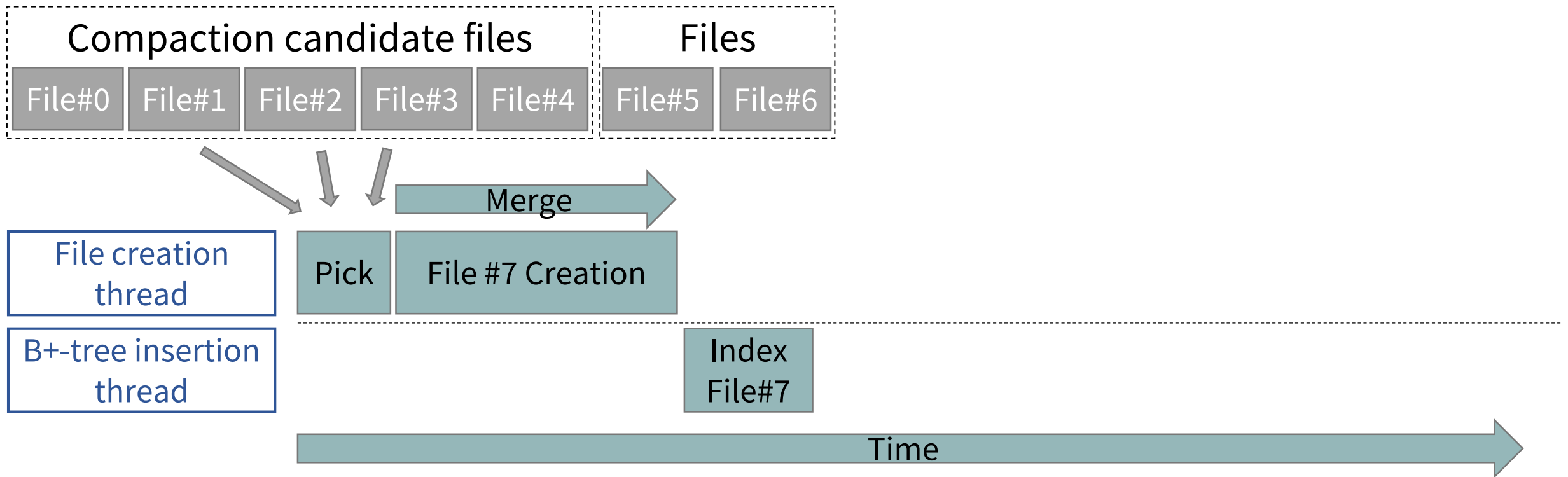
Compaction

- Compaction triggered when there are too many compaction candidate files



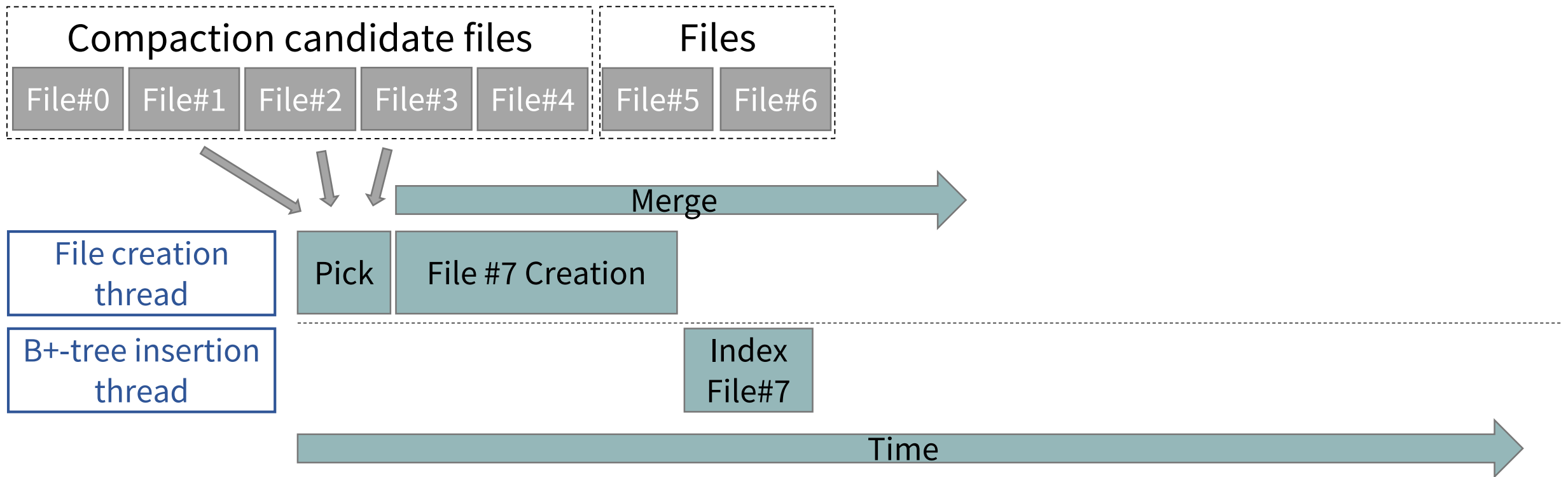
Compaction

- Compaction triggered when there are too many compaction candidate files



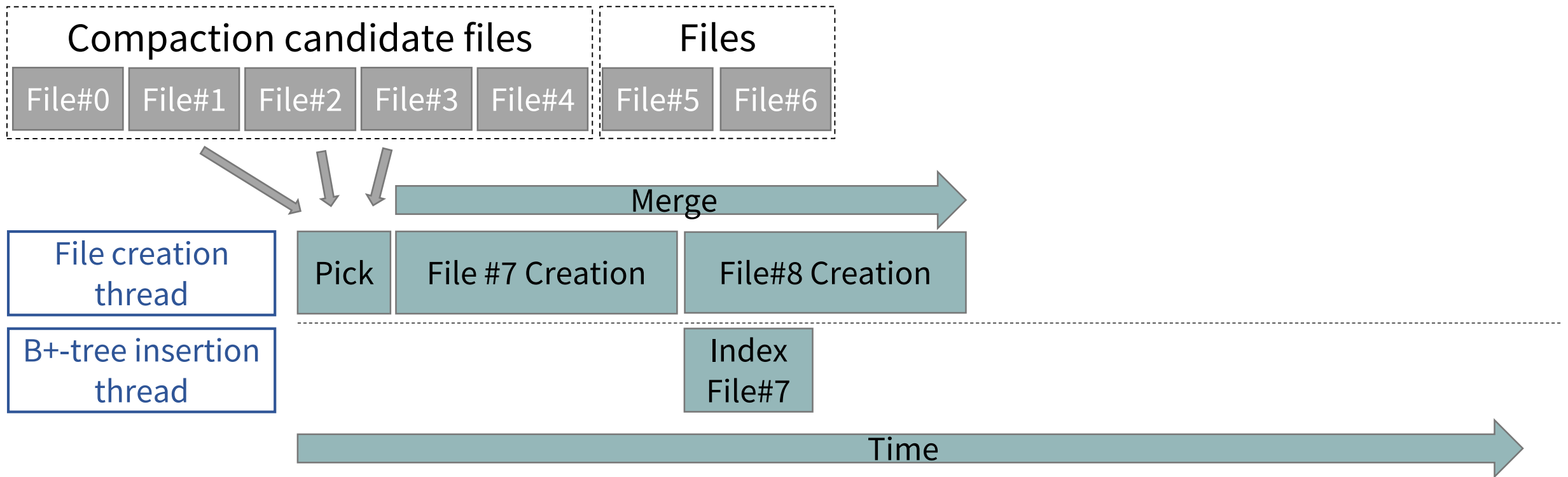
Compaction

- Compaction triggered when there are too many compaction candidate files



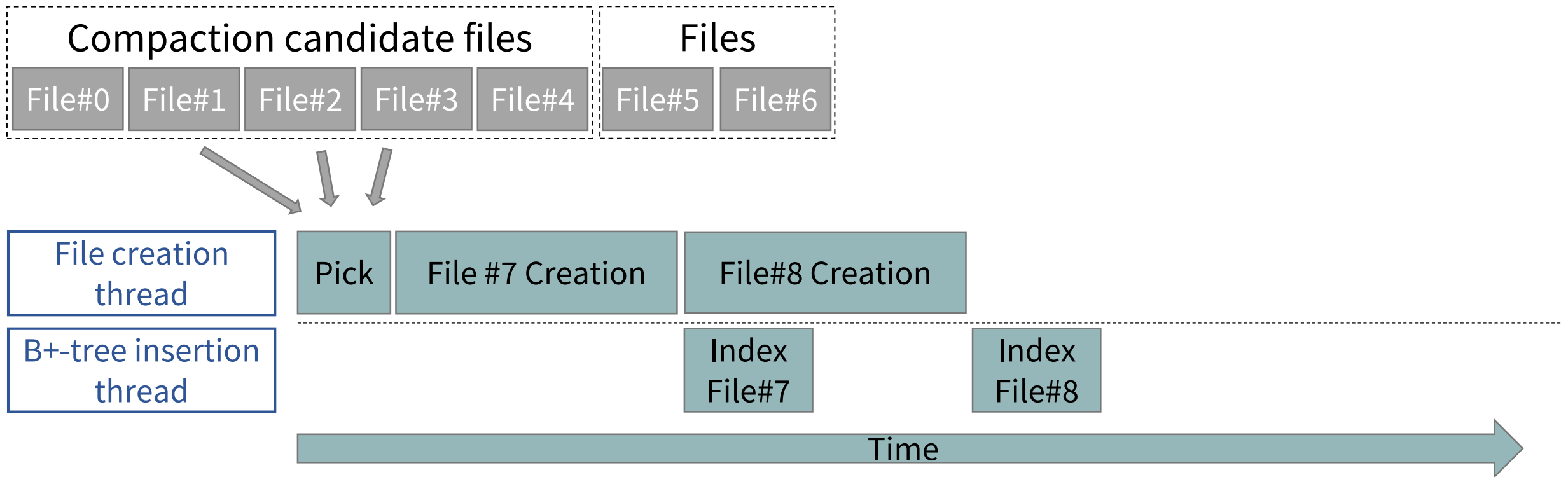
Compaction

- Compaction triggered when there are too many compaction candidate files



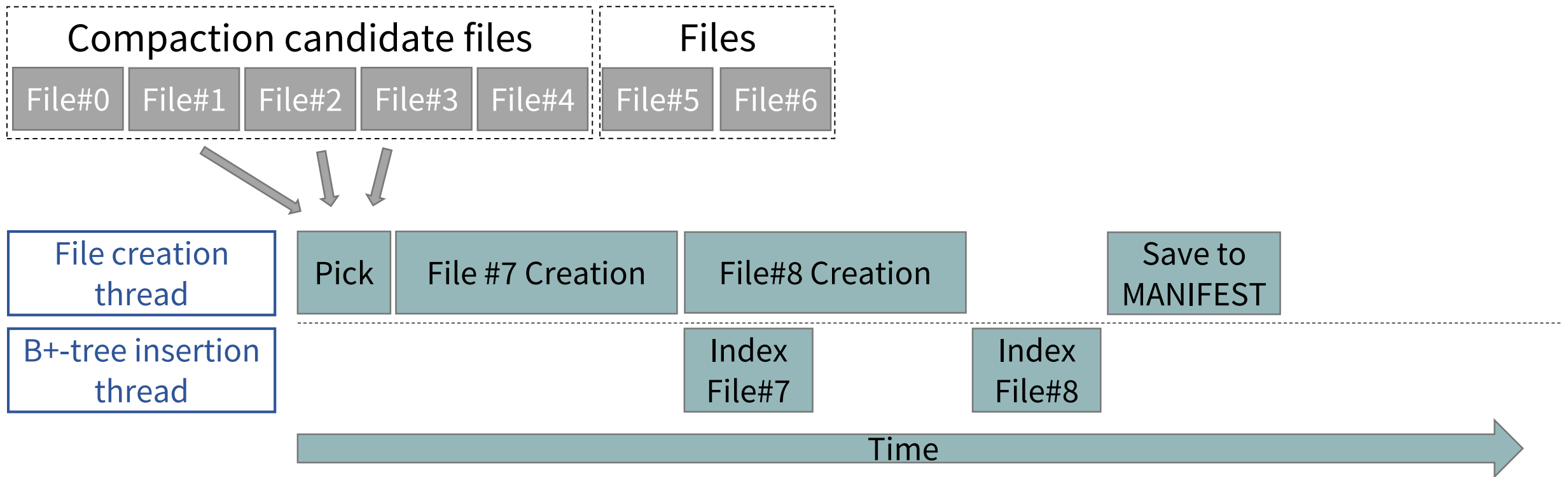
Compaction

- Compaction triggered when there are too many compaction candidate files



Compaction

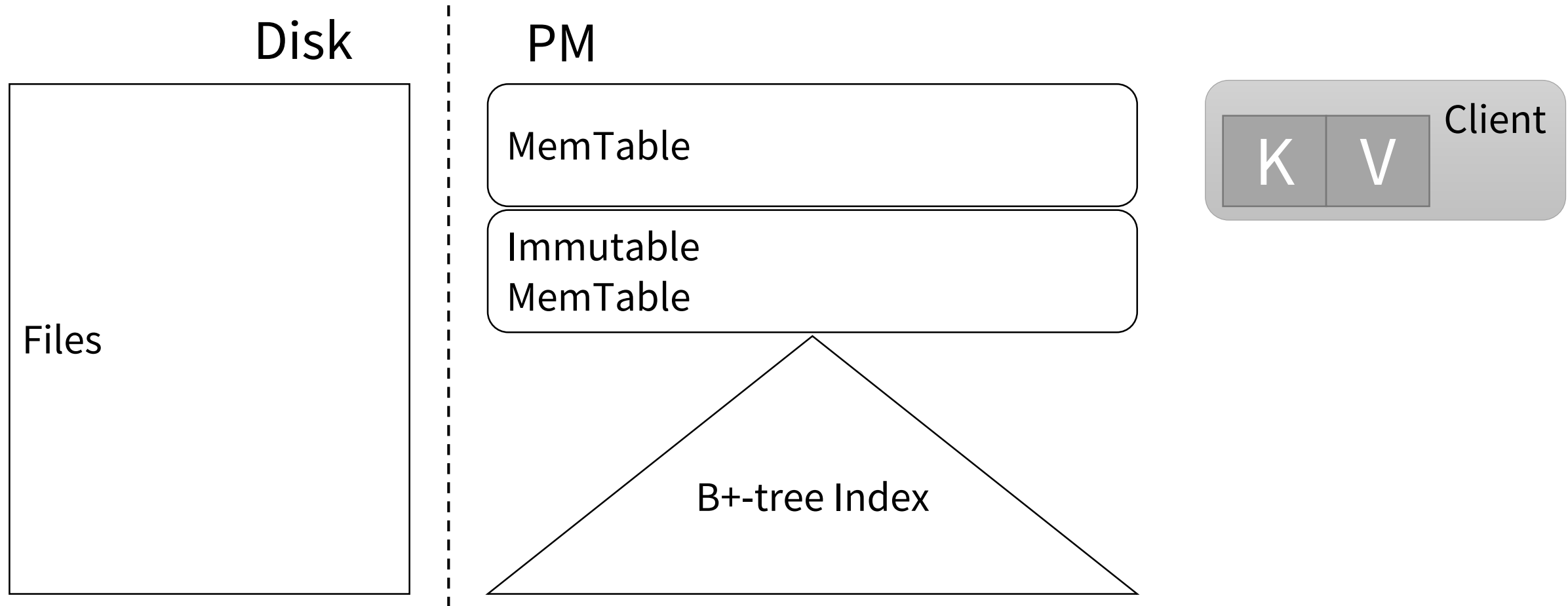
- Compaction triggered when there are too many compaction candidate files



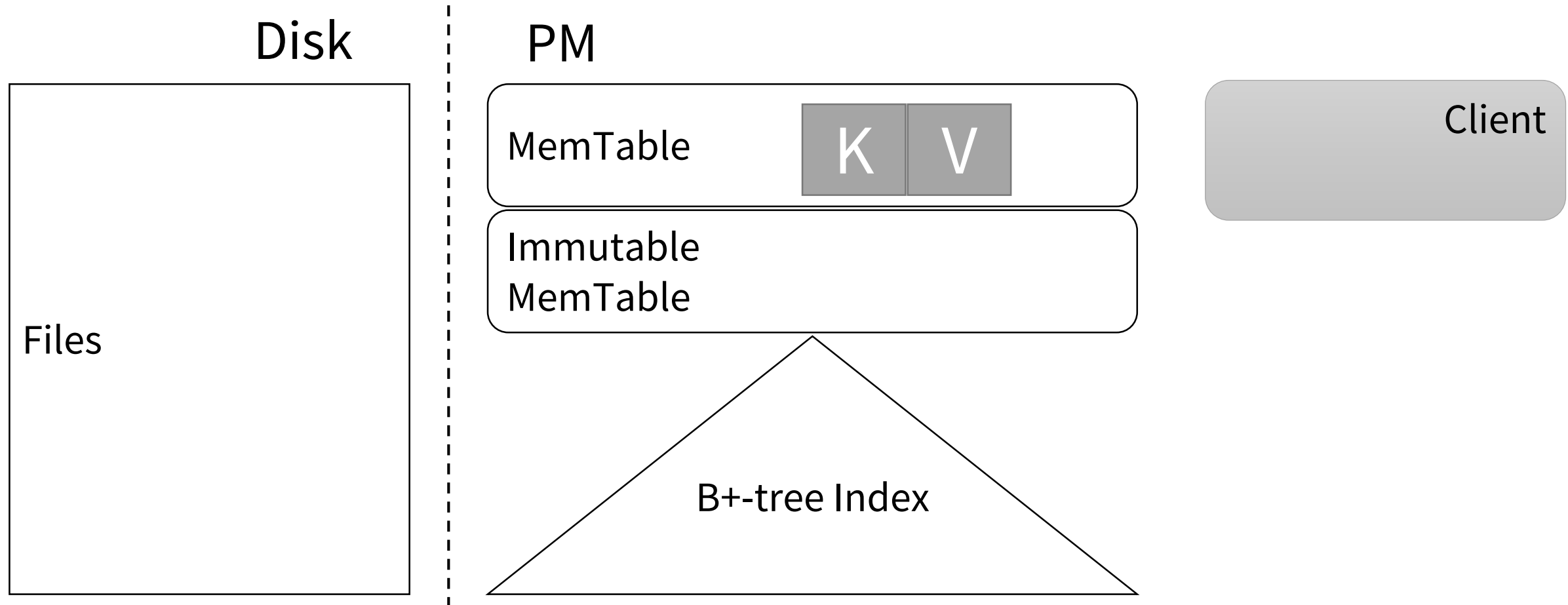
General operations

- Put
- Put if exists/Put if not exists
- Get
- Scan

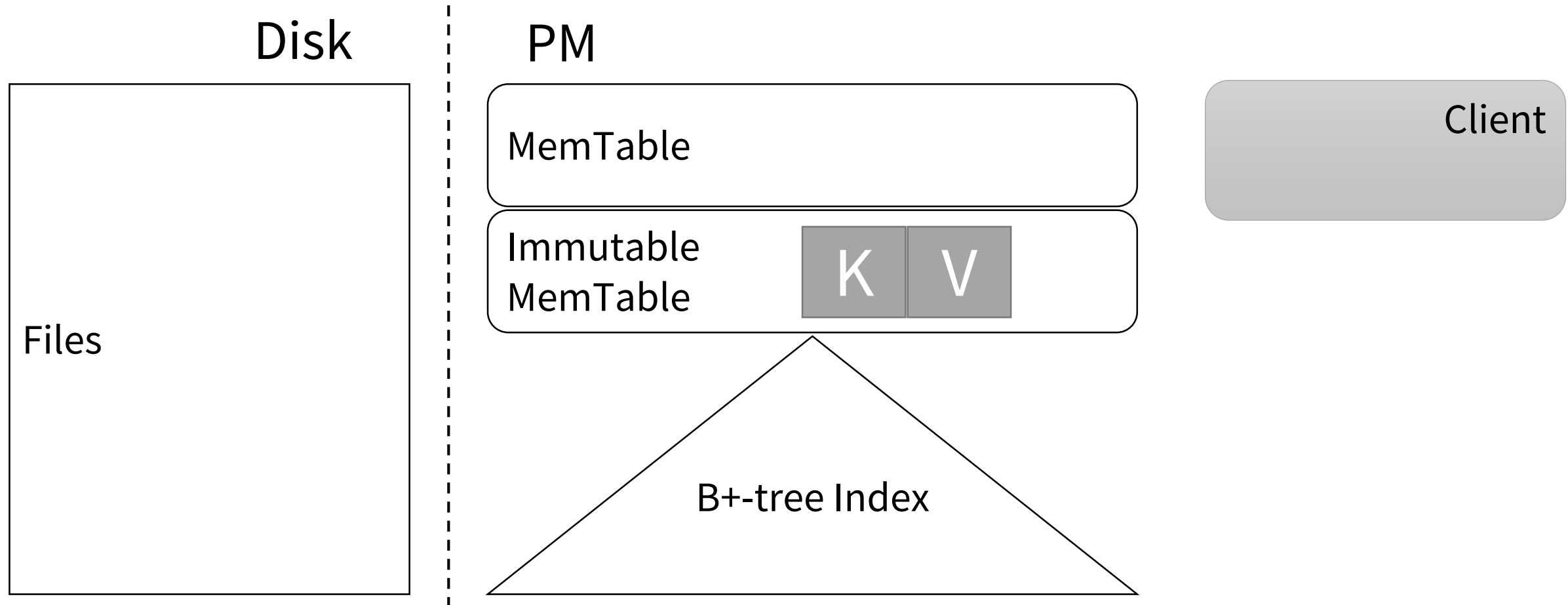
Put(key, value)



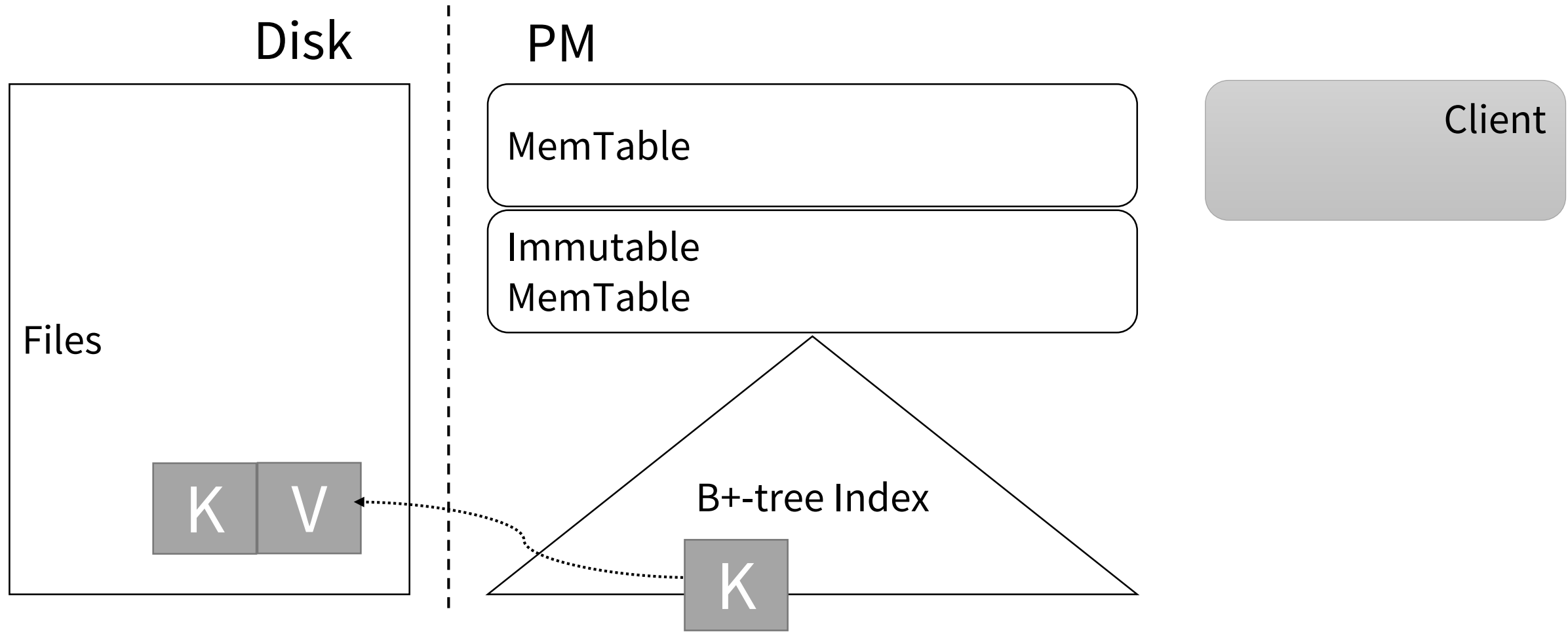
Put(key, value)



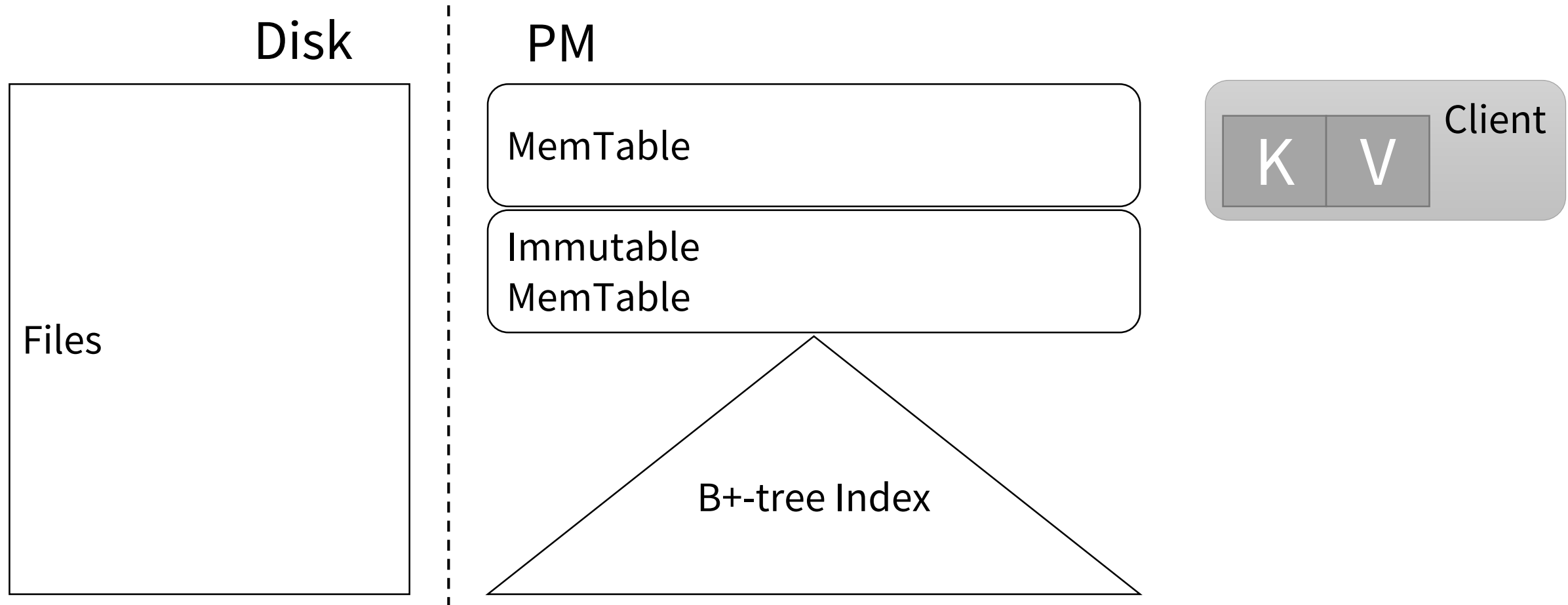
Put(key, value)



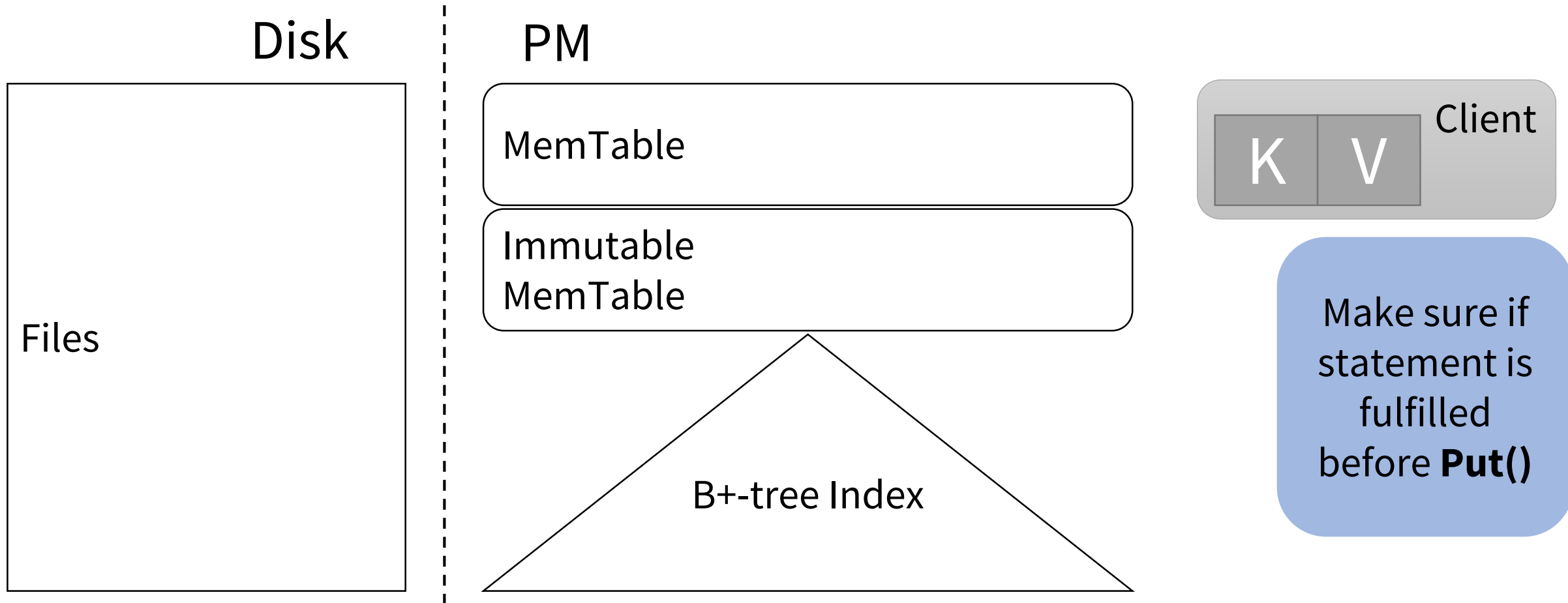
Put(key, value)



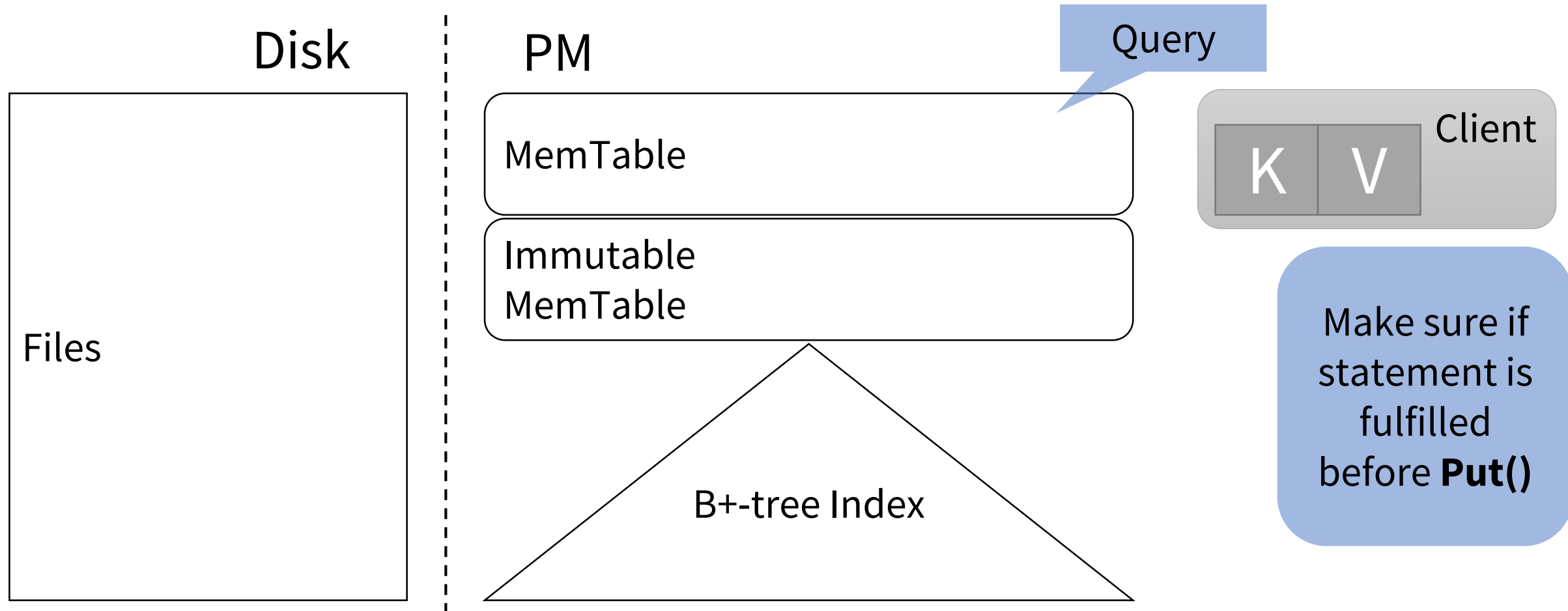
Put(key, value) if exists/if not exists



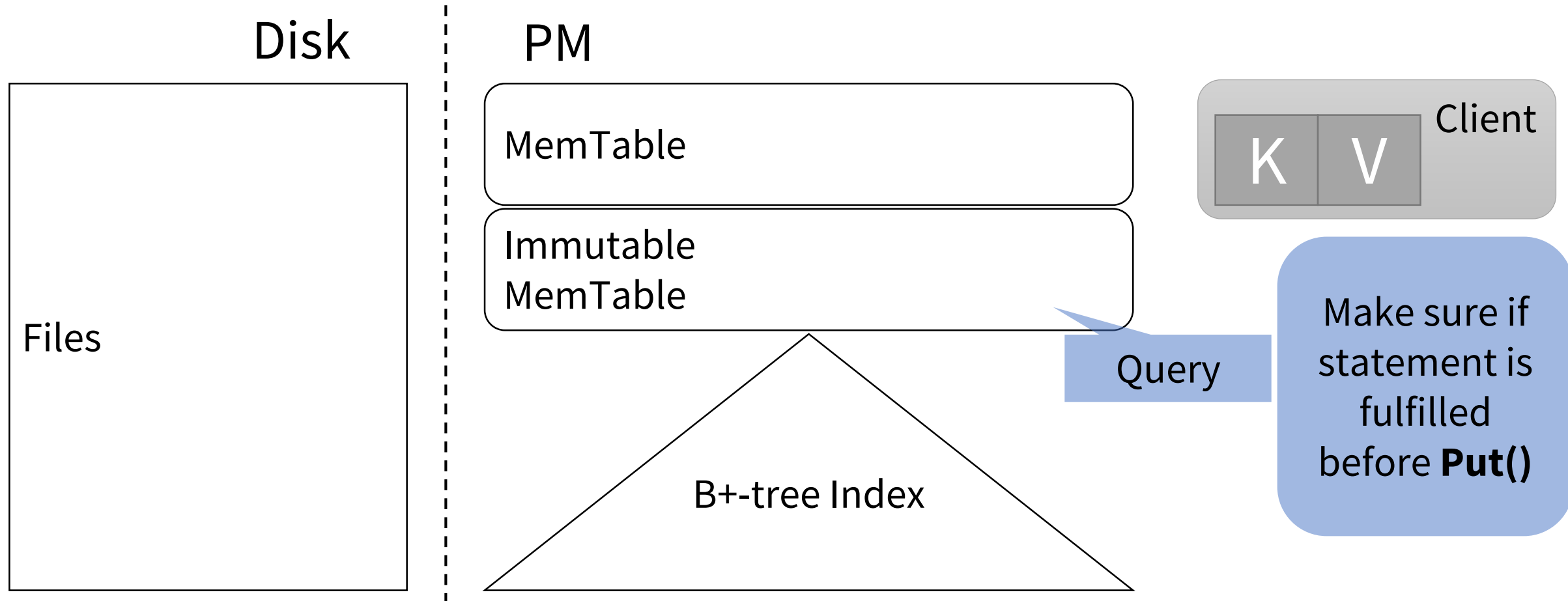
Put(key, value) if exists/if not exists



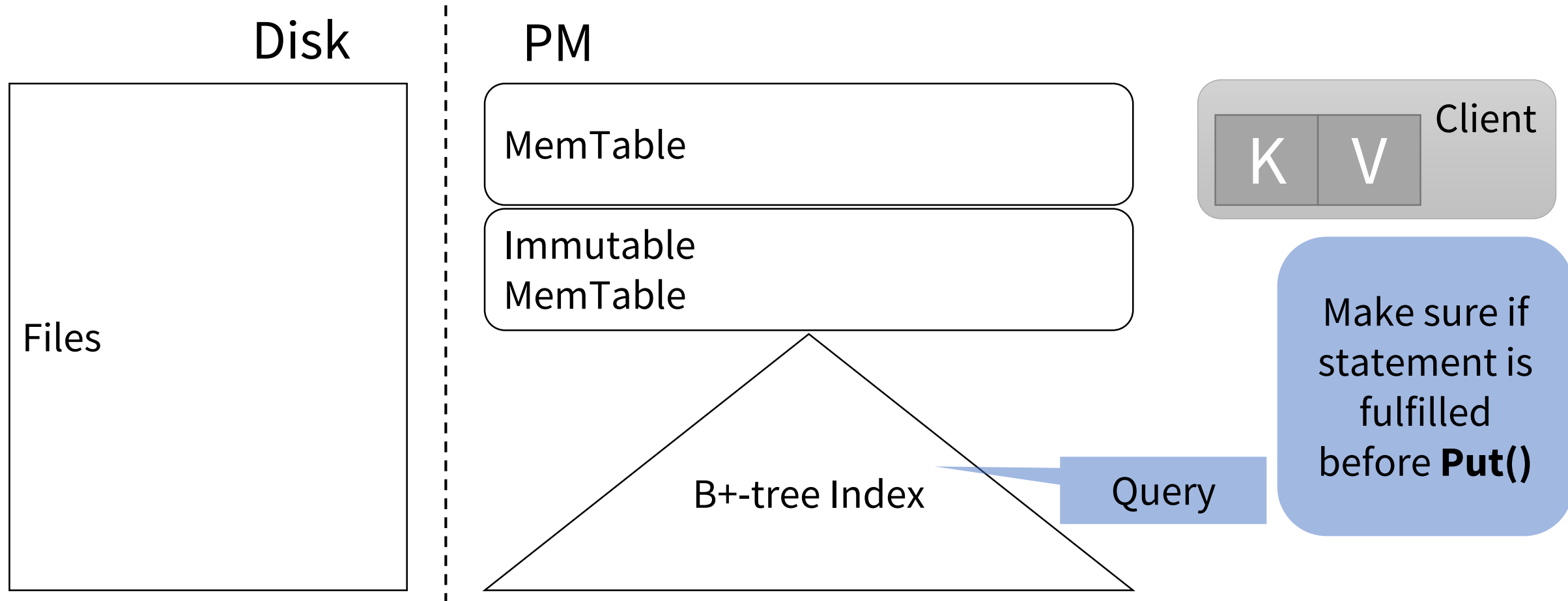
Put(key, value) if exists/if not exists



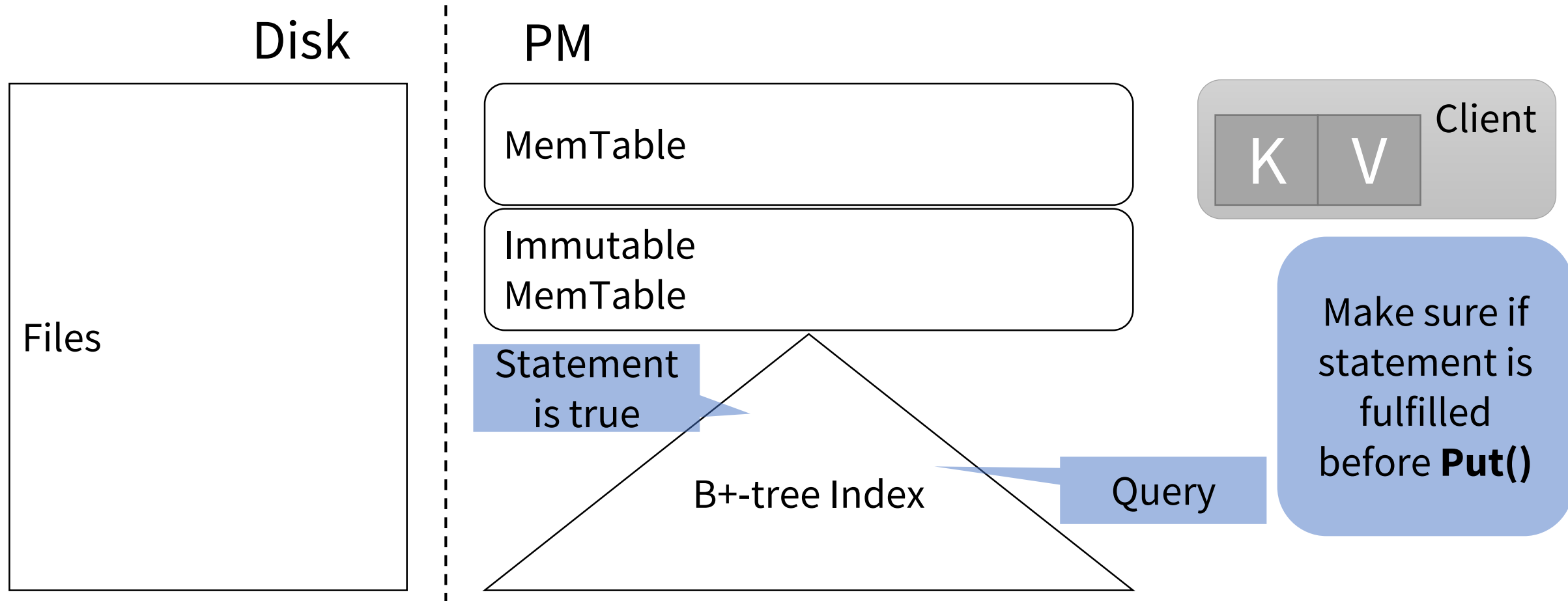
Put(key, value) if exists/if not exists



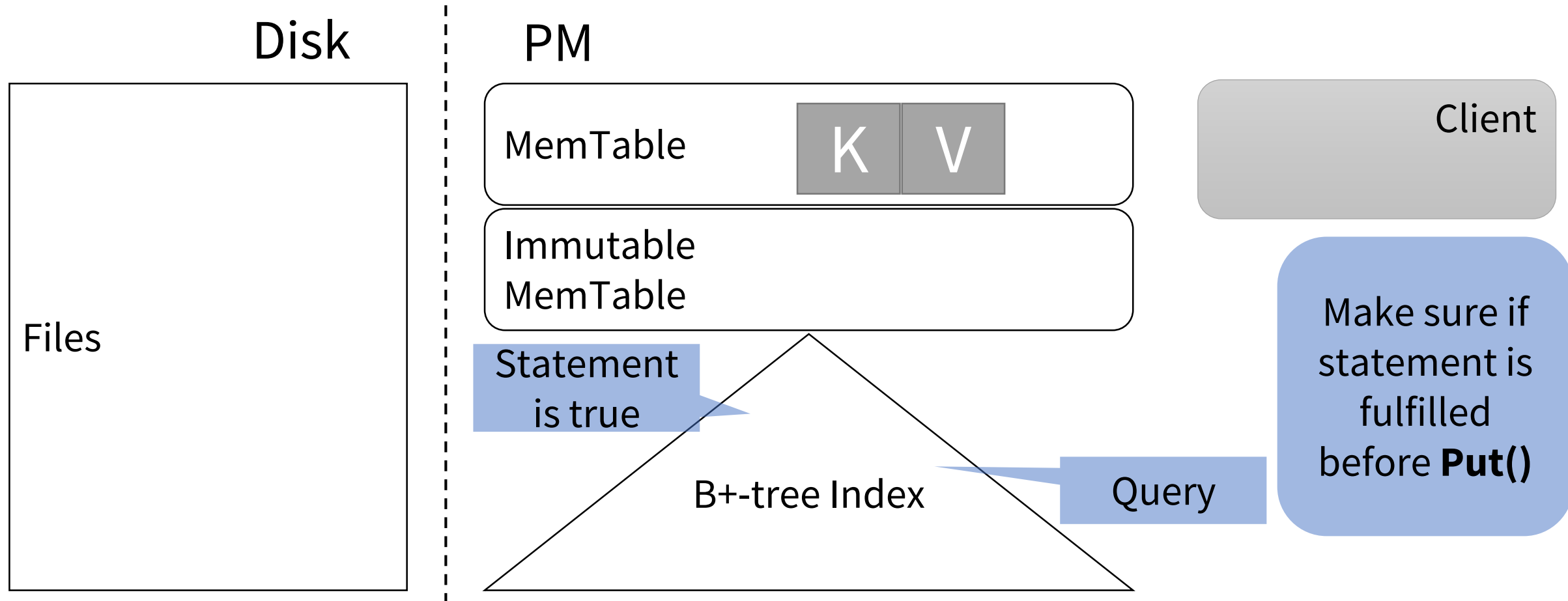
Put(key, value) if exists/if not exists



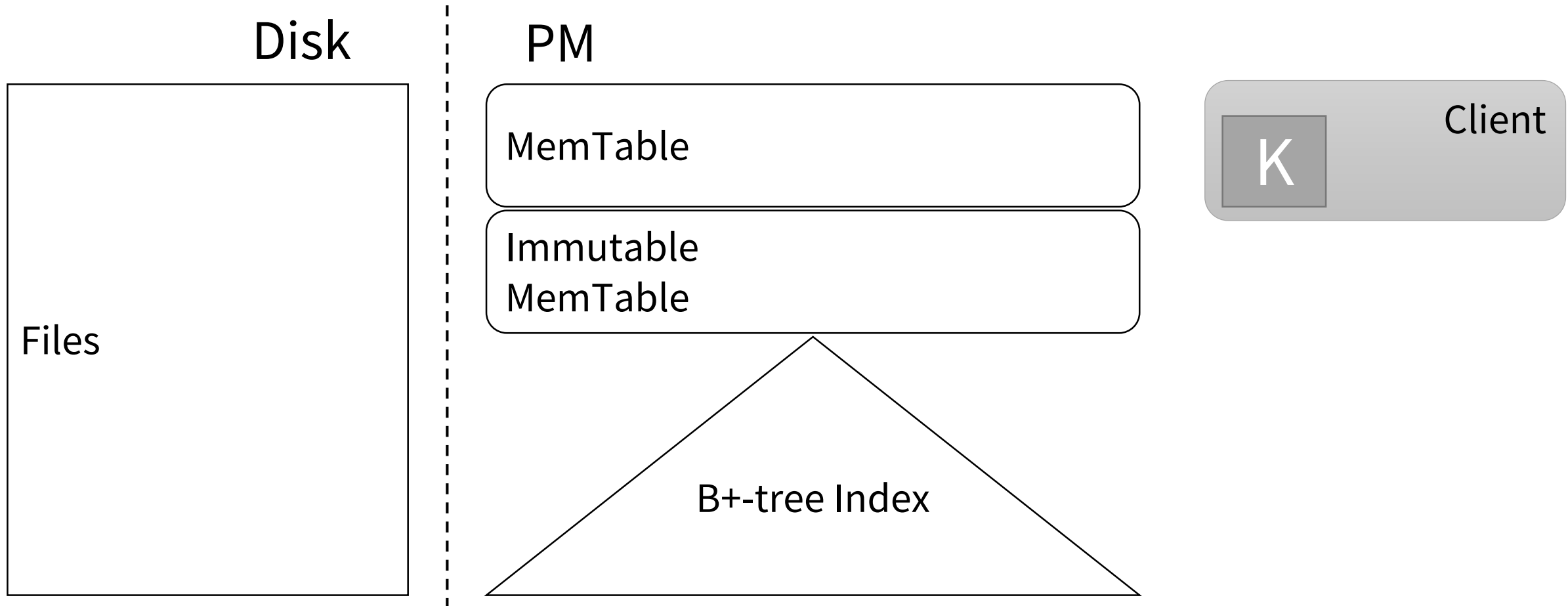
Put(key, value) if exists/if not exists



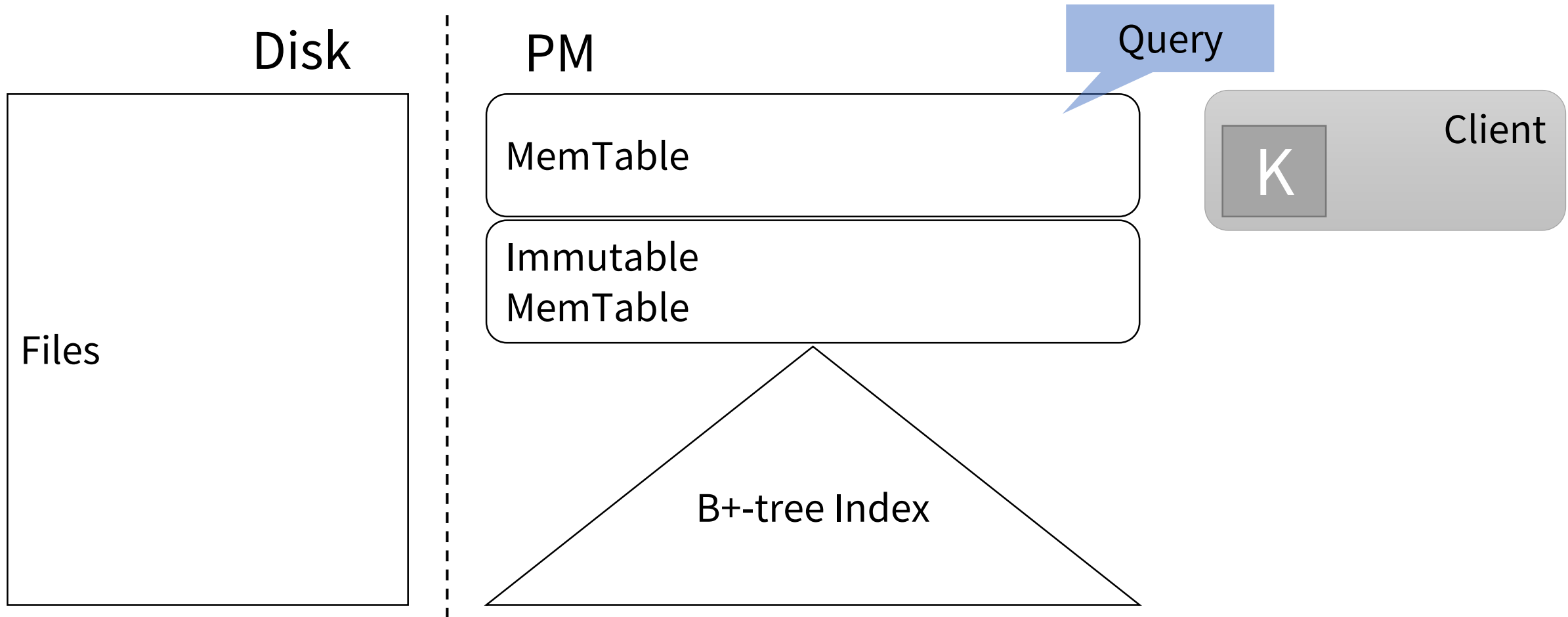
Put(key, value) if exists/if not exists



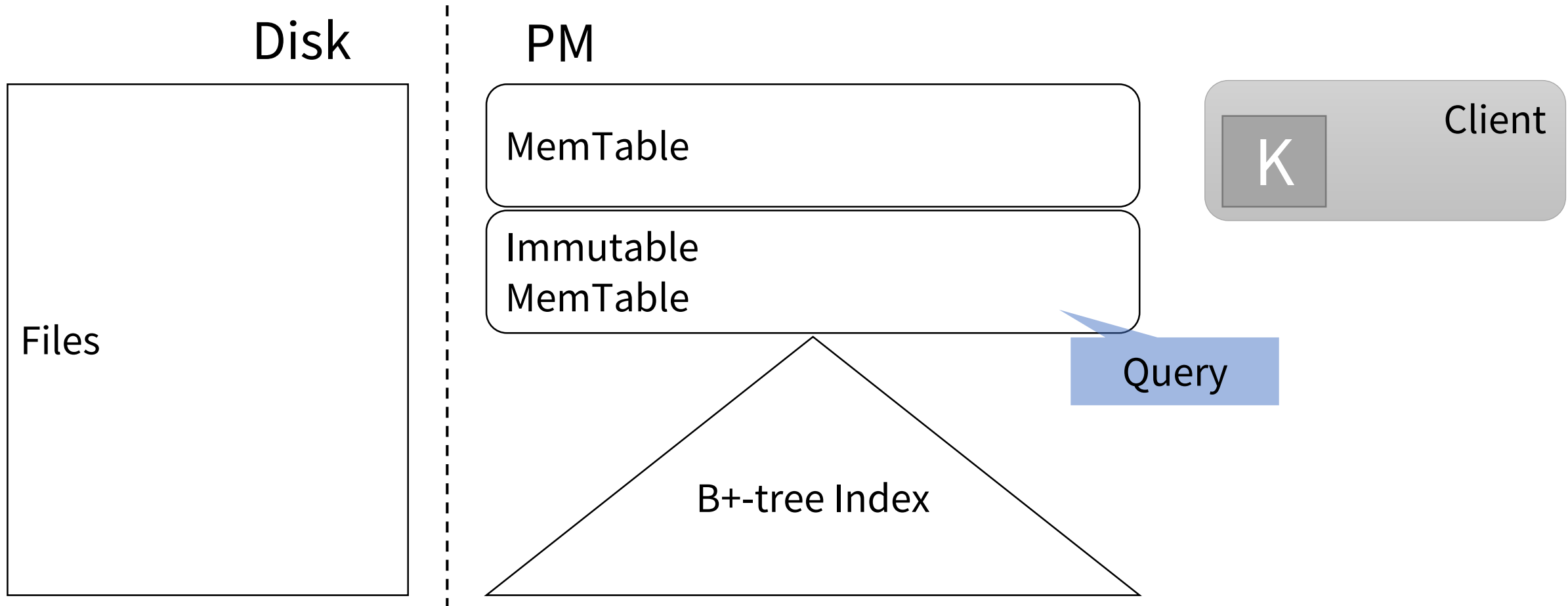
Get(key)



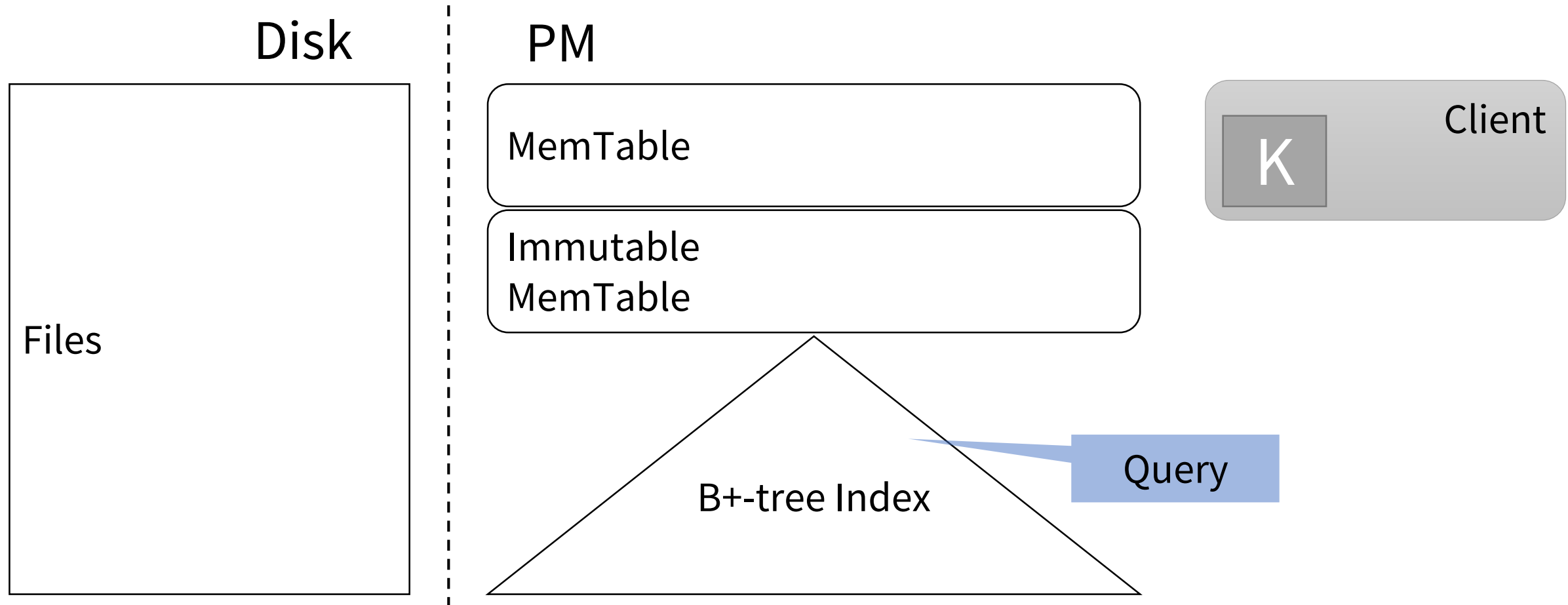
Get(key)



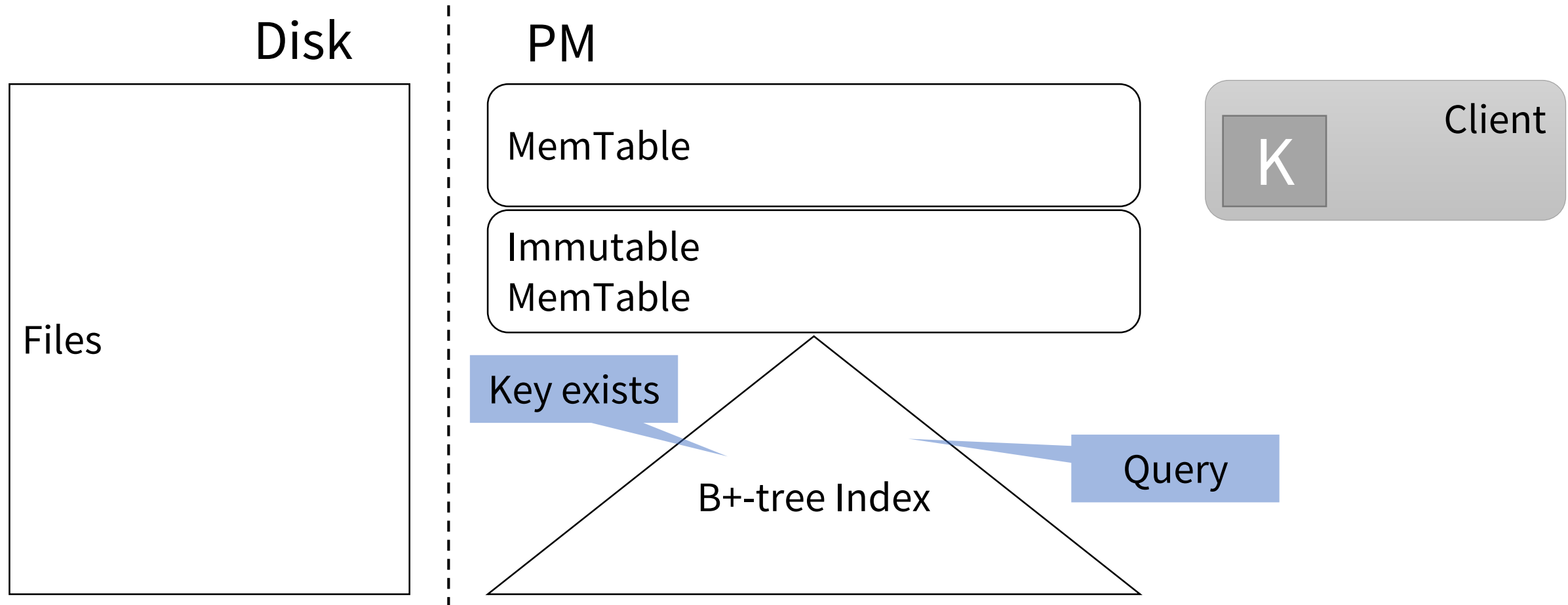
Get(key)



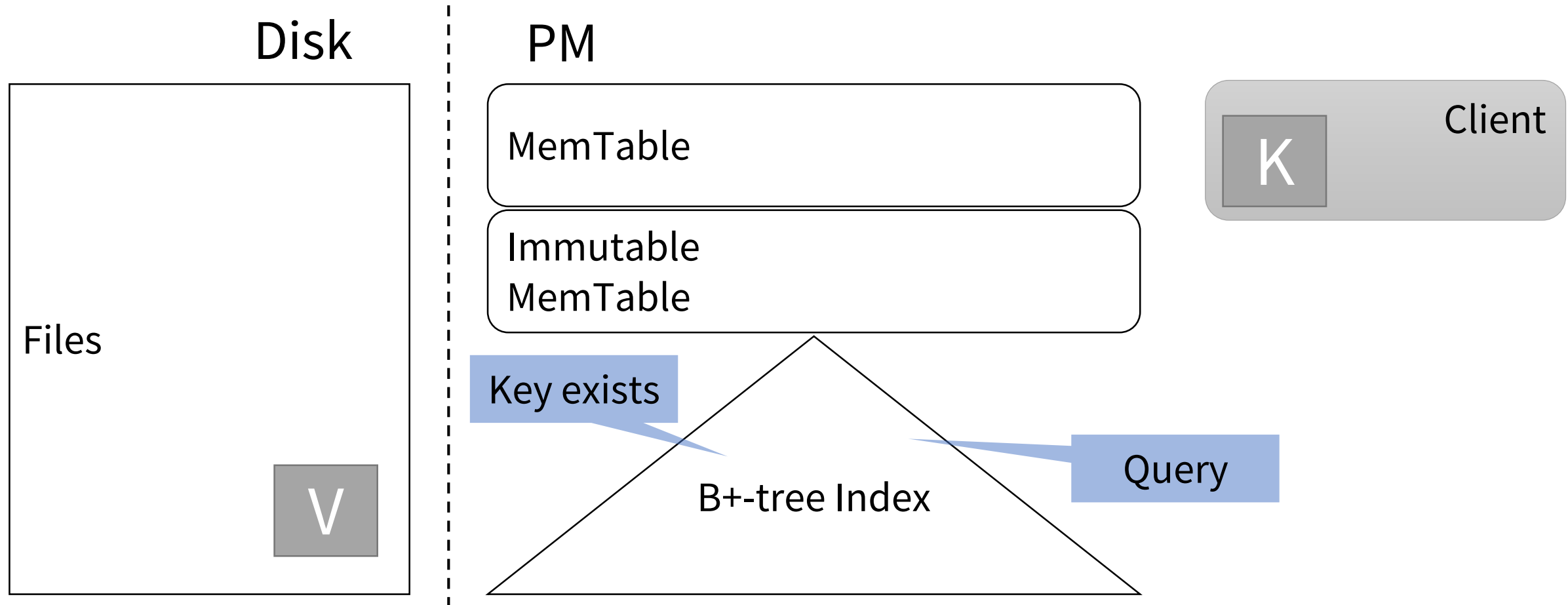
Get(key)



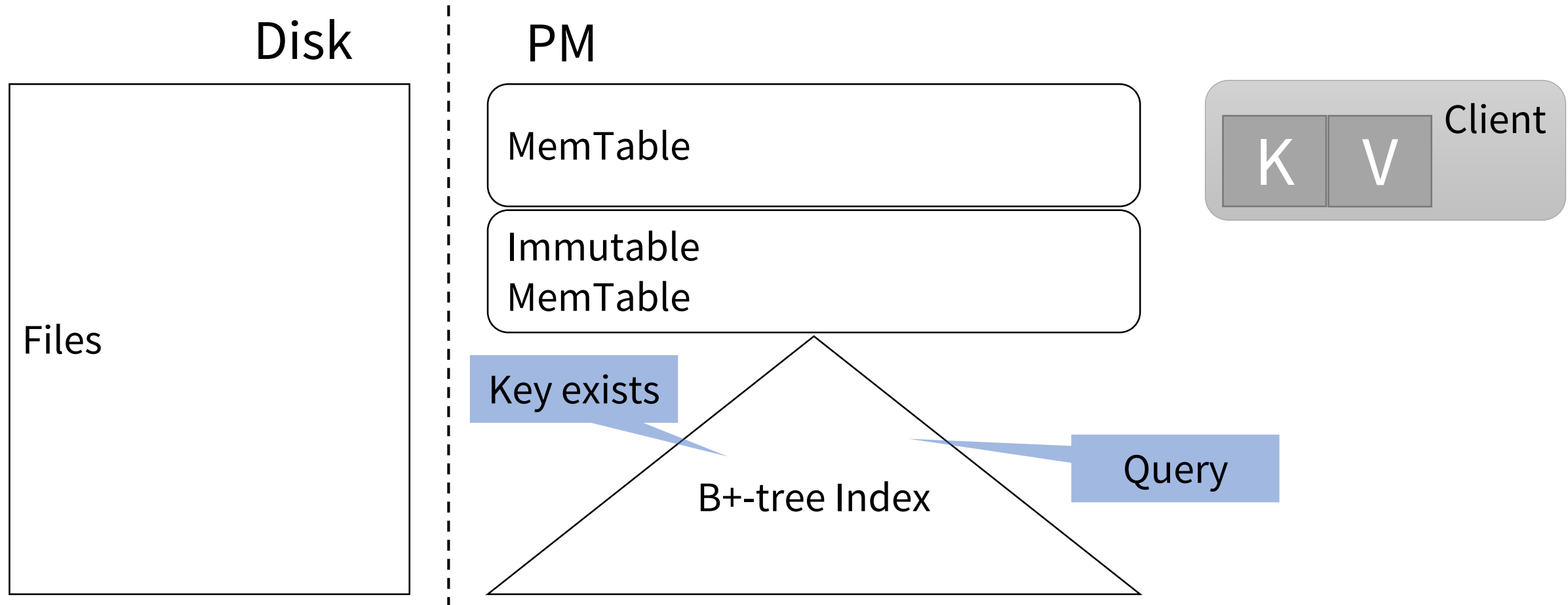
Get(key)



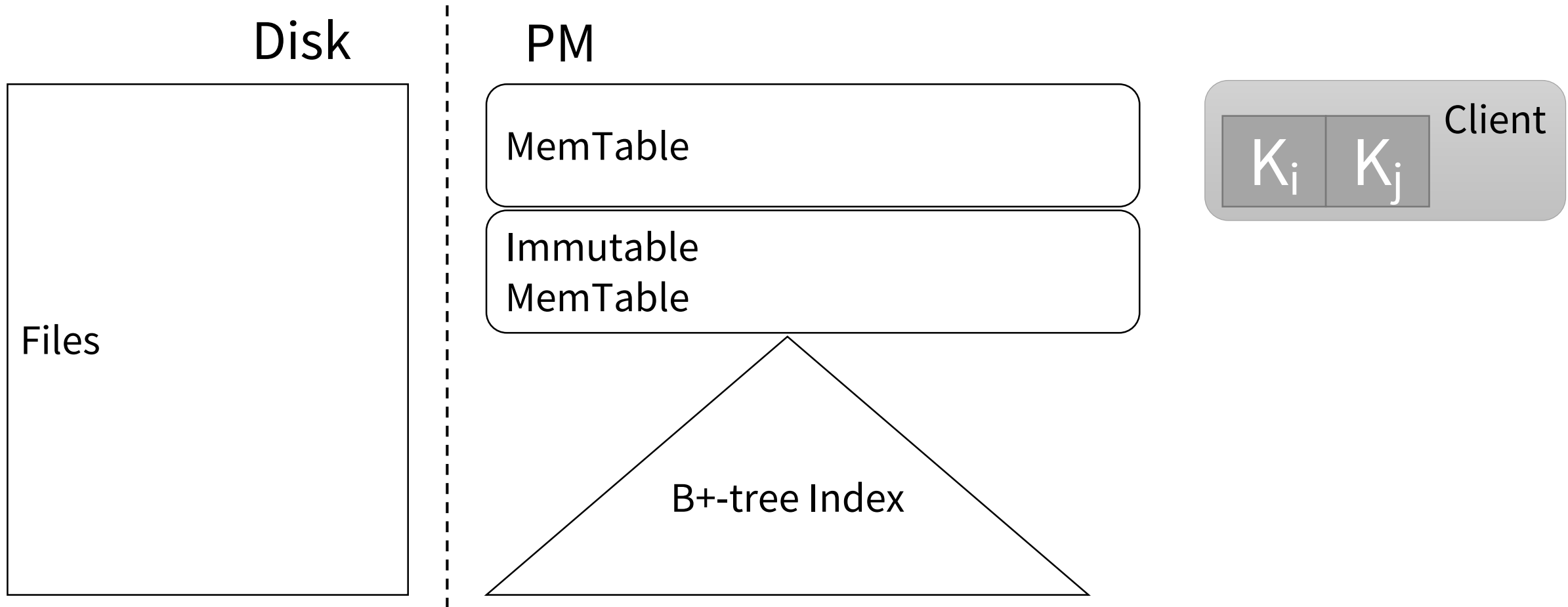
Get(key)



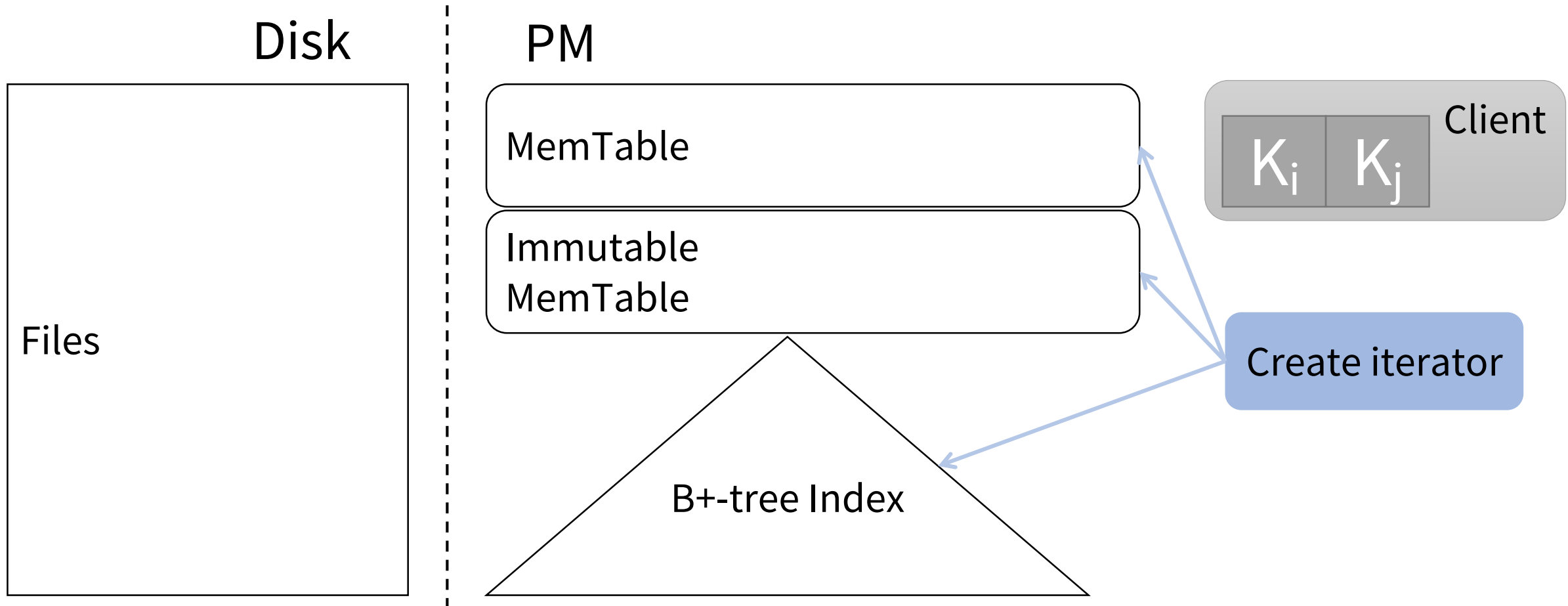
Get(key)



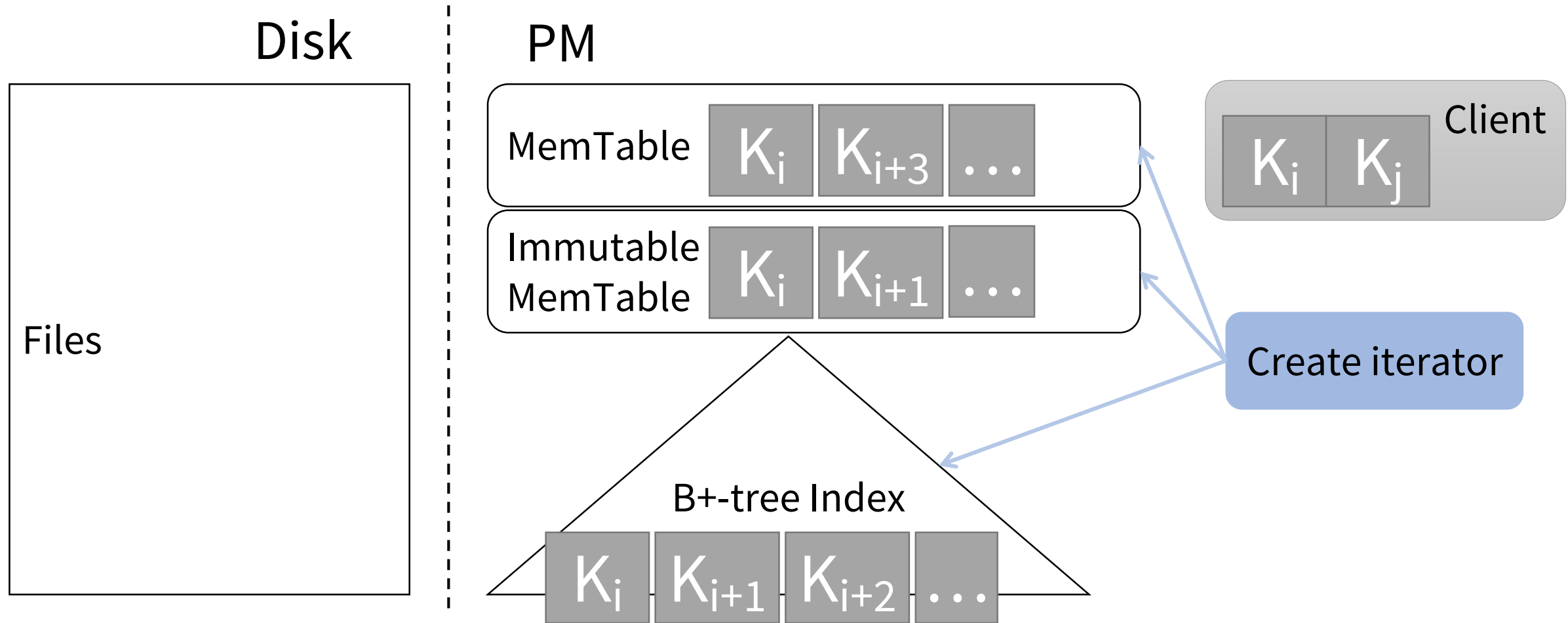
Scan(key_i , key_j)



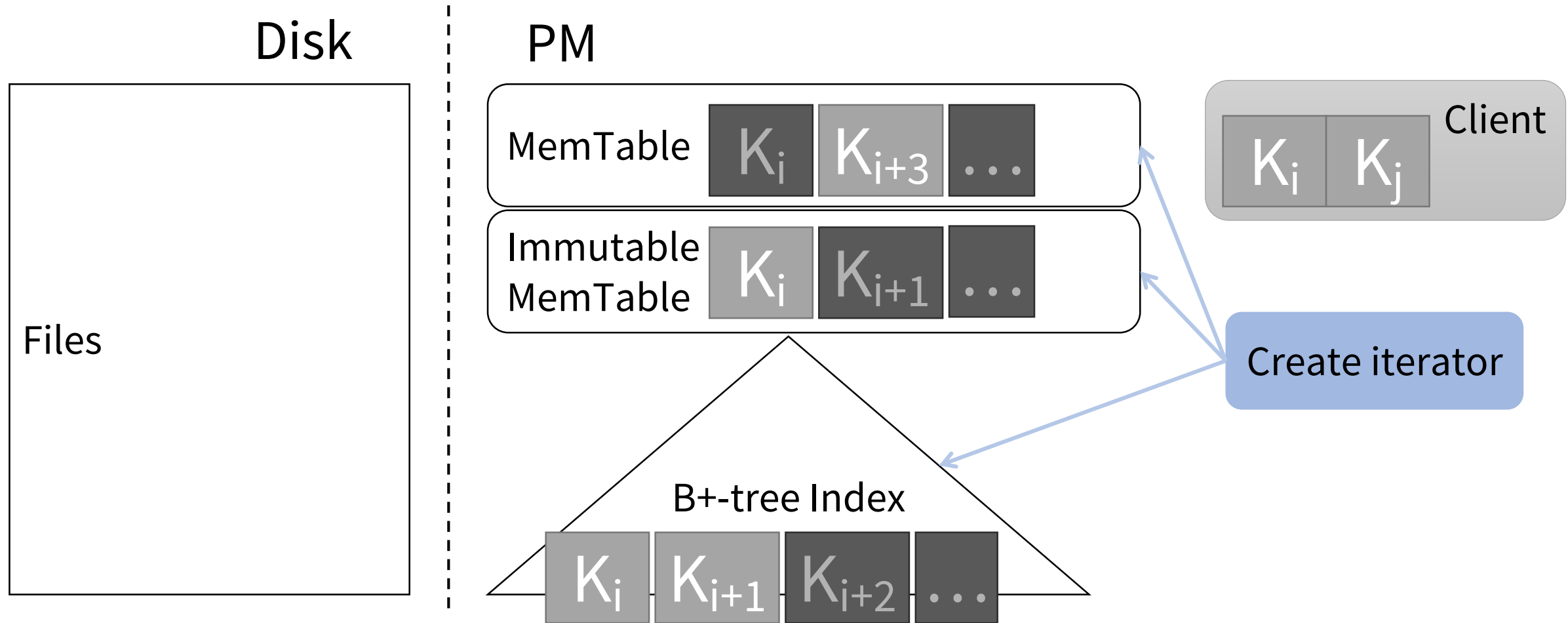
Scan(key_i , key_j)



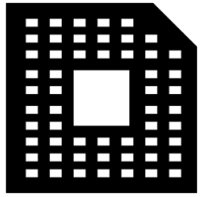
Scan(key_i , key_j)



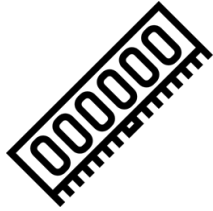
Scan(key_i , key_j)



Evaluation



Intel Xeon
E5-2640 v3



DRAM: **4GB**
Emulated PM: **7GB**



Intel SSD
DC S3520



Ubuntu 18.04
Kernel 4.15



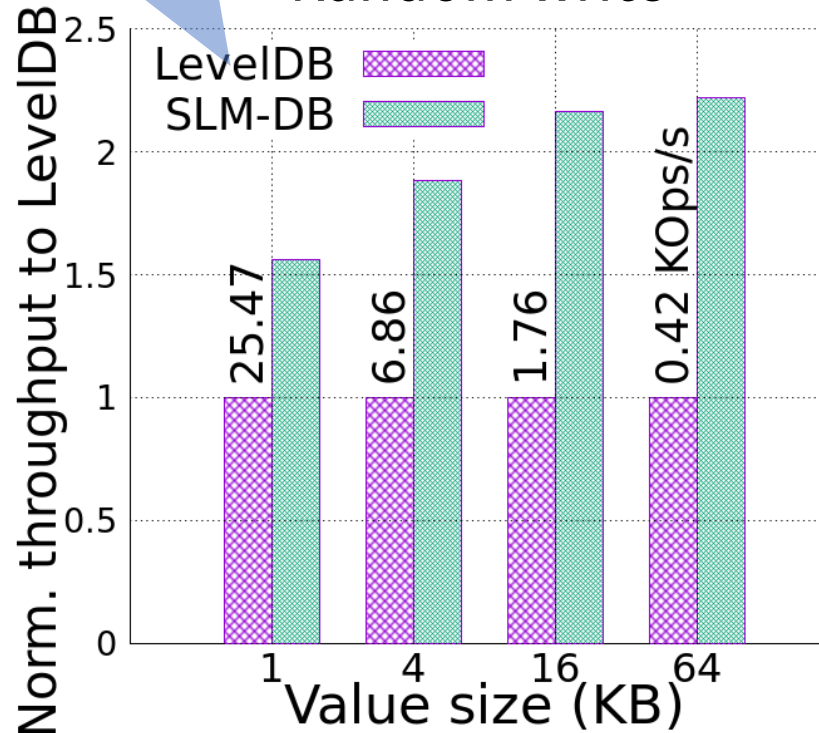
DB: **8GB/20GB**
Memtable: **64MB**

- PM write latency 500ns (5x of DRAM write latency)
- PM read latency & bandwidth same as DRAM's
- Intel's PMDK used to control PM pool

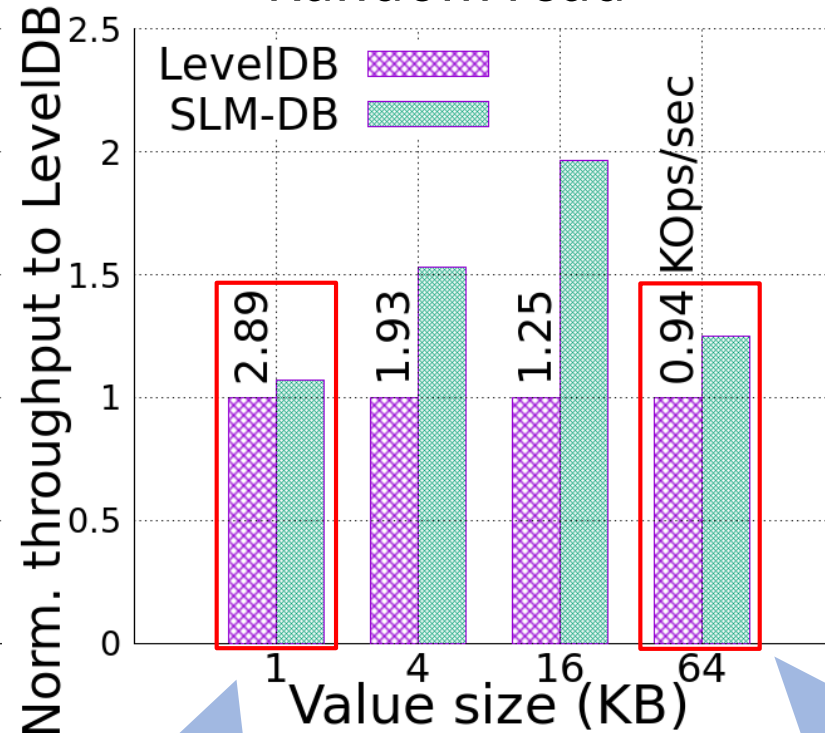
db_bench microbenchmark

Steady performance increase

Random write



Random read

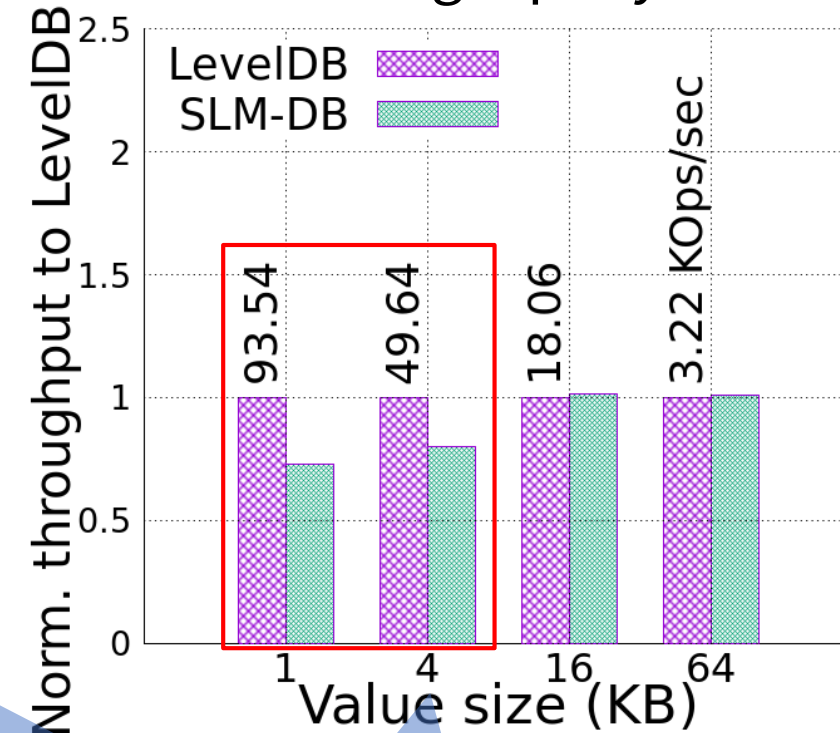


Low file locating overhead

Overhead amortized from large value size

Range size = 100

Range query



Low sequentiality

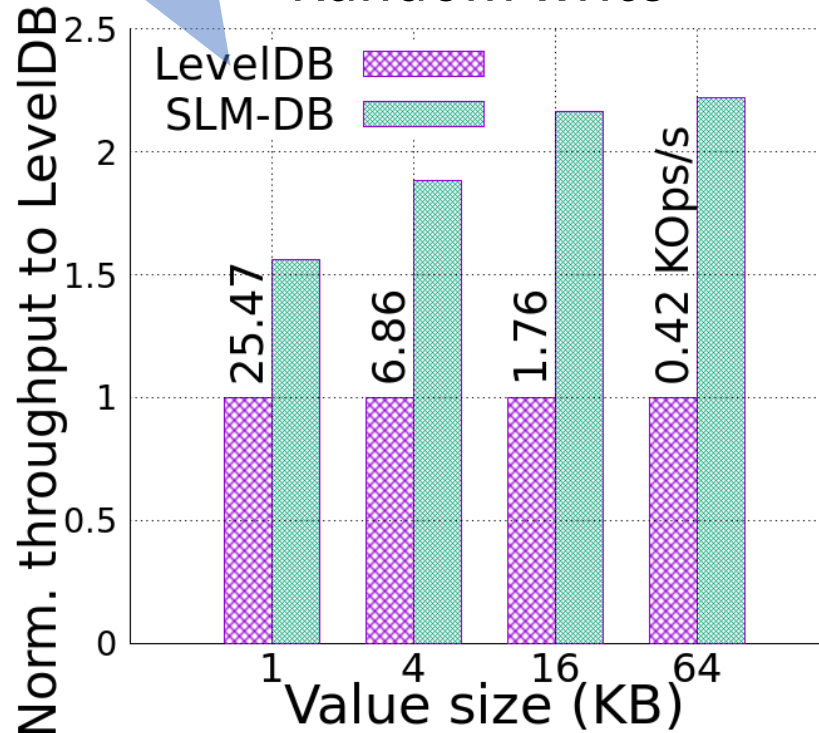
db_bench microbenchmark

Steady performance increase

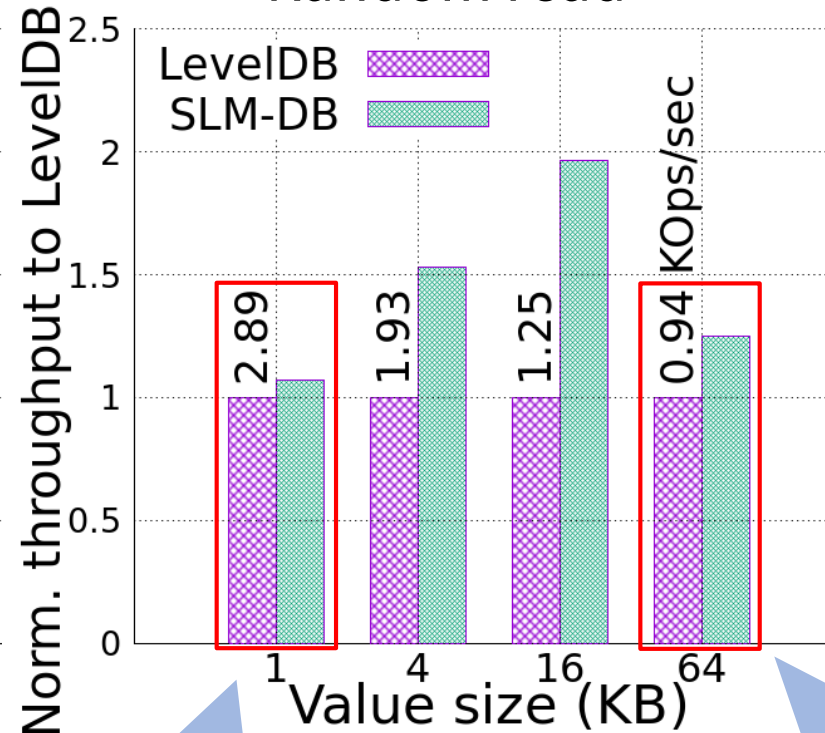
- ~2.56x less disk write amplification
- Max 700MB used in PM

Range size = 100

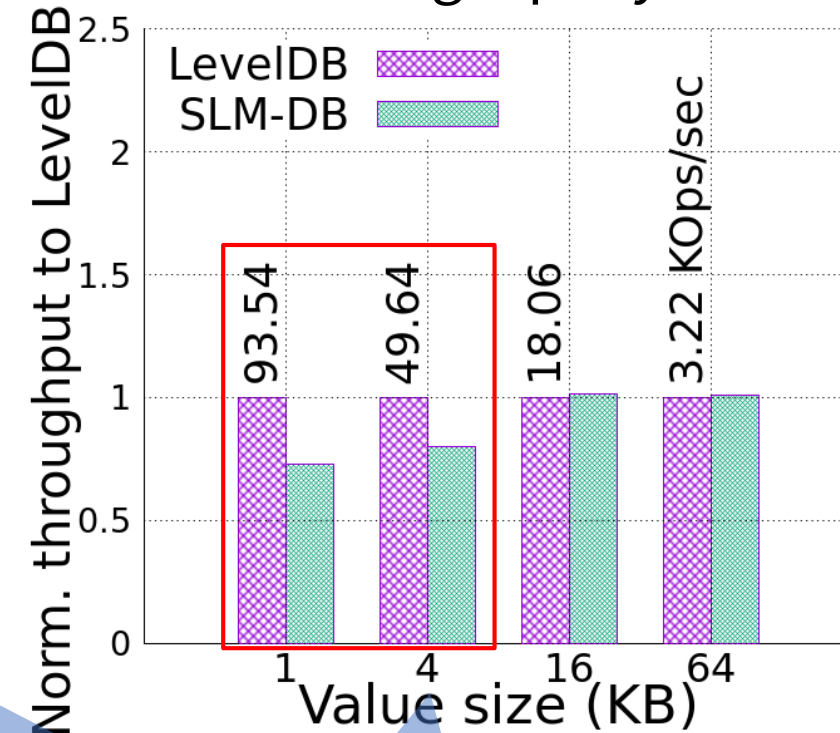
Random write



Random read



Range query



Low file locating overhead

Overhead amortized from large value size

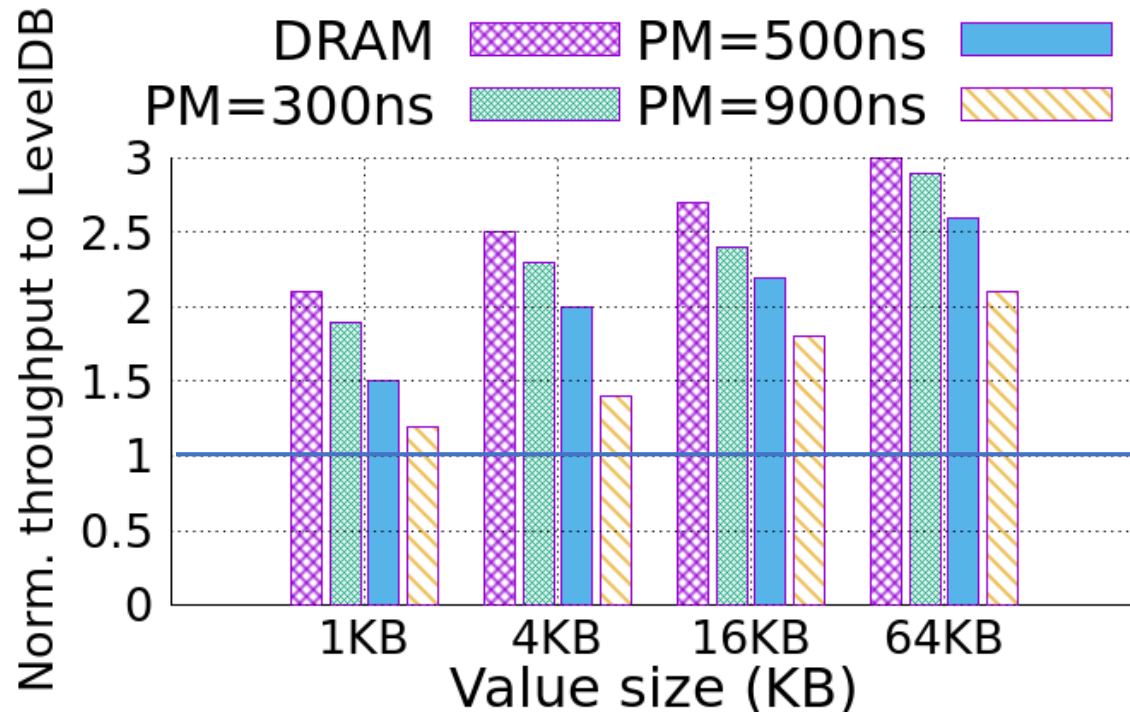
Low sequentiality

PM sensitivity

db_bench

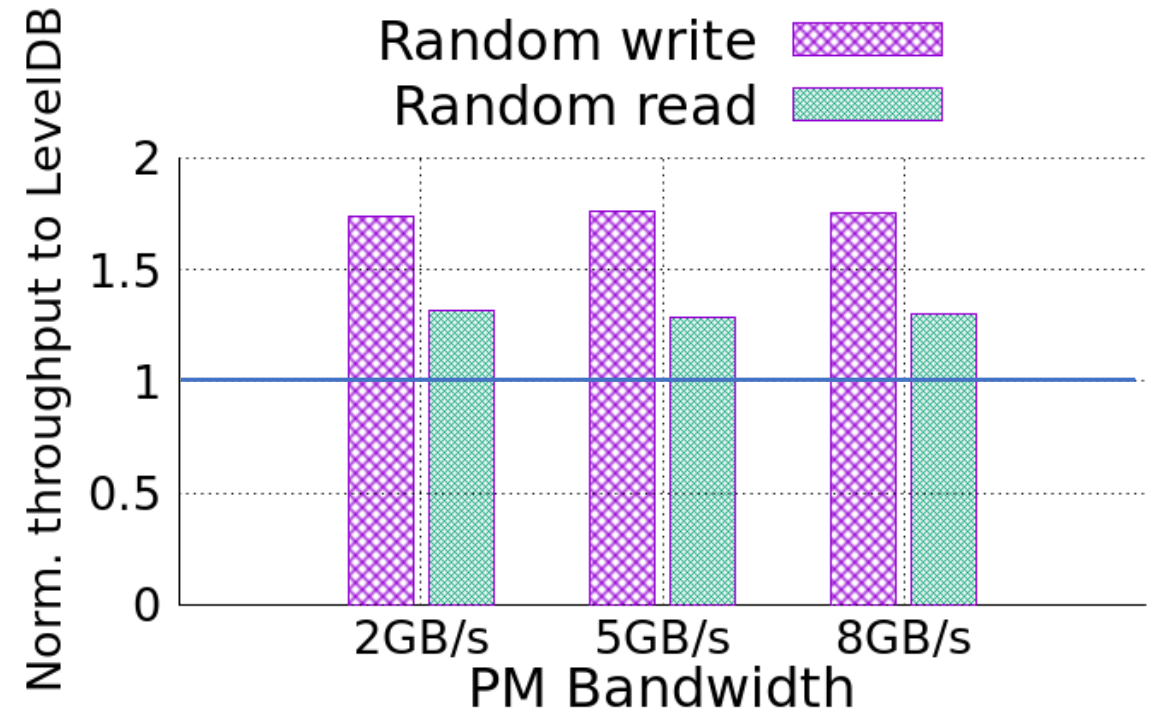
PM write latency sensitivity

Random write benchmark



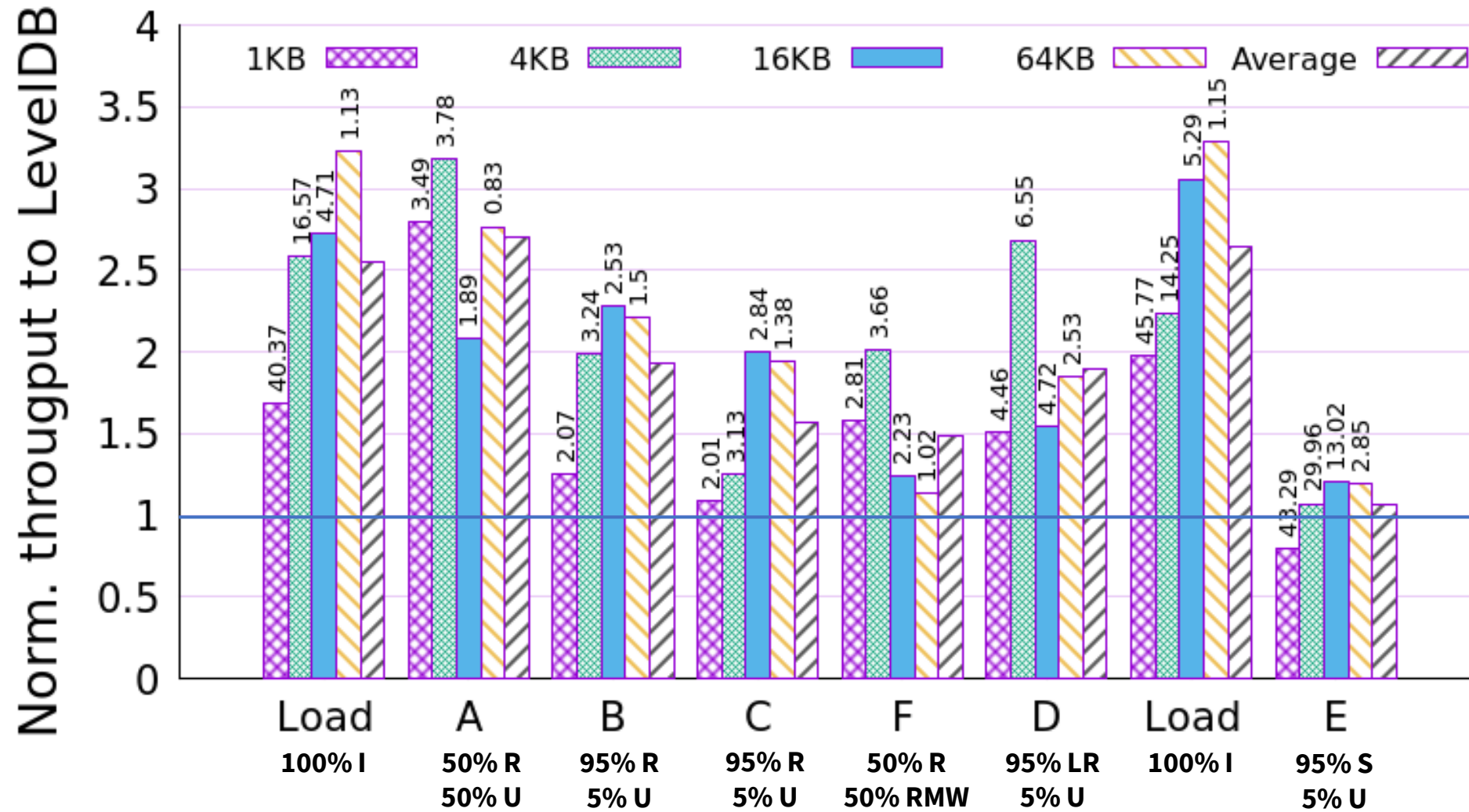
Emulated by inserting cpu pause after clflush()

PM bandwidth sensitivity

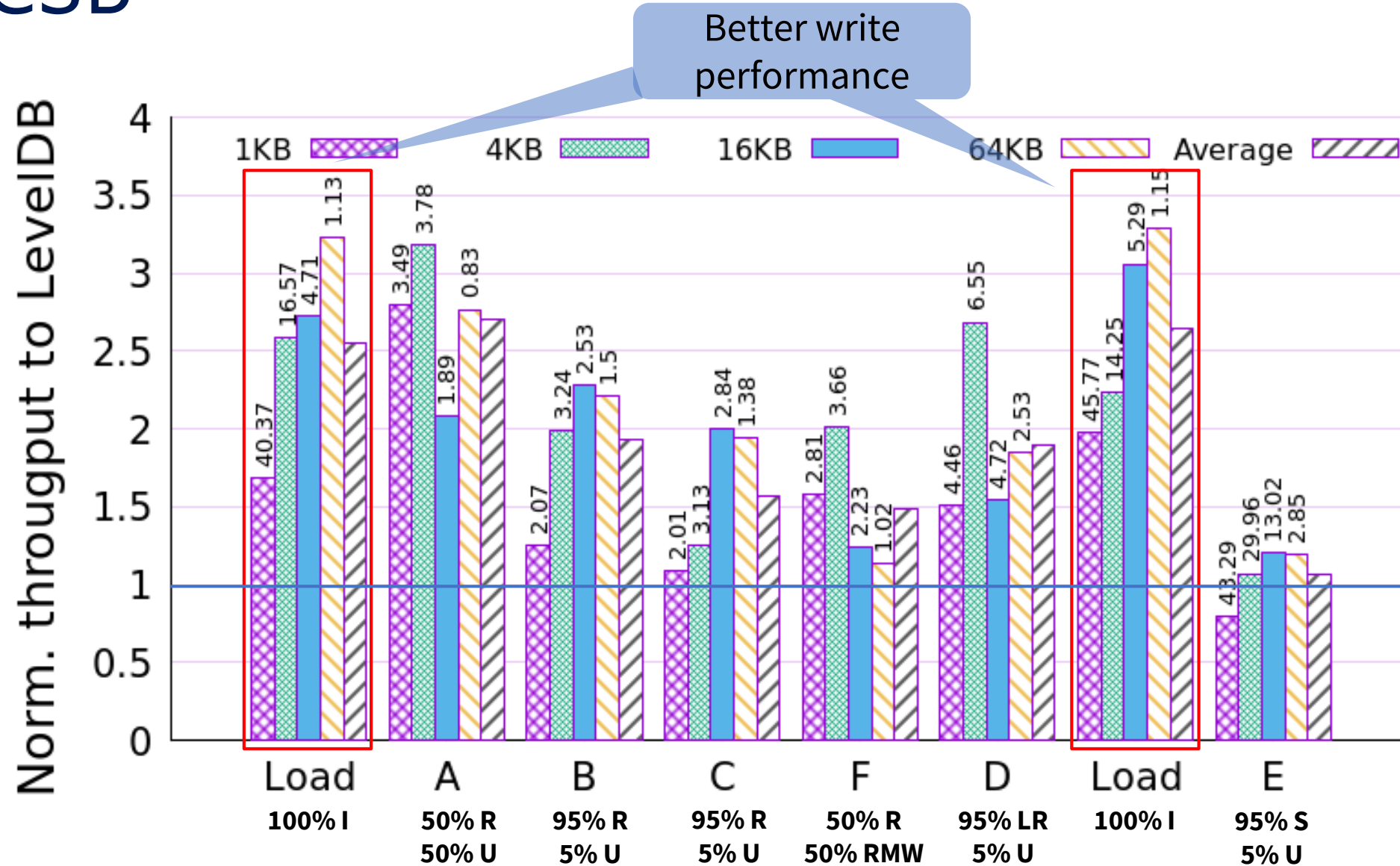


Emulated using Thermal Throttling

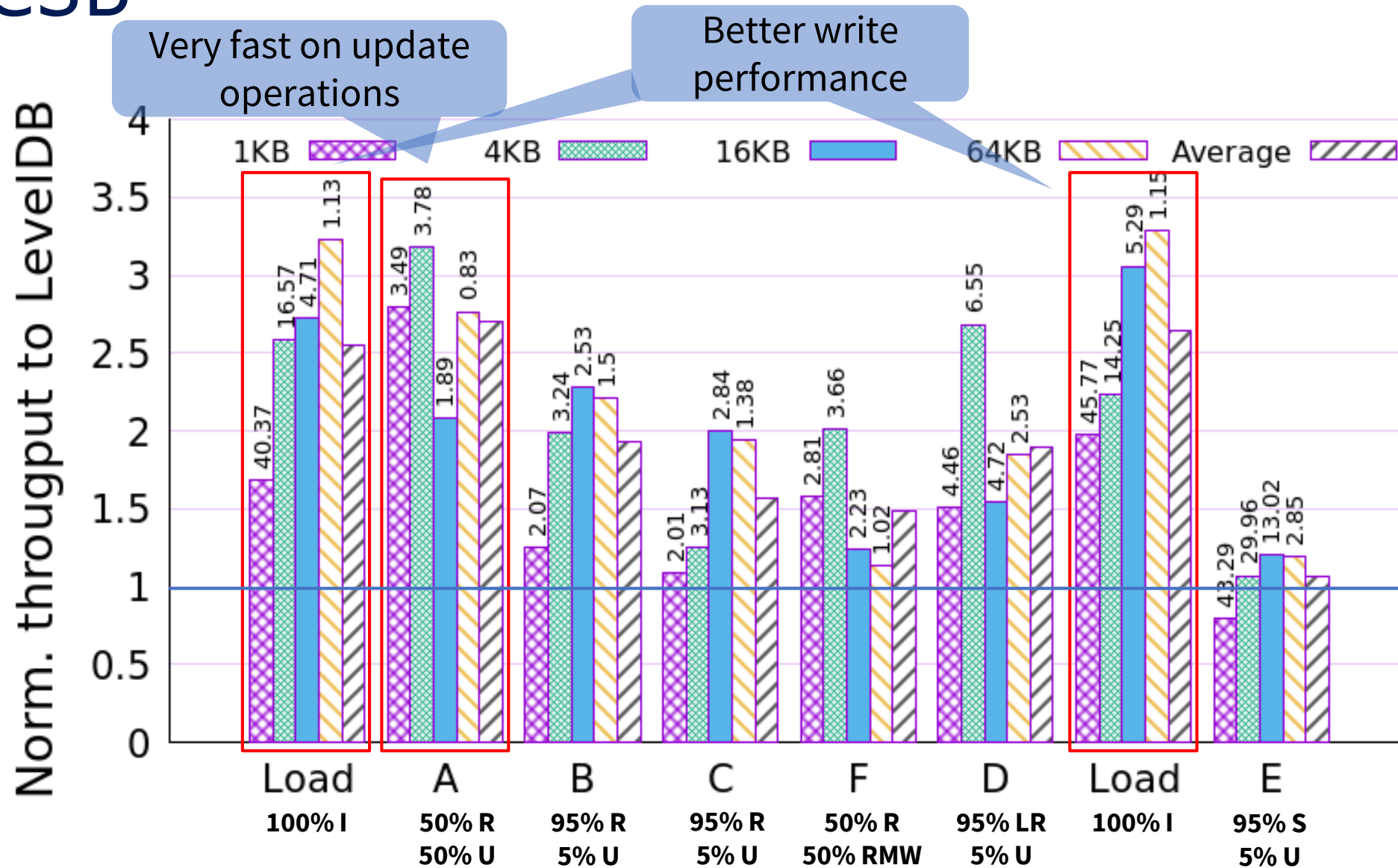
YCSB



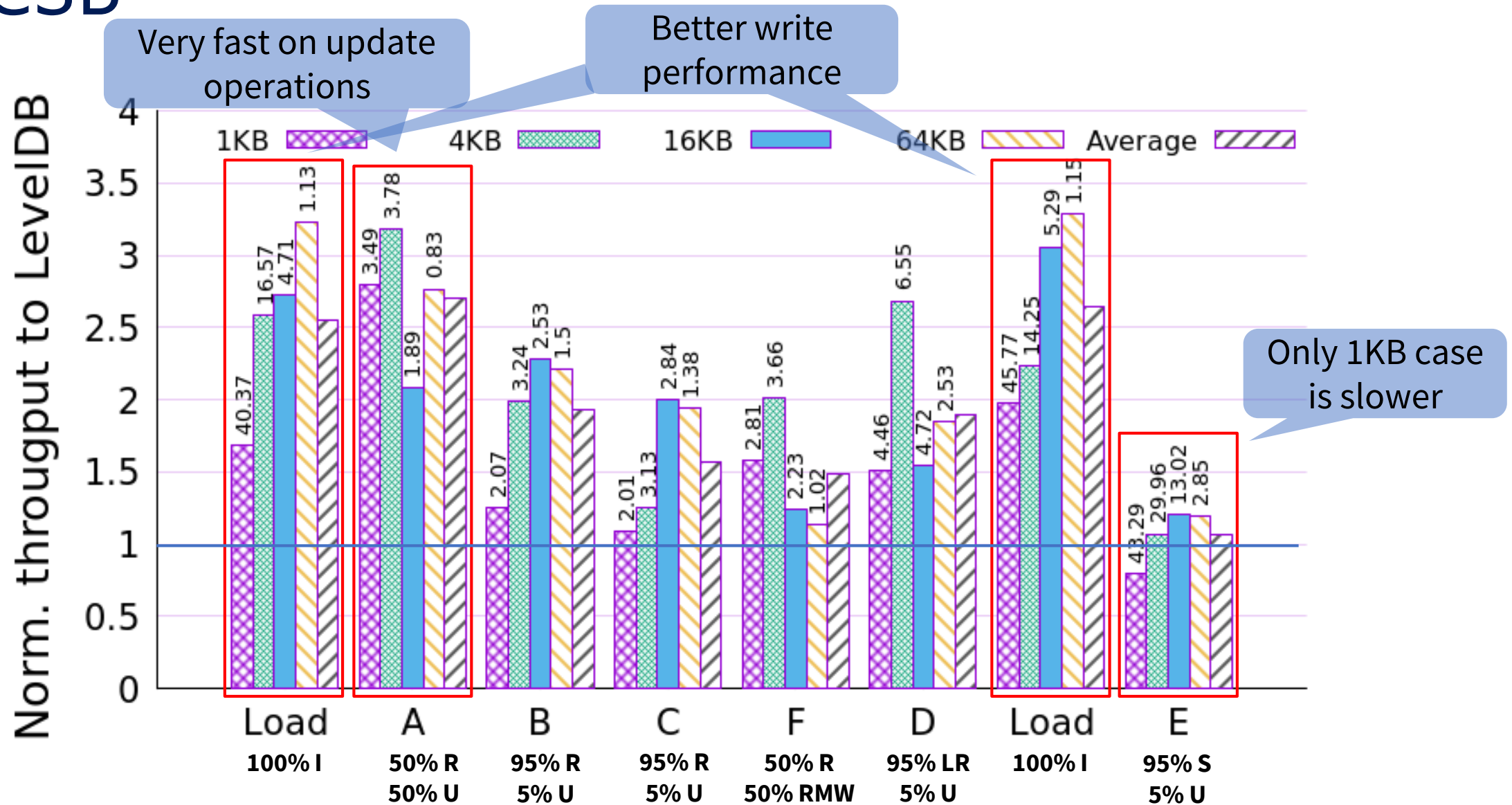
YCSB



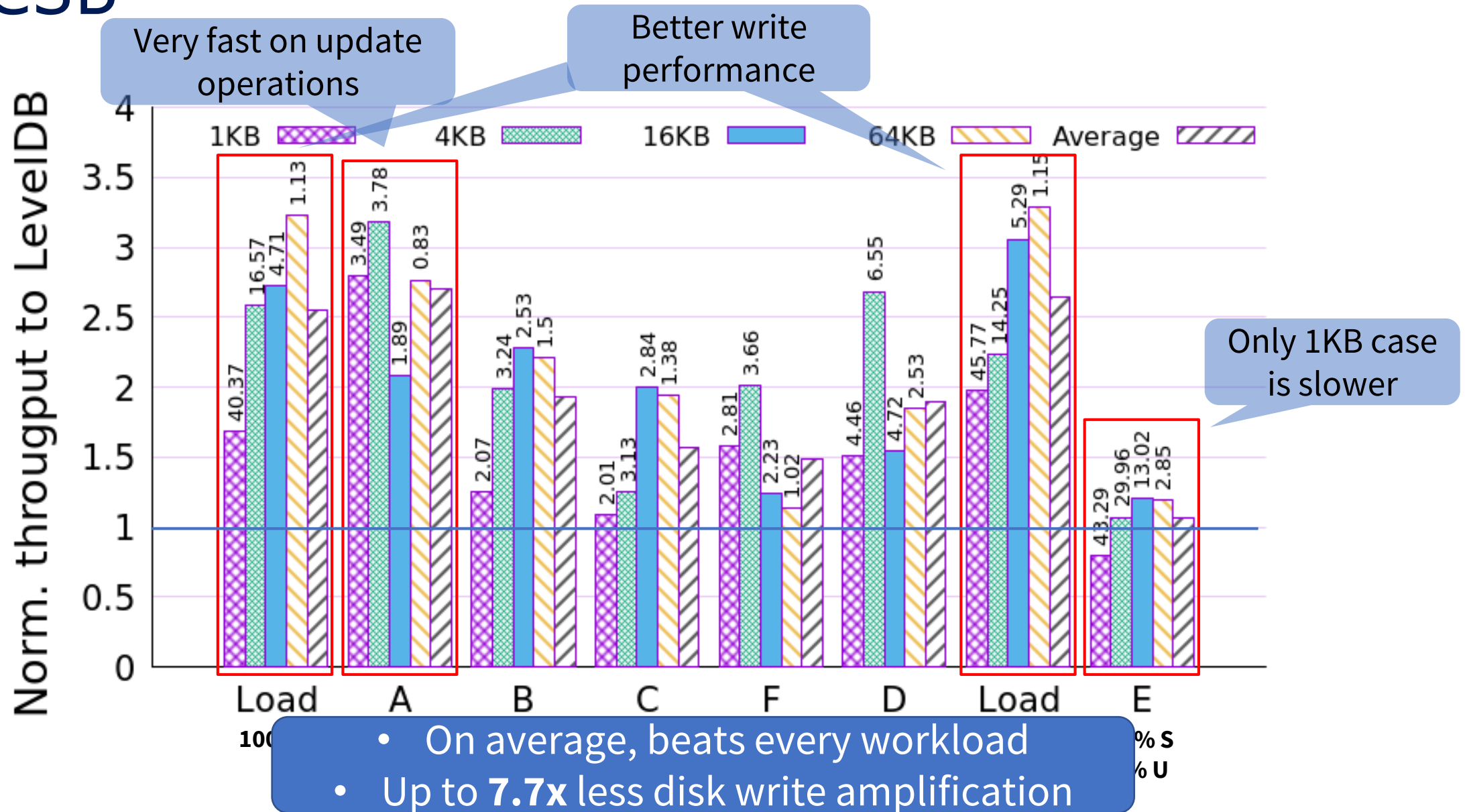
YCSB



YCSB



YCSB



Conclusion

- Novel design of Key-Value stores with Persistent Memory
- High write/read performance compared to LevelDB
- Comparable scan performance
- Low write amplification
- Near-optimal read amplification

Thanks!

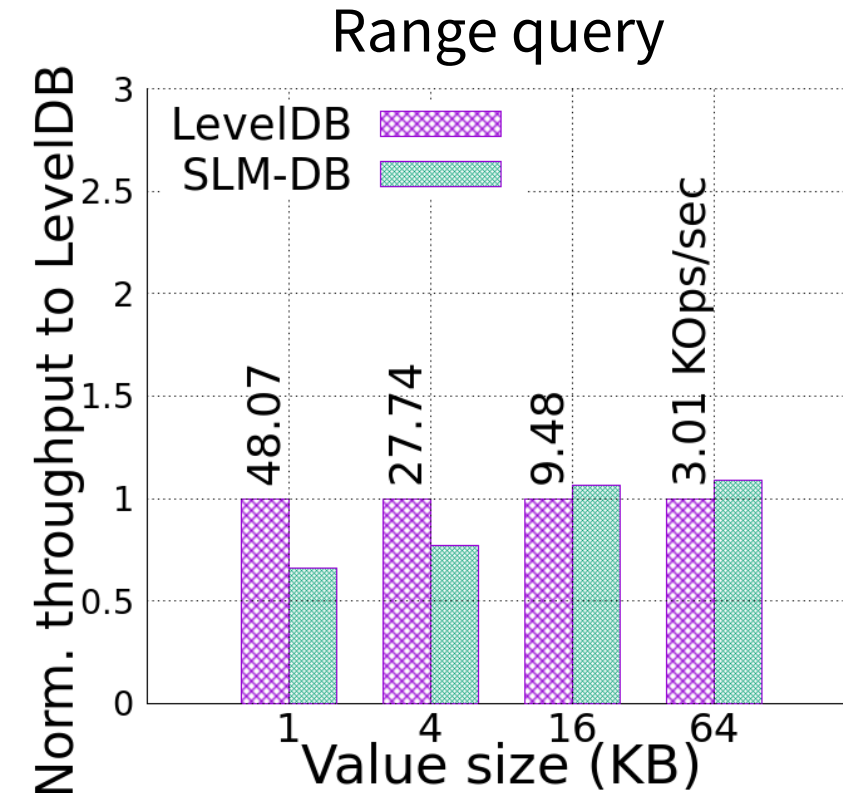
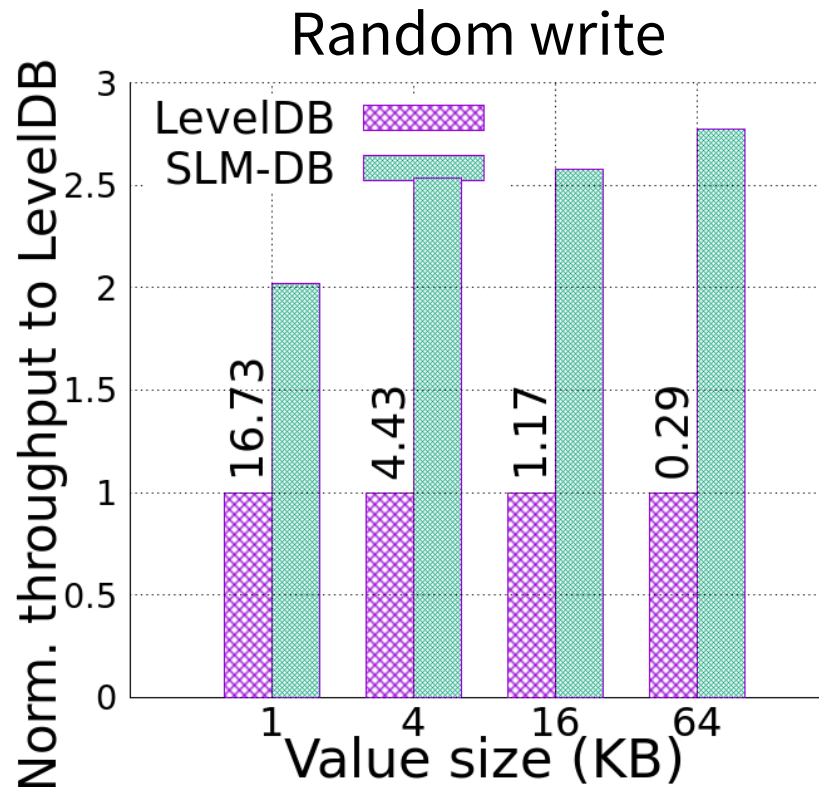
Questions?

SLM-DB: Single Level Merge Key-Value Store with Persistent Memory

Olzhas Kaiyrakhmet, Songyi Lee, Beomseok Nam, Sam H. Noh, Young-ri Choi

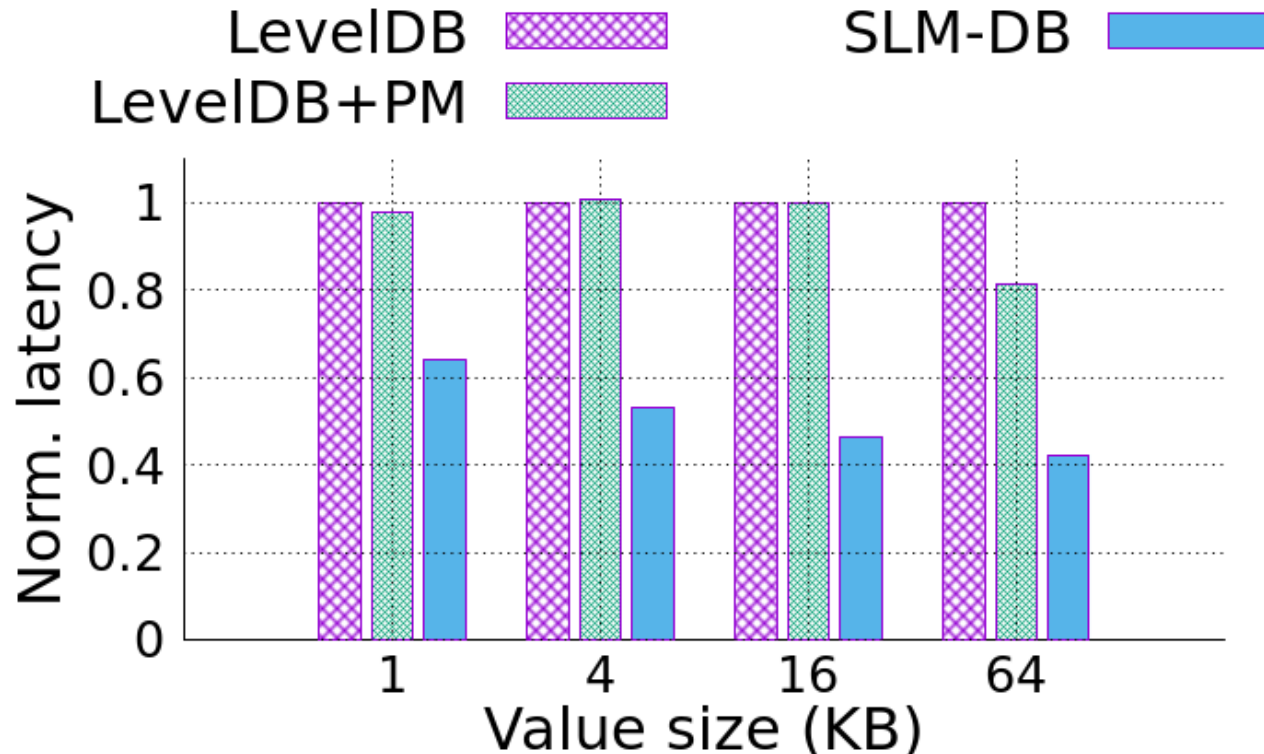


db_bench microbenchmark (20GB)

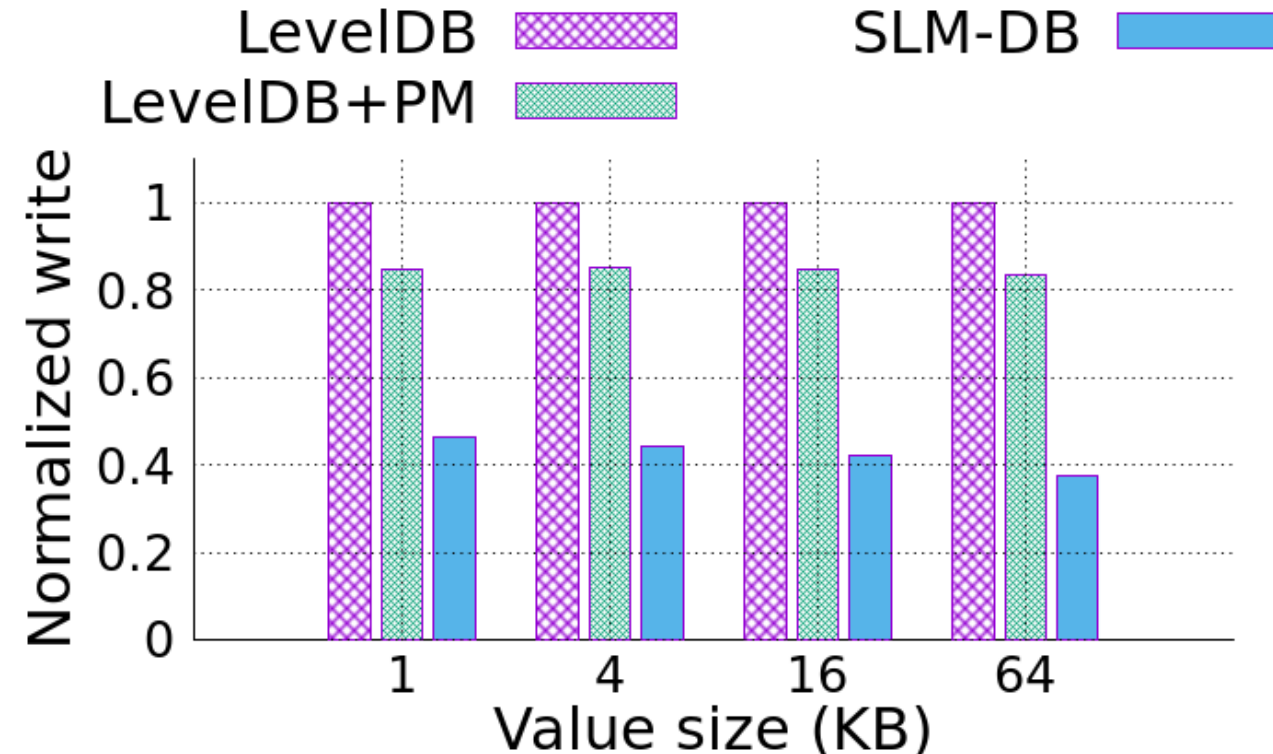


Effect of persistent MemTable

Random write performance

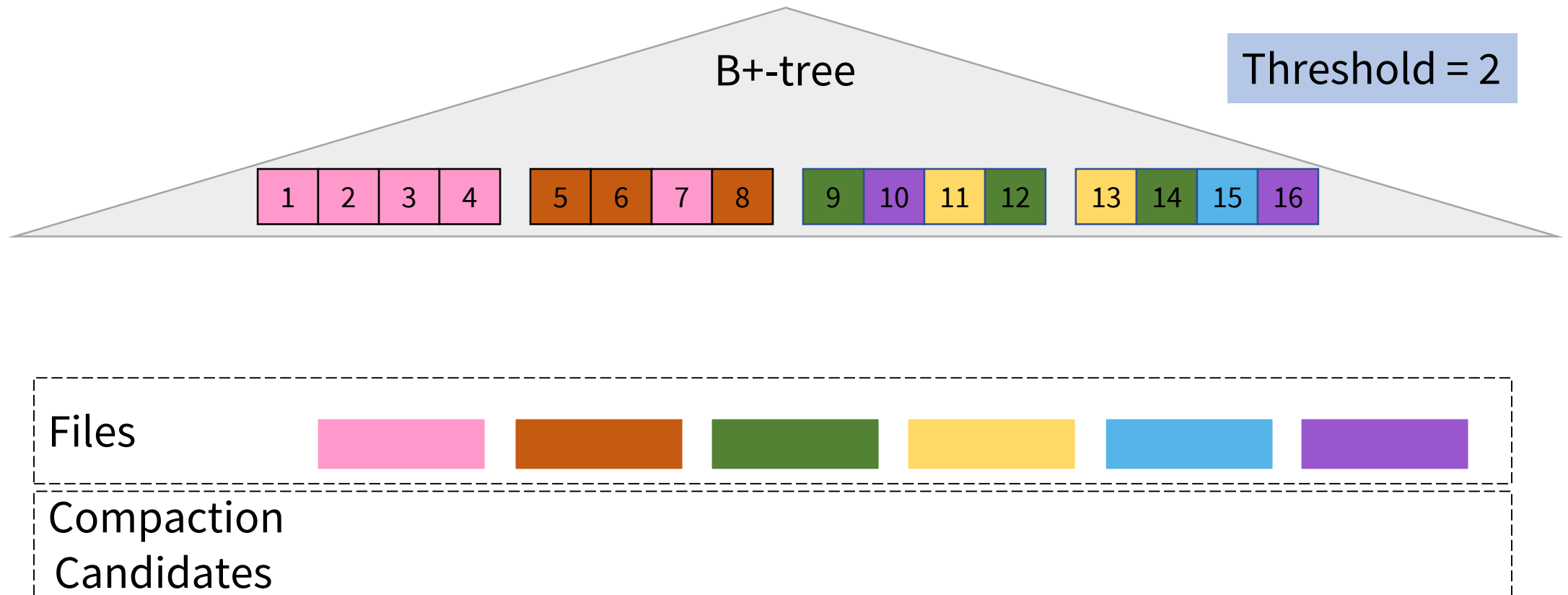


Total disk write



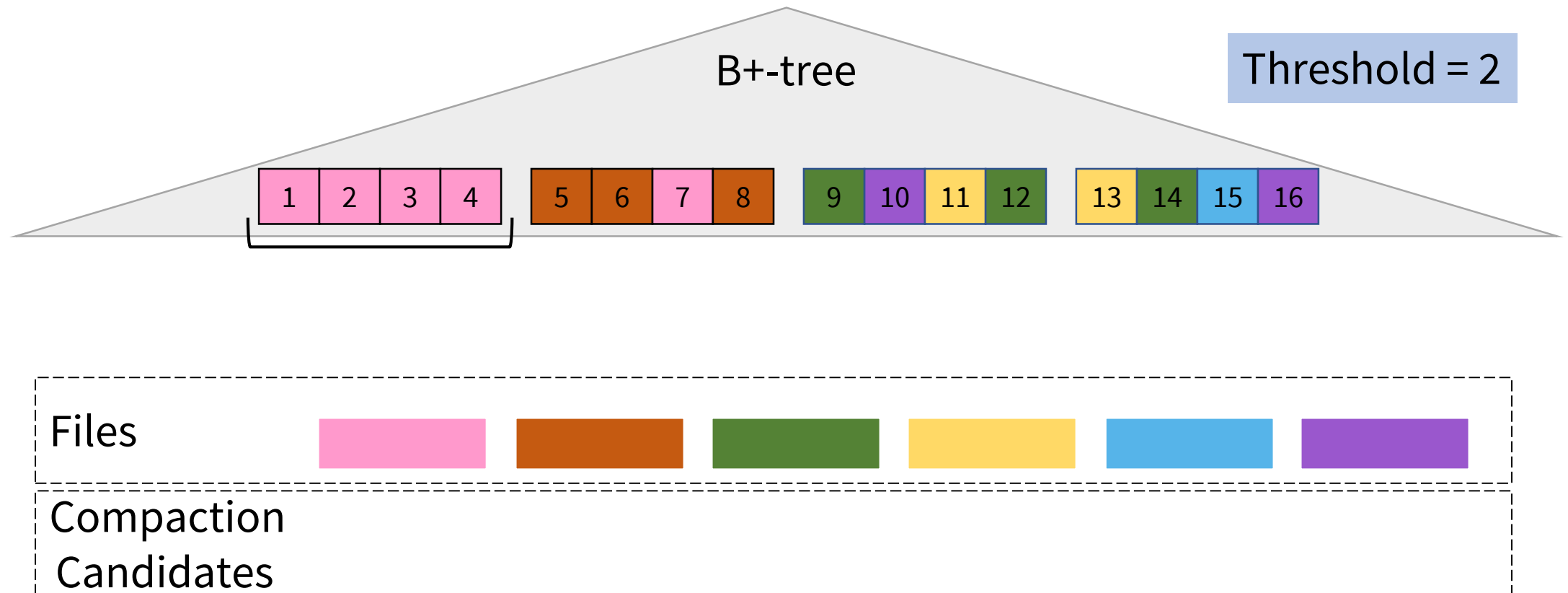
B+-tree leaf node scan

- To increase sequentiality of key-values with scans in round-robin fashion
- If the number of unique file accesses is above threshold, then add to candidates



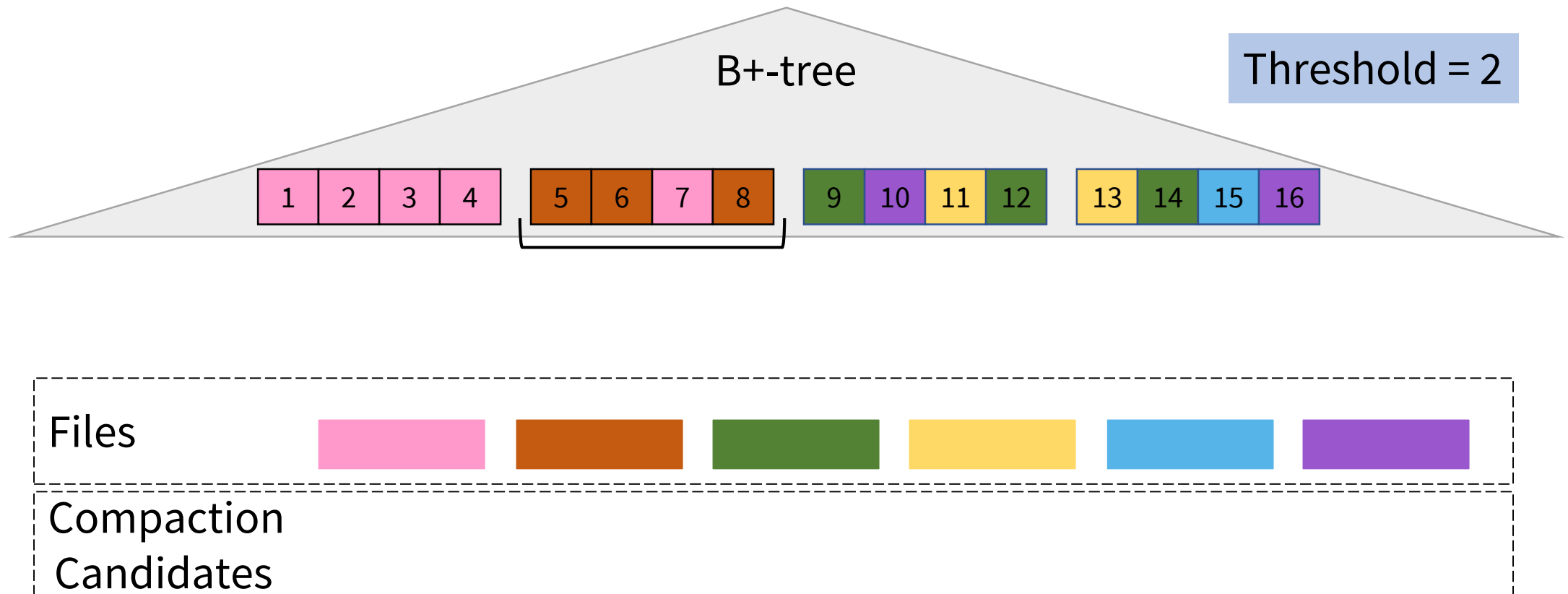
B+-tree leaf node scan

- To increase sequentiality of key-values with scans in round-robin fashion
- If the number of unique file accesses is above threshold, then add to candidates



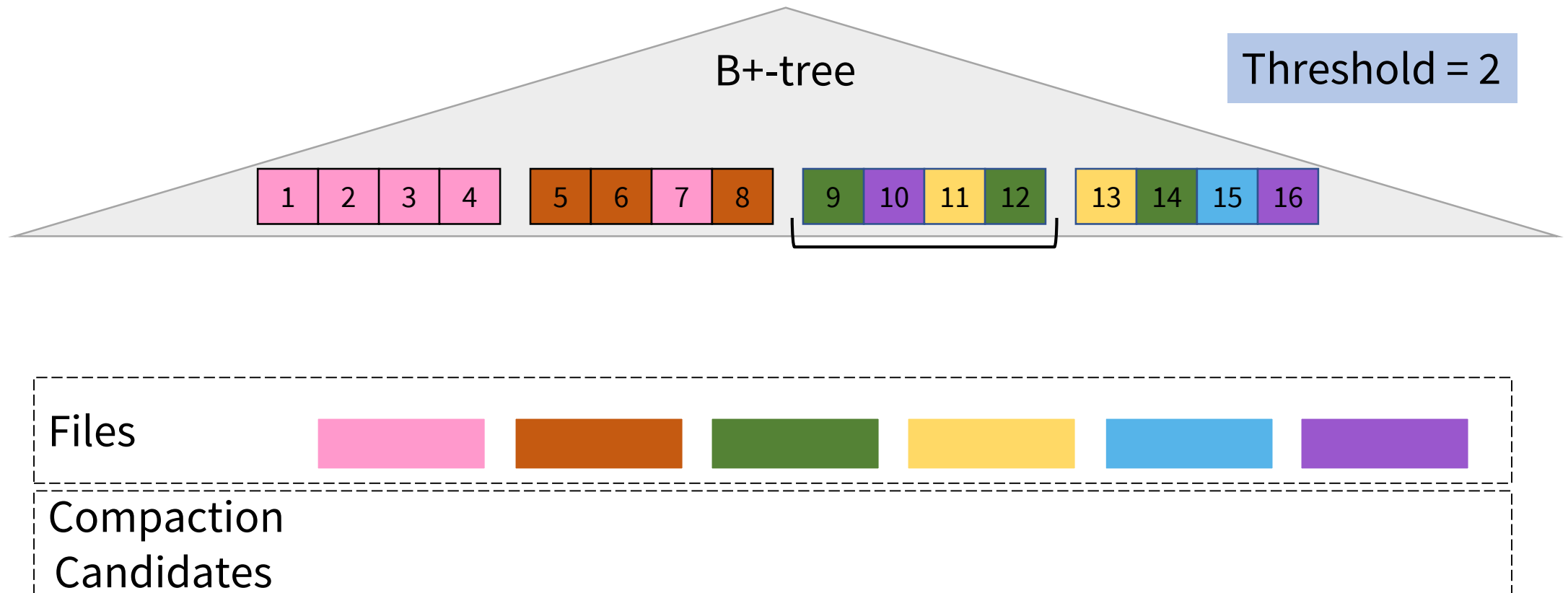
B+-tree leaf node scan

- To increase sequentiality of key-values with scans in round-robin fashion
- If the number of unique file accesses is above threshold, then add to candidates



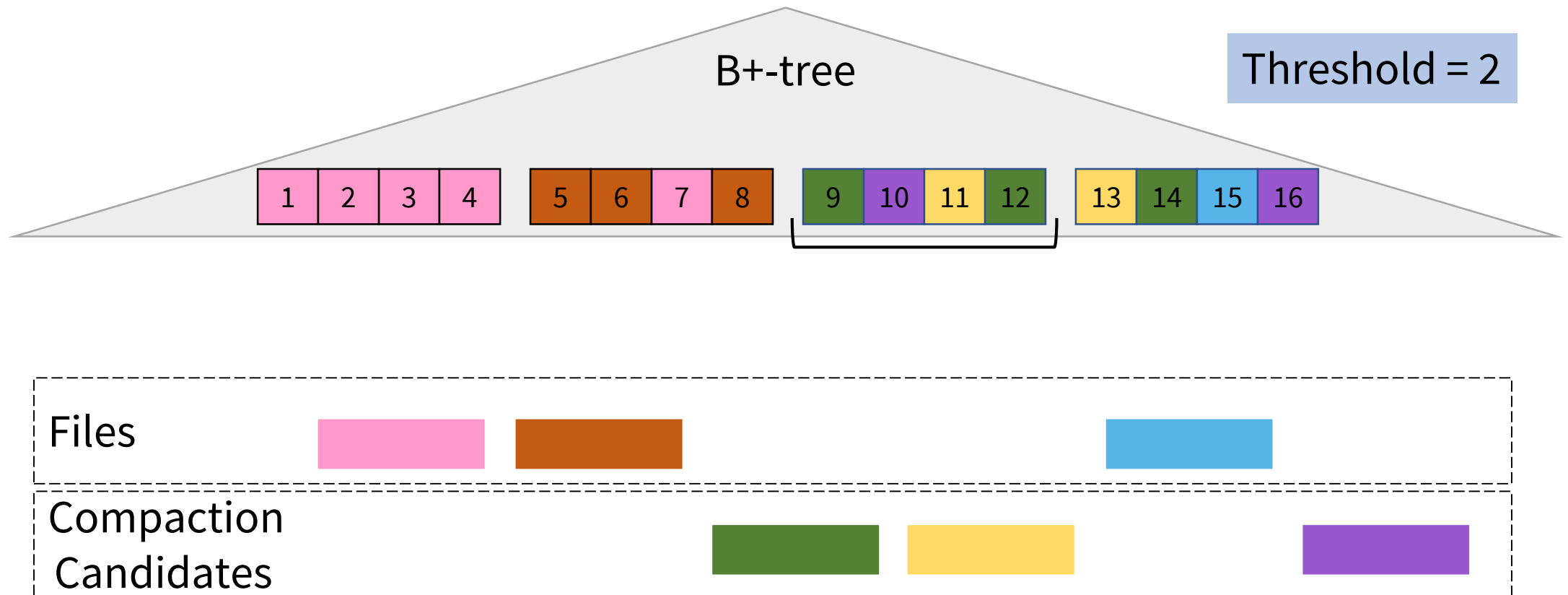
B+-tree leaf node scan

- To increase sequentiality of key-values with scans in round-robin fashion
- If the number of unique file accesses is above threshold, then add to candidates



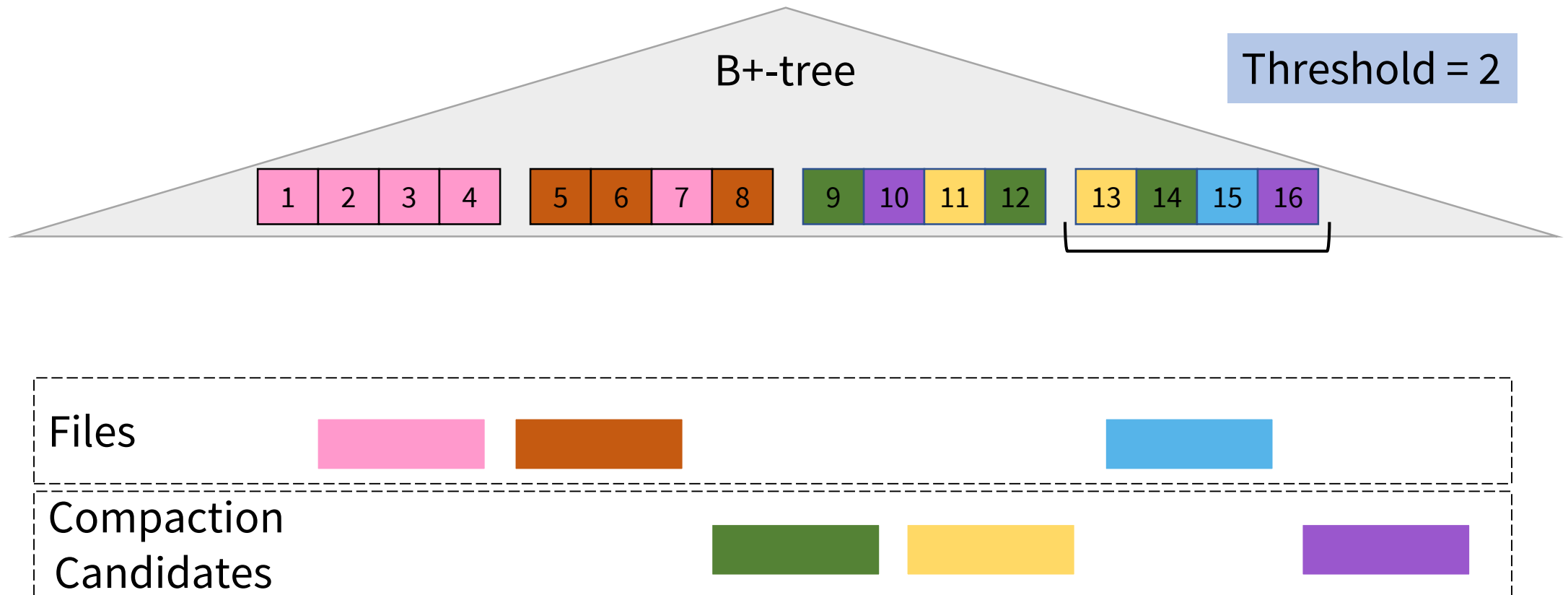
B+-tree leaf node scan

- To increase sequentiality of key-values with scans in round-robin fashion
- If the number of unique file accesses is above threshold, then add to candidates



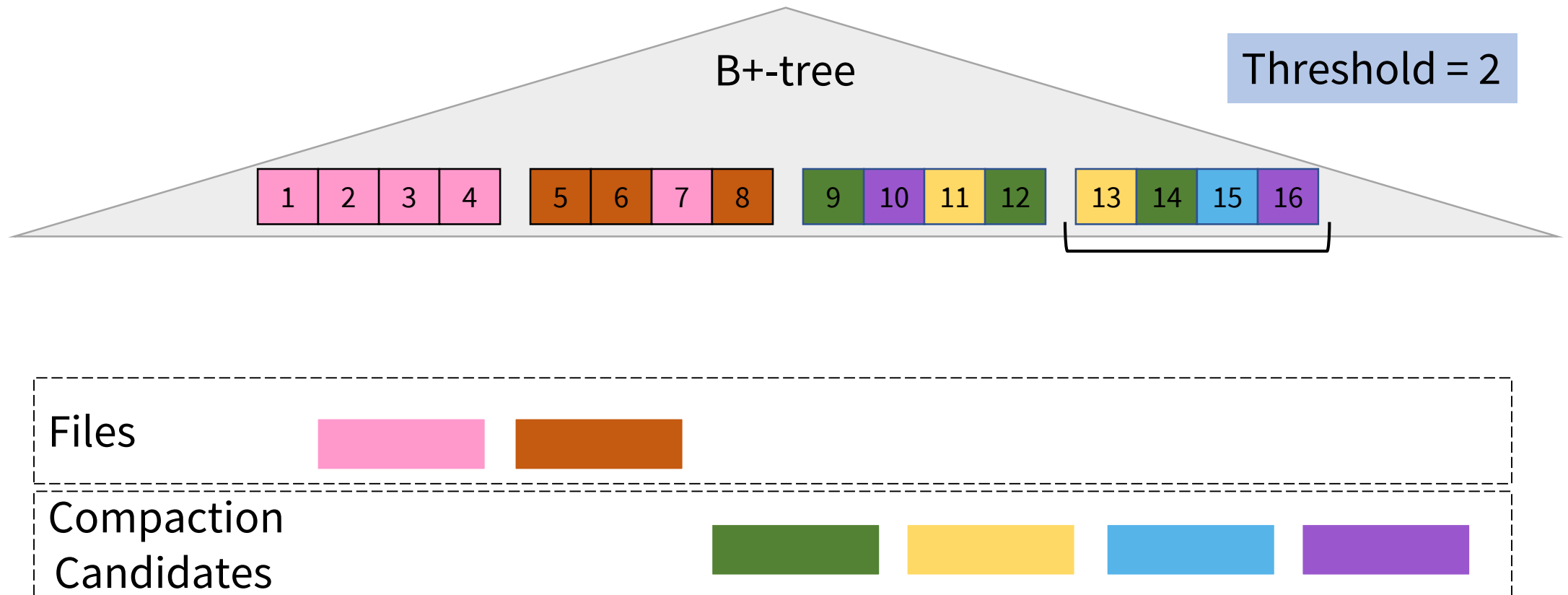
B+-tree leaf node scan

- To increase sequentiality of key-values with scans in round-robin fashion
- If the number of unique file accesses is above threshold, then add to candidates



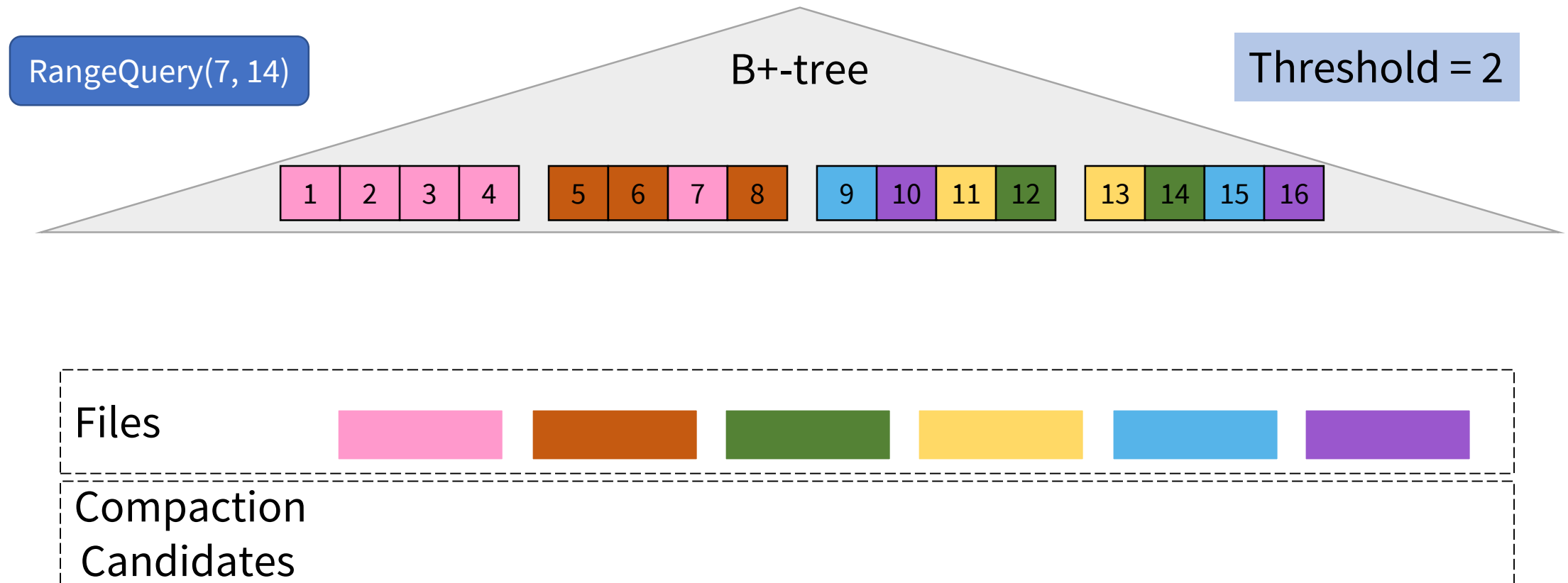
B+-tree leaf node scan

- To increase sequentiality of key-values with scans in round-robin fashion
- If the number of unique file accesses is above threshold, then add to candidates



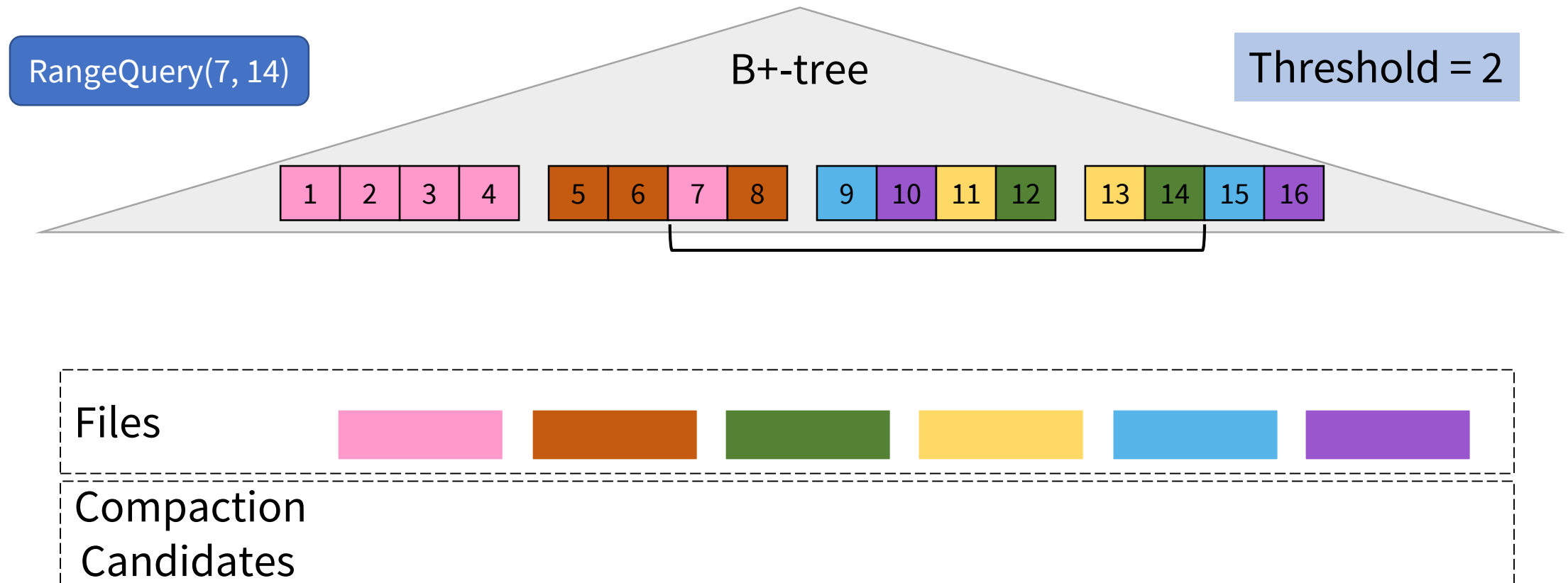
Degree of sequentiality per range query

- To increase sequentiality of key-values during range query operation
- If subrange max unique file accesses is above threshold, then add to candidates



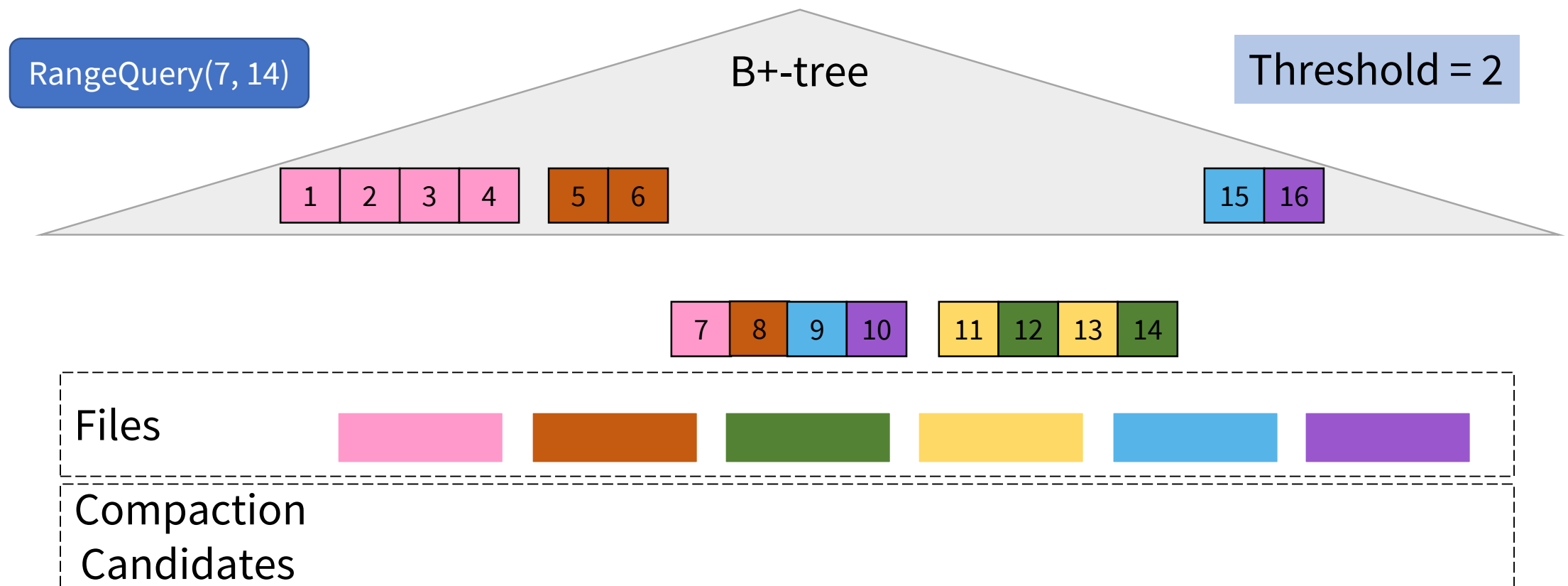
Degree of sequentiality per range query

- To increase sequentiality of key-values during range query operation
- If subrange max unique file accesses is above threshold, then add to candidates



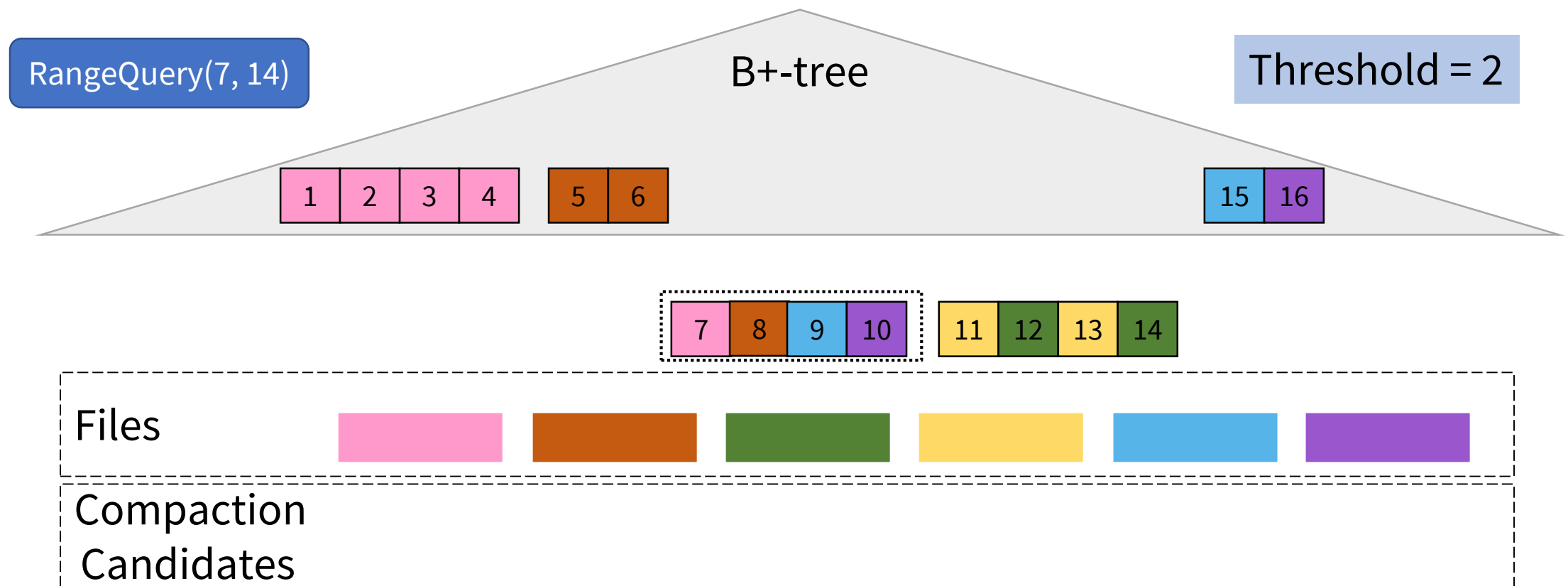
Degree of sequentiality per range query

- To increase sequentiality of key-values during range query operation
- If subrange max unique file accesses is above threshold, then add to candidates



Degree of sequentiality per range query

- To increase sequentiality of key-values during range query operation
- If subrange max unique file accesses is above threshold, then add to candidates



Degree of sequentiality per range query

- To increase sequentiality of key-values during range query operation
- If subrange max unique file accesses is above threshold, then add to candidates

