# ChameleonDB:
# a Key-value Store for Optane Persistent memory
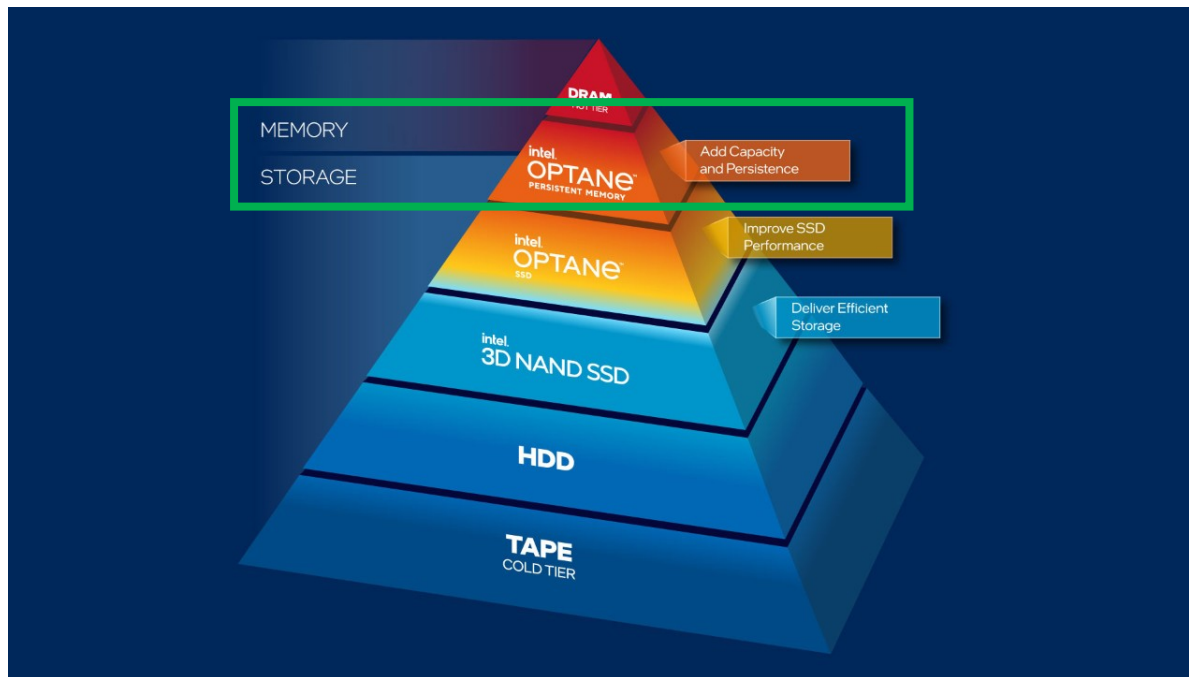
Wenhui Zhang, Xingsheng Zhao,

Song Jiang, and Hong Jiang
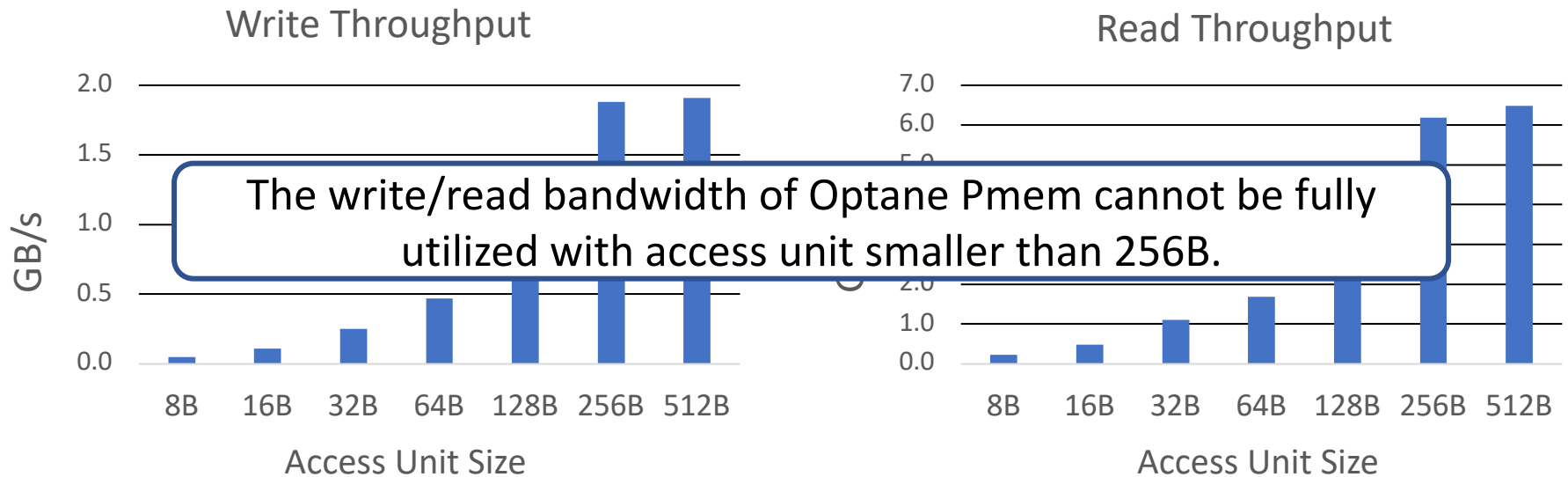

University of Texas at Arlington

# Why do we need a key-value store design for Optane persistent memory?



Because Optane Persistent Memory (Pmem) is very **different.**

It's different from **DRAM**, different from **traditional block devices**, and even different from **what was assumed about persistent memory** (slower, persistent DRAM).

# Optane Pmem is a block device with access unit as 256B.

**Write Throughput**

GB/s

2.0
1.5
1.0
0.5
0.0

8B  16B  32B  64B  128B  256B  512B

Access Unit Size

**Read Throughput**

7.0
6.0
5.0
2.0
1.0
0.0

8B  16B  32B  64B  128B  256B  512B

Access Unit Size

> The write/read bandwidth of Optane Pmem cannot be fully utilized with access unit smaller than 256B.
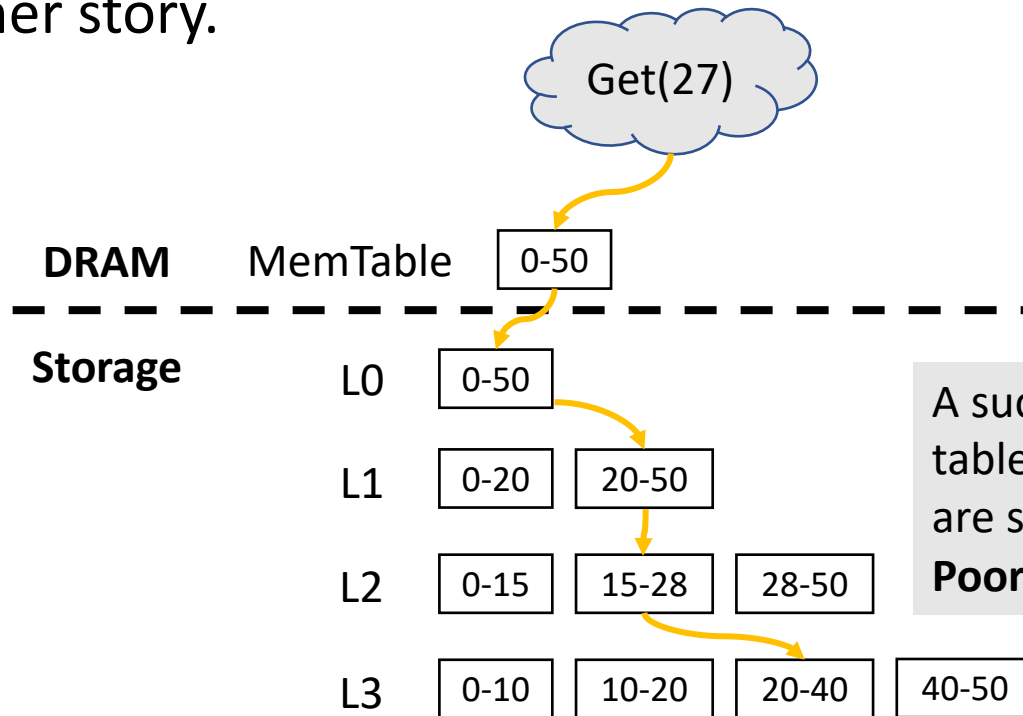
This property makes Optane Pmem different from what was assumed about persistent memory (cacheline as access unit).

KV store designs **employing small random writes** of persistent memory, including Level Hashing, CCEH, and FAST&FAIR, are **unable** to provide high **write performance** on Optane Pmem.

# Are LSM-tree based KV stores designed for block devices efficient for Optane Pmem?

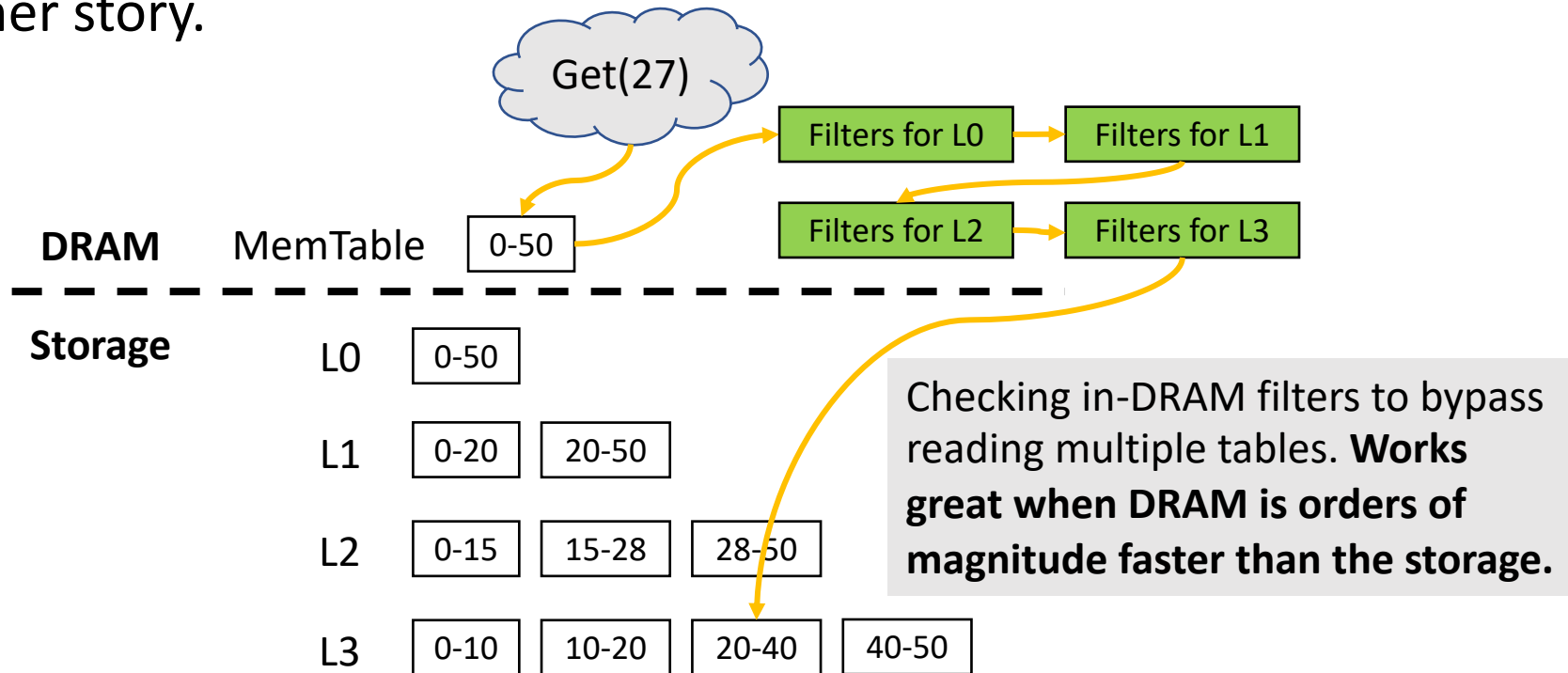They are efficient on write performance, but **read performance** is another story.

Get(27)

| | | |
|---|---|---|
| **DRAM** | MemTable | 0-50 |

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

**Storage**

L0    0-50

L1    0-20    20-50

L2    0-15    15-28    28-50

L3    0-10    10-20    20-40    40-50

A successful Get may need to read tables from multiple levels, which are stored in the slow block device. **Poor read performance.**
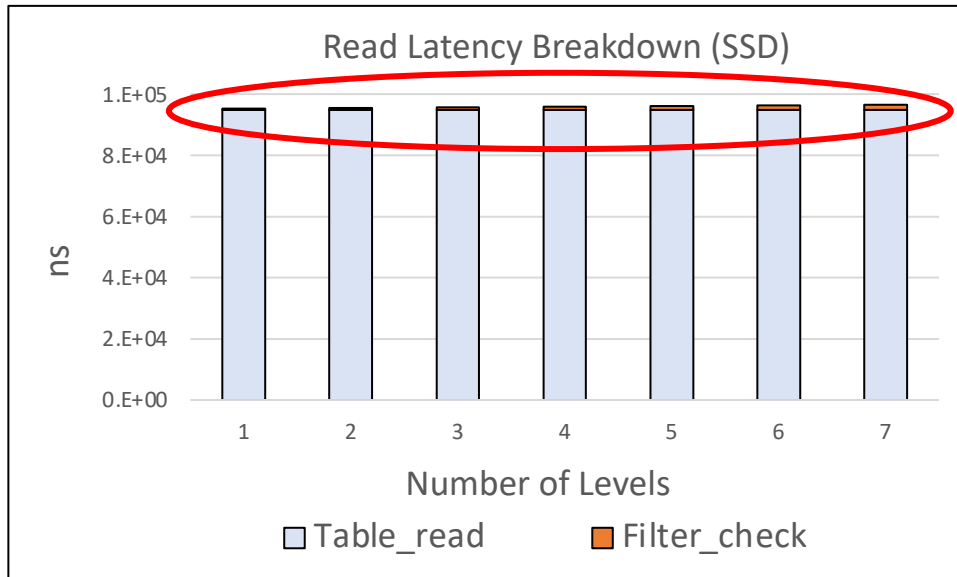
# Are LSM-tree based KV stores designed for block devices efficient for Optane Pmem?

They are efficient on write performance, but **read performance** is another story.



Get(27)

Filters for L0 → Filters for L1

Filters for L2 → Filters for L3

**DRAM** MemTable | 0-50 |

**Storage**

L0 | 0-50 |

L1 | 0-20 | 20-50 |

L2 | 0-15 | 15-28 | 28-50 |

L3 | 0-10 | 10-20 | 20-40 | 40-50 |

Checking in-DRAM filters to bypass reading multiple tables. **Works great when DRAM is orders of magnitude faster than the storage.**

# Are LSM-tree based KV stores designed for block devices efficient for Optane Pmem?

They are efficient on write performance, but **read performance** is another story.

When use **SSD** as storage, the time to checking filters in DRAM for multiple levels is **negligible**.



Read Latency Breakdown (SSD)

ns — Number of Levels

Table_read    Filter_check

Reading table from the slow storage (SSD) contributes to 99% of the read latency, while checking multiple in-DRAM filters is nearly negligible.
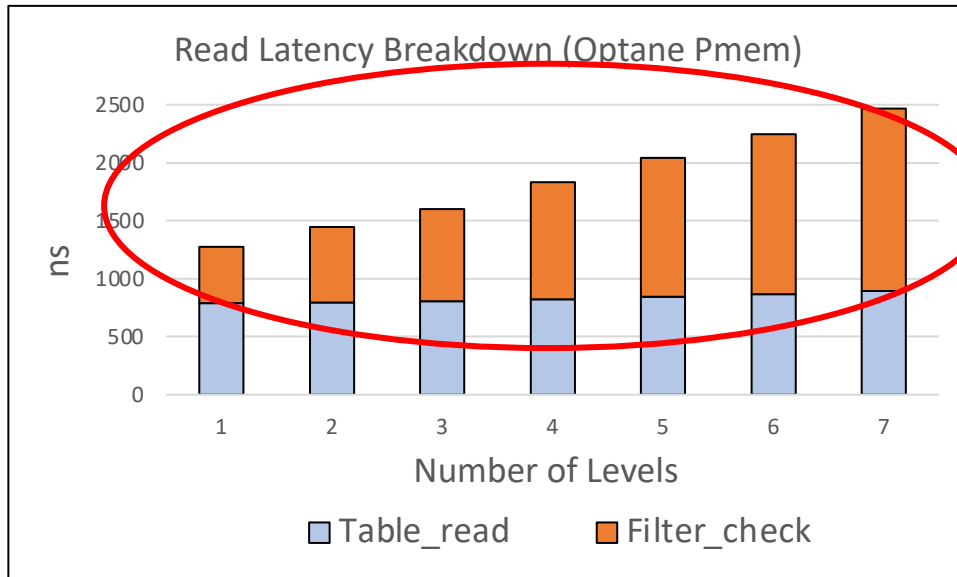
Read latency is stable with the multi-level structure.

# Are LSM-tree based KV stores designed for block devices efficient for Optane Pmem?

They are efficient on write performance, but **read performance** is another story.

When use **Optane Pmem**, whose latency is **~3x DRAM's**, the time to checking filters for multiple levels becomes **significant**.
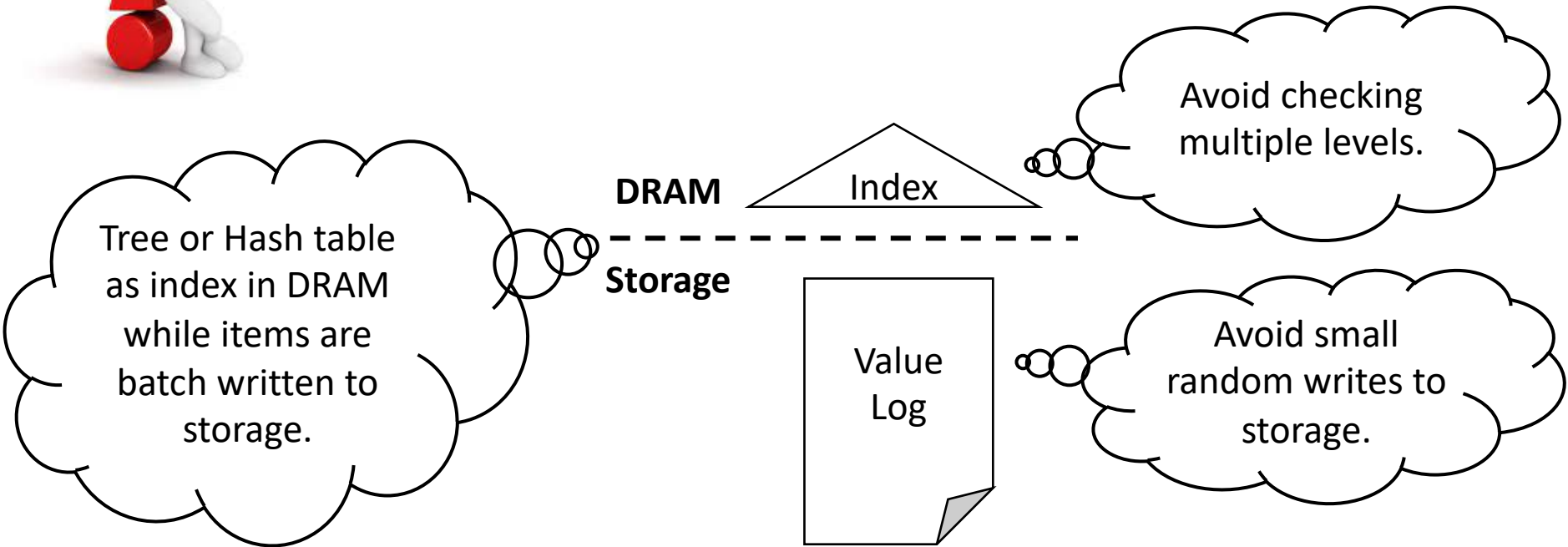


Read Latency Breakdown (Optane Pmem)

Checking filters for multiple levels contributes up to 63% of the read latency.

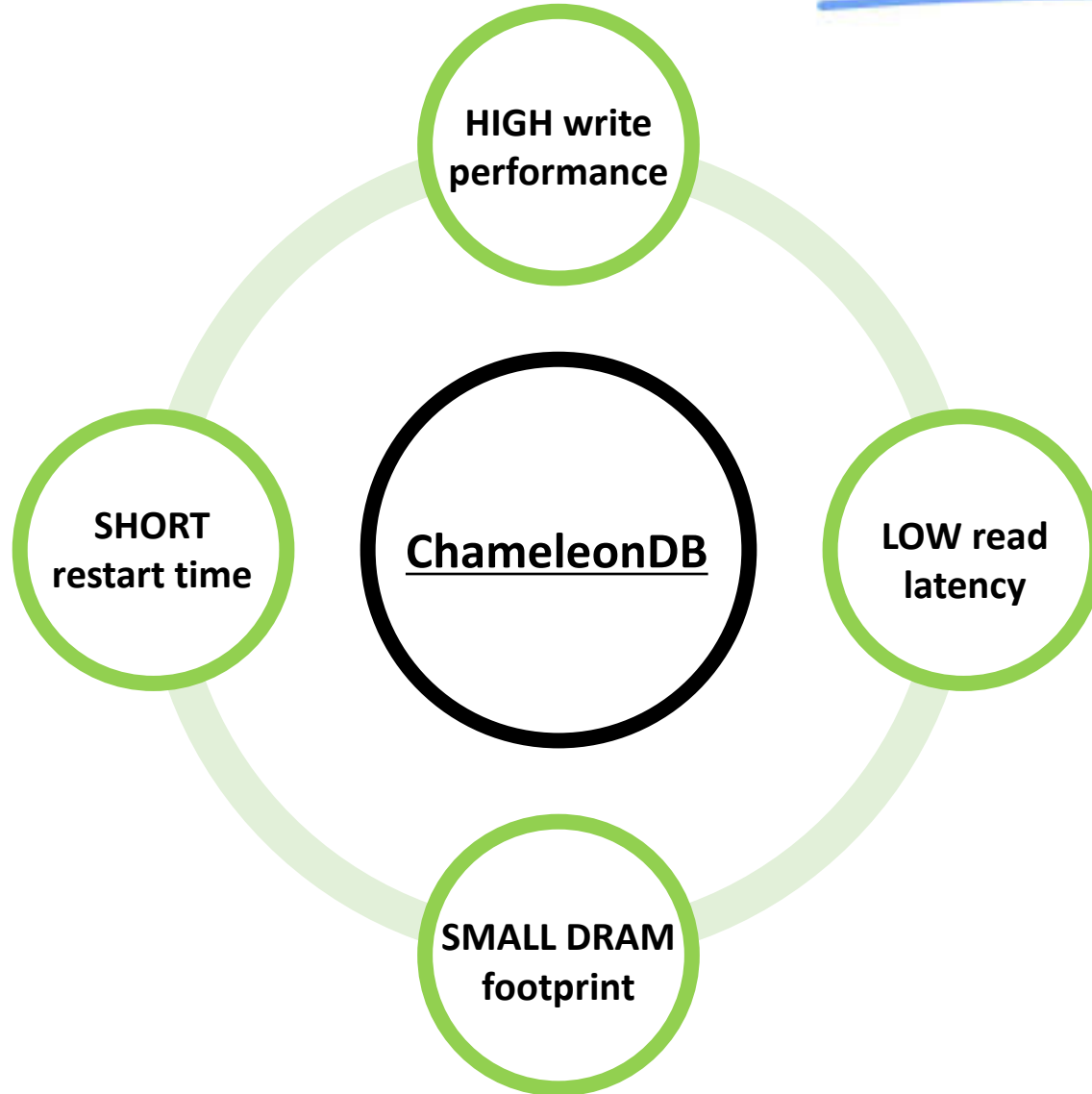**Multi-level structure** becomes a major **barrier** to achieving consistently low **read latency**.

# How about Log structure with in-DRAM index?

**Tree or Hash table as index in DRAM while items are batch written to storage.**

**DRAM** Index

**Avoid checking multiple levels.**

**Storage**

Value Log

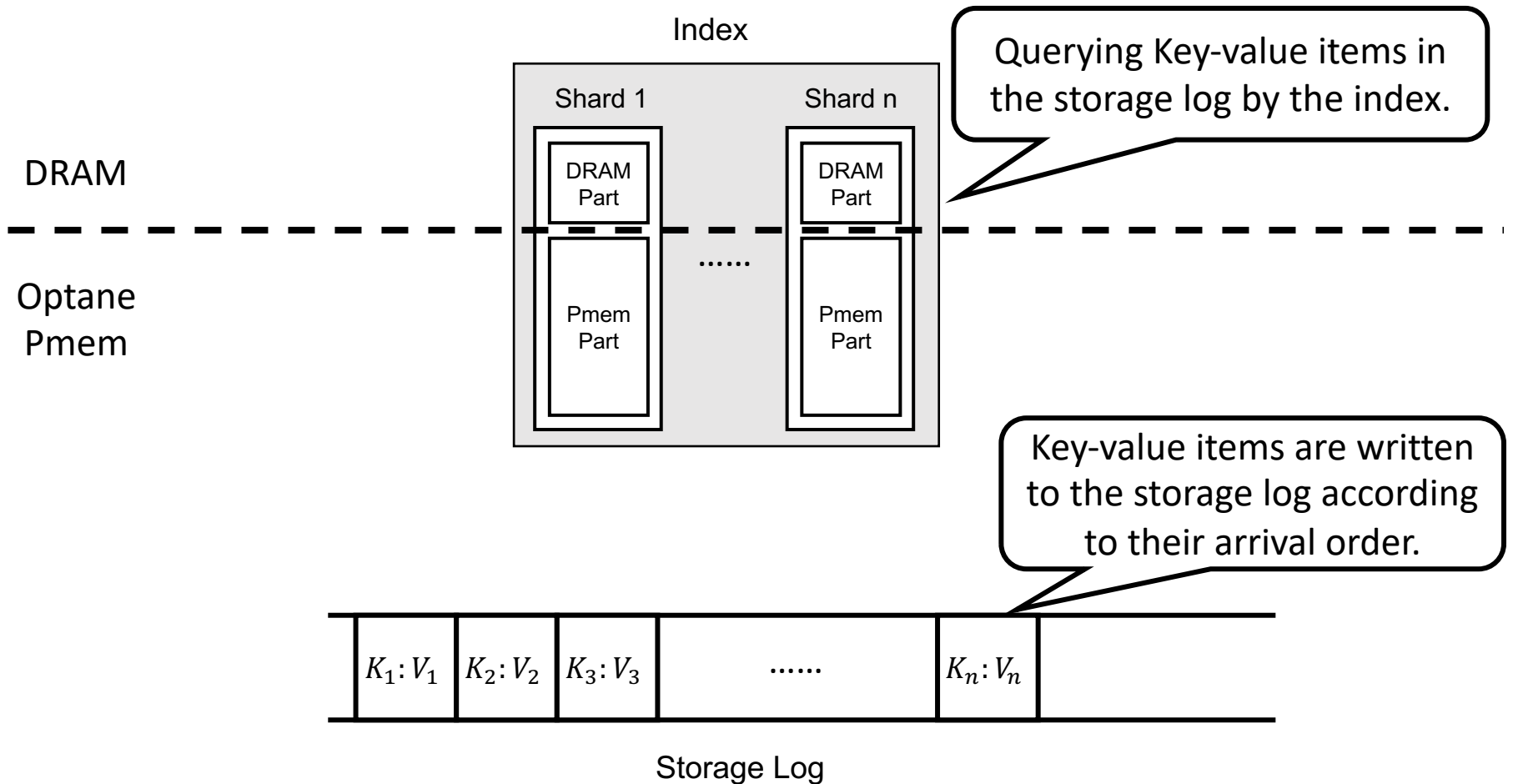**Avoid small random writes to storage.**

However, **DRAM footprint** for the index is considerably **large**, and recovering the index during a **restart** may take an **unacceptable long time**.
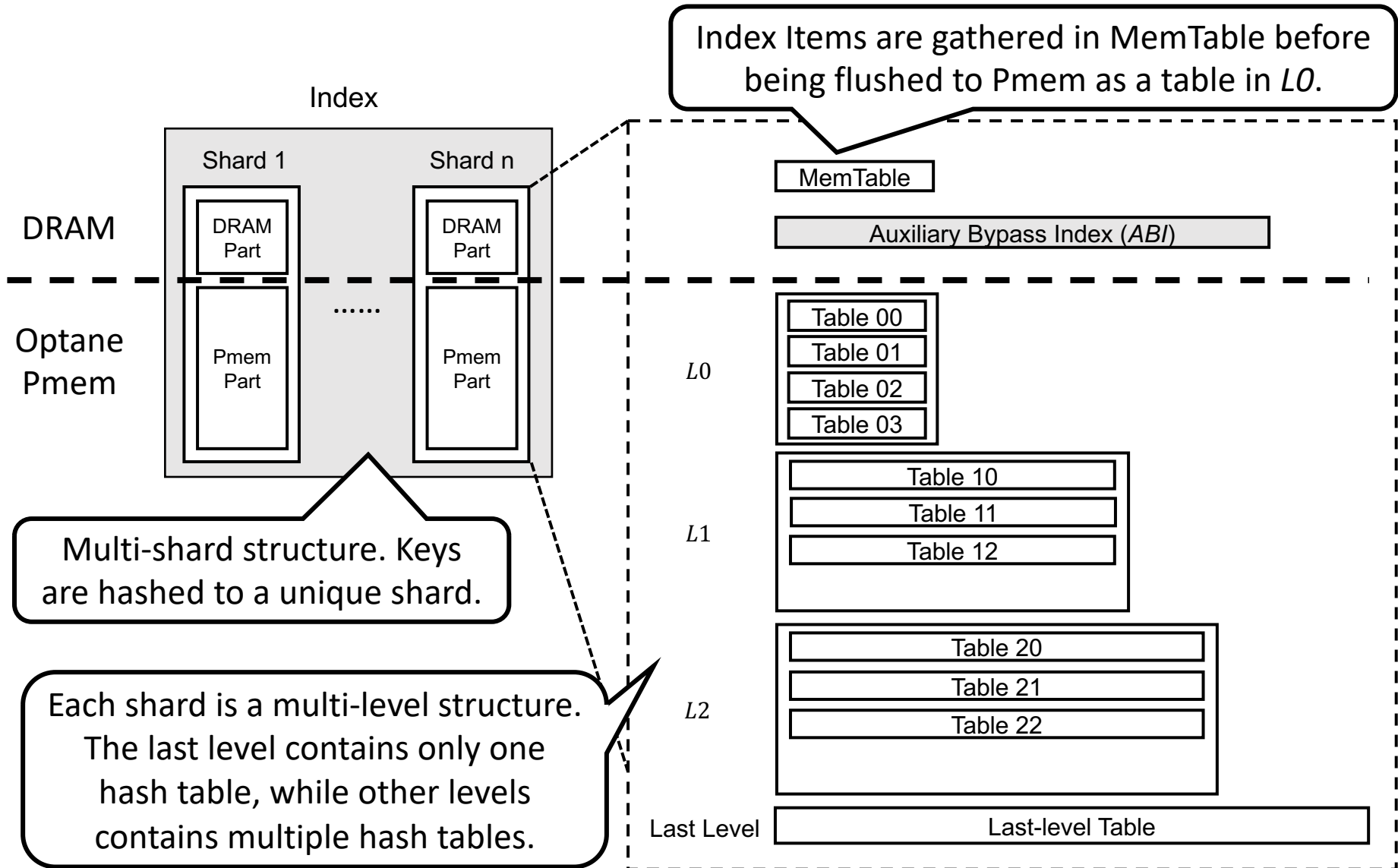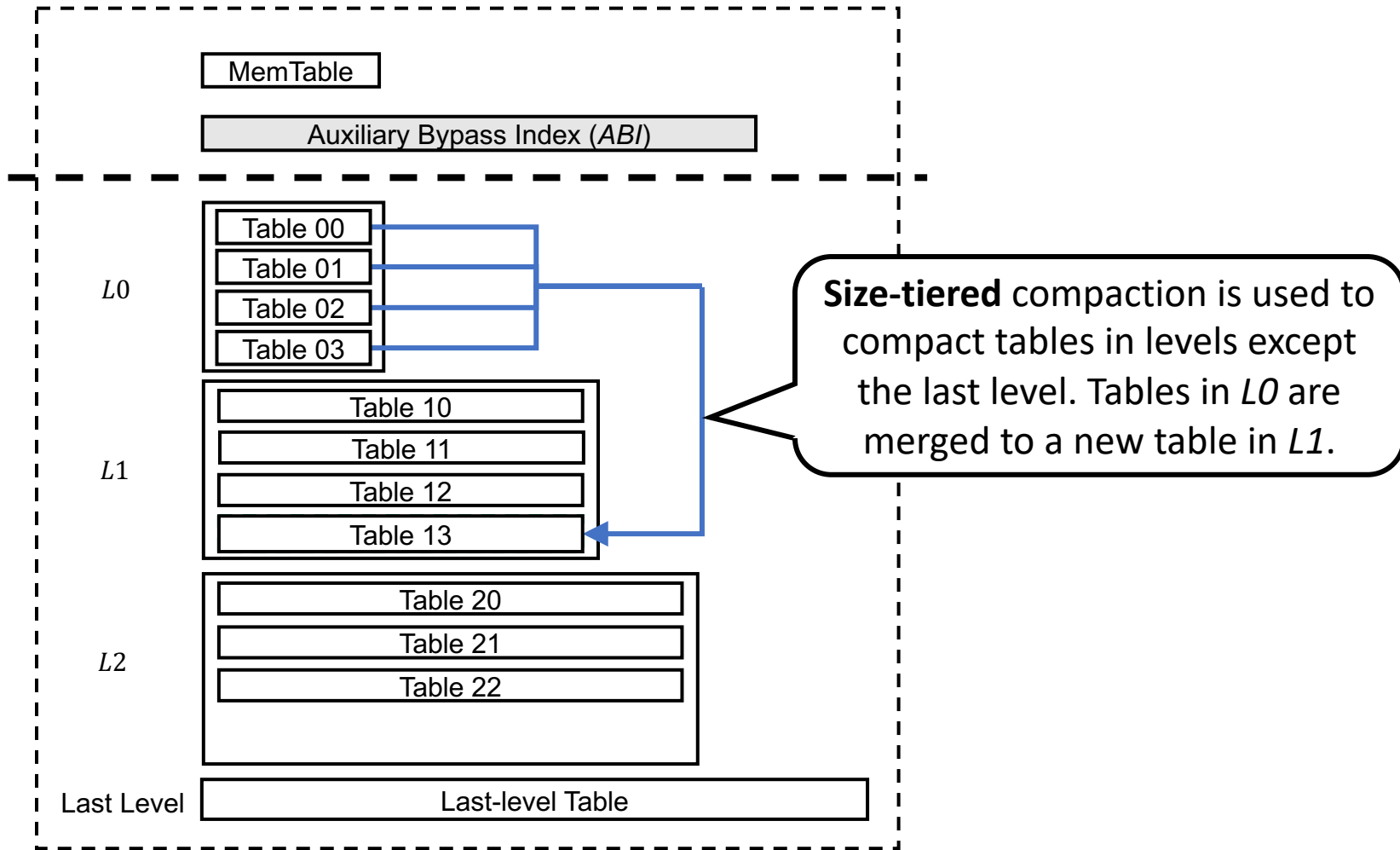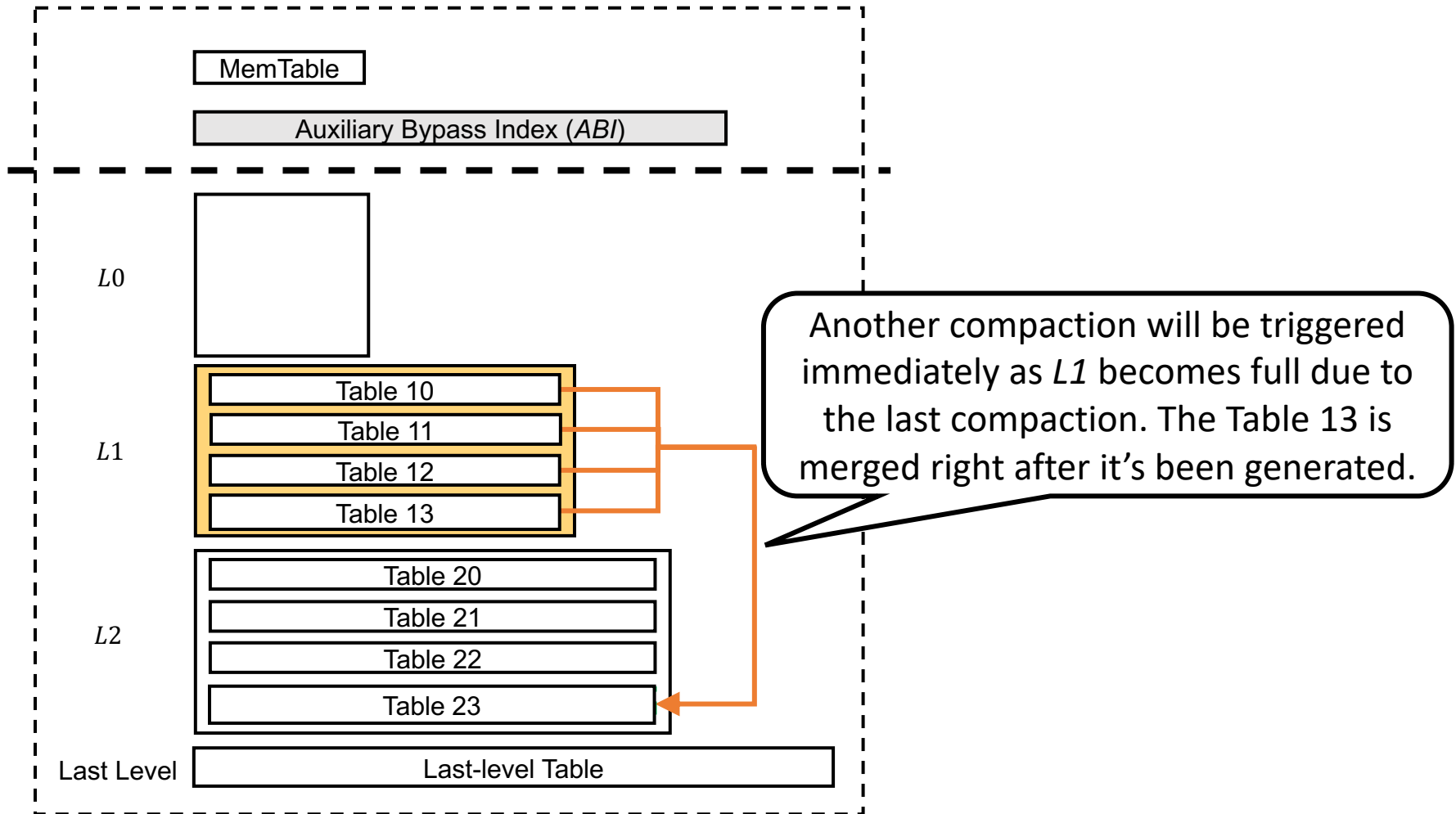
# Design Goals of ChameleonDB

# Structure of ChameleonDB

Index

Shard 1          Shard n

DRAM

DRAM Part      DRAM Part

...... 

Optane Pmem

Pmem Part      Pmem Part

Querying Key-value items in the storage log by the index.

Key-value items are written to the storage log according to their arrival order.

| $K_1:V_1$ | $K_2:V_2$ | $K_3:V_3$ | ...... | $K_n:V_n$ | |

Storage Log

# Structure of ChameleonDB

Index Items are gathered in MemTable before being flushed to Pmem as a table in *L0*.

Index

MemTable

DRAM

Shard 1          Shard n

DRAM Part        DRAM Part

......

Auxiliary Bypass Index (*ABI*)

Optane Pmem

Pmem Part        Pmem Part

*L0*
- Table 00
- Table 01
- Table 02
- Table 03

*L1*
- Table 10
- Table 11
- Table 12

Multi-shard structure. Keys are hashed to a unique shard.

*L2*
- Table 20
- Table 21
- Table 22

Each shard is a multi-level structure. The last level contains only one hash table, while other levels contains multiple hash tables.
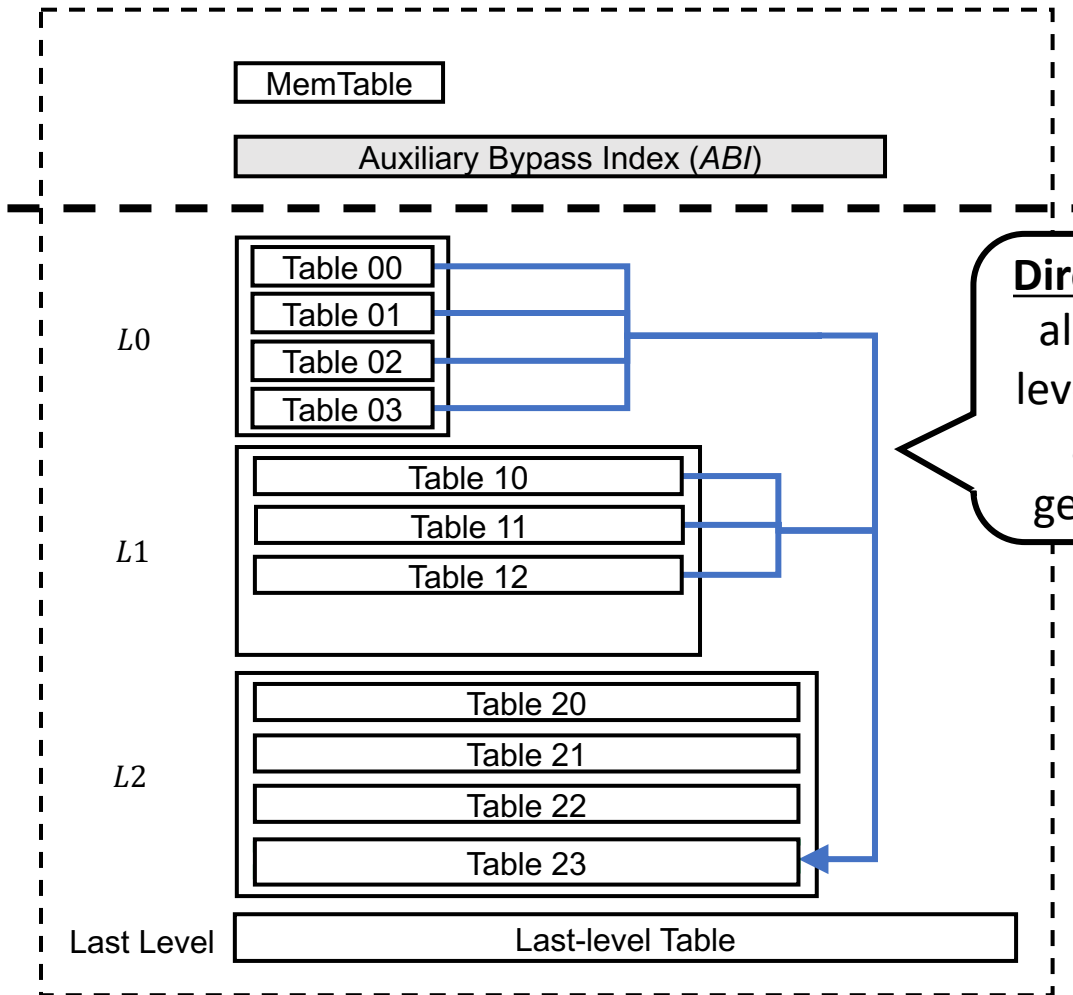
Last Level          Last-level Table

# Size-tiered Compaction in ChameleonDB

# Size-tiered Compaction in ChameleonDB

# Size-tiered Compaction in ChameleonDB

# Leveled Compaction in ChameleonDB

# Search items without Auxiliary Bypass Index

# Search items with Auxiliary Bypass Index

Get(key)

MemTable

Auxiliary Bypass Index (*ABI*)

L0
- Table 00
- Table 01
- Table 02

L1
- Table 10
- Table 11
- Table 12

L2
- Table 20
- Table 21
- Table 22

Last Level — Last-level Table

**ABI** is an in-DRAM hash table that holds all items in the upper levels.

Everytime ChameleonDB flushes a MemTable, it also inserts the items in the MemTable to the *ABI*. *ABI* will be cleared when a last-level compaction happened.

Use **ABI** to bypass checking upper levels one-by-one, and achieve consistently low read latency!

# Last level compacion with Auxiliary Bypass Index

# Recovery during restart



The MemTable should be recovered during a restart.

**DRAM**

MemTable

Auxiliary Bypass Index (*ABI*)

**Optane Pmem**

*L*0

Table 00
Table 01
Table 02

The *ABI* is not necessary to be recovered immediately during a restart as its items can be found by checking upper level tables that are persisted in Pmem.

*L*1

Table 10
Table 11
Table 12

*L*2

Table 20
Table 21
Table 22

Last Level    Last-level Table

# How does ChameleonDB achieve design goals?

➤ Batch KV items before writing them to the storage log, batch index in MemTable before flushing them to Optane Pmem, use multi-level structure to organize index
  $\Rightarrow$ avoid small random writes to Optane Pmem
    $\Rightarrow$ **high write performance**

➤ Use Auxiliary Bypass Index to accelerate Get operation
  $\Rightarrow$ avoid checking multiple levels one-by-one
    $\Rightarrow$ **consistently low read latency**

➤ Place only a portion of the index in DRAM
    $\Rightarrow$ **small DRAM footprint**

➤ Recovers only MemTable during a restart
    $\Rightarrow$ **short restart time**

# Experiment results

# Experiment results

# Experiment results

# Thank You!

wenhui.zhang@uta.edu