



# Tair: A Challenge of Performance in Intel® Optane™ DC Persistent Memory

路青霖

炸酱面组

Team with 李凯 胡伟东

# Tair<sup>[1]</sup>: Challenge of Performance

- High TPS in hotspot items (tens of millions)
  - Hotspot hashing mechanism
  - Limited in **Single node**
- Intel® Optane™ DC Persistent Memory Module
  - Performance<sup>[2]</sup>

Performance	PMEM	DRAM
Latency	Load	305ns/ <b>169ns(Seq)</b>
	Store	<b>94ns</b>
(MT, AVX-512)	Load	6.6GB/s
	Store	2.3GB/s

# Tair: Challenge of Performance

- Key-Value Engine

- Single node with **8G DRAM / 64G AEP** (Optane DC PMEM)

---

Intel(R) Xeon(R) Platinum 8269CY

---

**26 Cores (54 Threads)**

SSE4.2、AVX2、AVX-512

---

DRAM

---

**8GB**

---

Intel(R) Optane(R) DC PMEM

---

**64GB**

---

- Persistence
- Recovery

# Tair: Challenge of Performance

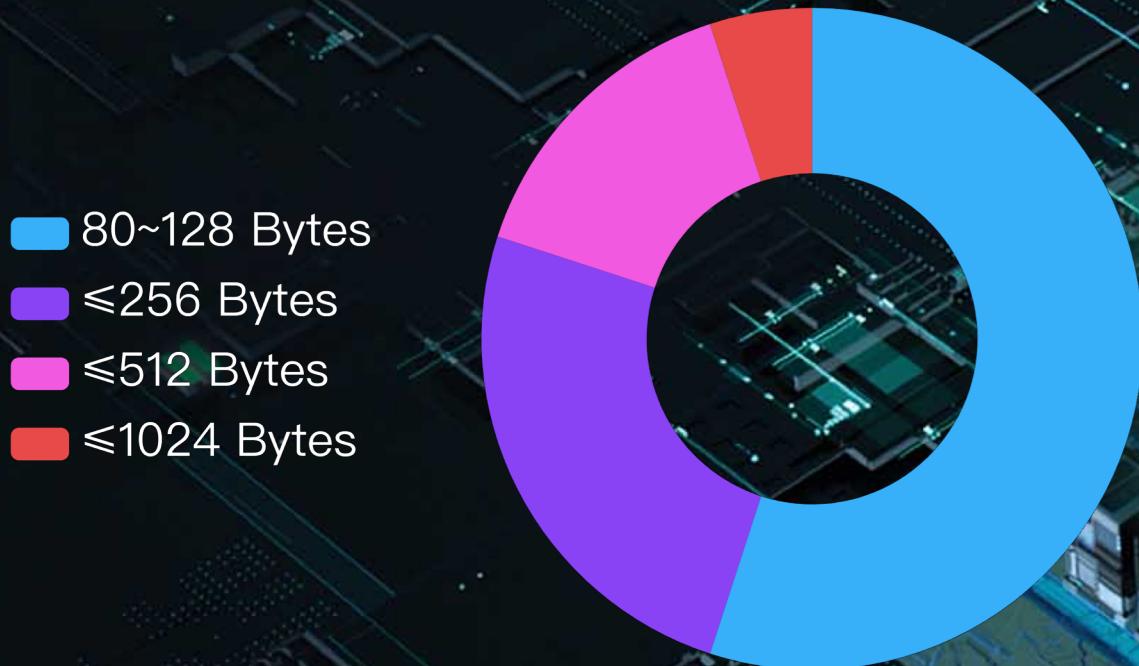
- Judgment
- Data
  - Key (**16B**)
  - Value (**80~1024B**, Max **50GB**)



<b>Correctness</b>	Concurrent Get & Set ( <b>16</b> Threads)
<b>Persistence</b>	Crash-Recover
<b>Performance</b>	<b>24M Set (16 Threads)</b> <b>24M Set (25%) + Get (75%, hotspot 75%)</b>

# Tair: Challenge of Performance

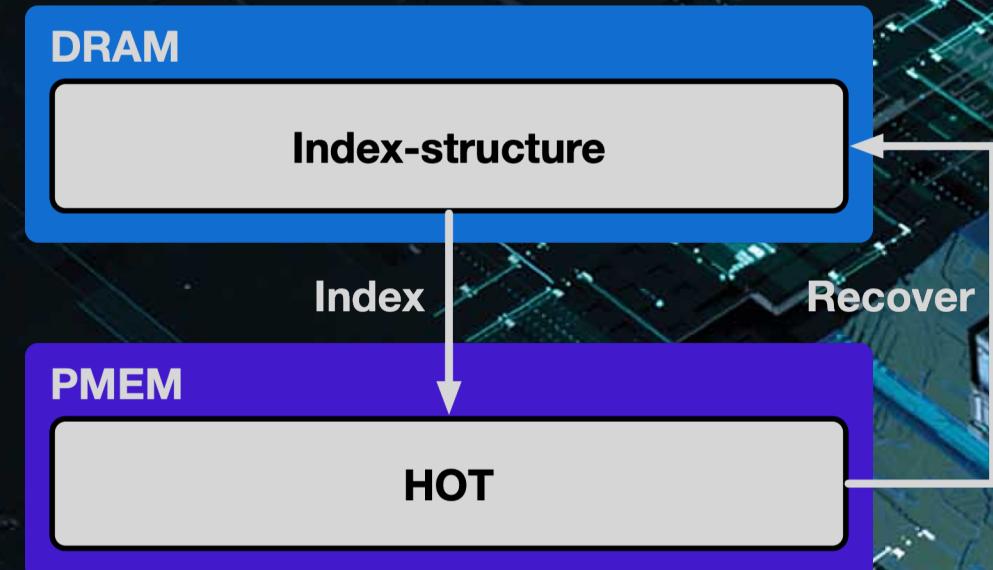
- **Index** under multi-threads
  - Concurrent consistency
  - Upper limit
- **Variable** value length
  - 80~1024 Bytes
  - Persistence
  - Collect & reuse (**25G**)
- Recovery
  - Flush and update
- Unleash potential performance of DC PMEM



# Overview Design

- What does KV-Engine do?
  - $f: Key \rightarrow Value$
- Heap-organized Table (HOT)<sup>[3]</sup>
  - **Unsorted**
- In-memory Index-structure
  - **Speedup** searching in HOT

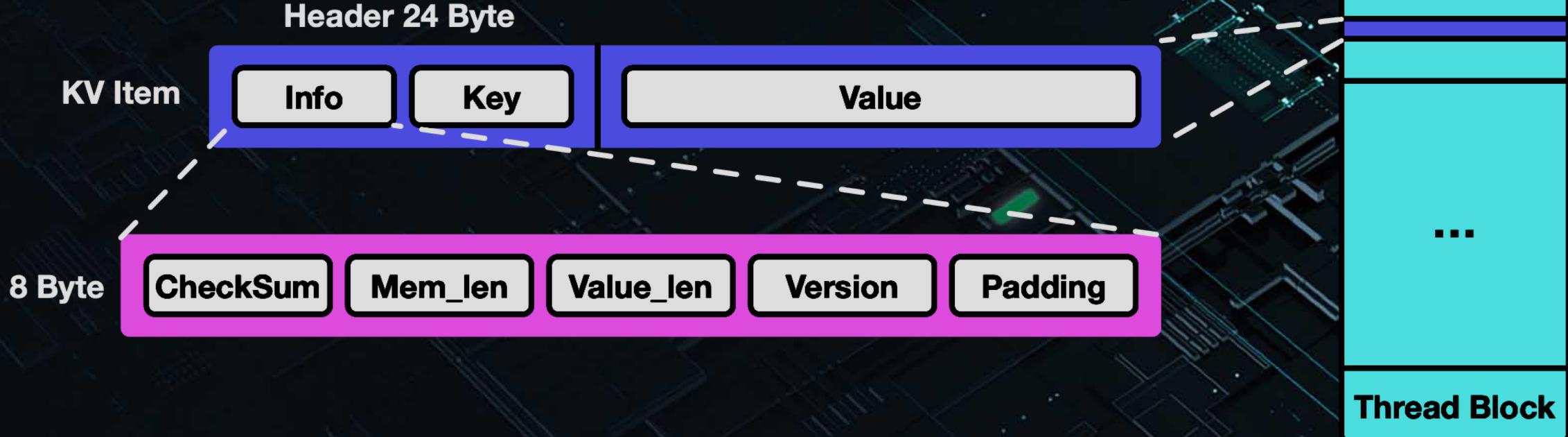
$$\left\{ \begin{array}{l} g: Key \rightarrow Index \\ h: Index \rightarrow Value \end{array} \right. \Rightarrow f(Key) = h(g(Key))$$



# Solution

- HOT – KV Item

- Concurrent / Sequential Write (Get) & Read (Set)
- Support recovery & reuse
- 16 Thread Blocks

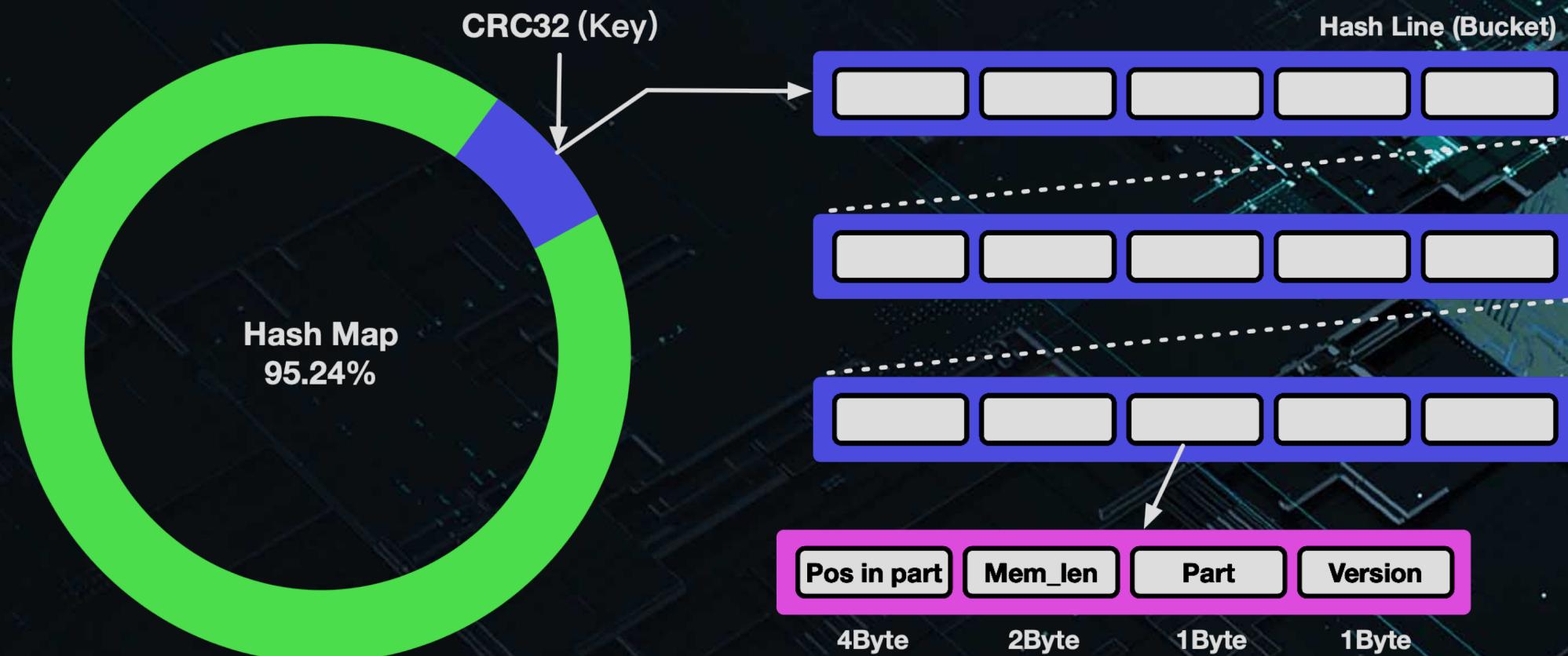


# Solution

- Index-structure - Hash map

➤  $g: Key \rightarrow Index$

$$\text{Redundancy } 5\%, \quad \alpha = \frac{N_{KV}}{N * N_{line}} = \frac{229037390}{48097859 * 5} = 95.24\%$$



# Solution

- Hash Map - Hash Line

- Hash Line  $X$  aligned by Cache Line (64Byte)  $C$ :

$$\text{minimize} \quad 64C - (16 + 8)X - 1$$

$$\text{subject to} \quad 64C - (16 + 8)X \geq 1$$

$$255 \geq X \geq 0$$

$X$  is an integer

$$C \geq 0$$

$C$  is an integer

$$\Rightarrow X = 5, C = 2$$

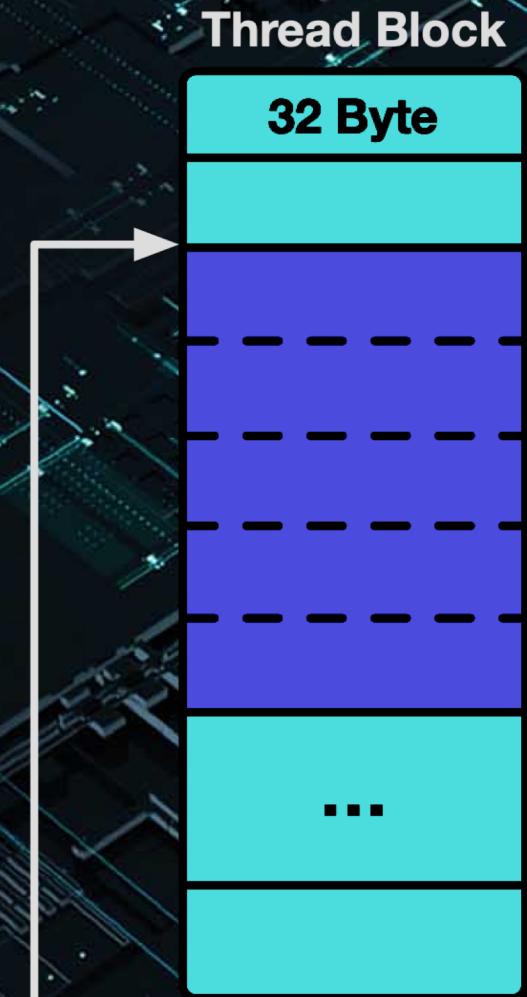
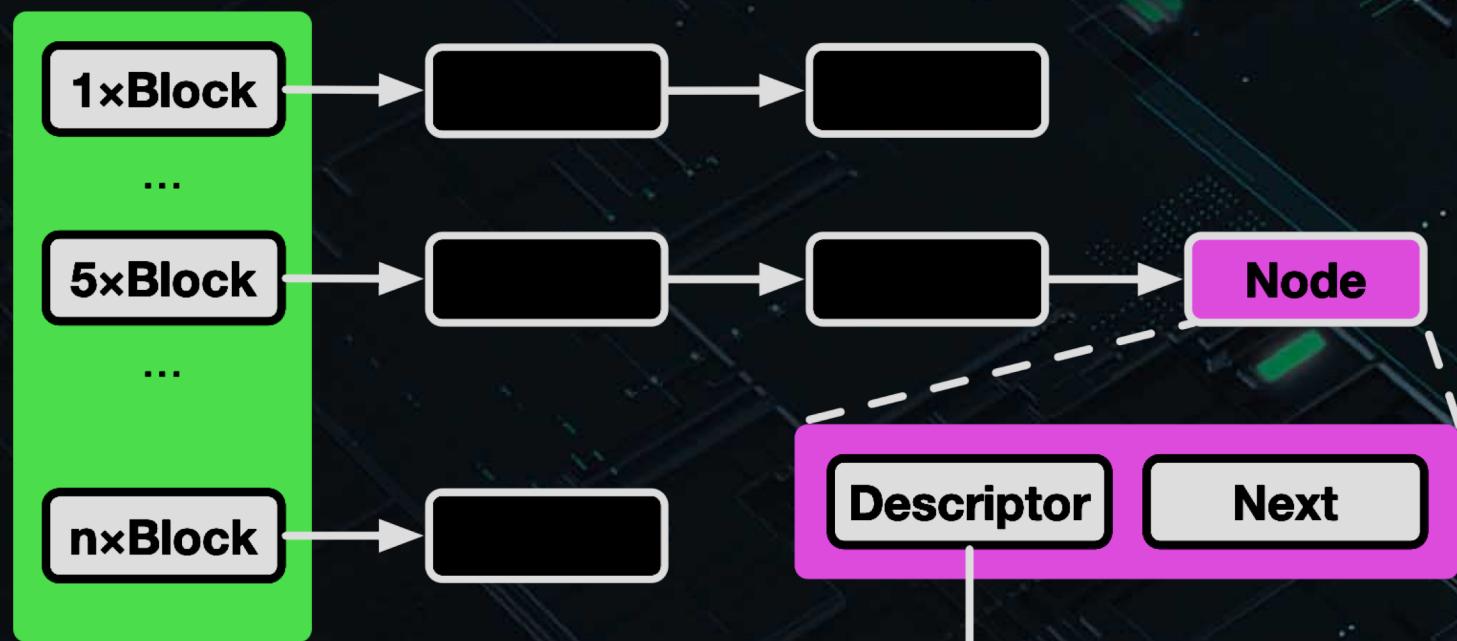
- Lock-free Insert

```
/* 128 Byte total */
HashLine:
    /* 80 Byte */
    Key    keys[5];
    /* 40 Byte */
    uint64 vals[5];
    /* 8 Byte */
    uint8  count;
    uint8  padding[7];
```

# Solution

- Allocator (Thread independent)

- Aligned (32 Byte) memory block
- Free-list (Page Arena)

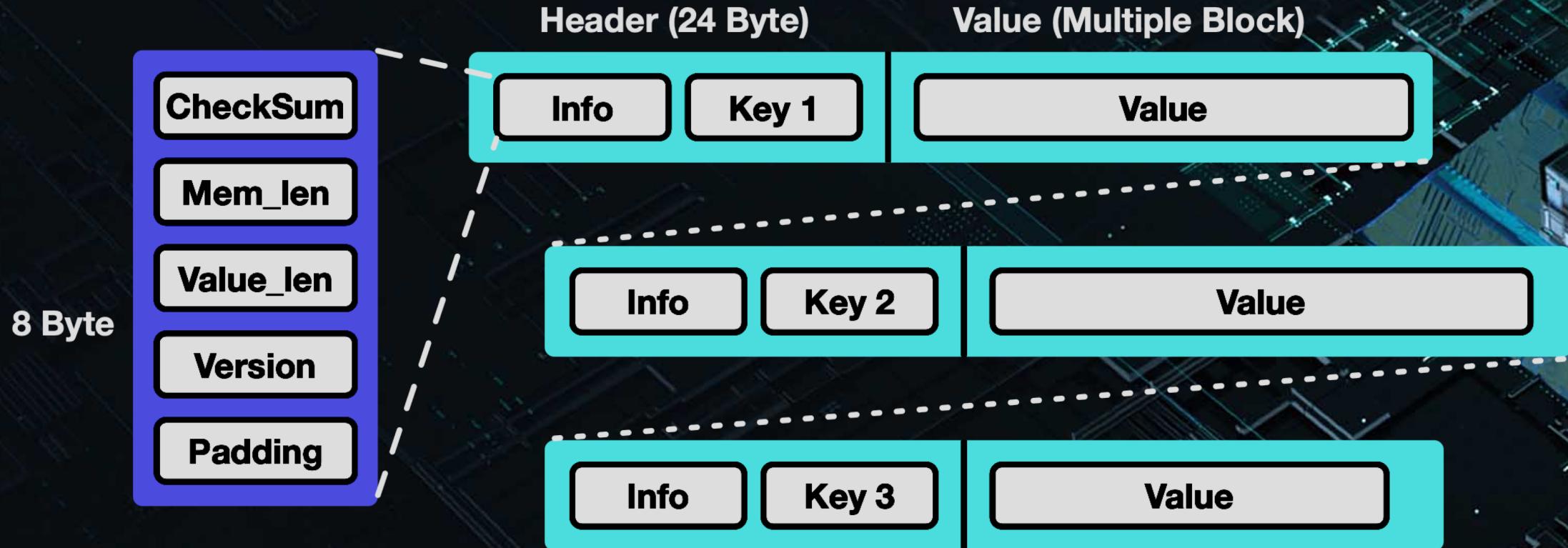


# Solution

- Recovery

- Consistence

- Atomic Write Unit (8 Byte)<sup>[3]</sup>

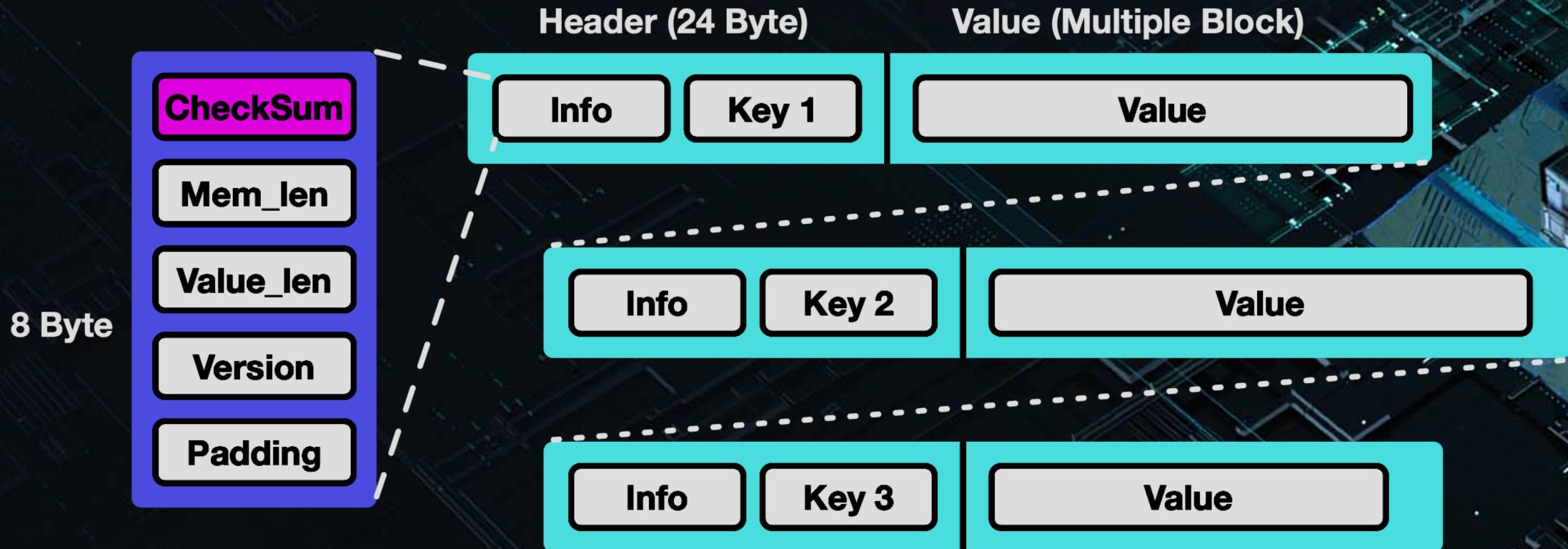


# Solution

- Recovery

- Consistence

- Atomic Write Unit (8 Byte)<sup>[3]</sup>



# Optimization

Score: 55

- False sharing
  - ThreadLocal - allocator offset & size
- Warm up
  - Page fault
  - pmem\_memset
  - <160s
- Get thread ID
  - Individual
  - Sequential

# Optimization

Score: 55

- False sharing
  - ThreadLocal - allocator offset & size
- Warm up
  - Page fault
  - pmem\_memset
  - <160s
- Get thread ID
  - Individual
  - Sequential

# Optimization

Score: 44

- False sharing
  - ThreadLocal - allocator offset & size
- Warm up
  - Page fault
  - pmem\_memset
  - <160s
- Get thread ID
  - Individual
  - Sequential

# Optimization

Score: 39

- False sharing
  - ThreadLocal - allocator offset & size
- Warm up
  - Page fault
  - pmem\_memset
  - <160s
- Get **thread ID**
  - Individual
  - Sequential

```
int64_t gettno() {
    static __thread int64_t tno = -1;
    static int64_t tno_seq = 0;
    if (UNLIKELY(-1 == tno)) {
        tno = __sync_fetch_and_add(&tno_seq, 1);
    }
    return tno;
}
```



# The Final Seconds

- Which part takes the most of time?

Score: 39

# The Final Seconds

- Which part takes the most of time?

Score: 39

➤ Hash Map

- Get
- Occupy
- Replace

➤ Heap Table

➤ Allocator

```
hash = CRC32(key);
for (i = 0; i < HASH_LINE_CNT; ++i) {

    /* Locate hash line */
    index = (hash + i) % HASH_LINE_CNT;

    /* find key in current hash line */
    hash_line[index].find(key);
    ...
}
```

Mod operation

Keys Comparing

# The Final Seconds

- Which part takes the most of time?

Score: 38

➤ Hash Map

- Get
- Occupy
- Replace

➤ Heap Table

➤ Allocator

```
index = CRC32(key) % HASH_LINE_CNT;
for (i = 0; i < HASH_LINE_CNT; ++i) {
    /* find key in current hash line */
    hash_line[index].find(key);
    ...
    /* next hash line */
    index++;
    index = UNLIKELY(index == HASH_LINE_CNT) ? 0 : index;
}

/* prefetch next keys to L0 cache */
_mm_prefetch(hash_line[index].keys, _MM_HINT_T0);
}
```

# The Final Seconds

- Which part takes the most of time?
  - Hash Map
  - **Heap Table**
    - Get
    - Put
  - Allocator

Score: 38

- MEMCPY
  - CRC of KV Item
- 
- rte\_memcpy
  - XXHASH3<sup>[4]</sup> (AVX2)

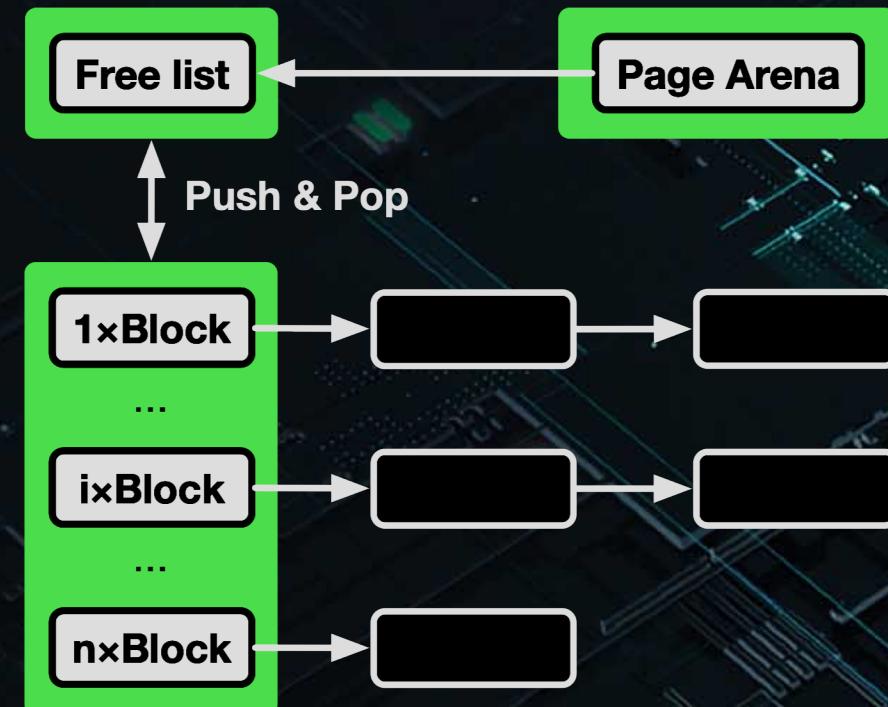
Hash Name	Width	Large Data Bandwidth	Small Data Velocity
XXH3	64	59.4 GB/s	133.1
XXH128	128	57.9 GB/s	118.1
CRC32C (hw)	32	13.0 GB/s	57.9

[4] Collet, Yann. "xxHash: Extremely Fast Hash Algorithm." (2016).

# The Final Seconds

- Which part takes the most of time?
  - Hash Map
  - Heap Table
  - Allocator

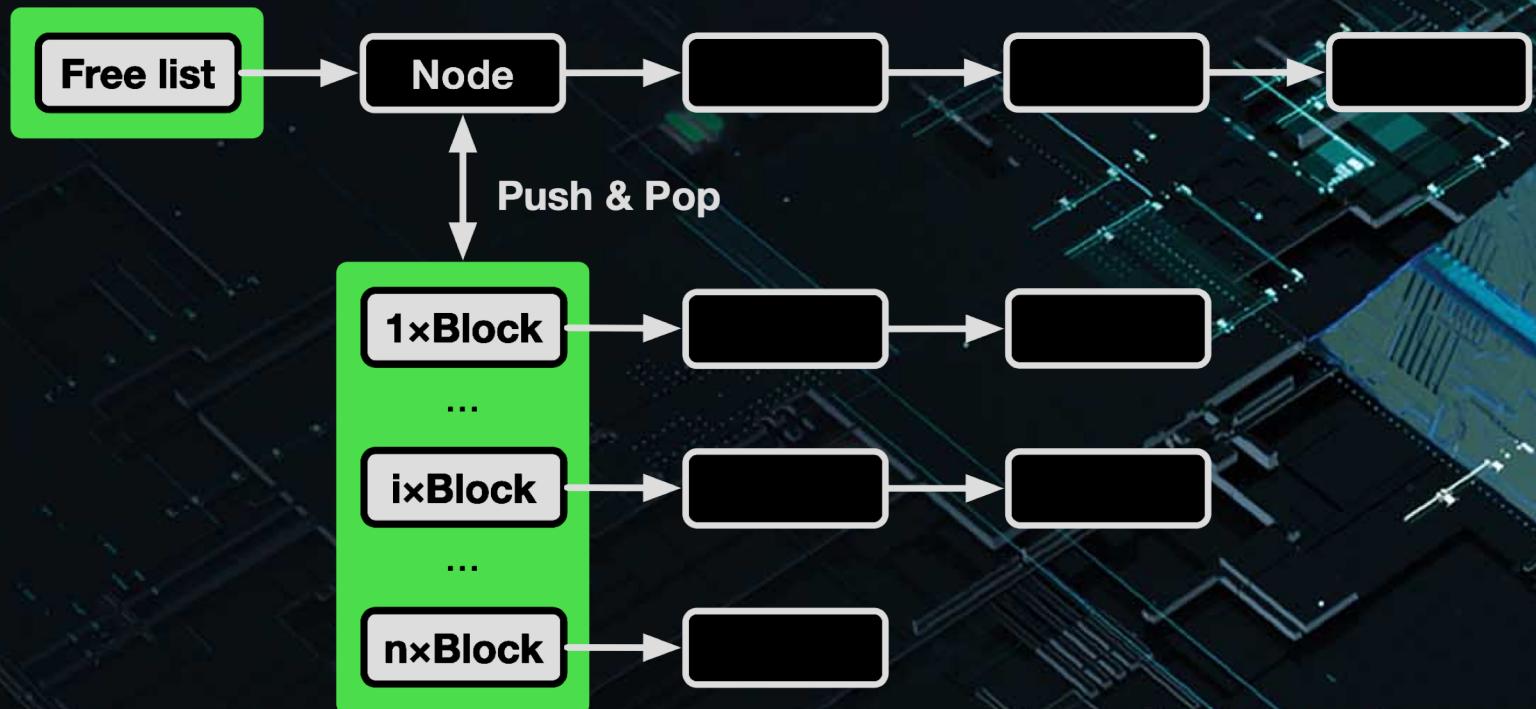
Score: 36



# The Final Seconds

- Which part takes the most of time?
  - Hash Map
  - Heap Table
  - Allocator

Score: 35.5



# Conclusion

- Hash Map (DRAM) + HOT (DC PMEM)
  - Thread blocks
  - Allocator
- Log & Profiling
- Value Cache
  - **Simplify** or not?
- 16 Threads => Thread pool
- **NOT** rely on the distribution of data

Score: 35.5 (25.5 + 10)



# Q & A

# Reference

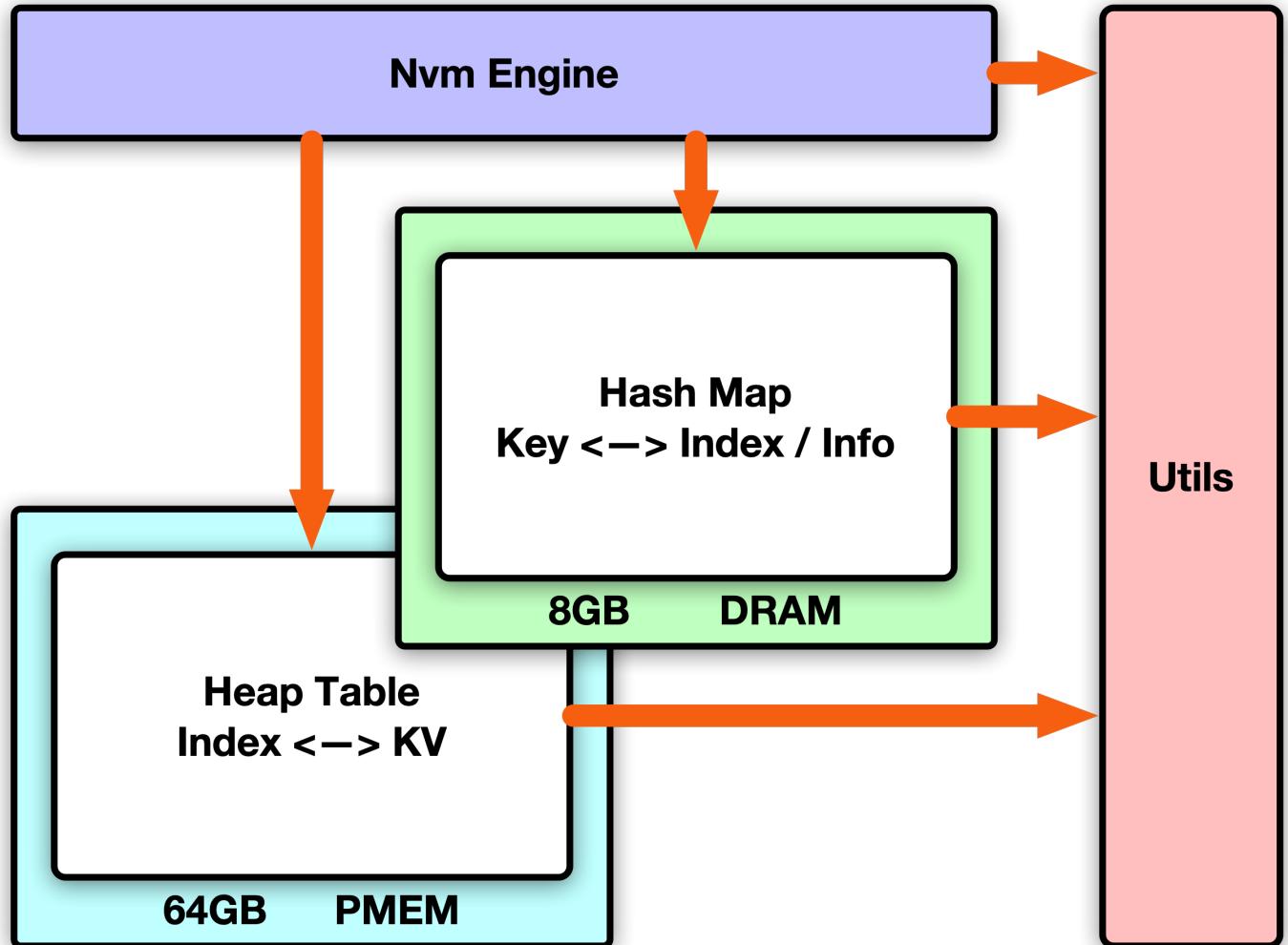
- [1] Duo long, Dao an, et al. 2017. TAIR, A distributed key-value storage system developed by Alibaba Group. <https://github.com/alibaba/tair>.
- [2] Izraelevitz, Joseph, et al. "Basic performance measurements of the intel optane DC persistent memory module." *arXiv preprint arXiv:1903.05714* (2019).
- [3] Islam, A. A. R., Narayanan, A., York, C., & Dai, D. A Performance Study of Optane Persistent Memory: From Indexing Data Structures' Perspective.
- [4] Collet, Yann. "xxHash: Extremely Fast Hash Algorithm." (2016).
- [5] Chen, Jiqiang, et al. "HotRing: A Hotspot-Aware In-Memory Key-Value Store." *18th {USENIX} Conference on File and Storage Technologies ({FAST} 20)*. 2020.

# Reference

- [1] Duo long, Dao an, et al. 2017. TAIR, A distributed key-value storage system developed by Alibaba Group. <https://github.com/alibaba/tair>.
- [2] Izraelevitz, Joseph, et al. "Basic performance measurements of the intel optane DC persistent memory module." *arXiv preprint arXiv:1903.05714* (2019).
- [3] Islam, A. A. R., Narayanan, A., York, C., & Dai, D. A Performance Study of Optane Persistent Memory: From Indexing Data Structures' Perspective.
- [4] Collet, Yann. "xxHash: Extremely Fast Hash Algorithm." (2016).
- [5] Chen, Jiqiang, et al. "HotRing: A Hotspot-Aware In-Memory Key-Value Store." *18th {USENIX} Conference on File and Storage Technologies ({FAST} 20)*. 2020.

# Solution

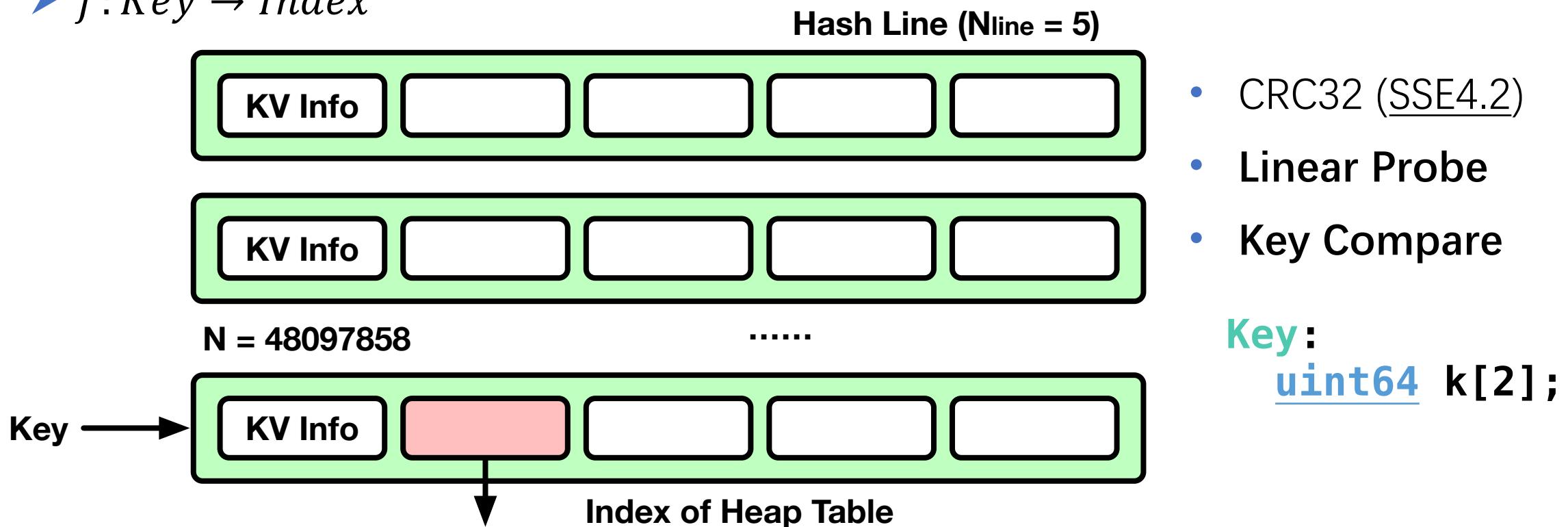
- NvmEngine
  - Hash Map (5.73GB)
    - Key to Index
  - Heap Table (64GB)
    - Index to KV
    - Recovery
  - Utils
    - Log Module
    - Profiling Module
    - System Call



# Solution

- Hash Map

➤  $f: Key \rightarrow Index$



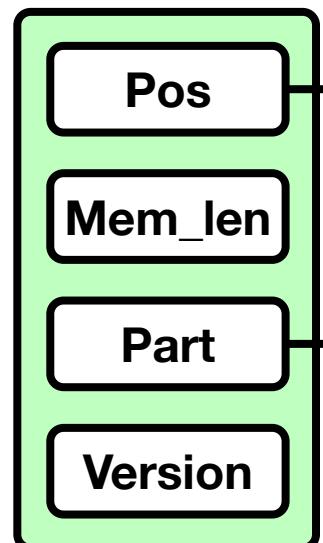
➤ Redundancy 5%,  $\alpha = \frac{N_{KV}}{N * N_{line}} = \frac{229037390}{48097859 * 5} = 95.24\%$

# Solution

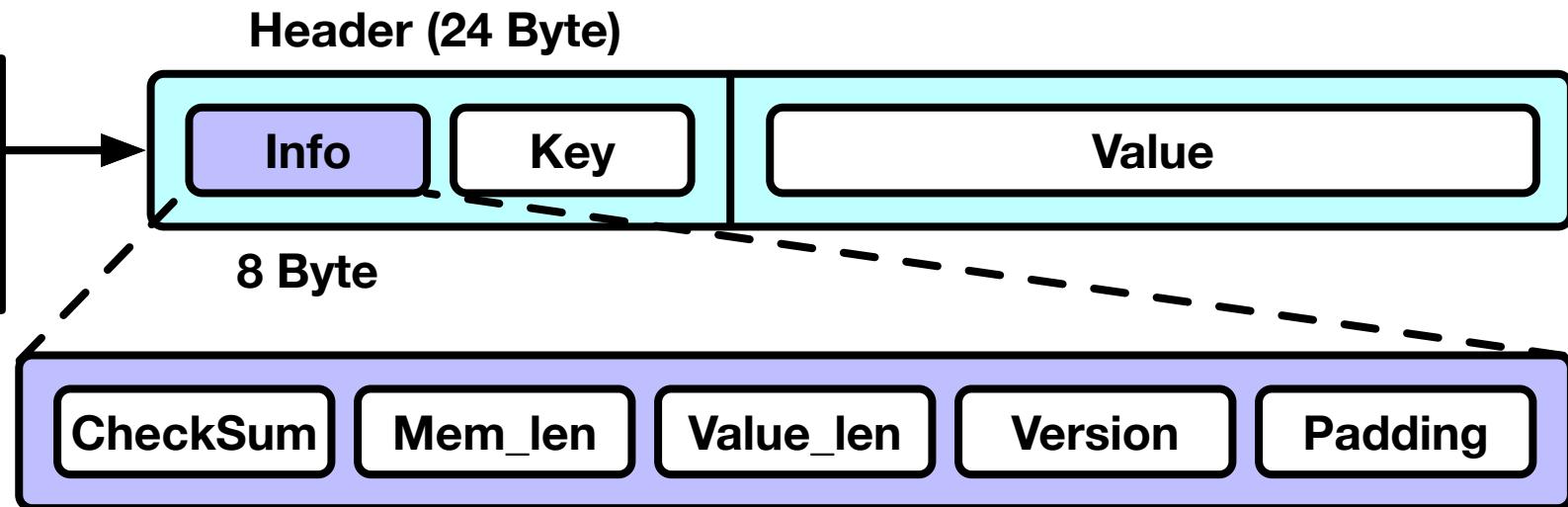
- **Heap Table**

- Concurrent / Sequential Write (**Get**) & Read (**Set**)
- Recovery
- Memory Reuse (**Allocator**)

**Index (8 Byte)**



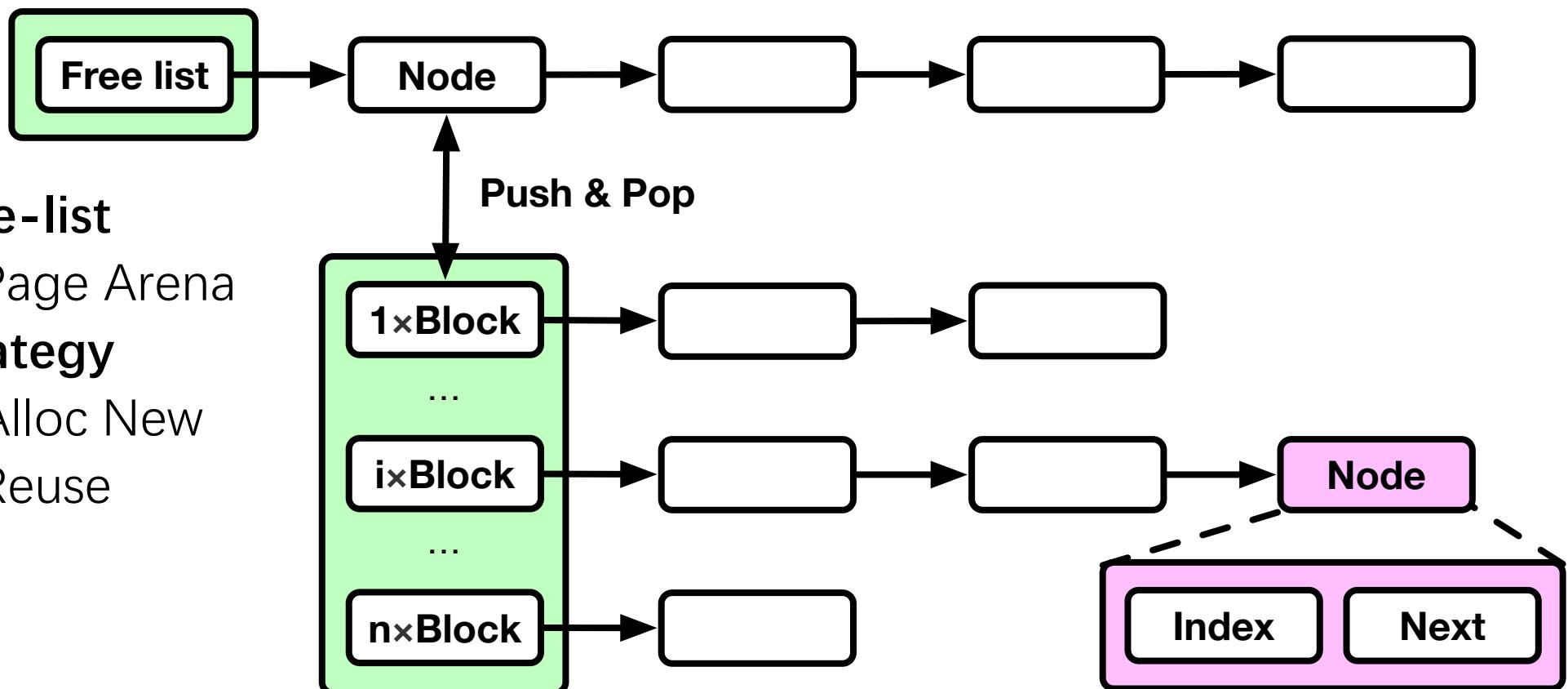
**KV Item**



# Solution

- **Allocator (Thread independent)**

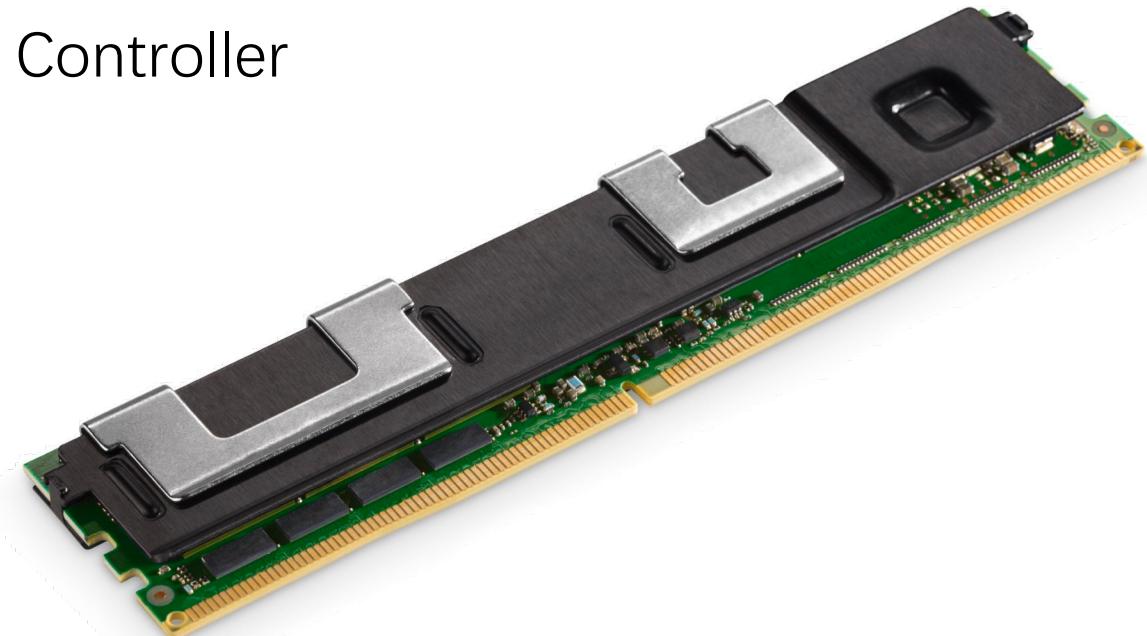
- Aligned (32 Byte) memory block



# What exactly is Optane DC PMEM?

- **Intel® Optane™ DC Persistent Memory Module**

- Capacity & Persistence
- Byte-granularity
- DIMM\* (DDR4) to CPU's Memory Controller
- Memory Mode
  - DRAM Cache
- **App-Direct Mode**
  - Persistent Device / Memory
  - File System
  - **Ext4 DAX**(Direct Access)
  - **User-space Persistent**



*DIMM: Dual In-line Memory Module*

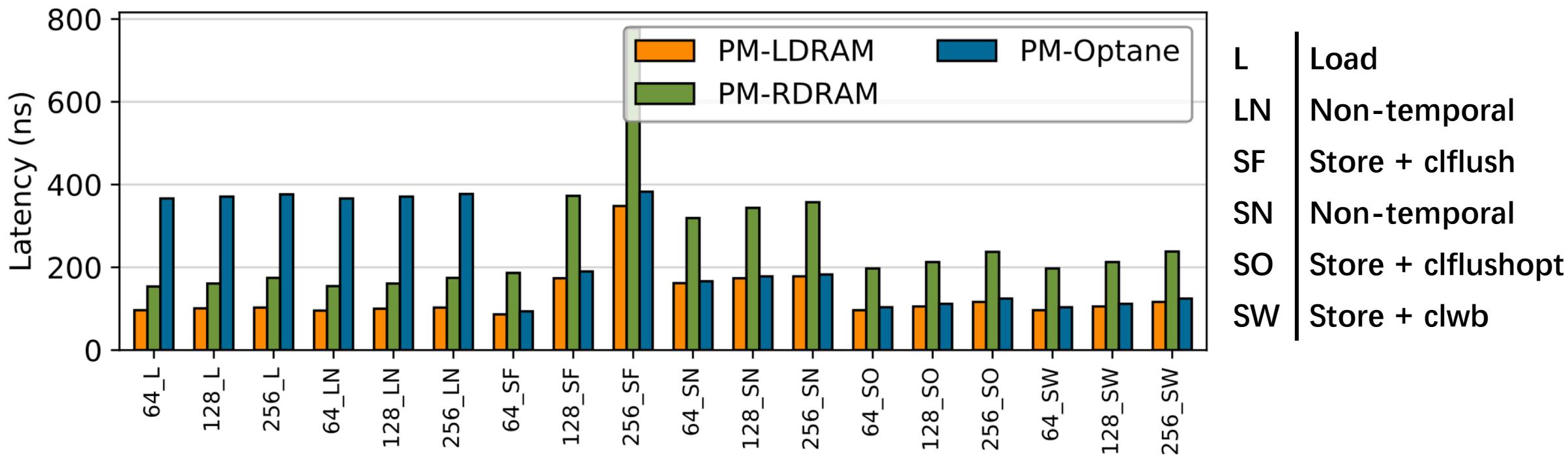
[2] Izraelevitz, Joseph, et al. "Basic performance measurements of the intel optane DC persistent memory module." *arXiv preprint arXiv:1903.05714* (2019).

# Intel® Optane™ DC Persistent Memory Module

- Performance<sup>[2]</sup>

- Load latency    **305ns** (Rand) / **169ns** (Seq)
- Store latency    **94ns**

DRAM **81ns** (Rand)  
DRAM **86ns**

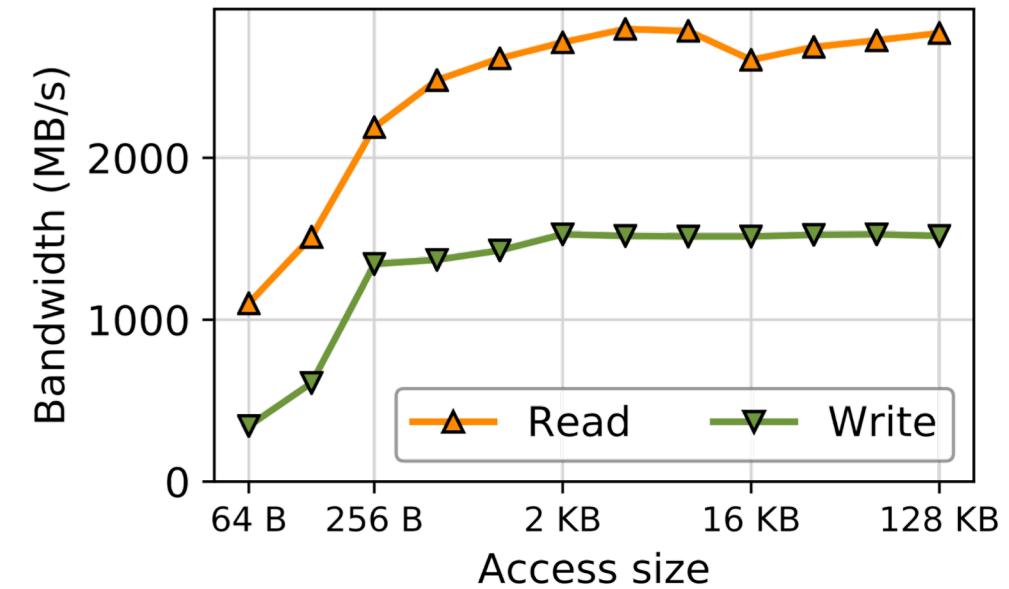
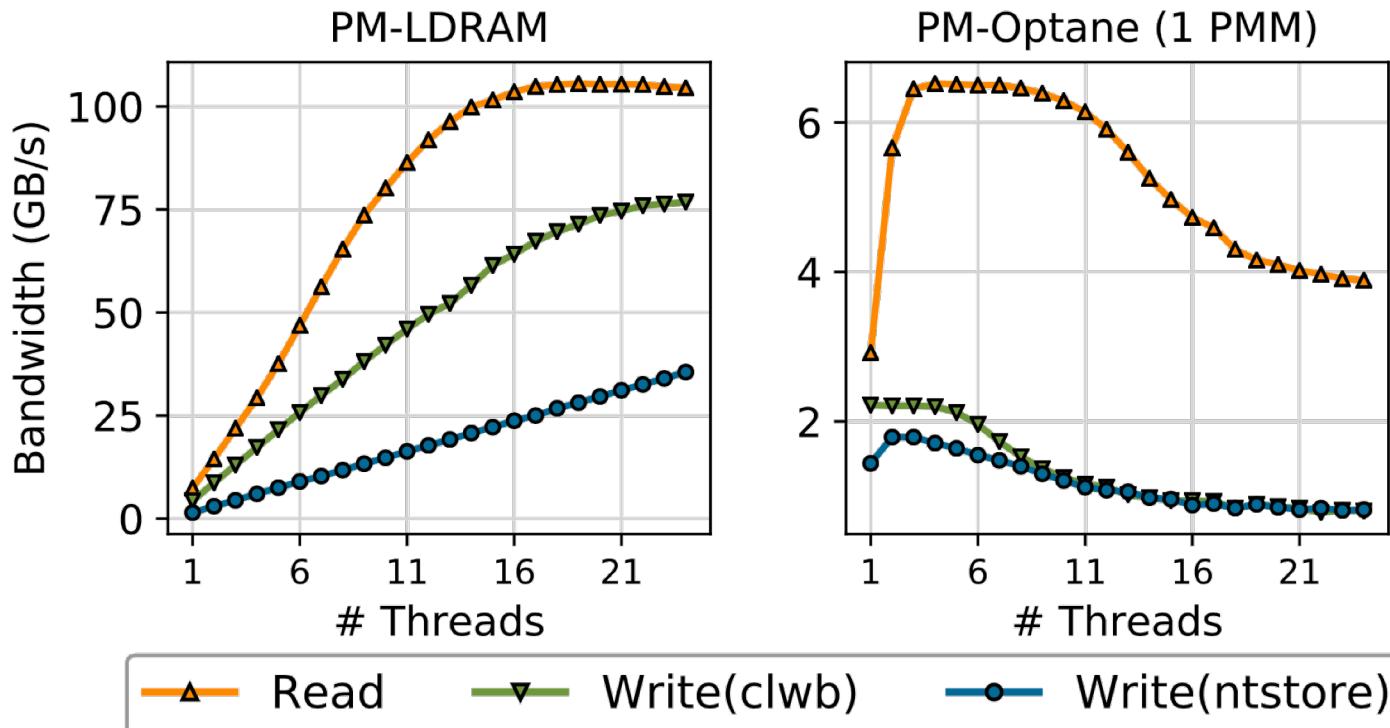


[2] Izraelevitz, Joseph, et al. "Basic performance measurements of the intel optane DC persistent memory module." *arXiv preprint arXiv:1903.05714* (2019).

# Intel® Optane™ DC Persistent Memory Module

- Performance<sup>[2]</sup>

- Max Bandwidth **6.6 GB/s** Read, **2.3 GB/s** Write (Multi-thread, **AVX-512**)



[2] Izraelevitz, Joseph, et al. "Basic performance measurements of the intel optane DC persistent memory module." *arXiv preprint arXiv:1903.05714* (2019).

# xxHash

- Extremely Fast Hash Algorithm<sup>[4]</sup>
  - Completed the **SMHasher** test suite

Hash Name	Width	Large Data Bandwidth	Small Data Velocity	Comment
XXH3	64	59.4 GB/s	133.1	AVX2 (optional)
XXH128	128	57.9 GB/s	118.1	AVX2 (optional)
<i>RAM Seq read</i>	N/A	28.0 GB/s	N/A	<i>for reference</i>
CRC32C (hw)	32	13.0 GB/s	57.9	SSE4.2
Murmur3	32	3.9 GB/s	56.1	
FNV64	64	1.2 GB/s	62.7	Poor avalanche properties
SHA1	160	0.8 GB/s	5.6	Cryptographic but broken
MD5	128	0.6 GB/s	7.8	Cryptographic but broken

[4] Collet, Yann. "xxHash: Extremely Fast Hash Algorithm." (2016).