

# ENDURANCE BENEFITS WITH PERSISTENT MEMORY

---



Sergey Vinogradov

*Software Development Engineer  
Intel*

Storage Performance Development Kit (SPDK)  
Persistent Memory Development Kit (PMDK)  
Intel® VTune™ Profiler

**Virtual Forum**

# AGENDA

1

Introduction

2

Endurance characteristics

3

Write Amplification in software

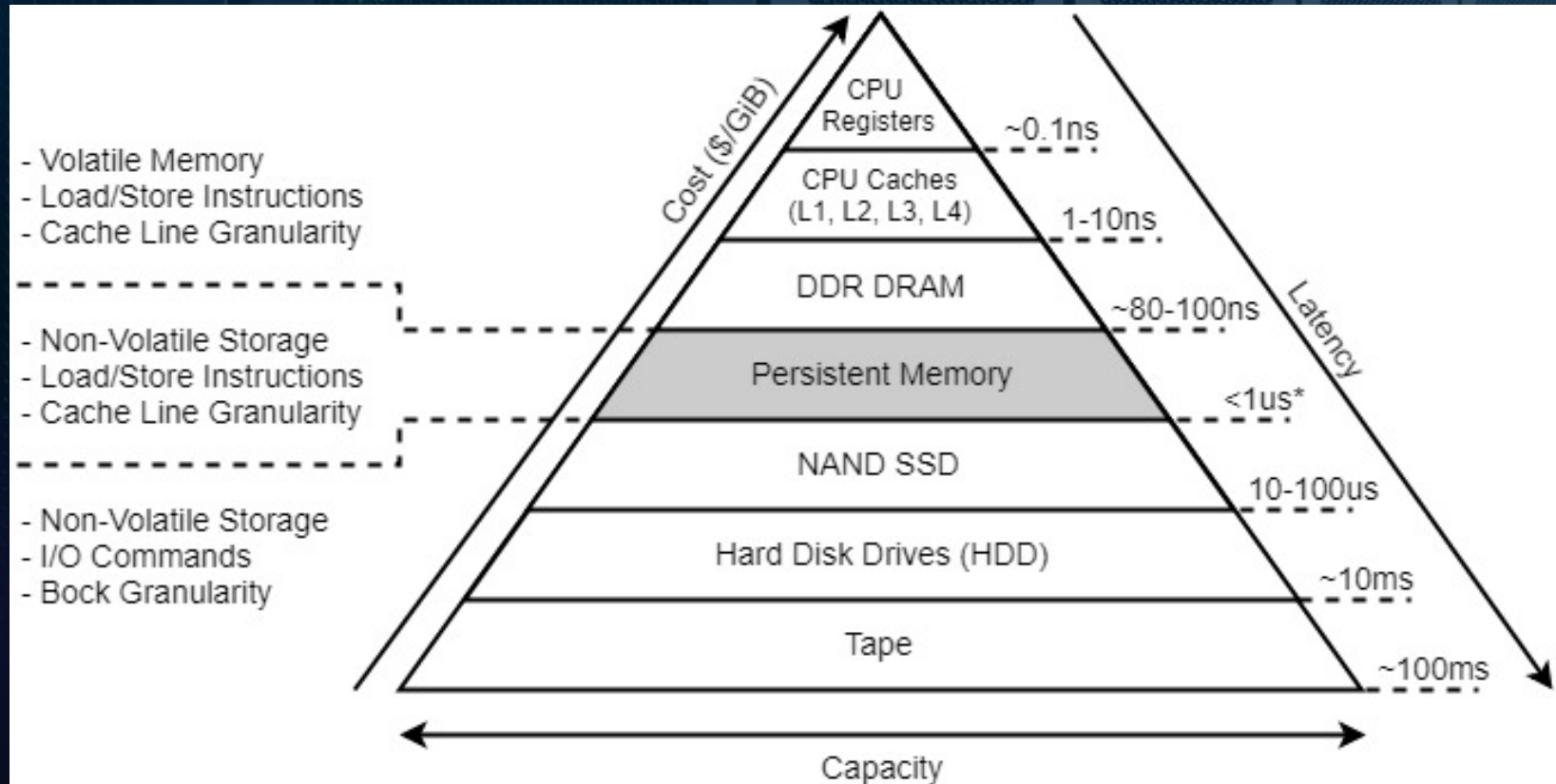
*Software modernization with Persistent memory*

4

Summary

# INTRODUCTION

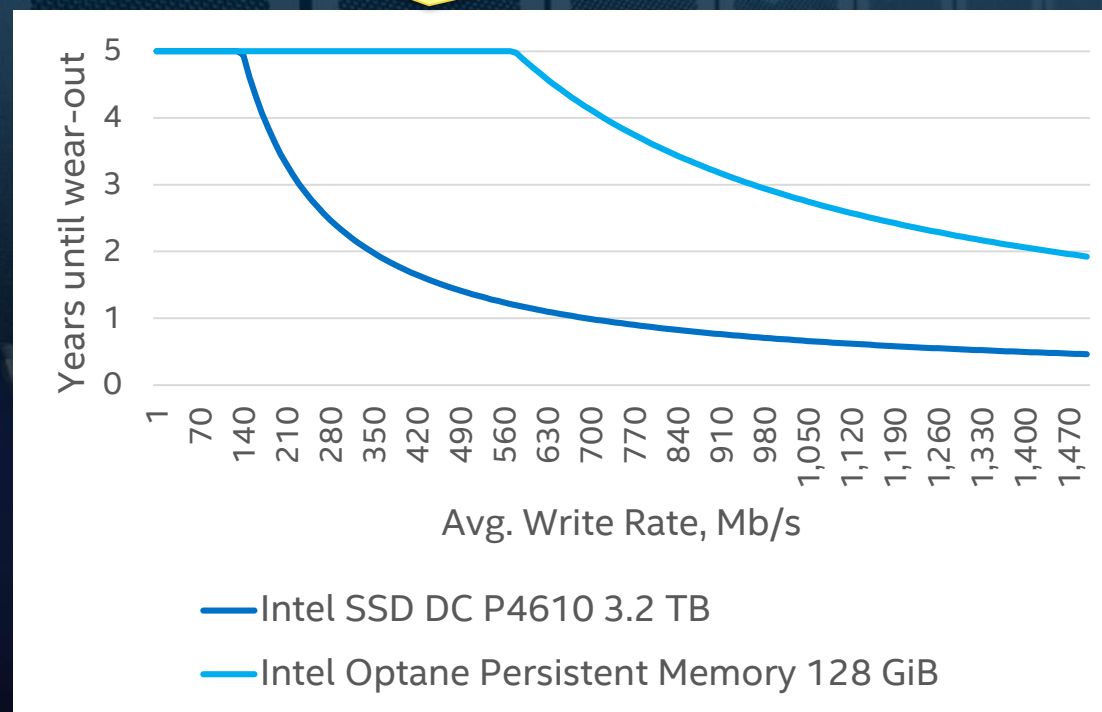
- Intel® Optane™ Persistent Memory has better performance characteristics comparing to NAND SDD. But... **What about endurance?**



# ENDURANCE CHARACTERISTICS OF PERSISTENT MEMORY

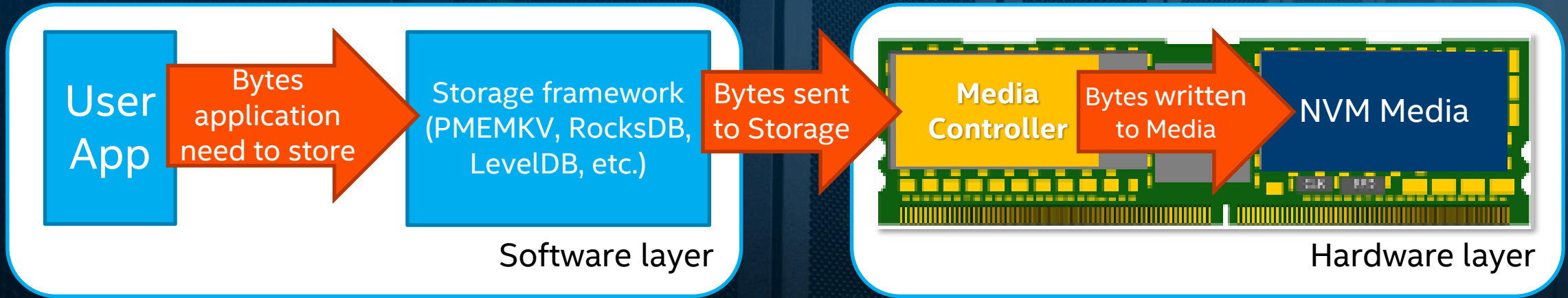
Persistent Memory survives better on write-intensive workloads.

| Params                 | SSD DC P4610 3.2 Tb | Optane Persistent Memory 128 Gib |             |
|------------------------|---------------------|----------------------------------|-------------|
| size                   | 3.2 Tb              | 128 Gb                           |             |
| Warranty               | 5 years             | 5 years                          |             |
| Endurance, 100% Writes | 21.85               | 64B                              | 256B        |
|                        |                     | 91 PBW                           | 292 PBW     |
| DWPD                   | 3.74                | 389.55                           | 1250        |
| Max Avg. Write Rate    | 138.5 MB/sec        | 577 MB/sec                       | 1851 MB/sec |



\* For Intel Optane Persistent Memory we take the worst case Endurance

# DEFINITION OF WRITE AMPLIFICATION (WA)



- Software write amplification caused by SW
  - Algorithm may duplicates data writes
  - Write granularity may cause Write Amplification
- Hardware write amplification caused by HW/FW

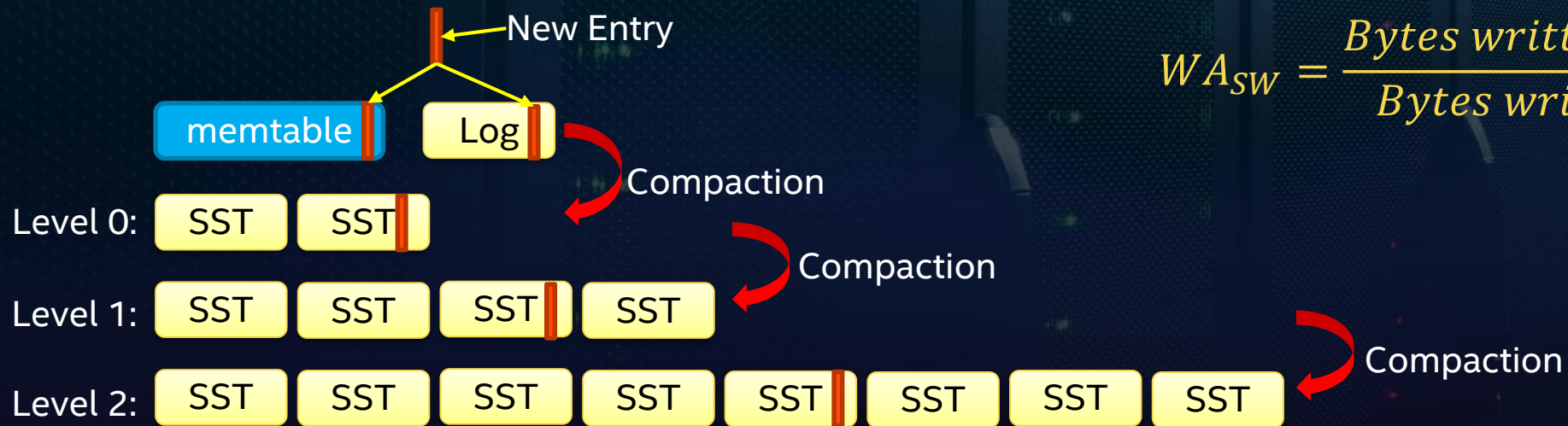
$$WA_{SW} = \frac{\text{Bytes sent to storage}}{\text{Bytes app want to store}}$$

$$WA_{HW} = \frac{\text{Bytes written to media}}{\text{Bytes sent to storage}}$$

$$WA = WA_{SW} \times WA_{HW}$$

# SOFTWARE WRITE AMPLIFICATION IN LSM TREE

- Log Structured Merge (LSM) tree is a key/value storage algorithm
  - Designed for block storage device
- Compaction causes Software Write Amplification

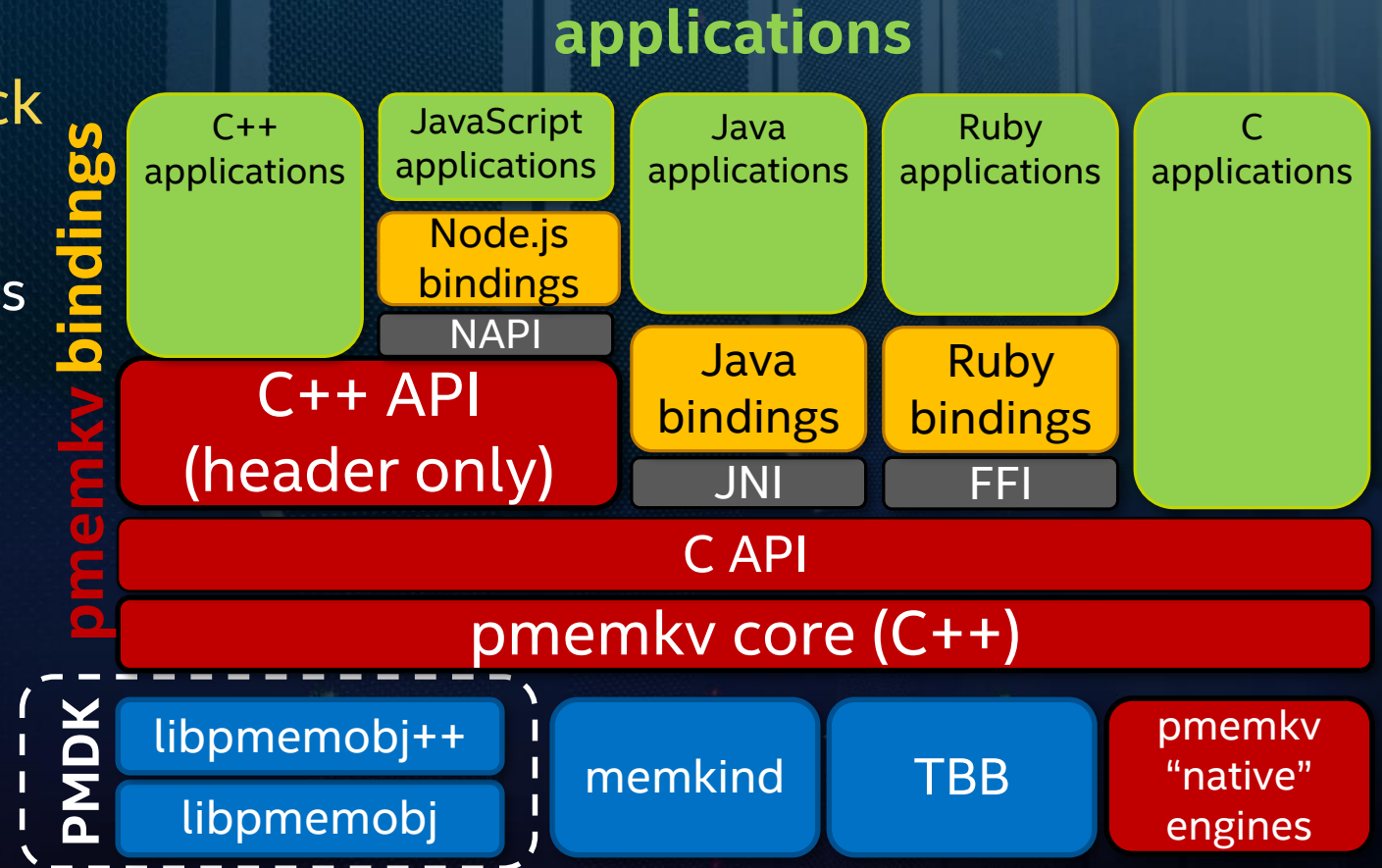


$$WA_{SW} = \frac{\text{Bytes written to the disk(SST)}}{\text{Bytes written to MemTable}}$$

# CAN WE OPTIMIZE SOFTWARE FOR PERSISTENT MEMORY?

We can compare PMEMKV with popular LSM-tree based frameworks

- LSM-tree designed for a block storage devices
  - Write optimized but the cost is write amplification
  - RocksDB and LevelDB are popular implementations
- **PMEMKV** is a key-value storage developed for Persistent Memory



# SOFTWARE WA: PMEMKV VS ROCKSDB VS LEVELDB

- **PMEMKV, RocksDB and LevelDB use the same benchmarking tool - db\_bench**
  - RocksDB and LevelDB use Persistent memory as a block storage device (Storage over AppDirect)
  - PMEMKV works in AppDirect mode (avoids page cache)
- **Selected 3 write-intensive benchmarks**
  - **fillseq** – insert entries to the empty storage in sequential order
  - **fillrandom** – insert entries to the empty storage in random order
  - **overwrite** – overwrite existing entries in prefilled storage in random order
- **Measured Write Amplification for 3 different data set sizes**
  - 10M entries, 20M entries, 30M entries



# METHODOLOGY TO MEASURE SOFTWARE WA

- **db\_bench** reports the **Raw Size**, which is cumulative size of entries inserted by the benchmark
- PMU uncore counters allow to monitor writes to Persistent Memory
  - **UNC\_M2M\_IMC\_WRITES.TO\_PMM**
  - **UNC\_M\_PMM\_WPQ\_INSERTS**
  - Counts number of cache lines sent to PMEM
- **EMON** is used to collect PMU events
  - Part of Intel® VTune™ Profiler
  - Perf can be used alternatively

$$WA_{sw} = \frac{\textit{Written to PMEM}}{\textit{Raw Size}}$$

*Written to PMEM* – measured by Uncore counter

*Raw Size* = #entries\* sizeof(entry)  
reported by benchmarks

# SOFTWARE WA: PMEMKV VS ROCKSDB VS LEVELDB

| Bench  | PMEMKV (csmmap)     |             | RocksDB             |               | LevelDB             |              |
|--|---------------------|-------------|---------------------|---------------|---------------------|--------------|
|  | Written to PMEM, Mb | WA          | Written to PMEM, Mb | WA            | Written to PMEM, Mb | WA           |
| <b>10 000 000 entries, Raw Size = 7782 Mb</b>  |                     |             |                     |               |                     |              |
| fillseq  | 32 378              | <b>4.16</b> | 35 068              | <b>4.51</b>   | 35 932              | <b>4.62</b>  |
| overwrite                                      | 20 143              | <b>2.59</b> | 1 853 610           | <b>238.19</b> | 271 942             | <b>34.95</b> |
| fillrandom                                     | 27 800              | <b>3.57</b> | 92 446              | <b>11.88</b>  | 243 860             | <b>31.34</b> |
| <b>20 000 000 entries, Raw Size = 15564 Mb</b> |                     |             |                     |               |                     |              |
| fillseq  | 64 294              | <b>4.13</b> | 70 439              | <b>4.53</b>   | 71 523              | <b>4.60</b>  |
| overwrite                                      | 40 285              | <b>2.59</b> | 2 191 380           | <b>140.80</b> | 599 998             | <b>38.55</b> |
| fillrandom                                     | 55 637              | <b>3.57</b> | 219 580             | <b>14.11</b>  | 551 925             | <b>35.46</b> |
| <b>40 000 000 entries, Raw Size = 31128 Mb</b> |                     |             |                     |               |                     |              |
| fillseq  | 128 196             | <b>4.12</b> | 141 099             | <b>4.53</b>   | 140 611             | <b>4.52</b>  |
| overwrite                                      | 80 575              | <b>2.59</b> | 2 478 590           | <b>79.63</b>  | 1 307 440           | <b>42.00</b> |
| fillrandom                                     | 111 377             | <b>3.58</b> | 588 285             | <b>18.90</b>  | 1 221 960           | <b>39.26</b> |

Key size = 16 bytes.  
Value size = 800 bytes

$$WA_{sw} = \frac{\text{Written to PMEM}}{\text{Raw Size}}$$

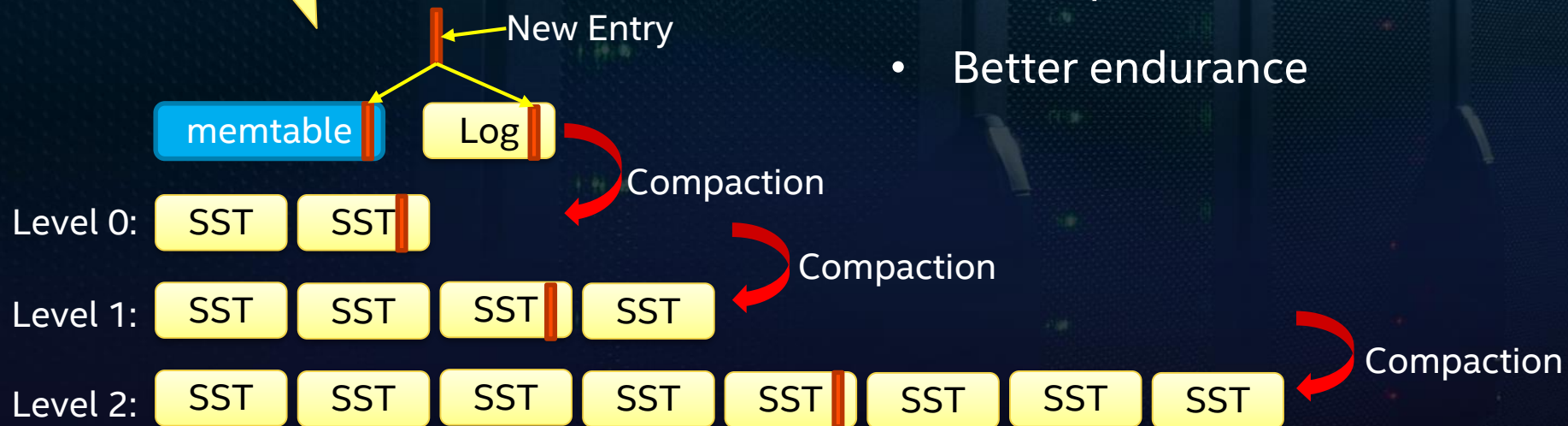
*Written to PMEM* – measured by Uncore counter

*Raw Size* = #entries \* sizeof(entry)  
reported by benchmarks

# POSSIBLE SOFTWARE MODERNIZATION (1/4)

No code changes.  
Use existing software!

- Run existing software on Persistent Memory
  - Storage over AppDirect mode
- Intel Optane Persistent Memory vs SSDs
  - Better performance
  - Better endurance



# POSSIBLE SOFTWARE MODERNIZATION (2/4)

Minor modifications.  
Use Libpmemlog as a memtable log.

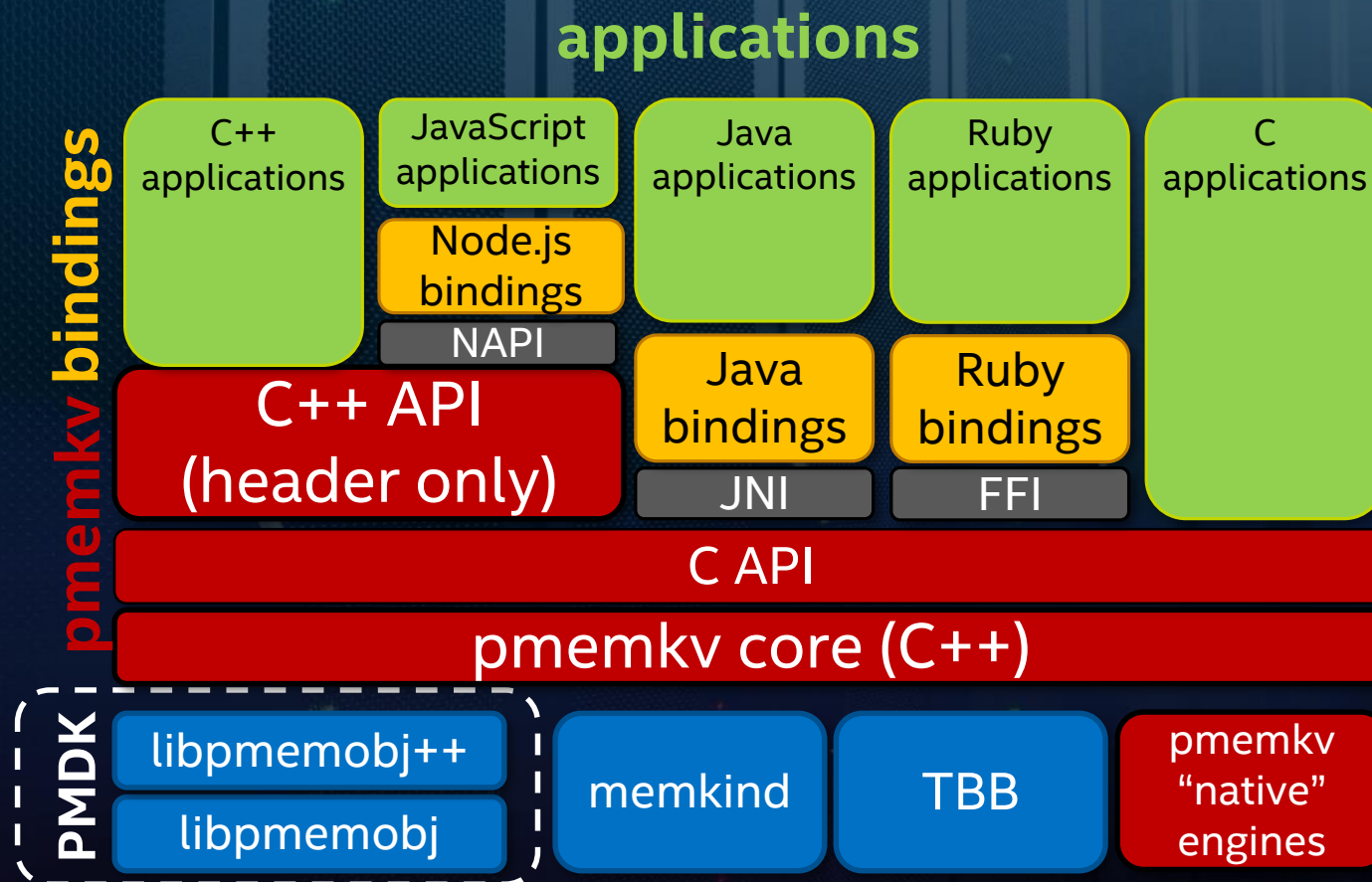
- **Libpmemlog** - persistent memory resident log file
  - Non-paged direct access
  - Improves insert performance
- **Still high WA due to compaction of SSTs**



# POSSIBLE SOFTWARE MODERNIZATION (3/4)

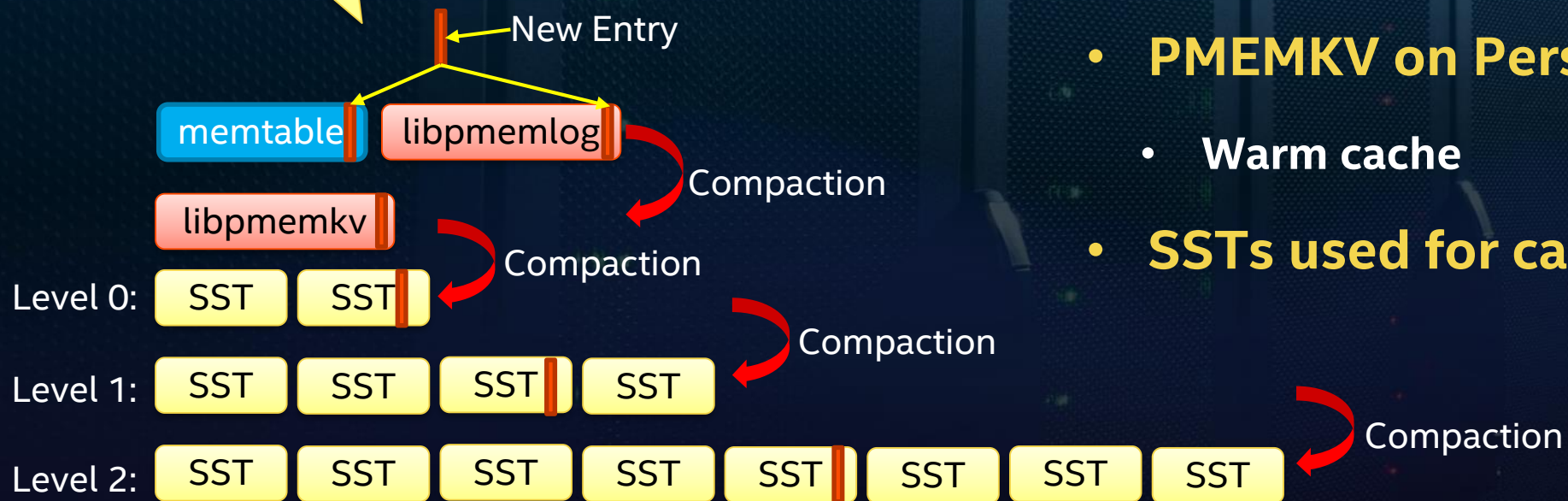
Replace existing key/value storage with PMEMKV

- Developed for Persistent Memory
- Pluggable engines
  - Native engines for persistent and volatile use-cases
- Native API for C/C++
- Easy integration with high-level language bindings



# POSSIBLE SOFTWARE MODERNIZATION (4/4)

Hybrid solution



- **Memtable on DRAM**
  - Fast hot cache
- **Libpmemlog persistent log**
- **PMEMKV on Persistent memory**
  - Warm cache
- **SSTs used for capacity**

# SUMMARY

- Intel Optane Persistent Memory is a new layer in a storage hierarchy
  - better endurance comparing to NAND SSD
  - better performance comparing to NAND SSD
- Intel Optane Persistent Memory might help to decrease TCO on write intensive workloads if SSD wear-out quickly
- PMDK provides set of libraries for easier code modernization and enabling Persistent Memory in applications



Storage Performance Development Kit (SPDK)  
Persistent Memory Development Kit (PMDK)  
Intel® VTune™ Profiler

**Virtual Forum**