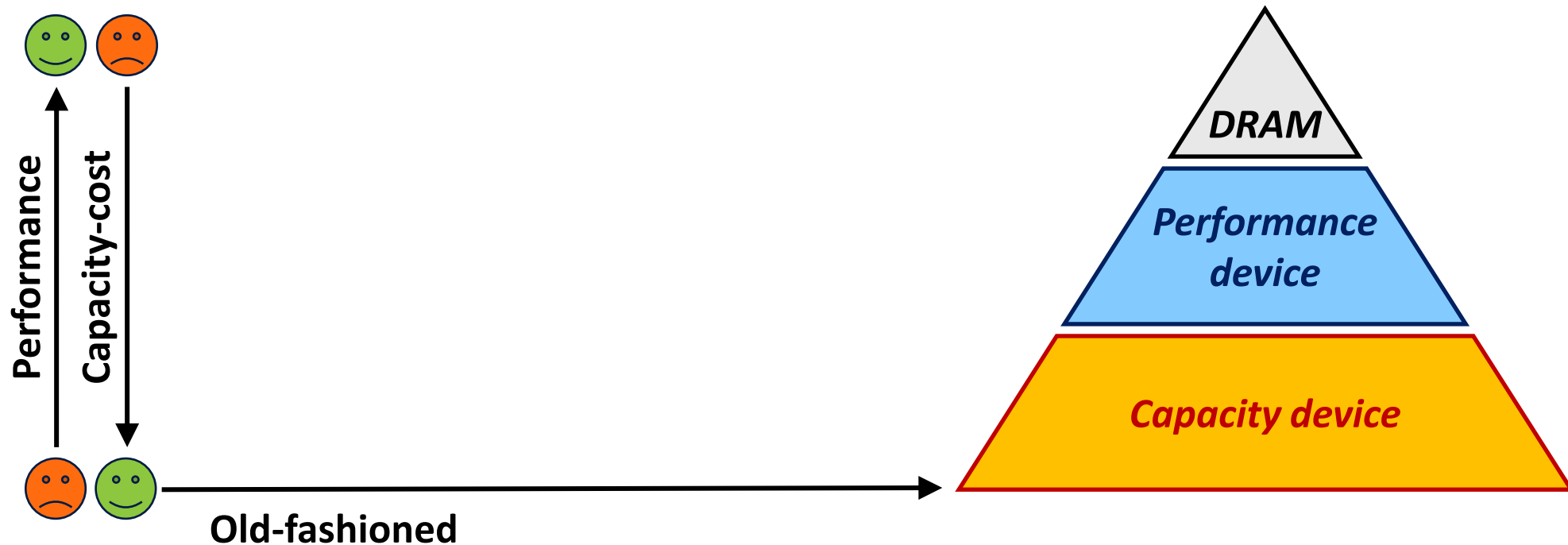# PRISM: Optimizing Key-Value Store for Modern Heterogeneous Storage Devices

*Yongju Song, Wook-Hee Kim, Sumit Kumar Monga,*
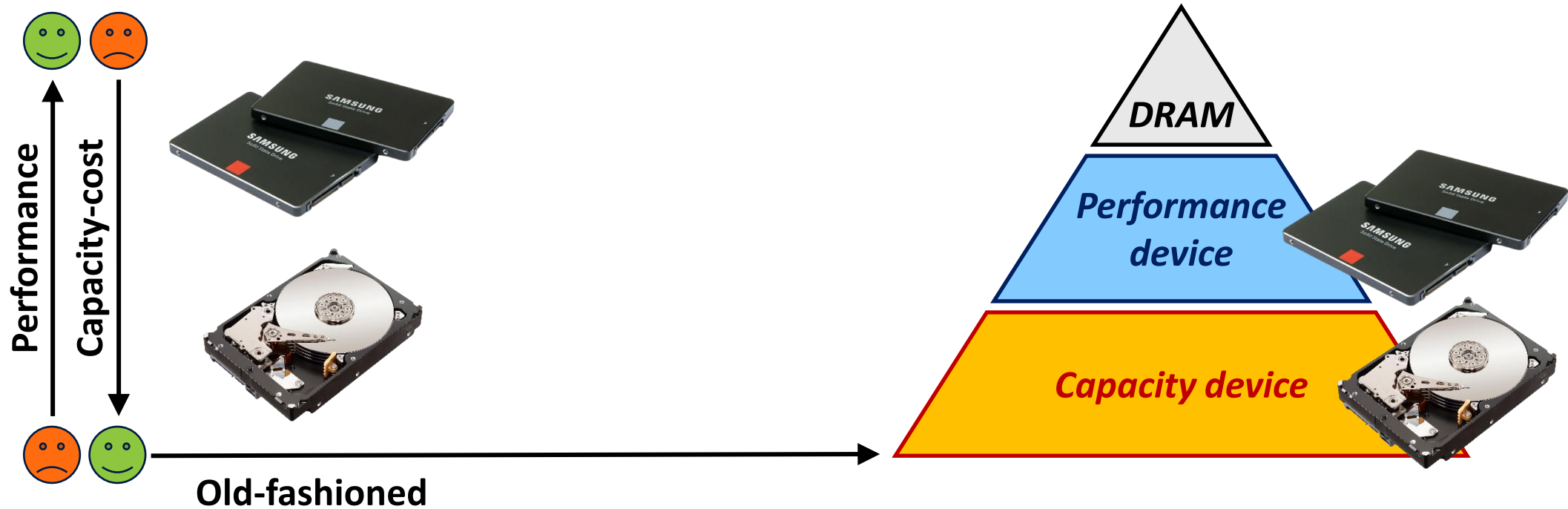
*Changwoo Min, and Young Ik Eom*

# Heterogeneous Storage Systems

A method for assigning different categories of data
to *various types of storage media* to *reduce overall storage costs*.

# Heterogeneous Storage Systems

A method for assigning different categories of data
to *various types of storage media* to *reduce overall storage costs*.
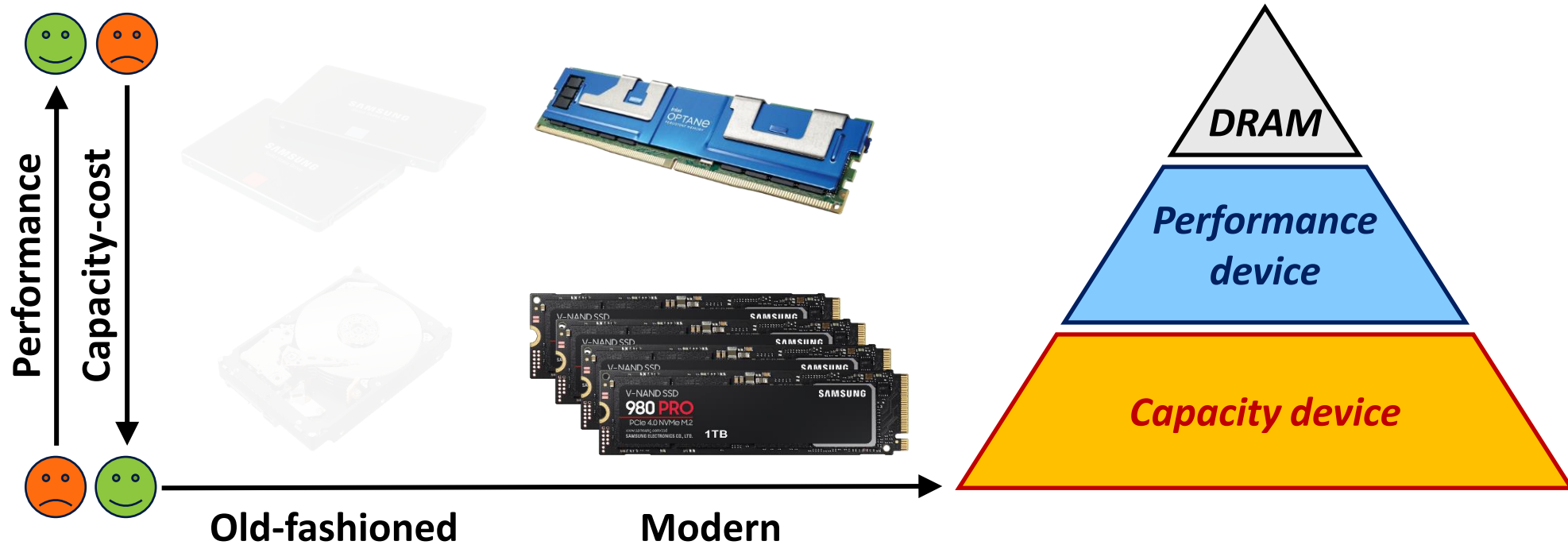


Old-fashioned

# Heterogeneous Storage Systems

A method for assigning different categories of data
to *various types of storage media* to *reduce overall storage costs*.



Performance

Capacity-cost

Old-fashioned        Modern

DRAM

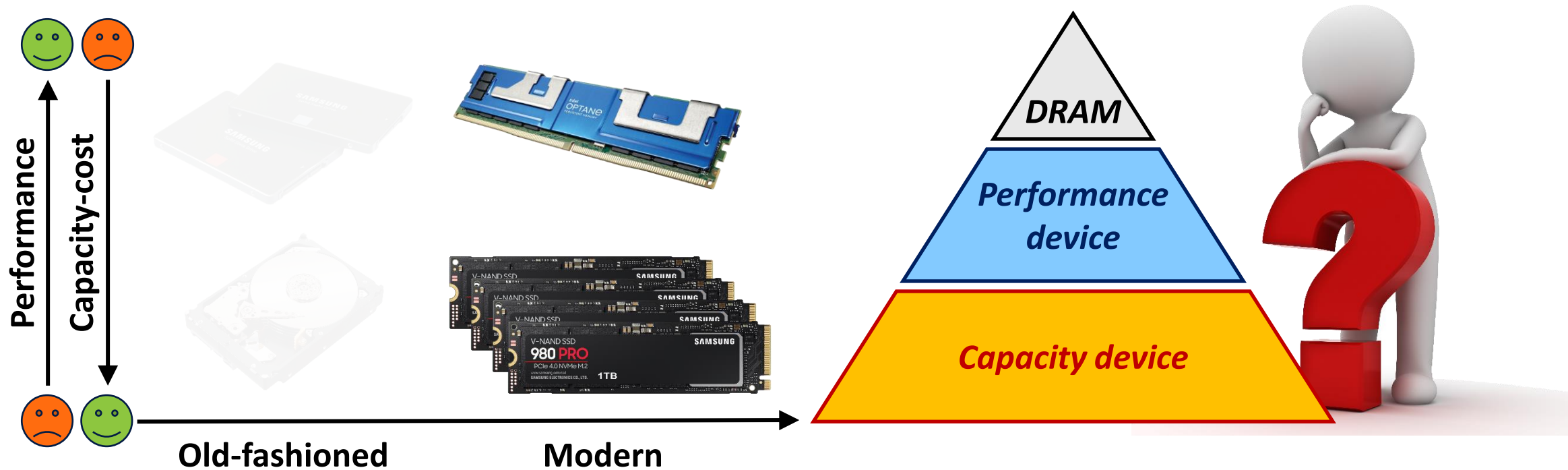*Performance device*

*Capacity device*

4

# Heterogeneous Storage Systems

A method for assigning different categories of data
to *various types of storage media* to *reduce overall storage costs*.



Performance / Capacity-cost

Old-fashioned          Modern

DRAM

*Performance device*

*Capacity device*

# Evolution of Storage Heterogeneity

| Specification | | Capacity | Cost | Performance | |
|---|---|---|---|---|---|
| Type | Model | GB | $/TB | Read Latency (usec) | Write Latency (usec) |
| DRAM | SK Hynix DRAM w/DDR4 | 16 | 5,427 | 0.08 | 0.08 |
| NVM | Intel Optane DCPMM w/DDR-T | 128 | 4,096 | 0.30 | 0.09 |
| NVM SSD | Intel Optane 905P w/PCIe 3 | 960 | 1,024 | 10 | 10 |
| Flash SSD | Samsung 980 Pro w/PCIe 4 | 1024 | 150 | 50 | 20 |
| Flash SSD | Samsung 980 w/PCIe 3 | 1024 | 100 | 60 | 20 |

# Evolution of Storage Heterogeneity

| Specification | | Capacity | Cost | Performance | |
|---|---|---|---|---|---|
| Type | Model | GB | $/TB | Read Latency (usec) | Write Latency (usec) |
| DRAM | SK Hynix DRAM w/DDR4 | 16 | 5,427 | 0.08 | 0.08 |
| NVM | Intel Optane DCPMM w/DDR-T | 128 | 4,096 | **0.30** | **0.09** |
| NVM SSD | Intel Optane 905P w/PCIe 3 | 960 | 1,024 | 10 | 10 |
| Flash SSD | Samsung 980 Pro w/PCIe 4 | **1024** | **150** | 50 | 20 |
| Flash SSD | Samsung 980 w/PCIe 3 | **1024** | **100** | 60 | 20 |

*Performance devices*

*Capacity devices*

# Evolution of Storage Heterogeneity

✓  ✓  ✓

| Specification | | Capacity | Cost | Performance | | | | Endurance |
|---|---|---|---|---|---|---|---|---|
| Type | Model | GB | $/TB | Read Latency (usec) | Write Latency (usec) | Read BW (GB/s) | Write BW (GB/s) | Warranty (PBW) |
| DRAM | SK Hynix DRAM w/DDR4 | 16 | 5,427 | 0.08 | 0.08 | 15 | 15 | ∞ |
| NVM | Intel Optane DCPMM w/DDR-T | 128 | 4,096 | **0.30** | **0.09** | 6.8 | 1.9 | **292** |
| NVM SSD | Intel Optane 905P w/PCIe 3 | 960 | 1,024 | 10 | 10 | 2.6 | 2.2 | 17.52 |
| Flash SSD | Samsung 980 Pro w/PCIe 4 | **1024** | **150** | 50 | 20 | **7** | **5** | 0.6 |
| Flash SSD | Samsung 980 w/PCIe 3 | **1024** | **100** | 60 | 20 | 3.5 | 3 | 0.6 |

# Evolution of Storage Heterogeneity

|  | | | | ✓ | | ✓ | ✓ | ✓ |
| Specification | | Capacity | Cost | Performance | | | | Endurance |
|---|---|---|---|---|---|---|---|---|
| Type | Model | GB | $/TB | Read Latency (usec) | Write Latency (usec) | Read BW (GB/s) | Write BW (GB/s) | Warranty (PBW) |
| DRAM | SK Hynix DRAM w/DDR4 | 16 | 5,427 | 0.08 | 0.08 | 15 | 15 | ∞ |
| NVM | Intel Optane DCPMM w/DDR-T | 128 | 4,096 | *0.30* | *0.09* | 6.8 | 1.9 | *292* |
| NVM SSD | Intel Optane 905P w/PCIe 3 | 960 | 1,024 | 10 | 10 | 2.6 | 2.2 | 17.52 |
| Flash SSD | Samsung 980 Pro w/PCIe 4 | *1024* | *150* | 50 | 20 | *7* | *5* | 0.6 |
| Flash SSD | Samsung 980 w/PCIe 3 | *1024* | *100* | 60 | 20 | 3.5 | 3 | 0.6 |

There is *No Clear Separation* between performance-/capacity- devices.

"The Storage Hierarchy is *Becoming a Jungle*." [CIDR'21, Dong Xie]

"The Storage Hierarchy is *Not a Hierarchy*." [FAST'21, Remzi H. Arpaci-Dusseau]

# Today's Storage Hierarchy

# Today's Storage Hierarchy

Placing hot data on NVM

- System can leverage the low latency of NVM but *suffer from its limited bandwidth*.
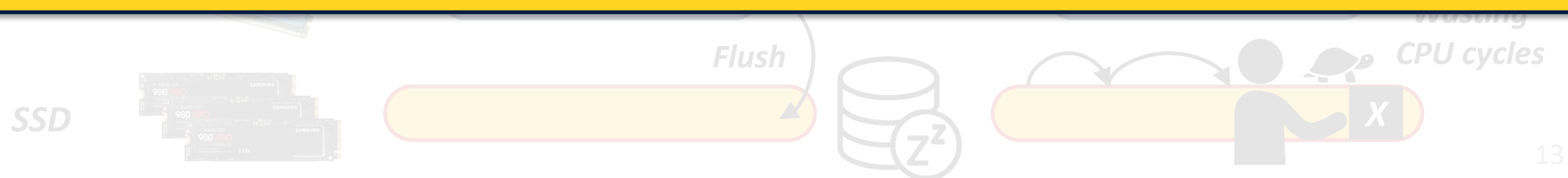
# Today's Storage Hierarchy

Placing hot data on NVM

- System can leverage the low latency of NVM but *suffer from its limited bandwidth*.

Traversing data layer by layer for handling read requests

- Inefficient traversal leads to *wasting CPU cycles*.

- Overall performance may be *bounded to the device with the lowest performance*.

# Today's Storage Hierarchy

Placing hot data on NVM

- System can leverage the low latency of NVM but suffer from its limited bandwidth.

Traversing data layer by layer for handling read requests

- Inefficient traversal leads to wasting CPU cycles.
- Overall performance may be bounded to the device with the lowest performance.

*How should we design a **Heterogeneous Storage System** in the **Modern Storage Landscape**?*

Flush

Wasting
CPU cycles

SSD

# Design Goals of *PRISM*

*Drawing the full potential of heterogeneous storage devices.*

*Minimizing the overhead of software stack for scalability*

*Providing a high level of crash consistency & concurrency*

# Overview of *PRISM*

# Overview of *PRISM*: Insert (k1, v1)

# Overview of *PRISM*: Insert (k1, v1)

**Background Reclamation of PWB**

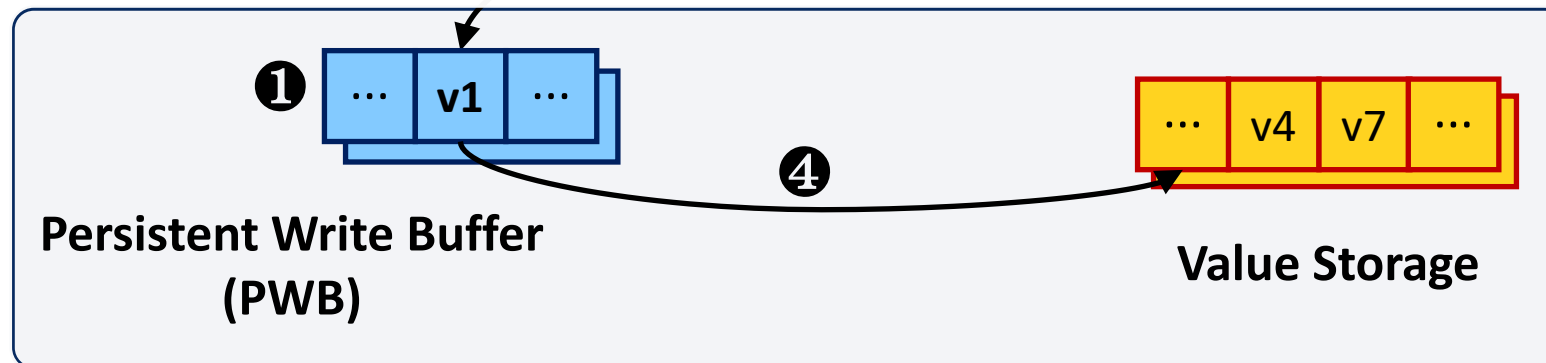- Preventing application threads from blocking

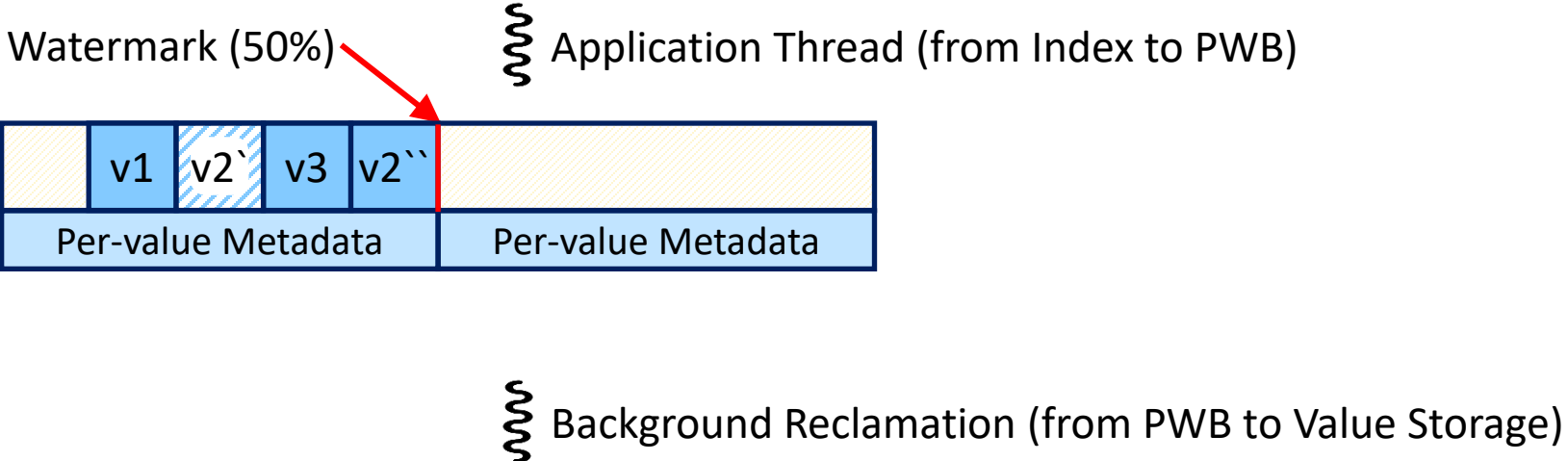**Asynchronous I/O batching**

- Achieving high bandwidth of SSD

# Asynchronous Bandwidth-Optimized WRITE

**Persistent Write Buffer (PWB)**

Watermark (50%)

Application Thread (from Index to PWB)

| | v1 | v2` | v3 | v2`` | |
|---|---|---|---|---|---|
| Per-value Metadata | | | | Per-value Metadata | |

Background Reclamation (from PWB to Value Storage)

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

**Value Storage**

Free Chunk N (512KB)

| Chunk 1 | ... | v1 | v2`` | v3 |
|---|---|---|---|---|
| | | Per-value Metadata | | |

Backward ptr

Value size

...

Append-only writes

18

# Asynchronous Bandwidth-Optimized WRITE

# Design Overview of PRISM: Lookup(k4)

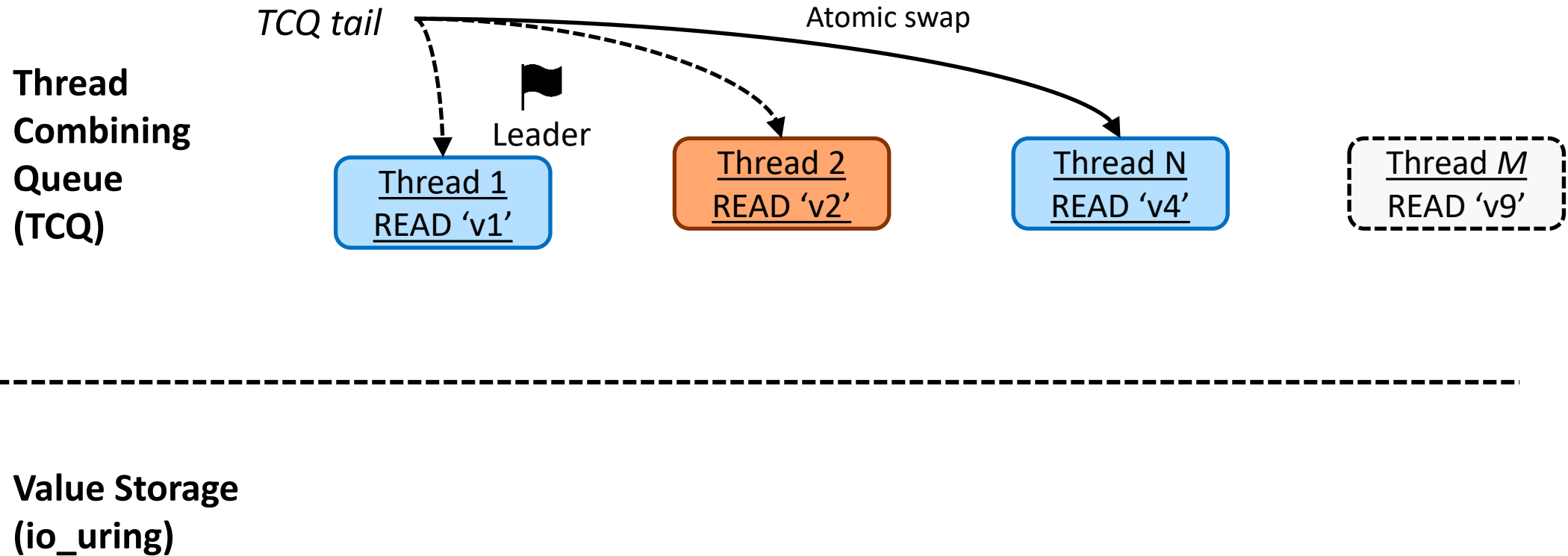# Design Overview of *PRISM*: Lookup(k4)

**Adjust the IO batch size for SSD reads according to thread concurrency**

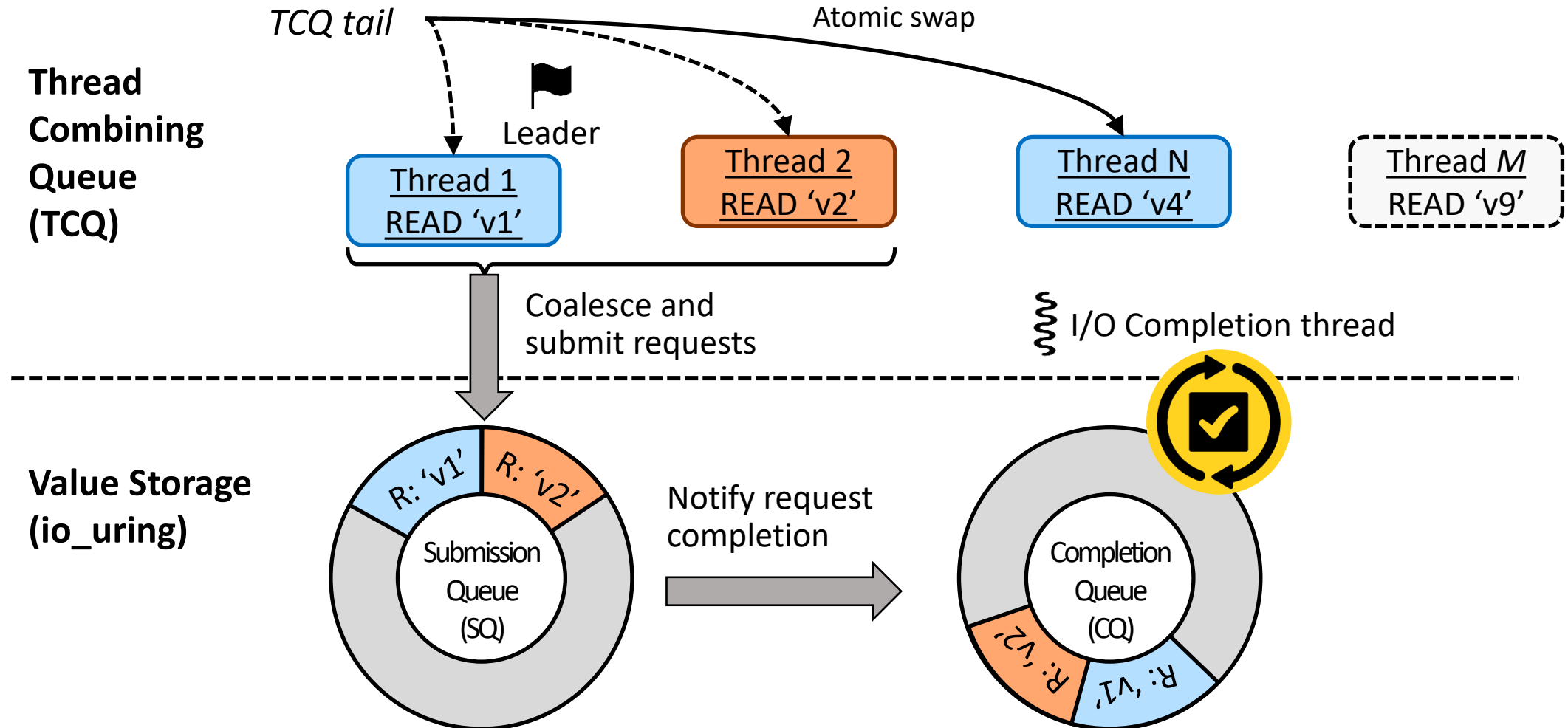**Combine reads from multiple threads to a single read operation**

- Aggressively utilizing the bandwidth and hide latency of SSD



**Value Storage**

*Opportunistic Thread Combining for READ*

# Opportunistic Thread Combining for READ

TCQ tail                                    Atomic swap

**Thread
Combining
Queue
(TCQ)**

🏴 Leader

| Thread 1 | Thread 2 | Thread N | Thread *M* |
|----------|----------|----------|------------|
| READ 'v1' | READ 'v2' | READ 'v4' | READ 'v9' |

**Value Storage
(io_uring)**

# Opportunistic Thread Combining for READ



TCQ tail

Atomic swap

**Thread Combining Queue (TCQ)**

Leader

Thread 1
READ 'v1'

Thread 2
READ 'v2'

Thread N
READ 'v4'

Thread *M*
READ 'v9'

Coalesce and submit requests

I/O Completion thread

**Value Storage (io_uring)**

R: 'v1'   R: 'v2'

Submission Queue (SQ)

Notify request completion

Completion Queue (CQ)

R: 'v2'   R: 'v1'

# Cross-media Crash Consistency



**Persistent Key Index**

**Heterogeneous Storage Index Table (HSIT)**

k1  k4  k7

... | ... | ... | ...

NVM

SSD

DRAM

v1` | v1`` | ...

**Persistent Write Buffer (PWB)**

... | v4 | v7 | ...

**Value Storage**

... | v7 | ...

**Scan-aware Value Cache (SVC)**

# Cross-media Crash Consistency

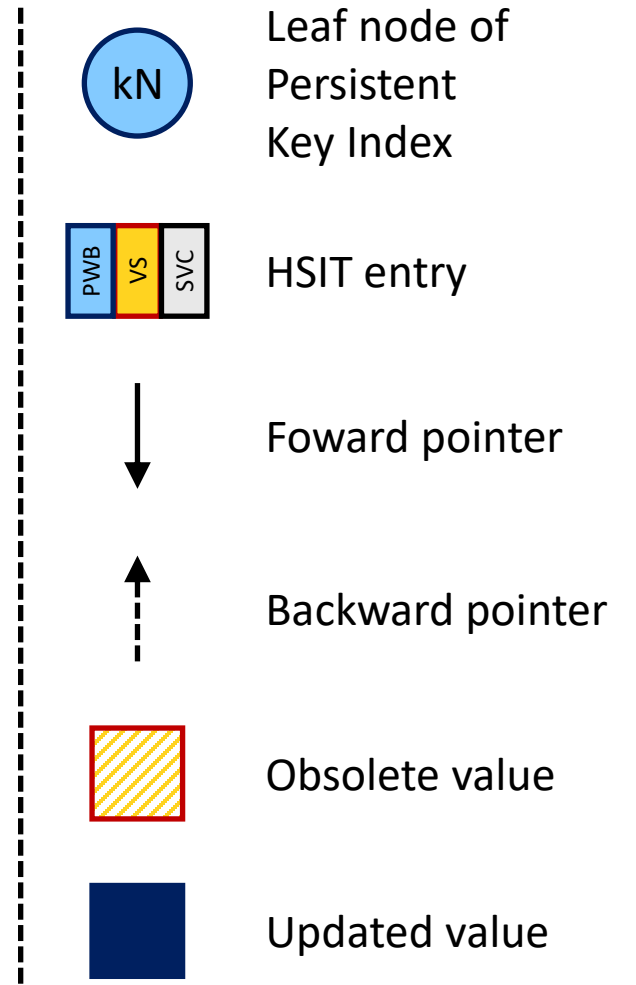## Components are scattered across multiple heterogeneous devices

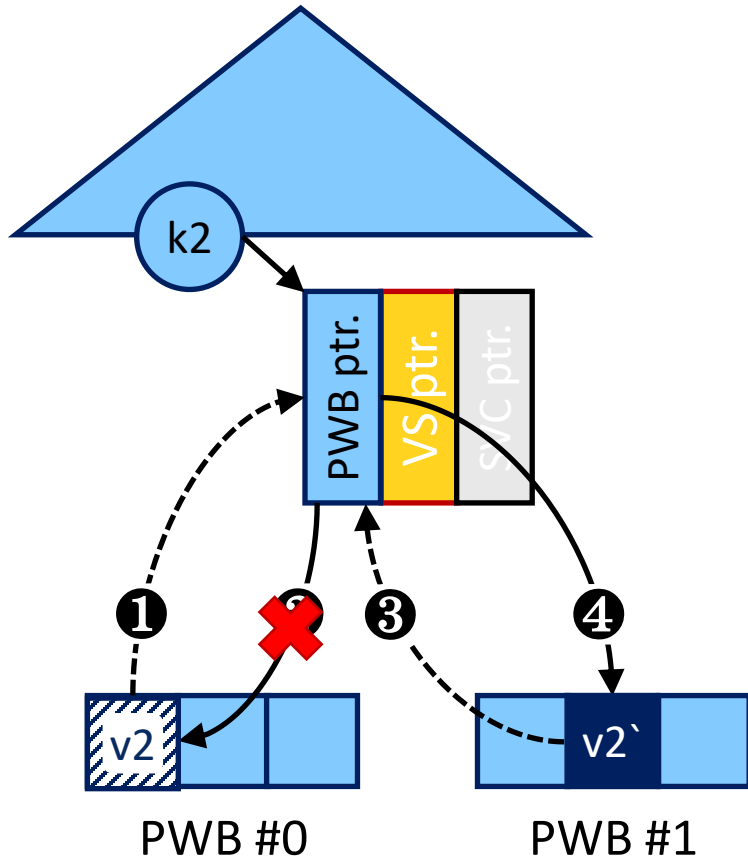- Lightweight crash consistency with Forward & Backward pointers

Cross-media Crash Consistency

Heterogeneous Storage Index Table (HSIT)

Persistent Write Buffer (PWB)

Value Storage

Scan-aware Value Cache (SVC)

# Crash Consistent Update of Values with HSIT



INSERT {k2, v2}

k2

PWB ptr.
VS ptr.
SVC ptr.

❶ ❷

v2

PWB #0    PWB #1

kN — Leaf node of Persistent Key Index

PWB VS SVC — HSIT entry

Foward pointer

Backward pointer

Obsolete value

Updated value

# Crash Consistent Update of Values with HSIT

# Crash Consistent Update of Values with HSIT



UPDATE {k2, v2} to {k2, v2`}

UPDATE {k3, v3} to {k3, v3`}

Leaf node of Persistent Key Index

HSIT entry

Foward pointer

Backward pointer

Obsolete value

Updated value

PWB #0     PWB #1

PWB #0     Value Storage #1

# Experimental Setup

## Hardware environment

- Two-socket Intel Xeon machine
- Each socket: 20 CPU cores,
  **Six 128GB Intel Optane DIMMSs**, and 96GB DRAM
- **Eight Samsung 980 PRO 1TB SSDs**
  with two NVMe RAID Controllers HighPoint SSD7103

## Competitors

- KVell: DRAM-SSD with up to 64 batched I/Os [SOSP'21]
- MatrixKV: DRAM-NVM-SSD [ATC'20]
- Allocated their hardware resources at the same cost levels

# Performance Comparison on YCSB

WRITE: Does not require level-compaction & Per-thread write buffer

READ: No need for traversing multiple levels & Efficient KV item caching

# Opportunistic Thread Combining

Prism opportunistically adjust the IO batch size
for read operations according to thread concurrency.

# In the paper…

Performance under other workloads

Performance impact of ..

- Number of SSDs
- Size of PWB/SVC
- Write amplification
- Garbage collection in Value Storage
- Individual techniques

Size of NVM space

Recovery

# Conclusion
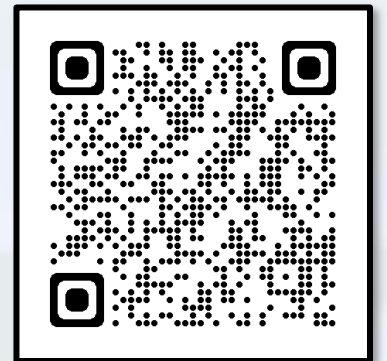
We answered the question:

*How should we design a Heterogeneous Storage System in the Modern Storage Landscape?*

- Synergistic Five Components
- Asynchronous Bandwidth-Optimized WRITE
- Opportunistic Thread Combining
- Cross-media Crash Consistency & Concurrency Control using Forward & Backward Pointers

# PRISM: Optimizing Key-Value Store for Modern Heterogeneous Storage Devices

*Yongju Song, Wook-Hee Kim, Sumit Kumar Monga,*

*Changwoo Min, and Young Ik Eom*

*Paper*