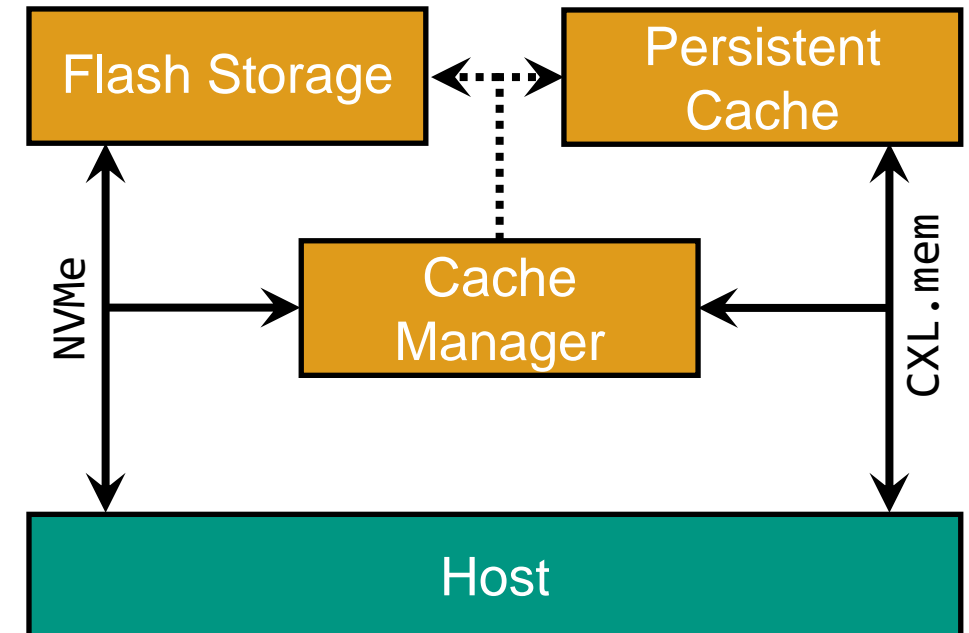# Fundamental OS Design Considerations for CXL-based Hybrid SSDs

**Daniel Habicht**, Yussuf Khalil, Lukas Werling, Thorsten Gröninger, and Frank Bellosa
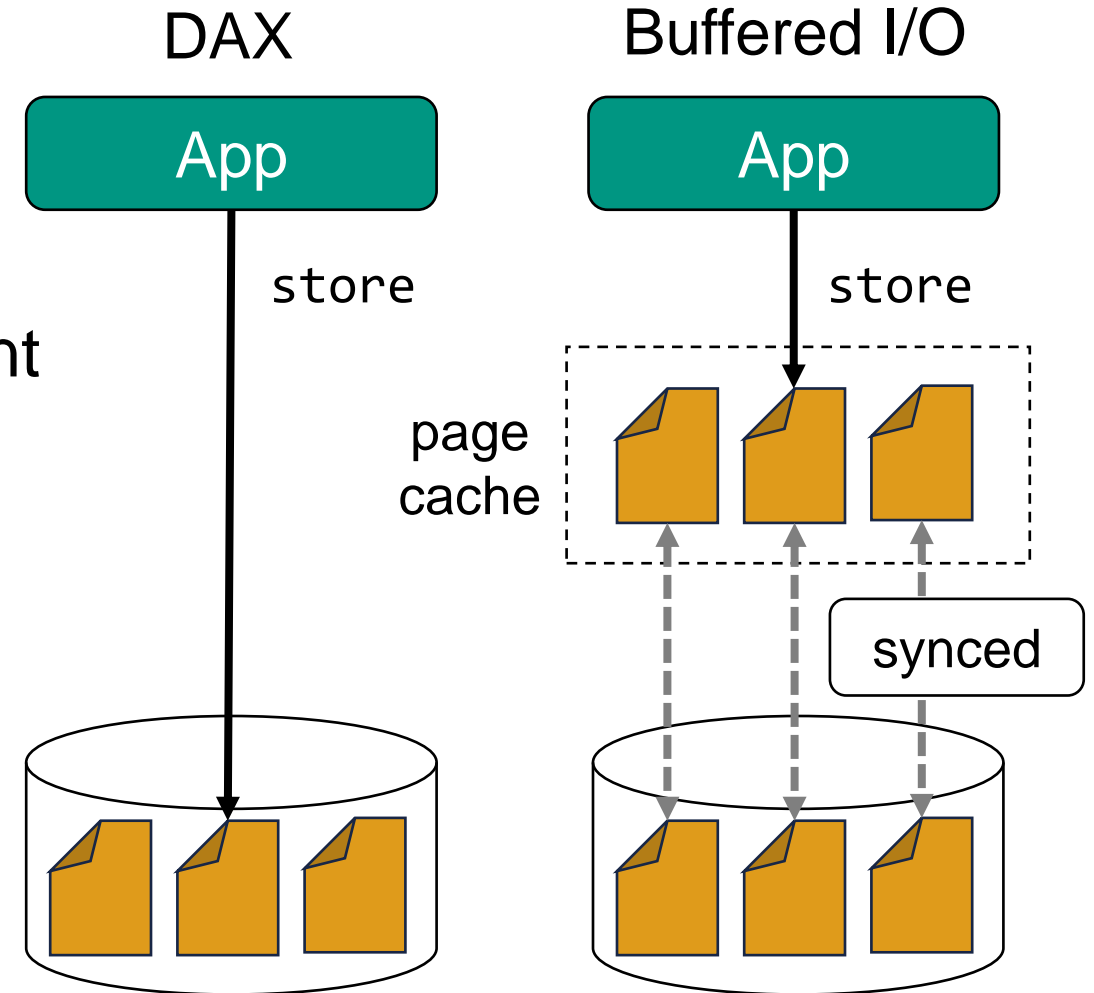
# Hybrid SSDs as Cost-Effective PM

✅ cost-effective (Flash ≫ cache)

❌ no uniform access due to cache

❌ existing OS abstractions unsuitable

Our contribution:
OS-centric hybrid SSD management

# Linux Direct Access (DAX)

- DAX = (volatile) page cache bypass

- Per-inode DAX flag
  - ❌ no fine-granular resource management
  - ❌ pressure on small on-device cache

- Assumes non-blocking access
  - ❌ CPU stalls on cache miss
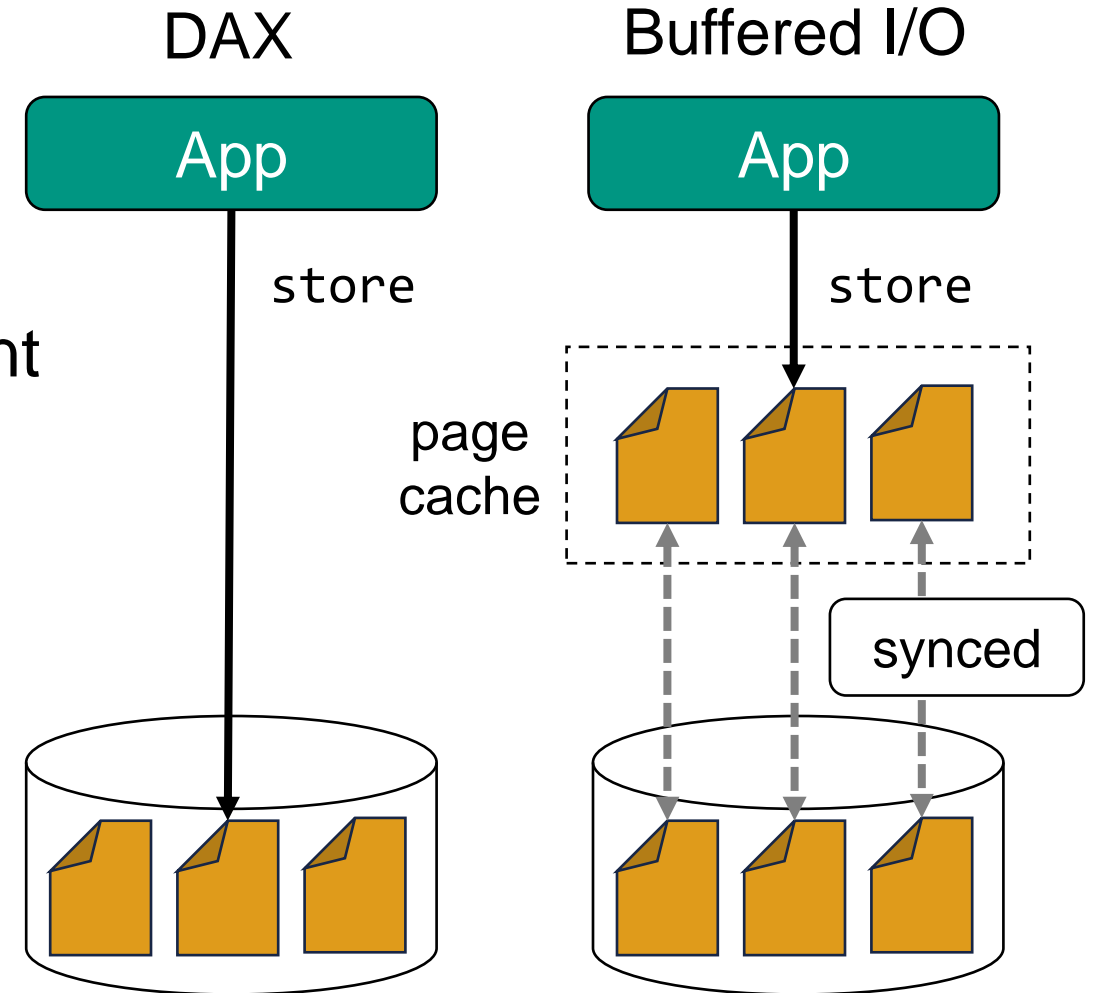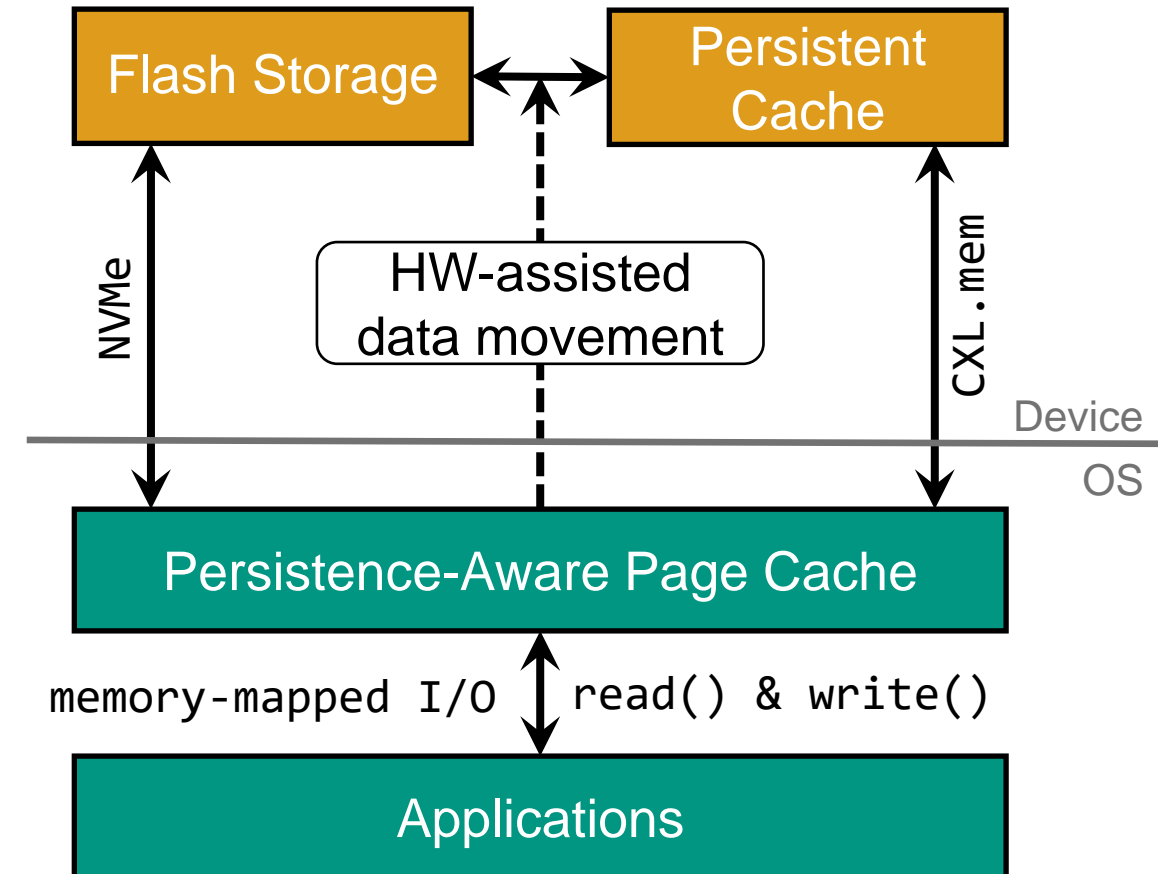
# Linux Direct Access (DAX)

- DAX = (volatile) page cache bypass

- Per-inode DAX flag
  - ❌ no fine-granular resource management
  - ❌ pressure on small on-device cache

- Assumes non-blocking access
  - ❌ CPU stalls on cache miss

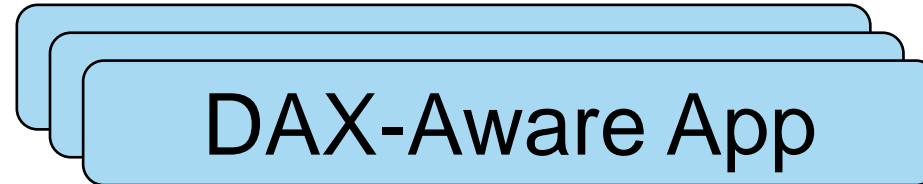➡ Existing DAX support unsuitable for hybrid SSDs

DAX

App

`store`

Buffered I/O

App

`store`

page cache

synced

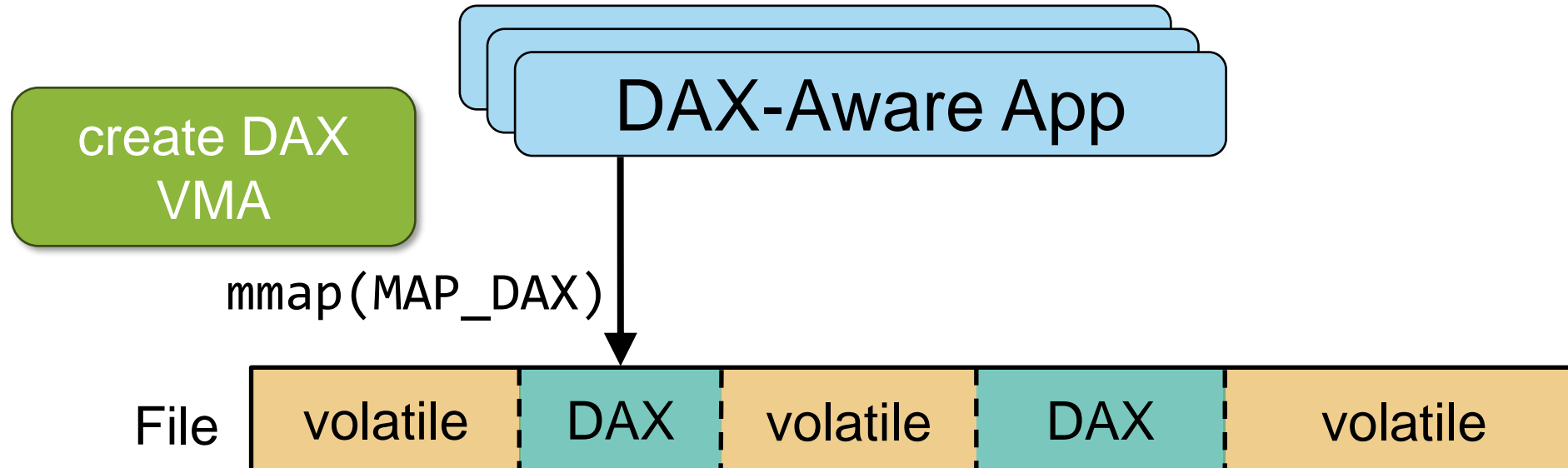# Our Solution: Persistence-Aware Page Cache

- **OS-centric cache management**
  - Host page tables reflect cache state
  - Hardware-assisted data movement
  - Expose resource management to apps

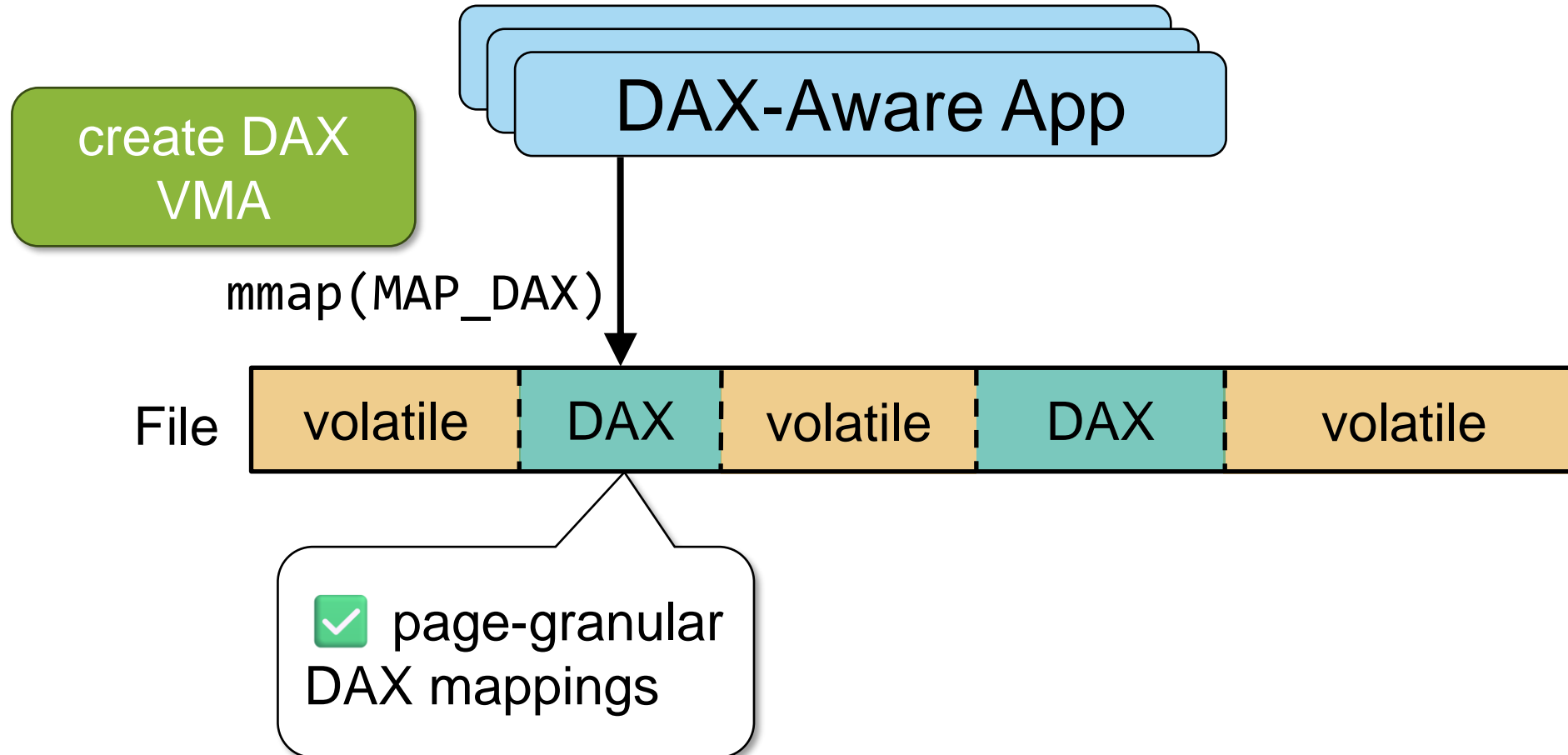- **Leverage persistence of DAX pages for lightweight `fsync()`**

# Fine-Granular Resource Management
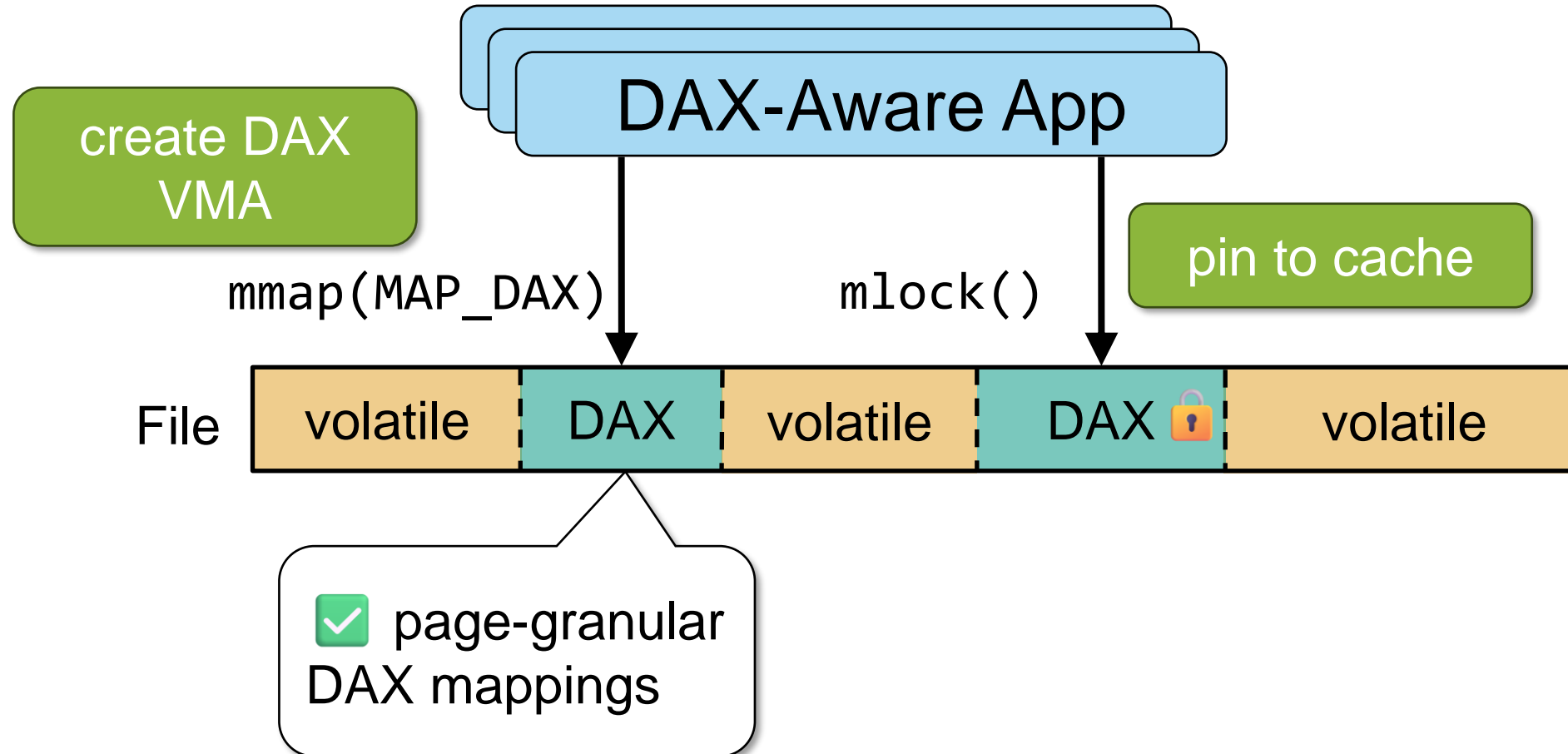
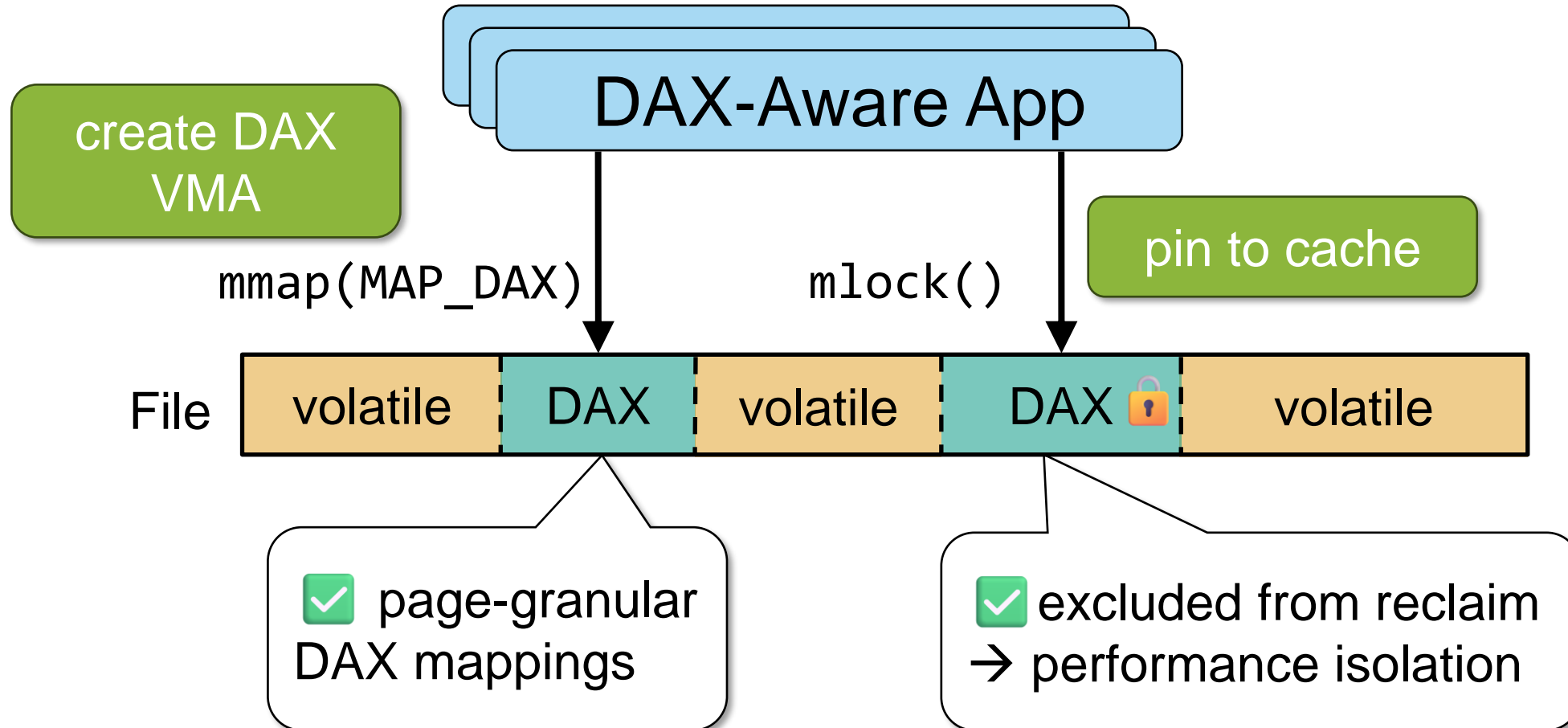# Fine-Granular Resource Management

# Fine-Granular Resource Management



create DAX VMA

DAX-Aware App

`mmap(MAP_DAX)`

File | volatile | DAX | volatile | DAX | volatile

✅ page-granular DAX mappings

# Fine-Granular Resource Management

create DAX VMA

DAX-Aware App

pin to cache

`mmap(MAP_DAX)`　　　　　`mlock()`

File | volatile | DAX | volatile | DAX 🔒 | volatile

✅ page-granular DAX mappings

# Fine-Granular Resource Management



create DAX VMA

pin to cache

DAX-Aware App

mmap(MAP_DAX)          mlock()

File | volatile | DAX | volatile | DAX 🔒 | volatile |

✅ page-granular DAX mappings

✅ excluded from reclaim
→ performance isolation

# DAX-Aware Page Cache Allocation

# DAX-Aware Page Cache Allocation



read()

0x0000

0x2000

0x4000

mmaped file

page cache

pgoff 3 DAX?

per-file rb tree of DAX VMAs

DAX page    volatile page

# DAX-Aware Page Cache Allocation



pgoff 3 DAX?

✅ index overlaps DAX VMA

per-file rb tree of DAX VMAs

read()

0x0000

0x2000

0x4000

X

mmaped file    page cache

DAX page    volatile page
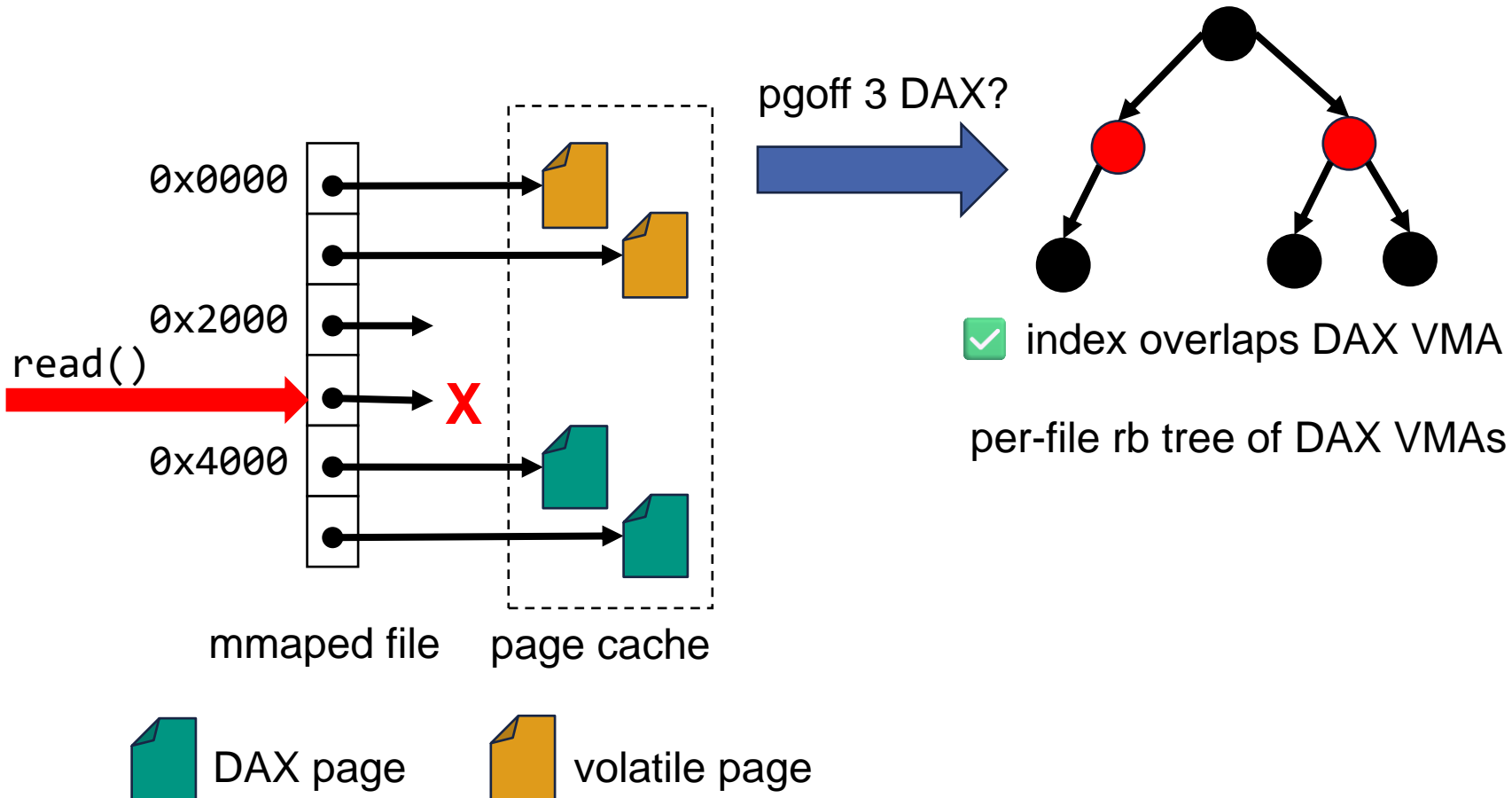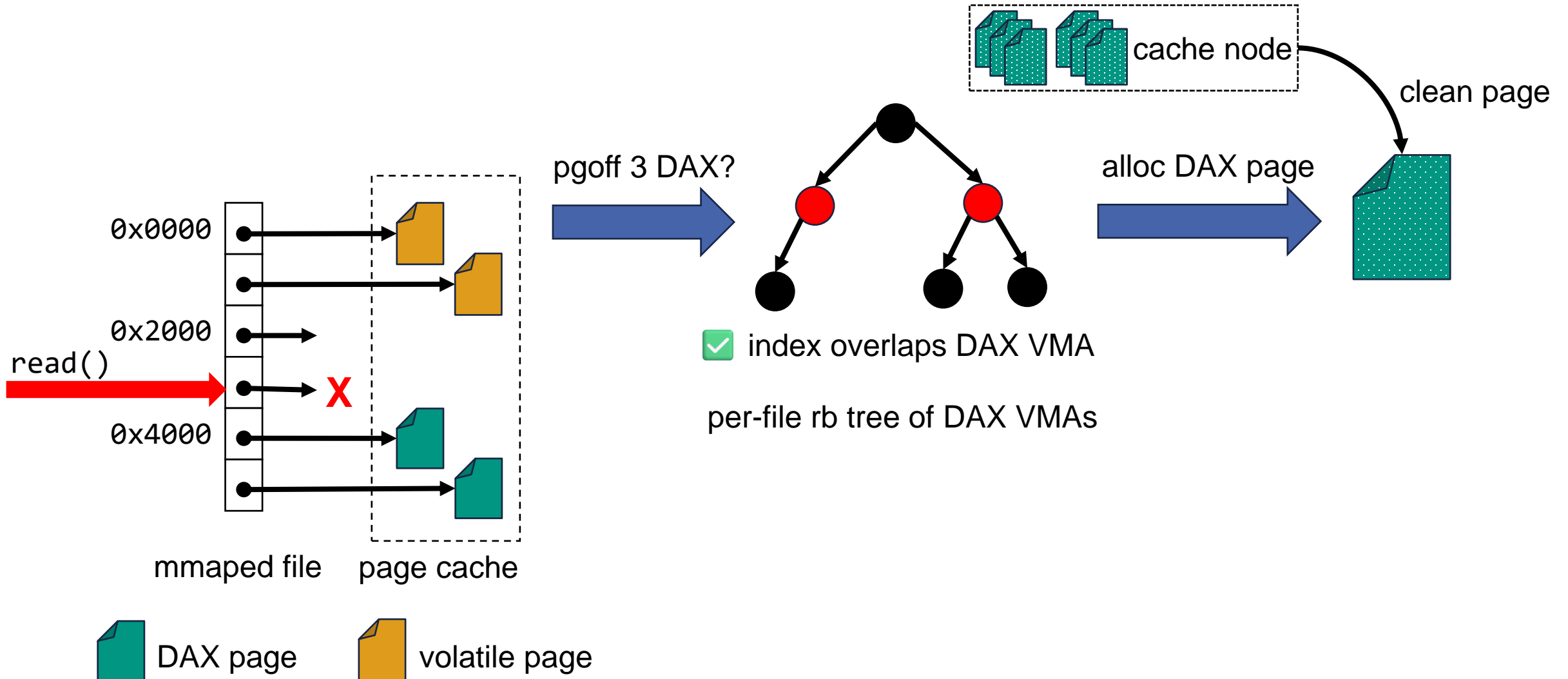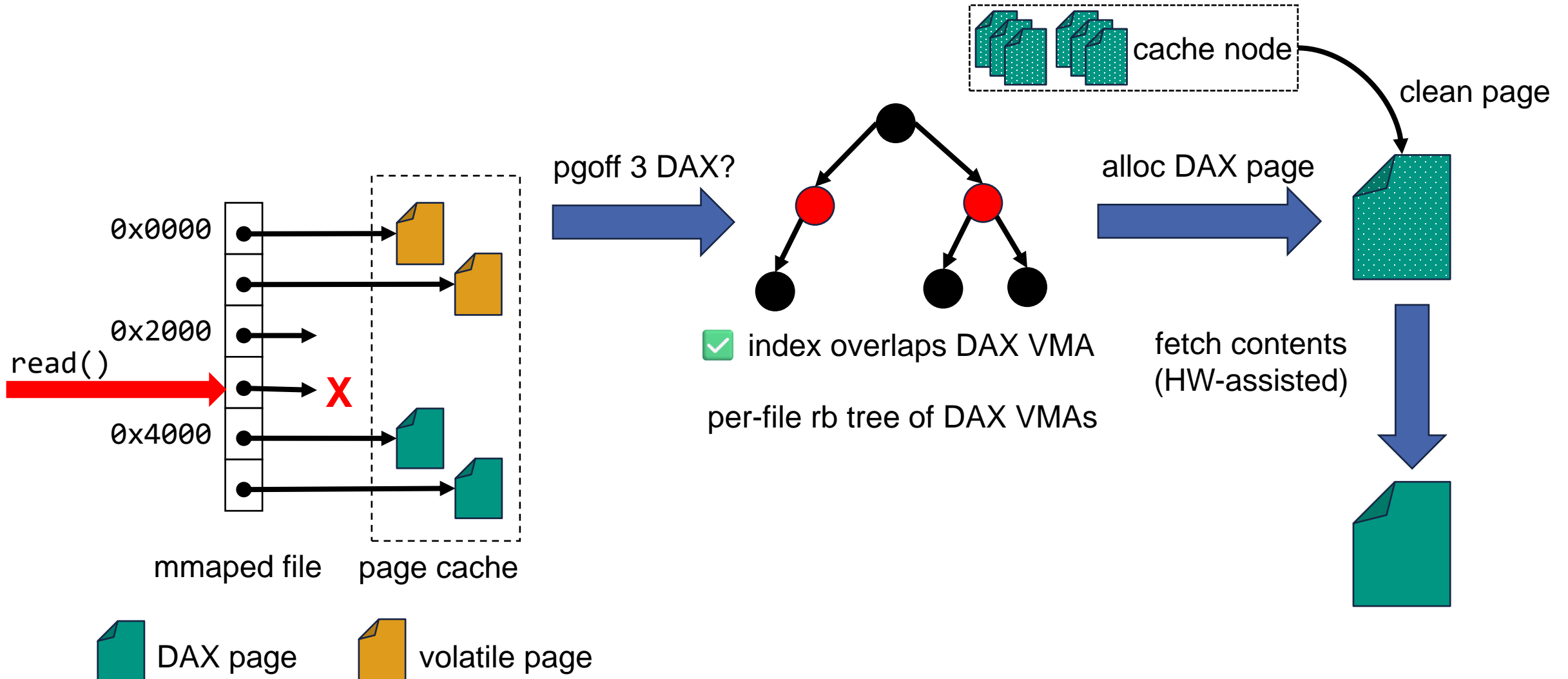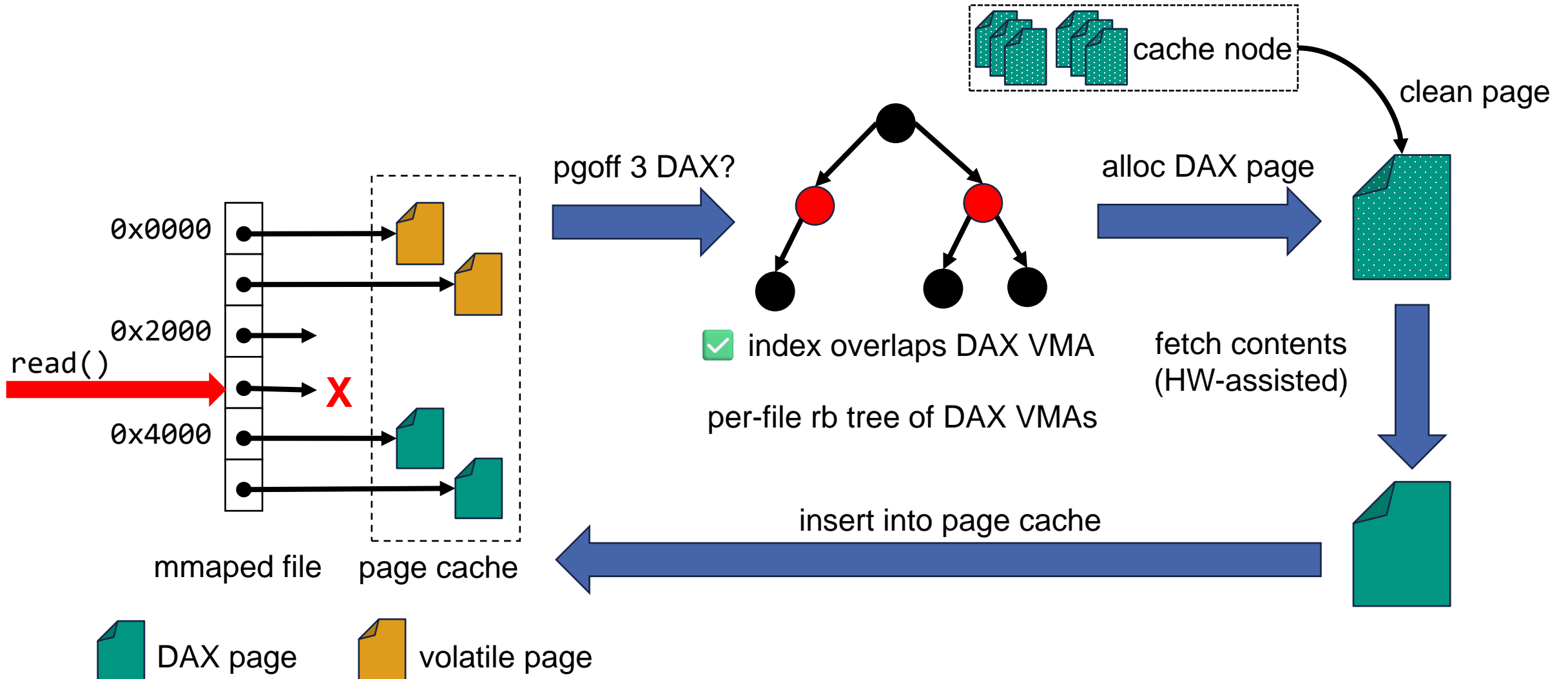
# DAX-Aware Page Cache Allocation

# DAX-Aware Page Cache Allocation



DAX page
volatile page

# DAX-Aware Page Cache Allocation



0x0000

0x2000

read()

0x4000

X

pgoff 3 DAX?

index overlaps DAX VMA

per-file rb tree of DAX VMAs

cache node

clean page

alloc DAX page

fetch contents
(HW-assisted)

insert into page cache

mmaped file    page cache

DAX page    volatile page

# DAX-Aware Page Cache Allocation



pgoff 3 DAX?

alloc DAX page

cache node

clean page

✅ index overlaps DAX VMA

per-file rb tree of DAX VMAs

fetch contents (HW-assisted)

insert into page cache

read()

0x0000

0x2000

0x4000

mmaped file

page cache

DAX page

volatile page

# Lightweight Synchronous Writeback

- Synchronous writeback critical for performance
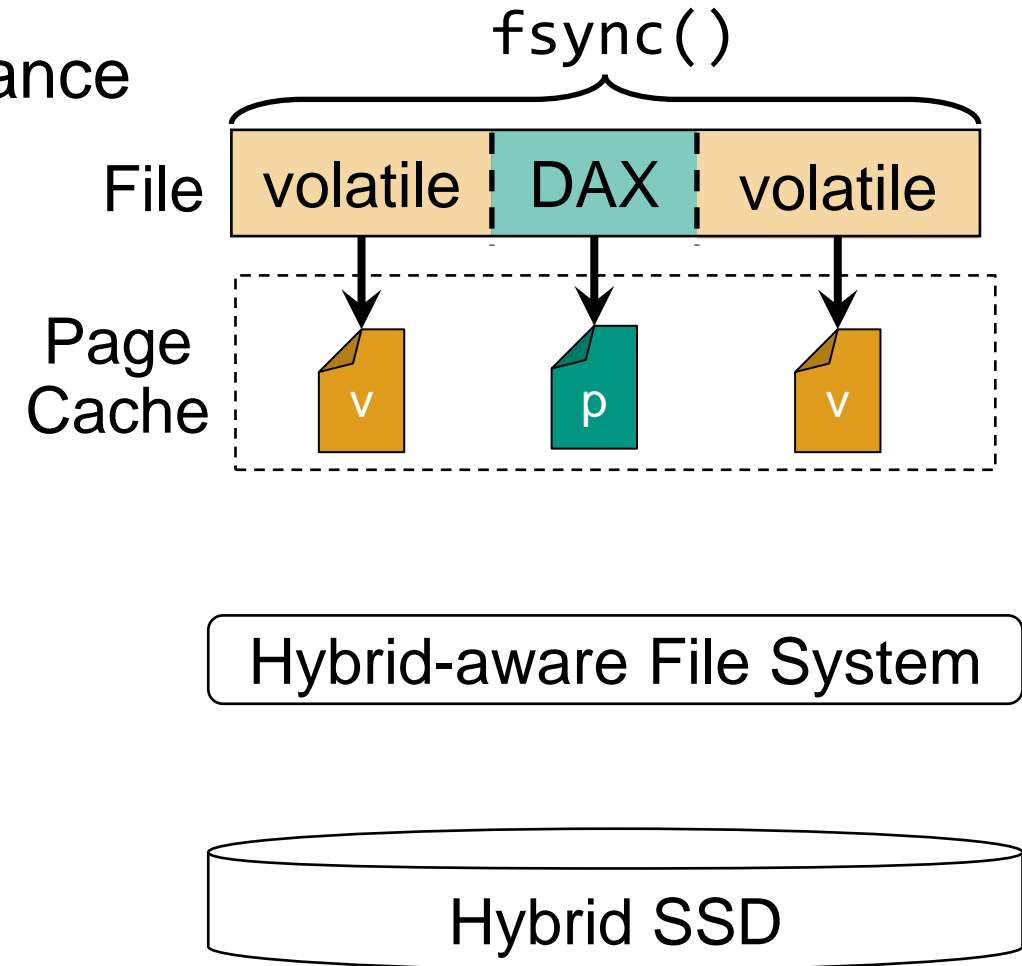  - On-device cache guarantees persistence
  → skip writeback of DAX pages
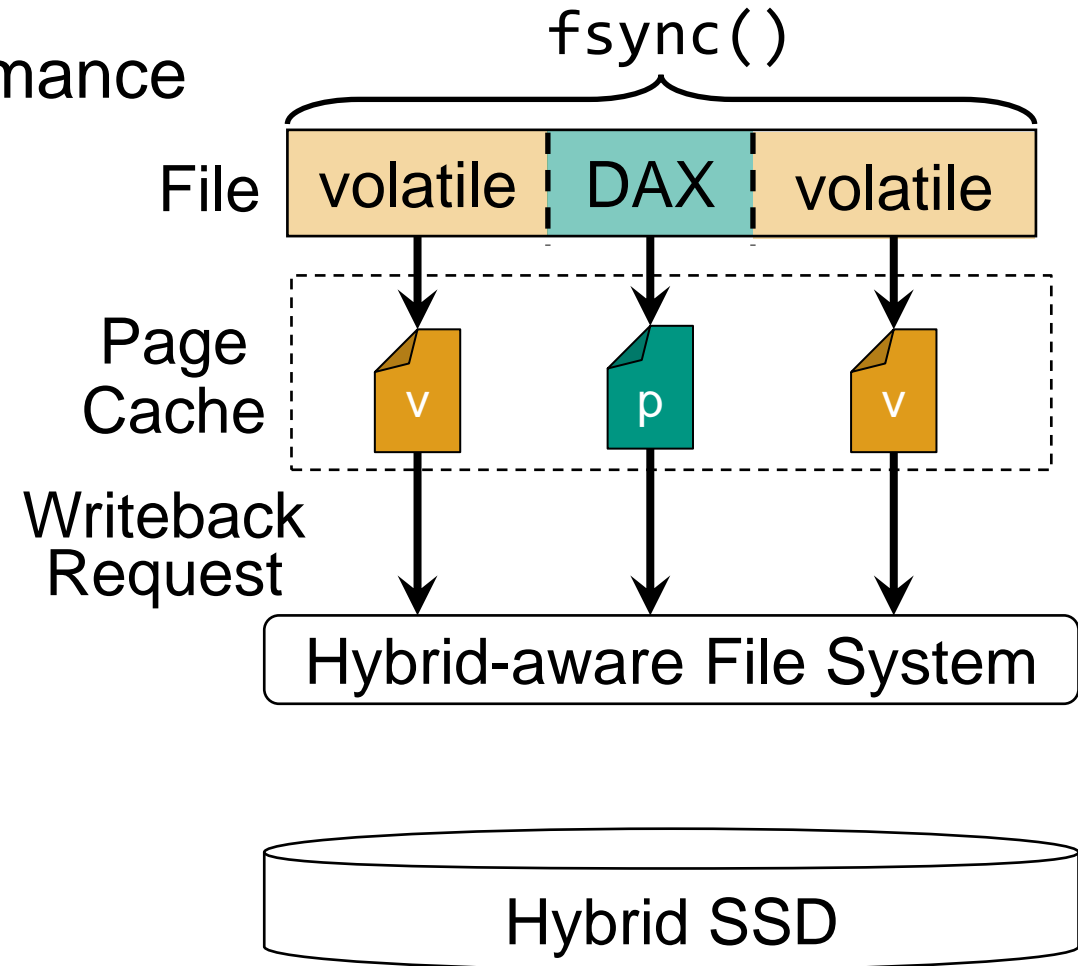  - DAX pages remain dirty

# Lightweight Synchronous Writeback

- Synchronous writeback critical for performance
  - On-device cache guarantees persistence
  → skip writeback of DAX pages
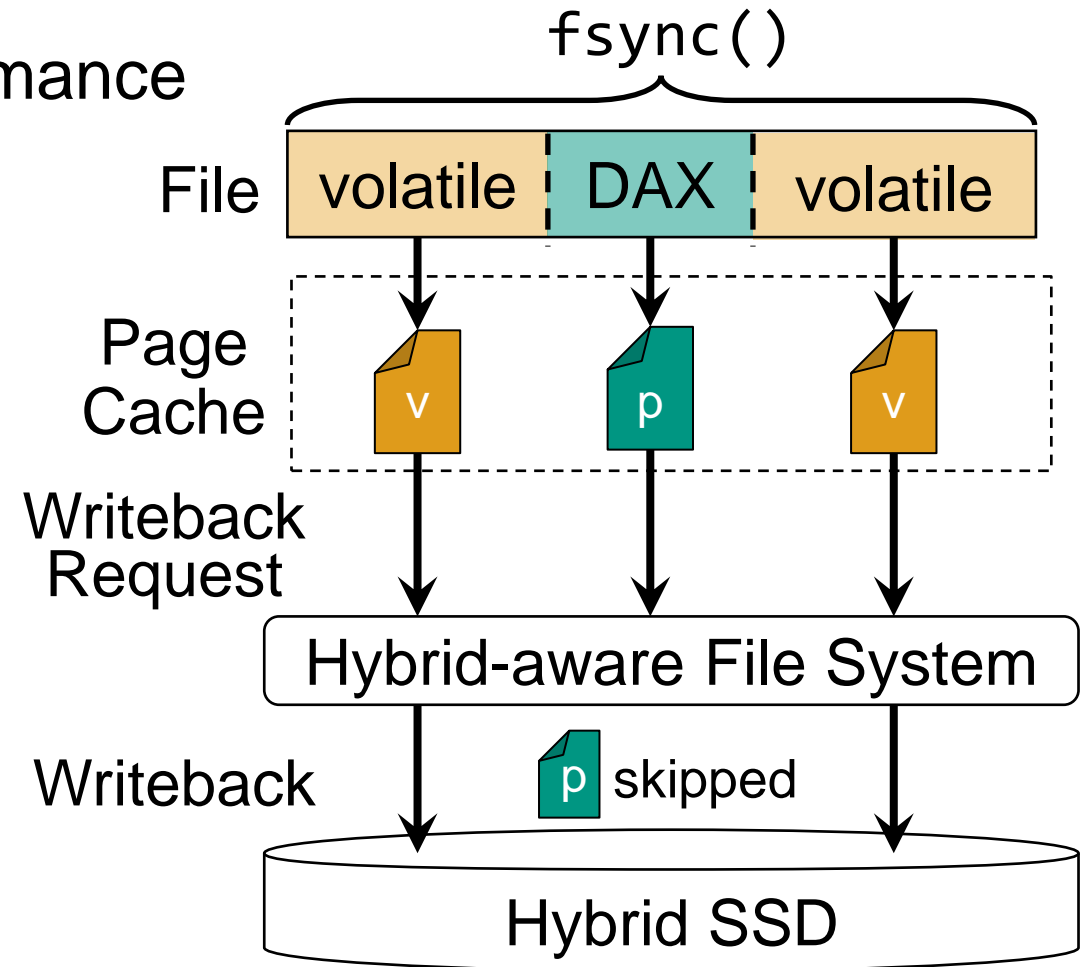  - DAX pages remain dirty

# Lightweight Synchronous Writeback

- Synchronous writeback critical for performance
  - On-device cache guarantees persistence
  → skip writeback of DAX pages
  - DAX pages remain dirty

# Lightweight Synchronous Writeback

- **Synchronous writeback critical for performance**
  - On-device cache guarantees persistence
    →skip writeback of DAX pages
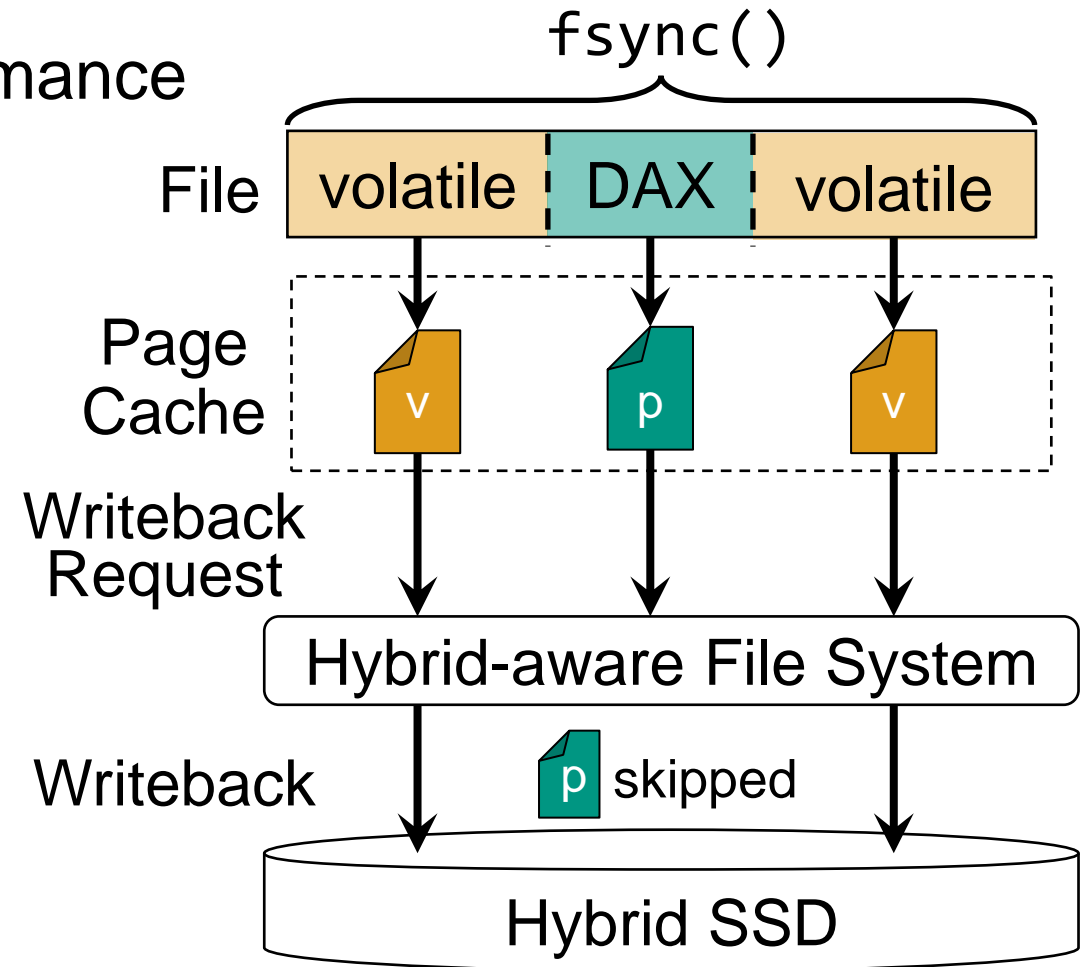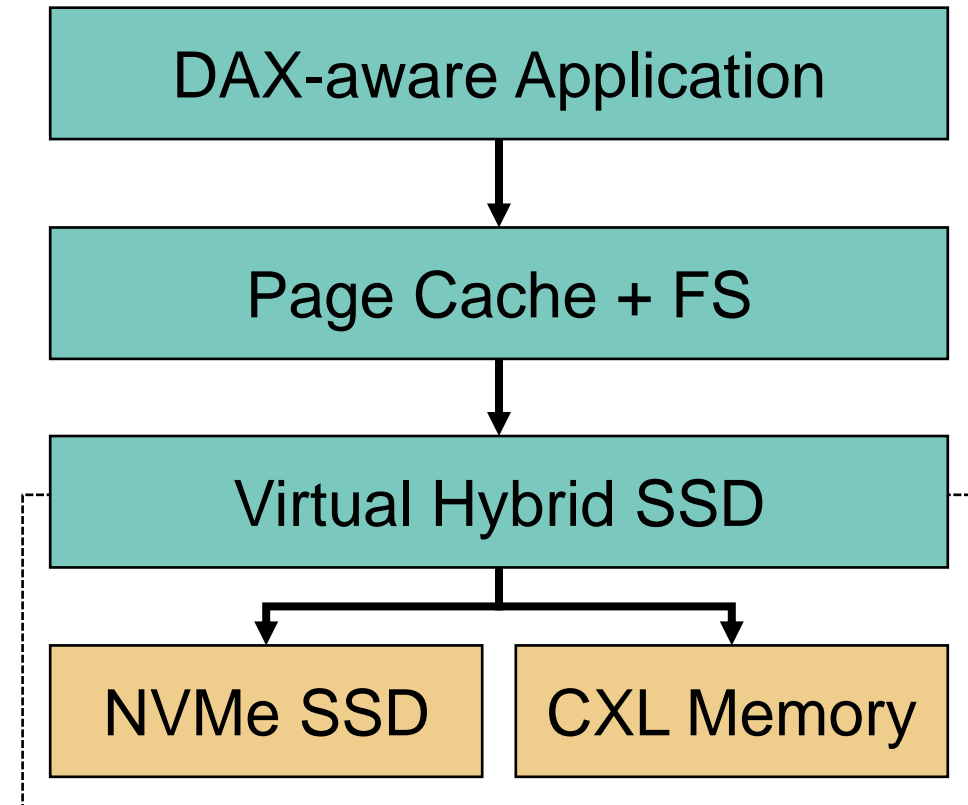  - DAX pages remain dirty

- **Asynchronous writeback unchanged**
  - Performance not critical
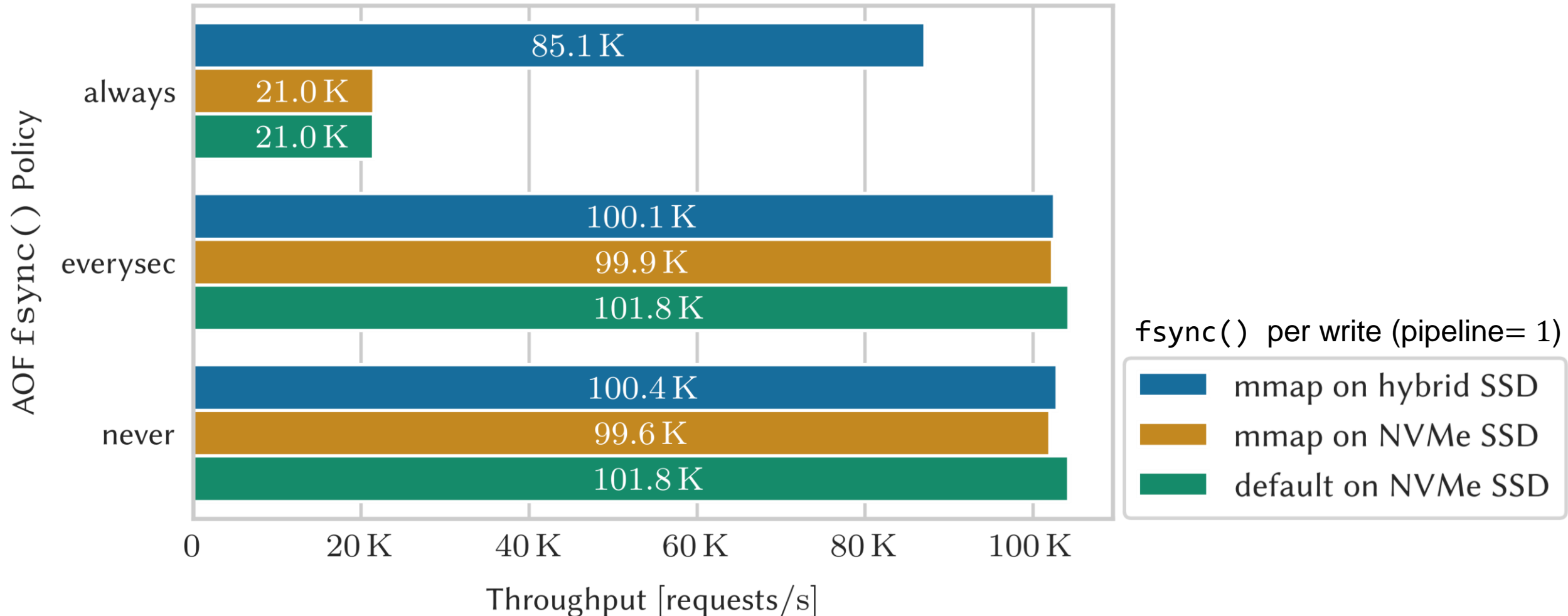  - Clean pages beneficial for reclaim

# Evaluation

- Emulated hybrid SSD (SSD + CXL mem)

- *Valkey* with Append-Only-File (AOF)
  - AOF writeback policy determines overhead
  - Write-only workload evaluated (worst-case)
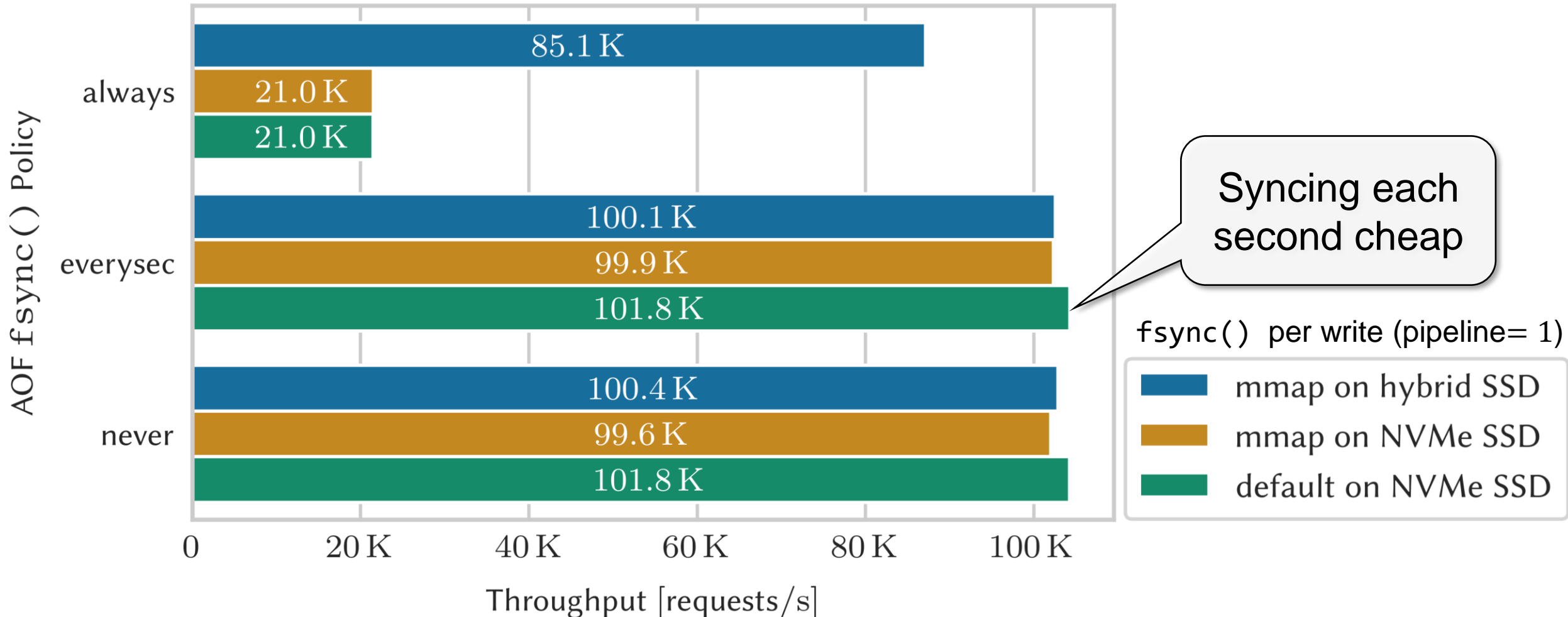  - mmap AOF backend for hybrid SSD

✅ <u>Throughput</u>, tail latencies, CPU and <u>energy efficiency</u> improved

```
┌─────────────────────────────┐
│   DAX-aware Application      │
└──────────────┬──────────────┘
               │
               ▼
┌─────────────────────────────┐
│      Page Cache + FS         │
└──────────────┬──────────────┘
               │
               ▼
┌─────────────────────────────┐
│     Virtual Hybrid SSD       │
└──────┬───────────────┬───────┘
       ▼               ▼
┌─────────────┐ ┌─────────────┐
│  NVMe SSD   │ │ CXL Memory  │
└─────────────┘ └─────────────┘
```
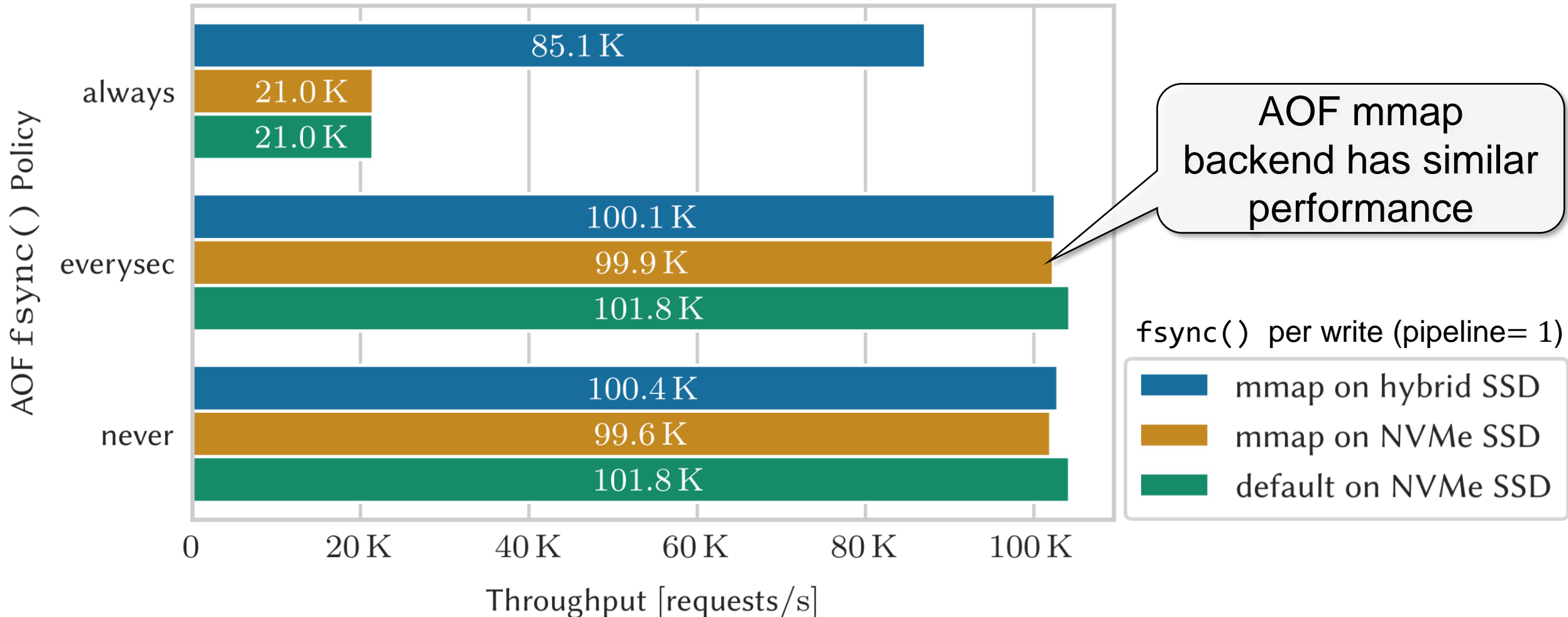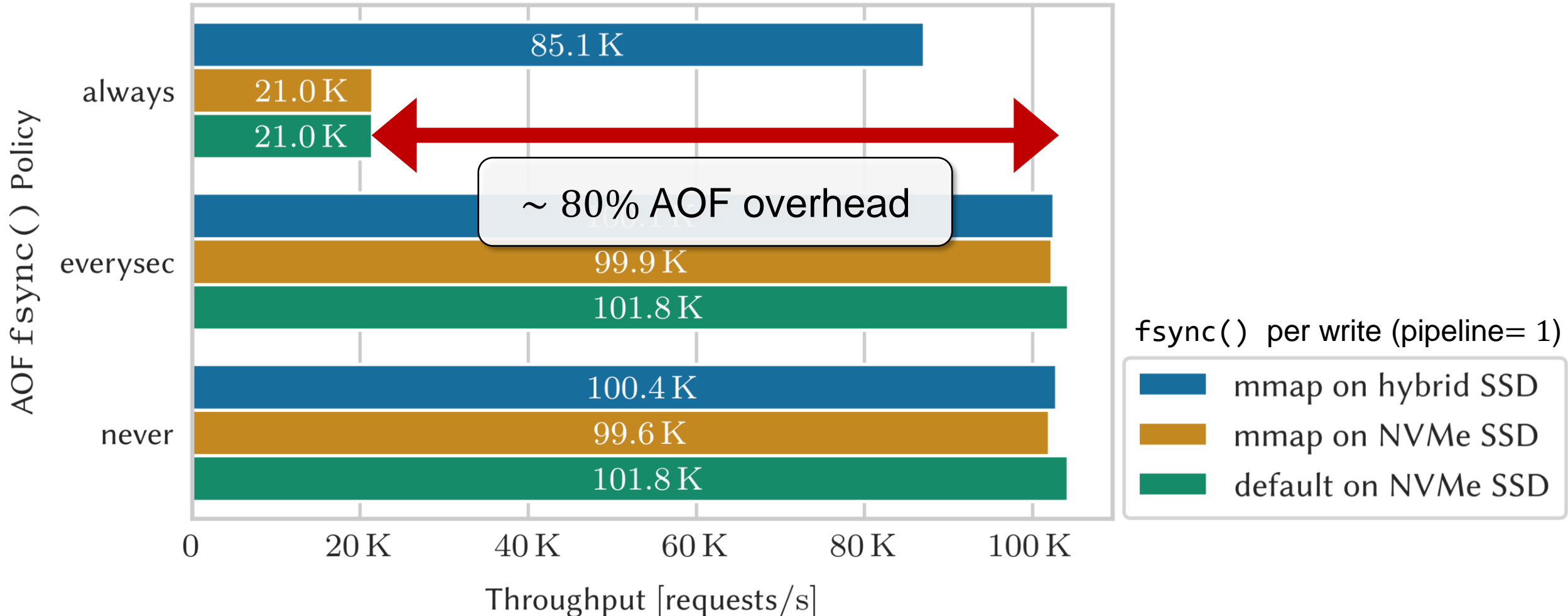
# Valkey AOF Throughput

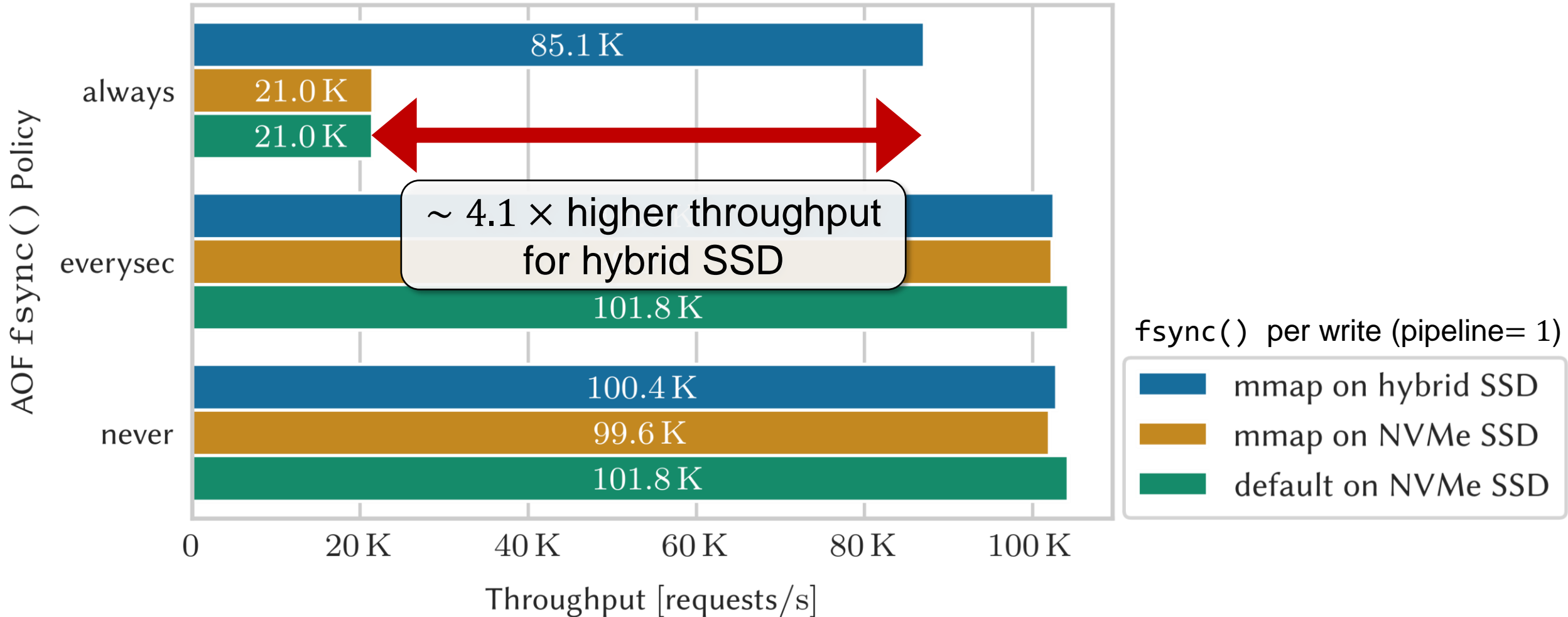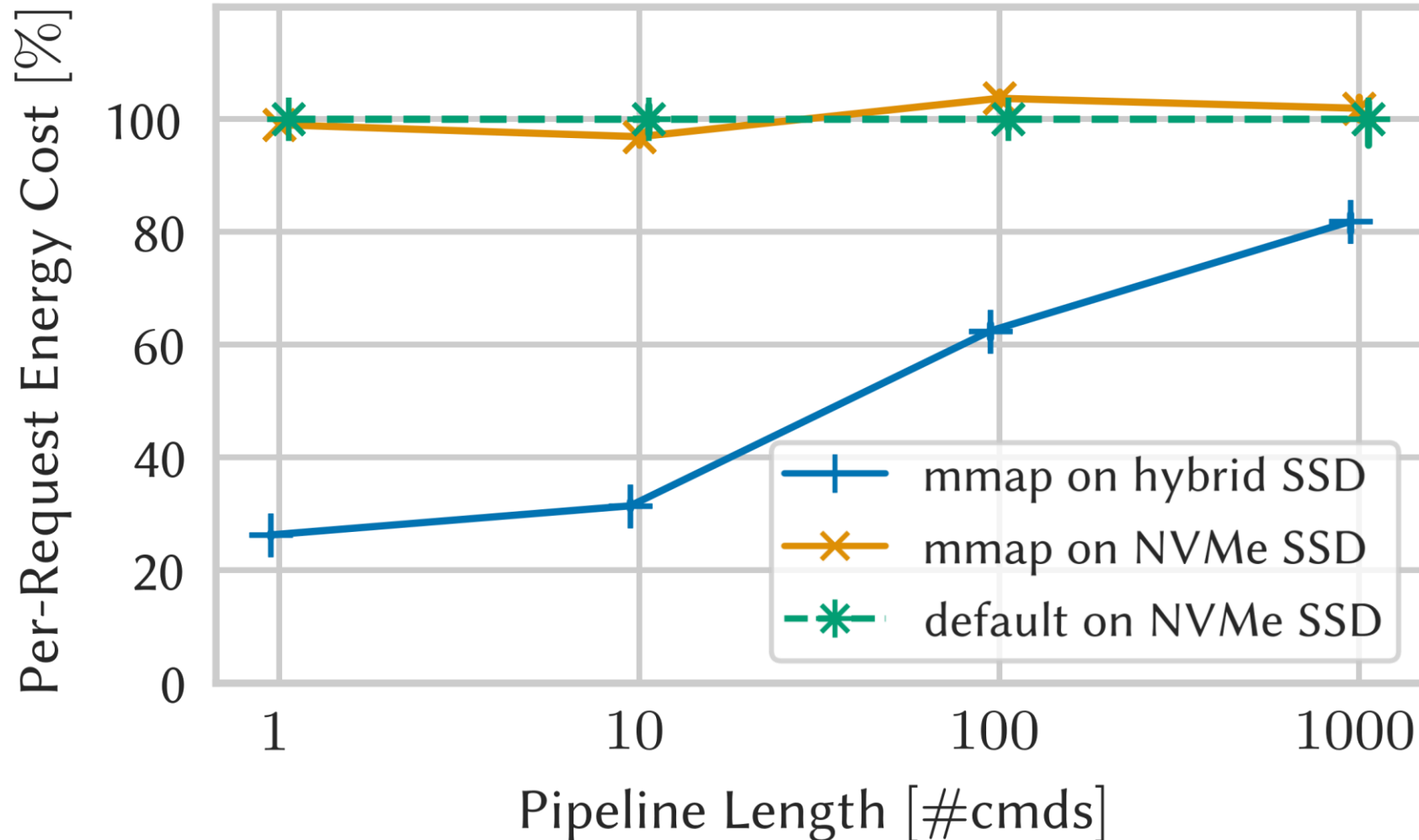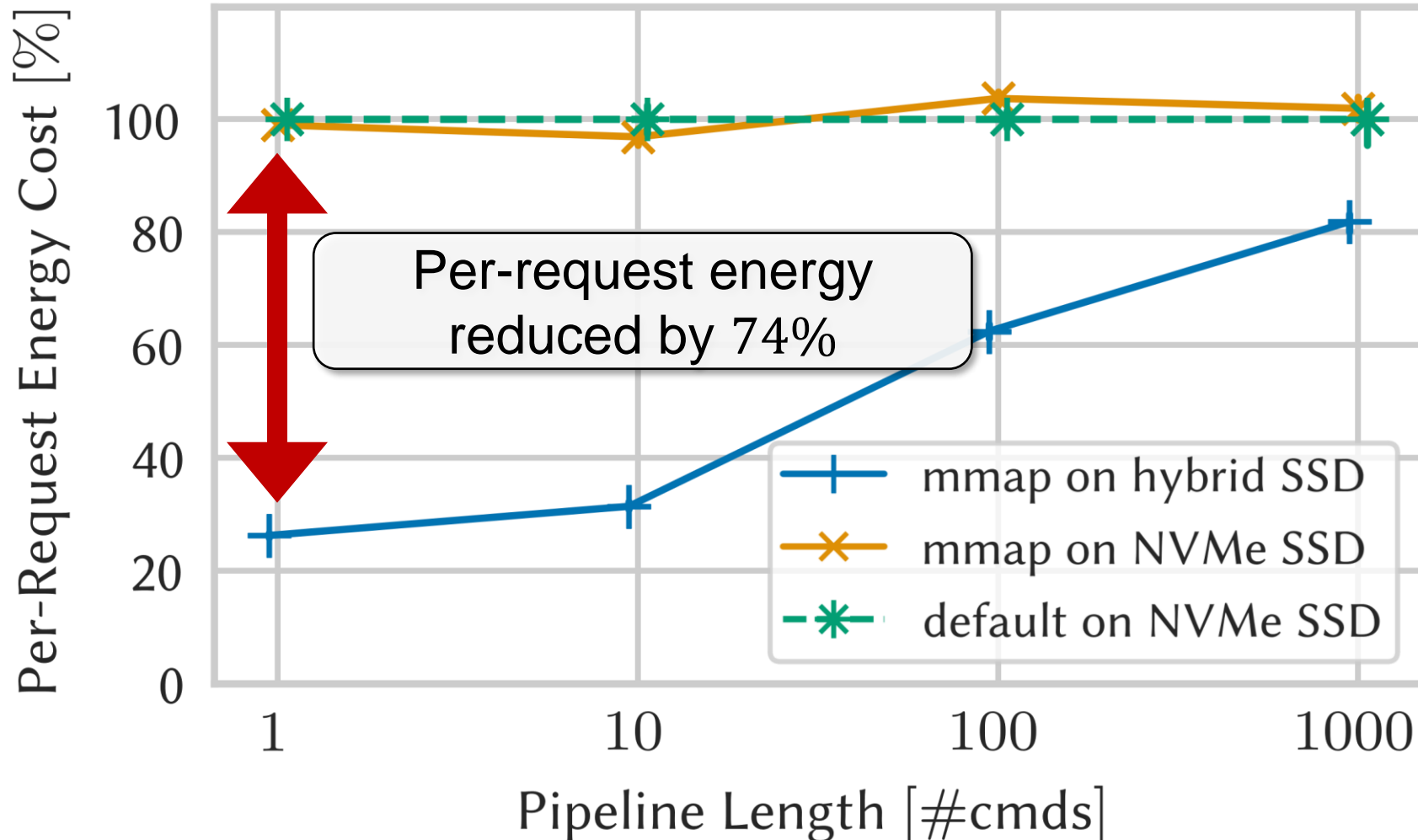# Valkey AOF Throughput

# Valkey AOF Throughput
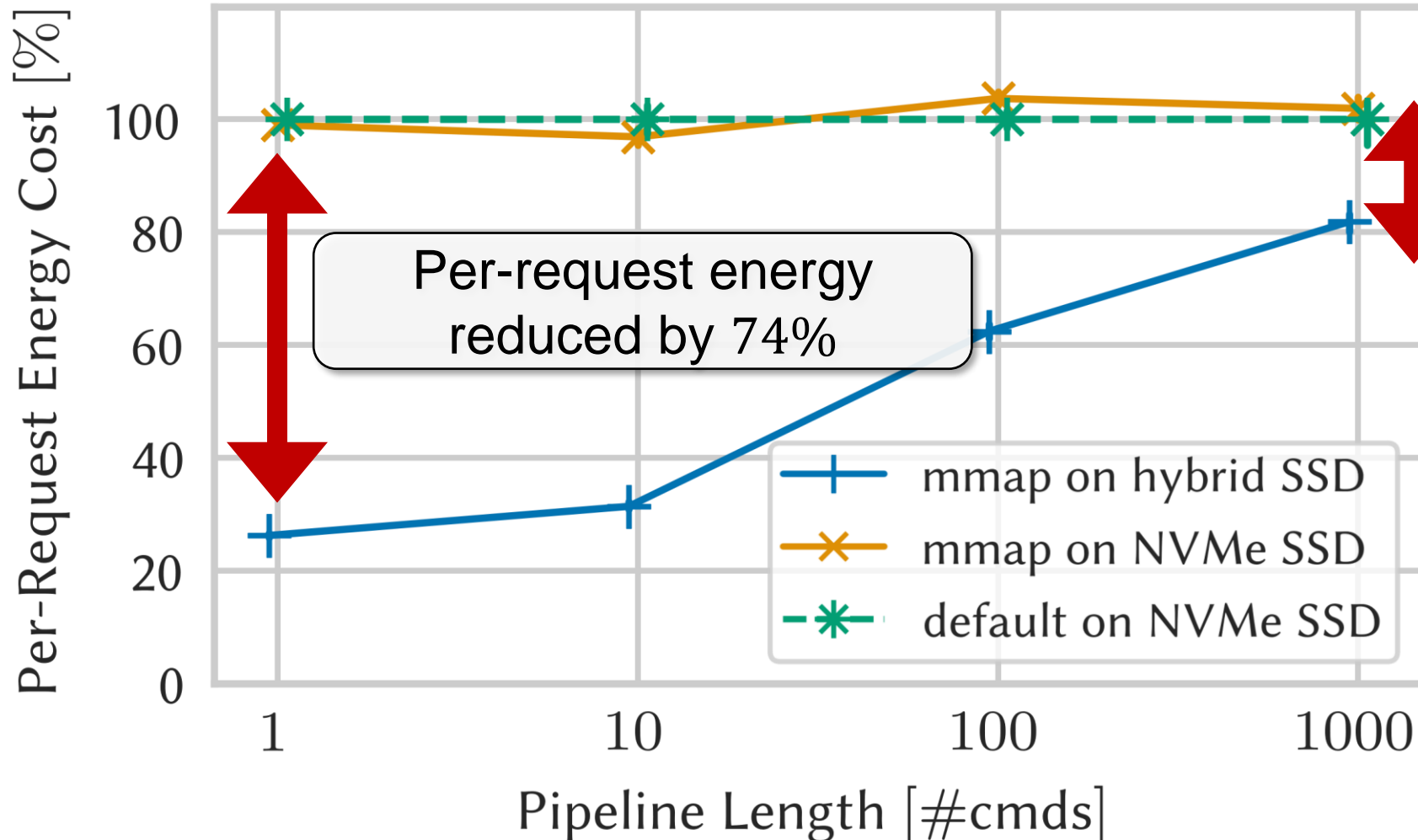
# Valkey AOF Throughput

# Valkey AOF Throughput

# Valkey Per-Request Energy Consumption

# Valkey Per-Request Energy Consumption



Per-request energy reduced by 74%

Per-Request Energy Cost [%] vs Pipeline Length [#cmds]

- mmap on hybrid SSD
- mmap on NVMe SSD
- default on NVMe SSD

# Valkey Per-Request Energy Consumption



Per-request energy reduced by 74%

Fewer `fsync()` calls → benefit of hybrid SSD decreases

Per-Request Energy Cost [%] vs Pipeline Length [#cmds]

- mmap on hybrid SSD
- mmap on NVMe SSD
- default on NVMe SSD

# Future Work

💡 Transparently establish of DAX mappings

💡 Study hardware design space for cache management

💡 Reevaluate on real-world hybrid SSDs

💡 Explore hybrid SSDs in consumer context

# Summary

- Hybrid SSD = NVMe + CXL.mem

- Existing OS abstractions unsuitable
  - Limited control over resource usage
  - CPU stalled on cache miss

- Our design:
  - Fine-granular resource management
  - Cache managed by OS

- Up to $4.1 \times$ higher *Valkey* throughput and 78% lower energy consumption