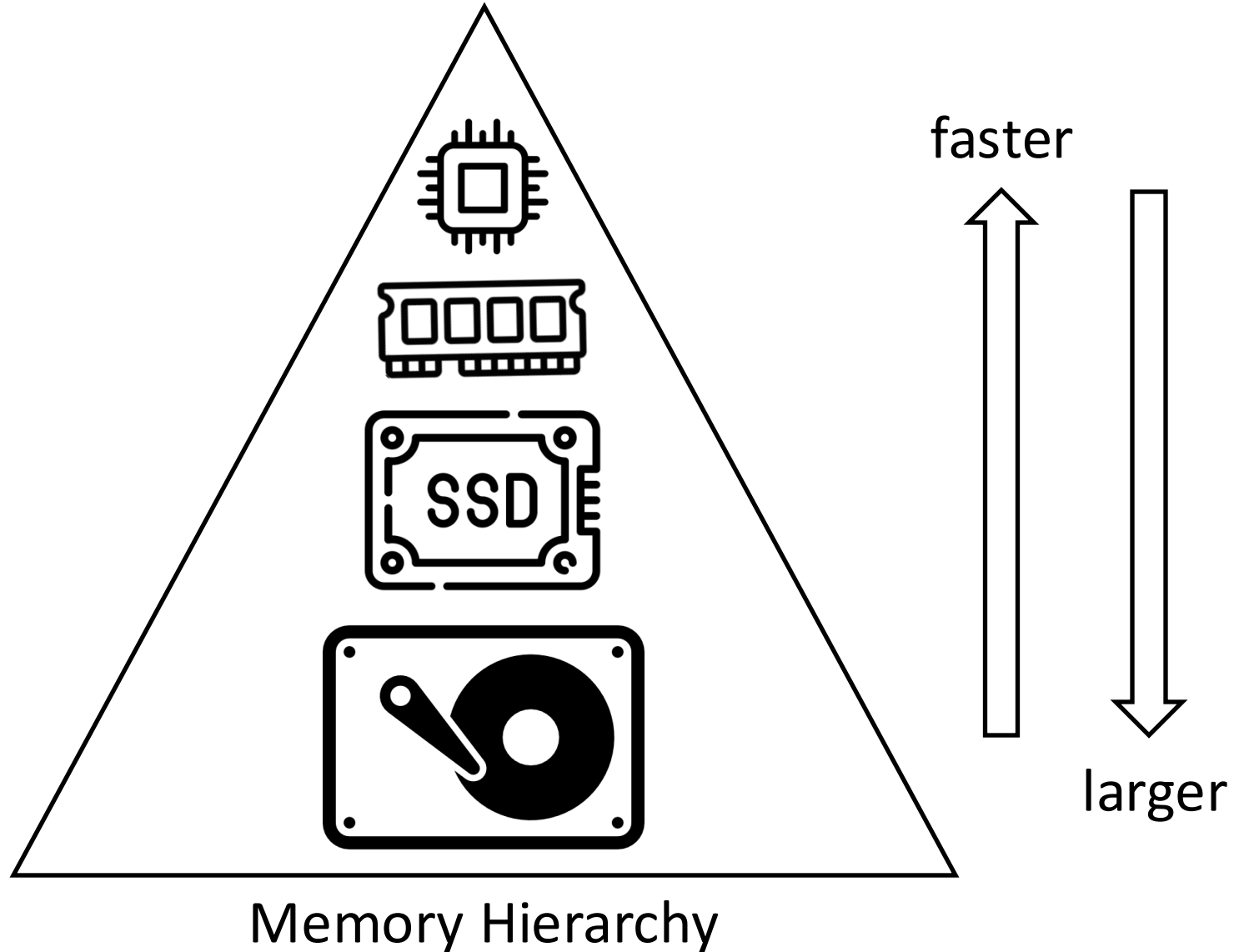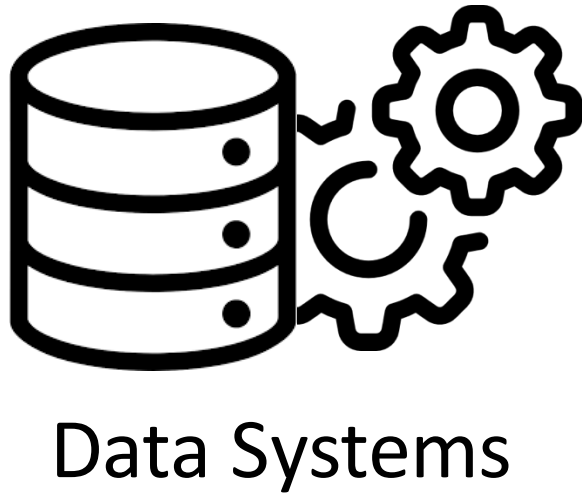# Asymmetry/Concurrency-Aware Bufferpool Manager for Modern Storage Devices
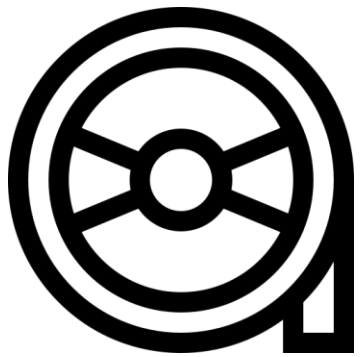
Tarikul Islam Papon

papon@bu.edu

Manos Athanassoulis

mathan@bu.edu

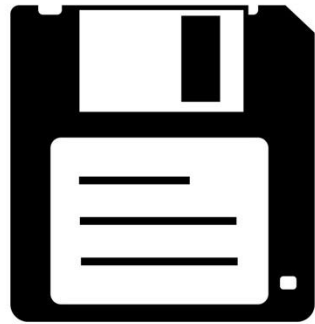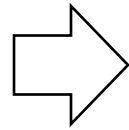# Data Systems & Hardware

Data Systems

Memory Hierarchy

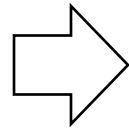faster

larger

# Evolution of Storage Devices
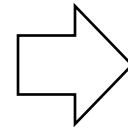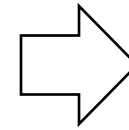
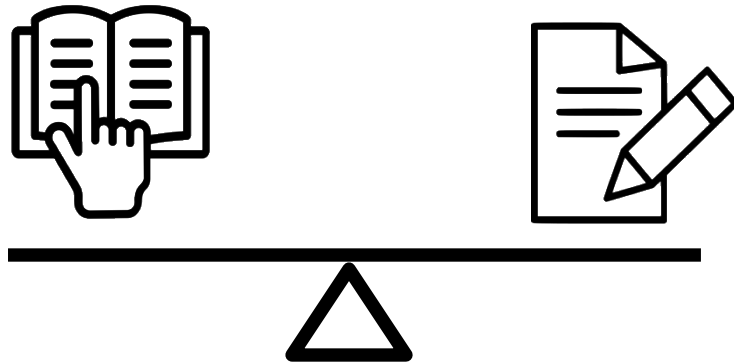Tape        Floppy        CD        HDD        SSD

# Hard Disk Drives



mechanical device

slow random access

**one block at a time**

**write latency ≈ read latency**

# Hard Disk Drives

**Symmetric** cost for **Read & Write** to disk

✓

**One I/O** at a time

✓

"Tape is Dead. Disk is Tape.

Flash is Disk."

- Jim Gray

# "Tape is Dead. Disk is Tape.
# Flash is Disk."

- Jim Gray

| Device | Size | Seq B/W | Time to read |
|--------|------|---------|--------------|
| HDD 1980 | 100 MB | 1.2 MB/s | ~ 1 min |
| HDD 2022 | 4 TB | 125 MB/s | ~ 9 hours |

# "Tape is Dead. Disk is Tape. Flash is Disk."

- Jim Gray

| Device | Size | Seq B/W | Time to read |
|--------|------|---------|--------------|
| HDD 1980 | 100 MB | 1.2 MB/s | ~ 1 min |
| HDD 2022 | 4 TB | 125 MB/s | ~ 9 hours |

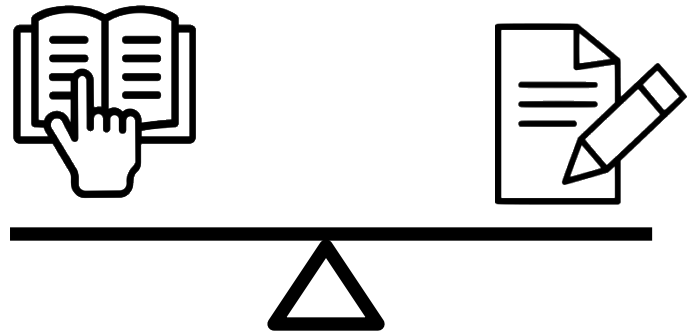HDDs are moving deeper in the memory hierarchy

# Solid State Drives



electronic device
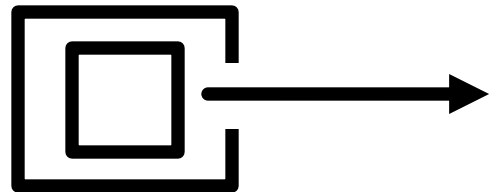
fast random access

**concurrent I/Os**

**write latency > read latency**

# HDD

# SSD

**Symmetric cost for Read & Write**

**Read/Write Asymmetry ($\alpha$)**

**One I/O at a time**

**Concurrency ($k$)**

# Concurrency

# Internals of an SSD

# Internals of an SSD



Parallelism at different levels (channel, chip, die, plane block, page)

# Read/Write Asymmetry

# Writes in SSD

Out-of-place updates cause invalidation

"*Erase before write*" approach



Plane

# Writes in SSD



Block 0

Block 1
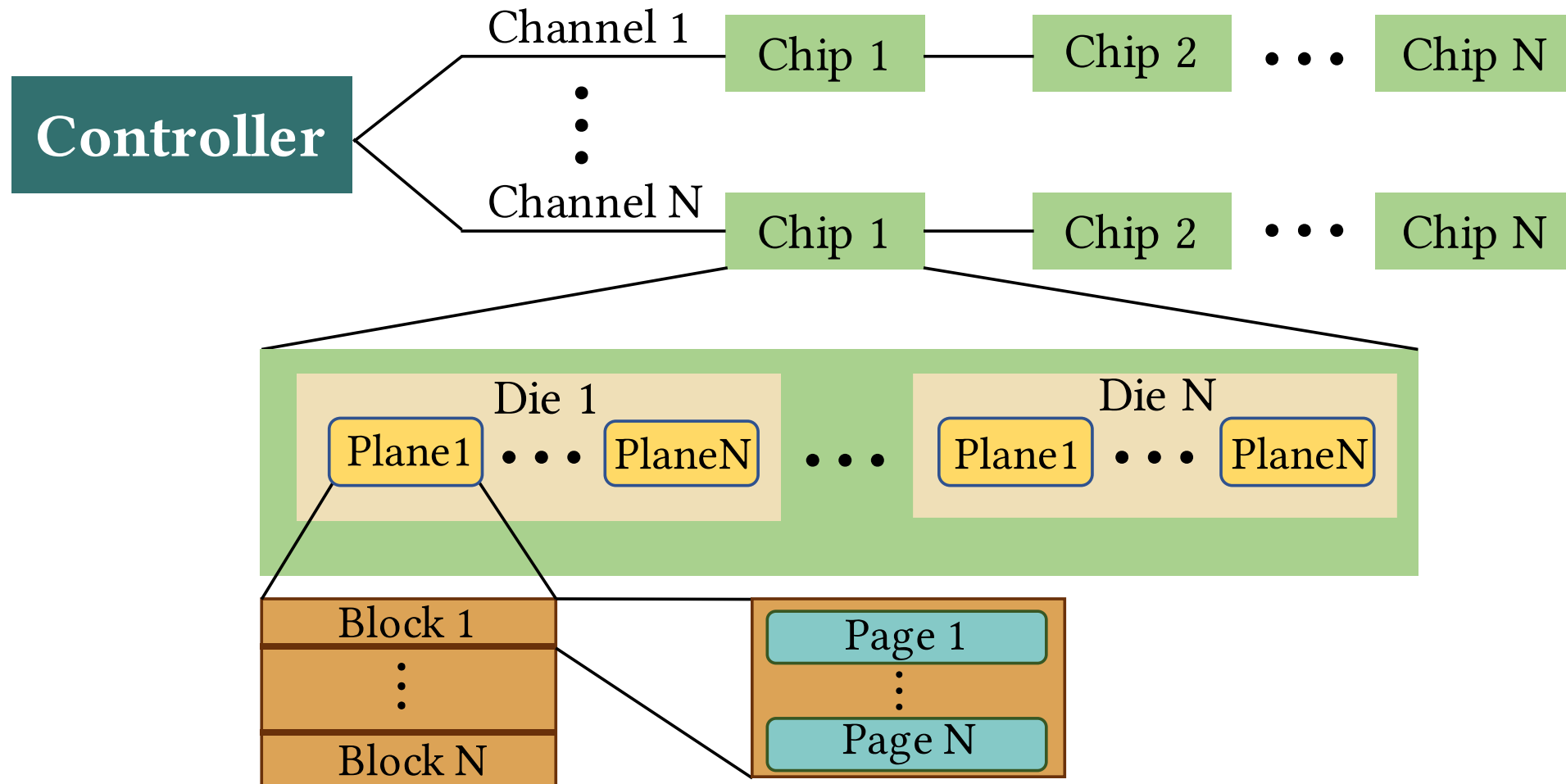
# Writes in SSD

| | | |
|---|---|---|
| A | B | C |
| D | E | F |
| G | H | Free |
| Free | Free | Free |

Block 0

| | | |
|---|---|---|
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |

Block 1

Writing in a free page isn't costly!

# Writes in SSD

Update

A, B, C, D



Block 0

Block 1

# Writes in SSD



Update

**A**, B, C, D

Block 0

Block 1

# Writes in SSD

Update

**A**, B, C, D

| | | |
|---|---|---|
| A | B | C |
| D | E | F |
| G | H | A' |
| Free | Free | Free |

Block 0

| | | |
|---|---|---|
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |

Block 1

# Writes in SSD



Update

A, **B**, C, D

Block 0

Block 1

# Writes in SSD

Update

A, **B**, C, D

| | | |
|---|---|---|
| A | B | C |
| D | E | F |
| G | H | A' |
| B' | Free | Free |

Block 0

| | | |
|---|---|---|
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |

Block 1

# Writes in SSD

Update

A, B, C, D

| | | |
|---|---|---|
| A | B | C |
| D | E | F |
| G | H | A' |
| B' | C' | D' |

Block 0

| | | |
|---|---|---|
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |

Block 1

Not all updates are costly!

23

# Writes in SSD

What if there is no space?



Block 0 ... Block N

# Writes in SSD

What if there is no space?

**Garbage Collection!**

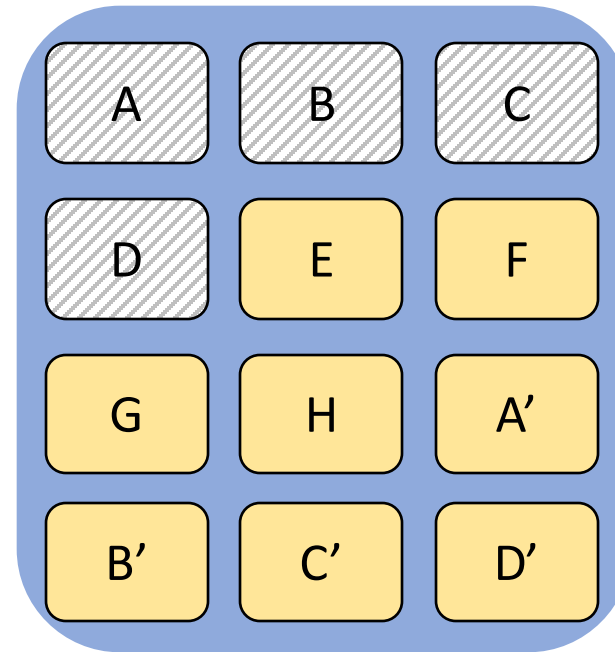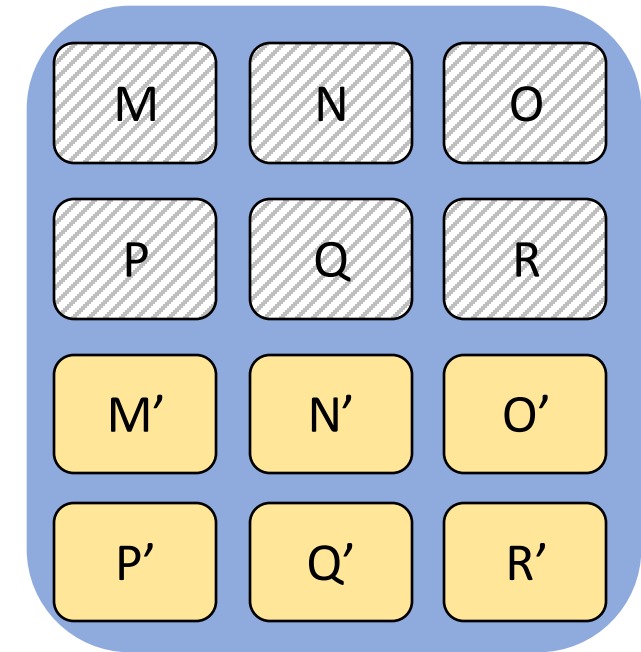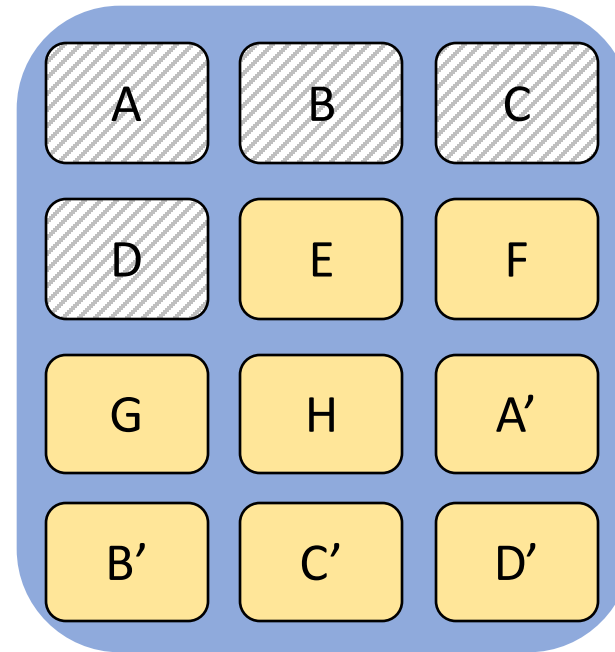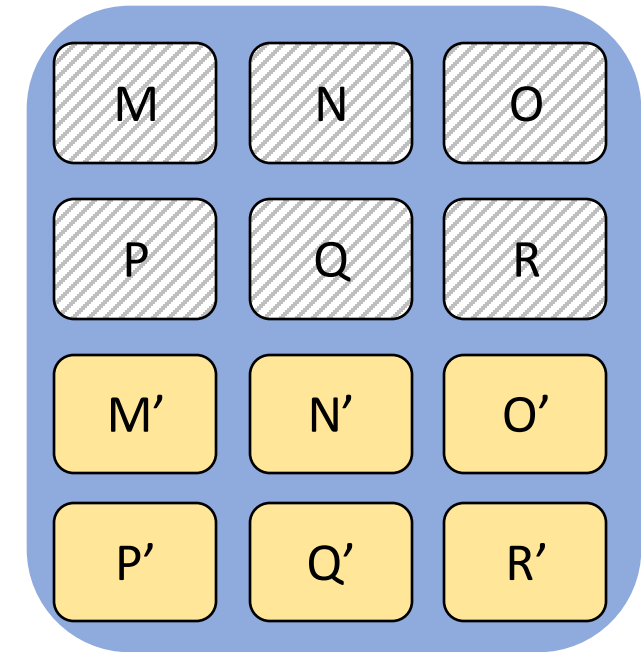| | | |
|---|---|---|
| A | B | C |
| D | E | F |
| G | H | A' |
| B' | C' | D' |

Block 0

| | | |
|---|---|---|
| M | N | O |
| P | Q | R |
| M' | N' | O' |
| P' | Q' | R' |

Block N

...

# Writes in SSD

What if there is no space?

**Garbage Collection!**

| | | |
|---|---|---|
| Erased | Erased | Erased |
| Erased | Erased | Erased |
| Erased | Erased | Erased |
| Erased | Erased | Erased |

| | | |
|---|---|---|
| Erased | Erased | Erased |
| Erased | Erased | Erased |
| Erased | Erased | Erased |
| Erased | Erased | Erased |

Block 0 ... Block N

Valid pages: | E | F | G | H | A' | B' | C' | D' | M' | N' | O' | P' | Q' | R' |

# Writes in SSD

What if there is no space?



**Garbage Collection!**

| | | |
|---|---|---|
| E | F | G |
| H | A' | B' |
| C' | D' | M' |
| N' | O' | P' |

| | | |
|---|---|---|
| Q' | R' | Free |
| Free | Free | Free |
| Free | Free | Free |
| Free | Free | Free |

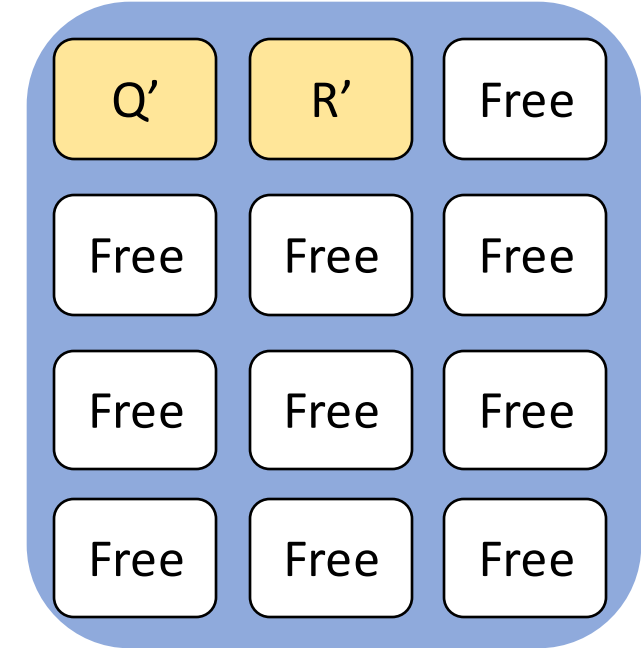Block 0                   …                   Block N

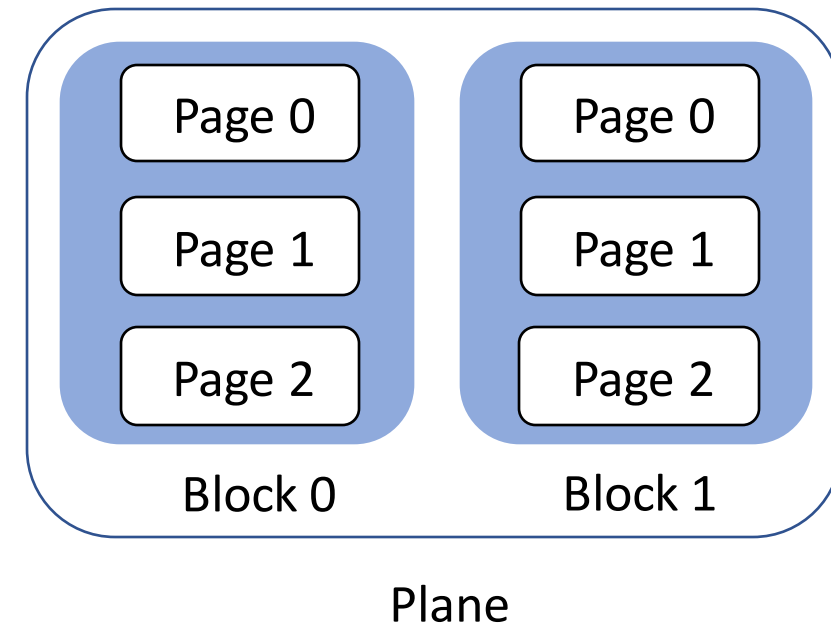Higher average update cost (due to GC) → *Read/Write asymmetry*

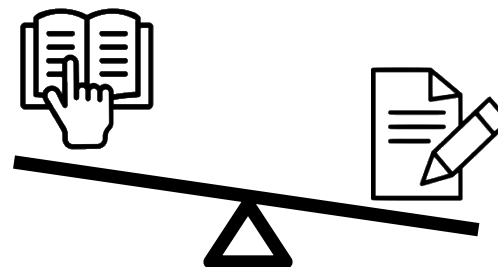# Read/Write Asymmetry

Out-of-place updates cause invalidation

"*Erase before write*" approach

Garbage Collection

Larger erase granularity
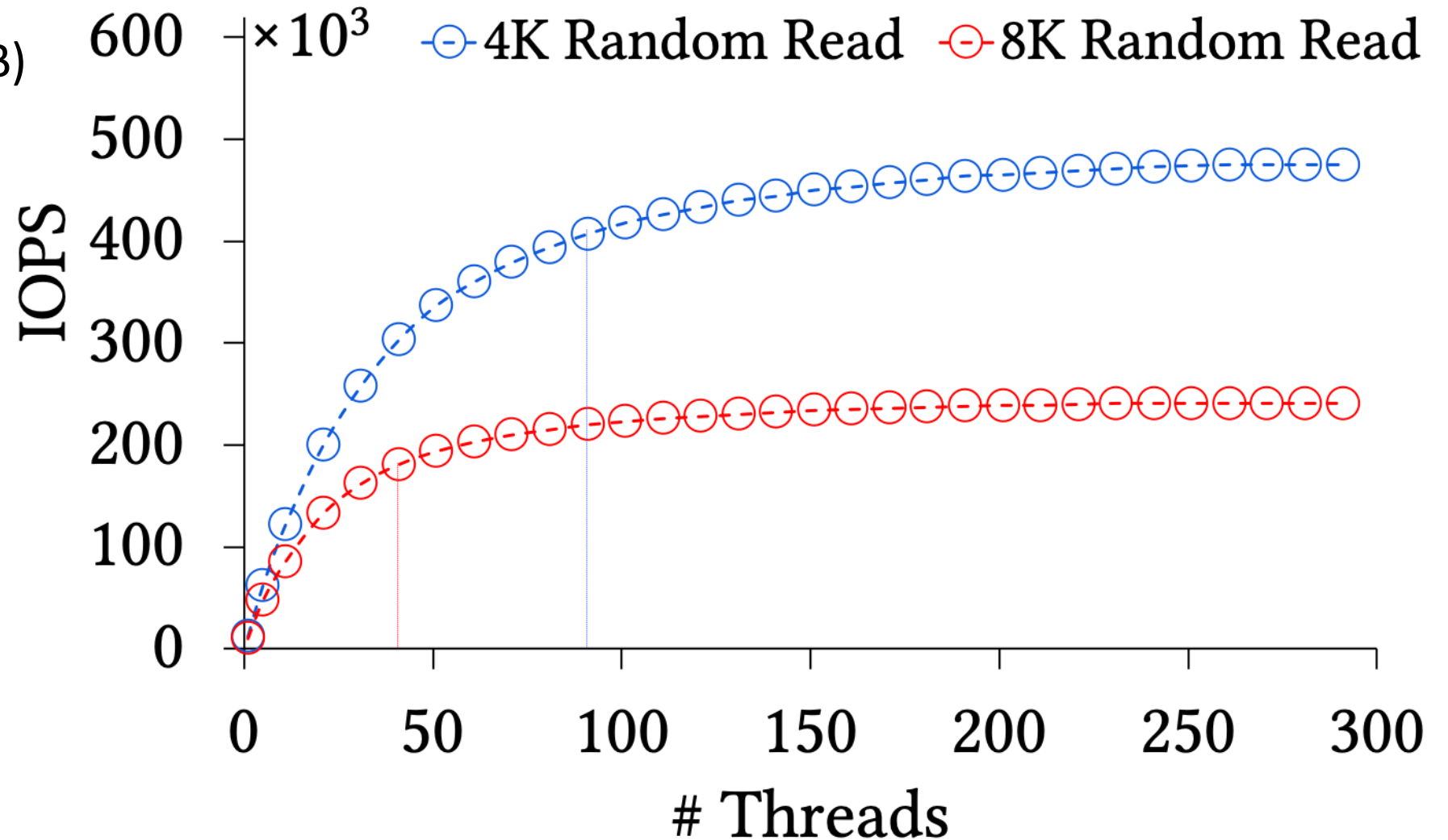
All these results in higher amortized write cost



Block 0 Block 1

Plane

# Read/Write Asymmetry - Example

| Device | Advertised Rand Read IOPS | Advertised Rand Write IOPS | Advertised Asymmetry |
|---|---|---|---|
| PCIe D5-P4320 | 427k | 36k | 11.9 |
| PCIe DC-P4500 | 626k | 51k | 12.3 |
| PCIe P4510 | 465k | 145k | 3.2 |
| SATA D3-S4610 | 92k | 28k | 3.3 |
| Optane P4800X | 550k | 500k | 1.1 |

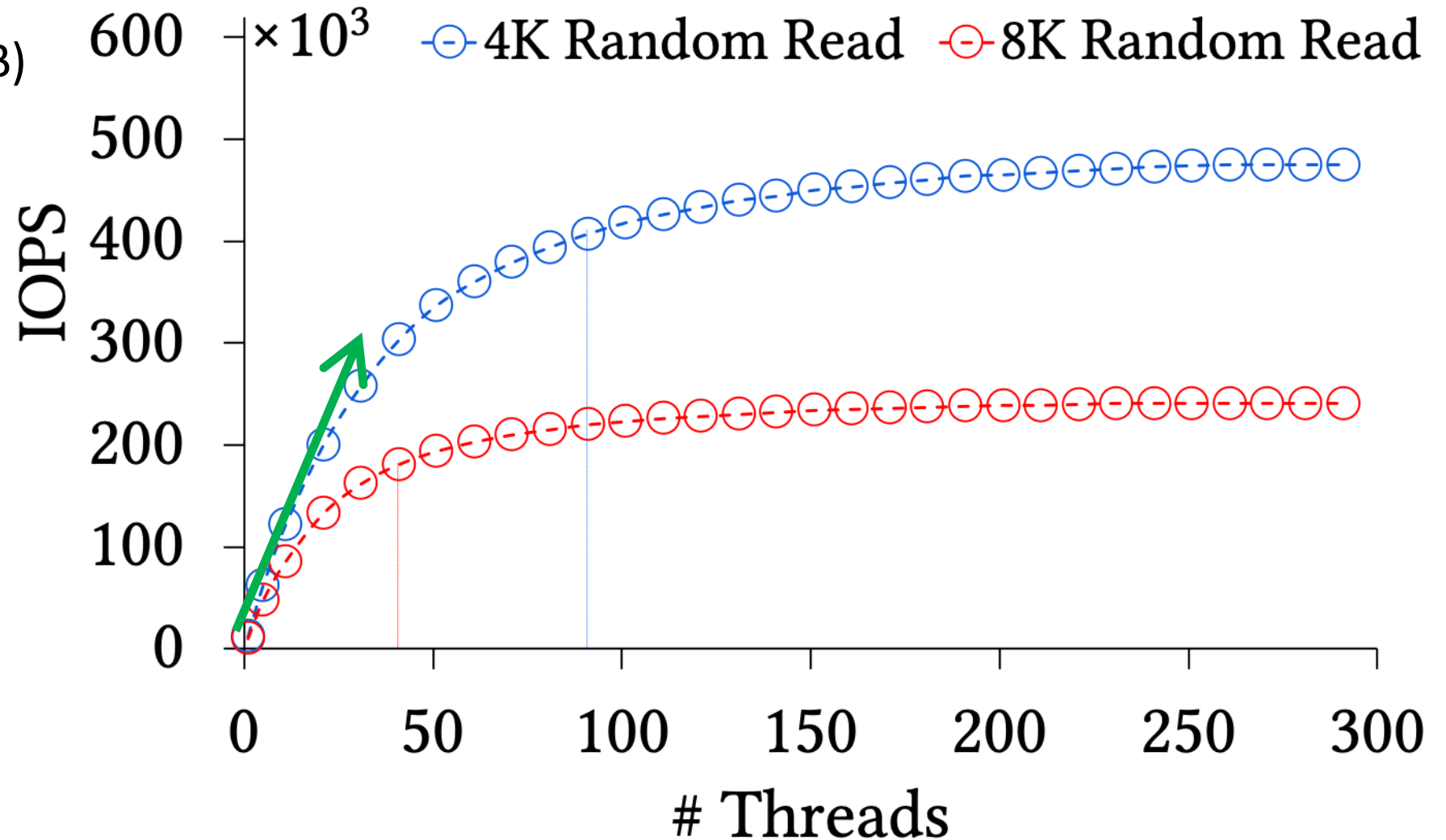# Quantifying Asymmetry & Concurrency

**Device**
PCIe SSD - P4510 (1TB)

# Quantifying Asymmetry & Concurrency

**Device**
PCIe SSD - P4510 (1TB)

# Quantifying Asymmetry & Concurrency

**Device**
PCIe SSD - P4510 (1TB)

# Quantifying Asymmetry & Concurrency

**Device**
PCIe SSD - P4510 (1TB)



read concurrency
($k_r$) = 80
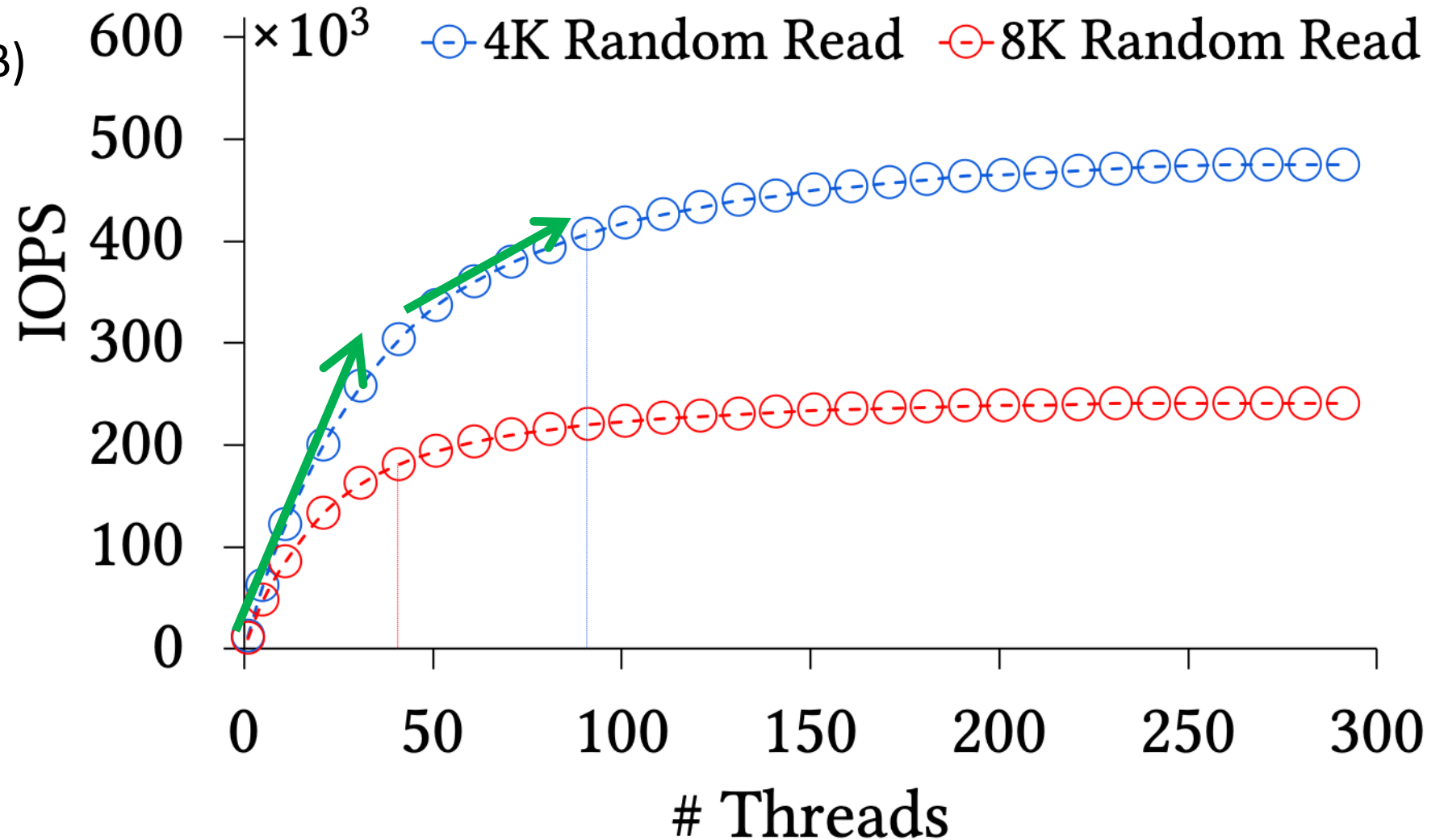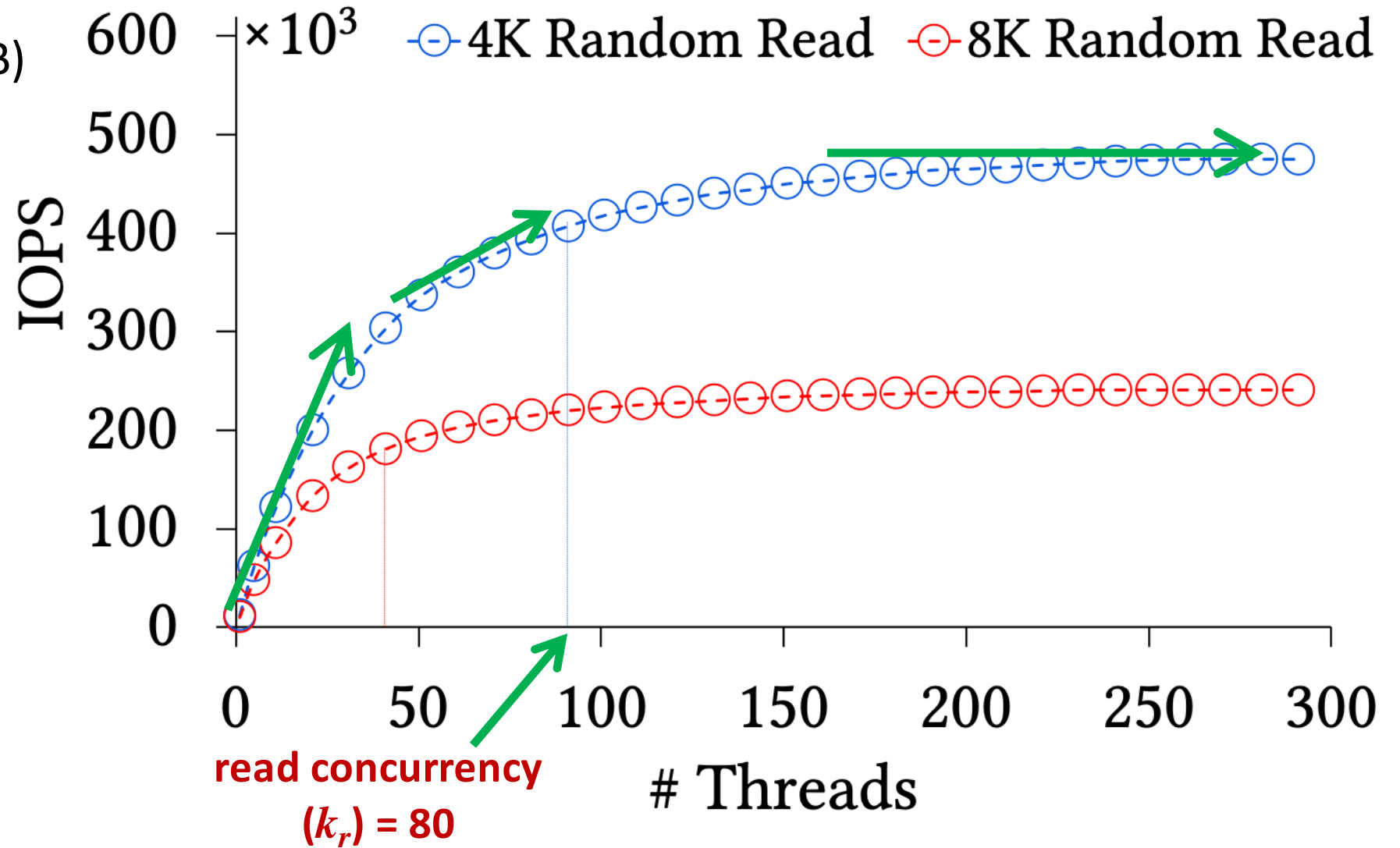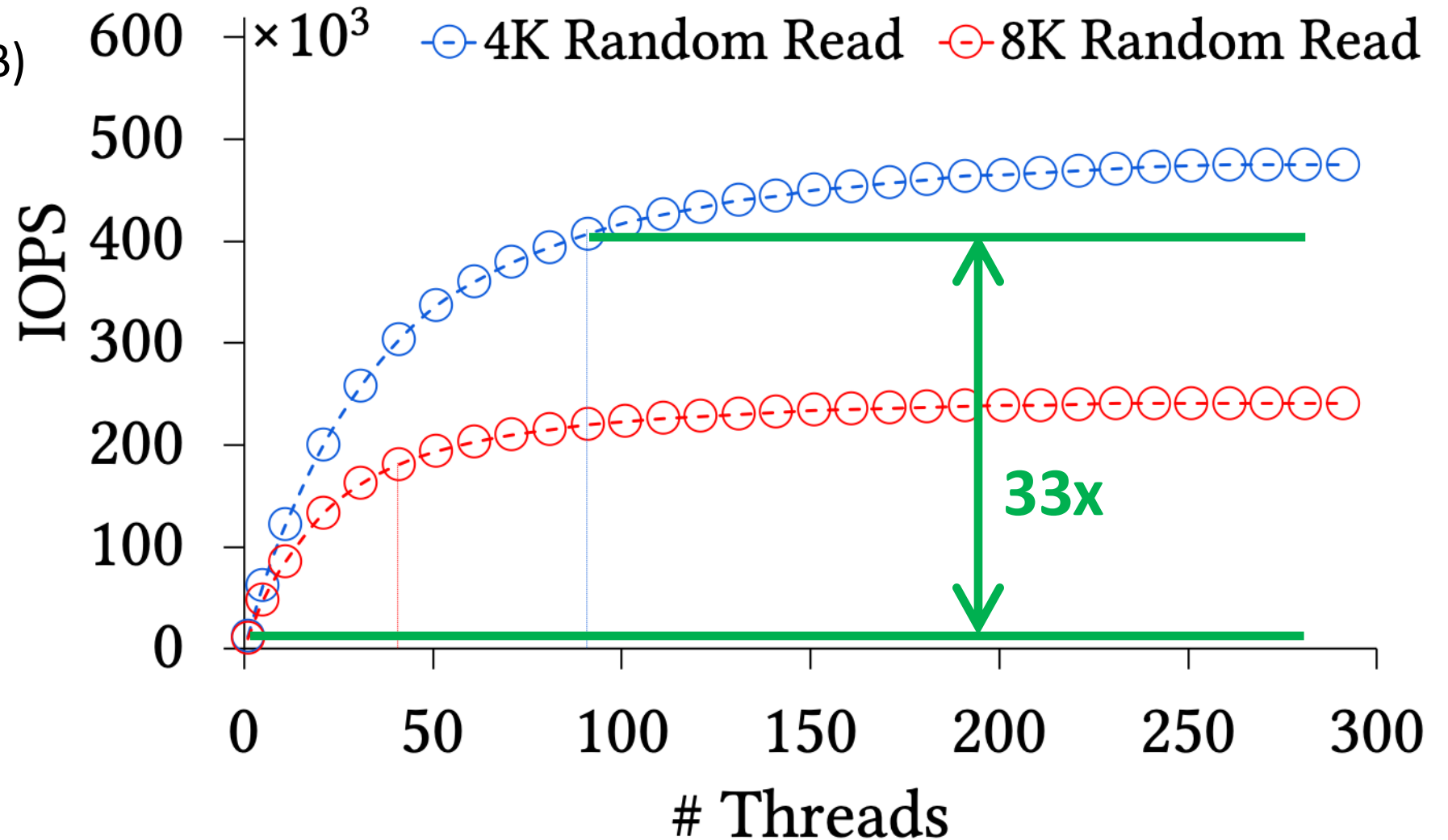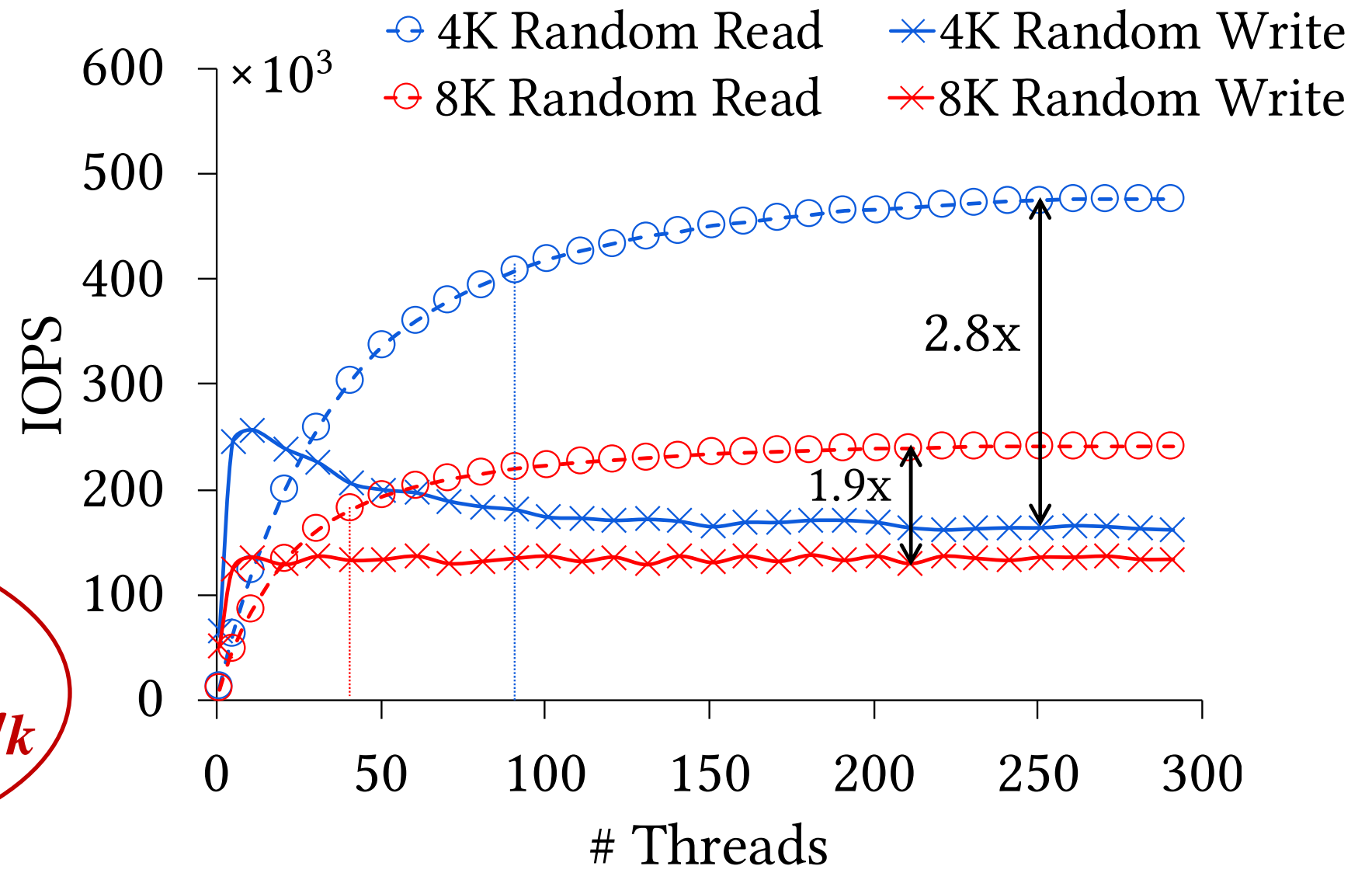
# Quantifying Asymmetry & Concurrency

**Device**
PCIe SSD - P4510 (1TB)

# Quantifying Asymmetry & Concurrency

**Device**
PCIe SSD - P4510 (1TB)

For 4K random read,

**Asymmetry**: 2.8

**Concurrency**: 80

**Yet, systems are not always tailored for $\alpha/k$**



- ⊖ 4K Random Read
- ✕ 4K Random Write
- ⊖ 8K Random Read
- ✕ 8K Random Write

IOPS ($\times 10^3$)

2.8x

1.9x

# Threads

# Empirical Asymmetry and Concurrency

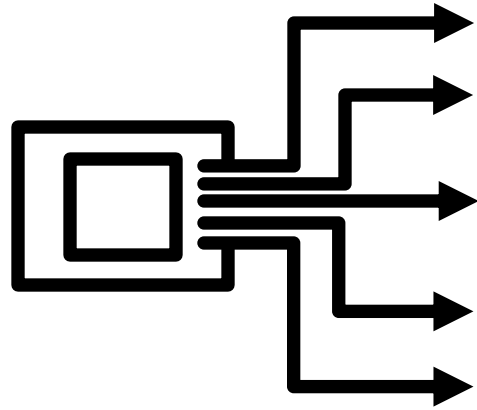| Device | $\alpha$ | $k_r$ | $k_w$ |
|---|---|---|---|
| Optane SSD | 1.1 | 6 | 5 |
| PCIe SSD | 2.8 | 80 | 8 |
| SATA SSD | 1.5 | 25 | 9 |
| Virtual SSD | 2.0 | 11 | 19 |

- "A Parametric I/O Model for Modern Storage Devices", DaMoN 2021
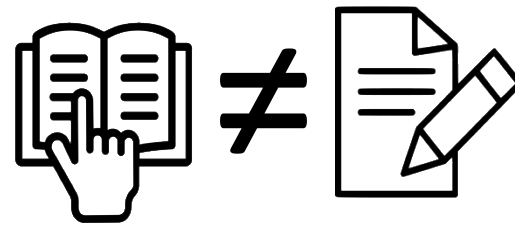
disc.bu.edu/papers/damon21-papon
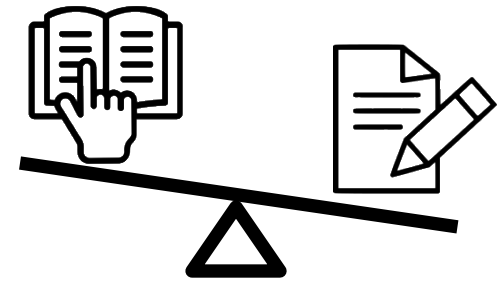
# Guidelines for Algorithm Design

Know Thy Device

Exploit concurrency
(with care)
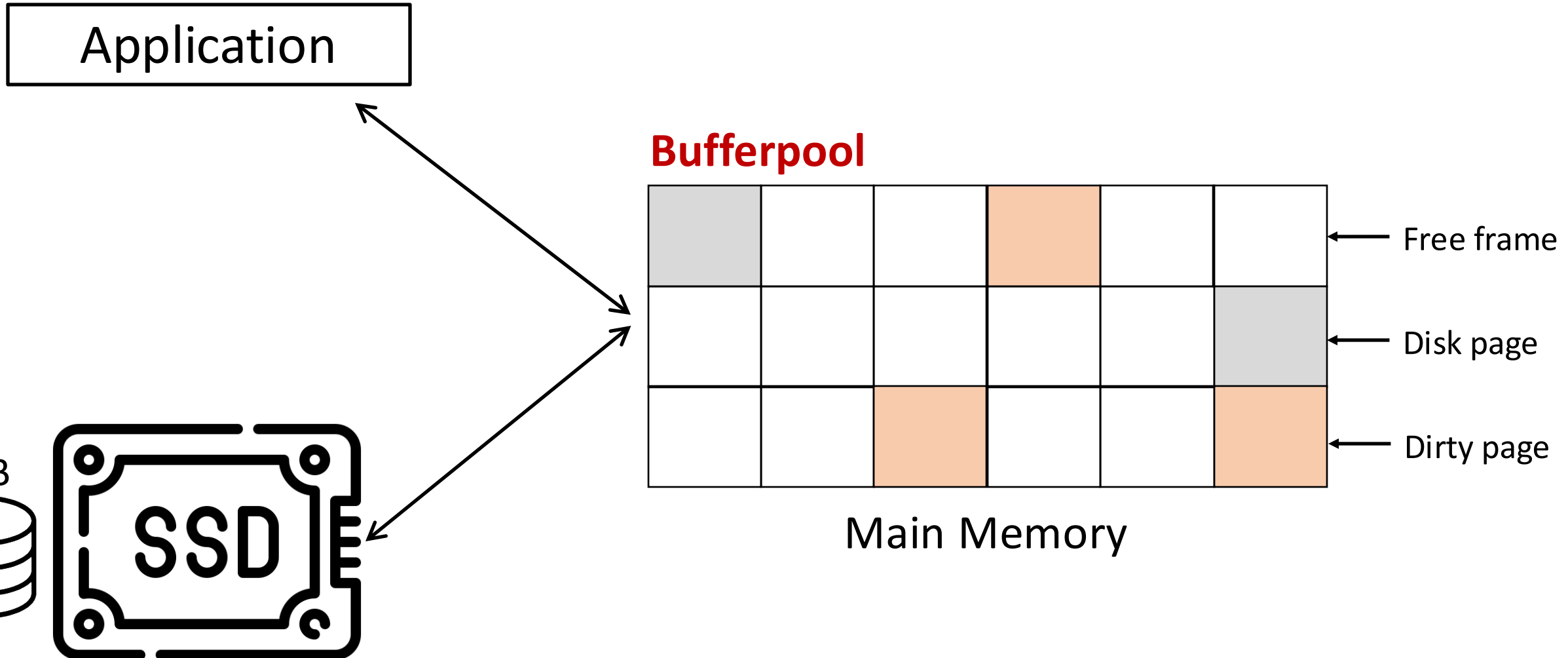
Treat read and
write differently.

asymmetry controls
performance

- "A Parametric I/O Model for Modern Storage Devices", DaMoN 2021
disc.bu.edu/papers/damon21-papon

# Bufferpool Manager &

# The Challenge
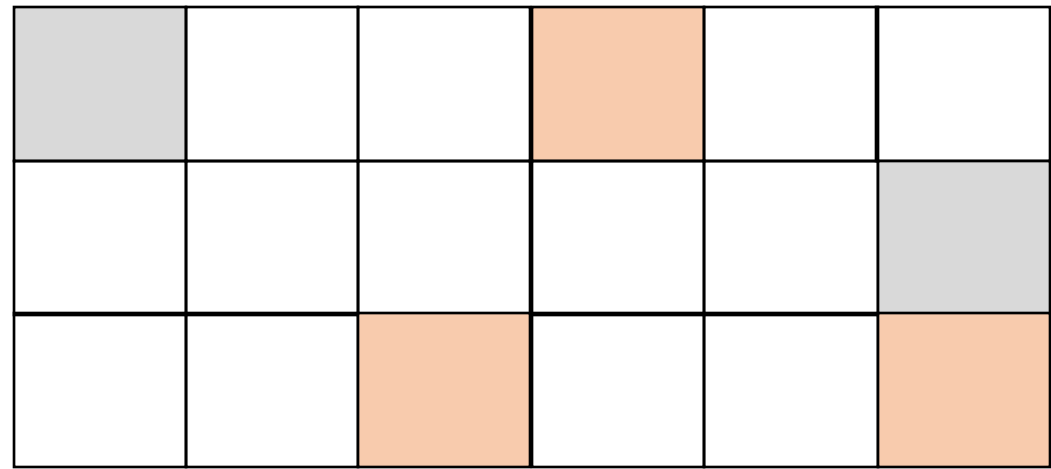
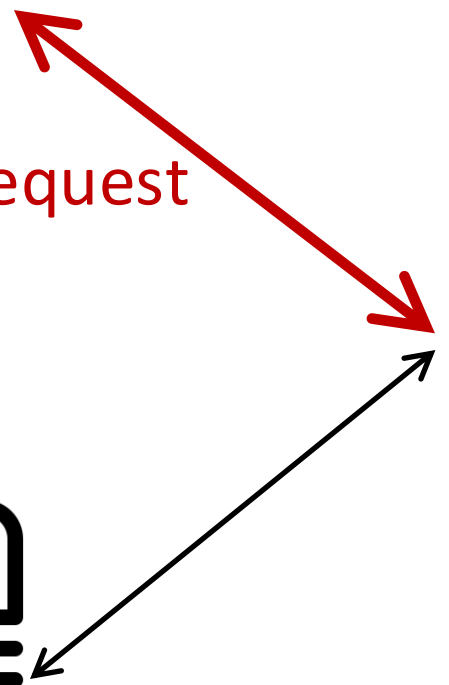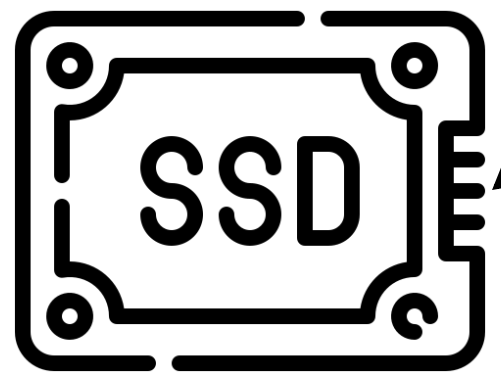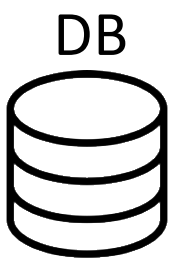# Bufferpool is Tightly Connected to Storage

# Bufferpool Manager

# Bufferpool Manager

Application

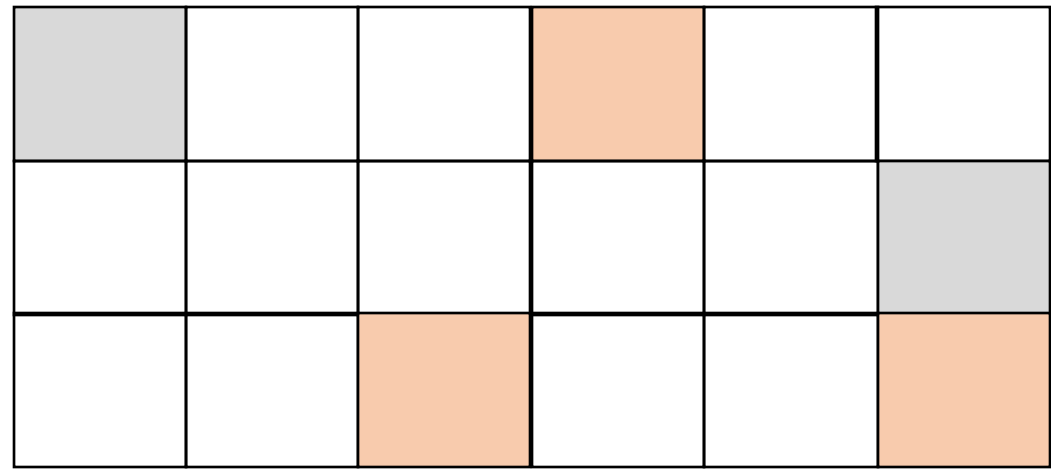**Bufferpool**

If page is not in BP,
fetch from disk

DB

SSD

Free frame

Disk page

Dirty page

Main Memory

# Bufferpool Manager

Application

**Bufferpool**



DB
SSD

Disk page

Dirty page

# Traditional Bufferpool Manager

Application

Page request

**Bufferpool**

Disk page

Dirty page

DB

SSD

If BP is full, one page is selected for eviction based on **page replacement policy**

# Traditional Bufferpool Manager

Application

**Bufferpool**



Disk page

Dirty page

DB

SSD

If the page is dirty, it is written back to disk

# Traditional Bufferpool Manager

Application

Page request comes

**Bufferpool**

Disk page

Dirty page

DB

SSD

Requested page is fetched in its place
**(exchanging one write for a read)**

# Popular Page Replacement Algorithms

LRU        (Most Popular)

LFU, FIFO     (Simple)

Clock Sweep   (Commercial)

CFLRU

LRU-WSR      Flash-Friendly

# CFLRU



Working region     Clean-first region

# CFLRU

Candidate for eviction

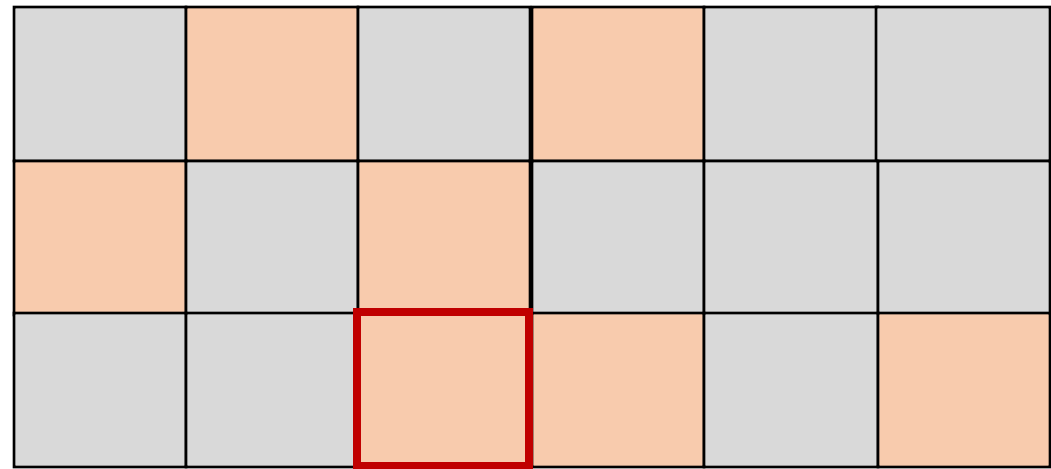| p6 | p5 | p4 | p3 | p2 | p1 |
|----|----|----|----|----|----|
| D | C | D | D | C | D |

Working region | Clean-first region

After Eviction:

Next Candidate for eviction

| p7 | p6 | p5 | p4 | p3 | p1 |
|----|----|----|----|----|----|
| C | D | C | D | D | D |

Working region | Clean-first region

# LRU-WSR

Cold flag

| p6 | p5 | p4 | p3 | p2 | p1 |
|----|----|----|----|----|----|
| D  | C  | D  | D  | C  | D  |
| 1  |    | 0  | 0  |    | 0  |

Cold flag NOT set!
This is be moved to front setting the cold flag

Cold flag

| p1 | p6 | p5 | p4 | p3 | p2 |
|----|----|----|----|----|----|
| D  | D  | C  | D  | D  | C  |
| 1  | 1  |    | 0  | 0  |    |

Candidate for eviction

After Eviction:

Cold flag

| p7 | p1 | p6 | p5 | p4 | p3 |
|----|----|----|----|----|----|
| C  | D  | D  | C  | D  | D  |
|    | 1  | 1  |    | 0  | 0  |

# LRU-WSR

p6  p5  p4  p3  p2  p1

| D | C | D | D | C | D |
|---|---|---|---|---|---|

Cold flag

| 1 | | 0 | 0 | | 1 |

Cold flag set!
**Candidate** for eviction

**After Eviction:**

p7  p6  p5  p4  p3  p2

| C | D | C | D | D | C |
|---|---|---|---|---|---|

Cold flag

| | 1 | | 0 | 0 | |

# The Challenges

- With write asymmetry, exchanging

  one write for one read is **NOT ideal**.

- Without exploiting concurrency,

  device remains vastly **underutilized**.

# Bufferpool Manager

## Eviction Policy

Which page to evict/write?
- LRU       – FIFO
- NRU       – 2Q
- Clock     – ARC
- Second Chance

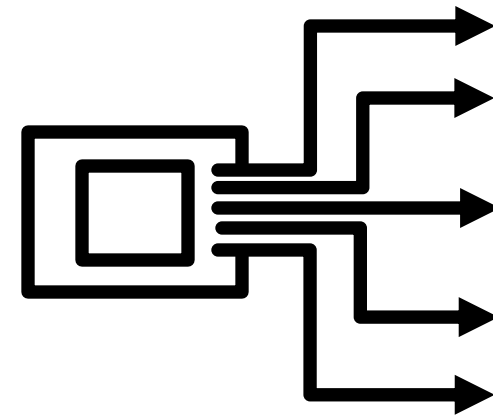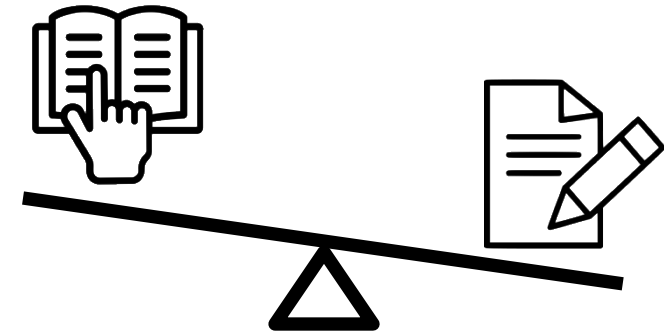- CFLRU     – CFLRU/C
- LRU-WSR   – CFLRU/E
- CCF-LRU   – DL-CFLRU/E

*Flash-friendly* policies

replacement policy

## Read-ahead Policy

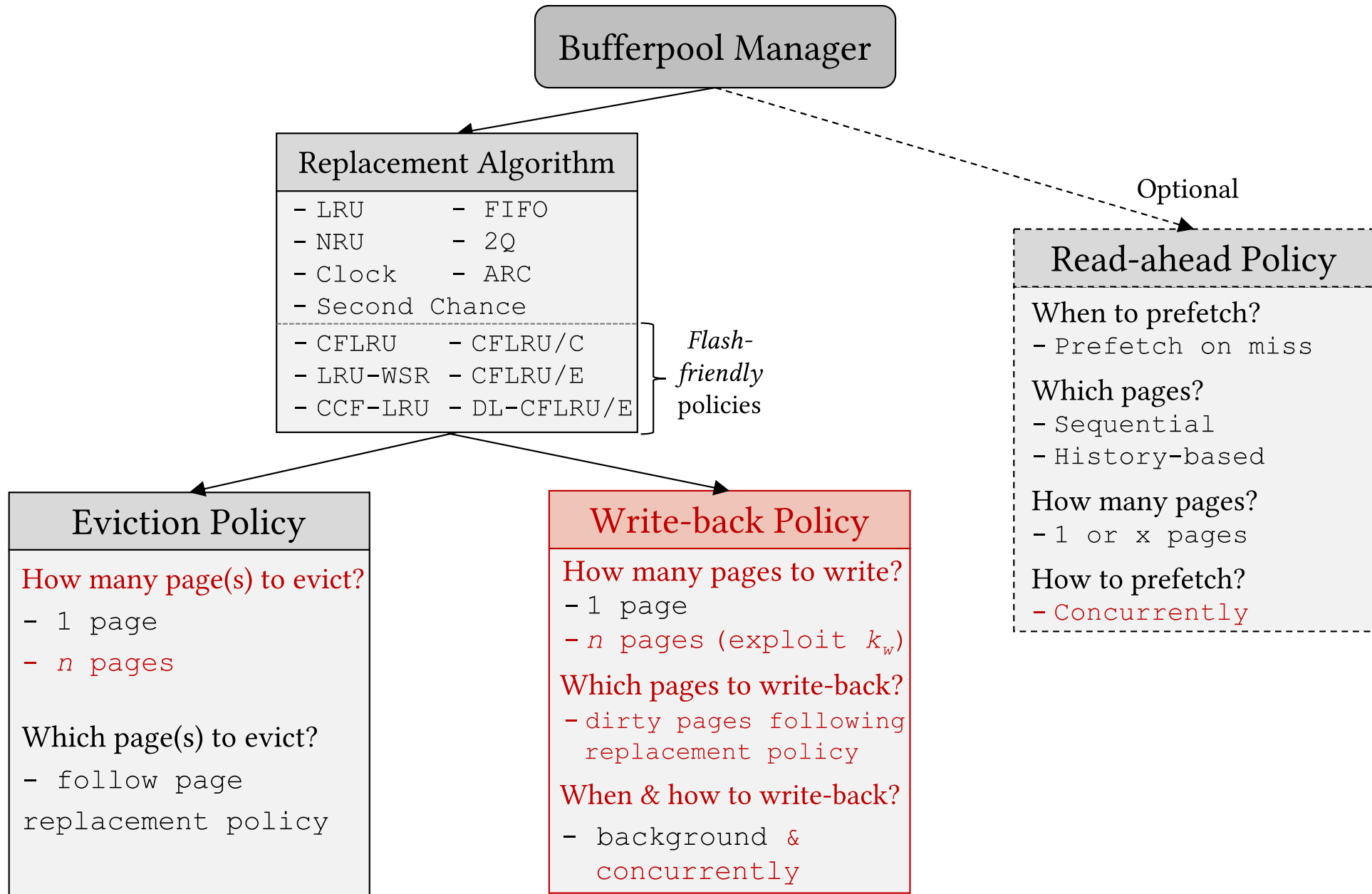Optional

When to prefetch?
- Prefetch on miss

Which pages?
- Sequential
- History-based

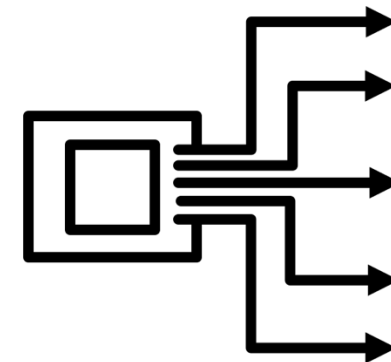How many pages?
- 1 or x pages

How to prefetch?
- Concurrently

# Asymmetry/Concurrency-Aware (`ACE`) Bufferpool Manager

# ACE Bufferpool Manager

Use device's properties

# $\mathrm{ACE}$ Bufferpool Manager

## ACE Bufferpool



**evict multiple pages**

**prefetch multiple pages**

**write back $k_w$ dirty pages**

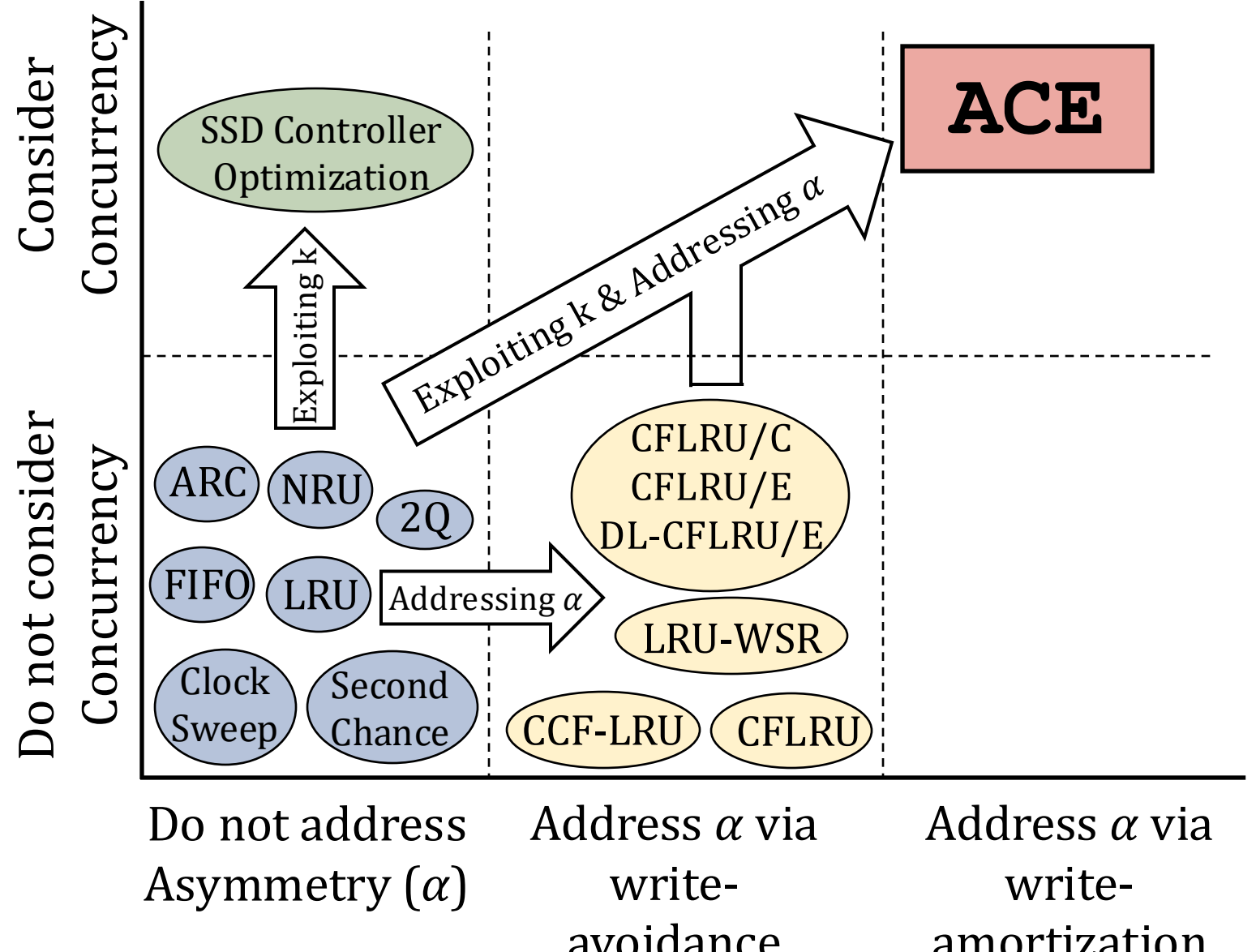DB

SSD

$\alpha, k_w, k_r$

✓ Can be integrated with **any** replacement algorithm

✓ **Any** prefetching technique can be used

# $\mathrm{ACE}$ Bufferpool Manager

Buffer Pool

| | | | |
|---|---|---|---|
| $p_1$ | $p_2$ | $p_4$ | $p_9$ |
| $p_5$ | $p_{12}$ | $p_{18}$ | $p_{10}$ |
| $p_{13}$ | $p_7$ | $p_{24}$ | $p_{21}$ |

$p$ Dirty page

$p$ Clean page

Evict $n_e$ pages

**Writer**

Concurrent *Device-aware* Writing

**Evictor**

... $p_2$ $p_1$ $p_5$ $p_9$

Candidates

**Reader**

Seq. Stream Prefetcher

History-based Prefetcher

Concurrently write back $n_w$ dirty pages

SSD

Parallelly prefetch $n_e$ - 1 pages

mru             lru

B | 6 | 2 | 3 | 5 | 7 | 4 | 9 |
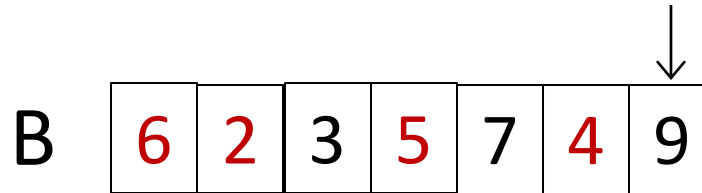
Let's assume: $k_w = 3$, LRU is the baseline replacement policy & red indicates dirty page

**Write request of page 8 comes**
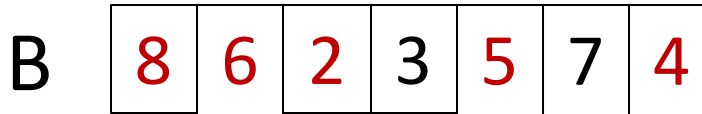
# An Example ($k_w = 3$)

**write page 8**

Candidate for eviction

↓

B | 6 | 2 | 3 | 5 | 7 | 4 | 9 |

Since candidate page is
clean, we simply evict 9
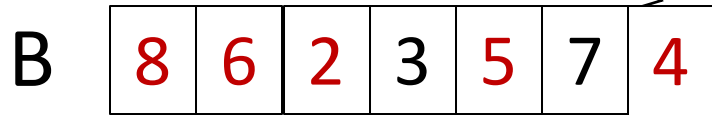
After eviction:

B | 8 | 6 | 2 | 3 | 5 | 7 | 4 |

**Write request of page 1 comes**

# An Example ($k_w = 3$)

**write page 1**

## LRU

Candidate

B | 8 | 6 | 2 | 3 | 5 | 7 | 4

After eviction:

B | 1 |

# An Example ($k_w = 3$)

**write page 1**

## LRU

## LRU+ACE(w/o PF)

Candidate

B | 8 | 6 | 2 | 3 | 5 | 7 | 4 |

| 8 | 6 | 2 | 3 | 5 | 7 | 4 |

After eviction:

B | 1 | 8 | 6 | 2 | 3 | 5 | 7 |

# An Example ($k_w = 3$)

**write page 1**

## LRU

**LRU+ACE(w/o PF)**

Candidate

B | 8 | 6 | 2 | 3 | 5 | 7 | 4 |

| 8 | 6 | 2 | 3 | 5 | 7 | 4 |

After eviction:

B | 1 | 8 | 6 | 2 | 3 | 5 | 7 |

4,5,2 concurrently written

4 evicted

**write page 1**

## LRU

## LRU+ACE(w/o PF)

B | 8 | 6 | 2 | 3 | 5 | 7 | 4 |

| 8 | 6 | 2 | 3 | 5 | 7 | |

After eviction:

After eviction:

B | 1 | 8 | 6 | 2 | 3 | 5 | 7 |

| 1 | 8 | 6 | 2 | 3 | 5 | 7 |

# An Example ($k_w = 3, n_e = 2$)

**write page 1**

## LRU

## LRU+ACE(w/o PF)

## LRU+ACE(w/PF)

Candidate

B | 8 | 6 | 2 | 3 | 5 | 7 | 4 |

| 8 | 6 | 2 | 3 | 5 | 7 | 4 |

| 8 | 6 | 2 | 3 | 5 | 7 | 4 |

After eviction:

After eviction:

B | 1 | 8 | 6 | 2 | 3 | 5 | 7 |

| 1 | 8 | 6 | 2 | 3 | 5 | 7 |

# An Example ($k_w = 3$, $n_e = 2$)

**write page 1**

## LRU

B

| 8 | 6 | 2 | 3 | 5 | 7 | 4 |

After eviction:

B

| 1 | 8 | 6 | 2 | 3 | 5 | 7 |

## LRU+ACE(w/o PF)

| 8 | 6 | 2 | 3 | 5 | 7 | 4 |

After eviction:

| 1 | 8 | 6 | 2 | 3 | 5 | 7 |

## LRU+ACE(w/PF)

eviction window

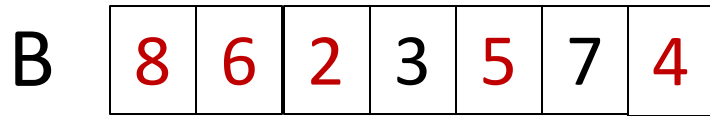| 8 | 6 | 2 | 3 | 5 | 7 | 4 |

4,5,2 concurrently written

4,7 evicted

# An Example ($k_w = 3$, $n_e = 2$)

**write page 1**

## LRU

B | 8 | 6 | 2 | 3 | 5 | 7 | 4 |

After eviction:

B | 1 | 8 | 6 | 2 | 3 | 5 | 7 |

## LRU+ACE(w/o PF)

| 8 | 6 | 2 | 3 | 5 | 7 | 4 |

After eviction:

| 1 | 8 | 6 | 2 | 3 | 5 | 7 |

## LRU+ACE(w/PF)

| 8 | 6 | 2 | 3 | 5 | | |

After eviction:

| 1 | 8 | 6 | 2 | 3 | 5 | 9 |

prefetched

# Experimental Evaluation



11.5

Clock Sweep
LRU
CFLRU
LRU-WSR

vs    their $\mathbb{ACE}$ counterparts

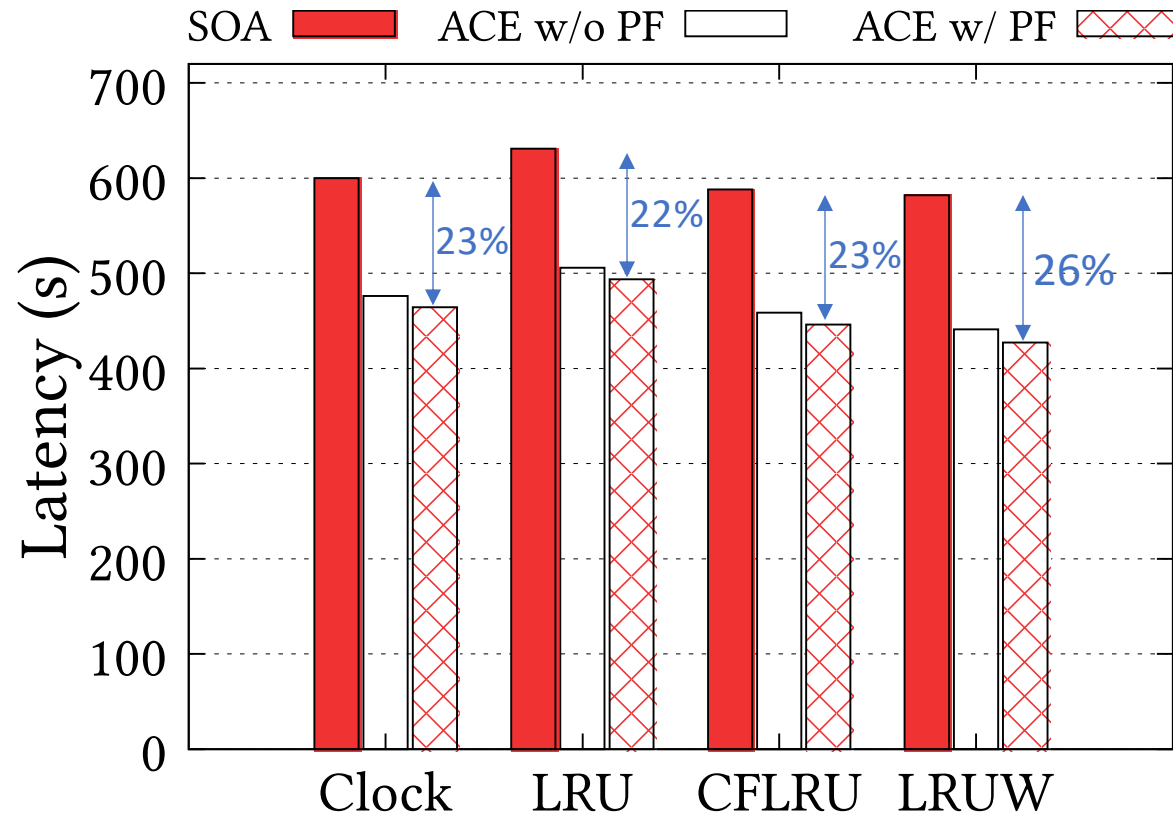| Device | $\alpha$ | $k_r$ | $k_w$ |
|---|---|---|---|
| Optane SSD | 1.1 | 6 | 5 |
| PCIe SSD | 2.8 | 80 | 8 |
| SATA SSD | 1.5 | 25 | 9 |
| Virtual SSD | 2.0 | 11 | 19 |

**Workload:**

synthesized traces

TPC-C benchmark

# ACE Improves Runtime

**Device: PCIe SSD**

$\alpha = 2.8$, $k_w = 8$



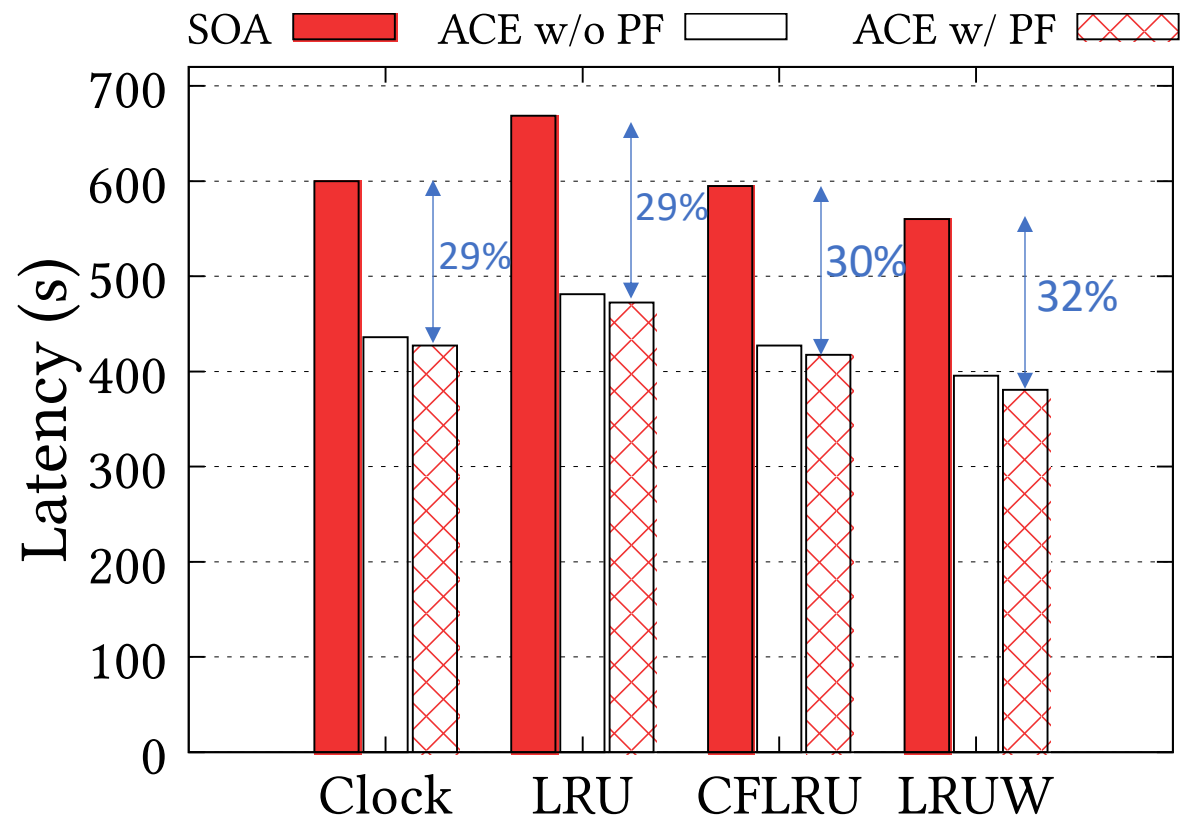ACE improves runtime by 22-26%

Negligible increase in buffer miss (<0.009%)

**Benefit comes at no cost**

# Higher Gain for Write-Heavy Workload

**Device: PCIe SSD**

$\alpha = 2.8$, $k_w = 8$



Write-intensive workloads have

higher benefit (up to 32%)

# Impact of R/W Ratio & Asymmetry


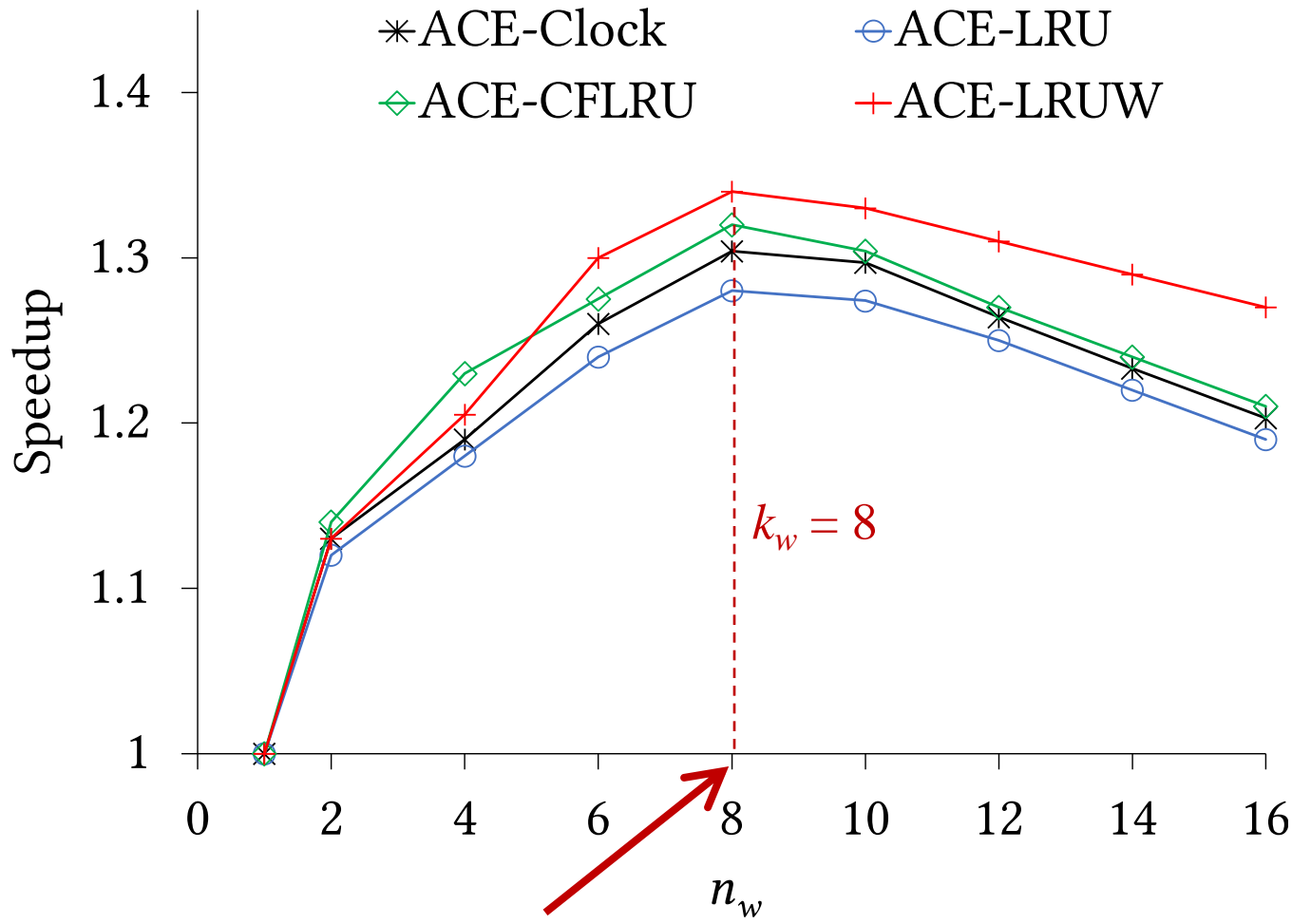
more writes, more speedup

higher asymmetry, higher speedup

good benefit even for low asymmetry
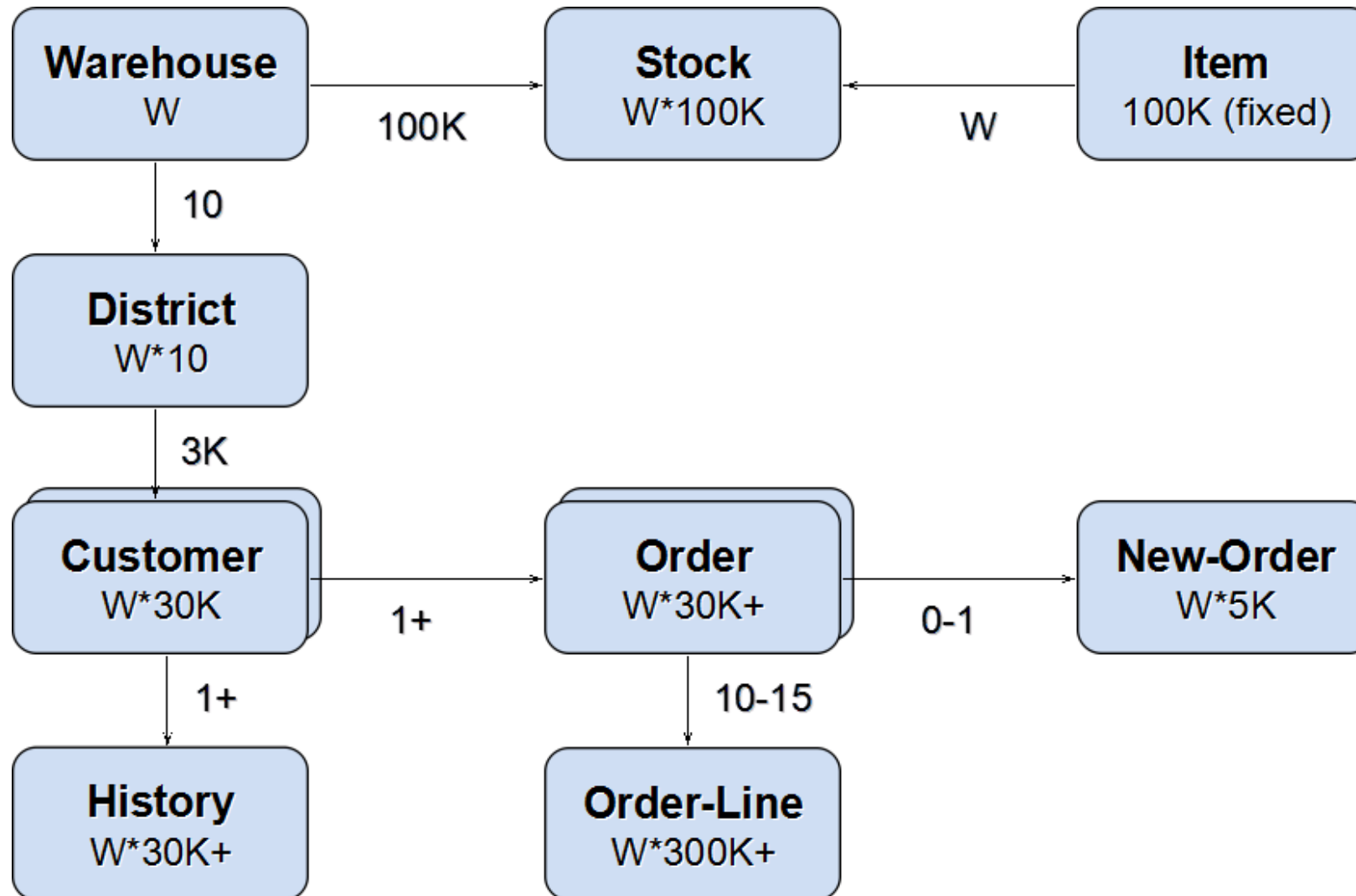
# Impact of #Concurrent I/Os



**Device: PCIe SSD**
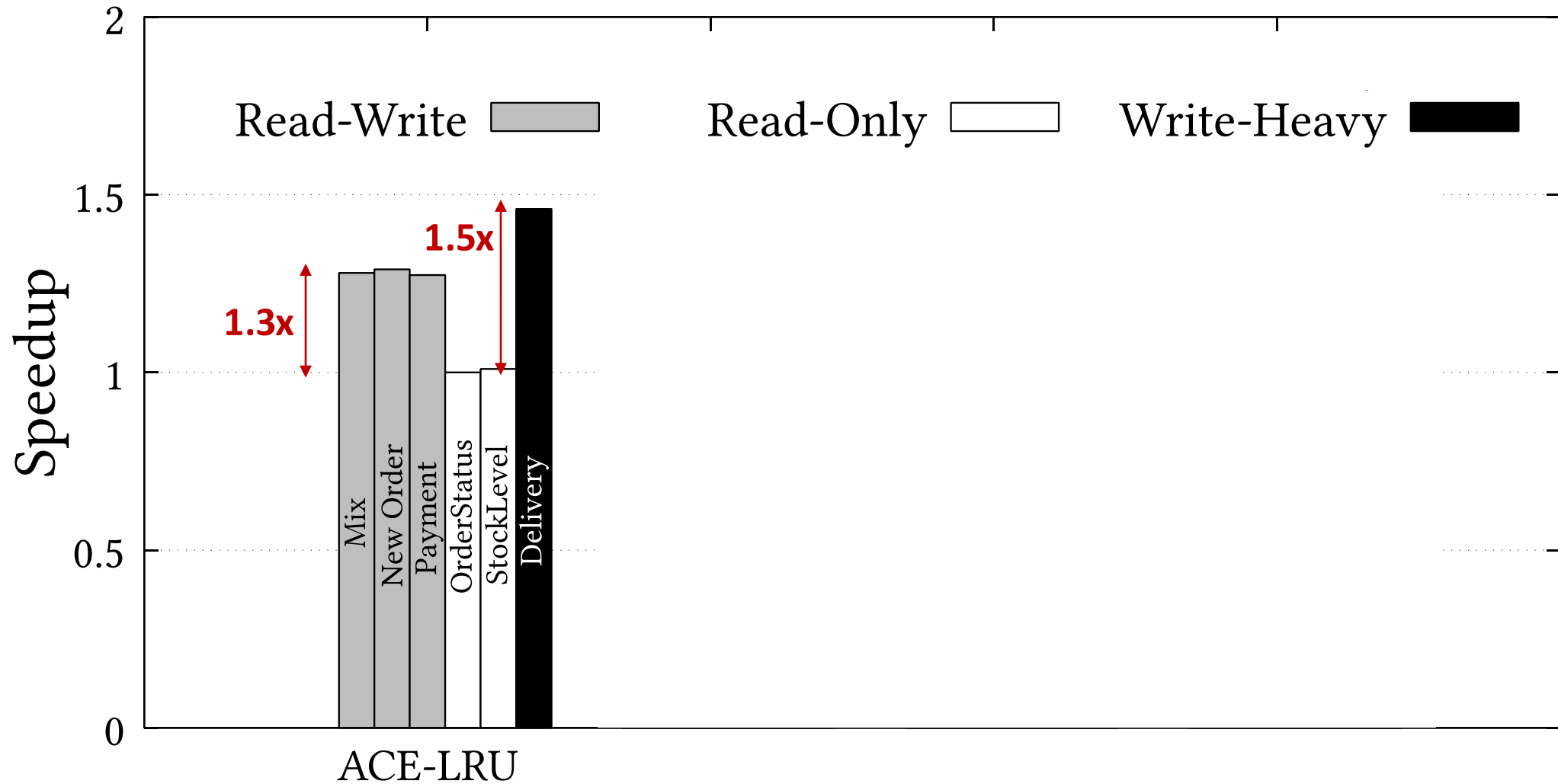
$\alpha = 2.8$, $k_w = 8$

**Highest speedup when optimal concurrency is used**

Legend:
- ✳ ACE-Clock
- ○ ACE-LRU
- ◇ ACE-CFLRU
- + ACE-LRUW

$k_w = 8$

$n_w$

Mixed Skewed Trace
(r/w: 50/50, locality
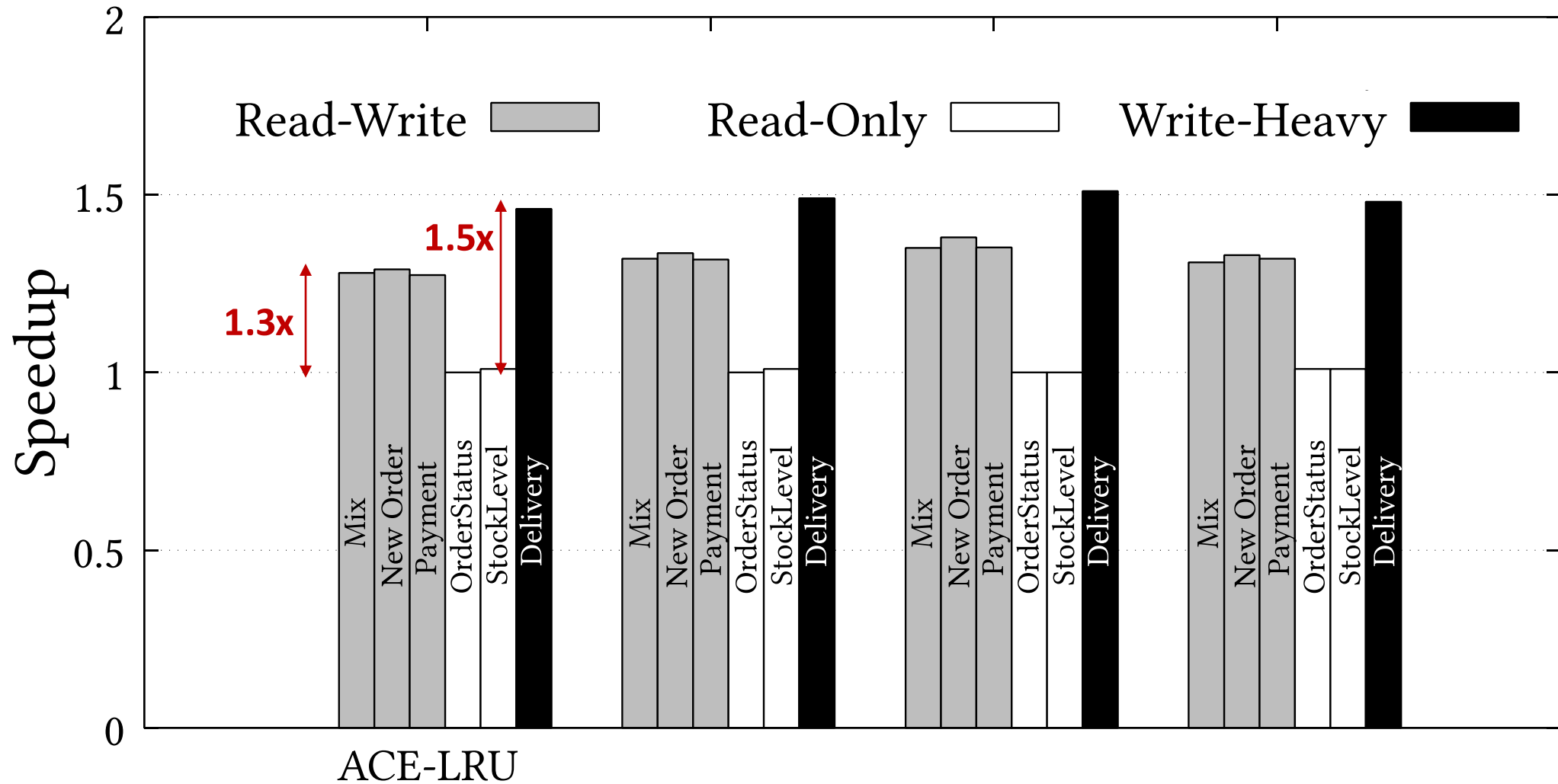
# Experimental Evaluation (TPC-C)

# Experimental Evaluation (TPC-C)



**ACE Achieves 1.3x for mixed TPC-C**

# Experimental Evaluation (TPC-C)



**ACE Achieves 1.3x for mixed TPC-C**

# Summary

ACE Bufferpool



DB

SSD

☑ Decoupled eviction and write-back mechanism

☑ `ACE` works with **any** page replacement policy

☑ **Any** prefetching technique can be used

☑ With low engineering effort, **any** DBMS

bufferpool can benefit from this approach

# Conclusion & Future Work

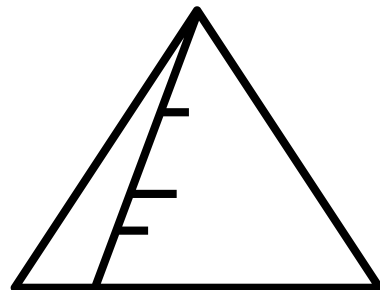Make *asymmetry and concurrency* part of *algorithm design*

*...* not simply an engineering optimization

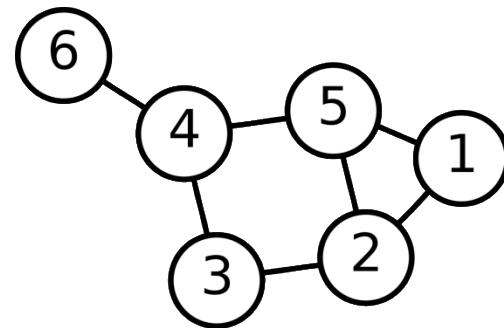Build algorithms/data structures for storage devices

with asymmetry $\alpha$ and concurrency $k$

Stay Tuned!

index structures

graph traversal

**SIGMOD '24**

**CAVE: Concurrency-Aware Graph Processing on SSDs**

disc.bu.edu/papers/sigmod24-cave

# Thank You!

[disc.bu.edu/papers/icde23-papon](disc.bu.edu/papers/icde23-papon)