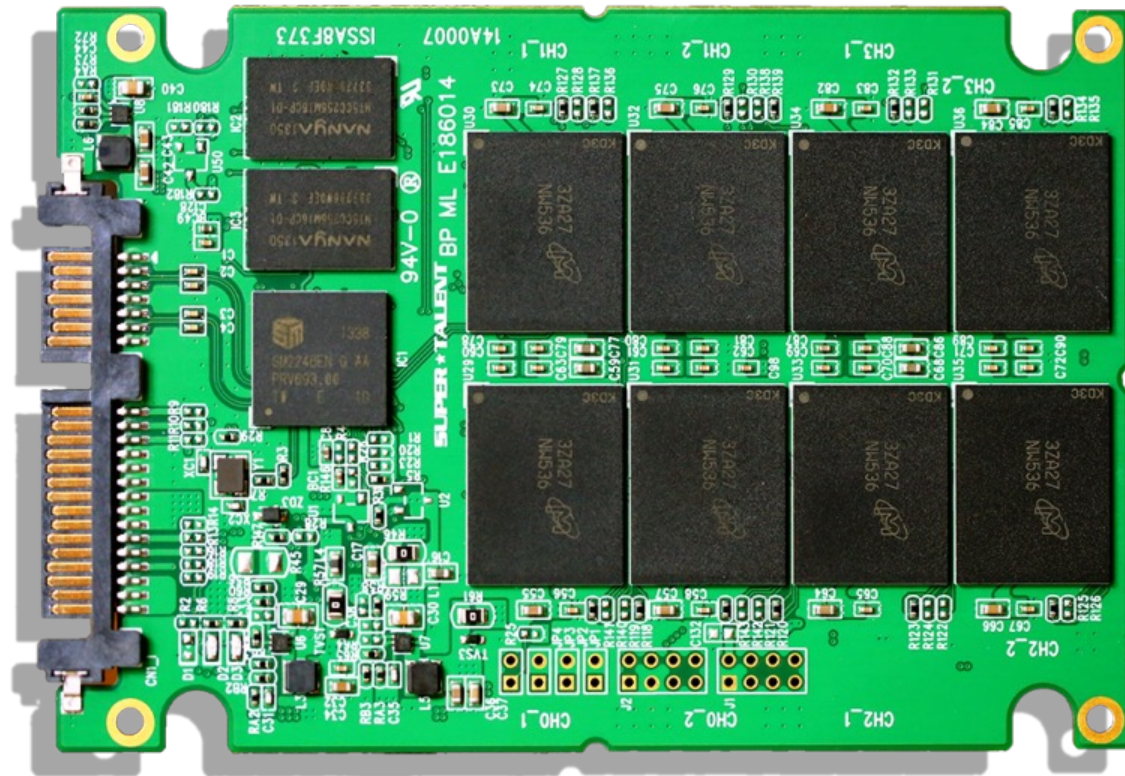


ACEing the Bufferpool Management Paradigm for Modern Storage Devices

Tarikul Islam Papon

Manos Athanassoulis

Solid State Drives



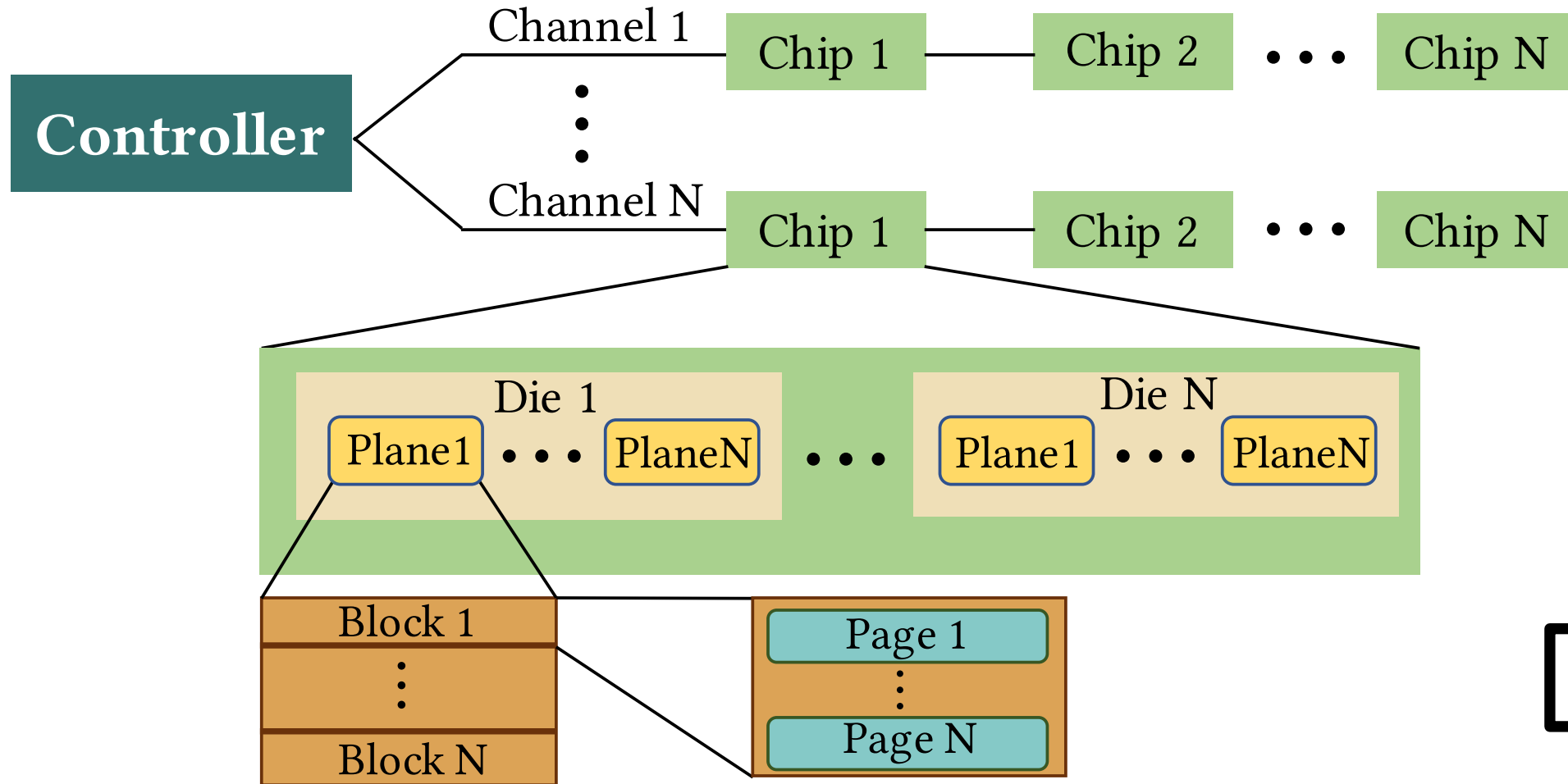
electronic device

fast random access

concurrent I/Os

write latency > read latency

Concurrency



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

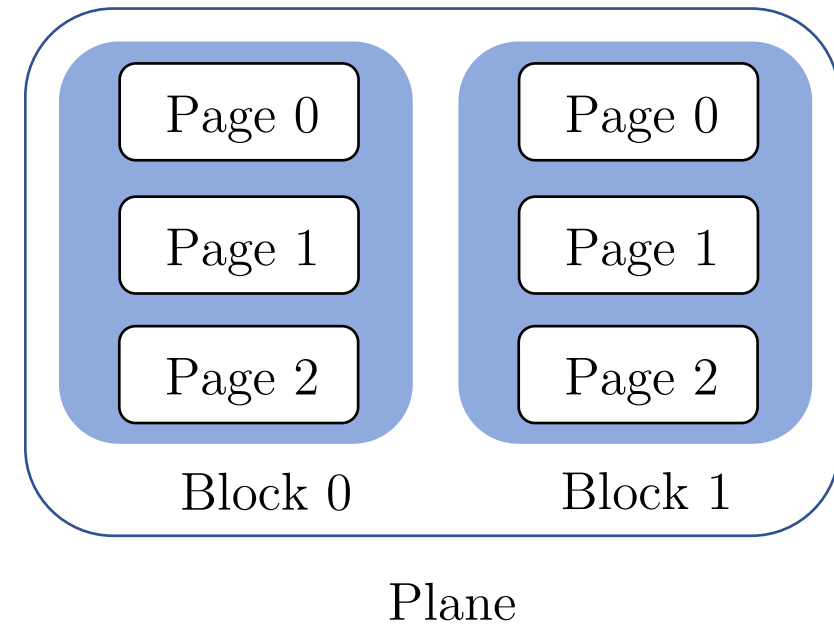
Read/Write Asymmetry

Out-of-place updates cause invalidation

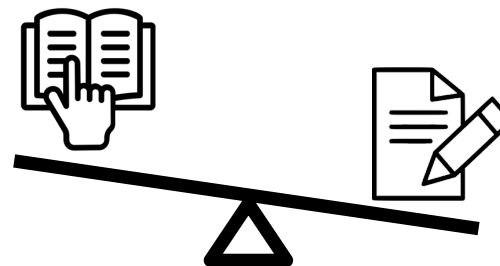
“Erase before write” approach

Garbage Collection

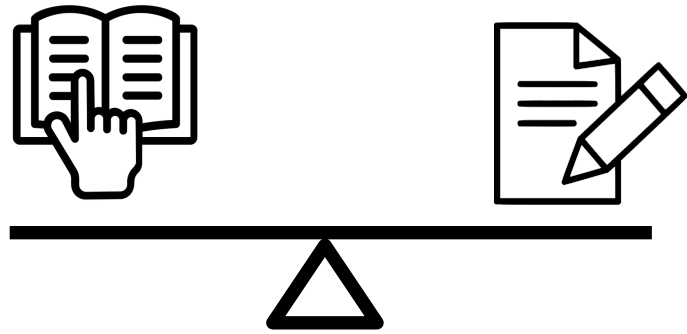
Larger erase granularity



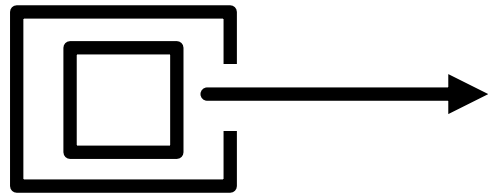
All these results in higher amortized write cost



HDD

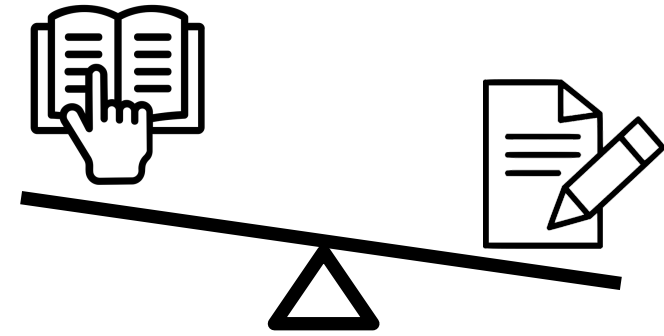


Symmetric cost for Read & Write

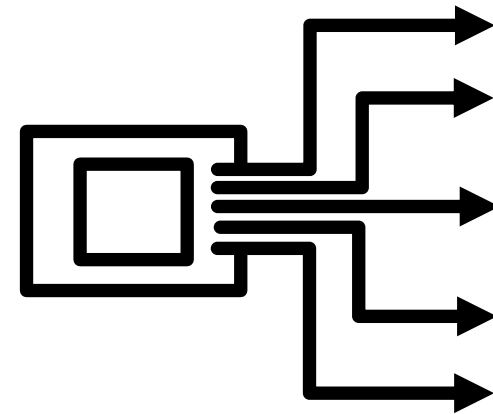


One I/O at a time

SSD



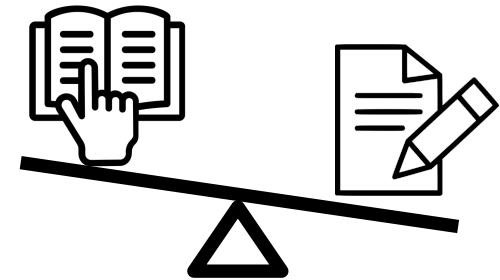
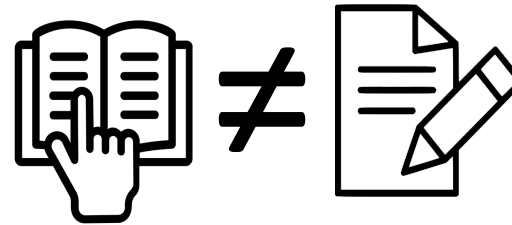
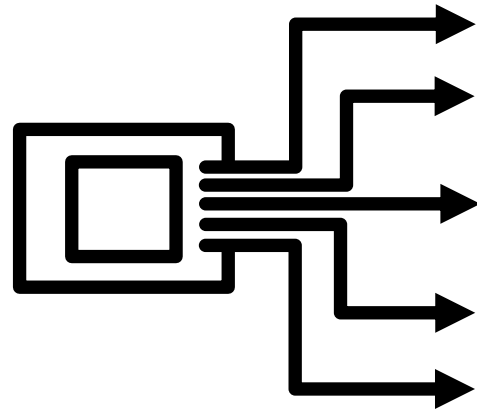
Read/Write Asymmetry (α)



Concurrency (k)

“A Parametric I/O Model for Modern Storage Devices”

DaMoN 2021



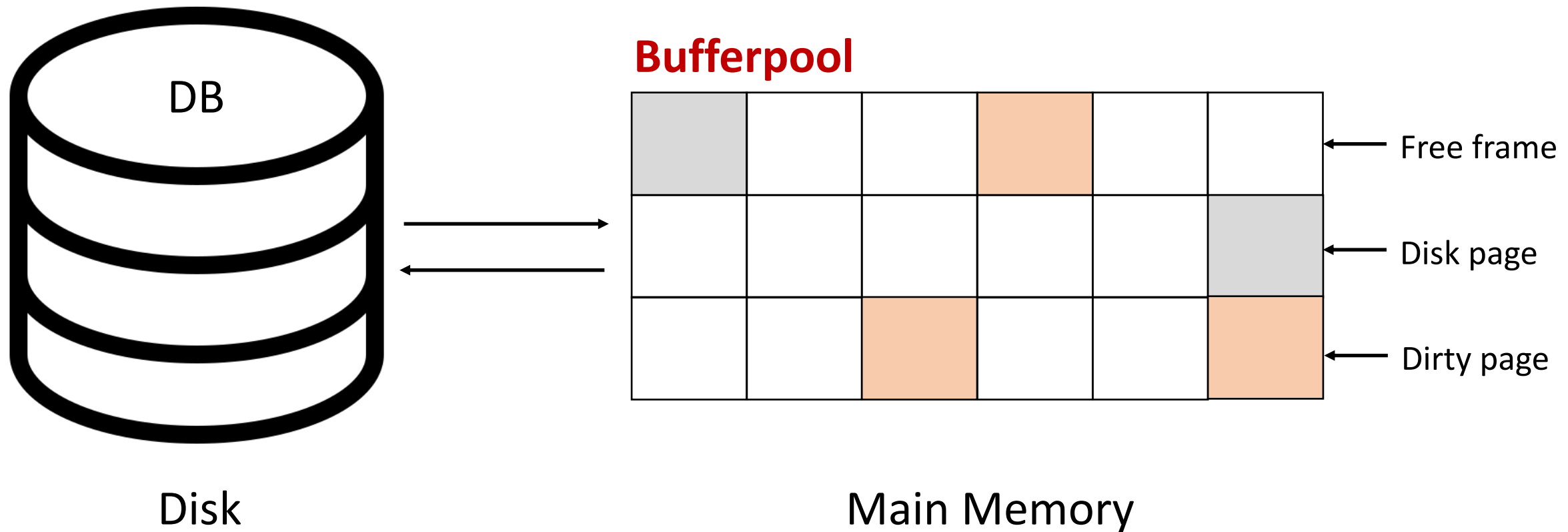
know Thy Device

exploit k_r and k_w
 (with care)
 read concurrency ←
 write concurrency

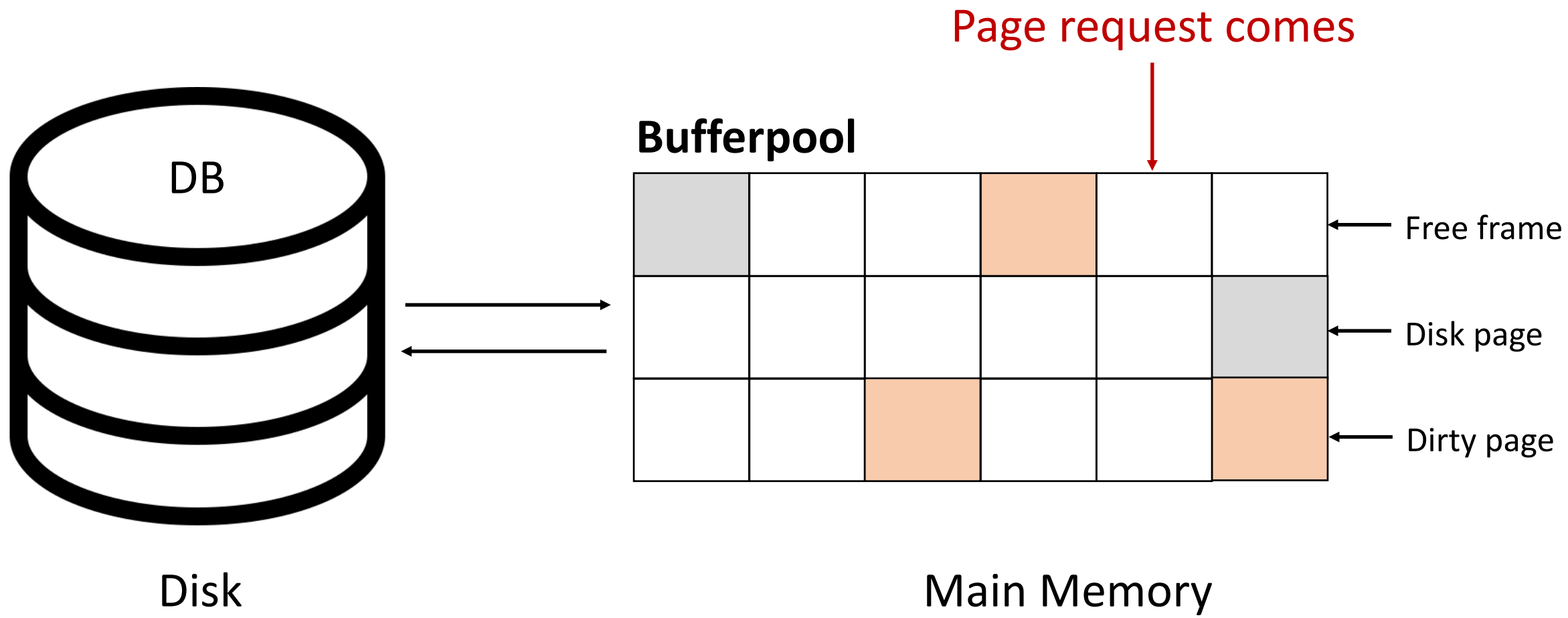
treat read and write differently

asymmetry (α) controls performance

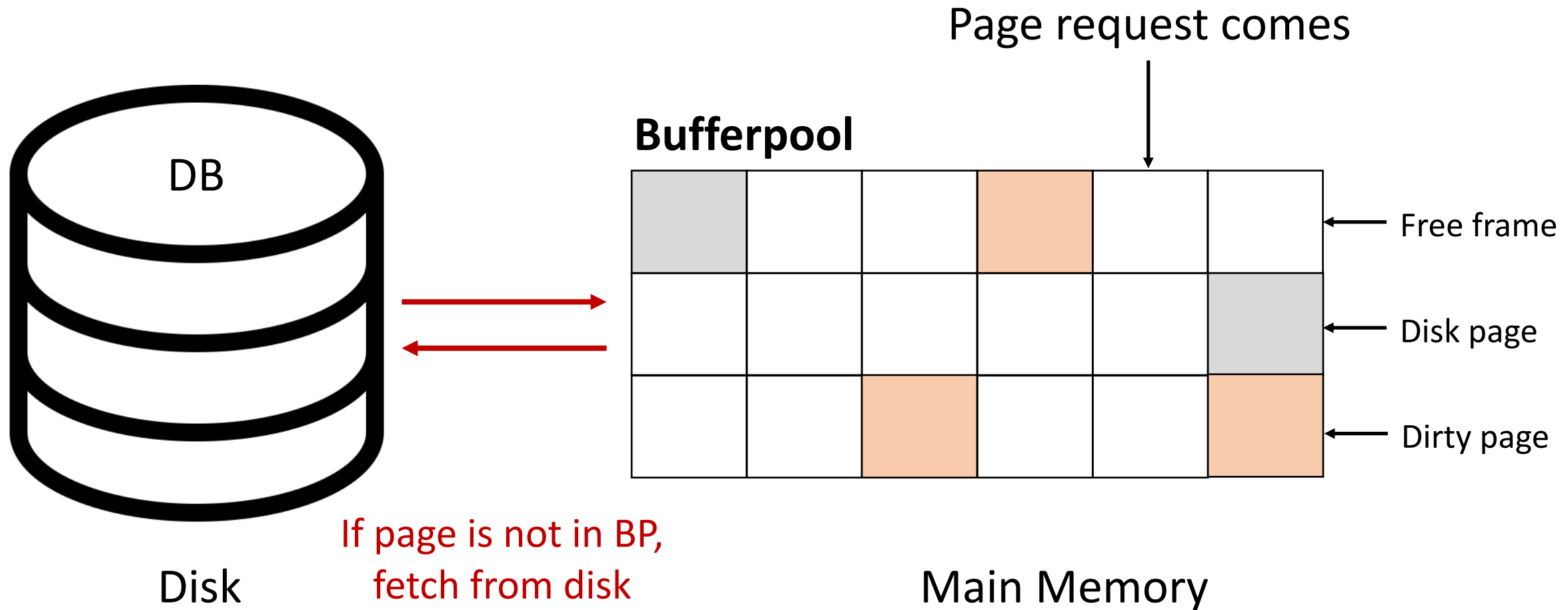
Bufferpool is Tightly Connected to Storage



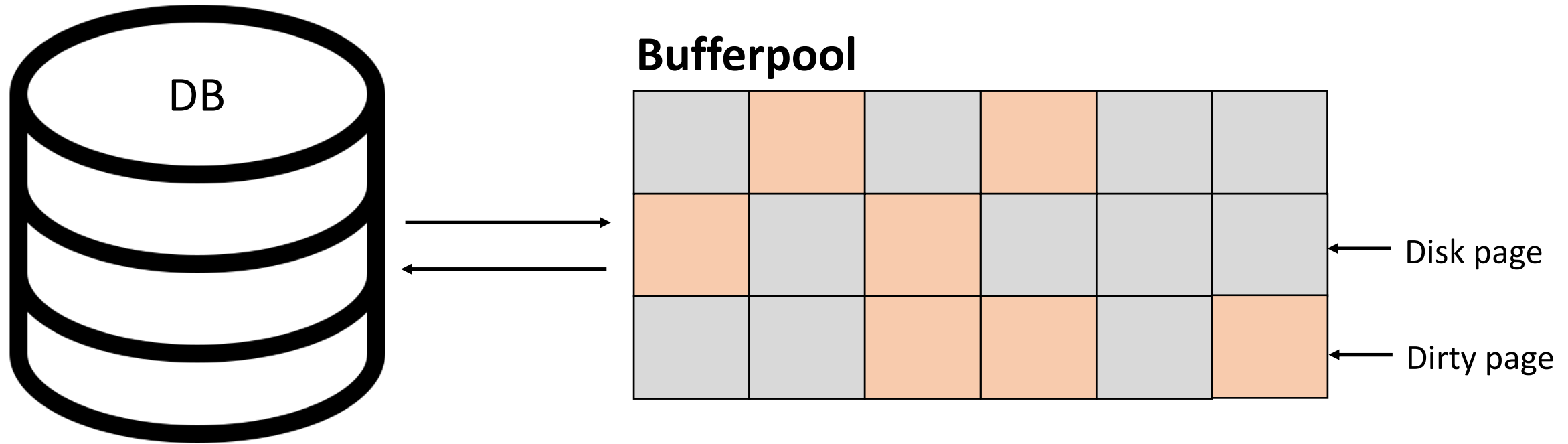
Traditional Bufferpool Manager



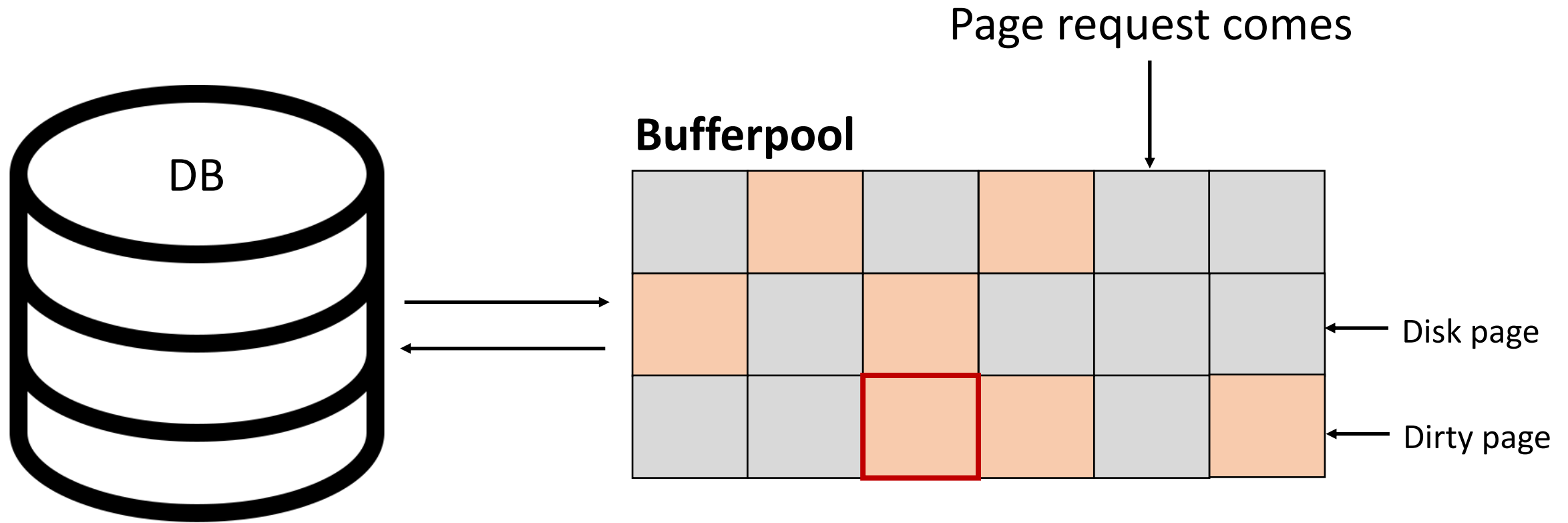
Traditional Bufferpool Manager



Traditional Bufferpool Manager

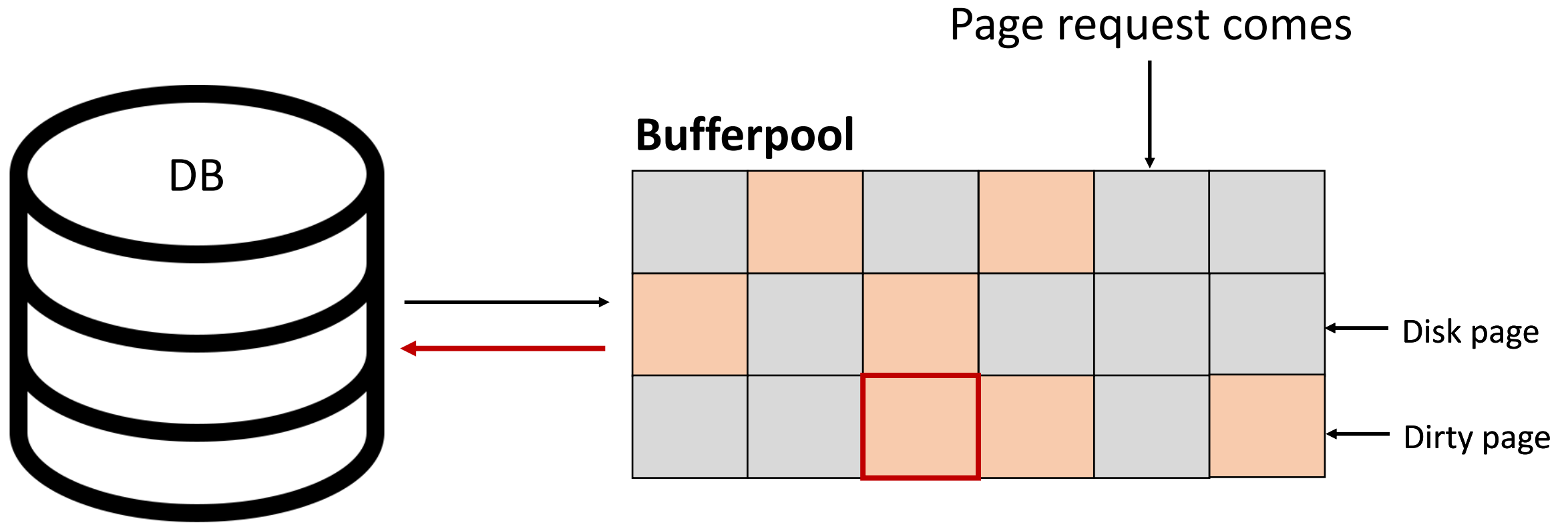


Traditional Bufferpool Manager



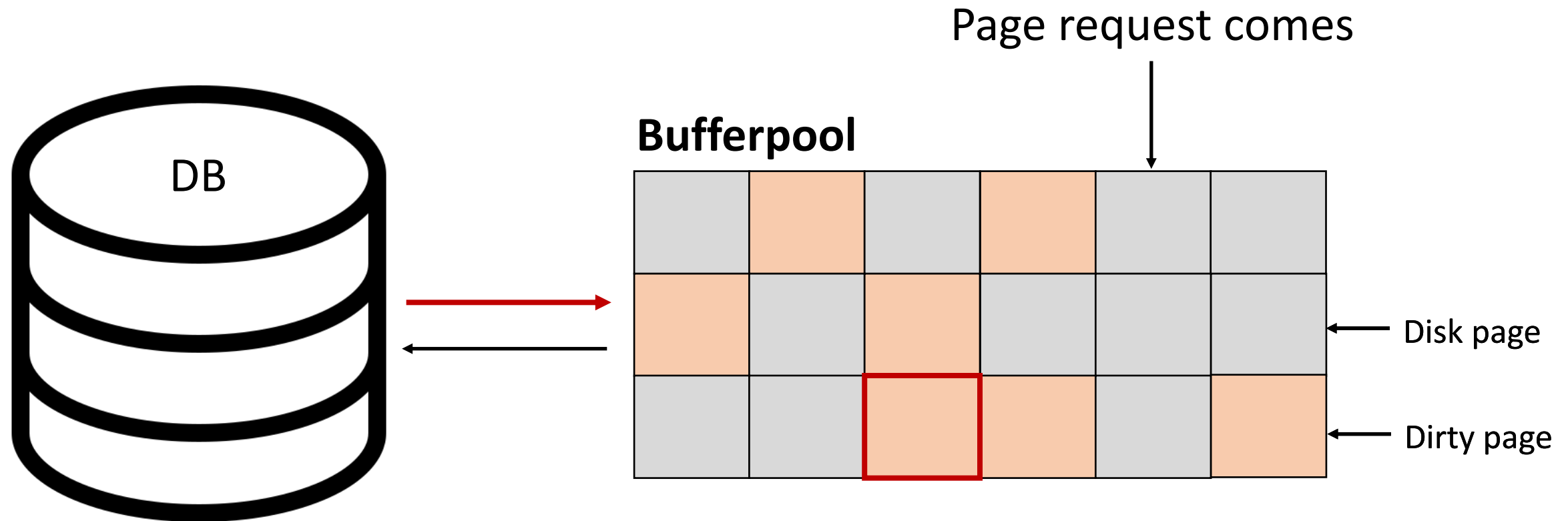
If BP is full, one page is selected for eviction based on **page replacement policy**

Traditional Bufferpool Manager



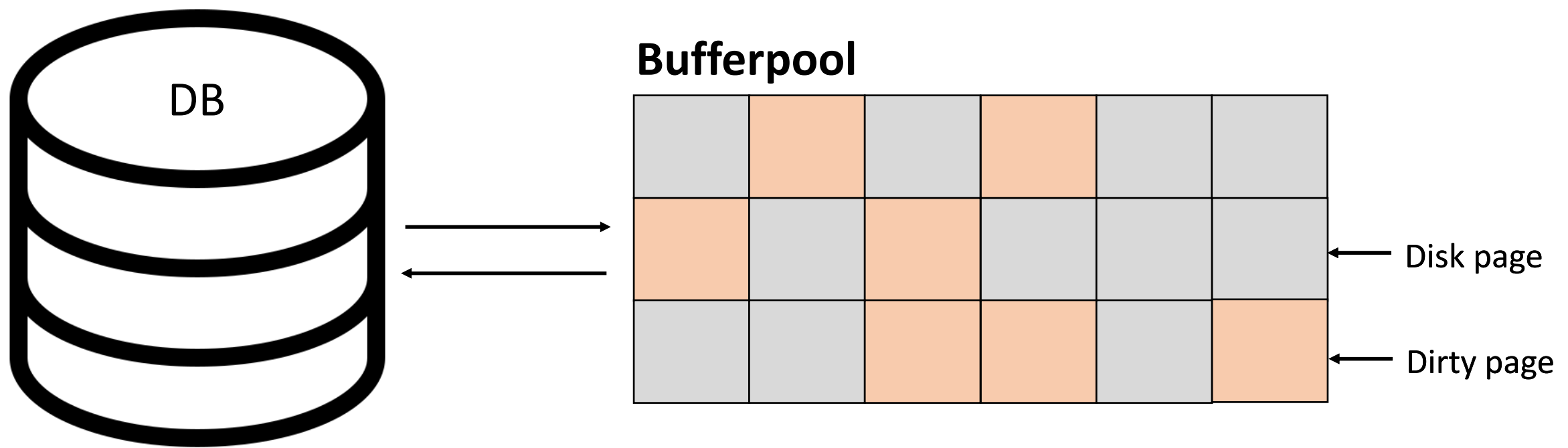
If the page is dirty, it is written back to disk

Traditional Bufferpool Manager



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

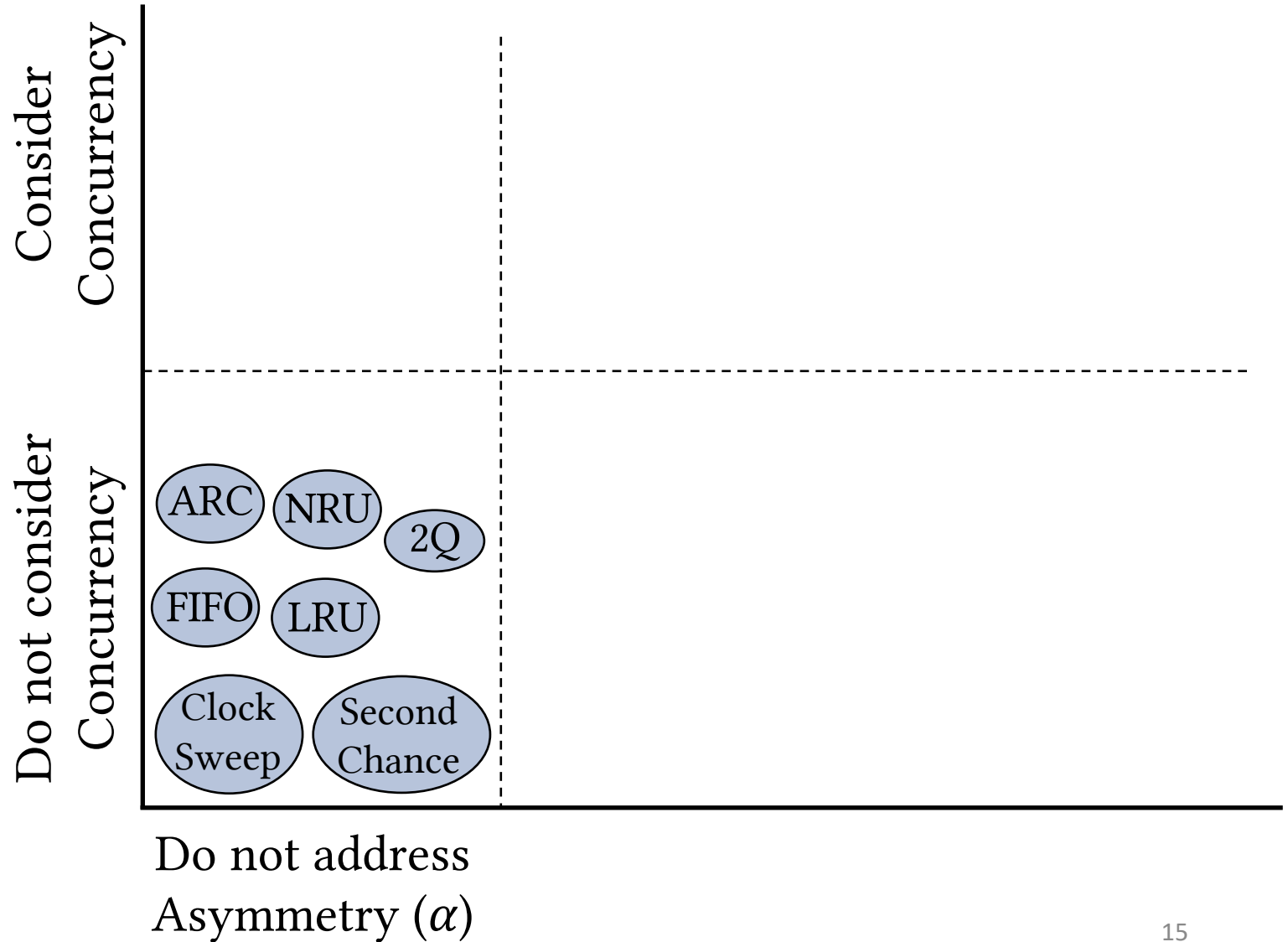
Traditional Bufferpool Manager



Is this Optimal?

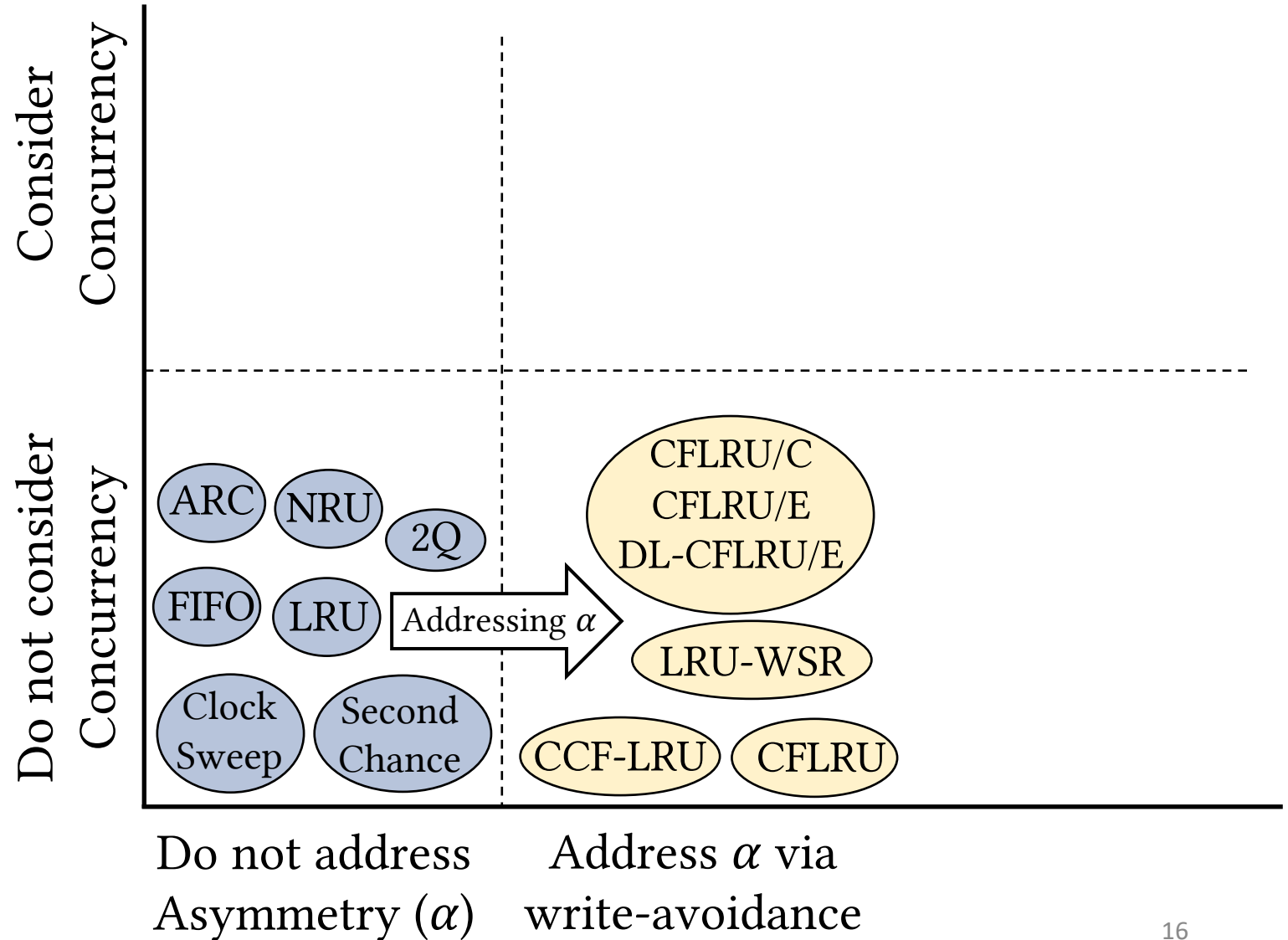
The Challenge

- With write asymmetry, it is **NOT** fair to exchange one write for one read.



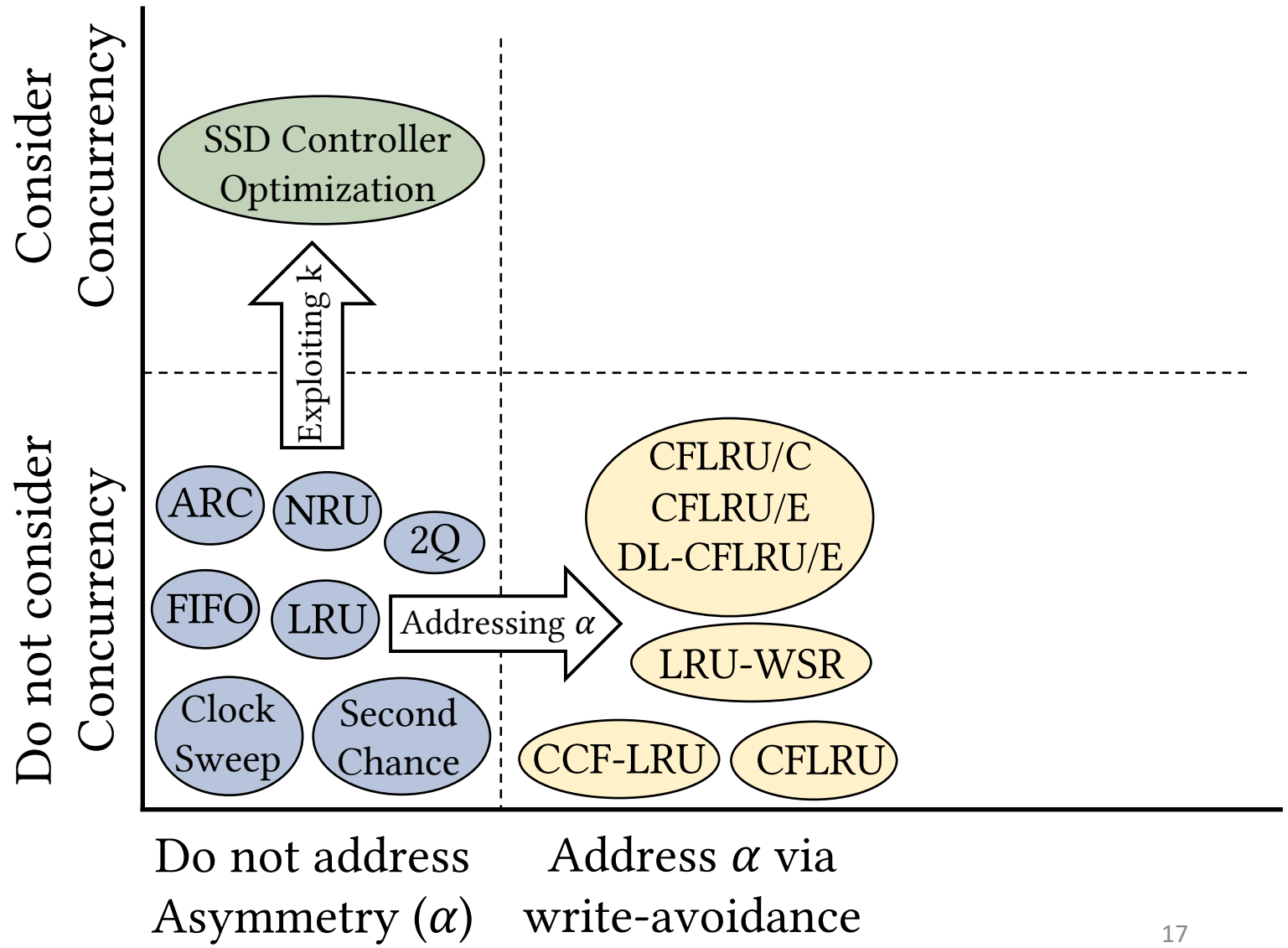
The Challenge

- With write asymmetry, it is **NOT** fair to exchange one write for one read.



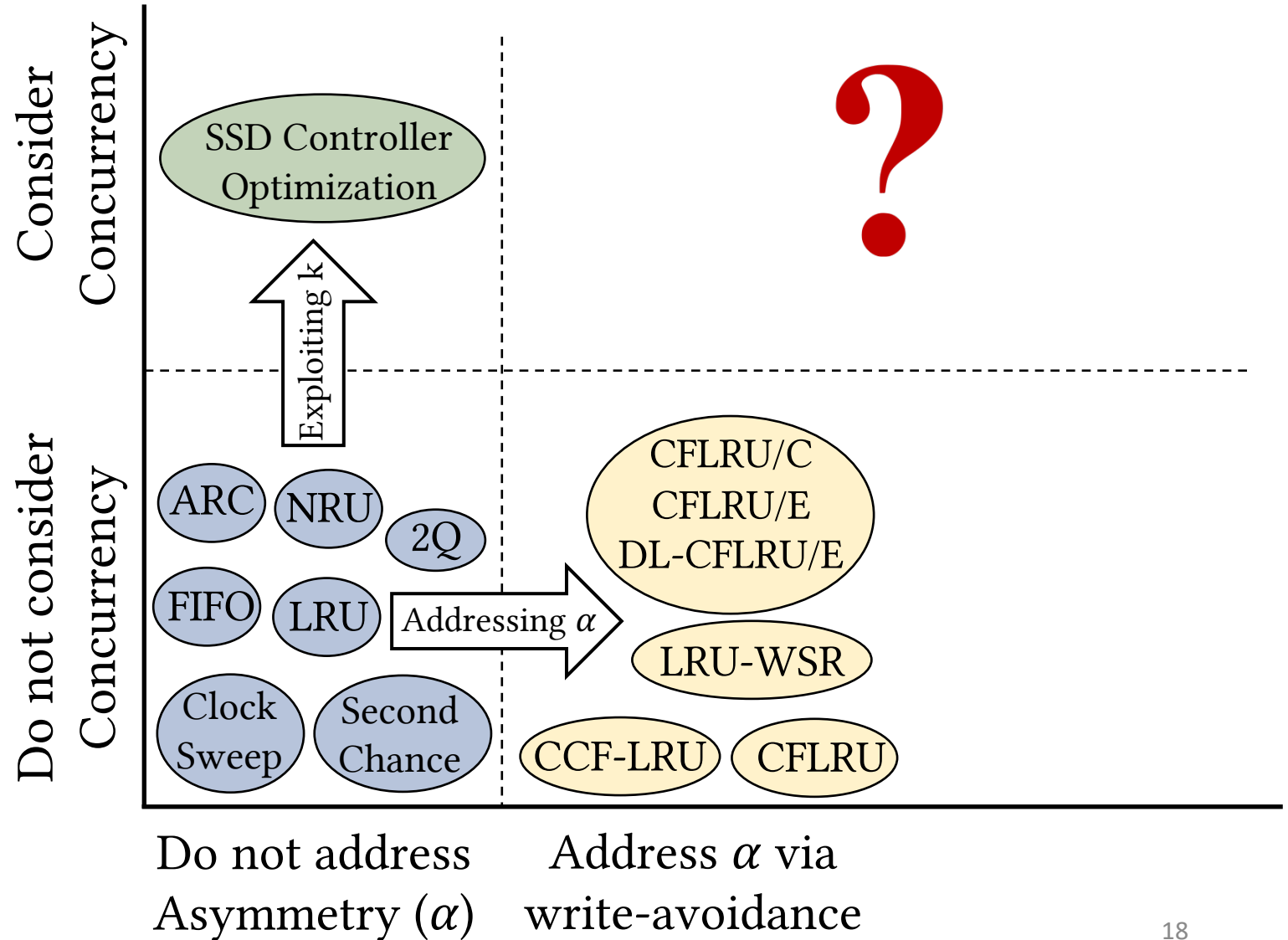
The Challenge

- With write asymmetry, it is **NOT** fair to exchange one write for one read.
- Do not expressly utilize the device concurrency.



The Challenge

- With write asymmetry, it is **NOT** fair to exchange one write for one read.
- Do not expressly utilize the device concurrency.



The Challenge

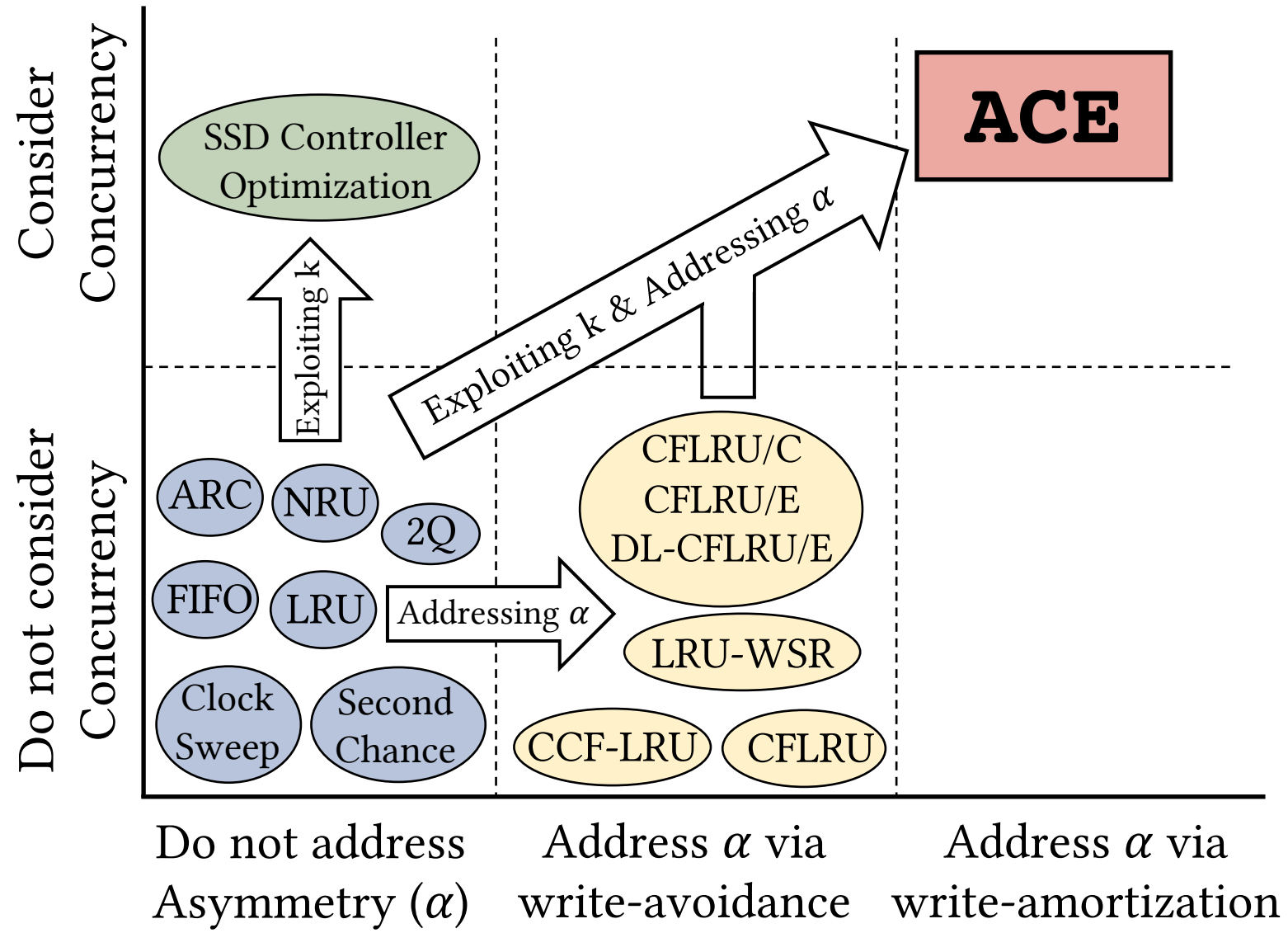
device under-utilization

poor end-to-end performance

high deployment cost

The Solution

- device under-utilization ✓
- poor end-to-end performance ✓
- high deployment cost ✓
- ease of integration** ✓



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*

Optional

Read-ahead Policy

When to prefetch?

- Prefetch on miss

Which pages?

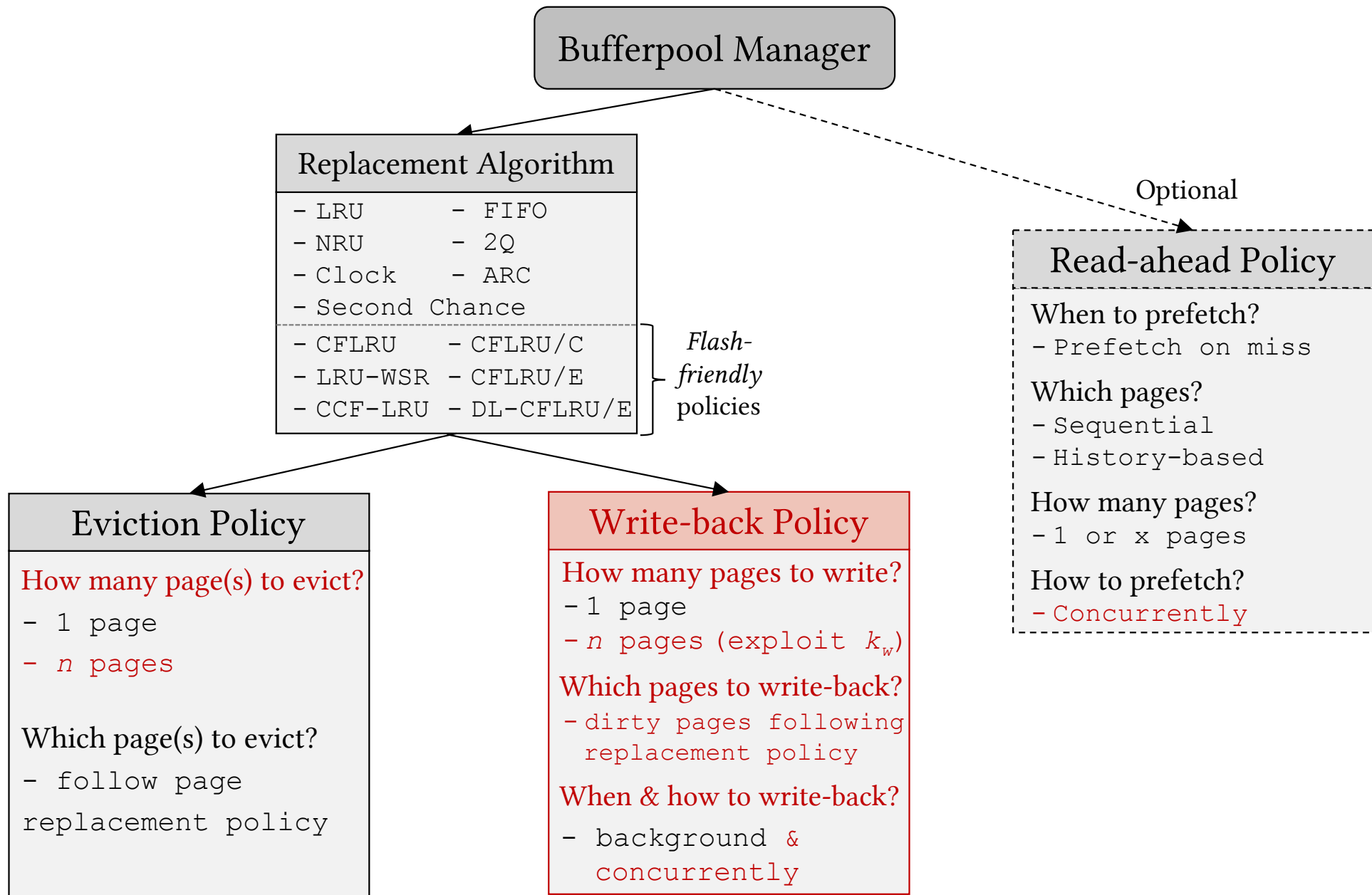
- Sequential
- History-based

How many pages?

- 1 or x pages

How to prefetch?

- **Concurrently**

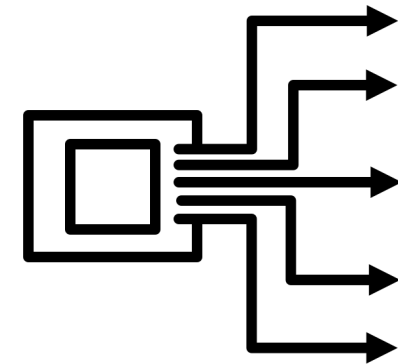
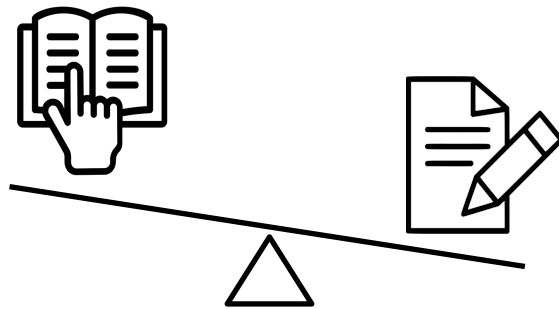


Asymmetry/Concurrency-Aware (ACE) Bufferpool Manager

ACE Bufferpool Manager



Use device's properties



ACE Bufferpool Manager



$1 \leq n_e \leq$ read concurrency (k_r)

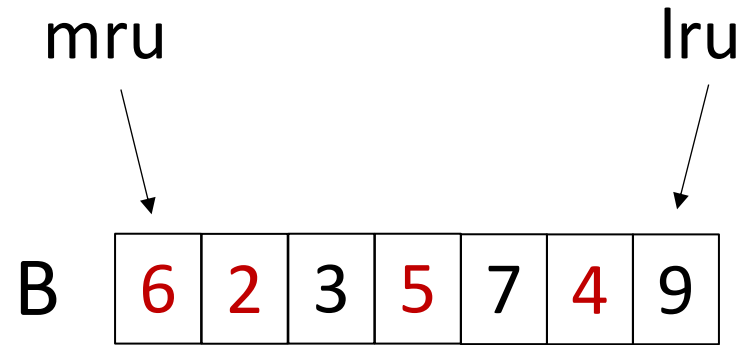
n_w = device's write concurrency (k_w)

write n_w **dirty pages** concurrently

evict n_e **pages**

prefetch $n_e - 1$ **pages** concurrently

An Example



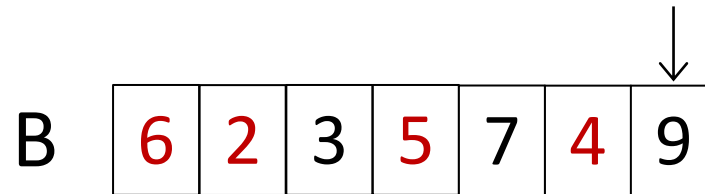
Let's assume: $k_w = 3$, LRU is the baseline replacement policy & **red** indicates dirty page

Write request of page 8 comes

An Example ($k_w = 3$)

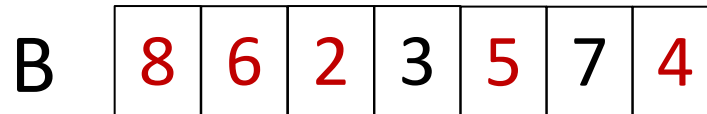
write page 8

Candidate for eviction



Since candidate page is clean, we simply evict 9

After eviction:



Write request of page 1 comes

An Example ($k_w = 3$)

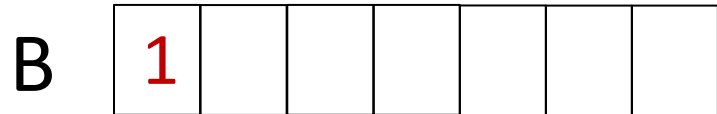
write page 1

LRU

Candidate



After eviction:



An Example ($k_w = 3$)

write page 1

LRU

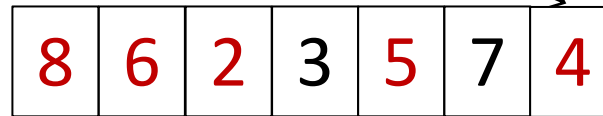


After eviction:



LRU+ACE (w/o PF)

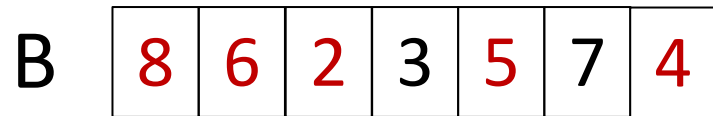
Candidate



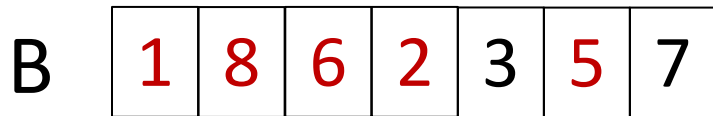
An Example ($k_w = 3$)

write page 1

LRU

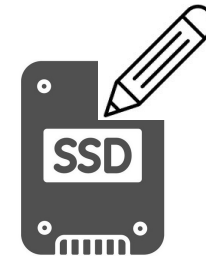
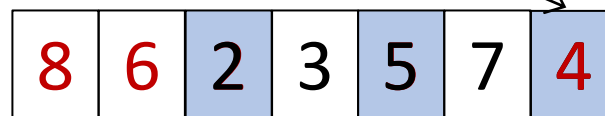


After eviction:



LRU+ACE (w/o PF)

Candidate



4,5,2 concurrently written
4 evicted

An Example ($k_w = 3$)

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

After eviction:

1	8	6	2	3	5	7
---	---	---	---	---	---	---

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

write page 1

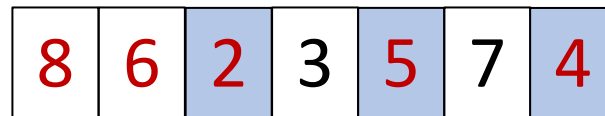
LRU



After eviction:



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



After eviction:



Candidate



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

write page 1

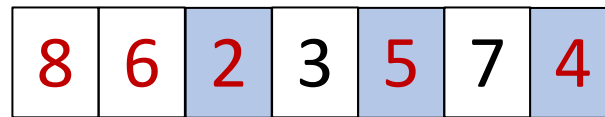
LRU



After eviction:



LRU+ACE (w/o PF)

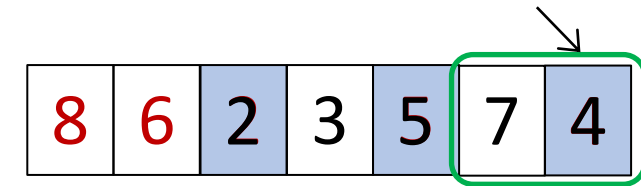


After eviction:



LRU+ACE (w/ PF)

eviction window



4,5,2 concurrently written
4,7 evicted

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

write page 1

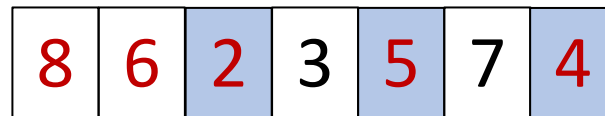
LRU



After eviction:



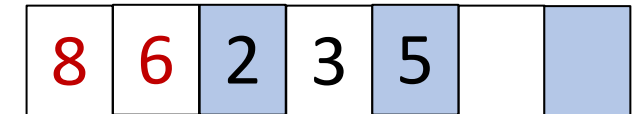
LRU+ACE (w/o PF)



After eviction:



LRU+ACE (w/PF)



After eviction:



↑
prefetched

Experimental Evaluation



Clock Sweep

LRU

CFLRU

LRU-WSR

vs their ACE counterparts

Device	α	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:

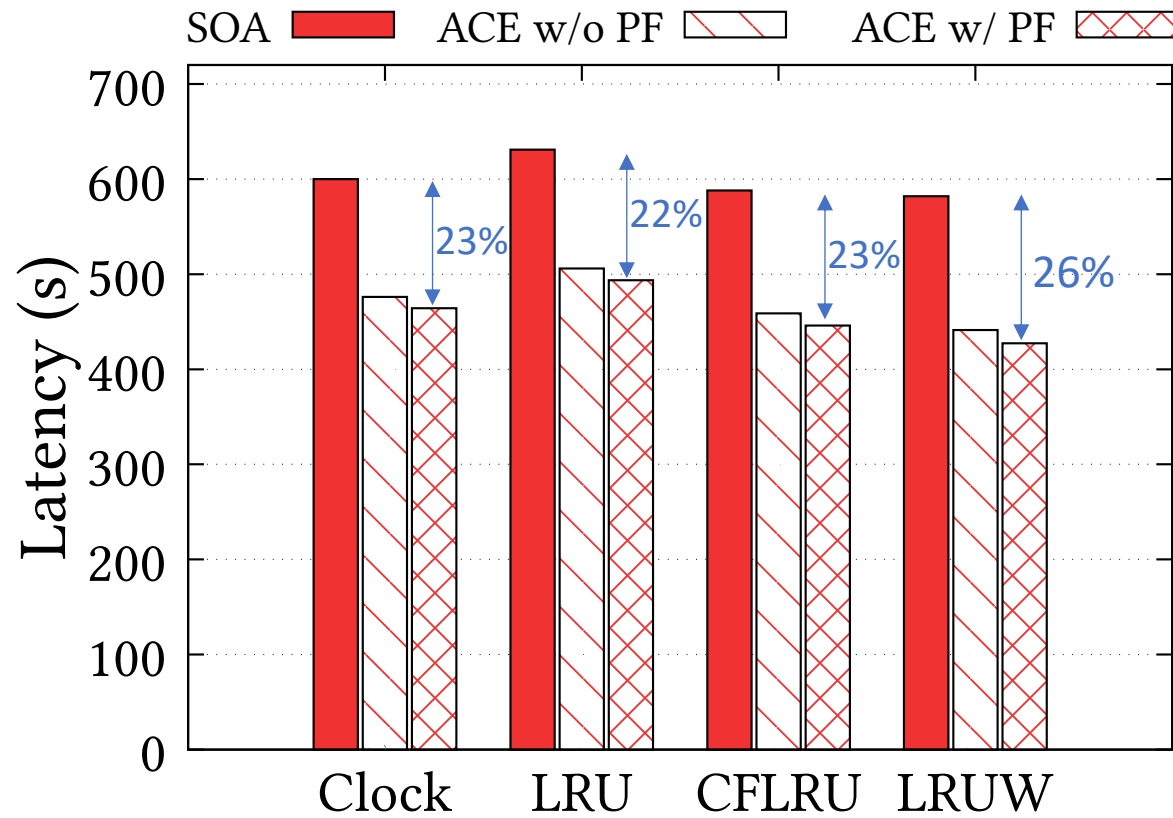
synthesized traces

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%)

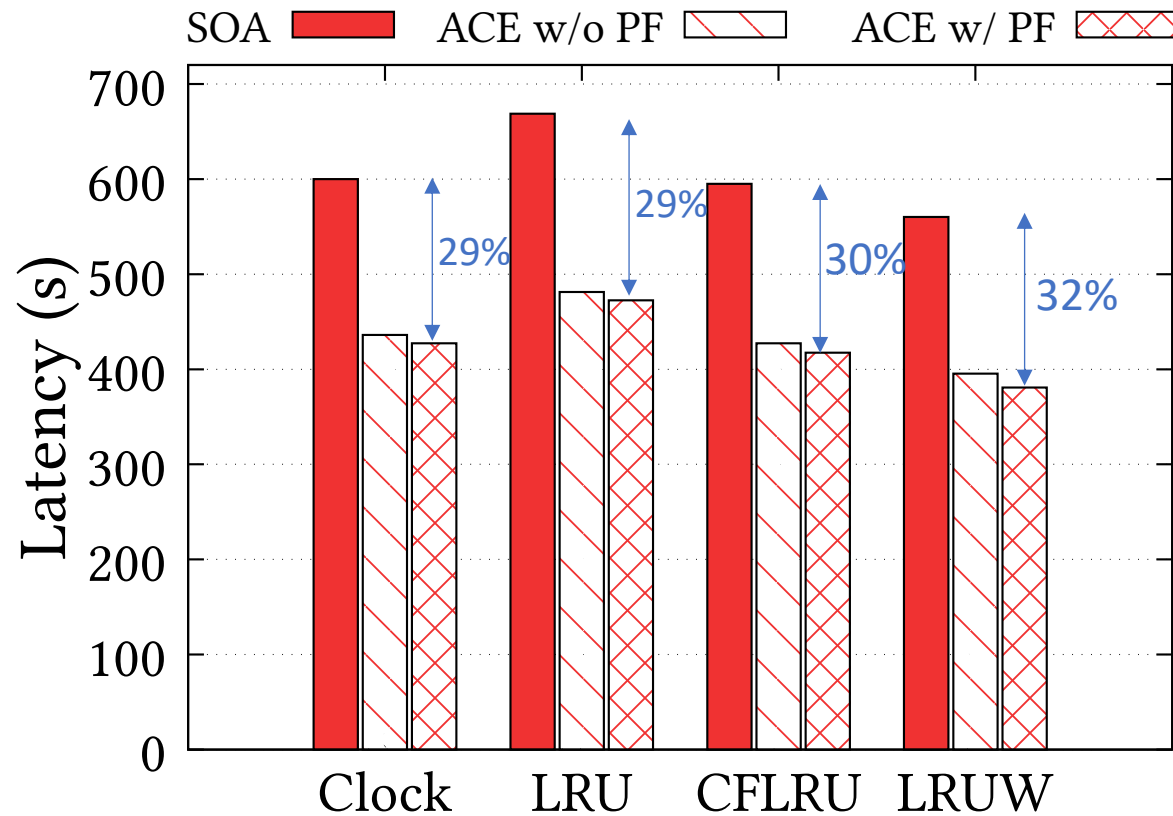
Benefit comes at no cost

Mixed Skewed Trace
(r/w: 50/50, locality 90/10)

Higher Gain for Write-Heavy Workload

Device: PCIe SSD

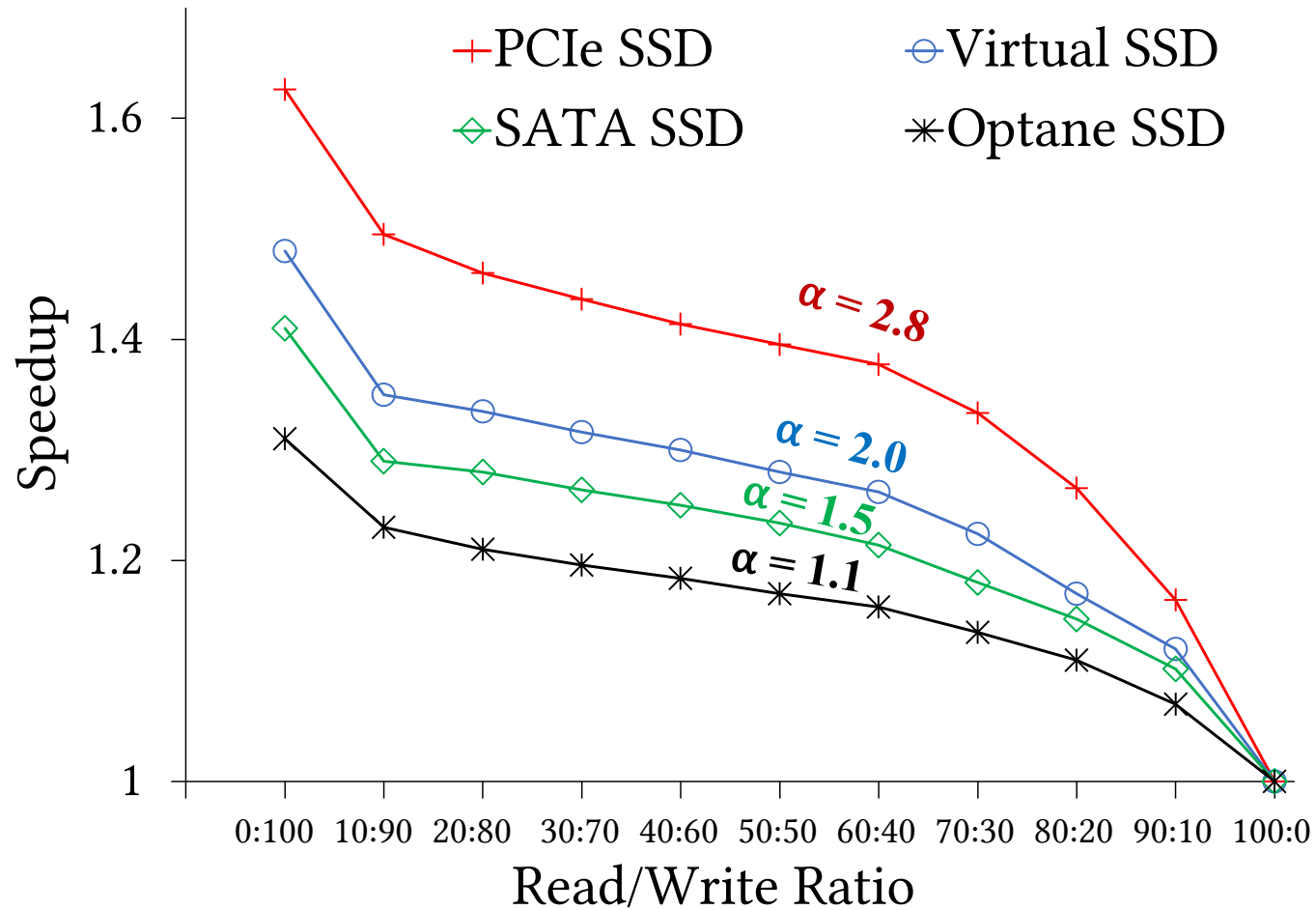
$\alpha = 2.8, k_w = 8$



Mixed Skewed Trace
(r/w: 50/50, locality 90/10)

Write-intensive workloads have
higher benefit (up to 32%)

Impact of R/W Ratio & Asymmetry



more writes, more speedup

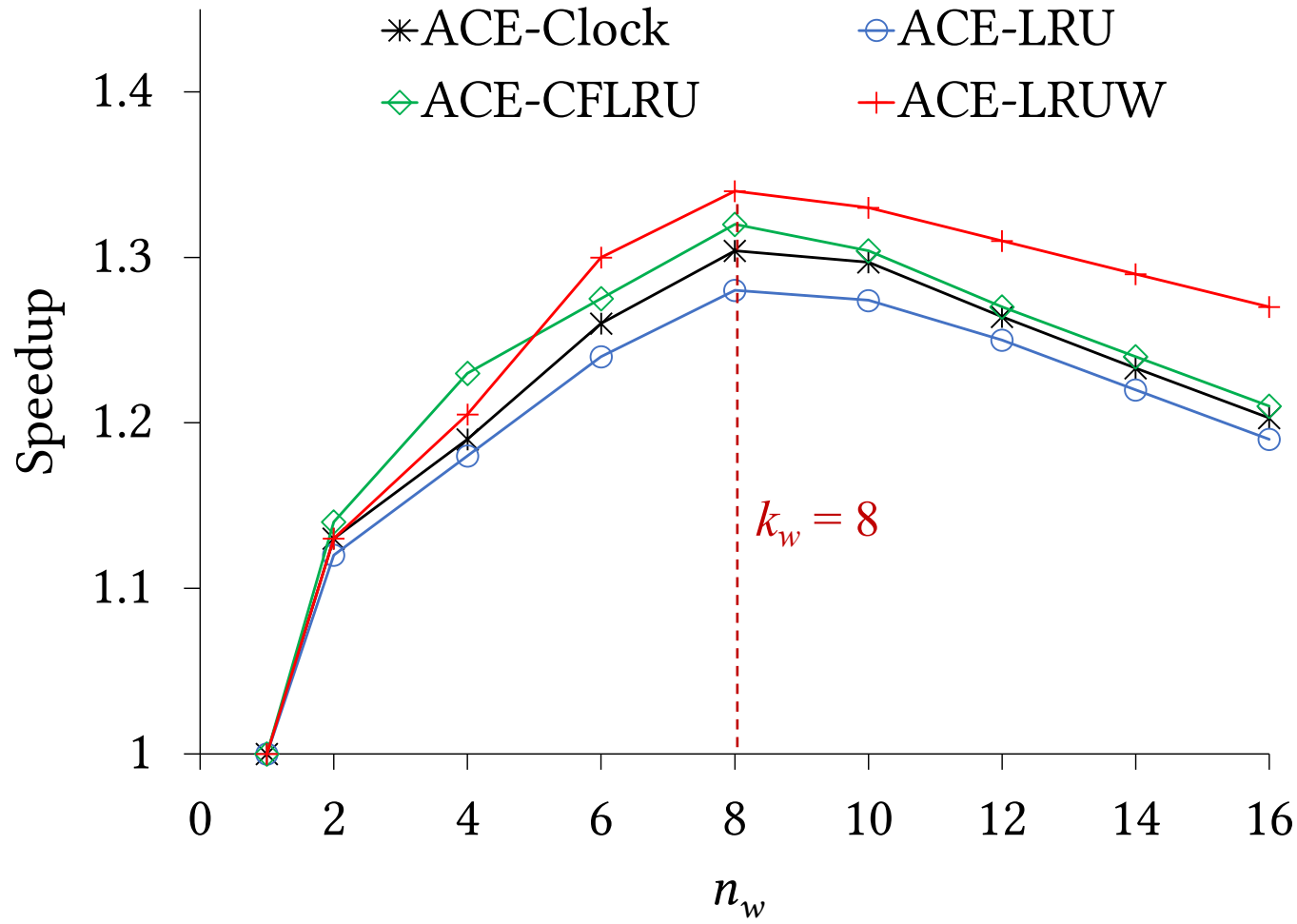
higher asymmetry, higher speedup

good benefit even for low asymmetry

Impact of #Concurrent I/Os

Device: PCIe SSD

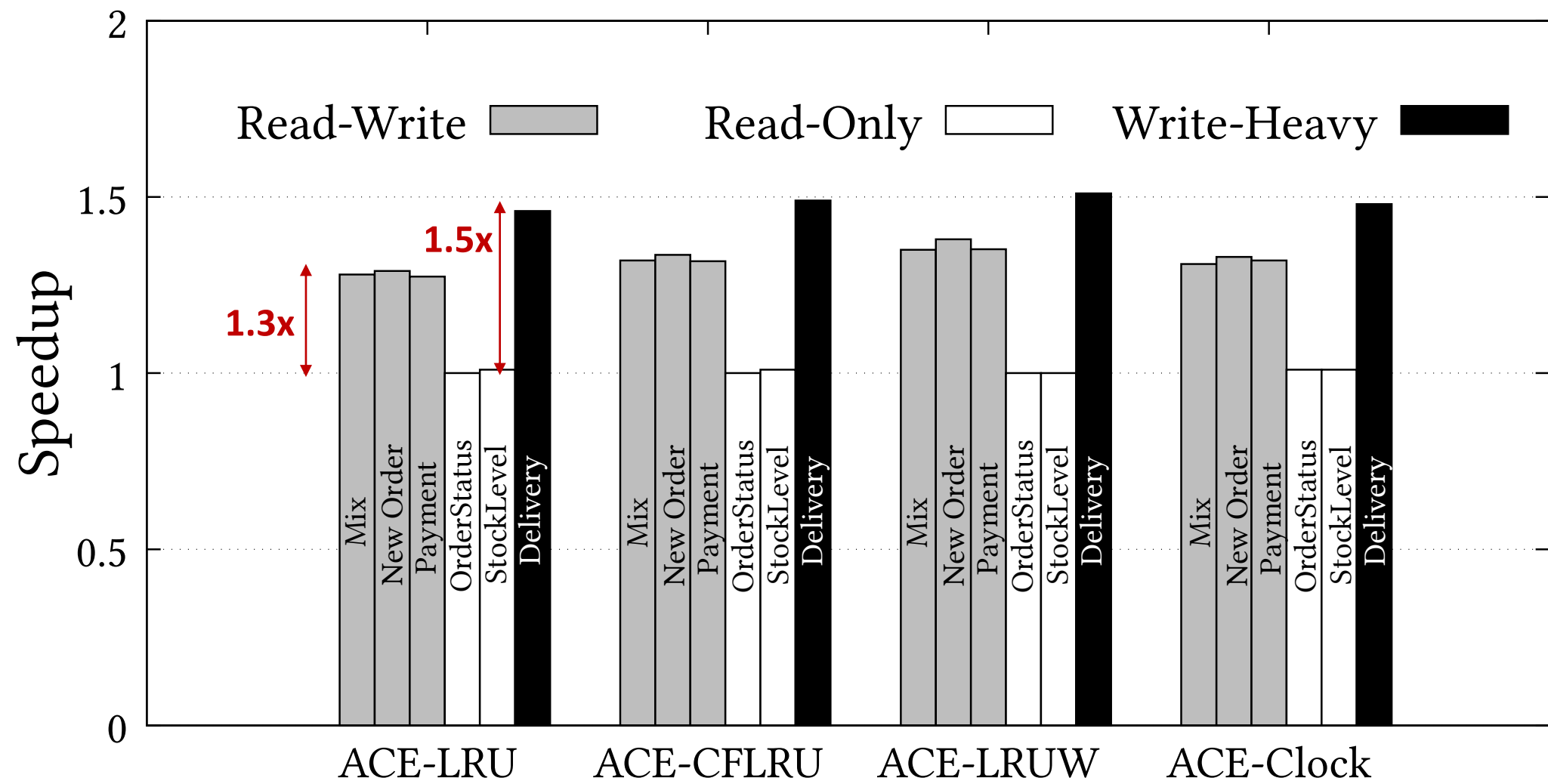
$\alpha = 2.8, k_w = 8$



Highest speedup when optimal concurrency is used

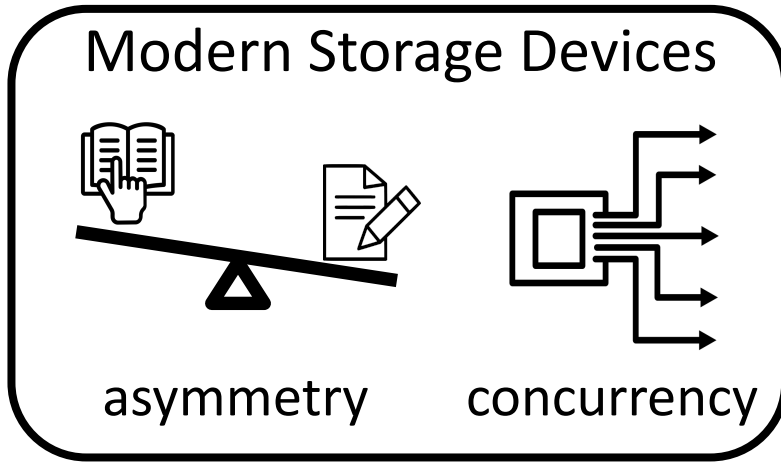
Mixed Skewed Trace
(r/w: 50/50, locality 90/10)

Experimental Evaluation (TPC-C)



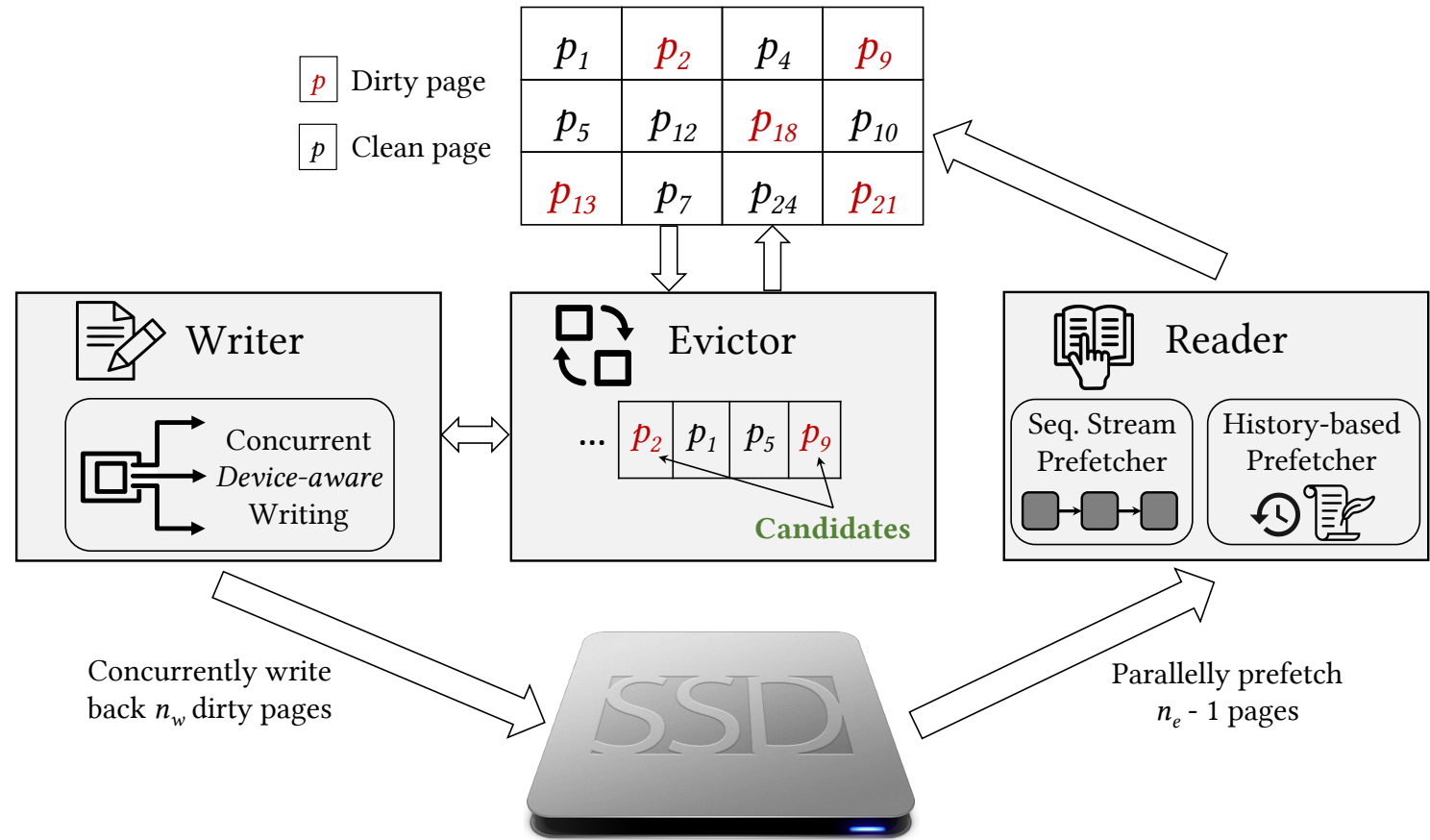
ACE Achieves 1.3x for mixed TPC-C

Summary



decoupled eviction and write-back mechanism
 can be integrated with **any** replacement policy
 benefit comes with no penalty

ACE Bufferpool



Thank You!

papon@bu.edu

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

Measuring Asymmetry & Concurrency

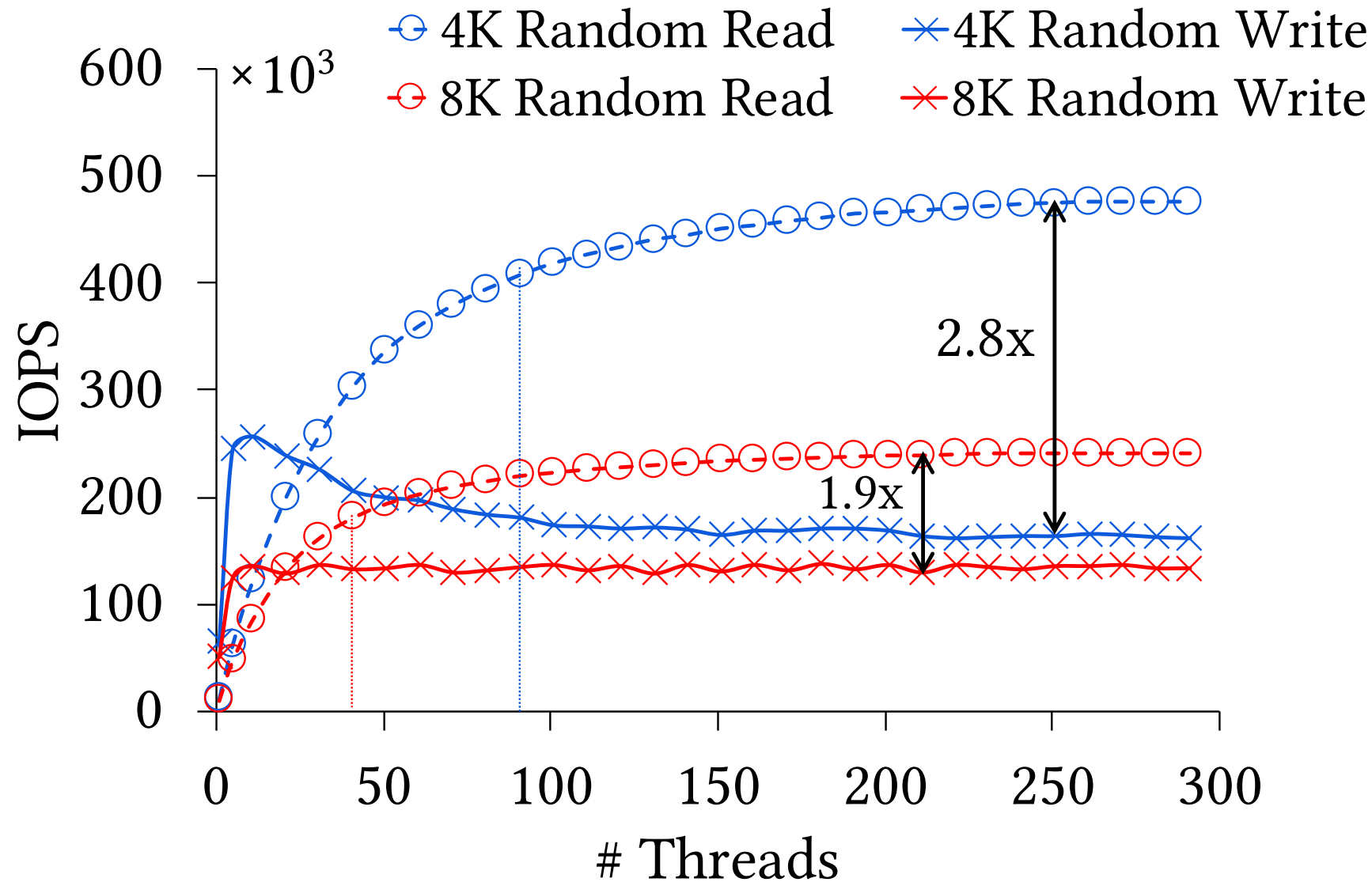
Device

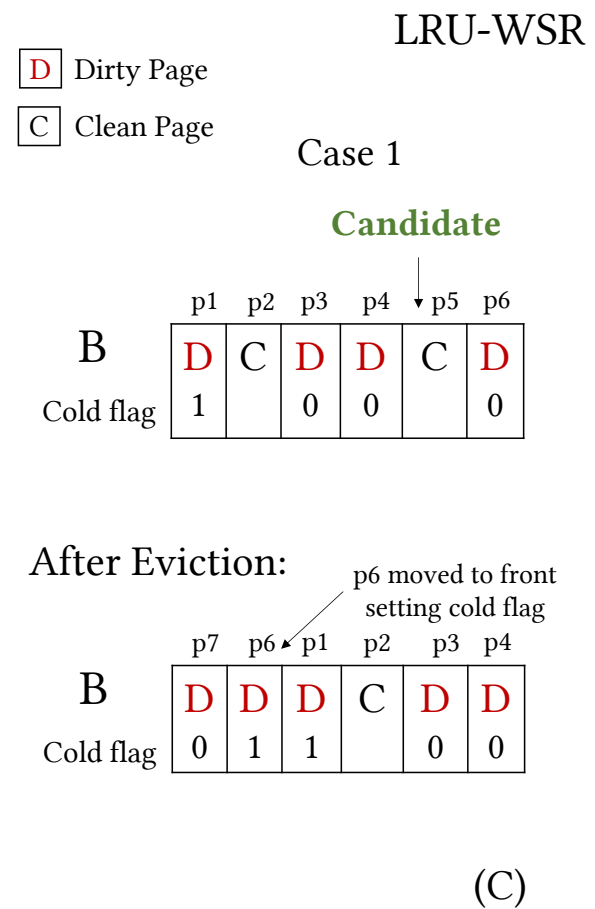
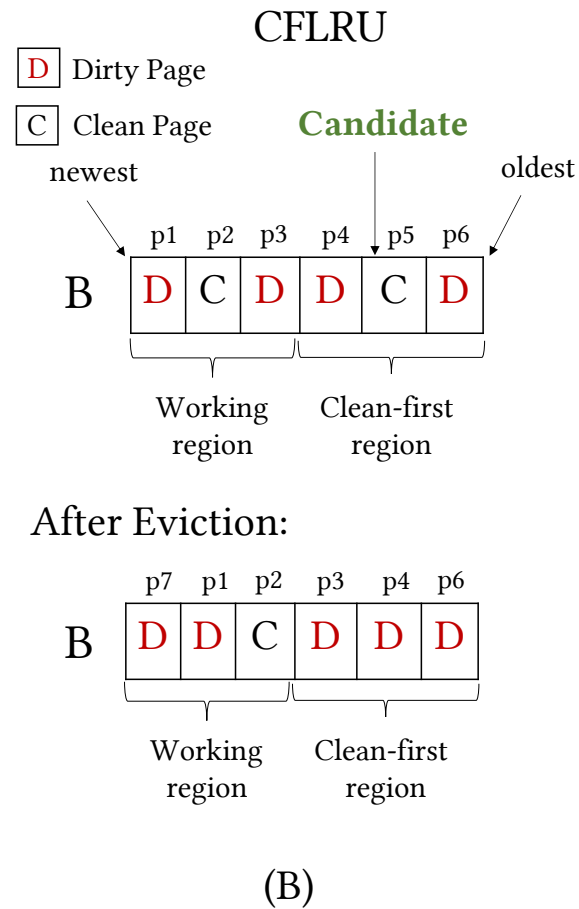
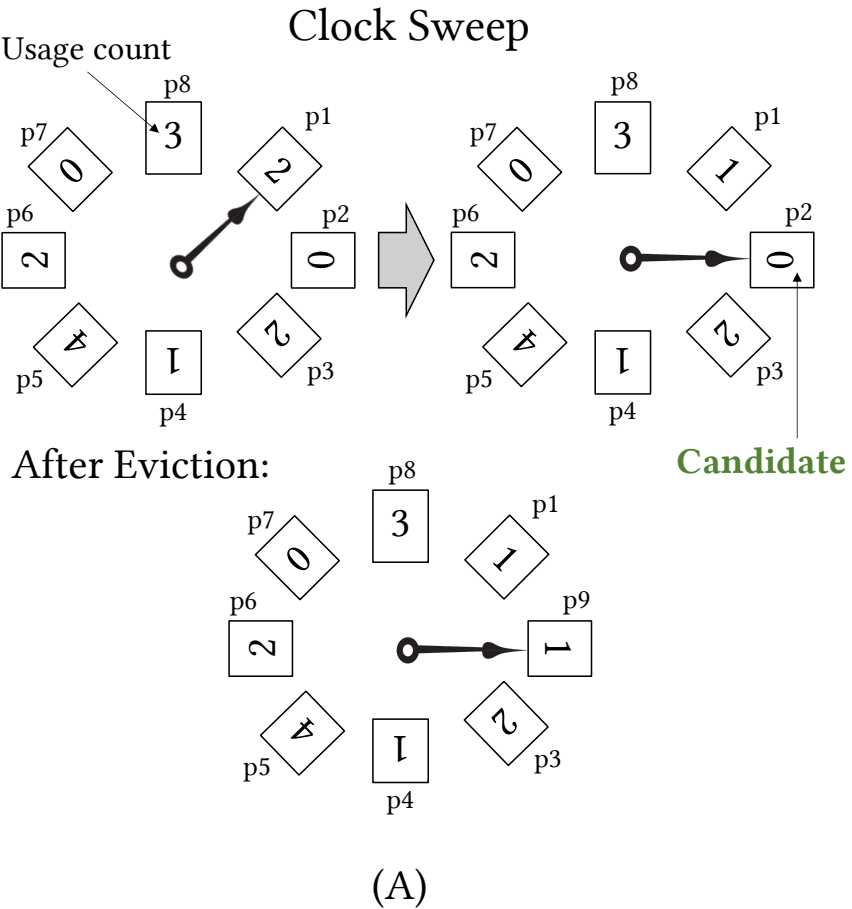
PCIe SSD - P4510 (1TB)

For 4K random read,

Asymmetry: 2.8

Concurrency: 80





```

Input:  $P, n_w, n_e$  is pf_enabled
1 //  $P$  is the accessed page
2 //  $n_w$  is the maximum effective write concurrency ( $n_w = k_w$ )
3 //  $n_e$  is the number of concurrent reads when prefetching is enabled
4 // is_pf_enabled determines if prefetching is enabled or not
5 if  $P$  in bufpool then
6 | return  $P$ 
7 else
8 // miss! need to bring  $P$  from disk
9 if bufpool not full then
10 | if is_pf_enabled == true then
11 | | // reads  $P$  and prefetches up to  $n_e - 1$  pages from disk
12 | | (depending on available slots)
13 | | - prefetch_pages ( $P, n_e - 1$ )
14 | | else
15 | | | - read  $P$  from disk
16 | | end if
17 | else
18 | |  $top\_page$  = replacement_policy.get_one_page_to_evict()
19 | | if  $top\_page$  is clean then
20 | | | // follow classical approach if page is clean
21 | | | - drop  $top\_page$  from bufpool
22 | | | - read  $P$  from disk
23 | | | else
24 | | | //  $top\_page$  is dirty. concurrently write  $n_w$  dirty pages
25 | | | //  $P_{wb}$  is a vector containing the candidate dirty pages
26 | | | -  $P_{wb}$  = populate_pages_to_writeback()
27 | | | - issue  $\|length(P_{wb})\|$  concurrent writes,  $\forall p \in P_{wb}$ 
28 | | | - mark  $\|length(P_{wb})\|$  pages as clean,  $\forall p \in P_{wb}$ 
29 | | | if is_pf_enabled == true then
30 | | | | // evict  $n_e$  pages
31 | | | | // pages written and to be evicted can be different
32 | | | | //  $P_{ev}$  is a vector containing the pages to evict
33 | | | |  $P_{ev}$  = replacement_policy.get_n_pages_to_evict()
34 | | | | - drop  $\|length(P_{ev})\|$  pages from bufpool,  $\forall p \in P_{ev}$ 
35 | | | | // Now, prefetch
36 | | | | - prefetch_pages ( $P, n_e - 1$ )
37 | | | | - empty  $P_{ev}$ 
38 | | | | else
39 | | | | | // evict 1 page
40 | | | | | - drop  $top\_page$  from bufpool
41 | | | | | - read  $P$  from disk
42 | | | | end if
43 | | | else
44 | | | | - empty  $P_{wb}$ 
45 | | | end if
46 | | end if
47 | end if
48 end if

```

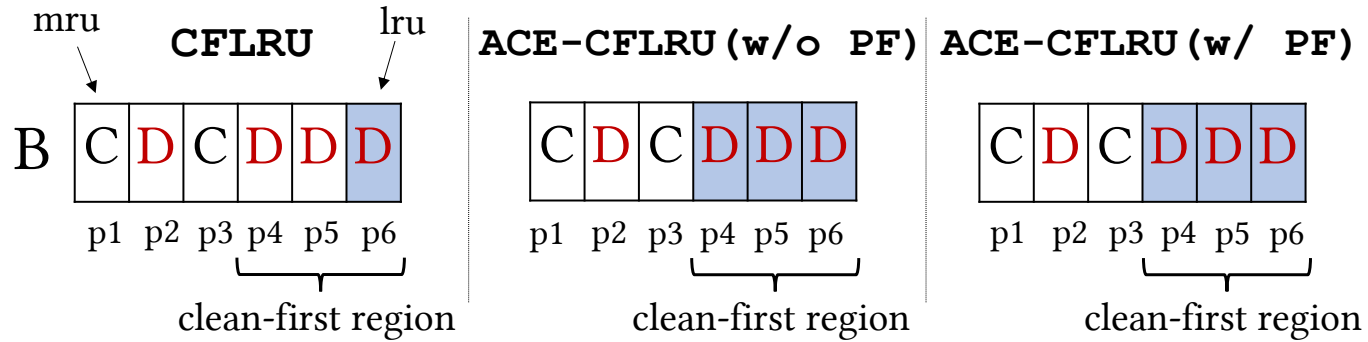
```

1 Procedure populate_pages_to_writeback()
2 | // follow the underlying page replacement policy to generate  $P_{wb}$ 
3 | - select next  $n_w$  dirty pages based on the underlying page replacement policy
4 | - return this vector

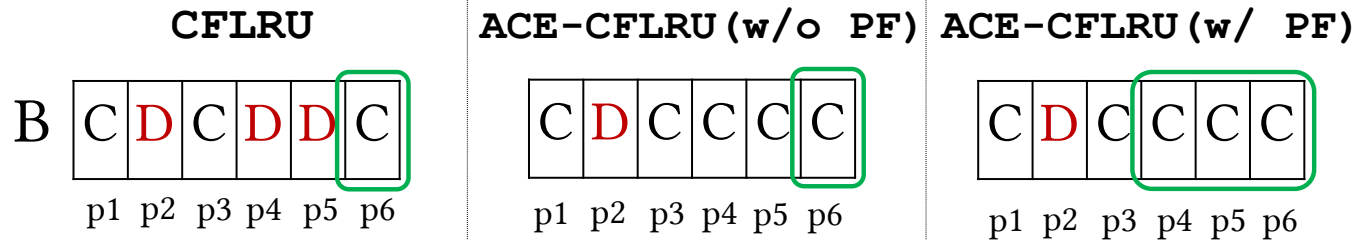
1 Procedure prefetch_pages(page  $P$ , int  $x$ )
2 | if  $P$  in Sequential_Table then
3 | | // start of a sequential stream!
4 | | // read  $P$  and the next  $x$  pages concurrently
5 | | - prefetch_sequential ( $P$ )
6 | | else
7 | | | // use the history based prefetcher
8 | | | // read  $P$  and  $x$  pages (selected by prefetcher) concurrently
9 | | | - prefetch_history ( $P$ )
10 | | end if
11 | /* note that  $P$  should be placed in the most recently used position
12 | | in the bufpool whereas other pages should be placed in the
13 | | least recently used positions */
14 | - place these  $x + 1$  pages into bufpool

```

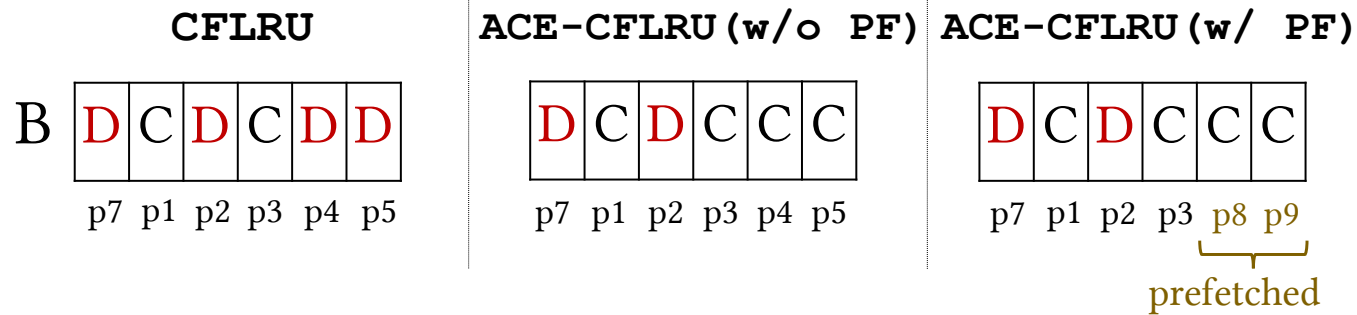
Initial State (Before 'write p7')



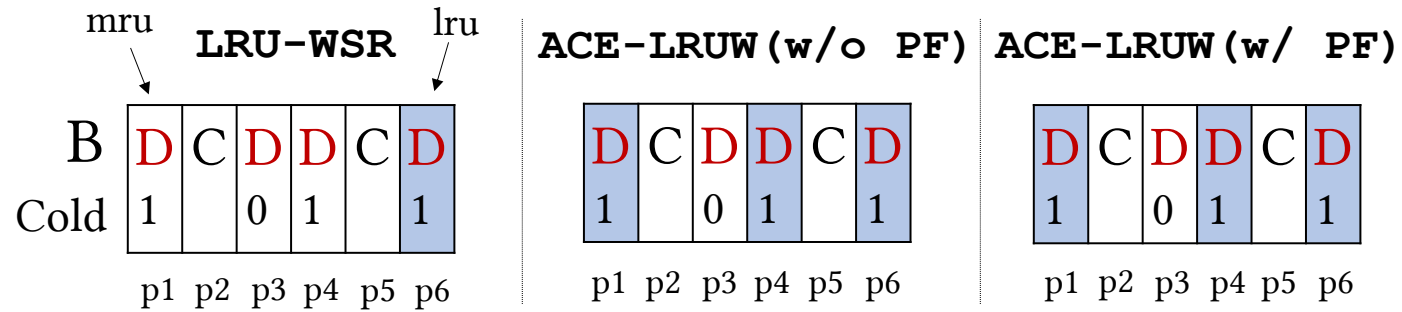
Intermediate State



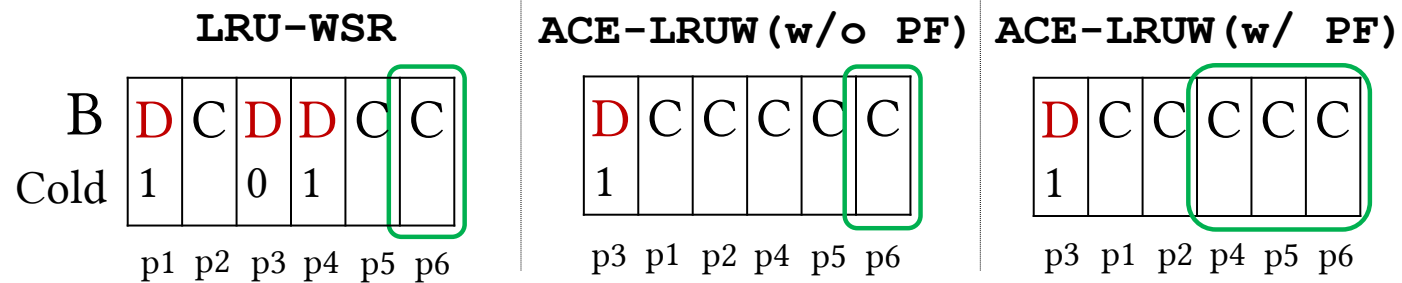
Final State (After 'write p7')



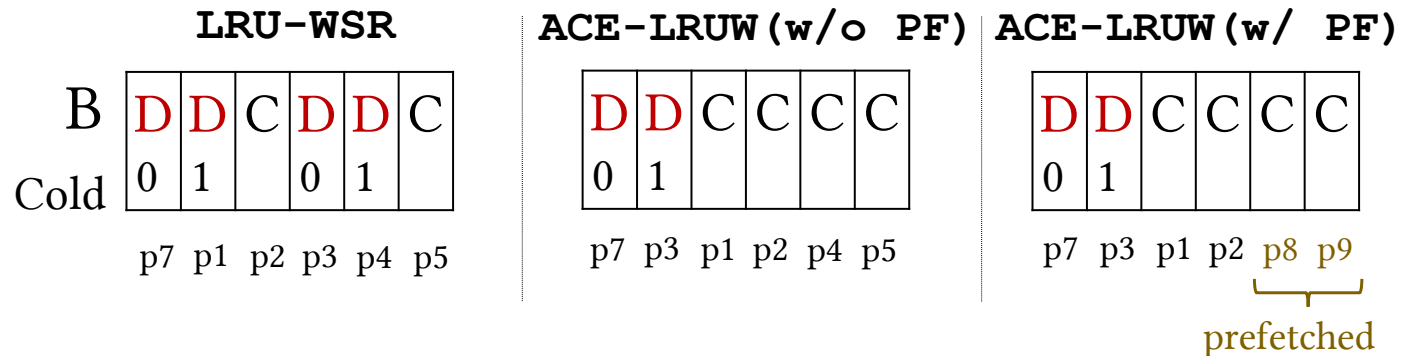
Initial State (Before 'write p7')



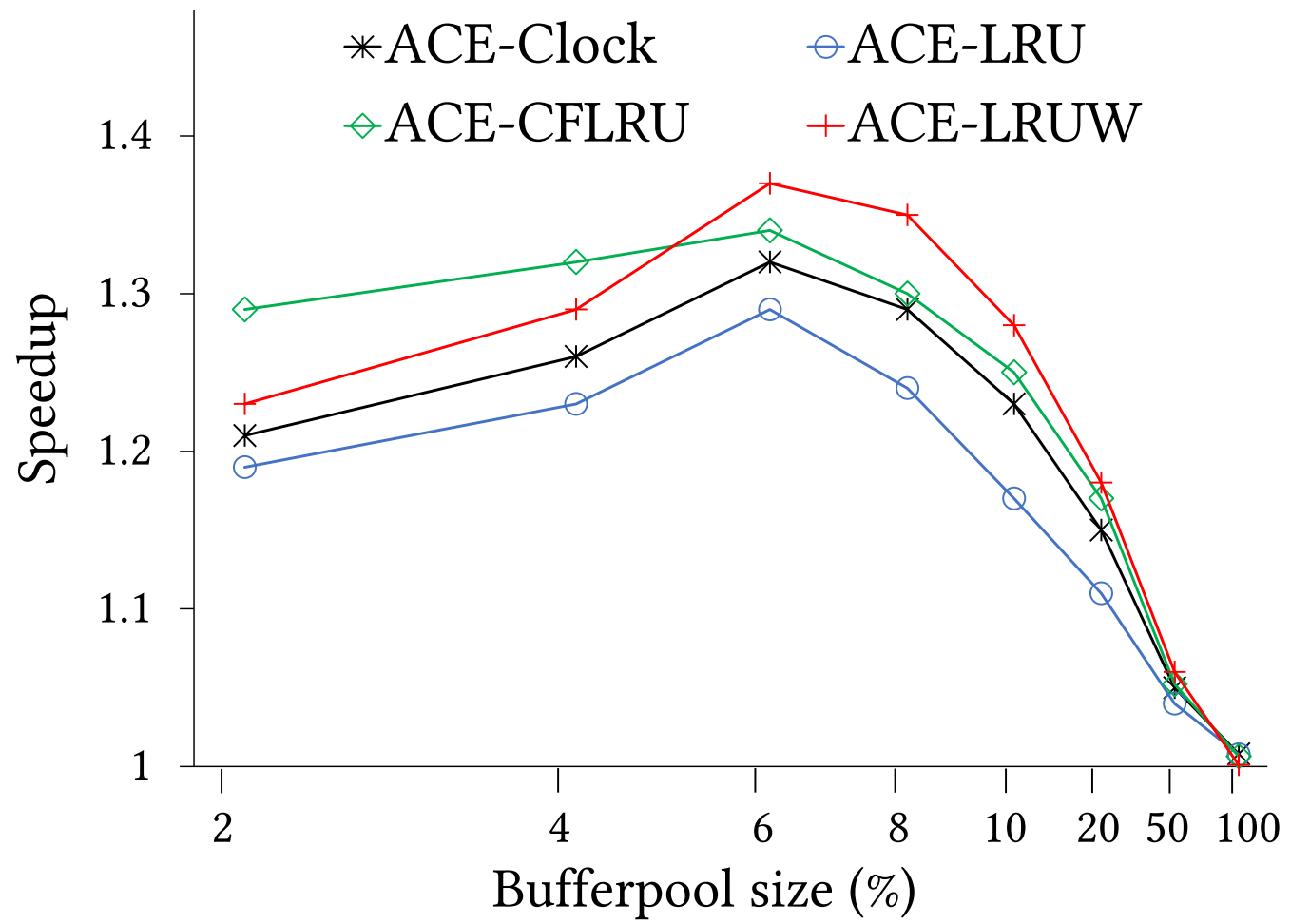
Intermediate State



Final State (After 'write p7')



Experimental Evaluation



ACE performs well under memory pressure

Impact on #writes

