

# CS 561: Data Systems Architectures

# Class 16

### Asymmetry and Concurrency Aware Storage Management

Manos Athanassoulis

https://bu-disc.github.io/CS561/

# Data Systems & Hardware



Memory Hierarchy

# Hardware Trends

**Evolution of Storage Technology** 



# Solid State Drives





# Goal: Developing Hardware-Aware Data Systems



SSD-Aware Systems [IEEE ICDE '24]

# Goal: Developing Hardware-Aware Data Systems



SSD-Aware Systems [IEEE ICDE '24]



# **SSD** Concurrency



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



Block 0



Block 0

Block 1

Writing in a free page isn't costly!









Block 0



Block 0



Not all updates are costly!

#### What if there is no space?



. . .

Block 0

Block N

#### What if there is no space?



**Garbage Collection!** 





Block 0

. . .

Block N



#### What if there is no space?



**Garbage Collection!** 



. . .

Block 0

Block N

# Read/Write Asymmetry in SSD



Higher average update cost (due to GC)  $\rightarrow$  *Read/Write asymmetry* 













# Impact of the File System



# Impact of the File System



# Impact of the File System



### Empirical Asymmetry and Concurrency

		4KB			8KB	
Devices	α	k <sub>r</sub>	$k_w$	α	k <sub>r</sub>	$k_w$
Optane SSD	1.1	6	5	1.0	4	4
PCIe SSD (with FS)	2.8	80	8	1.9	40	7
PCIe SSD (w/o FS)	3.0	16	6	3.0	15	4
SATA SSD	1.5	25	9	1.3	21	5

DaMoN@SIGMOD 2021

#### Guidelines for System Design in SSDs



# Goal: Developing Hardware-Aware Data Systems



SSD-Aware Systems [IEEE ICDE '24]

**IEEE ICDE 2023** 

# **Bufferpool is Tightly Connected to Storage**



#### IEEE ICDE 2023










### Traditional Bufferpool Manager



# The Challenges

With write asymmetry, exchanging

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

Without exploiting concurrency,

device remains vastly **underutilized**.







# Asymmetry/Concurrency-Aware (ACE) Bufferpool Manager

### ACE Bufferpool Manager



### Use device's properties









## ACE Bufferpool Manager



### An Example



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

Write request of page 8 comes

# An Example ( $k_w = 3$ )

write page 8 B  $\begin{bmatrix} 6 & 2 & 3 & 5 & 7 & 4 & 9 \end{bmatrix}$ Candidate for eviction  $\downarrow$ 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



SSD

ഀ൬൬ഀ

After eviction:





After eviction:





An Example 
$$(k_w = 3)$$

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



LRU

write page 1

After eviction:



After eviction:



more clean pages

An Example (
$$k_w = 3$$
)  
LRU+ACE (w/o PF) LRU+ACE (w/PF)  
Candidate

8

6

3

2

5

7

4



LRU

write page 1



After eviction:

В	1	8	6	2	3	5	7
	-	•		-	U		•

After eviction:

An Example 
$$(k_w = 3, n_e = 2)$$
  
write page 1  
LRU  
B 8 6 2 3 5 7 4  
After eviction:  
P 1 9 6 3 2 5 7 7  
An Example  $(k_w = 3, n_e = 2)$   
LRU+ACE  $(w/o PF)$  LRU+ACE  $(w/PF)$   
eviction window  
8 6 2 3 5 7 4  
After eviction:  
D 1 9 6 3 2 5 7  
An Example  $(k_w = 3, n_e = 2)$ 

R	1	8	6	2	3	5	7	
D	<b>_</b>	0	U	Ζ	5	5	/	

4,5,2 concurrently written 4,7 evicted

An Example 
$$(k_w = 3, n_e = 2)$$
  
write page 1  
LRU LRU+ACE (w/o PF) LRU+ACE (w/PF)







After eviction:

В	1	8	6	2	3	5	7
		_			-	-	-

After eviction:

After eviction:



# **Experimental Evaluation**



Device	α	k <sub>r</sub>	$k_w$
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<sub>w</sub> = 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<sub>w</sub> = 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: PCle SSD** 

IEEE ICDE 2023

$$\alpha$$
 = 2.8, k<sub>w</sub> = 8

### Highest speedup when optimal concurrency is used

# Experimental Evaluation (TPC-C)



# Experimental Evaluation (TPC-C)

TPC-C consists of 5 transactions

NewOrder (45%) R/W Mix

Payment (43%) R/W Mix

OrderStatus (4%) R-only

StockLevel (4%) R-only

Delivery (4%) W-heavy



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



**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

### Goal: Developing Hardware-Aware Data Systems



SSD-Aware Systems [IEEE ICDE '24]

# Rise of Large Graphs

Graphs are everywhere!



Social Network

**Physical Science** 

Transportation Network

Machine Learning

Real-world graphs often have more than a billion nodes

### Processing Large Graphs







**Distributed Systems** 

Single-node in-memory systems Single-node out-of-core systems

### Out of Core Systems







Data partitioning

Improve memory & disk locality

Reduce random I/O

**Designed for HDDs** 

# Our Goal

- Optimize for storage-based workload
- Focus on **traversal** operations
- Utilize efficient SSD concurrency by parallelizing independent I/Os
- Maintain **core** algorithm properties

### Concurrency-Aware Graph (V, E) Manager

### CAVE

# Concurrent Graph Algorithms

- Parallel Breadth-First Search
- Parallel pseudo Depth-First Search
- Parallel Weakly Connected Components
- Parallel PageRank
- Parallel Random Walk

# Parallel BFS





processing in progress



yet to be processed

Each iteration involves

- 1. processing a list of vertices aka the **frontier**
- 2. accessing the neighbors of each vertex
- 3. updating vertex values
- determining which vertices should be visited in the next iteration


# Parallel BFS





processing in progress



Each iteration involves

- 1. processing a list of vertices aka the **frontier**
- 2. accessing the neighbors of each vertex
- 3. updating vertex values
- determining which vertices should be visited in the next iteration











## **Experimental Evaluation**

Dataset	Description	#Nodes	#Edges	Diameter	Size
FS	Friendster Social Network	65M	1.8B	32	32 GB
TW	Twitter Social Network	53M	2B	18	28 GB
RN	RoadNet Network of PA	1M	1.5M	786	47 MB
LJ	LiveJournal Social Network	5M	69M	16	1 GB
YT	YouTube Social Network	1.1M	3M	20	39 MB
SD	Synthetic data	50M	1.25B	6	20 GB

6 datasets

3 devices Optane SSD ( $k_r = 6$ ) PCIe SSD ( $k_r = 80$ ) SATA SSD ( $k_r = 25$ )

Approaches Used:

GraphChi, GridGraph, Mosaic, CAVE, CAVE\_blocked

## CAVE's Preprocessing is Efficient

System	Preprocessing Time (s)		Data File Size (GB)		
System	Dataset: FS	Dataset: TW	Dataset: FS	Dataset: TW	
GraphChi	819	784	8.3	8.4	
GridGraph	55	86	84	75	
Mosaic	469	370	27	17	
CAVE	52	49	14	13	

### CAVE Performs Efficient PBFS



### CAVE Performs Efficient PBFS



## CAVE Performs Efficient PBFS



Both CAVE implementations outperforms GridGraph, Mosaic and GraphChi

# CAVE Utilizes Concurrent I/O

Dataset: FS



SATA SSD ( $k_r = 25$ ) PCIe SSD ( $k_r = 80$ ) Optane SSD ( $k_r = 6$ )

#### **Device gets saturated at** *optimal concurrency*

## Conclusion

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

