# CSE462/562: Database Systems (Spring 23) Lecture 3: Buffer Management 2/7/2023





User applications		
DBMS	SQL Parser/API	
	Query Execution	
	File Organization/Access Methods	
	Buffer Management	
	Disk space/File management	
	Operating System	

Hardware devices







#### How does database access data pages?

- Data pages are stored in disk file
  - suppose we want to count how many rows there are in a database table
    - need to scan all pages
  - page must be loaded into memory before any computation happens



#### How does database access data pages?

- Data pages are stored in disk file
  - suppose we want to count how many rows there are in a database table
    - need to scan all pages
  - page must be loaded into memory before any computation happens



#### How does database access data pages?

- Data pages are stored in disk file
  - suppose we want to count how many rows there are in a database table
    - need to scan all pages
  - page must be loaded into memory before any computation happens
- What if we want to scan the data file for multiple passes?
  - Option 1: read/write the entire page on demand before reading/writing the integer <- very slow
  - Option 2: read all data pages into memory at the beginning <- not scalable
    - May not fit in memory
    - What to do on modify?
      - Immediately write back? Or Flush when program shutsdown?
      - Data persistence?
- Solution: buffer pool

# Buffer management in DBMS

- Buffer manager manages a fixed-size pool of in-memory page frames which
  - are of the same size as the data pages (e.g., 4KB)



- Handling page request
  - Suppose we want to read/write a page in the file with page number = 100



- Handling page request
  - Suppose we want to read/write a page in the file with page number = 100



- Handling page request
  - Suppose we want to read/write a page in the file with page number = 100



- Handling page request
  - Suppose we want to read/write a page in the file with page number = 100



- Handling page request
  - Suppose we want to read/write a page in the file with page number = 100



- Handling page request
  - Suppose we want to read/write a page in the file with page number = 100



*Cost: 1 I/O* 

- Handling page request
  - Suppose we want to read the same page again (pid = 100)



- Handling page request
  - Suppose we want to read the same page again (pid = 100)



- Handling page request
  - Suppose we want to read the same page again (pid = 100)



*Cost: 0 I/O* 

• How to implement line 2?

2 if pid exists in some buffer frame i:

- Need to store the page numbers, but where?
  - For each buffer frame, we maintain a metadata structure which includes **pid**.





CSE462/562 (Spring 2023): Lecture 3



- Practical consideration for hash tables
  - DBMS usually has its own hash tables implementation for buffer manager -- why?
    - memory constraints, efficiency, concurrency control, ...



- Practical consideration for hash tables
  - For Project 2: feel free to use libraries (e.g., absl::flat\_hash\_map, std::unordered\_map)
    - Tips for time and memory efficiency: avoid rehashing
      - Set the initial bucket count K >= m / max\_load\_factor()



- What if we run out of buffer frames?
  - e.g., we are scanning a table with N = 100 pages, but buffer pool size m = 10



- What if we run out of buffer frames?
  - Buffer eviction: choose a victim to remove from the buffer pool
    - Several possible policies (more on this later)



- What if we run out of buffer frames?
  - Buffer eviction: choose a victim to remove from the buffer pool
    - Several possible policies (more on this later)



- What if we run out of buffer frames?
  - Buffer eviction: choose a victim to remove from the buffer pool
    - Several possible policies (more on this later)



#### Page requested for writes

- Potential problem with page eviction?
  - What if the evicted page is modified? (e.g., UPDATE A SET x = x + 10 WHERE ...)
    - We must write modified page back before eviction



#### Page requested for writes

- Potential problem with page eviction?
  - What if the evicted page is modified? (e.g., UPDATE A SET x = x + 10 WHERE ...)
    - We must write modified page back before eviction



#### Page requested for writes

- Potential problem with page eviction?
  - What if the evicted page is modified? (e.g., UPDATE A SET x = x + 10 WHERE ...)
    - We must write modified page back before eviction



- Problems with concurrency
  - One thread reading a block while the other tries to evict it

T1: char \* frame = BufMgr.HandlePageRequest(110) // &frames[0]



- Problems with concurrency
  - One thread reading a block while the other tries to evict it

T1: char \* f1 = BufMgr.HandlePageRequest(110) // &frames[0] *f1 now contains a wrong page for T1* T2: char \* f2 = BufMgr.HandlePageRequest(99) // &frames[0]



- Solution: introducing a buffer pin count per buffer frame
  - Upon page request, pin count++
  - Upon page release, pin count--
  - Never evict a page with pin count > 0



- Solution: introducing a buffer pin count per buffer frame
  - Upon page request, pin count++
  - Upon page release, pin count--
- T1: BufferId b1 = BufMgr.PinPage(110, &f1) // b1 = 0
- Never evict a page with pin count > 0



- Solution: introducing a buffer pin count per buffer frame
  - Upon page request, pin count++
  - Upon page release, pin count--
  - Never evict a page with pin count > 0

T1: BufferId b1 = BufMgr.PinPage(110, &f1) // b1 = 0

T2: BufferId b2 = BufMgr.PinPage(99, &f2) // b2 = 1



- Solution: introducing a buffer pin count per buffer frame
  - Upon page request, pin count++
  - Upon page release, pin count--
  - Never evict a page with pin count > 0

T1: BufferId b1 = BufMgr.PinPage(110, &f1) // b1 = 0

T2: BufferId b2 = BufMgr.PinPage(99, &f2) 
$$// b2 = 1$$

T1: BufMgr.UnpinPage(b1)



- Solution: introducing a buffer pin count per buffer frame
  - Upon page request, pin count++
  - Upon page release, pin count--
  - Never evict a page with pin count > 0

T1: BufferId b1 = BufMgr.PinPage(110, &f1)// b1 = 0T2: BufferId b2 = BufMgr.PinPage(99, &f2)// b2 = 1

T1: BufMgr.UnpinPage(b1)



# **Eviction policy**

- How do we choose a victim for eviction?
  - Randomly? The one with the lowest buffer ID that is not pinned? (Inefficient!)



# **Eviction policy**

- Eviction policy (aka replacement policy)
  - An algorithm for choosing unpinned frames when there's no free frame
    - It can have huge impacts on the # of I/Os, depending on the access pattern
  - Many common choices:
    - Least recently used (LRU)
    - Most recently used (MRU)
    - Clock
    - Database workload specific policies
    - ...

- Least Recently Used (LRU)
  - for each page in buffer pool, the order of the pages were last unpinned
  - replace the frame which has the oldest (earliest) time
  - very common policy: intuitive and simple
    - Works well for repeated accesses to popular pages -> typical transactional workload
- Example: P stands for PinPage, and U stands for UnpinPage, m = 3
  - P(1), P(2), P(3), U(2), U(3), P(4), U(1), P(3), U(3)



- Least Recently Used (LRU)
  - for each page in buffer pool, the order of the pages were last unpinned
  - replace the frame which has the oldest (earliest) time
  - very common policy: intuitive and simple
    - Works well for repeated accesses to popular pages -> typical transactional workload
- Example: P stands for PinPage, and U stands for UnpinPage, m = 3
  - P(1), P(2), P(3), U(2), U(3), P(4), U(1), P(3), U(3)



- Least Recently Used (LRU)
  - for each page in buffer pool, the order of the pages were last unpinned
  - replace the frame which has the oldest (earliest) time
  - very common policy: intuitive and simple
    - Works well for repeated accesses to popular pages -> typical transactional workload
- Example: P stands for PinPage, and U stands for UnpinPage, m = 3
  - P(1), P(2), P(3), U(2), U(3), P(4), U(1), P(3), U(3)



LRU list:



- Least Recently Used (LRU)
  - for each page in buffer pool, the order of the pages were last unpinned
  - replace the frame which has the oldest (earliest) time
  - very common policy: intuitive and simple
    - Works well for repeated accesses to popular pages -> typical transactional workload
- Example: P stands for PinPage, and U stands for UnpinPage, m = 3

• P(1), P(2), P(3), U(2), U(3), P(4), U(1), P(3), U(3)



How to implement in practice?

Exercise: how to remove a node in the middle of LRU list when there's a buffer hit?

- Least Recently Used (LRU)
  - for each page in buffer pool, the order of the pages were last unpinned
  - replace the frame which has the oldest (earliest) time
  - very common policy: intuitive and simple
    - Works well for repeated accesses to popular pages -> typical transactional workload
- Example: P stands for PinPage, and U stands for UnpinPage, m = 3

• P(1), P(2), P(3), U(2), U(3), P(4), U(1), P(3), U(3)

1 victim for eviction



- Least Recently Used (LRU)
  - for each page in buffer pool, the order of the pages that were last unpinned
  - replace the frame which has the oldest (earliest) time
  - very common policy: intuitive and simple
    - Works well for repeated accesses to popular pages -> typical transactional workload
- Problems?
  - Sequential flooding:
    - # buffer frames < # pages in file means every existing page in the buffer gets evicted
    - Prevents buffer hit for other transactions working on other files
- DB may know the access pattern before hand so that it can adapt its replacement policies
  - e.g., using a small ring buffer for sequential scan to avoid flooding the entire buffer pool

- Approximate LRU
- Each buffer frame has a clock bit
  - Set upon page pinned or unpinned (why?)
- When we need an eviction, move the clock hand
  - Ignore any page that is still pinned
  - Otherwise
    - If bit is set, clear it
    - If bit is clear, evict it
    - i.e., second chance



- Approximate LRU
- Each buffer frame has a clock bit
  - Set upon page pinned or unpinned (why?)
- When we need an eviction, move the clock hand
  - Ignore any page that is still pinned
  - Otherwise
    - If bit is set, clear it
    - If bit is clear, evict it
    - i.e., second chance



- Approximate LRU
- Each buffer frame has a clock bit
  - Set upon page pinned or unpinned (why?)
- When we need an eviction, move the clock hand
  - Ignore any page that is still pinned
  - Otherwise
    - If bit is set, clear it
    - If bit is clear, evict it
    - i.e., second chance
- Why this might be faster and easier to implement than LRU?
  - Hint: put the clock bit into the buffer meta structures
    - scan buffer meta structures instead



- Approximate LRU
- Each buffer frame has a clock bit
  - Set upon page pinned or unpinned (why?)
- When we need an eviction, move the clock hand
  - Ignore any page that is still pinned
  - Otherwise
    - If bit is set, clear it
    - If bit is clear, evict it
    - i.e., second chance
- Alternative: third/fourth/... chance
  - allowing clock counters up to 2/3/...



#### **Buffer flush**

- When are dirty pages written back to disk?
  - When evicted
  - During shutdown
  - Forced flush: flushing certain dirty pages to disk
    - when data need to be persisted for data consistency
    - only unpinned page may be flushed
    - other constraints apply (discussed later this semester)

OS does disk space & buffer management as well: why not let OS manage these tasks?

- Some limitations, e.g., files can't span disks.
- Buffer management in DBMS requires ability to:
  - pin a page in buffer pool, force a page to disk & order writes (important for implementing CC, concurrency control, & recovery)
  - adjust *eviction policy,* and prefetch pages based on access patterns in typical DB operations.

#### Summary

- Buffer management in DBMS
  - Buffer manager implementation
  - Eviction policy
- Next lecture
  - Data storage layout
- Project 2 released today (2/7/2023)
  - Due on Thursday, 2/16/2023, 1:00 AM EST
  - Write-up due on Saturday, 2/18/2023, 1:00 AM EST
  - Solution code for project 1 will be released on Piazza by the end of today.
- HW1 released today (not graded; no need to submit)
  - Solution to HW 1 will be released in a week