

# Buffer Pools



Lecture #05



Database Systems  
15-445/15-645  
Fall 2018

AP

Andy Pavlo  
Computer Science  
Carnegie Mellon Univ.

# UPCOMING DATABASE EVENTS

---

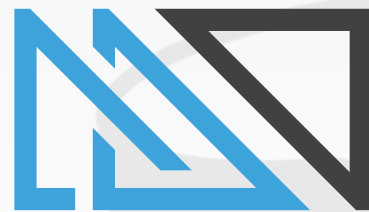
## Relational AI Talk

- Wednesday Sep 12<sup>th</sup> @ 4:00pm
- GHC 8102

**relationalAI**

## MapD Talk

- Thursday Sept 20<sup>th</sup> @ 12pm
- CIC 4<sup>th</sup> Floor



**MAPD**

# DATABASE STORAGE

---

**Problem #1:** How the DBMS represents the database in files on disk.

**Problem #2:** How the DBMS manages its memory and move data back-and-forth from disk.

# DATABASE STORAGE

---

## **Spatial Control:**

- Where to write pages on disk.
- The goal is to keep pages that are used together often as physically close together as possible on disk.

## **Temporal Control:**

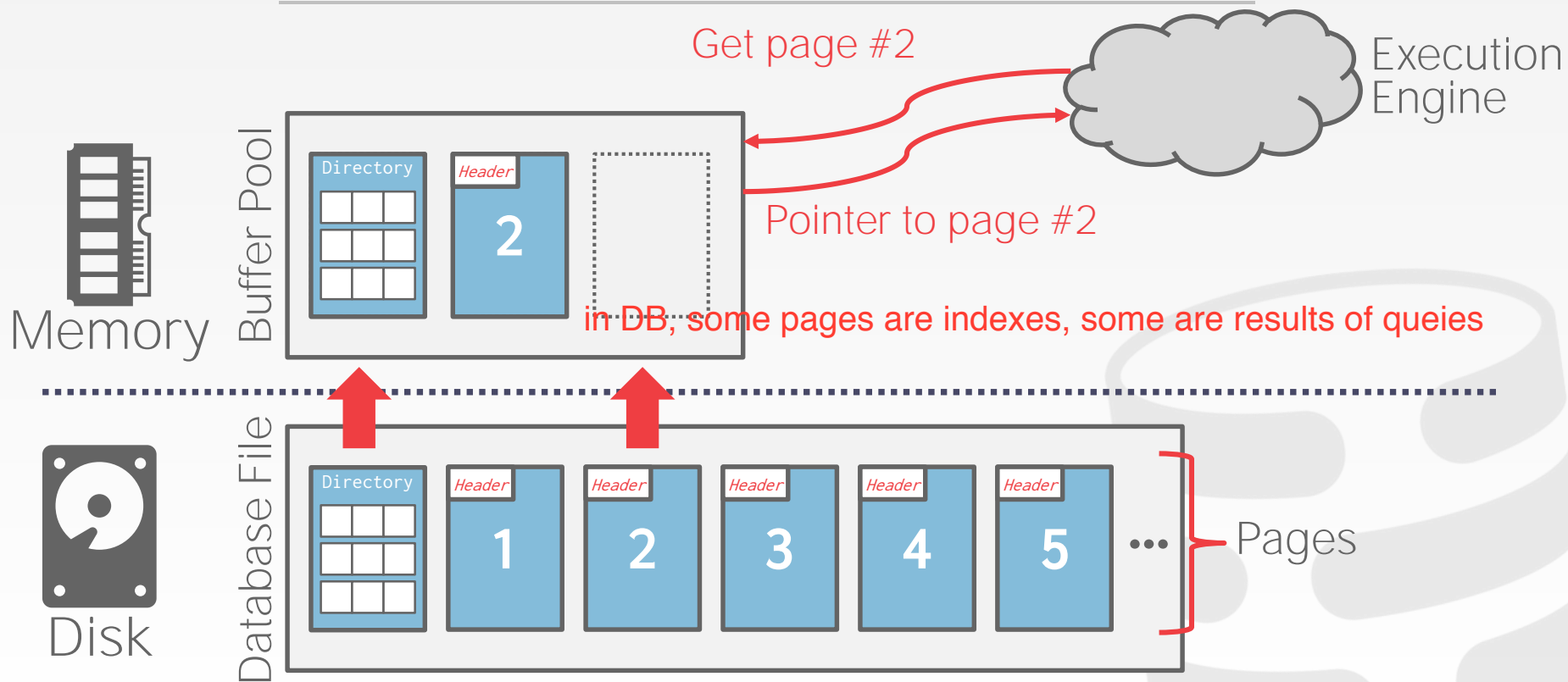
- When to read pages into memory, and when to write them to disk.
- The goal is minimize the number of stalls from having to read data from disk.

Buffer pool: maintain all the pages  
that read from the files on disk

The minimal granularity of OS reading and writing is page  
OS treats all pages all the same

5

## DISK-ORIENTED DBMS



# TODAY'S AGENDA

---

Buffer Pool Manager

Replacement Policies

Allocation Policies

Other Memory Pools



# BUFFER POOL ORGANIZATION

Memory region organized as an array of fixed-size pages.

An array entry is called a frame.

When the DBMS requests a page, an exact copy is placed into one of these frames.

Buffer  
Pool



On-Disk File

File -> Pages -> Slots -> Tuples in a slot

-> Frames (in buffer pool)

To check if one page is in the buffer pool, otherwise go find it from disk

## BUFFER POOL META-DATA

The page table keeps track of pages that are currently in memory.

Also maintains additional meta-data per page:

→ **Dirty Flag**

→ **Pin/Reference Counter**

Page  
Table

page1
page3

Buffer  
Pool

page1
page3
frame3
frame4

---

page1	page2	page3	page4
-------	-------	-------	-------

On-Disk File



# BUFFER POOL META-DATA

The page table keeps track of pages that are currently in memory.

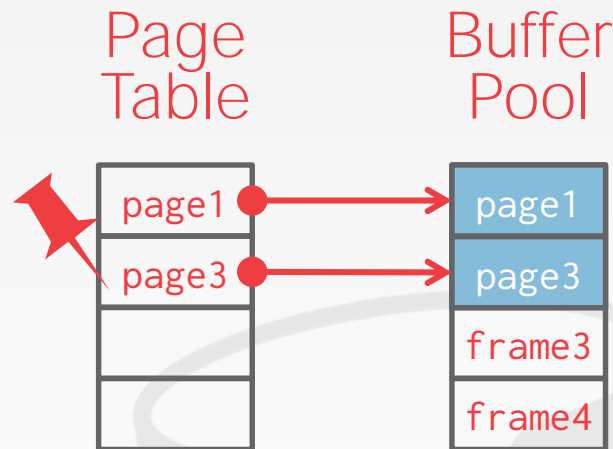
Also maintains additional meta-data per page:

→ **Dirty Flag**

→ **Pin/Reference Counter** prevent buffer pool removing this page if it is being referring

anytime when accessing to a page in buffer pool,  
we need to pin it in the page table

metadata of page in the buffer pool will be removed once this page get evicted from memory



On-Disk File

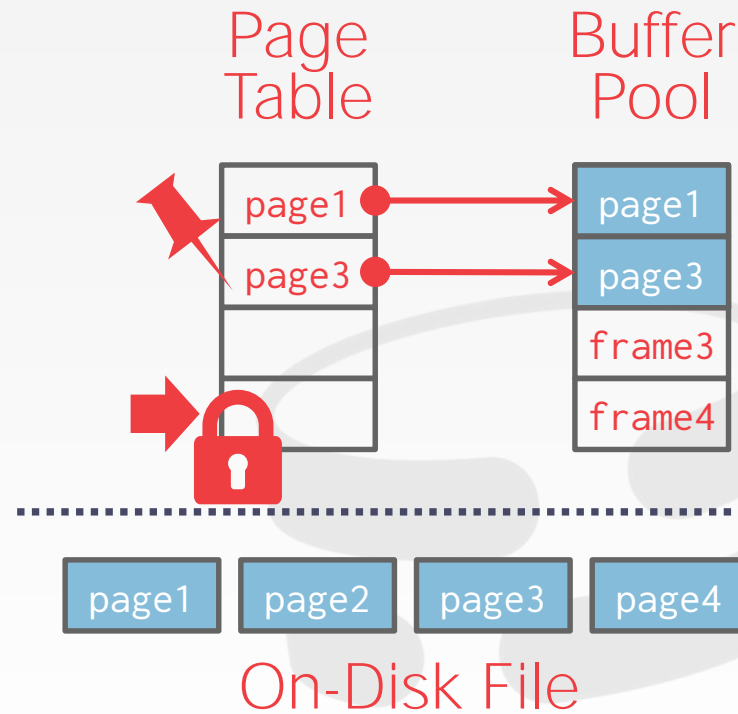
# BUFFER POOL META-DATA

The page table keeps track of pages that are currently in memory.

Also maintains additional meta-data per page:

→ **Dirty Flag**

→ **Pin/Reference Counter**

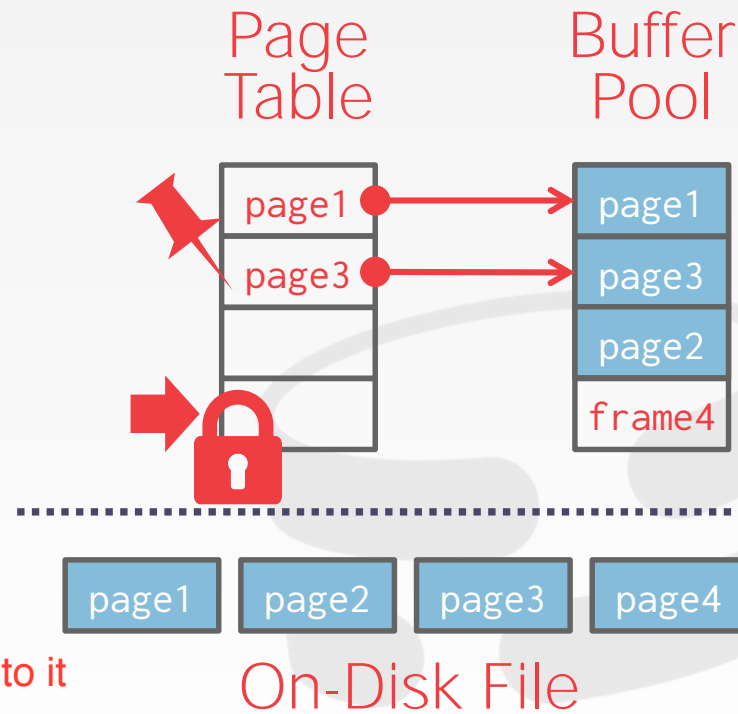


# BUFFER POOL META-DATA

The page table keeps track of pages that are currently in memory.

Also maintains additional meta-data per page:

- **Dirty Flag**
- **Pin/Reference Counter**



Latch the location in the hashtable and then do changes to it

# BUFFER POOL META-DATA

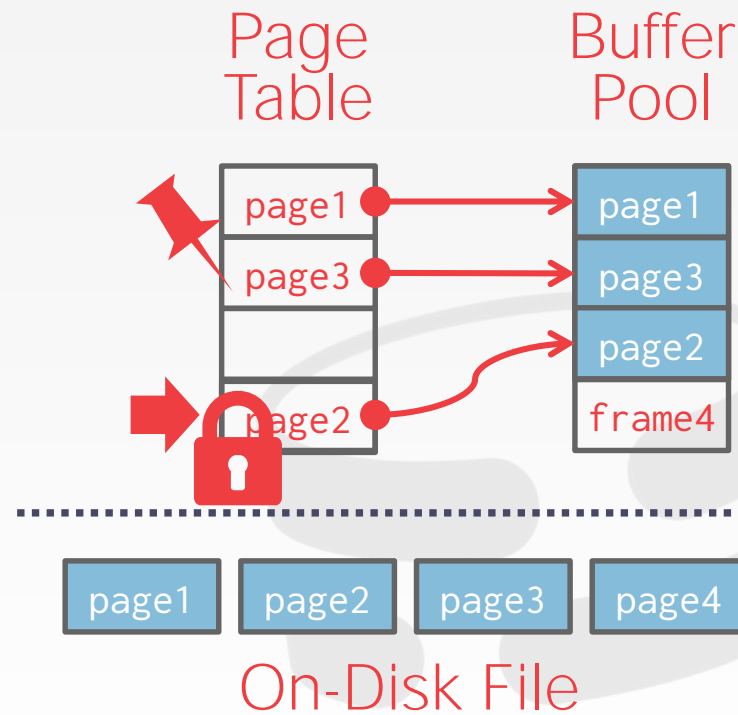
The page table keeps track of pages that are currently in memory.

Also maintains additional meta-data per page:

- **Dirty Flag**
- **Pin/Reference Counter**

read / write: pin it in the table

write: maintain dirty flag to  
let changes written back to the disk



# LOCKS VS. LATCHES

## Locks:

database or tuple

- Protects the database's **logical contents** from other transactions.
- Held for transaction duration.
- Need to be able to rollback changes.

## Latches:

IO

- Protects the critical sections of the DBMS's internal data structure from other threads.
- Held for operation duration.
- Do not need to be able to rollback changes.

← **Mutex**

# PAGE TABLE VS. PAGE DIRECTORY

page id -> page location on db

The **page directory** is the mapping from page ids to page locations in the database files. **run when restart db** **need to store on non-volatile place**  
→ All changes must be recorded on disk to allow the DBMS to find on restart.

**in memory**  
The **page table** is the mapping from page ids to a copy of the page in buffer pool frames.

→ This is an in-memory data structure that does not need to be stored on disk.

do not need to store

One buffer pool with multi threads accessing would cause bottleneck

## MULTIPLE BUFFER POOLS

The DBMS does not always have a single buffer pool for the entire system.

- Multiple buffer pool instances
- Per-database buffer pool
- Per-page type buffer pool

partition memory to several buffer pools and same data use same latch (simply by hash or mod to map page id to buffer pool)

Helps reduce latch contention and improve locality.

Again, OS does not know anything in a file but bytes, but DBMS can tell if this page contains indexes or data and do specific things on that



# PRE-FETCHING

which OS cannot do

The DBMS can also prefetch pages based on a query plan.

- Sequential Scans
- Index Scans

Buffer Pool



Disk Pages





# PRE-FETCHING

The DBMS can also prefetch pages based on a query plan.

- Sequential Scans
- Index Scans

Buffer Pool



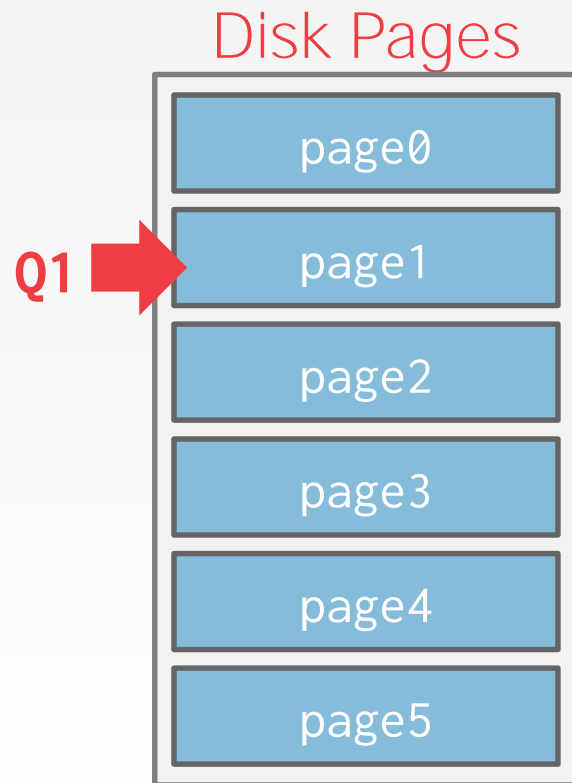
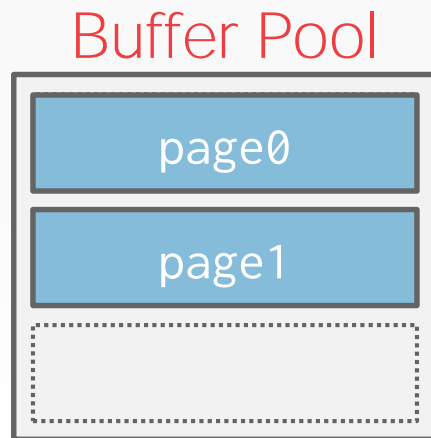
Disk Pages



# PRE-FETCHING

The DBMS can also prefetch pages based on a query plan.

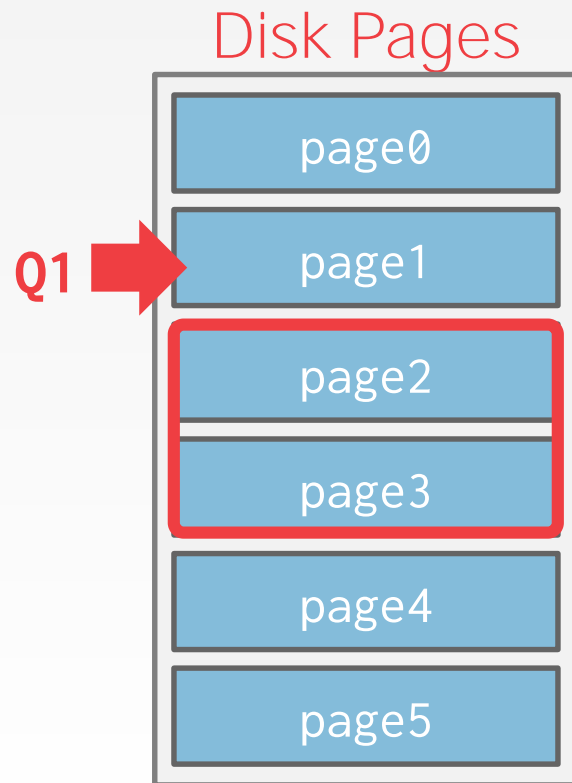
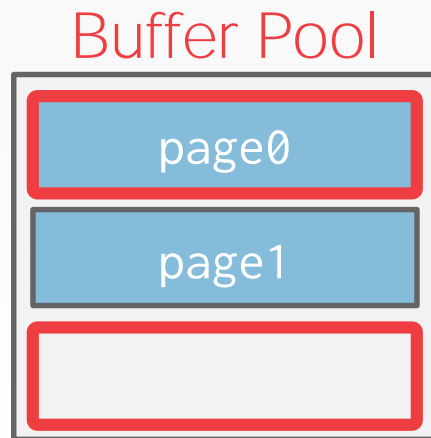
- Sequential Scans
- Index Scans



# PRE-FETCHING

The DBMS can also prefetch pages based on a query plan.

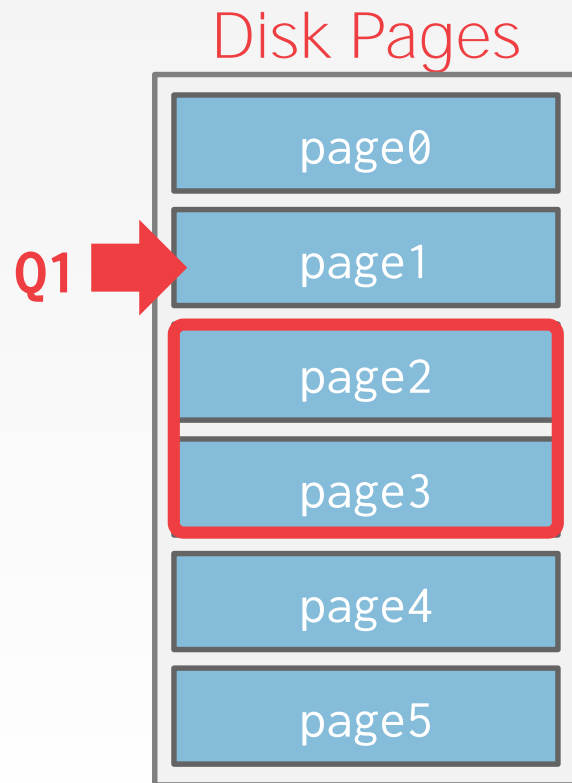
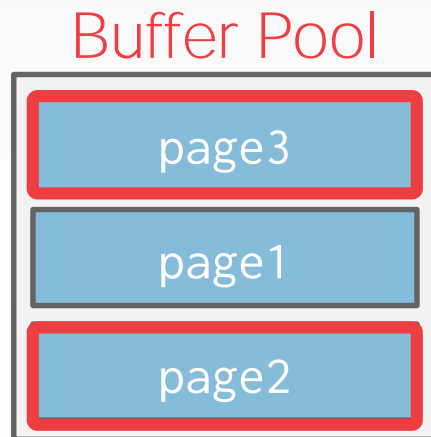
- Sequential Scans
- Index Scans



# PRE-FETCHING

The DBMS can also prefetch pages based on a query plan.

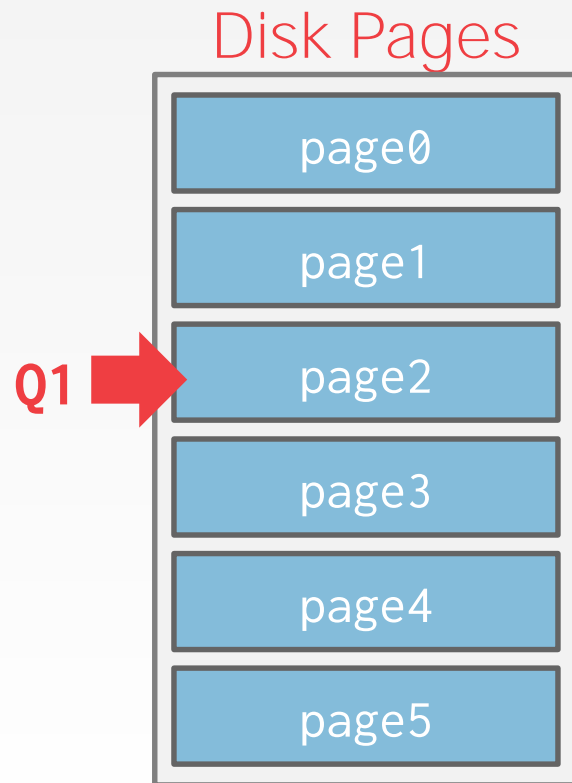
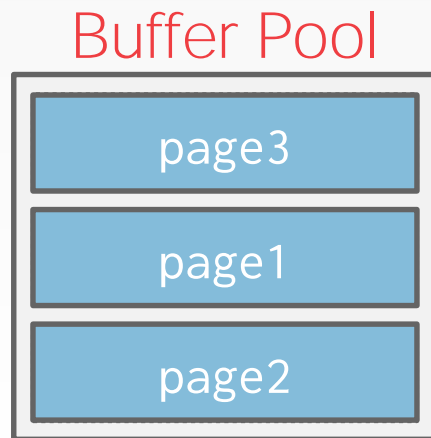
- Sequential Scans
- Index Scans



# PRE-FETCHING

The DBMS can also prefetch pages based on a query plan.

- Sequential Scans
- Index Scans



# PRE-FETCHING

The DBMS can also prefetch pages based on a query plan.

- Sequential Scans
- Index Scans

## Buffer Pool

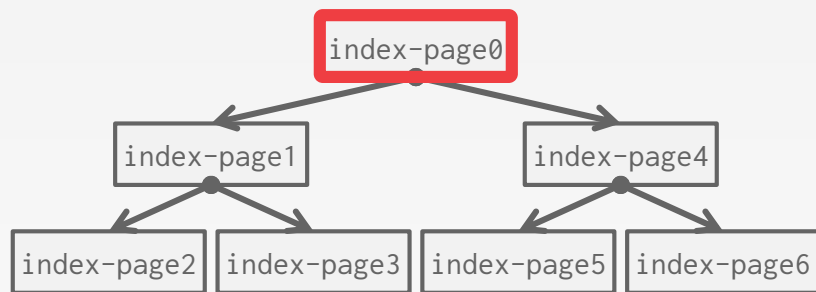


## Disk Pages



because not all scans are sequential scan, indexes scan is more common

## PRE-FETCHING



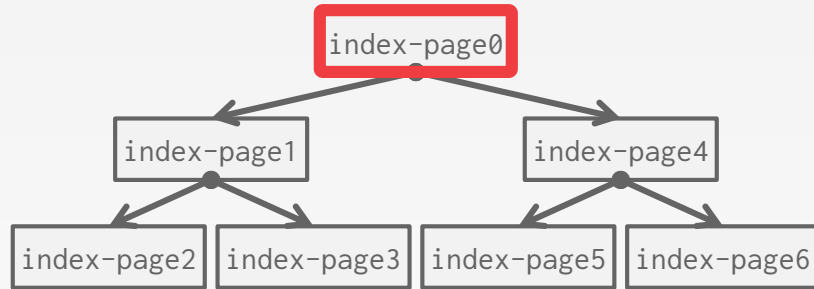
### Buffer Pool



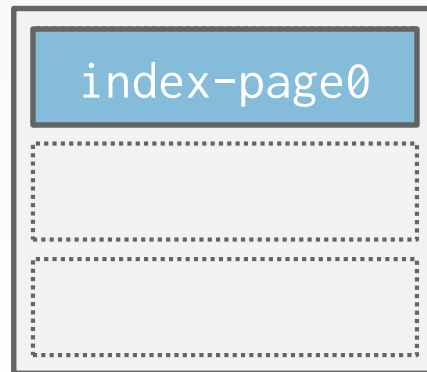
### Disk Pages



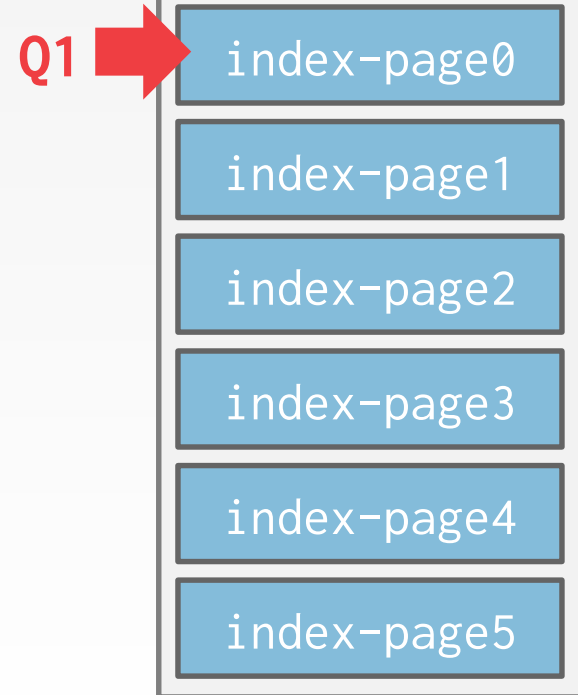
# PRE-FETCHING



## Buffer Pool

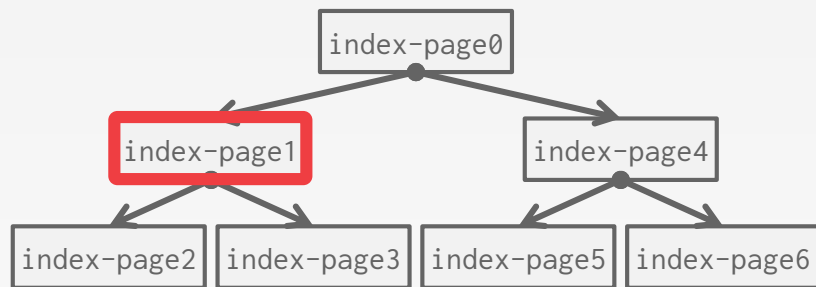


## Disk Pages

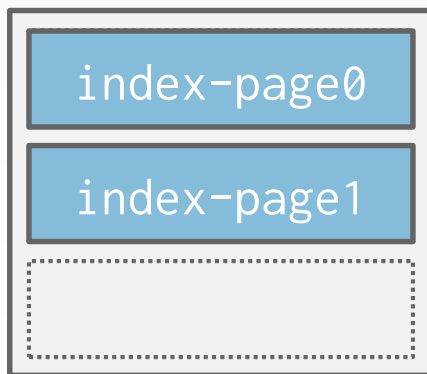




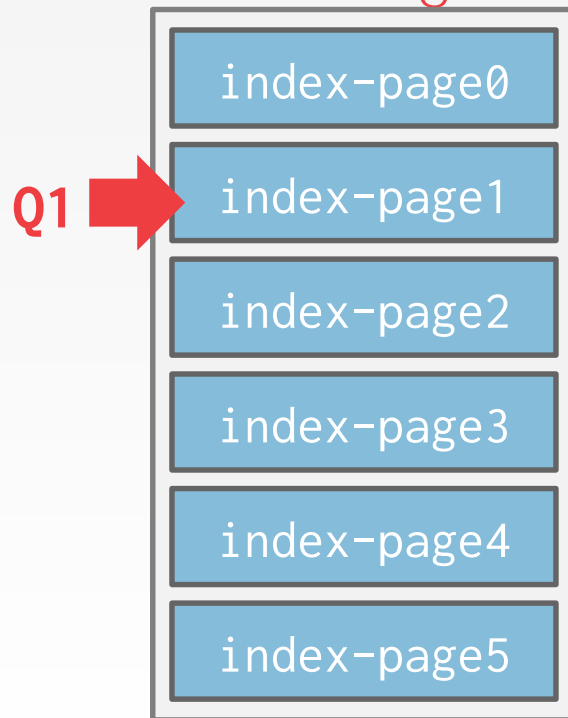
# PRE-FETCHING



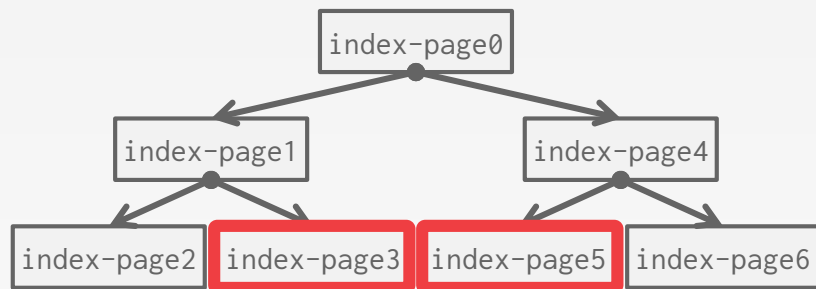
## Buffer Pool



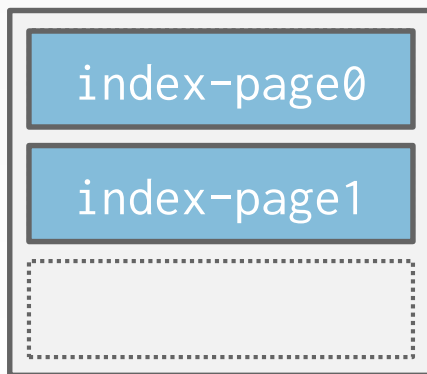
## Disk Pages



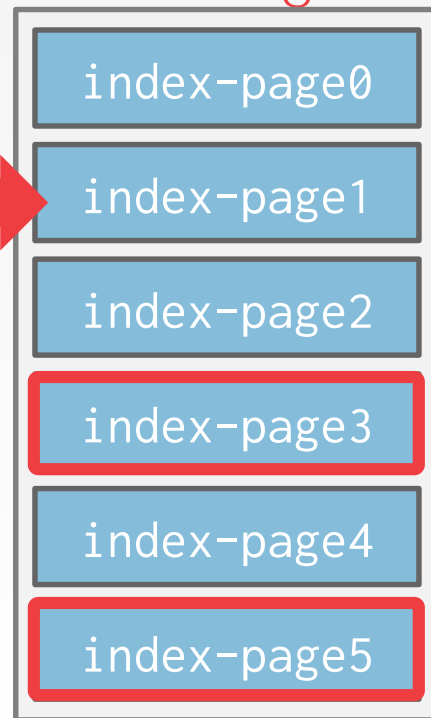
# PRE-FETCHING



## Buffer Pool



## Disk Pages



# SCAN SHARING

---

Queries are able to **reuse data** retrieved from storage or operator computations.

→ **This is different from result caching.**

Allow multiple queries to attach to a single cursor that scans a table.

→ Queries do not have to be exactly the same.

→ Can also share intermediate results.

# SCAN SHARING

---

If a query starts a scan and if there one already doing this, then the DBMS will attach to the second query's cursor.

→ The DBMS keeps track of where the second query joined with the first so that it can finish the scan when it reaches the end of the data structure.

Fully supported in IBM DB2 and MSSQL.

Oracle only supports cursor sharing for identical queries.



image each query has a cursor to traverse pages

## SCAN SHARING

**Q1** `SELECT SUM(val) FROM A`

Buffer Pool



Disk Pages



# SCAN SHARING

**Q1** `SELECT SUM(val) FROM A`

Buffer Pool



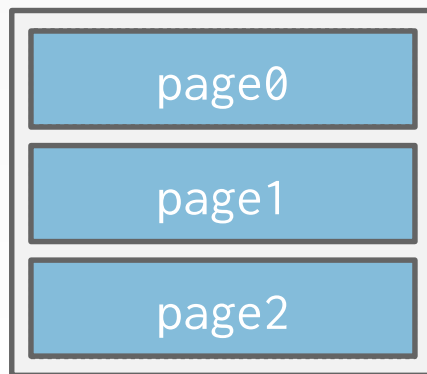
Disk Pages



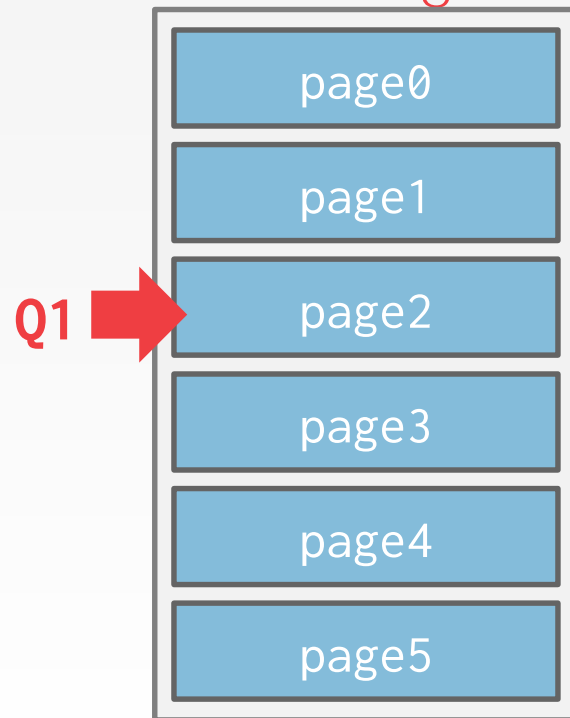
# SCAN SHARING

**Q1** `SELECT SUM(val) FROM A`

Buffer Pool



Disk Pages



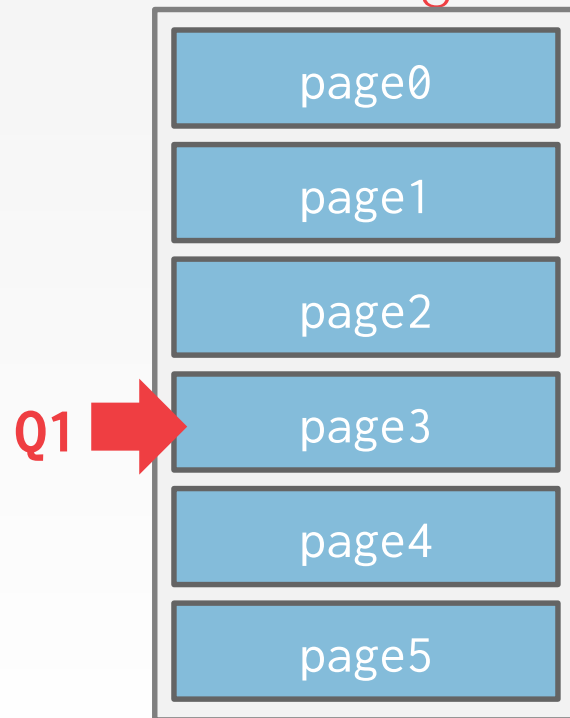
# SCAN SHARING

**Q1** `SELECT SUM(val) FROM A`

Buffer Pool



Disk Pages





# SCAN SHARING

**Q1** `SELECT SUM(val) FROM A`

Buffer Pool



Disk Pages



At this moment, page 0 is the last victim  
page moved out from buffer pool

## SCAN SHARING

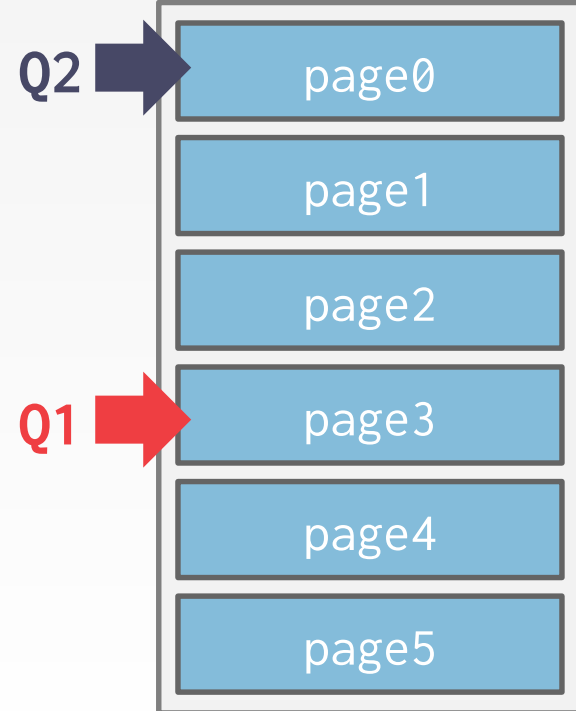
**Q1** `SELECT SUM(val) FROM A`

**Q2** `SELECT AVG(val) FROM A`

### Buffer Pool



### Disk Pages



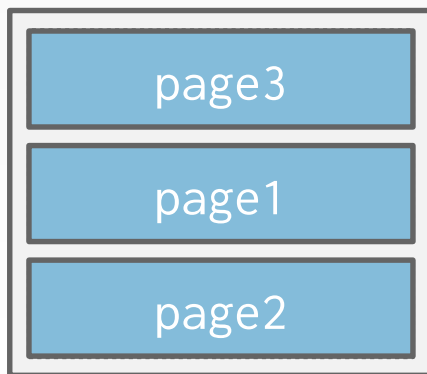
so to avoid extra disk reads, DBMS can move q2 and q1 cursor together,  
to make them read same data pages. After q1 finishes, move pages q2 needs  
to buffer pool

## SCAN SHARING

**Q1** `SELECT SUM(val) FROM A`

**Q2** `SELECT AVG(val) FROM A`

### Buffer Pool



**Q2 Q1** →

### Disk Pages



# SCAN SHARING

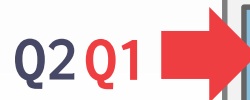
**Q1** `SELECT SUM(val) FROM A`

**Q2** `SELECT AVG(val) FROM A`

## Buffer Pool



## Disk Pages



# SCAN SHARING

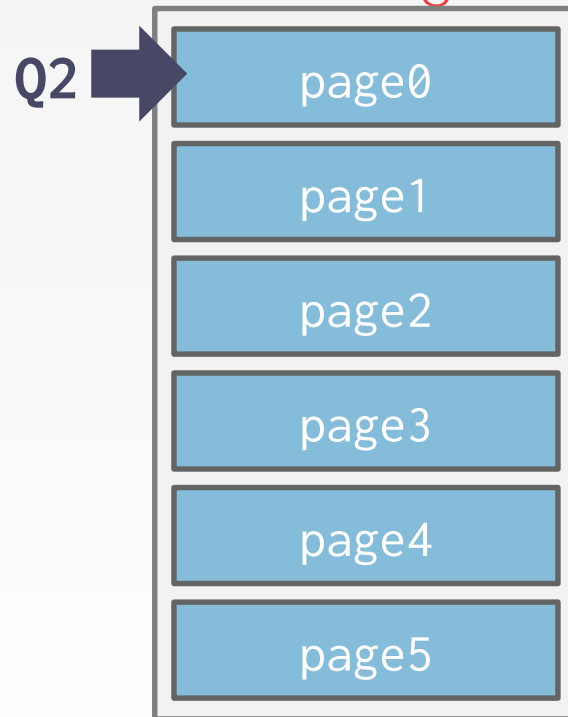
**Q1** `SELECT SUM(val) FROM A`

**Q2** `SELECT AVG(val) FROM A`

## Buffer Pool



## Disk Pages

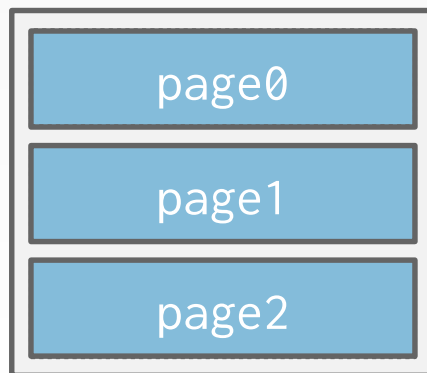


# SCAN SHARING

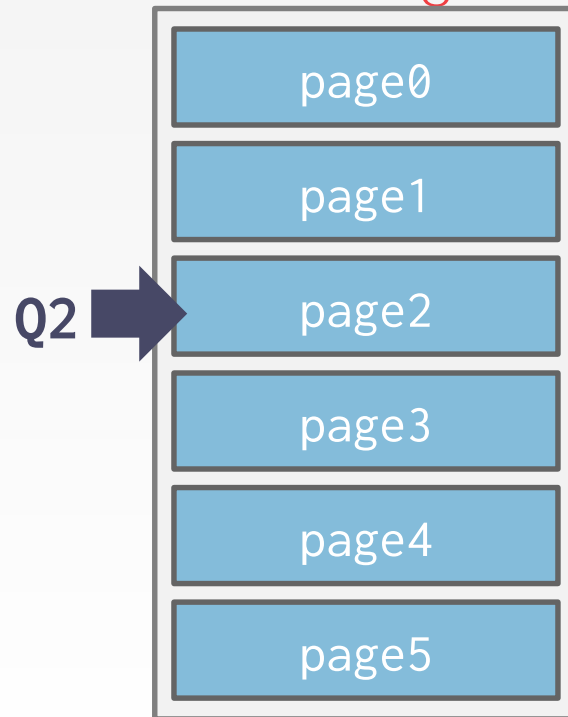
**Q1** `SELECT SUM(val) FROM A`

**Q2** `SELECT AVG(val) FROM A`

## Buffer Pool



## Disk Pages



# BUFFER POOL BYPASS

---

The sequential scan operator will not store fetched pages in the buffer pool to avoid overhead.

- Memory is local to running query.
- Works well if operator needs to read a large sequence of pages that are contiguous on disk.

Called "Light Scans" in Informix.

The Informix logo, featuring the word "Informix" in a bold, sans-serif font. The "i" is red and the "x" is blue with a stylized, multi-colored tail. A registered trademark symbol (®) is to the right of the "x". The logo is overlaid on a large, faint, light gray background image of a database cylinder.

# OS PAGE CACHE

---

Most disk operations go through the OS API.

Unless you tell it not to, the OS maintains its own filesystem cache.

Most DBMSs use direct I/O (O\_DIRECT) to bypass the OS's cache.

- Redundant copies of pages.
- Different eviction policies.



# BUFFER REPLACEMENT POLICIES

---

When the DBMS needs to free up a frame to make room for a new page, it must decide which page to evict from the buffer pool.

Goals:

- Correctness
- Accuracy
- Speed
- Meta-data overhead



# LEAST-RECENTLY USED

---

Maintain a timestamp of when each page was last accessed.

When the DBMS needs to evict a page, select the one with the oldest timestamp.

→ Keep the pages in sorted order to reduce the search time on eviction.

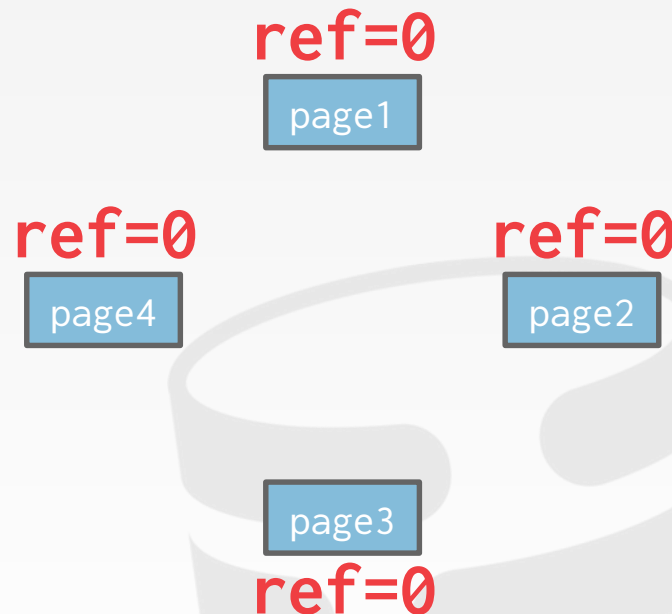
# CLOCK

Approximation of LRU without needing a separate timestamp per page.

- Each page has a reference bit.
- When a page is accessed, set to 1.

Organize the pages in a circular buffer with a "clock hand":

- Upon sweeping, check if a page's bit is set to 1.
- If yes, set to zero. If no, then evict.



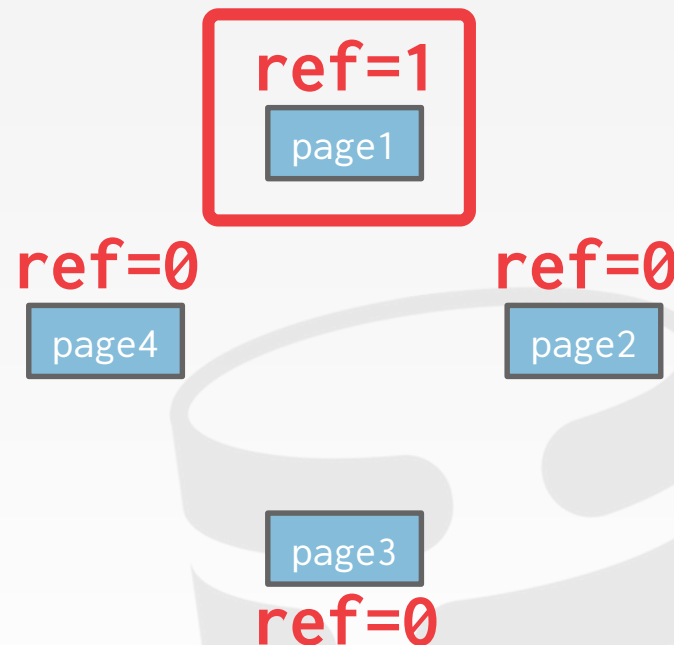
# CLOCK

Approximation of LRU without needing a separate timestamp per page.

- Each page has a reference bit.
- When a page is accessed, set to 1.

Organize the pages in a circular buffer with a "clock hand":

- Upon sweeping, check if a page's bit is set to 1.
- If yes, set to zero. If no, then evict.



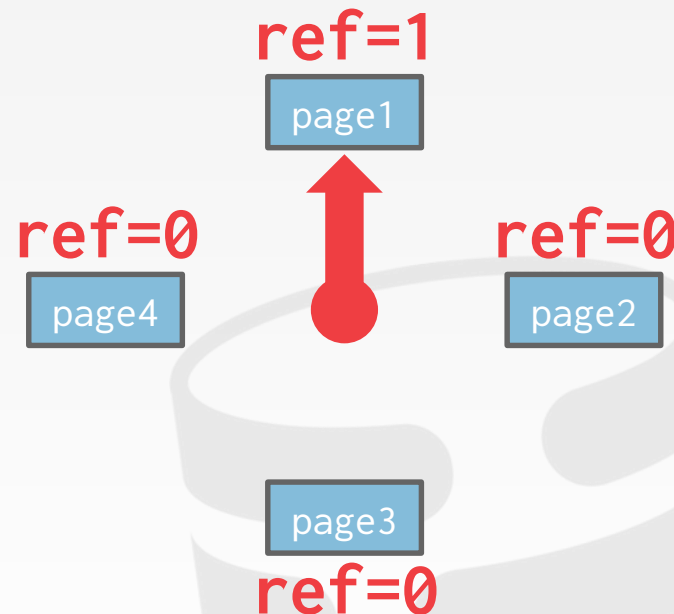
# CLOCK

Approximation of LRU without needing a separate timestamp per page.

- Each page has a reference bit.
- When a page is accessed, set to 1.

Organize the pages in a circular buffer with a "clock hand":

- Upon sweeping, check if a page's bit is set to 1.
- If yes, set to zero. If no, then evict.



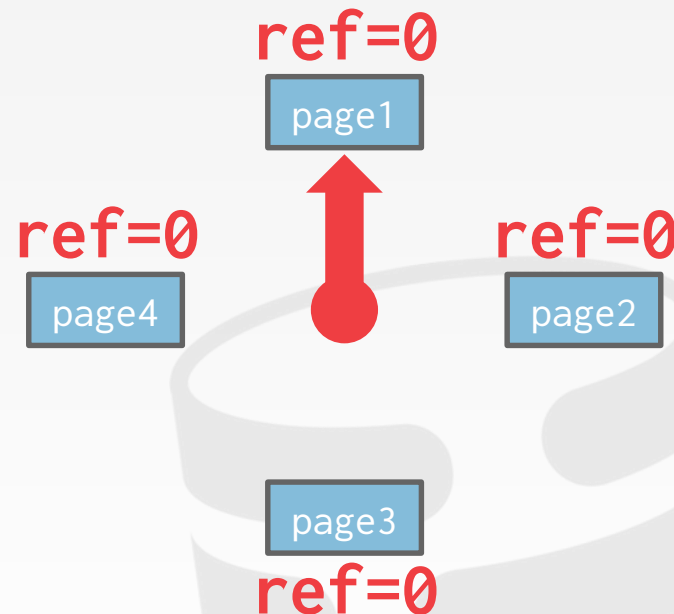
# CLOCK

Approximation of LRU without needing a separate timestamp per page.

- Each page has a reference bit.
- When a page is accessed, set to 1.

Organize the pages in a circular buffer with a "clock hand":

- Upon sweeping, check if a page's bit is set to 1.
- If yes, set to zero. If no, then evict.



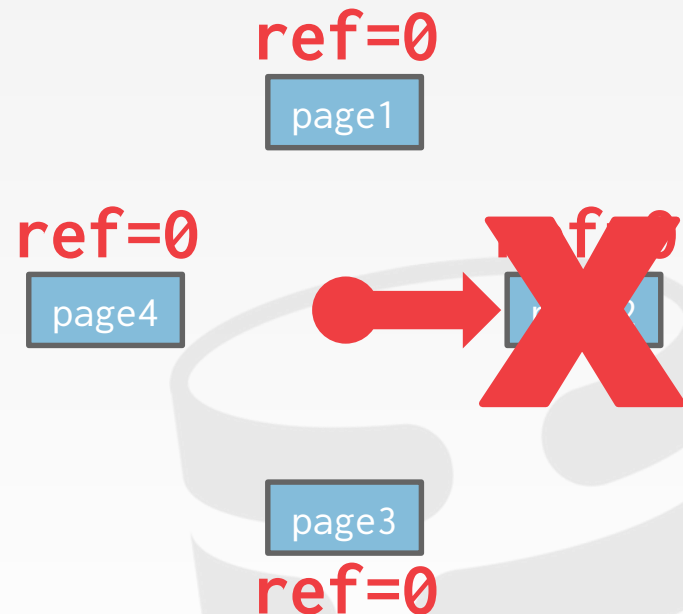
# CLOCK

Approximation of LRU without needing a separate timestamp per page.

- Each page has a reference bit.
- When a page is accessed, set to 1.

Organize the pages in a circular buffer with a "clock hand":

- Upon sweeping, check if a page's bit is set to 1.
- If yes, set to zero. If no, then evict.



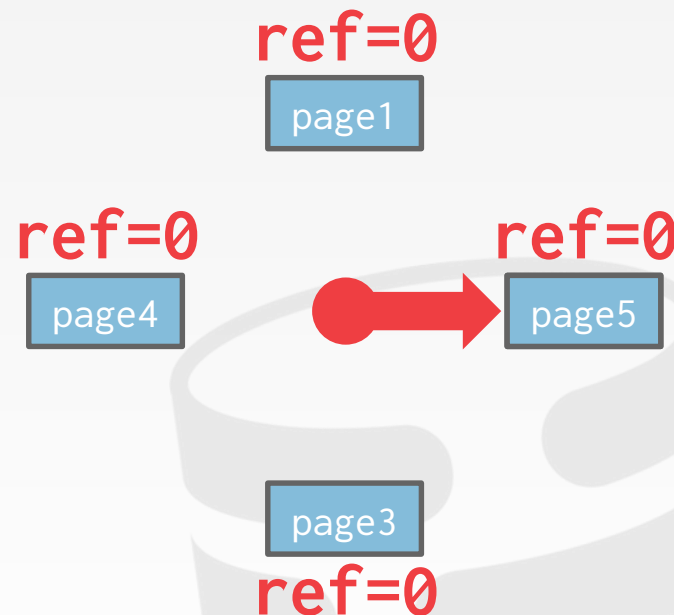
# CLOCK

Approximation of LRU without needing a separate timestamp per page.

- Each page has a reference bit.
- When a page is accessed, set to 1.

Organize the pages in a circular buffer with a "clock hand":

- Upon sweeping, check if a page's bit is set to 1.
- If yes, set to zero. If no, then evict.





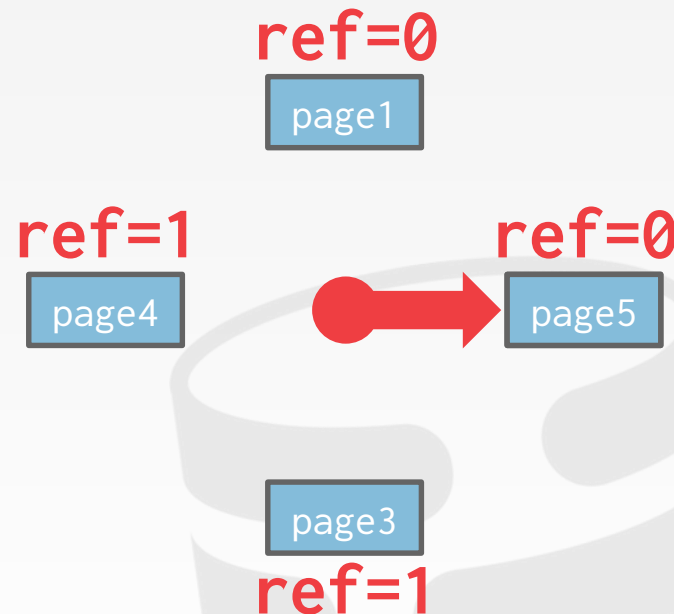
# CLOCK

Approximation of LRU without needing a separate timestamp per page.

- Each page has a reference bit.
- When a page is accessed, set to 1.

Organize the pages in a circular buffer with a "clock hand":

- Upon sweeping, check if a page's bit is set to 1.
- If yes, set to zero. If no, then evict.



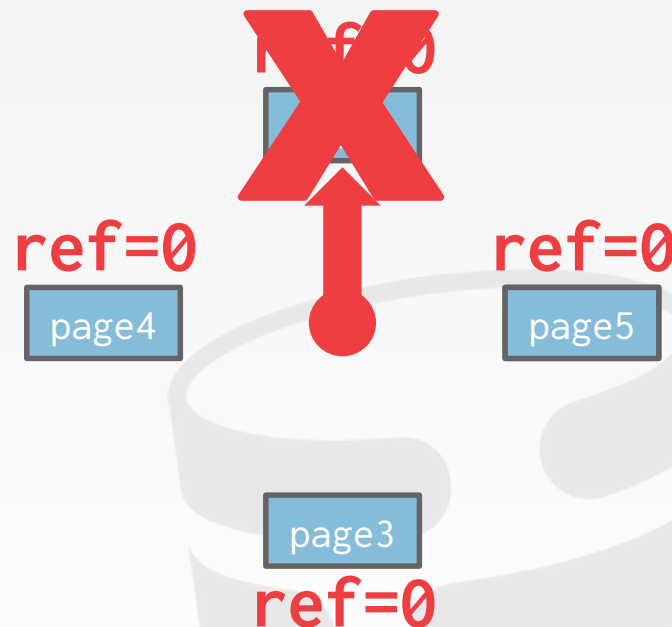
# CLOCK

Approximation of LRU without needing a separate timestamp per page.

- Each page has a reference bit.
- When a page is accessed, set to 1.

Organize the pages in a circular buffer with a "clock hand":

- Upon sweeping, check if a page's bit is set to 1.
- If yes, set to zero. If no, then evict.



# PROBLEMS

---

LRU and CLOCK replacement policies are susceptible to sequential flooding.

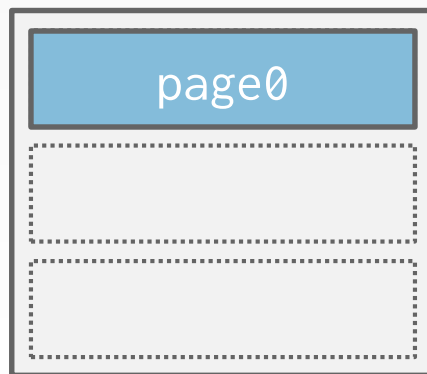
- A query performs a sequential scan that reads every page.
- This pollutes the buffer pool with pages that are read once and then never again.

The most recently used page is actually the most unneeded page.

# SEQUENTIAL FLOODING

**Q1** `SELECT * FROM A WHERE id = 1`

Buffer Pool



Disk Pages



# SEQUENTIAL FLOODING

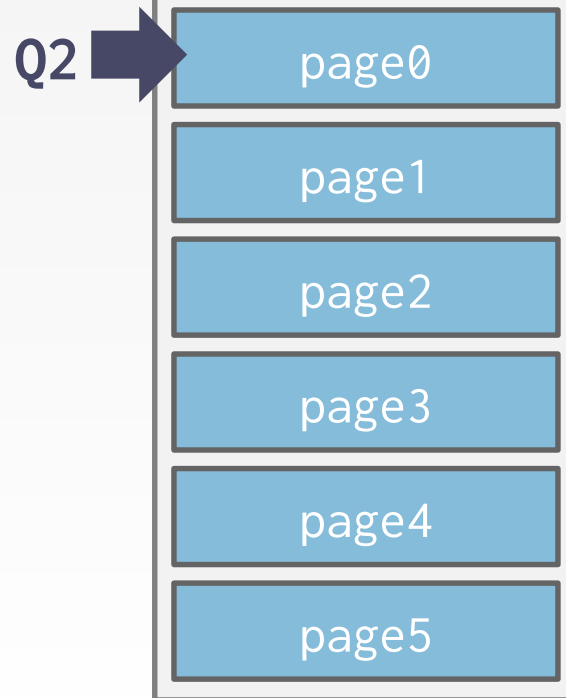
**Q1** `SELECT * FROM A WHERE id = 1`

**Q2** `SELECT AVG(val) FROM A`

Buffer Pool



Disk Pages

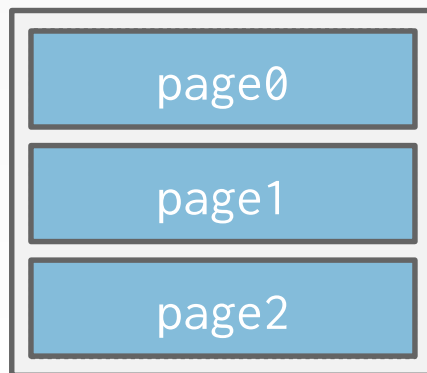


# SEQUENTIAL FLOODING

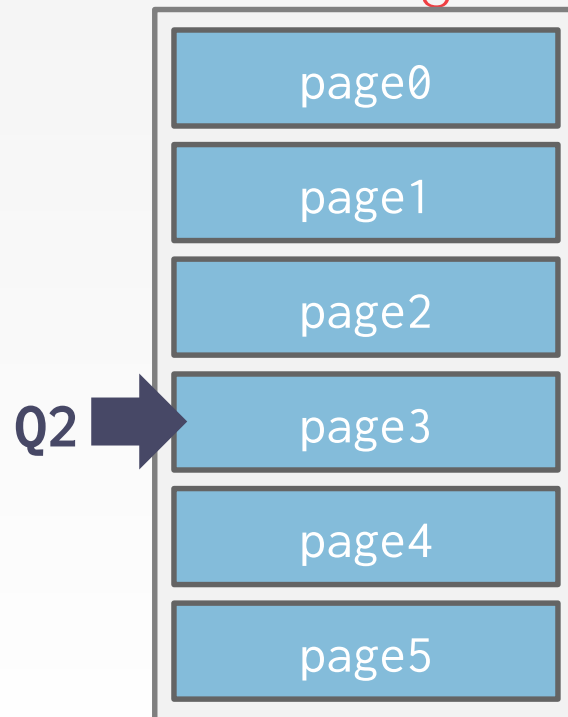
**Q1** `SELECT * FROM A WHERE id = 1`

**Q2** `SELECT AVG(val) FROM A`

## Buffer Pool



## Disk Pages



# SEQUENTIAL FLOODING

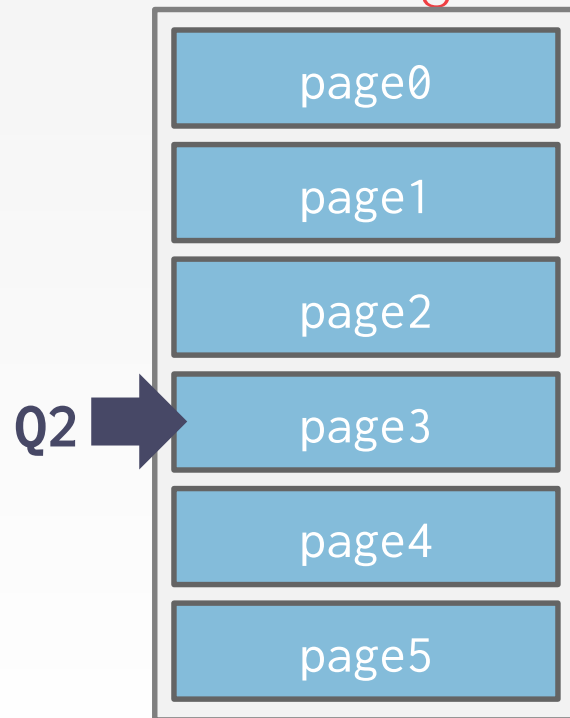
**Q1** `SELECT * FROM A WHERE id = 1`

**Q2** `SELECT AVG(val) FROM A`

## Buffer Pool



## Disk Pages



# SEQUENTIAL FLOODING

**Q1** `SELECT * FROM A WHERE id = 1`

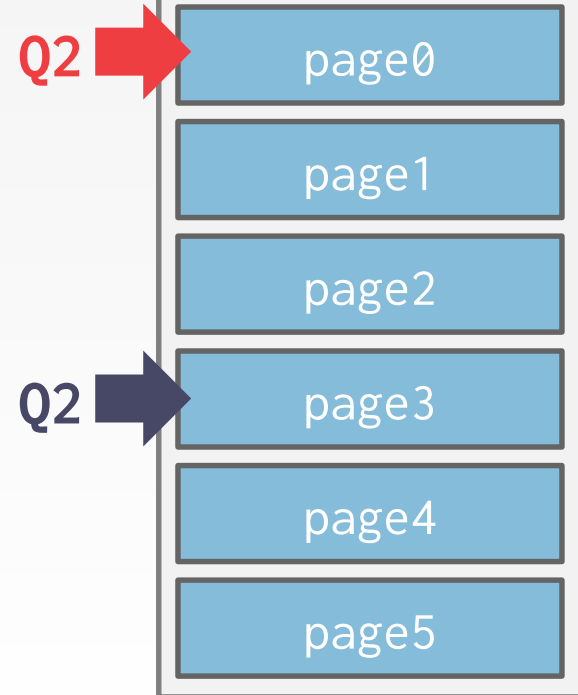
**Q2** `SELECT AVG(val) FROM A`

**Q3** `SELECT * FROM A WHERE id = 1`

## Buffer Pool



## Disk Pages





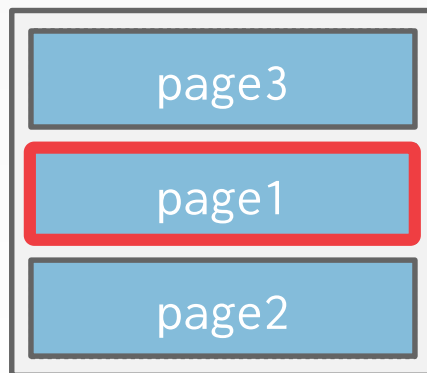
# SEQUENTIAL FLOODING

**Q1** `SELECT * FROM A WHERE id = 1`

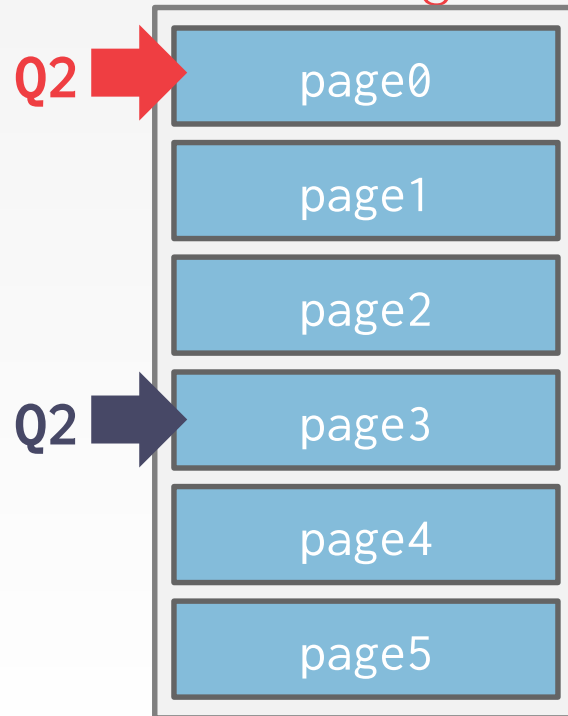
**Q2** `SELECT AVG(val) FROM A`

**Q3** `SELECT * FROM A WHERE id = 1`

## Buffer Pool



## Disk Pages



## BETTER POLICIES: LRU-K

---

Take into account history of the last  $K$  references as timestamps and compute the interval between subsequent accesses.

The DBMS then uses this history to estimate the next time that page is going to be accessed.

# BETTER POLICIES: LOCALIZATION

---

The DBMS chooses which pages to evict on a per txn/query basis. This minimizes the pollution of the buffer pool from each query.

→ Keep track of the pages that a query has accessed.

Example: Postgres maintains a small ring buffer that is private to the query.

# BETTER POLICIES: PRIORITY HINTS

---

The DBMS knows what the context of each page during query execution.

It can provide hints to the buffer pool on whether a page is important or not.

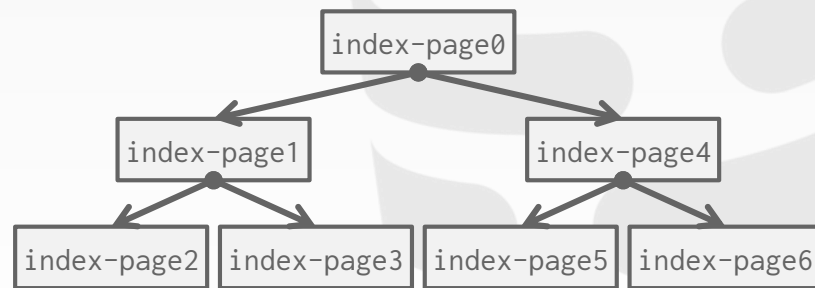


# BETTER POLICIES: PRIORITY HINTS

The DBMS knows what the context of each page during query execution.

It can provide hints to the buffer pool on whether a page is important or not.

**Q1** `INSERT INTO A VALUES (id++)`

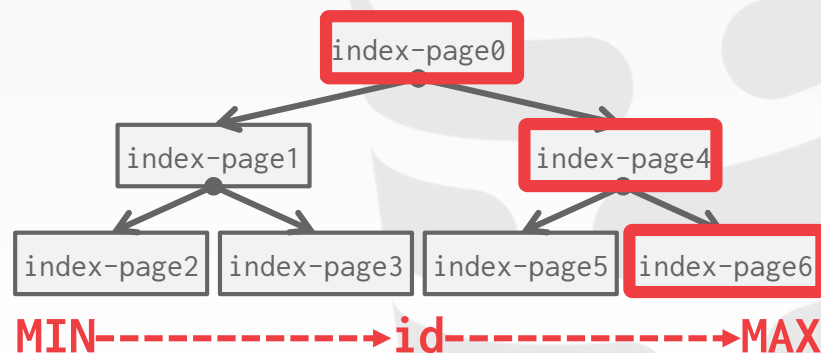


# BETTER POLICIES: PRIORITY HINTS

The DBMS knows what the context of each page during query execution.

It can provide hints to the buffer pool on whether a page is important or not.

**Q1** `INSERT INTO A VALUES (id++)`



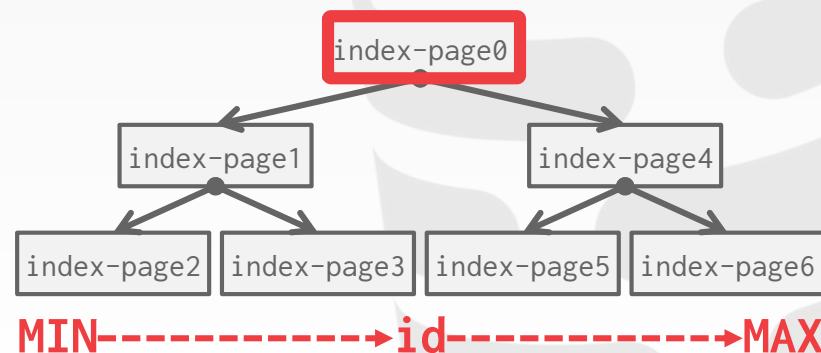
# BETTER POLICIES: PRIORITY HINTS

The DBMS knows what the context of each page during query execution.

It can provide hints to the buffer pool on whether a page is important or not.

**Q1** INSERT INTO A VALUES (*id++*)

**Q2** SELECT \* FROM A WHERE id = ?



# DIRTY PAGES

---

**FAST:** If a page in the buffer pool is not dirty, then the DBMS can simply "drop" it.

**SLOW:** If a page is dirty, then the DBMS must write back to disk to ensure that its changes are persisted.

Trade-off between fast evictions versus dirty writing pages that will not be read again in the future.



# BACKGROUND WRITING

---

The DBMS can periodically walk through the page table and write dirty pages to disk.

When a dirty page is safely written, the DBMS can either evict the page or just unset the dirty flag.

Need to be careful that we don't write dirty pages before their log records have been written...

# ALLOCATION POLICIES

---

## **Global Policies:**

→ Make decisions for all active txns.

## **Local Policies:**

→ Allocate frames to a specific txn without considering the behavior of concurrent txns.

→ Still need to support sharing pages.

# OTHER MEMORY POOLS

---

The DBMS needs memory for things other than just tuples and indexes.

These other memory pools may not always be backed by disk. Depends on implementation.

- Sorting + Join Buffers
- Query Caches
- Maintenance Buffers
- Log Buffers
- Dictionary Caches



# CONCLUSION

---

The DBMS can manage that sweet, sweet memory better than the OS.

Leverage the semantics about the query plan to make better decisions:

- Evictions
- Allocations
- Pre-fetching



# PROJECT #1

---

You will build the first component of your storage manager.

- Extendible Hash Table
- LRU Replacement Policy
- Buffer Pool Manager

All of the projects are based on SQLite, but you will not be able to use your storage manager just yet after this first project.



Due Date:  
Wed Sept 26<sup>th</sup> @ 11:59pm

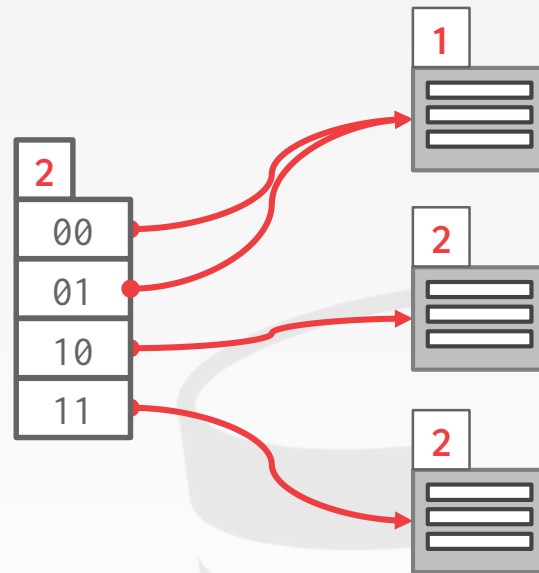
# TASK #1 – EXTENDIBLE HASH TABLE

Build a thread-safe extendible hash table.

- Use unordered buckets to store key/value pairs.
- You must support growing table size.
- You do not need to support shrinking.

General Hints:

- You can use **std::hash** and **std::mutex**.



## TASK #2 – LRU REPLACEMENT POLICY

---

Build a data structure that tracks the usage of **Page** objects in the buffer pool using the least-recently used policy.

General Hints:

→ Your **LRUReplacer** does not need to worry about the "pinned" status of a **Page**.

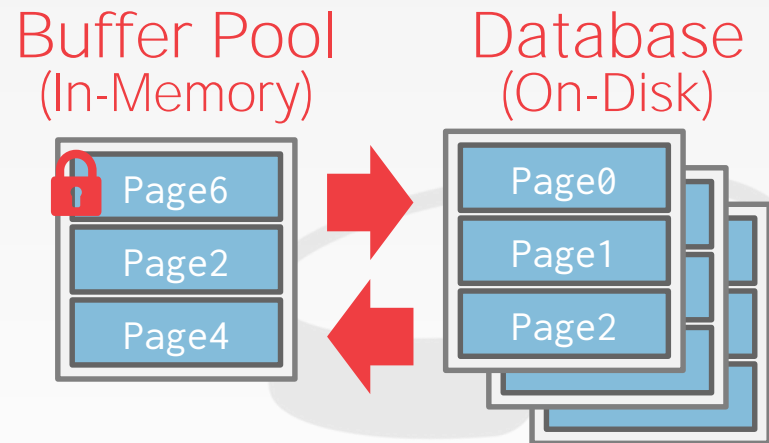
# TASK #3 – BUFFER POOL MANAGER

Combine your hash table and LRU replacer together to manage the allocation of pages.

- Need to maintain an internal data structures of allocated + free pages.
- We will provide you components to read/write data from disk.

## General Hints:

- Make sure you get the order of operations correct when pinning.





# GETTING STARTED

---

Download the source code from the project webpage.

Make sure you can build it on your machine.

- We've test it on Andrew machines, OSX, and Linux.
- It should compile on Windows 10 w/ Ubuntu, but we haven't tried it.

# THINGS TO NOTE

---

Do **not** change any file other than the six that you have to hand in.

The projects are cumulative.

We will **not** be providing solutions.

Post your questions on Piazza or come to our office hours. We will **not** help you debug.

# PLAGIARISM WARNING

---

Your project implementation must be your own work.

- You may not copy source code from other groups or the web.
- Do not publish your implementation on Github.

Plagiarism will not be tolerated.  
See [CMU's Policy on Academic Integrity](#) for additional information.



# NEXT CLASS

---

HASH TABLES!

