



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Ack: Prof. Jignesh Patel @ CMU
Prof. Andy Pavlo @CMU

CSC3170

6: Memory & Disk I/O Management

Chenhao Ma

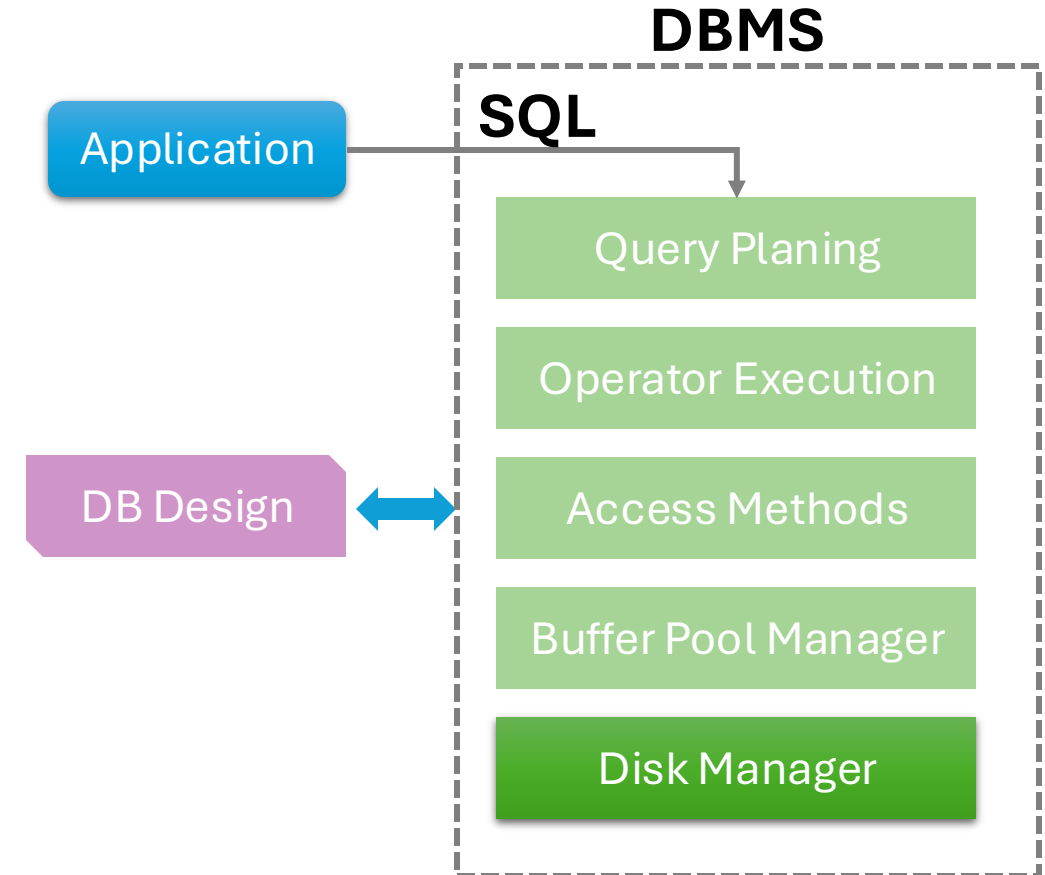
School of Data Science

The Chinese University of Hong Kong, Shenzhen

Last Lecture

- **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

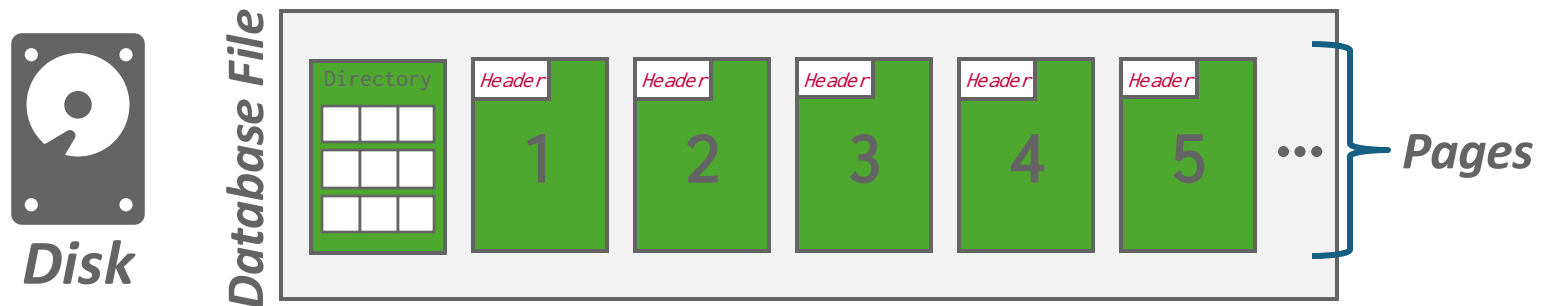
- **Spatial Control:**

- 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.

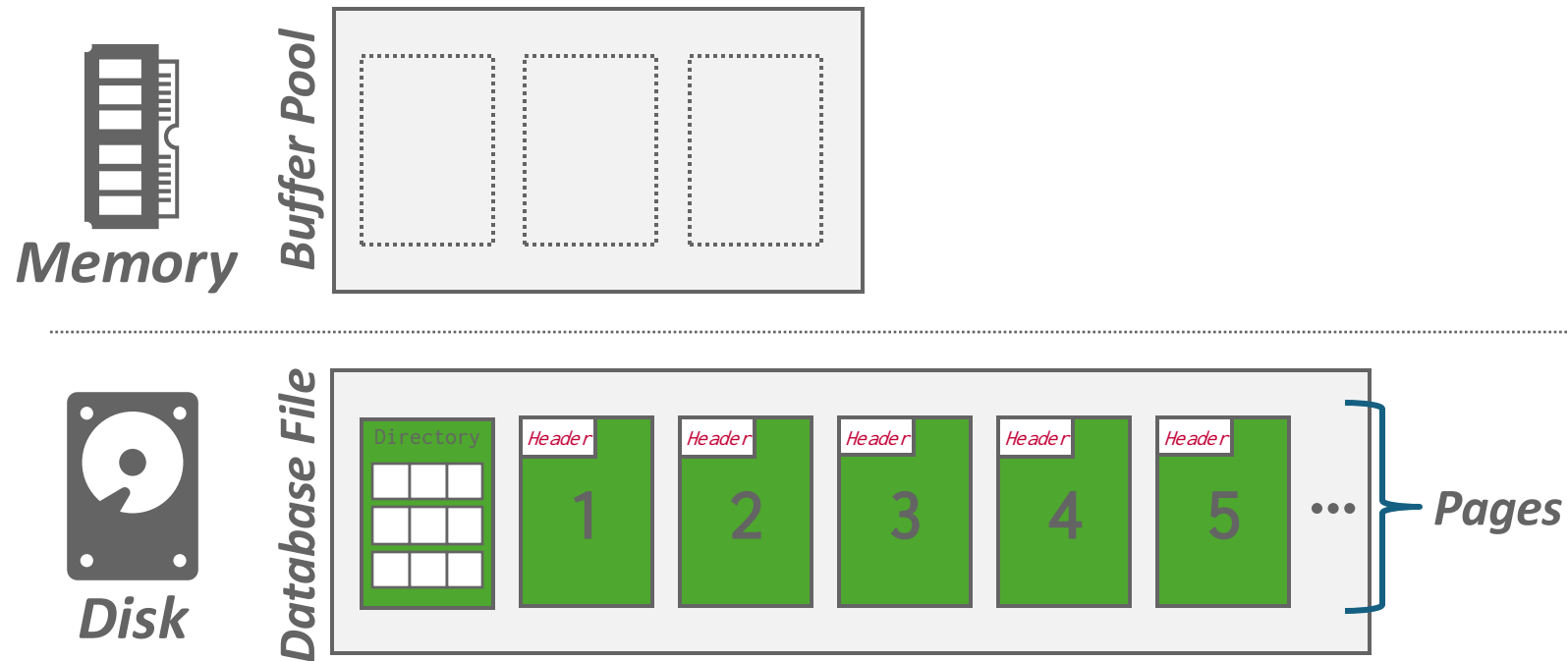
- **Temporal Control:**

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

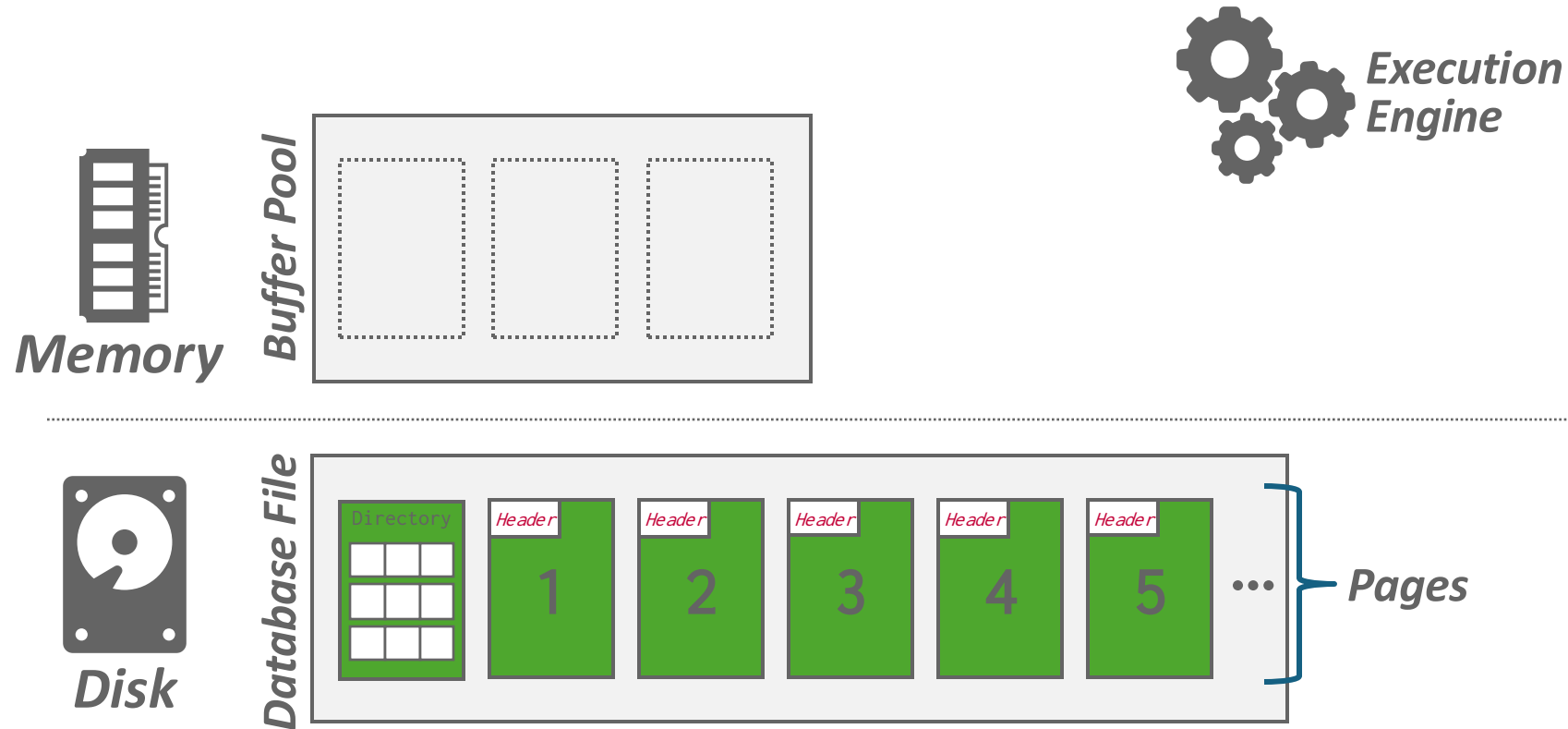
Disk-oriented DBMS



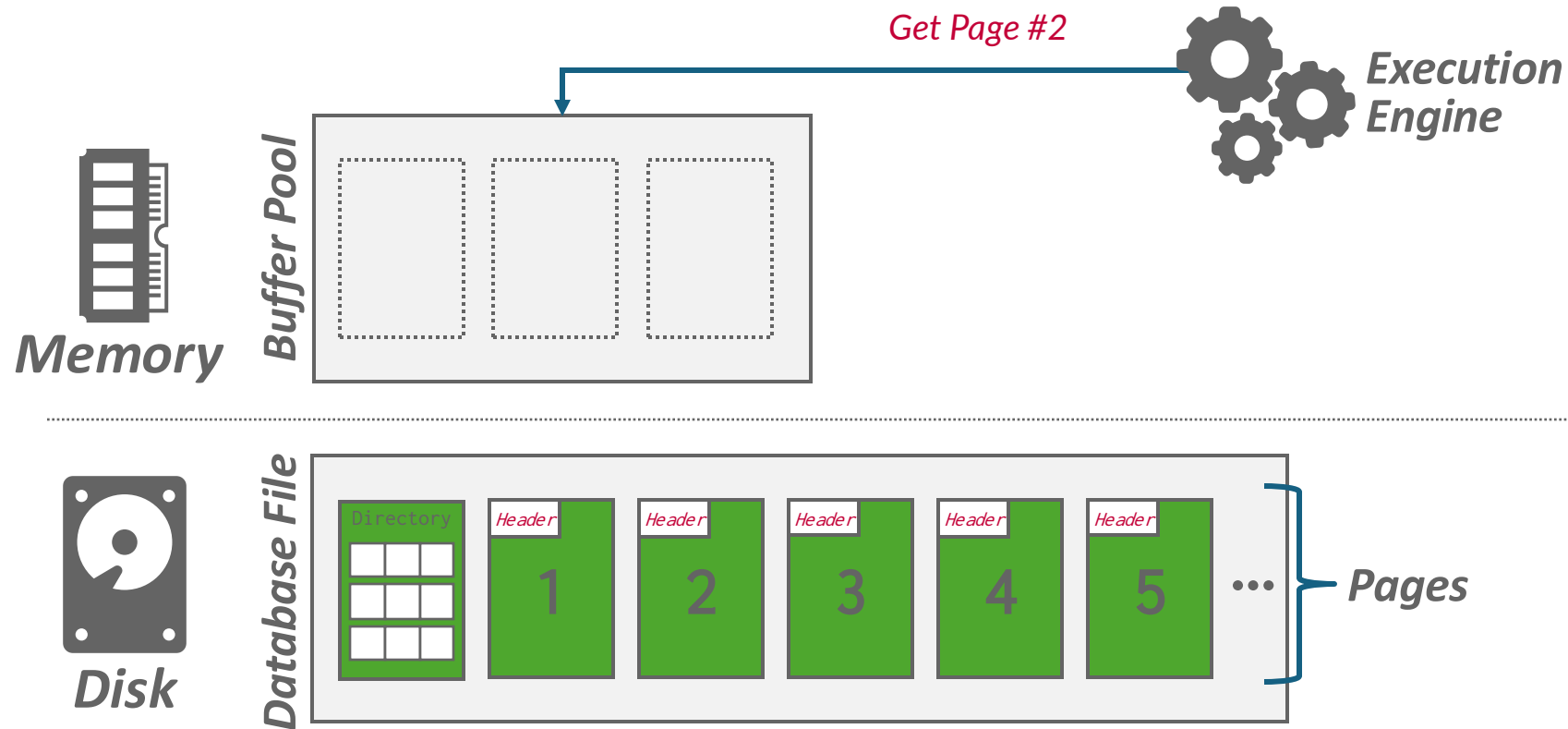
Disk-oriented DBMS



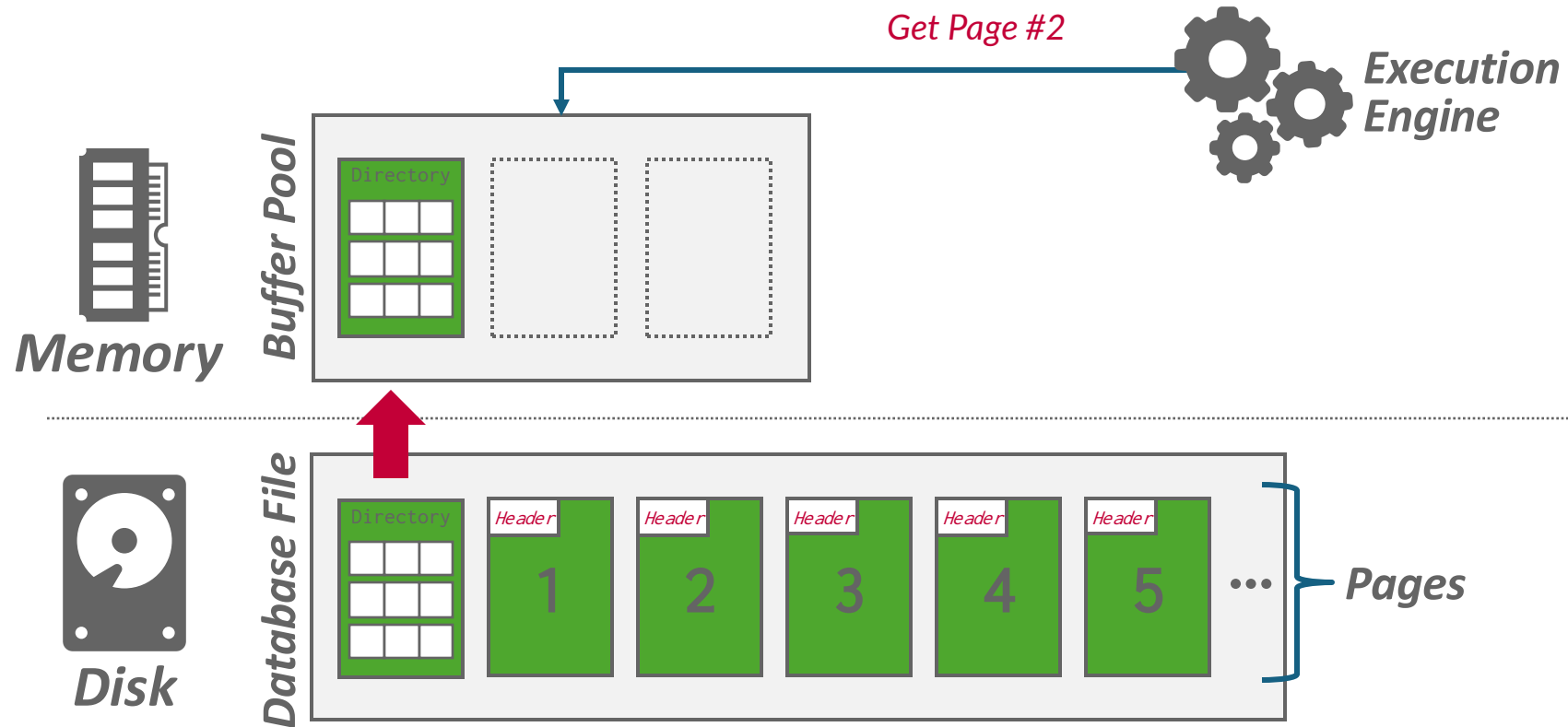
Disk-oriented DBMS



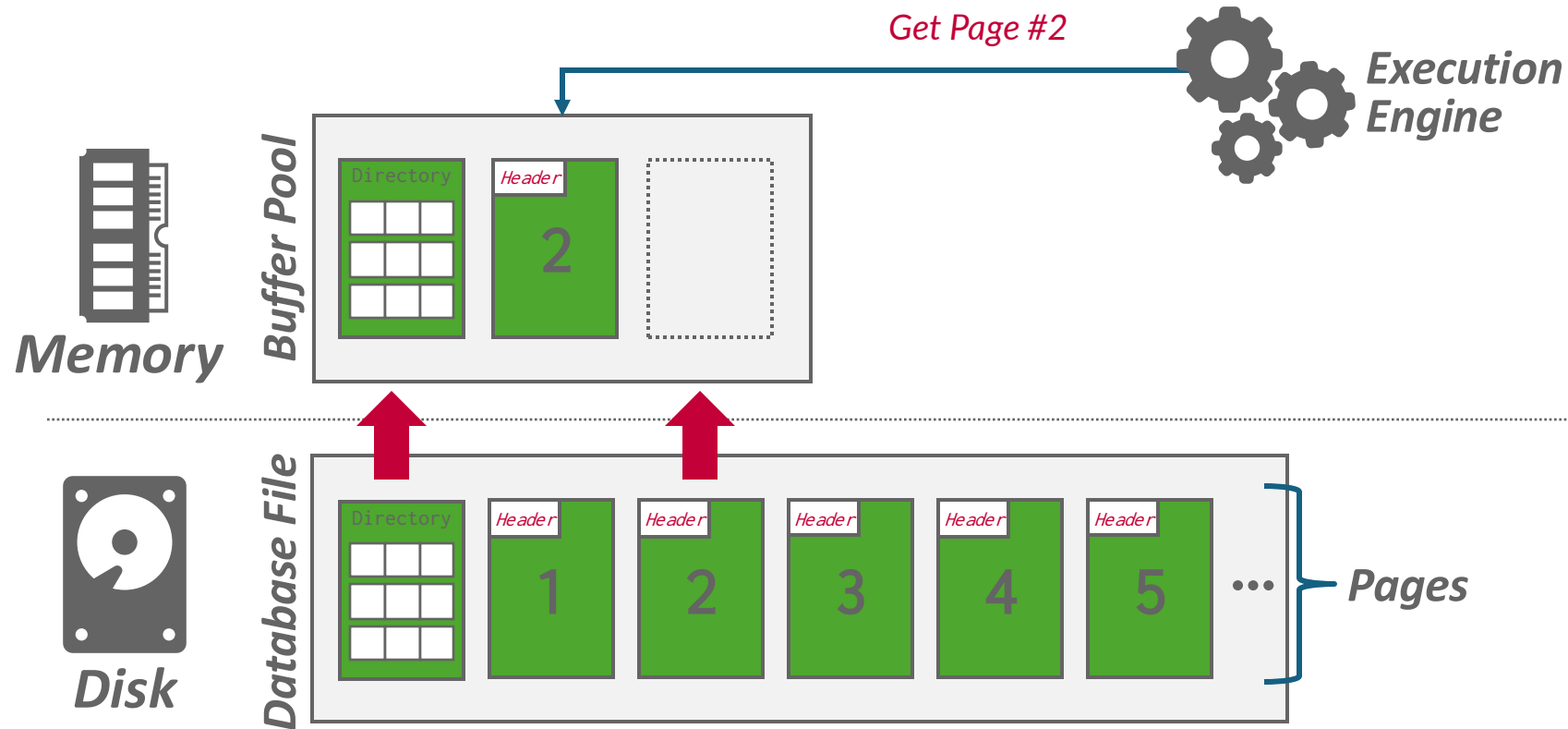
Disk-oriented DBMS



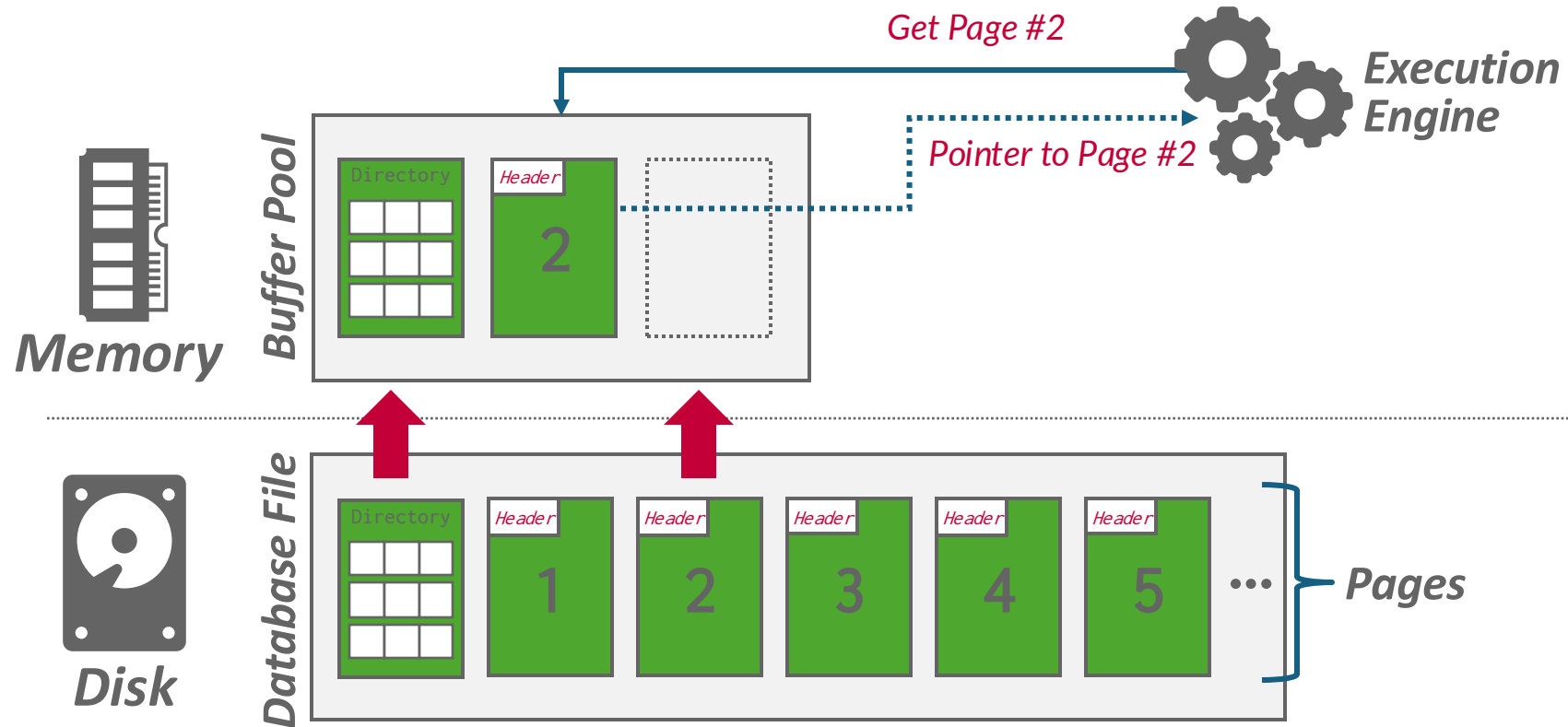
Disk-oriented DBMS



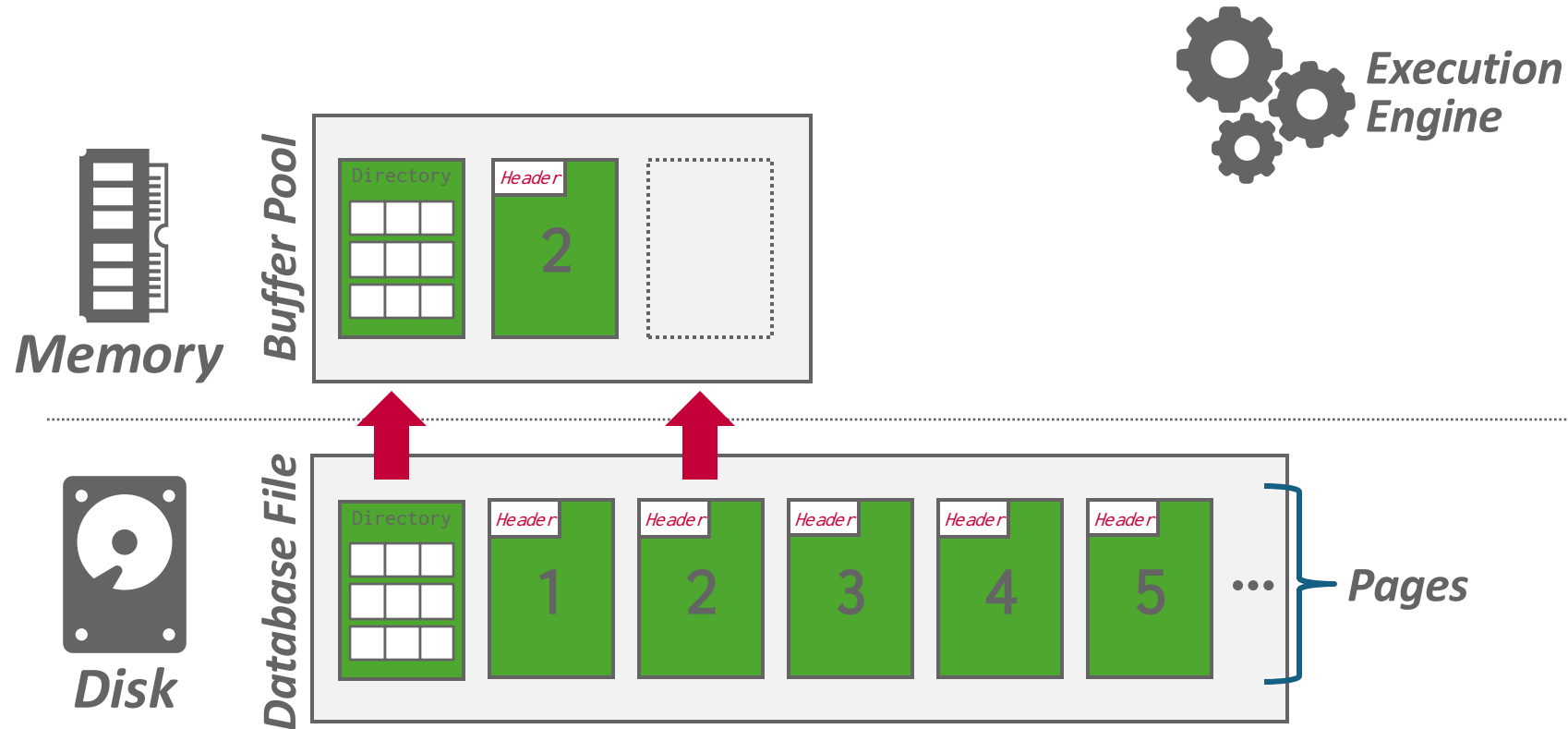
Disk-oriented DBMS



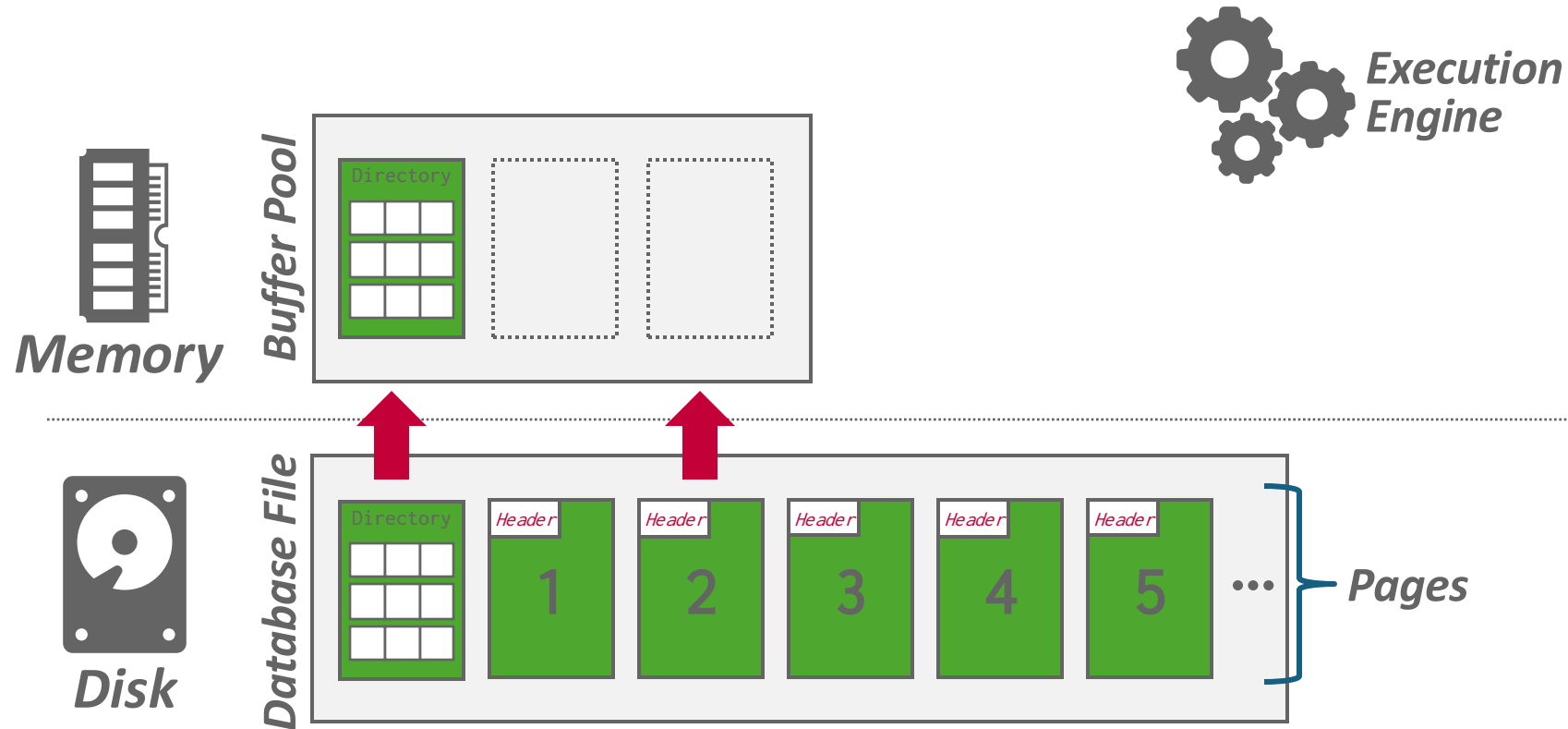
Disk-oriented DBMS



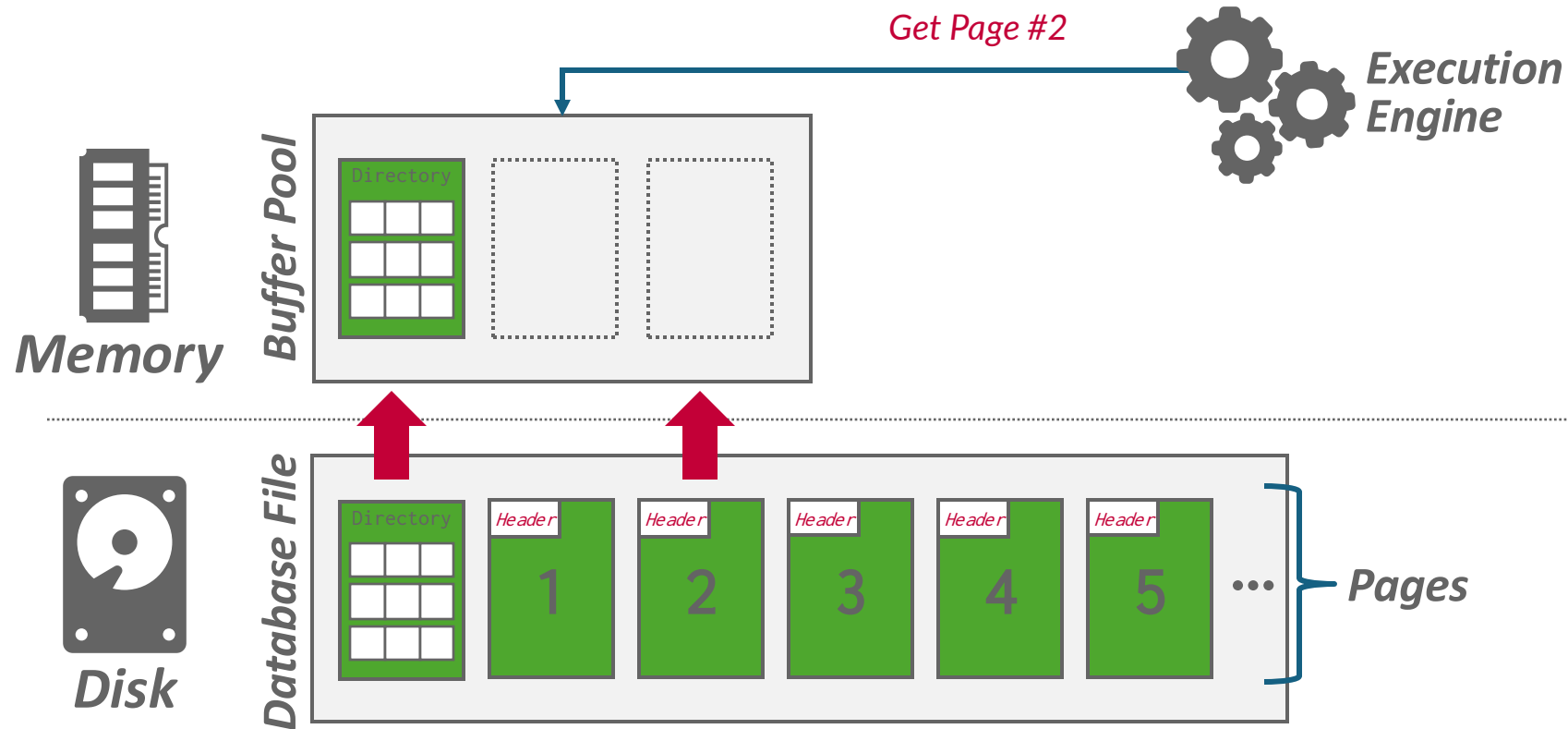
Disk-oriented DBMS



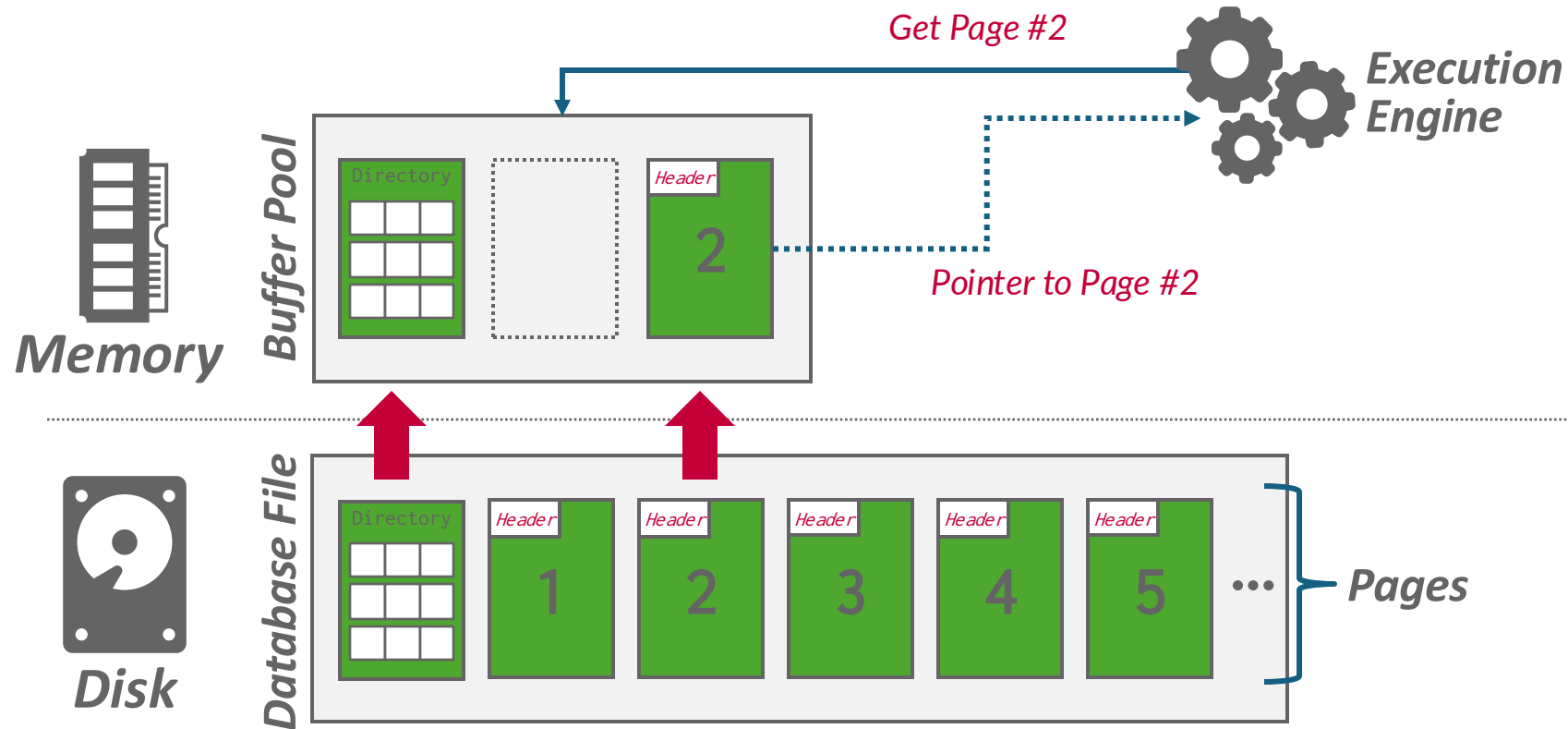
Disk-oriented DBMS



Disk-oriented DBMS



Disk-oriented DBMS



This Lecture

- Buffer Pool Manager
- Disk I/O Scheduling
- Replacement Policies
- Other Memory Pools

Buffer Pool Manager

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.
- Dirty pages are buffered and **not** written to disk immediately
 - Write-Back Cache

*Buffer
Pool*

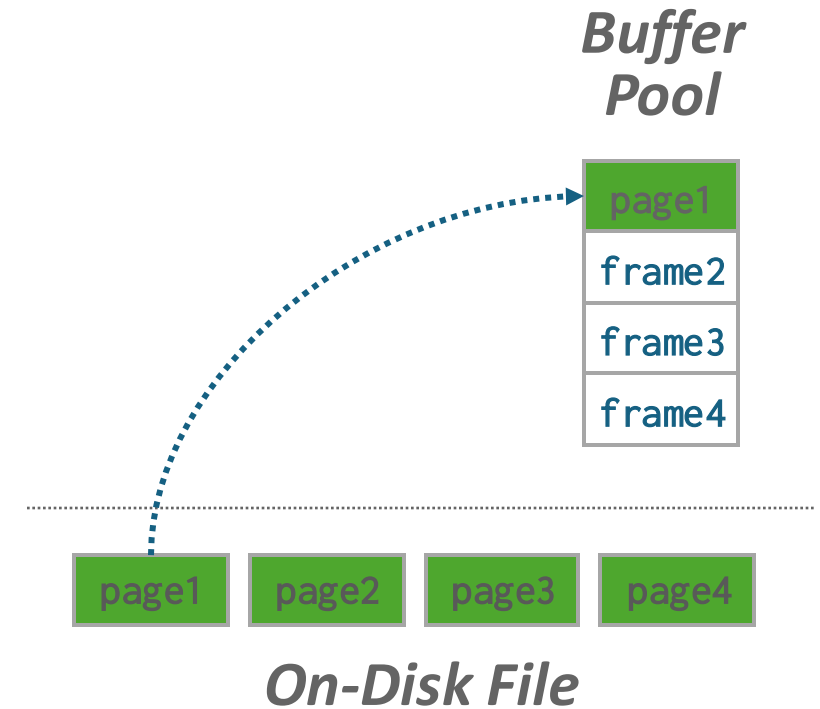
frame1
frame2
frame3
frame4

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

On-Disk File

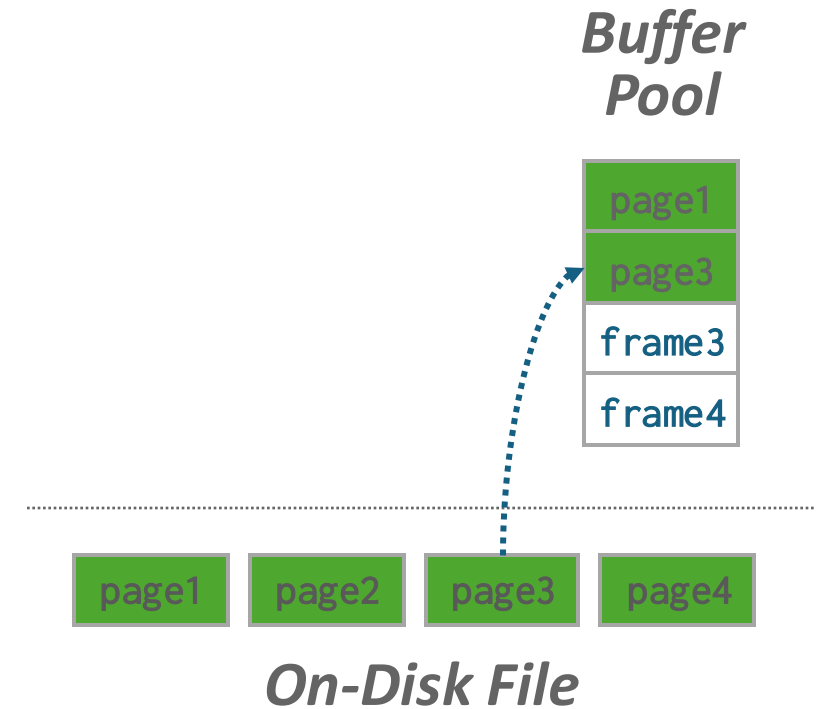
Buffer Pool Organization

- 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.
- Dirty pages are buffered and **not** written to disk immediately
 - Write-Back Cache



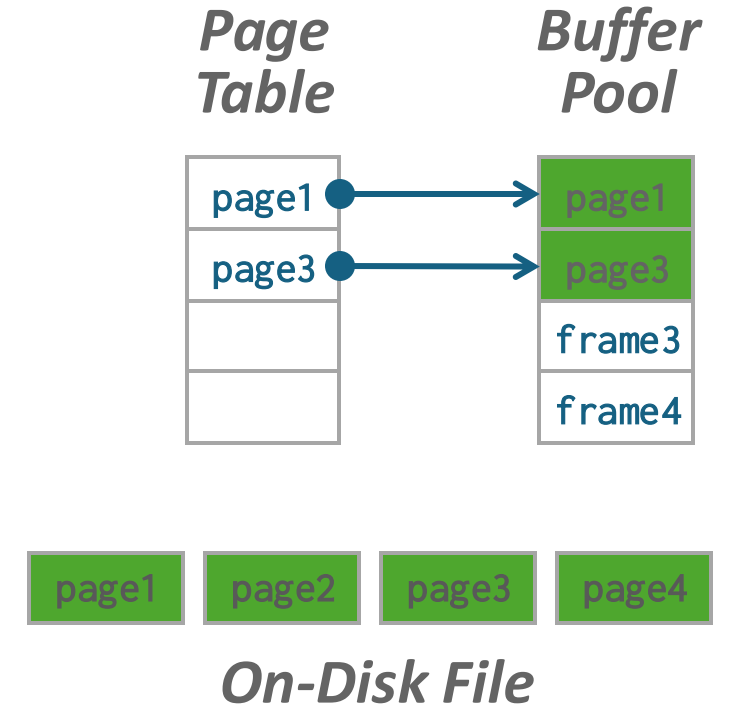
Buffer Pool Organization

- 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.
- Dirty pages are buffered and **not** written to disk immediately
 - Write-Back Cache



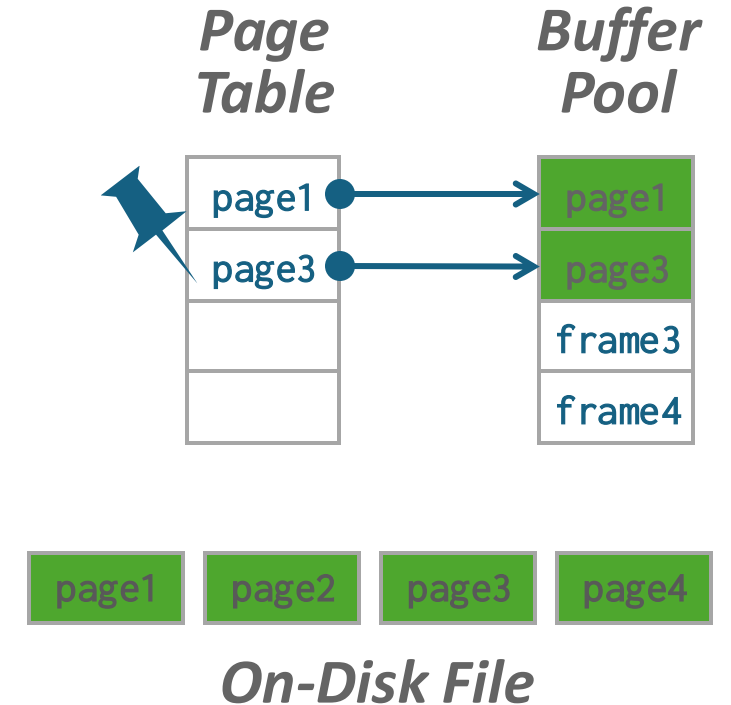
Buffer Pool Meta-data

- The **page table** keeps track of pages that are currently in memory.
 - Usually a fixed-size hash table protected with latches to ensure thread-safe access.
- Additional meta-data per page:
 - **Dirty Flag**
 - **Pin/Reference Counter**
 - **Access Tracking Information**



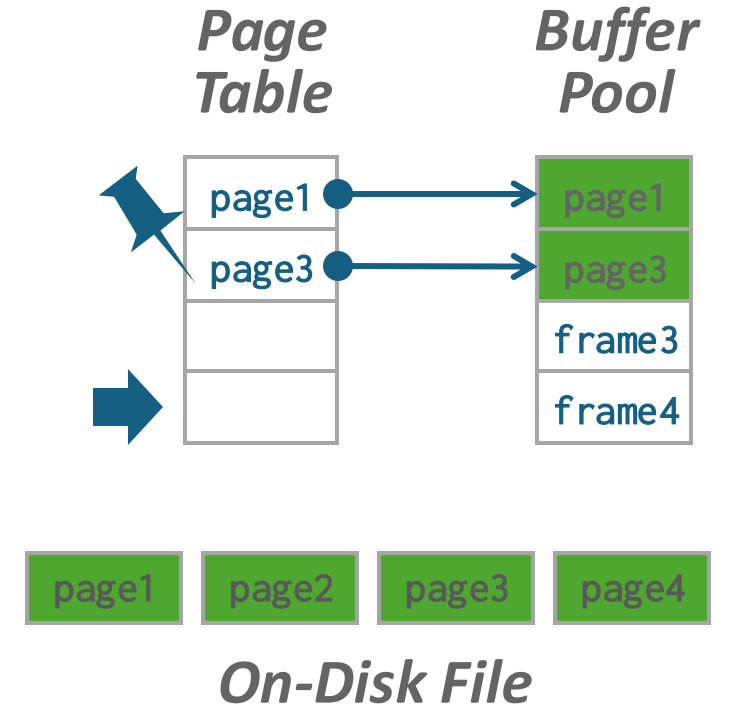
Buffer Pool Meta-data

- The **page table** keeps track of pages that are currently in memory.
 - Usually a fixed-size hash table protected with latches to ensure thread-safe access.
- Additional meta-data per page:
 - **Dirty Flag**
 - **Pin/Reference Counter**
 - **Access Tracking Information**



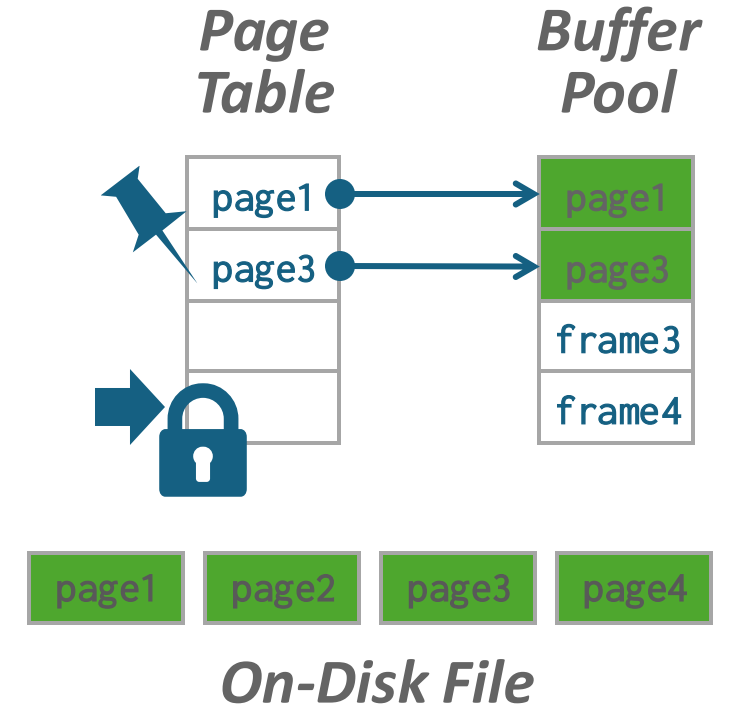
Buffer Pool Meta-data

- The **page table** keeps track of pages that are currently in memory.
 - Usually a fixed-size hash table protected with latches to ensure thread-safe access.
- Additional meta-data per page:
 - **Dirty Flag**
 - **Pin/Reference Counter**
 - **Access Tracking Information**



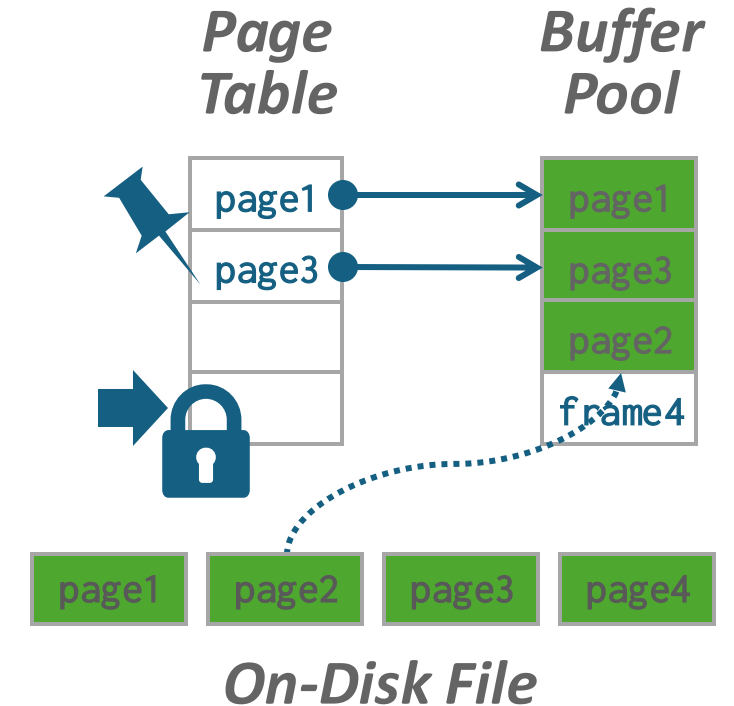
Buffer Pool Meta-data

- The **page table** keeps track of pages that are currently in memory.
 - Usually a fixed-size hash table protected with latches to ensure thread-safe access.
- Additional meta-data per page:
 - **Dirty Flag**
 - **Pin/Reference Counter**
 - **Access Tracking Information**



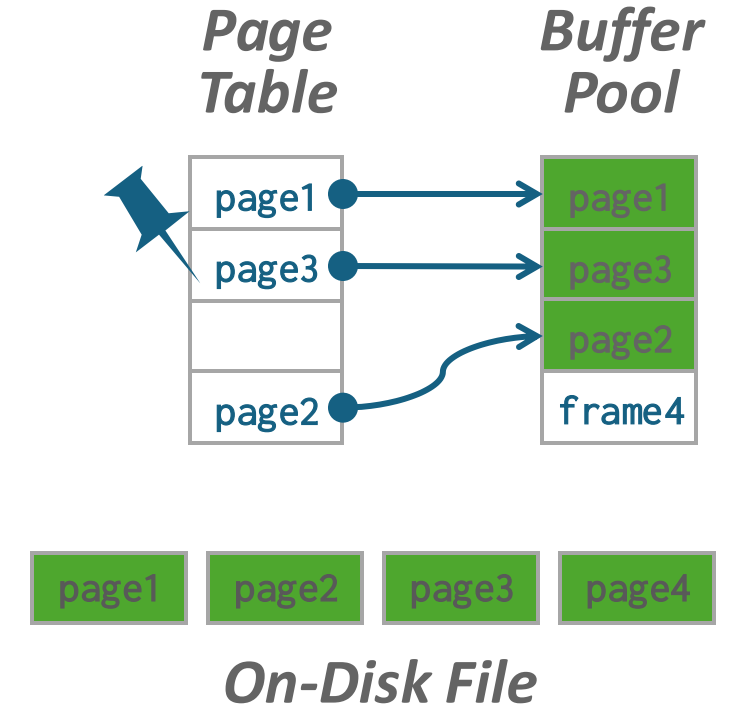
Buffer Pool Meta-data

- The **page table** keeps track of pages that are currently in memory.
 - Usually a fixed-size hash table protected with latches to ensure thread-safe access.
- Additional meta-data per page:
 - Dirty Flag
 - Pin/Reference Counter
 - Access Tracking Information



Buffer Pool Meta-data

- The **page table** keeps track of pages that are currently in memory.
 - Usually a fixed-size hash table protected with latches to ensure thread-safe access.
- Additional meta-data per page:
 - **Dirty Flag**
 - **Pin/Reference Counter**
 - **Access Tracking Information**



Locks vs. Latches

- **Locks:**

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

- **Latches:**

- 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.

Locks vs. Latches

- **Locks:**

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

- **Latches:**

- 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 Tables vs. Page Directory

- The **page directory** is the mapping from page ids to page locations in the database files.
 - All changes must be recorded on disk to allow the DBMS to find on restart.
- 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.

Allocation Policies

- **Global Policies:**
 - Make decisions for all active queries.
- **Local Policies:**
 - Allocate frames to a specific queries without considering the behavior of concurrent queries.
 - Still need to support sharing pages.

Buffer Pool Optimizations

Buffer Pool Optimizations

- Multiple Buffer Pools
- Pre-Fetching
- Scan Sharing
- Buffer Pool Bypass

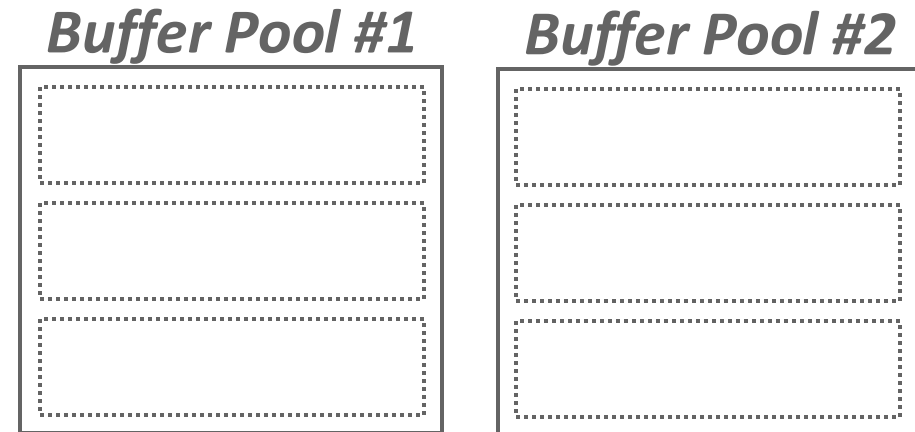
Multiple Buffer Pools

- 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
- Partitioning memory across multiple pools helps reduce latch contention and improve locality.



Multiple Buffer Pools

- **Approach #1: Object Id**
 - Embed an object identifier in record ids and then maintain a mapping from objects to specific buffer pools.
- **Approach #2: Hashing**
 - Hash the page id to select which buffer pool to access.



Multiple Buffer Pools

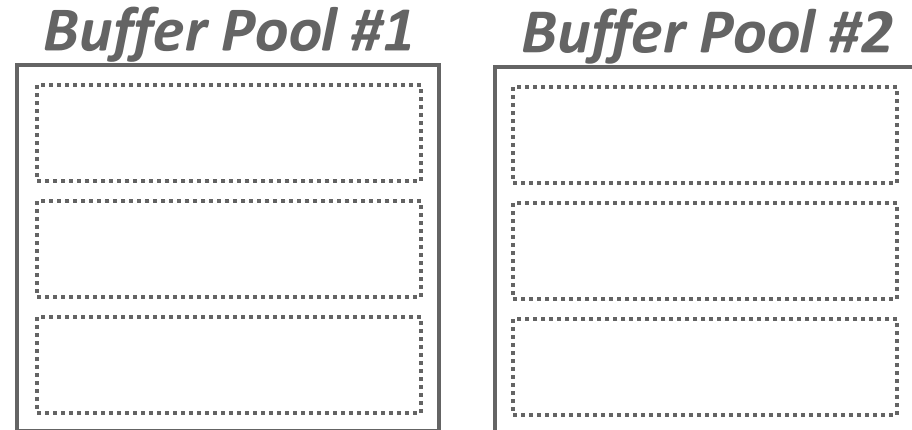
- **Approach #1: Object Id**

- Embed an object identifier in record ids and then maintain a mapping from objects to specific buffer pools.

Q1 GET RECORD #123

- **Approach #2: Hashing**

- Hash the page id to select which buffer pool to access.



Multiple Buffer Pools

- **Approach #1: Object Id**

- Embed an object identifier in record ids and then maintain a mapping from objects to specific buffer pools.

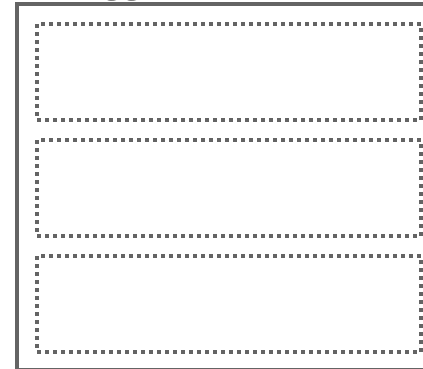
Q1 GET RECORD #123

<ObjectId, PageId, SlotNum>

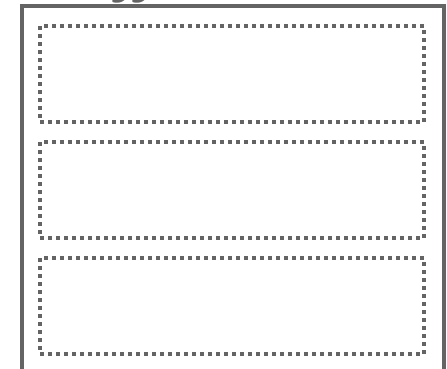
- **Approach #2: Hashing**

- Hash the page id to select which buffer pool to access.

Buffer Pool #1



Buffer Pool #2



Multiple Buffer Pools

- **Approach #1: Object Id**

- Embed an object identifier in record ids and then maintain a mapping from objects to specific buffer pools.

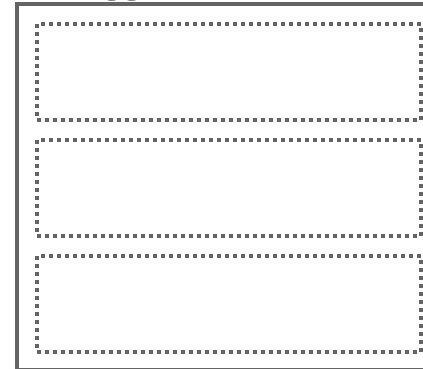
Q1 GET RECORD #123

<ObjectId, PageId, SlotNum>

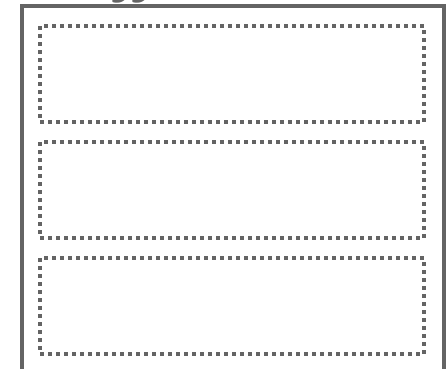
- **Approach #2: Hashing**

- Hash the page id to select which buffer pool to access.

Buffer Pool #1

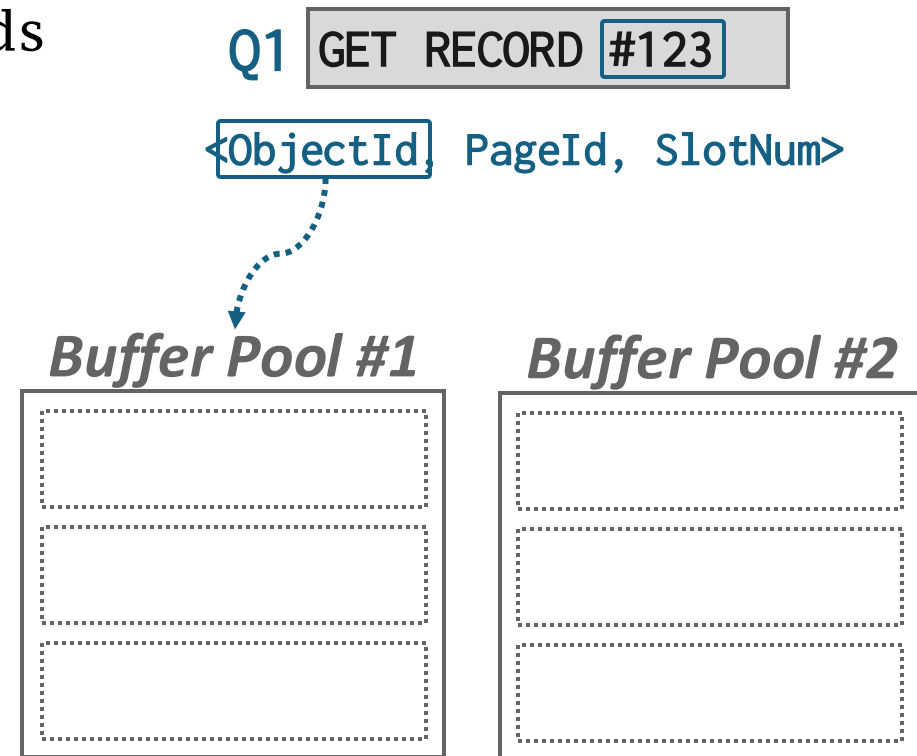


Buffer Pool #2



Multiple Buffer Pools

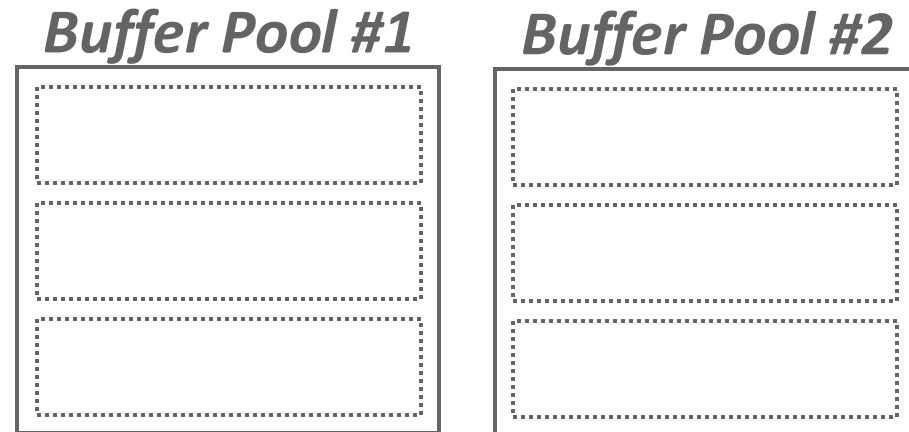
- **Approach #1: Object Id**
 - Embed an object identifier in record ids and then maintain a mapping from objects to specific buffer pools.
- **Approach #2: Hashing**
 - Hash the page id to select which buffer pool to access.



Multiple Buffer Pools

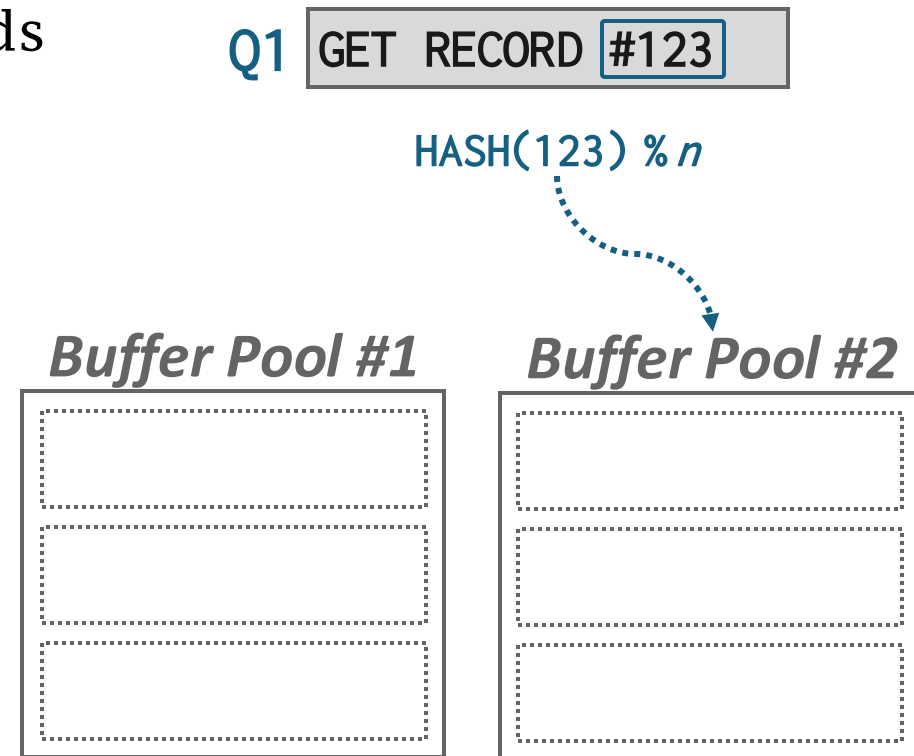
- **Approach #1: Object Id**
 - Embed an object identifier in record ids and then maintain a mapping from objects to specific buffer pools.
- **Approach #2: Hashing**
 - Hash the page id to select which buffer pool to access.

Q1 GET RECORD #123



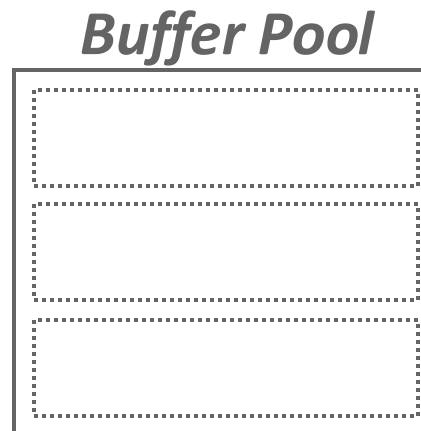
Multiple Buffer Pools

- **Approach #1: Object Id**
 - Embed an object identifier in record ids and then maintain a mapping from objects to specific buffer pools.
- **Approach #2: Hashing**
 - Hash the page id to select which buffer pool to access.



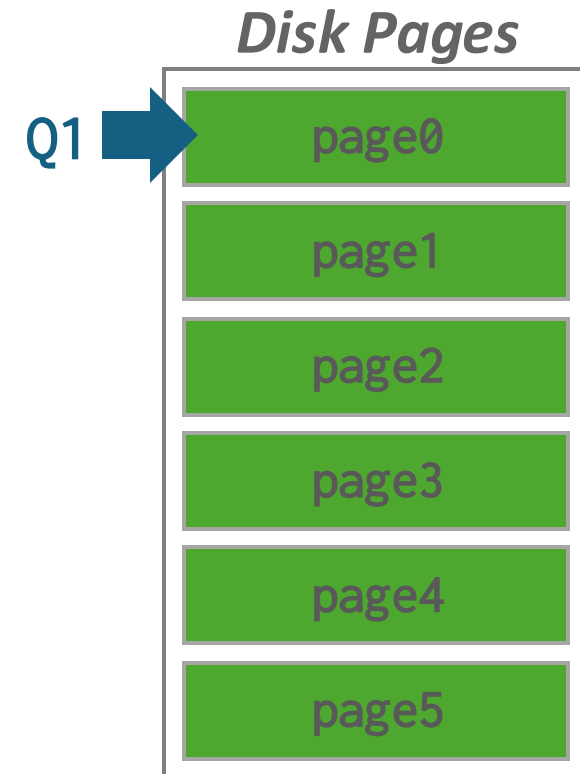
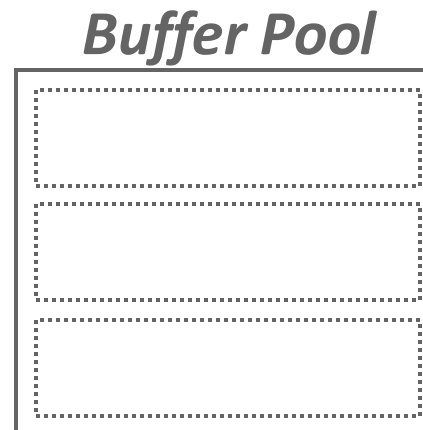
Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans



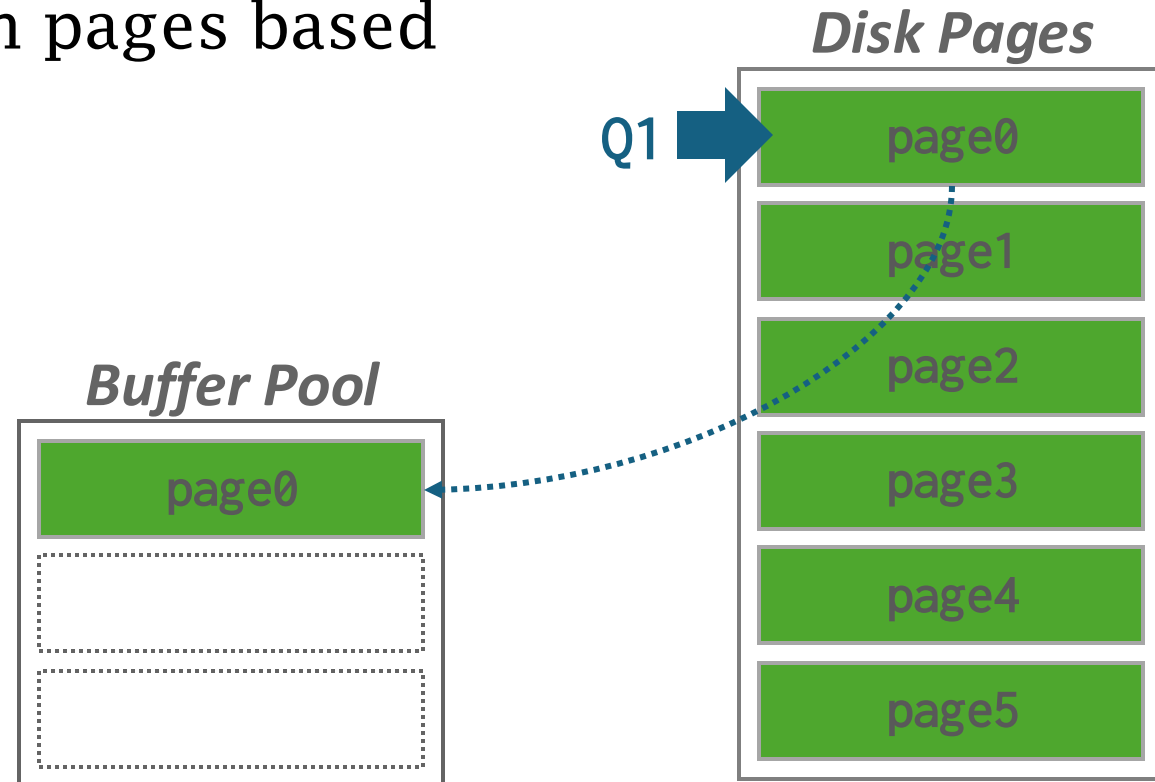
Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans



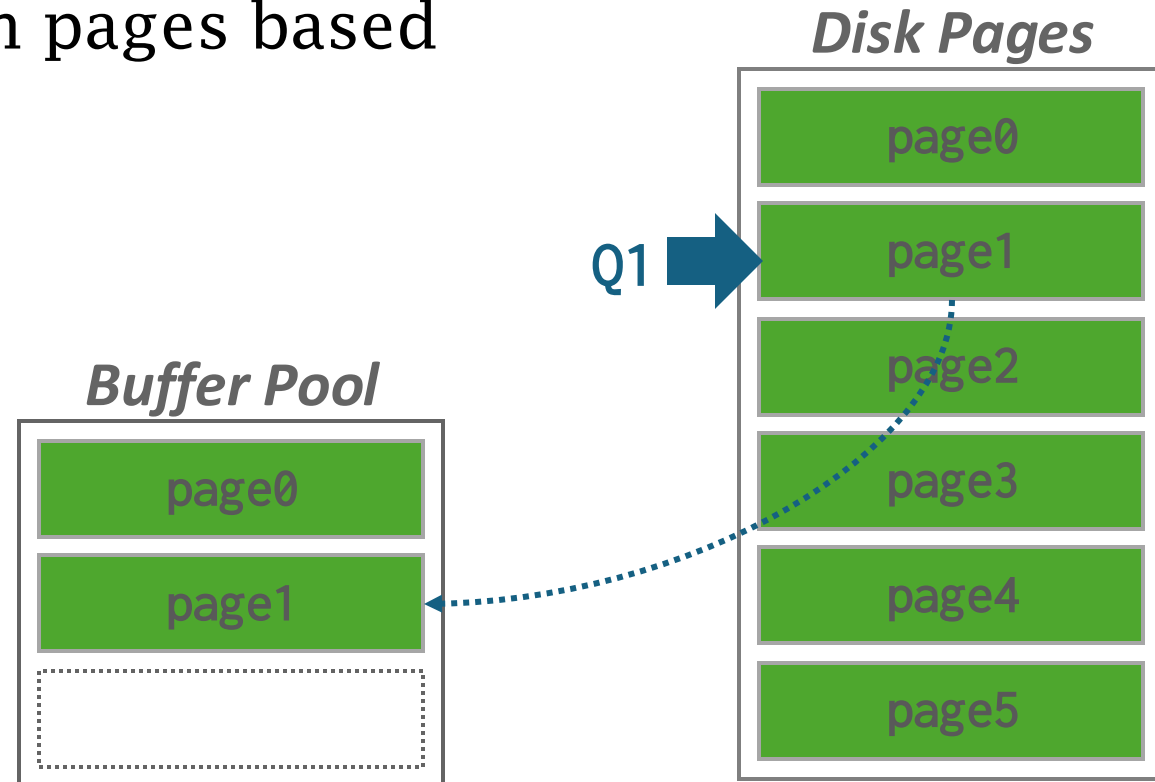
Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans



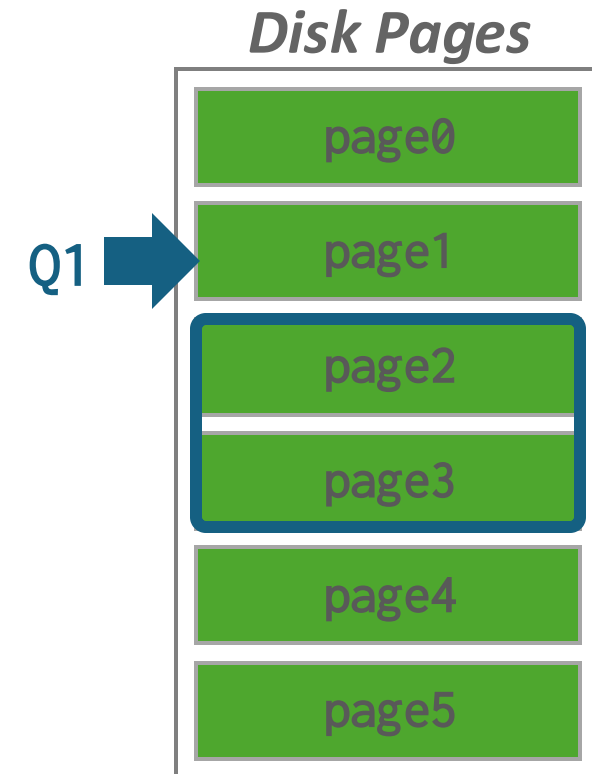
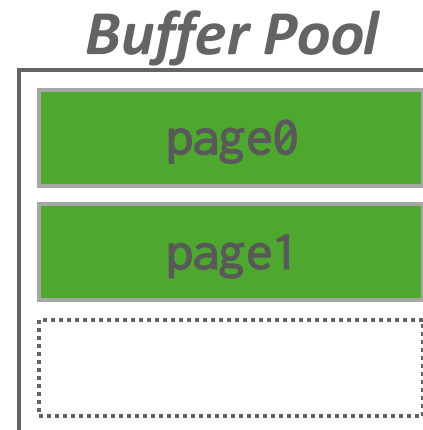
Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans



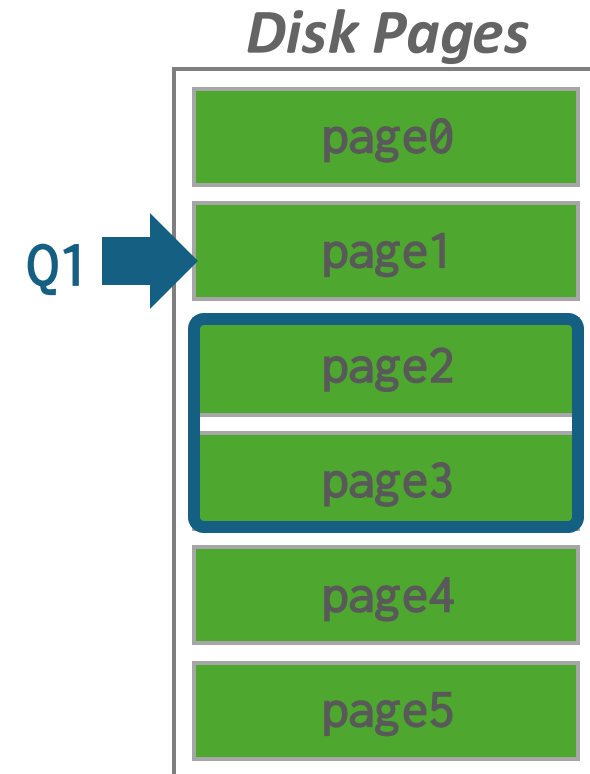
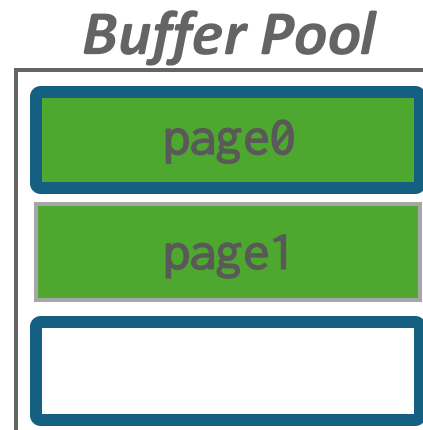
Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans



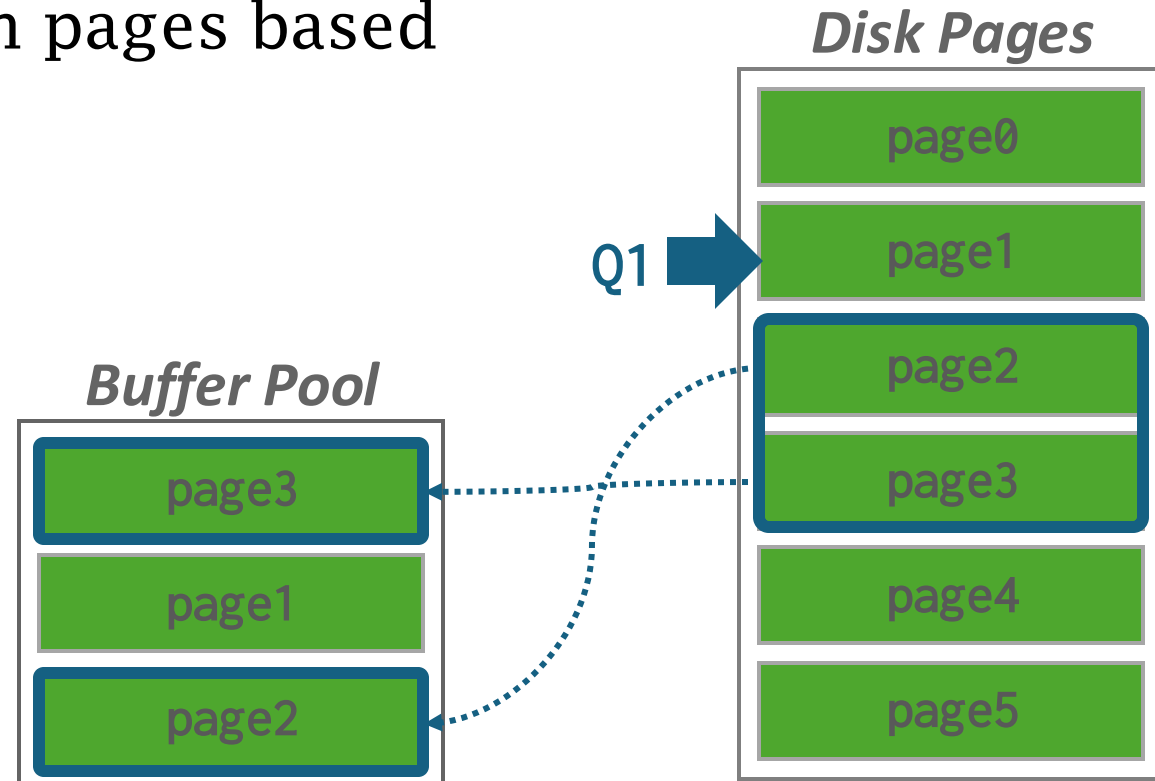
Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans



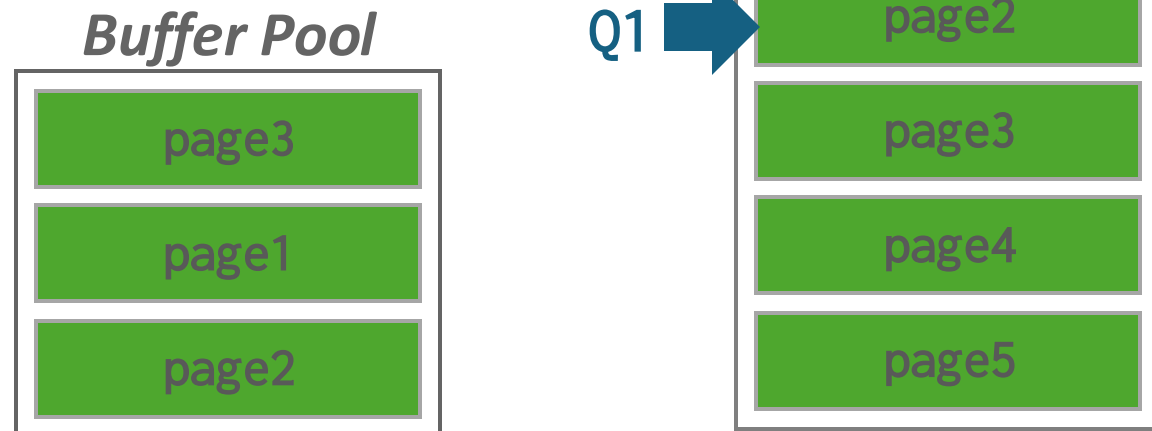
Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans



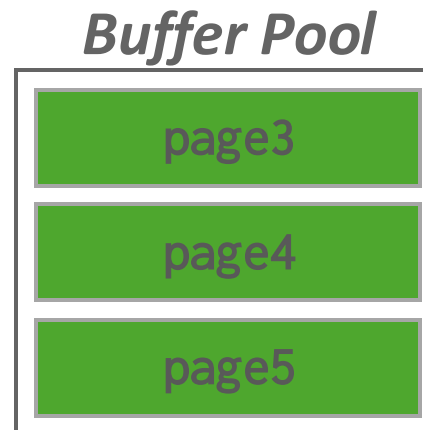
Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans

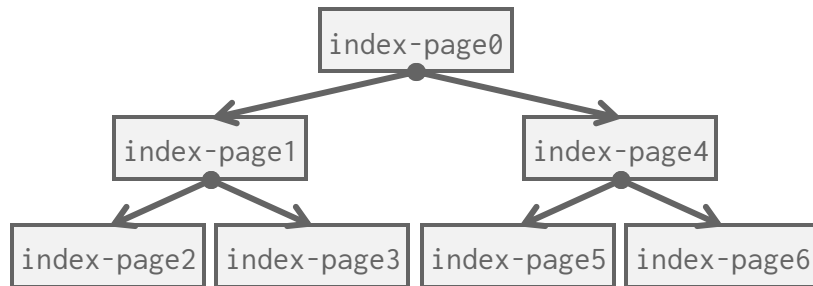


Pre-fetching

- The DBMS can also prefetch pages based on a query plan.
 - Sequential Scans
 - Index Scans



Pre-fetching



Q1

```
SELECT * FROM A
WHERE val BETWEEN 100 AND 250
```

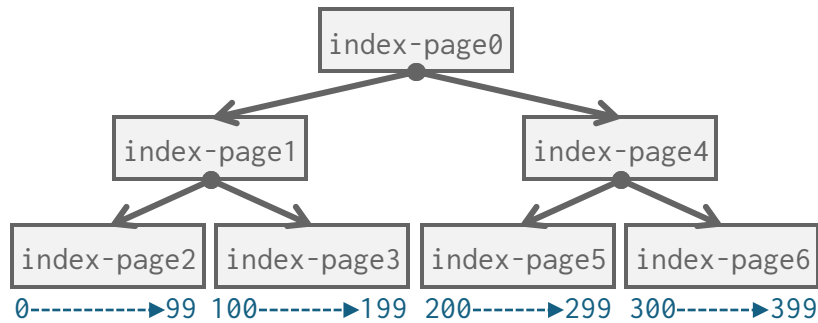
Buffer Pool



Disk Pages



Pre-fetching



Q1

```
SELECT * FROM A
WHERE val BETWEEN 100 AND 250
```

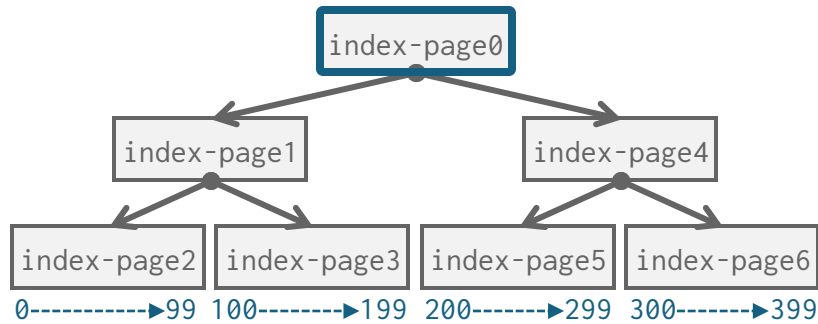
Buffer Pool



Disk Pages



Pre-fetching



Q1

```
SELECT * FROM A
WHERE val BETWEEN 100 AND 250
```

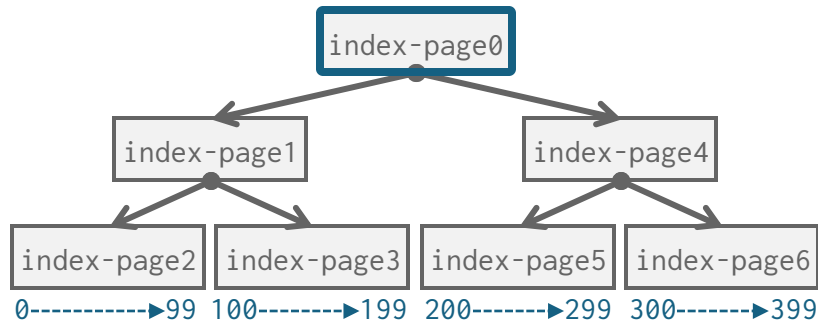
Buffer Pool



Disk Pages



Pre-fetching



Q1

```
SELECT * FROM A
WHERE val BETWEEN 100 AND 250
```

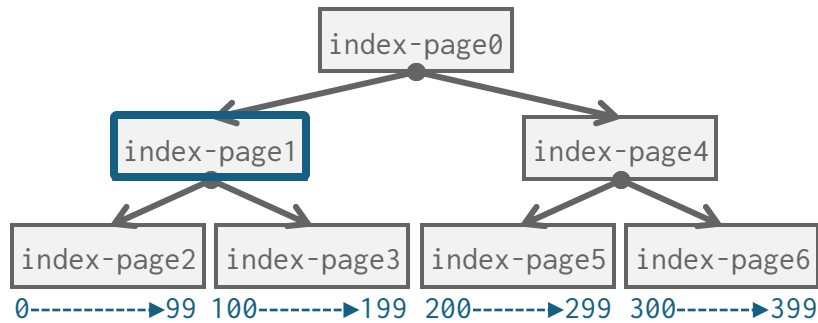
Buffer Pool



Disk Pages



Pre-fetching



Q1

```
SELECT * FROM A
WHERE val BETWEEN 100 AND 250
```

Buffer Pool

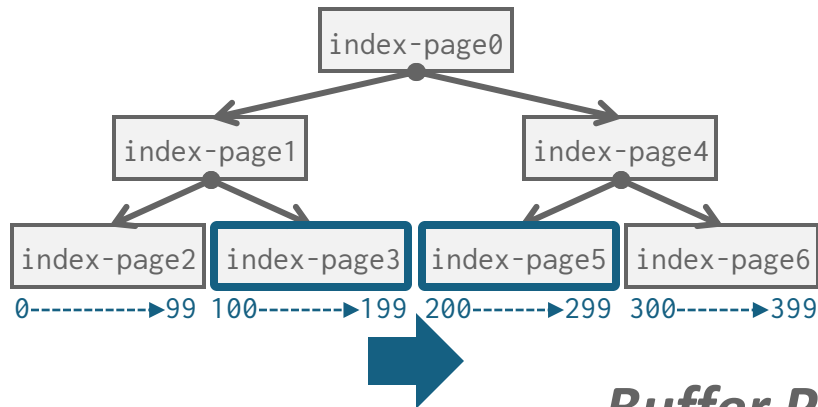


Disk Pages

Q1 →



Pre-fetching



Q1

```
SELECT * FROM A
WHERE val BETWEEN 100 AND 250
```

Buffer Pool

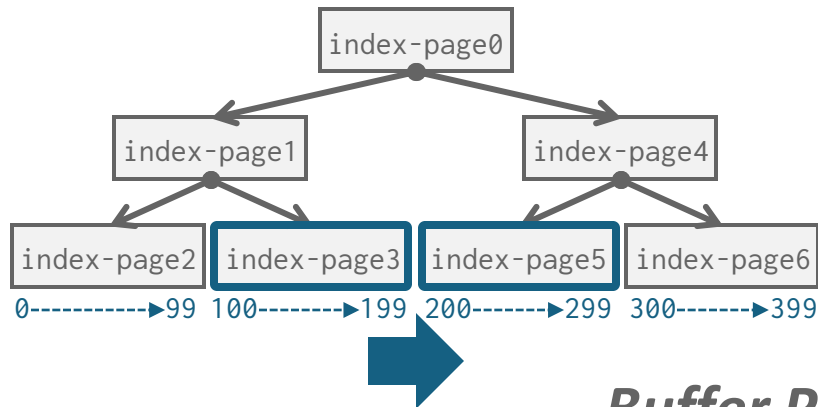


Disk Pages

Q1 →

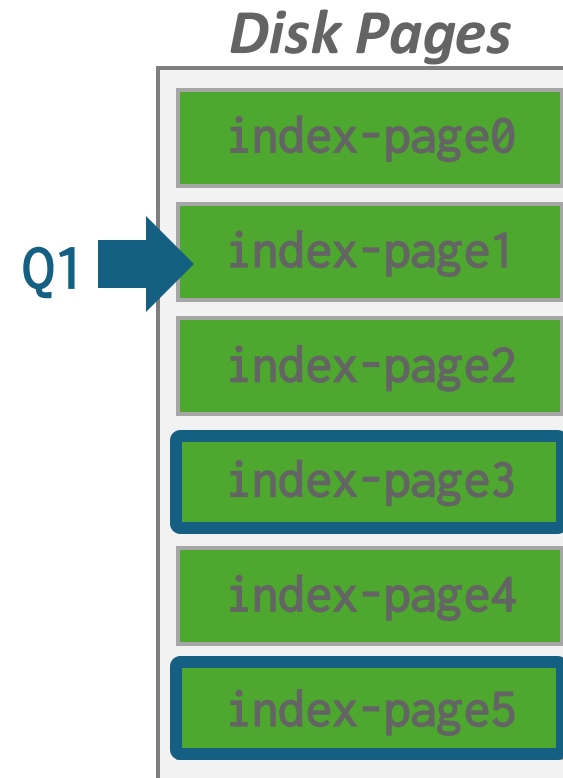
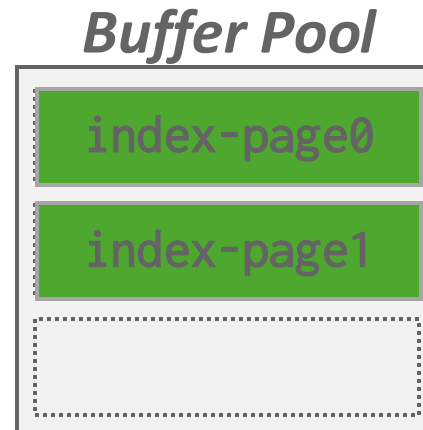


Pre-fetching



Q1

```
SELECT * FROM A
WHERE val BETWEEN 100 AND 250
```



Scan Sharing

- Queries can reuse data retrieved from storage or operator computations.
 - Also called *synchronized scans*.
 - 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 the same.
 - Can also share intermediate results.

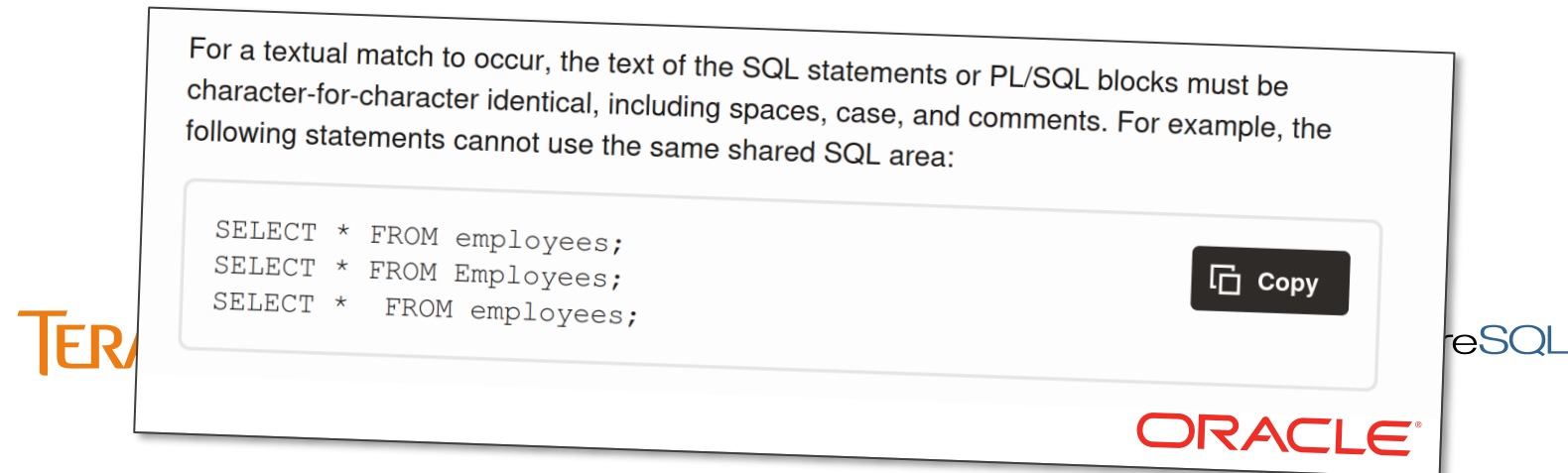
Scan Sharing

- If a query wants to scan a table and another query is already doing this, then the DBMS will attach the second query's cursor to the existing cursor.
- Examples:
 - Fully supported in DB2, MSSQL, Teradata, and Postgres.
 - Oracle only supports [cursor sharing](#) for identical queries.



Scan Sharing

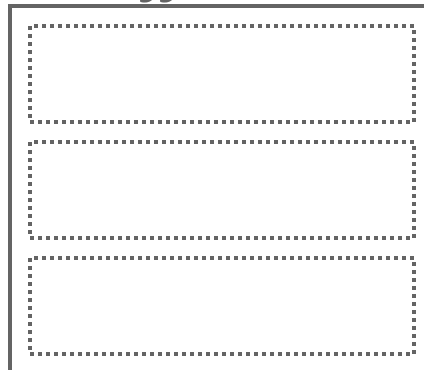
- If a query wants to scan a table and another query is already doing this, then the DBMS will attach the second query's cursor to the existing cursor.
- Examples:
 - Fully supported in DB2, MSSQL, Teradata, and Postgres.
 - Oracle only supports [cursor sharing](#) for identical queries.



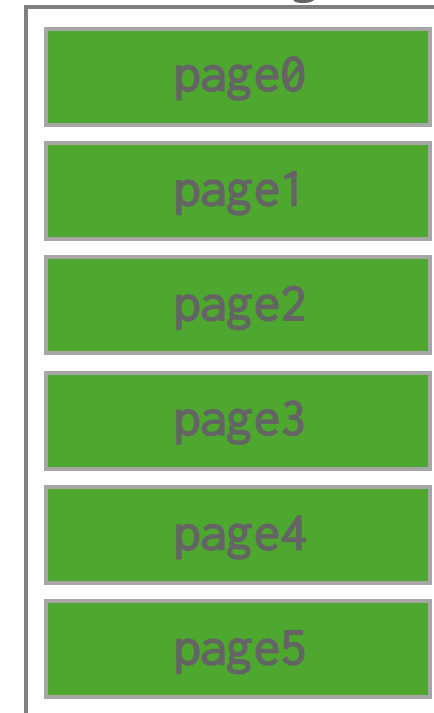
Scan Sharing

Q1 `SELECT SUM(val) FROM A`

Buffer Pool



Disk Pages

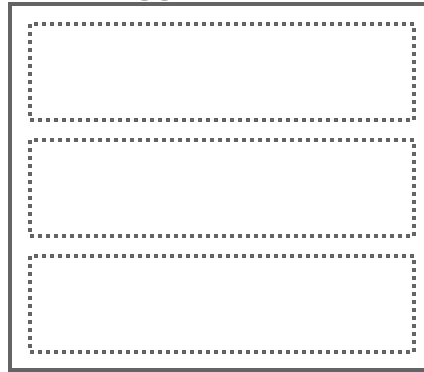


Scan Sharing

Q1

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

Buffer Pool



Disk Pages



Scan Sharing

Q1

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

Buffer Pool

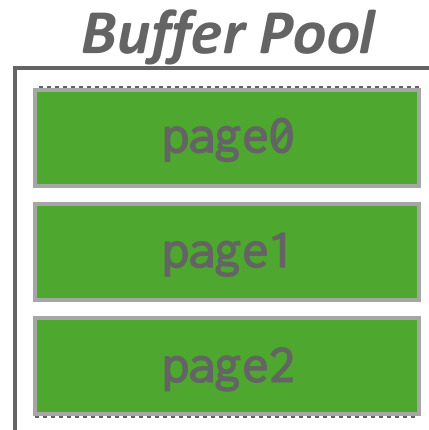


Disk Pages



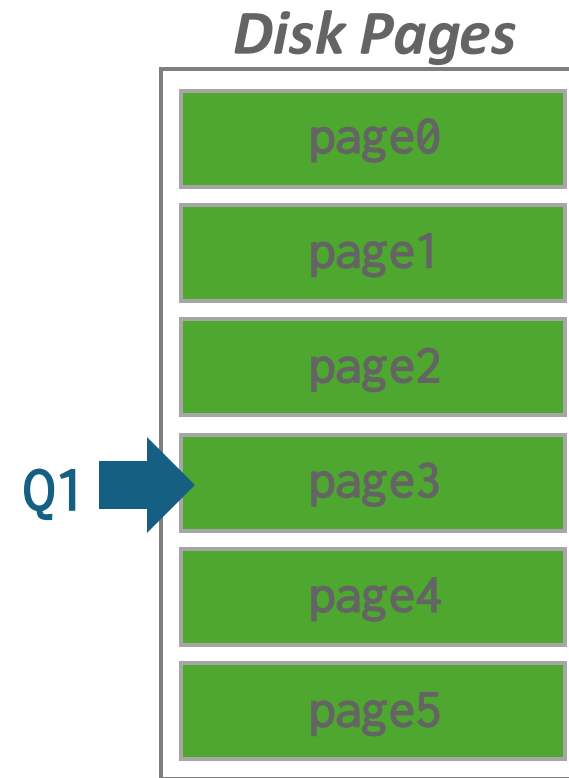
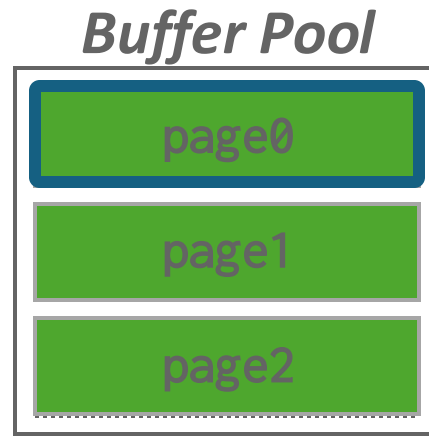
Scan Sharing

Q1 `SELECT SUM(val) FROM A`



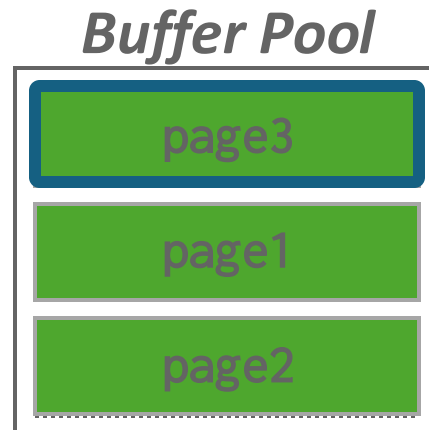
Scan Sharing

Q1 `SELECT SUM(val) FROM A`



Scan Sharing

Q1 `SELECT SUM(val) FROM A`



Scan Sharing

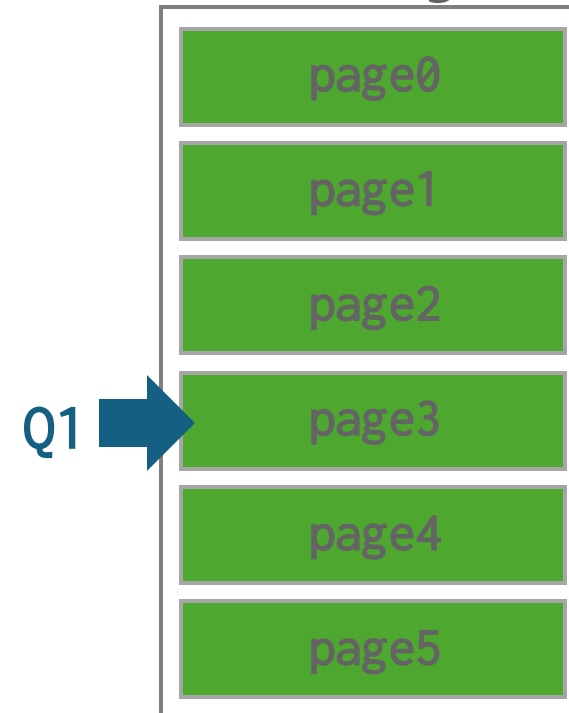
Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



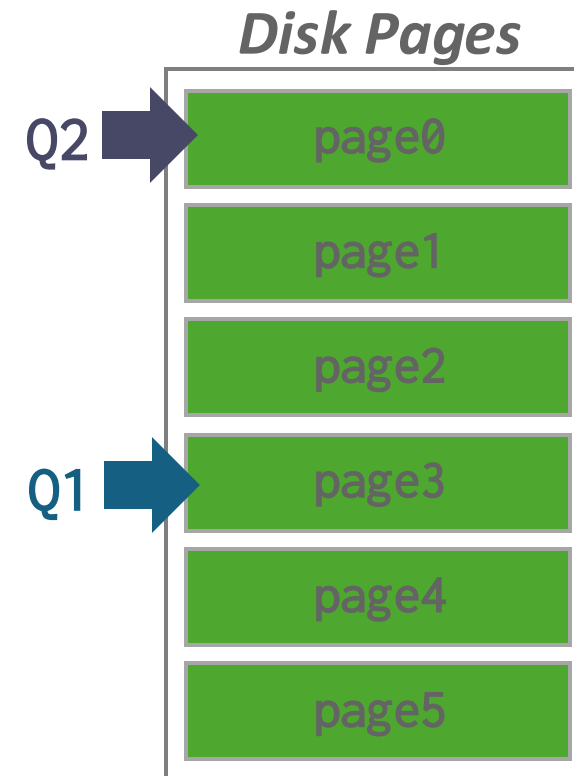
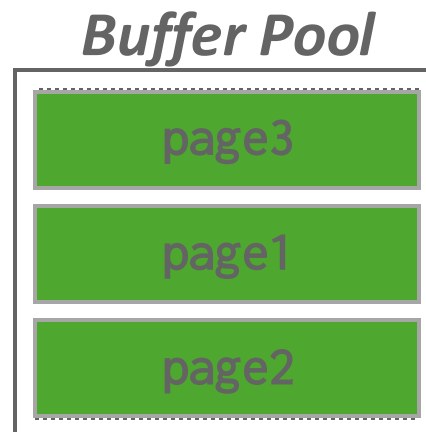
Disk Pages



Scan Sharing

Q1 `SELECT SUM(val) FROM A`

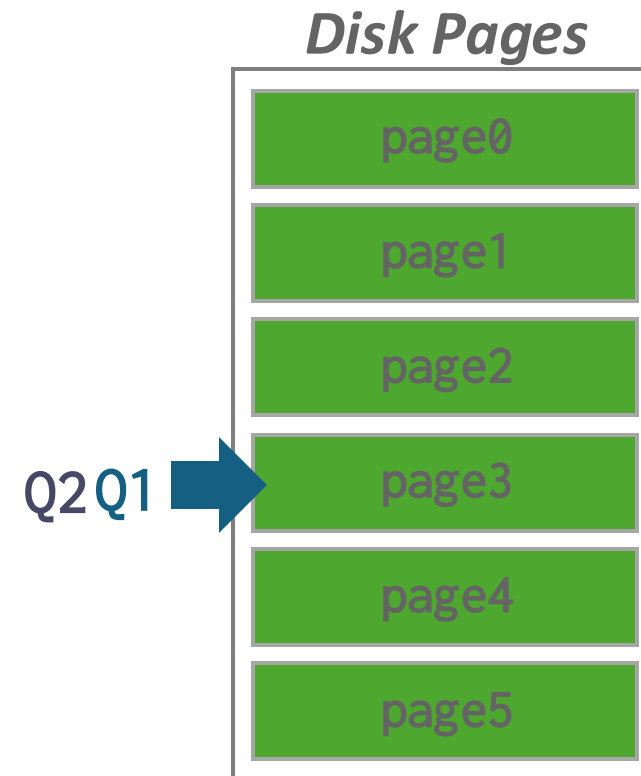
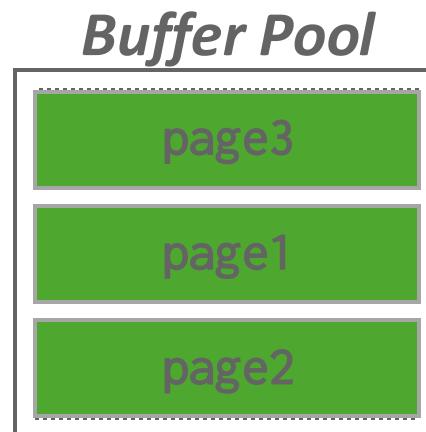
Q2 `SELECT AVG(val) FROM A`



Scan Sharing

Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A`



Scan Sharing

Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



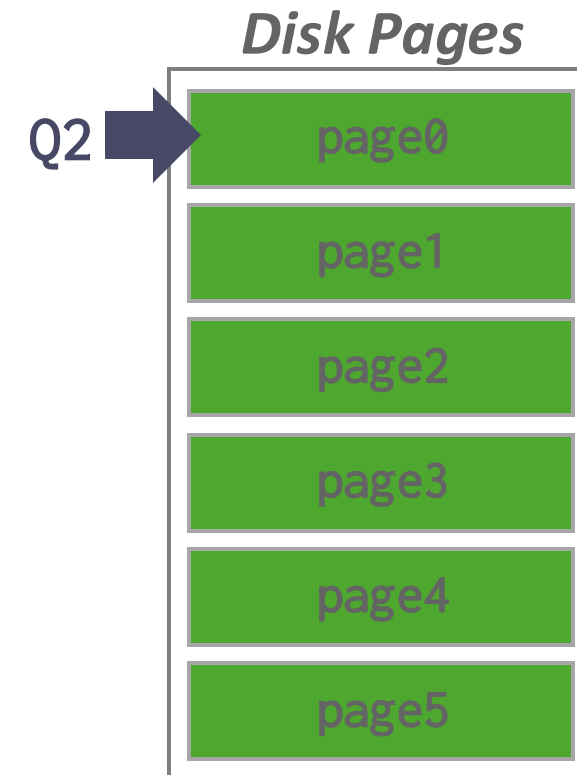
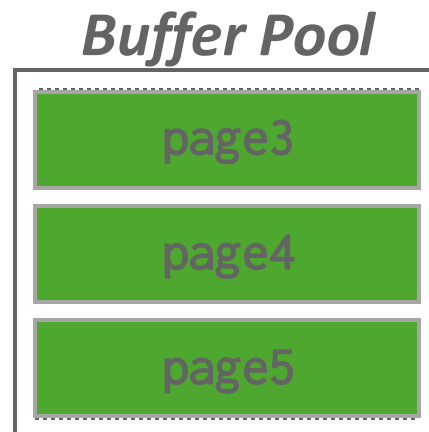
Disk Pages



Scan Sharing

Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A`



Scan Sharing

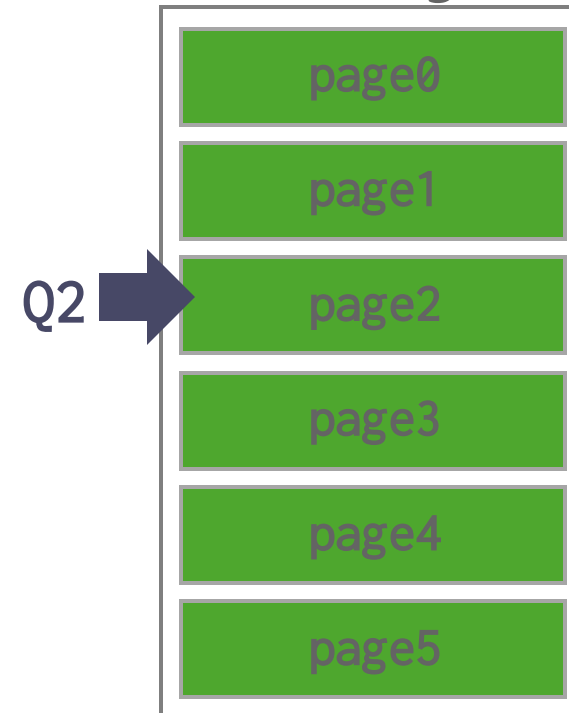
Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



Disk Pages



Scan Sharing

Q1

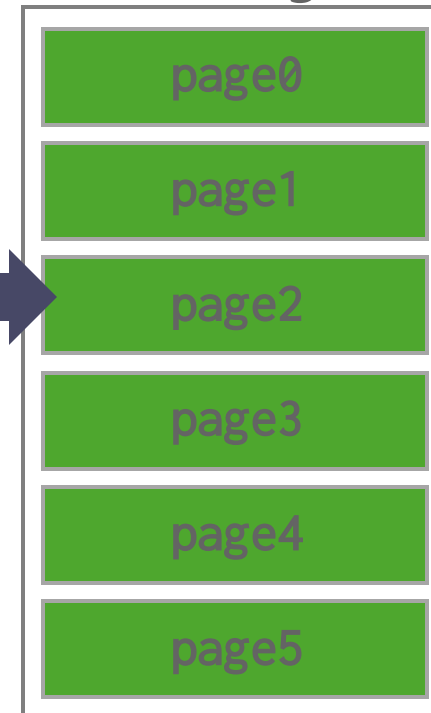
```
SELECT SUM(val) FROM A
```

Buffer Pool



Disk Pages

Q2 →



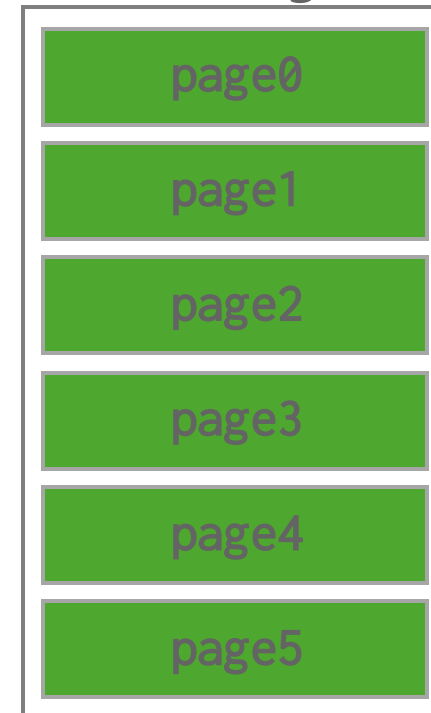
Scan Sharing

Q1 `SELECT SUM(val) FROM A`

Buffer Pool

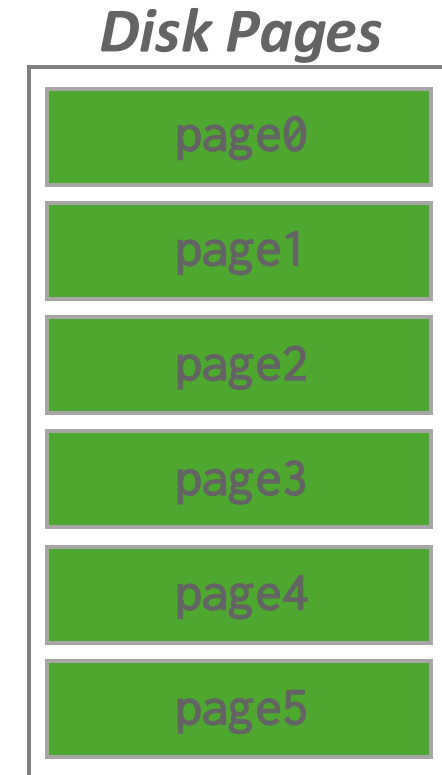


Disk Pages



Continuous Scan Sharing

- Instead of trying to be clever, the DBMS continuously scans the database files repeatedly.
 - One continuous cursor per table.
 - Queries “hop” on board the cursor while it is running and then disconnect once they have enough data.
- Not viable if you pay per IOP.
Only done in academic prototypes.



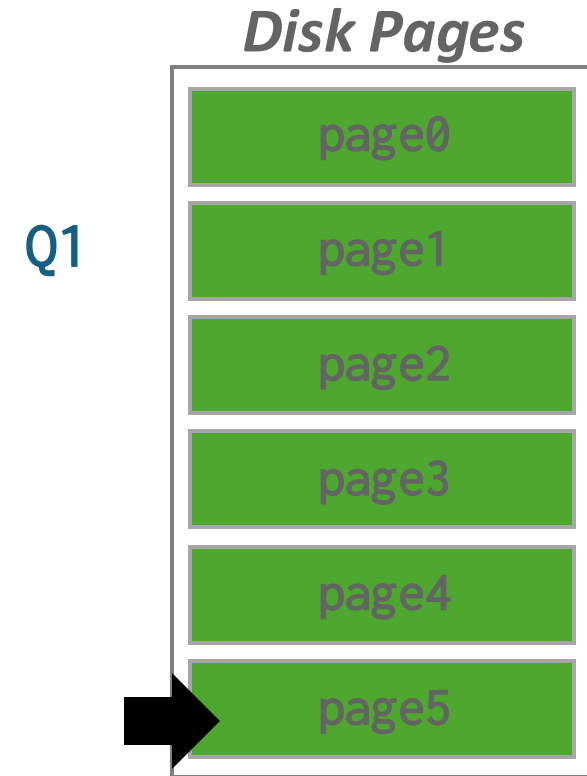
Continuous Scan Sharing

- Instead of trying to be clever, the DBMS continuously scans the database files repeatedly.
 - One continuous cursor per table.
 - Queries “hop” on board the cursor while it is running and then disconnect once they have enough data.
- Not viable if you pay per IOP.
Only done in academic prototypes.



Continuous Scan Sharing

- Instead of trying to be clever, the DBMS continuously scans the database files repeatedly.
 - One continuous cursor per table.
 - Queries “hop” on board the cursor while it is running and then disconnect once they have enough data.
- Not viable if you pay per IOP.
Only done in academic prototypes.



Continuous Scan Sharing

- Instead of trying to be clever, the DBMS continuously scans the database files repeatedly.
 - One continuous cursor per table.
 - Queries “hop” on board the cursor while it is running and then disconnect once they have enough data.
- Not viable if you pay per IOP.
Only done in academic prototypes.



Continuous Scan Sharing

- Instead of trying to be clever, the DBMS continuously scans the database files repeatedly.
 - One continuous cursor per table.
 - Queries “hop” on board the cursor while it is running and then disconnect once they have enough data.
- Not viable if you pay per IOP.
Only done in academic prototypes.



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.
 - Can also be used for temporary data (sorting, joins).
- Called “Light Scans” in Informix.



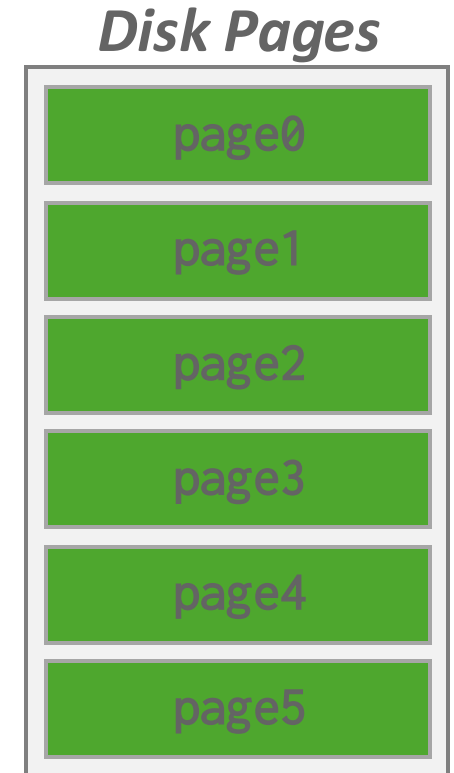
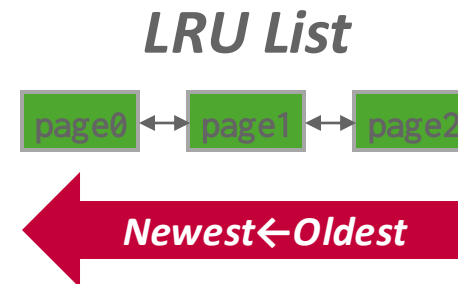
Buffer Replacement Policies

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
 - Metadata overhead

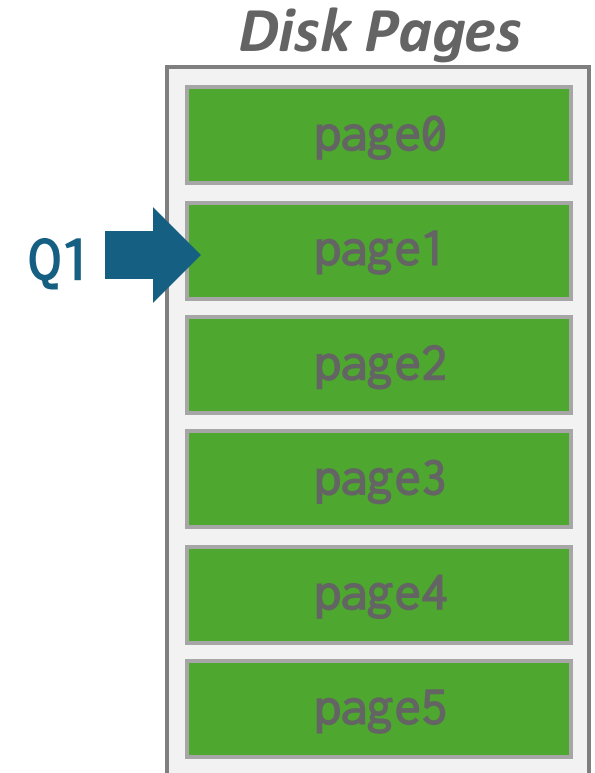
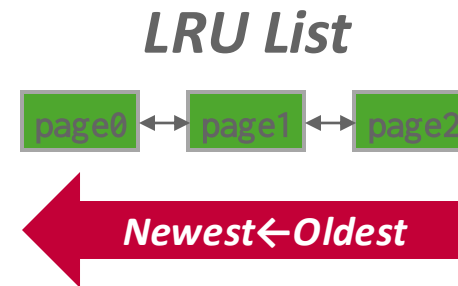
Least-Recently Used

- Maintain a single 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.



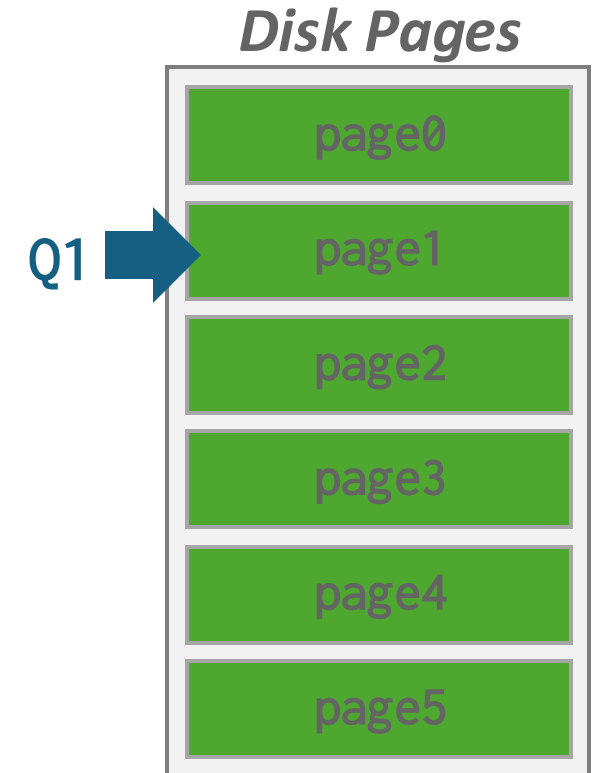
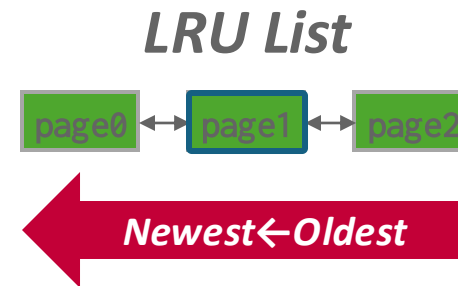
Least-Recently Used

- Maintain a single 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.



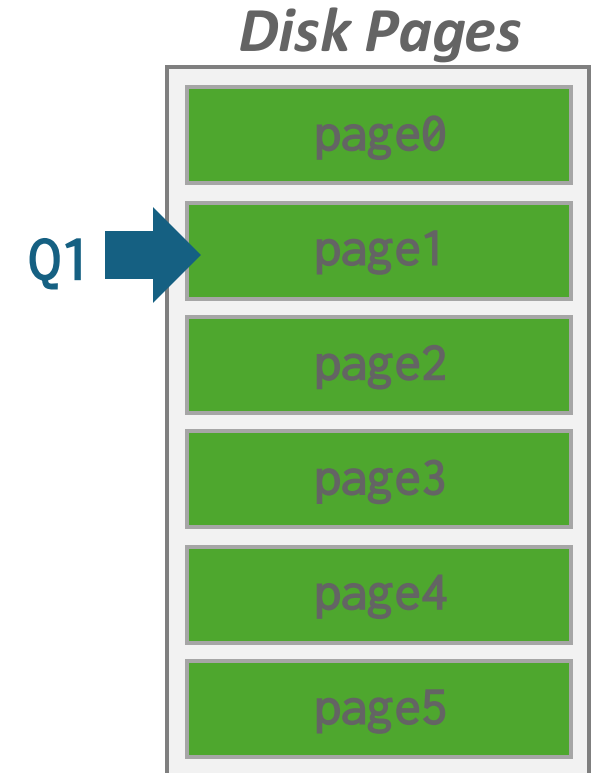
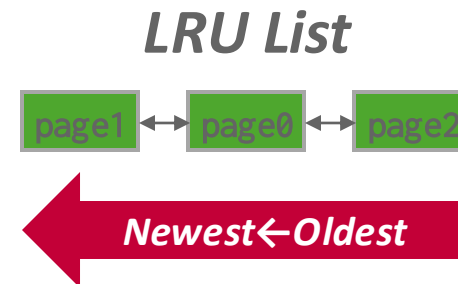
Least-Recently Used

- Maintain a single 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.



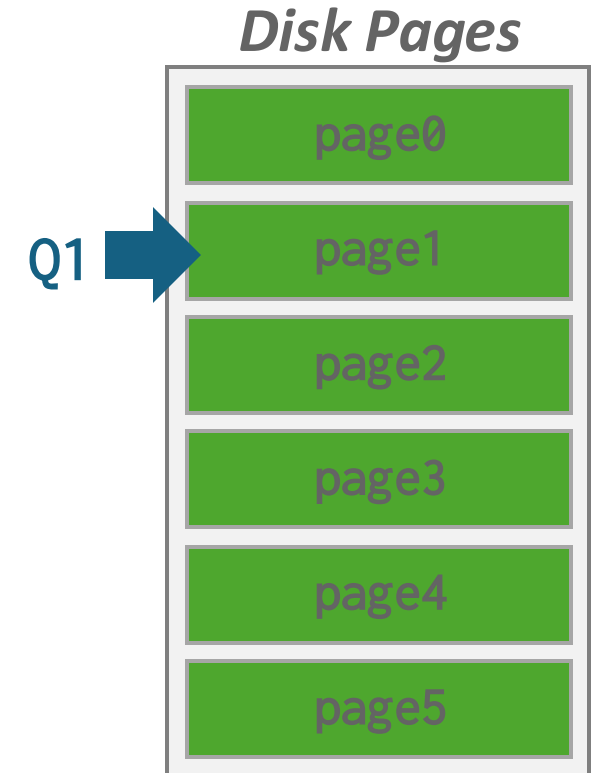
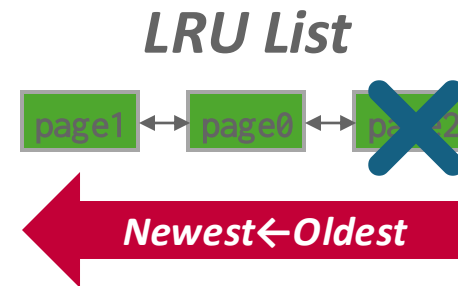
Least-Recently Used

- Maintain a single 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.



Least-Recently Used

- Maintain a single 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.

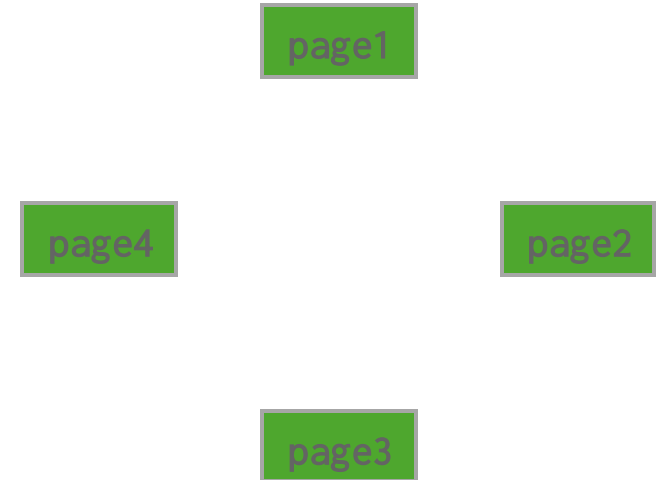


Clock

- Approximation of LRU that does not need 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.

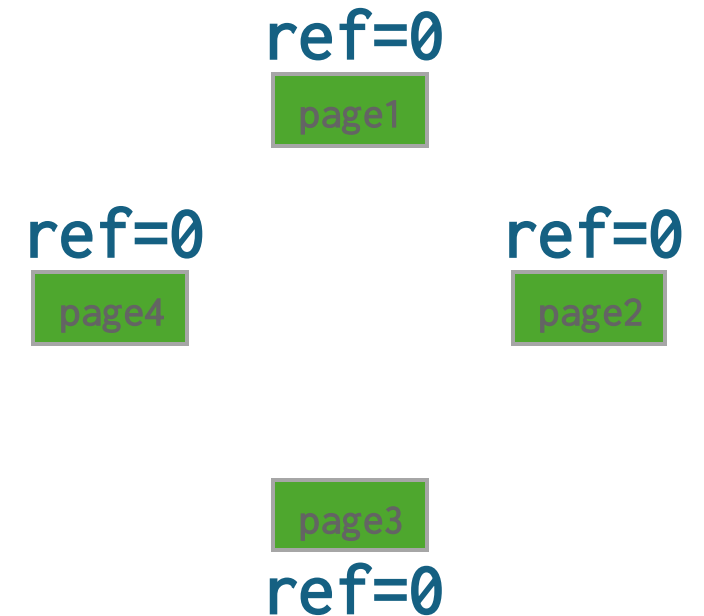
Clock

- Approximation of LRU that does not need 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.



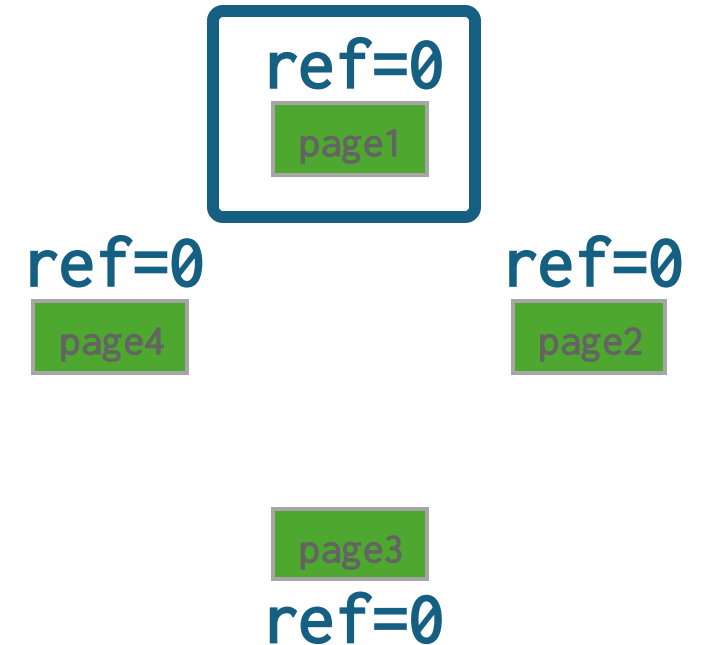
Clock

- Approximation of LRU that does not need 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.



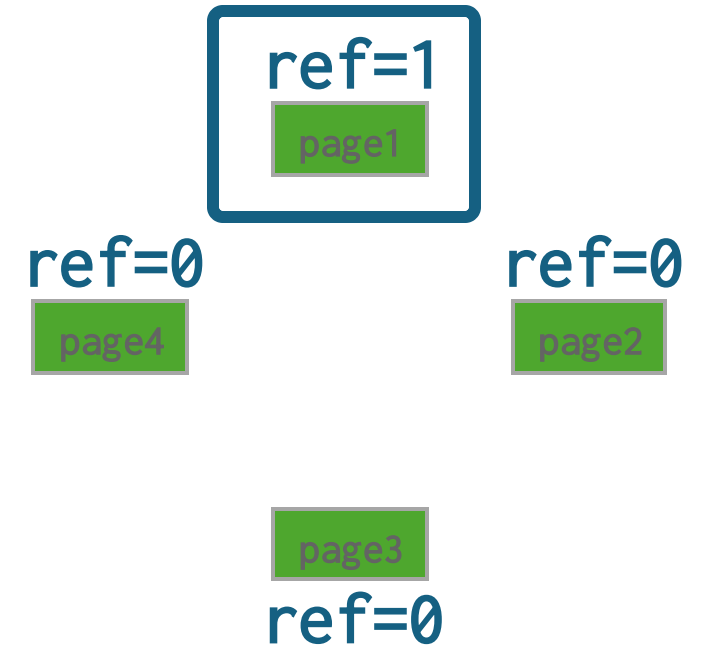
Clock

- Approximation of LRU that does not need 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.



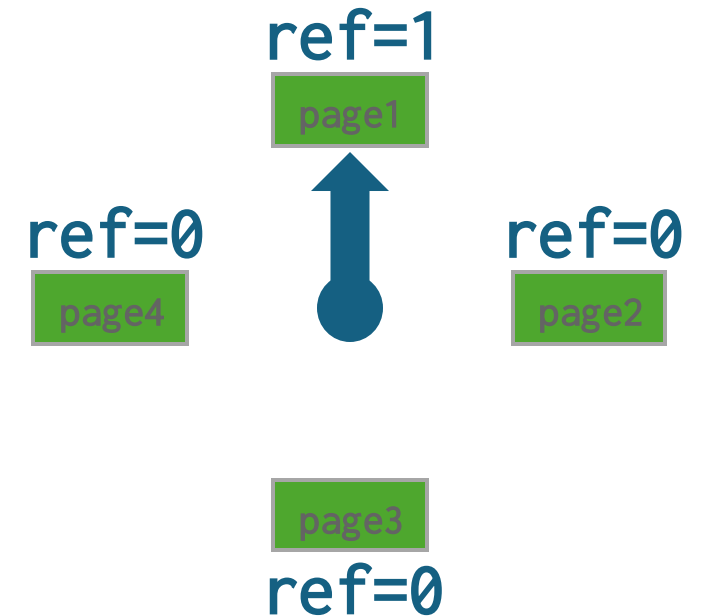
Clock

- Approximation of LRU that does not need 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.



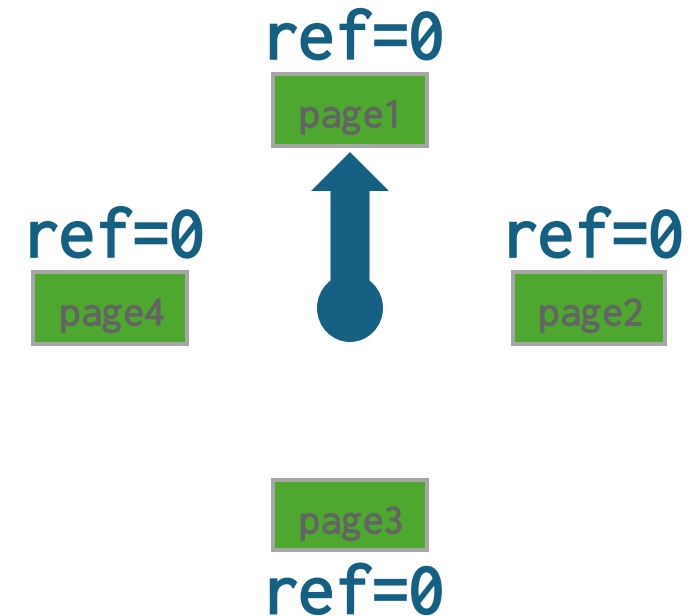
Clock

- Approximation of LRU that does not need 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.



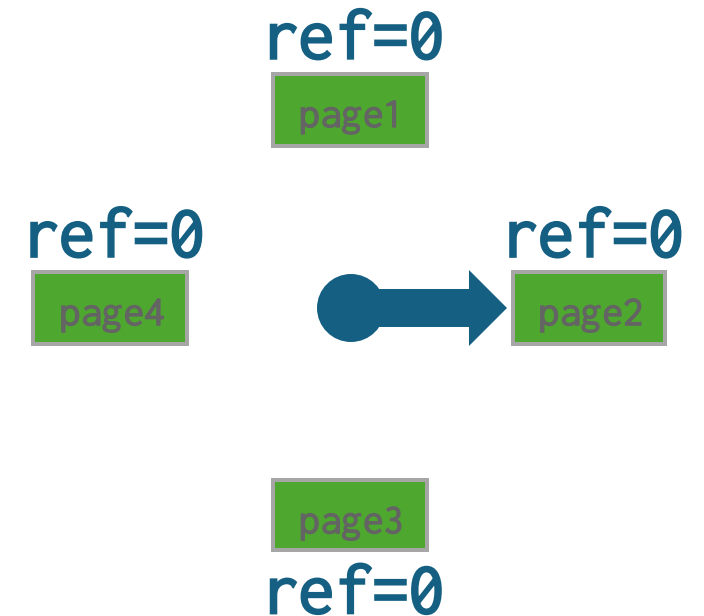
Clock

- Approximation of LRU that does not need 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.



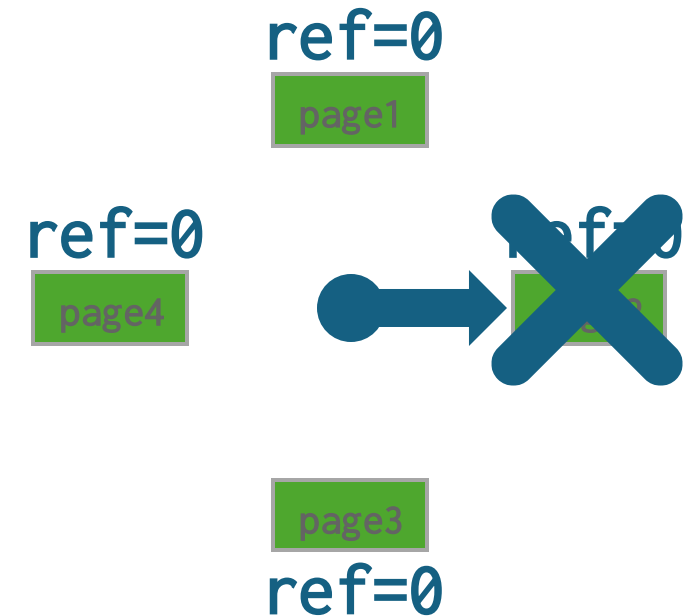
Clock

- Approximation of LRU that does not need 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.



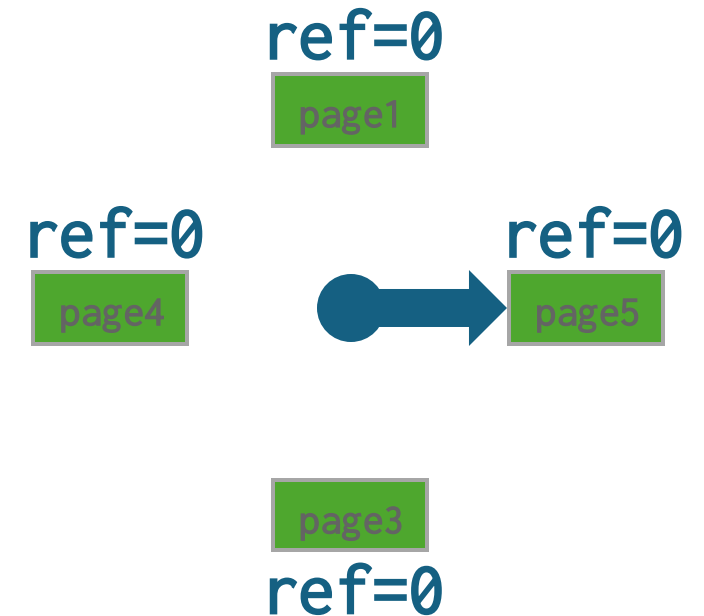
Clock

- Approximation of LRU that does not need 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.



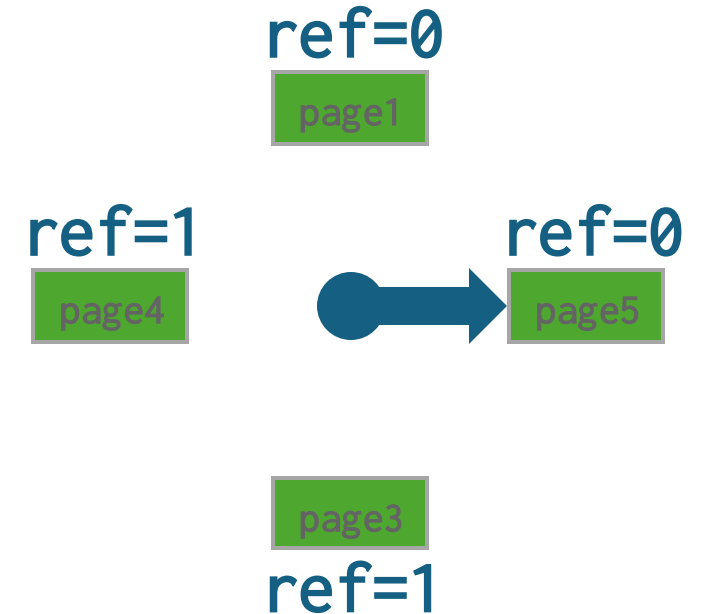
Clock

- Approximation of LRU that does not need 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.



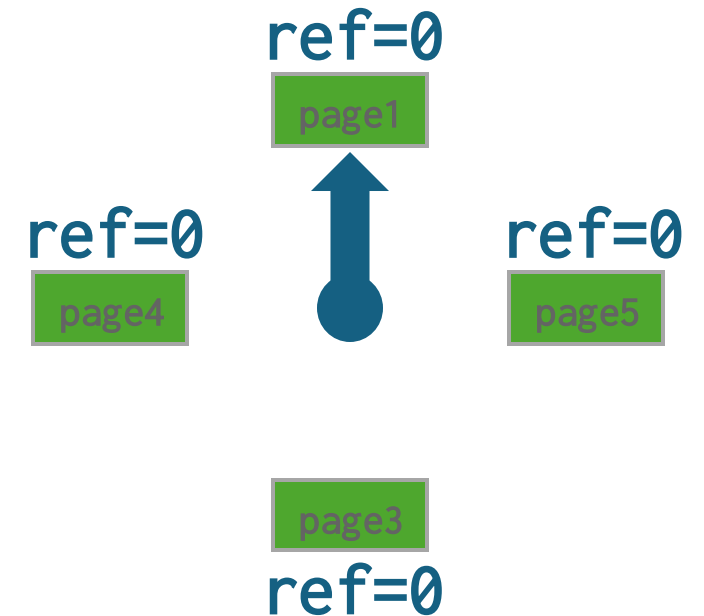
Clock

- Approximation of LRU that does not need 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.



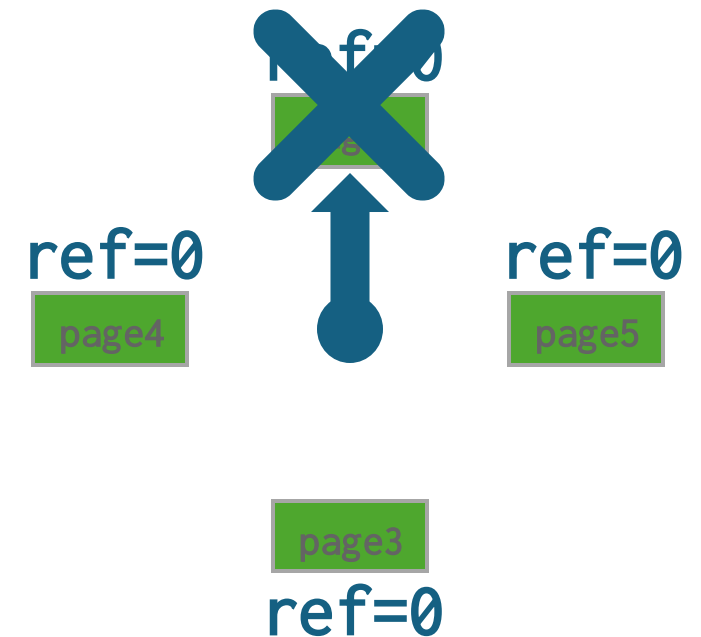
Clock

- Approximation of LRU that does not need 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.



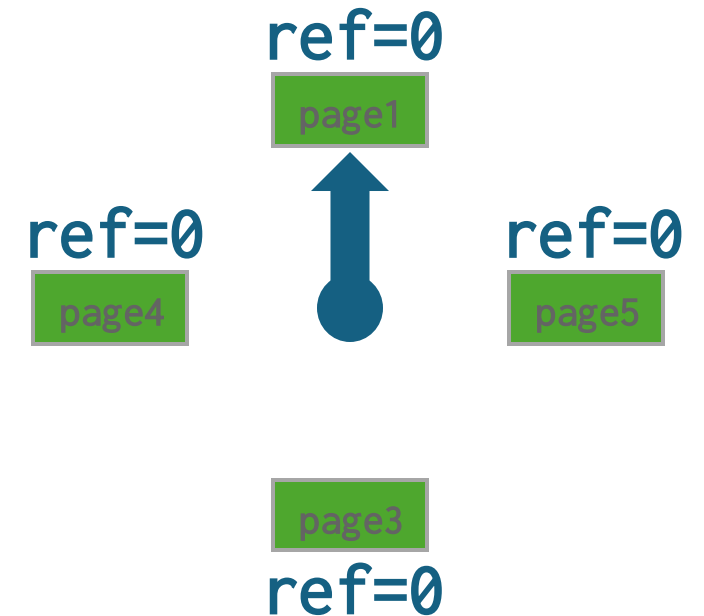
Clock

- Approximation of LRU that does not need 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 that does not need 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.



Observation

- LRU + 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.
 - In OLAP workloads, the *most recently used* page is often the best page to evict.
- LRU + CLOCK only tracks when a page was last accessed, but not how often a page is accessed.

Sequential Flooding

Q1 `SELECT * FROM A WHERE primaryKey = 1`

Buffer Pool

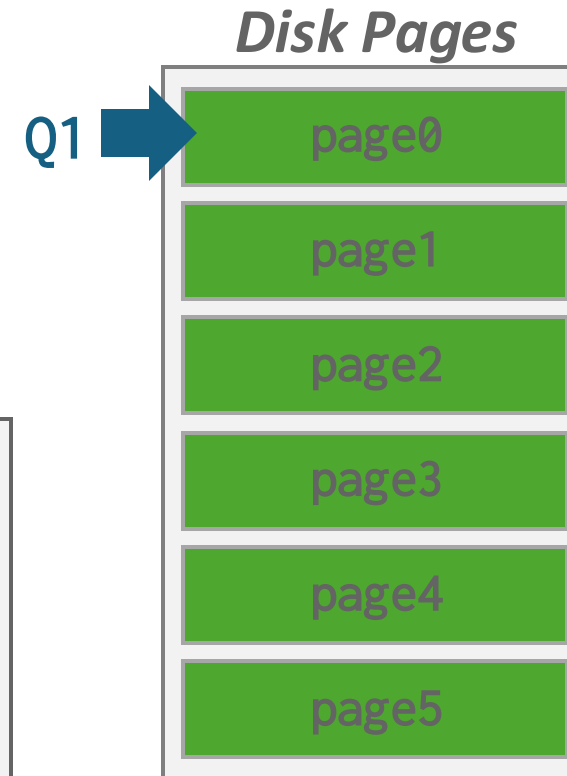
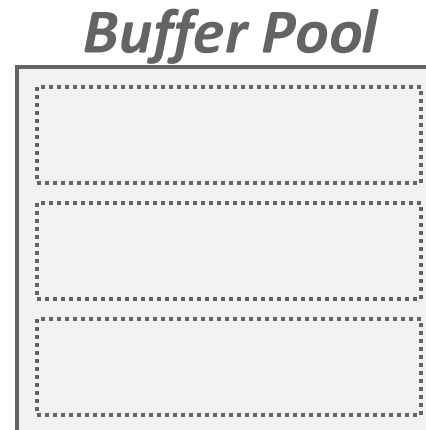


Disk Pages



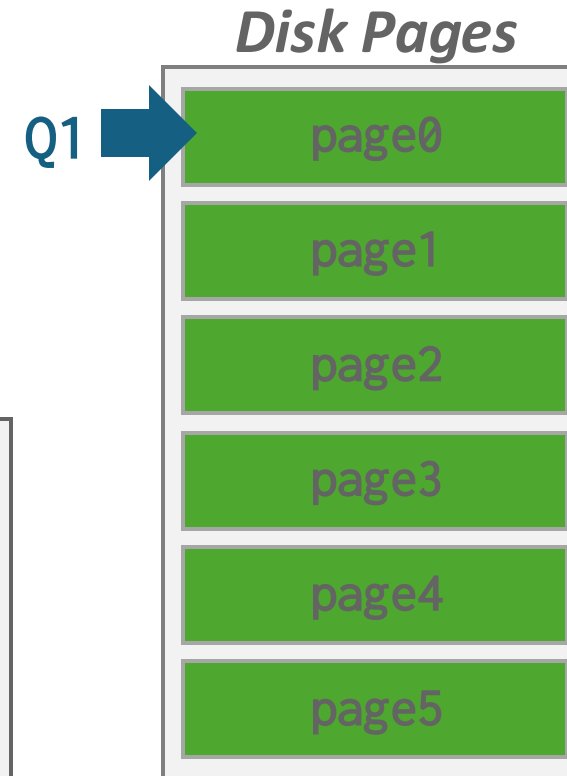
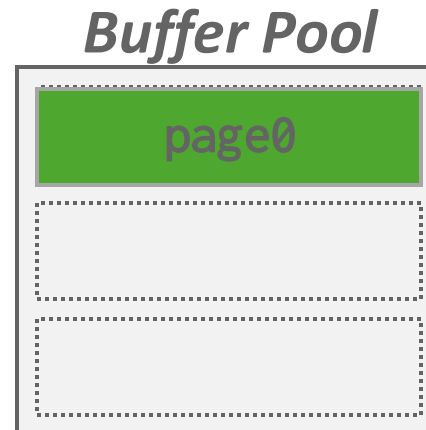
Sequential Flooding

Q1 `SELECT * FROM A WHERE primaryKey = 1`



Sequential Flooding

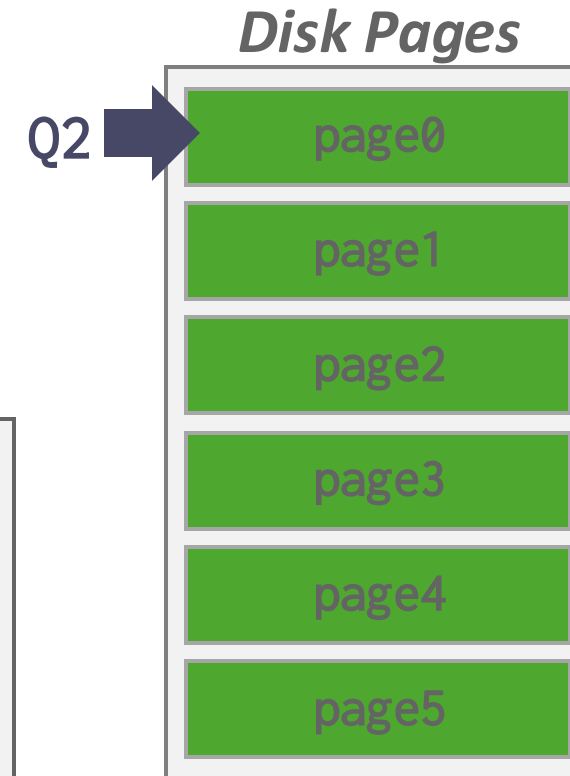
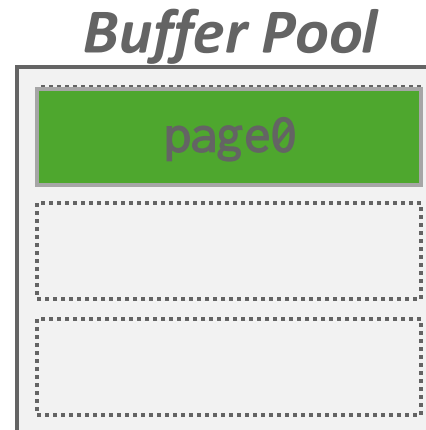
Q1 `SELECT * FROM A WHERE primaryKey = 1`



Sequential Flooding

Q1 `SELECT * FROM A WHERE primaryKey = 1`

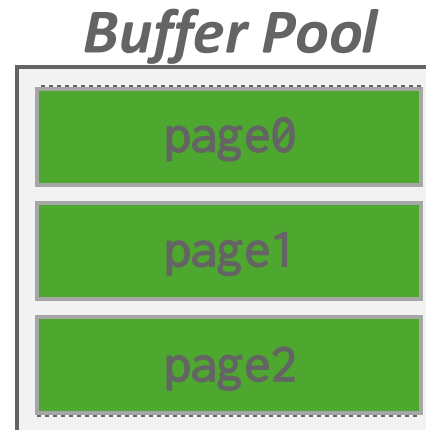
Q2 `SELECT AVG(val) FROM A`



Sequential Flooding

Q1 `SELECT * FROM A WHERE primaryKey = 1`

Q2 `SELECT AVG(val) FROM A`



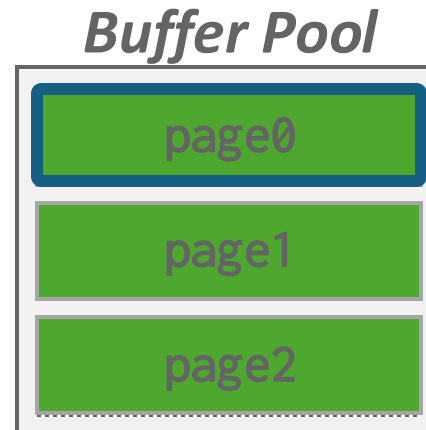
Q2 →



Sequential Flooding

Q1 `SELECT * FROM A WHERE primaryKey = 1`

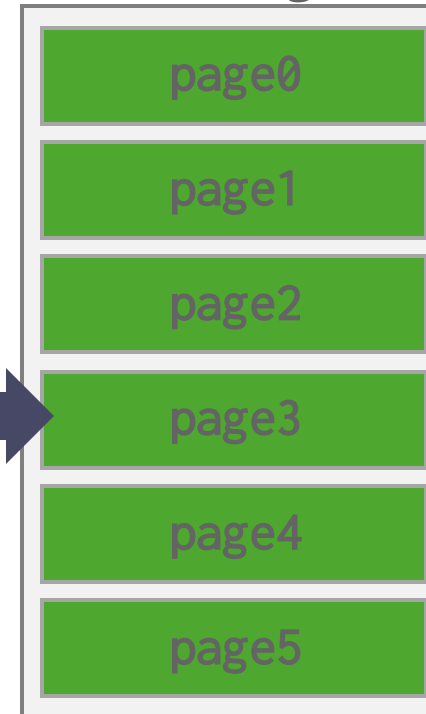
Q2 `SELECT AVG(val) FROM A`



Q2



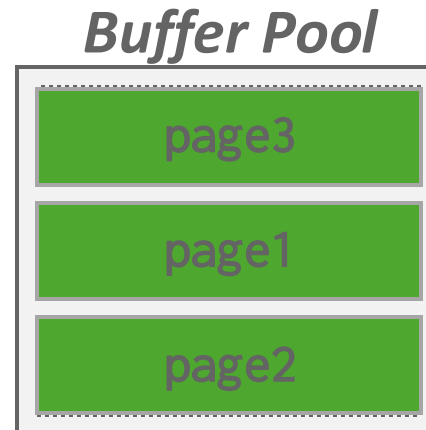
Disk Pages



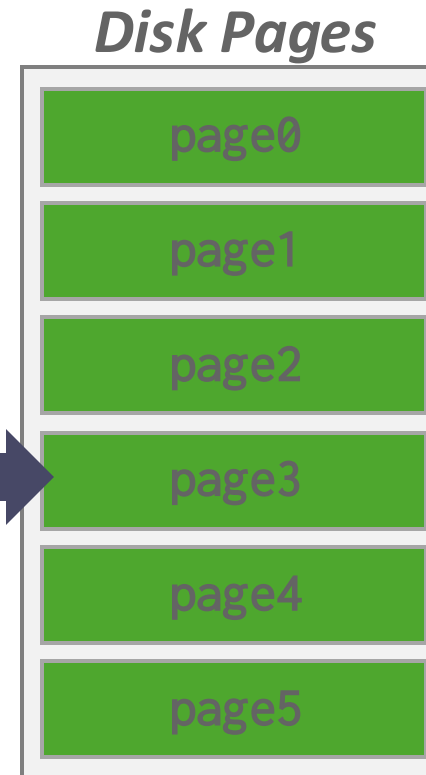
Sequential Flooding

Q1 `SELECT * FROM A WHERE primaryKey = 1`

Q2 `SELECT AVG(val) FROM A`



Q2 →

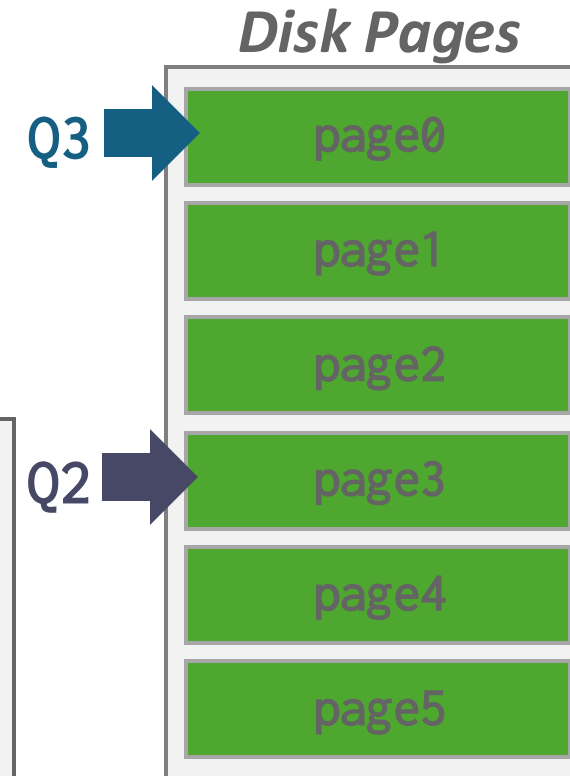
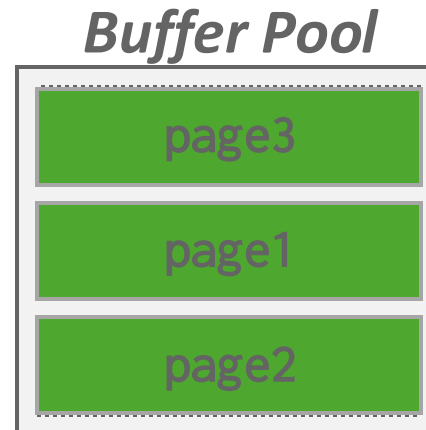


Sequential Flooding

Q1 SELECT * FROM A WHERE primKey = 1

Q2 SELECT AVG(val) FROM A

Q3 SELECT * FROM A WHERE primKey = 1

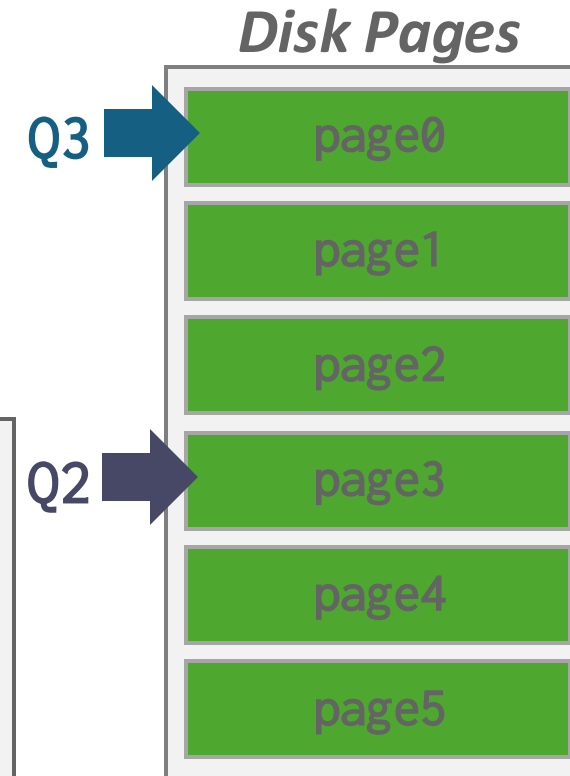
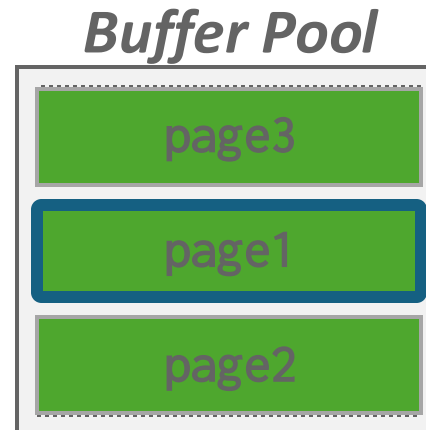


Sequential Flooding

Q1 SELECT * FROM A WHERE primKey = 1

Q2 SELECT AVG(val) FROM A

Q3 SELECT * FROM A WHERE primKey = 1



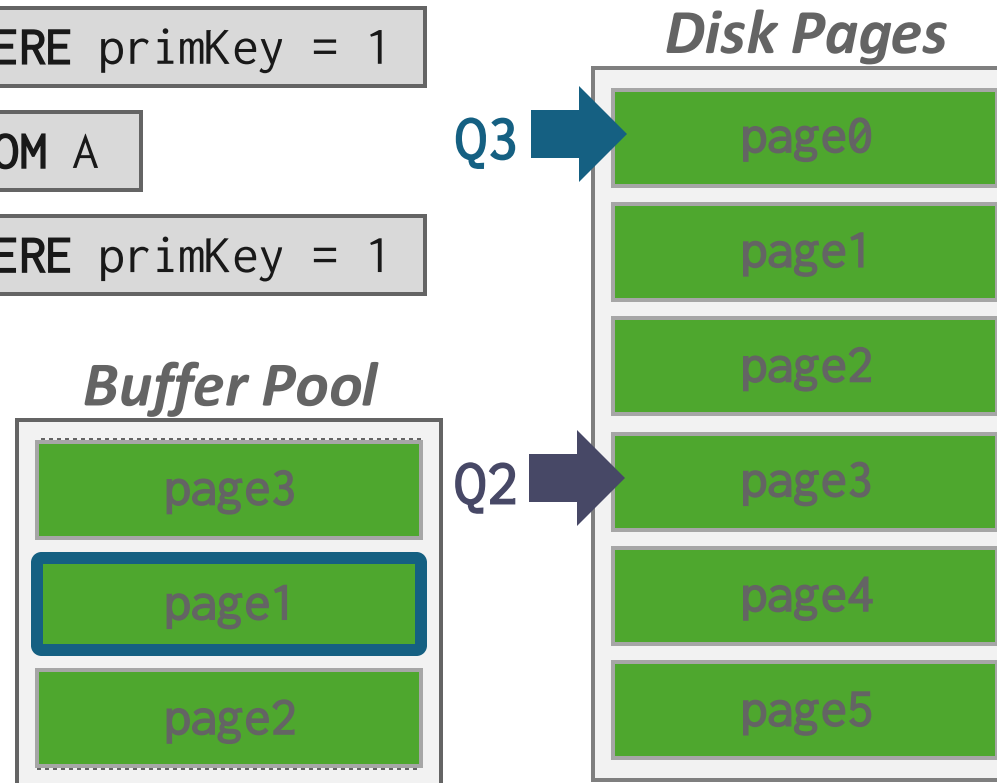
Sequential Flooding

Q1 `SELECT * FROM A WHERE primaryKey = 1`

Q2 `SELECT AVG(val) FROM A`

Q3 `SELECT * FROM A WHERE primaryKey = 1`

Sequential flooding can occur when a table is scanned multiple times within one query, such as with a Nested-Blocks Join.



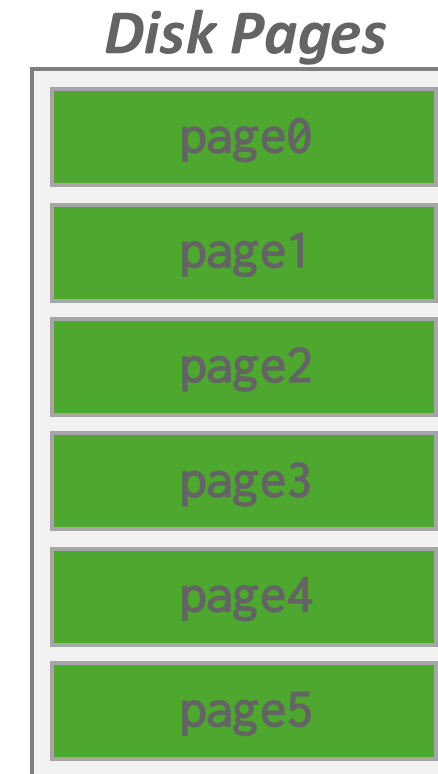
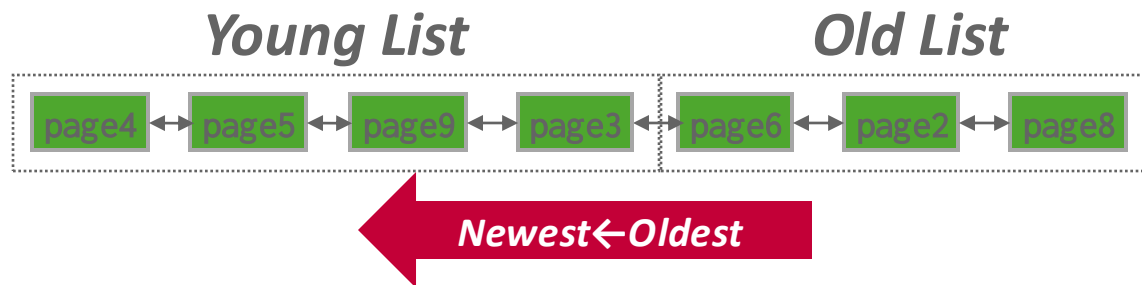
Better Policies: LRU- k

- Track the last k references to each page and compute the interval between subsequent accesses.
 - Can get fancy with distinguishing between reference types.
- The DBMS then uses this history to estimate the next time that page is going to be accessed.
 - Replace the page with the oldest “ k -th” access.
 - A balance between recency and frequency of access.
 - Maintain an ephemeral in-memory cache for recently evicted pages to prevent them from always being evicted.



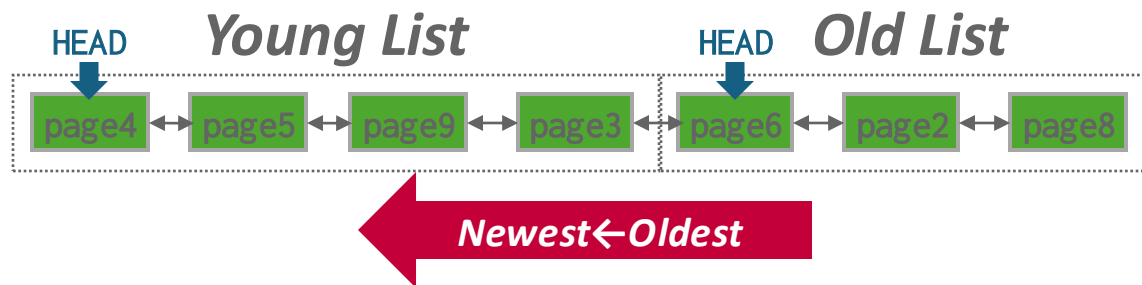
MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



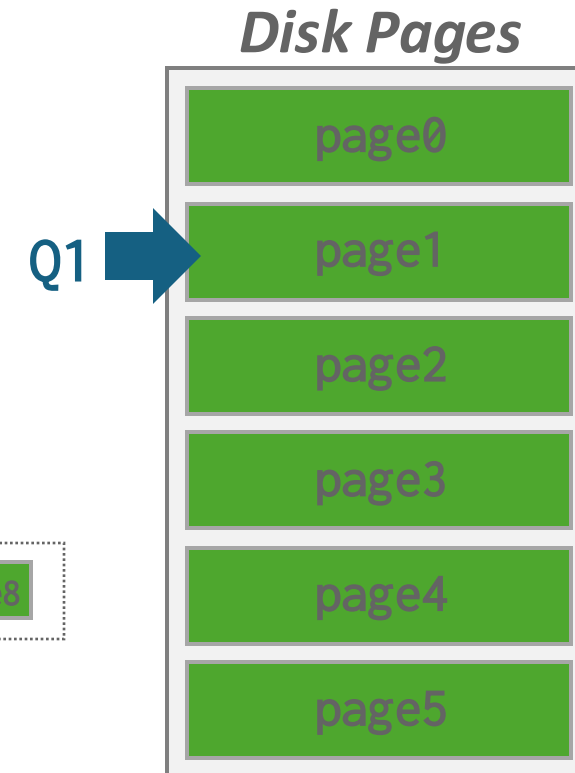
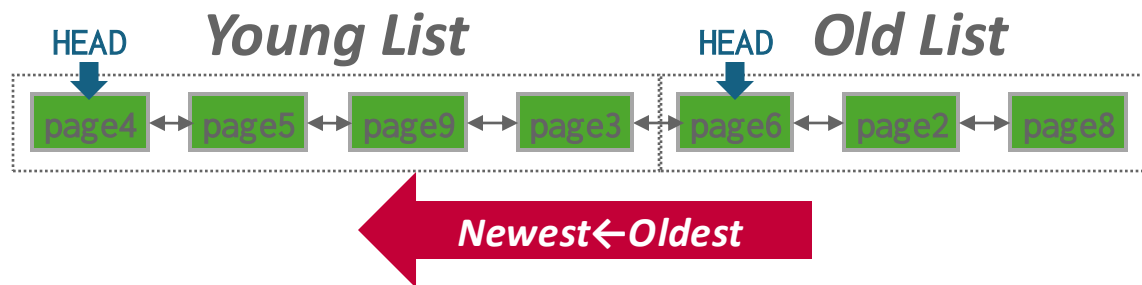
MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



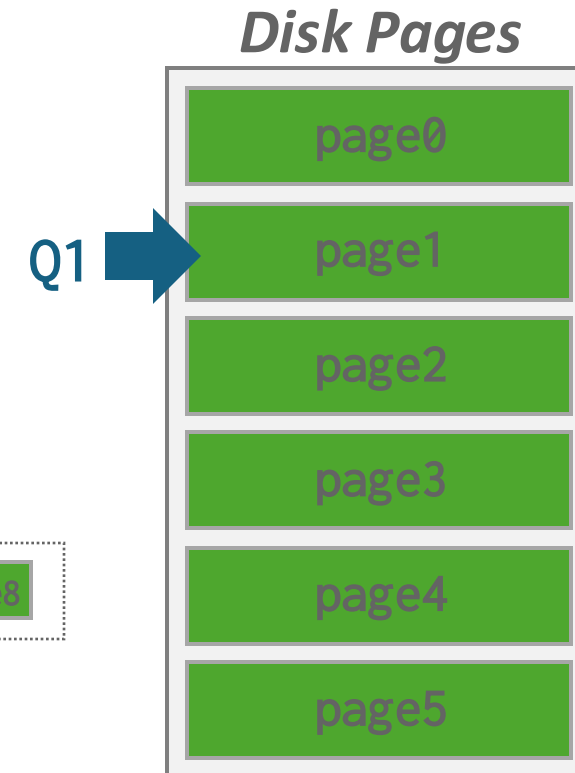
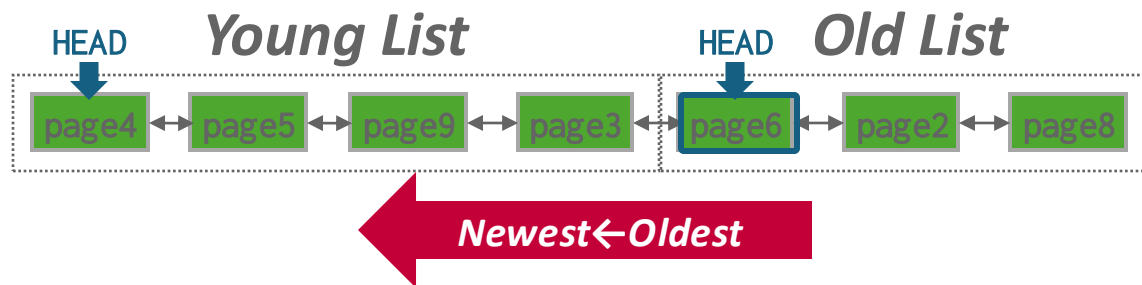
MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



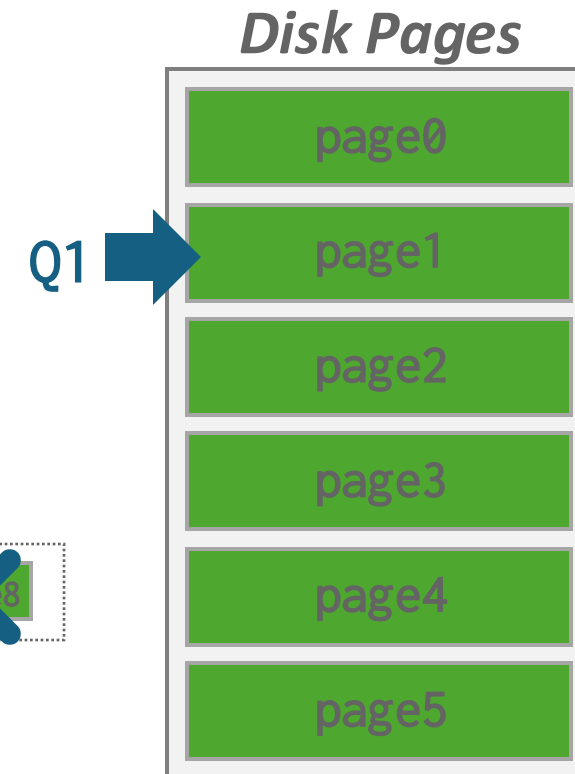
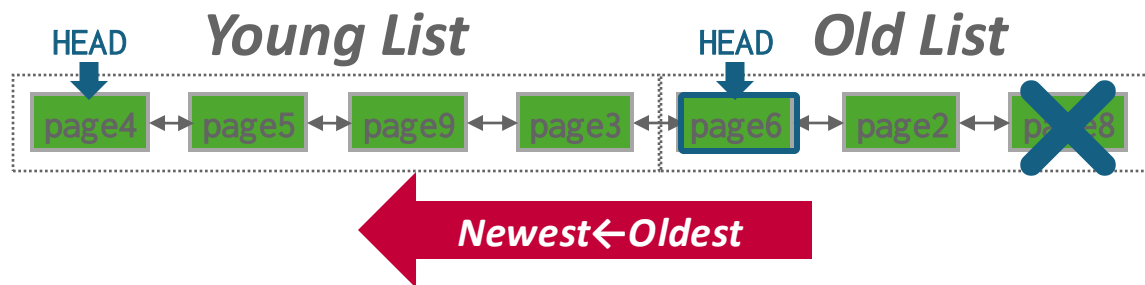
MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



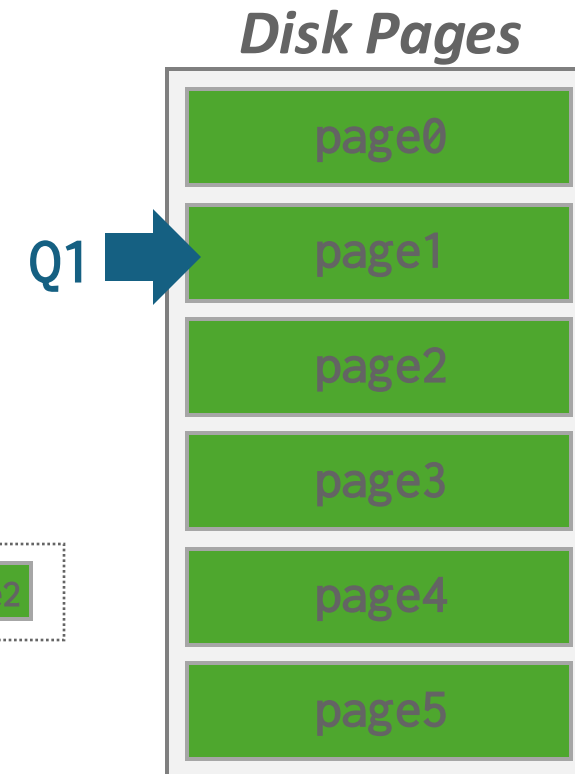
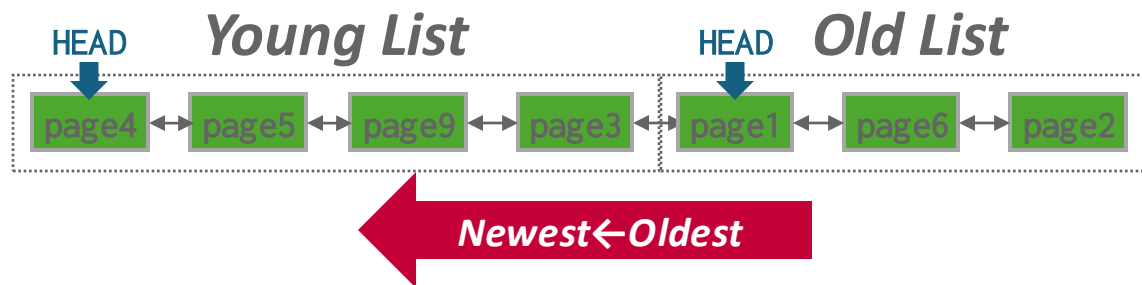
MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



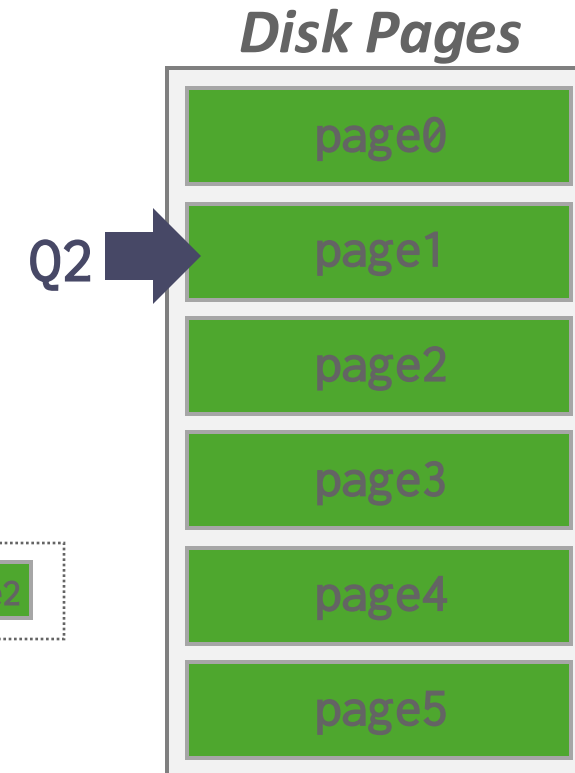
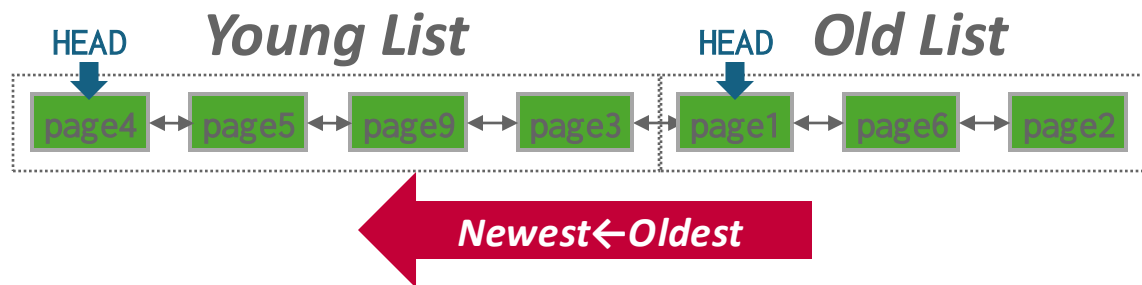
MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



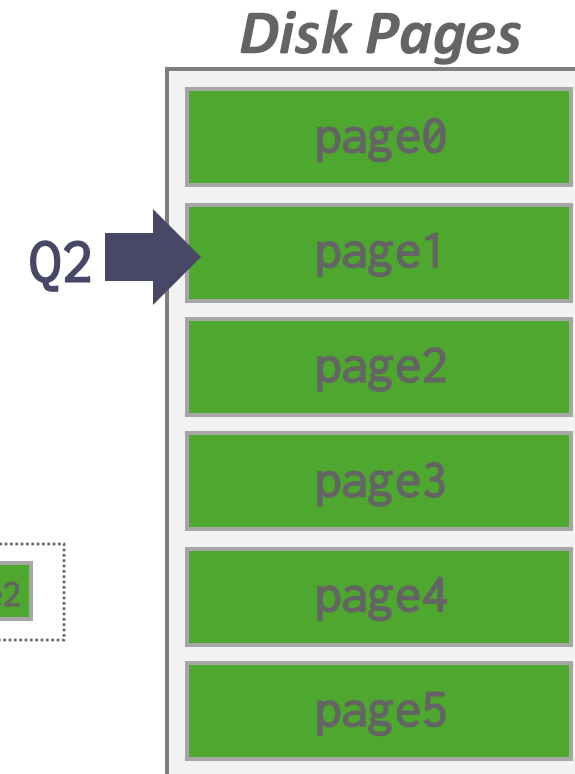
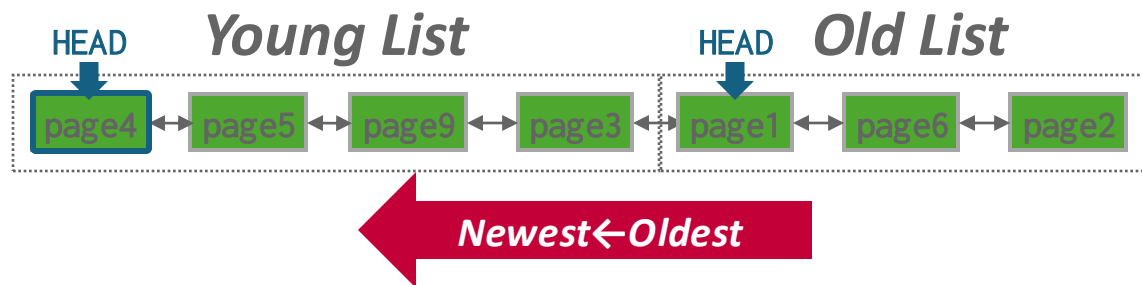
MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



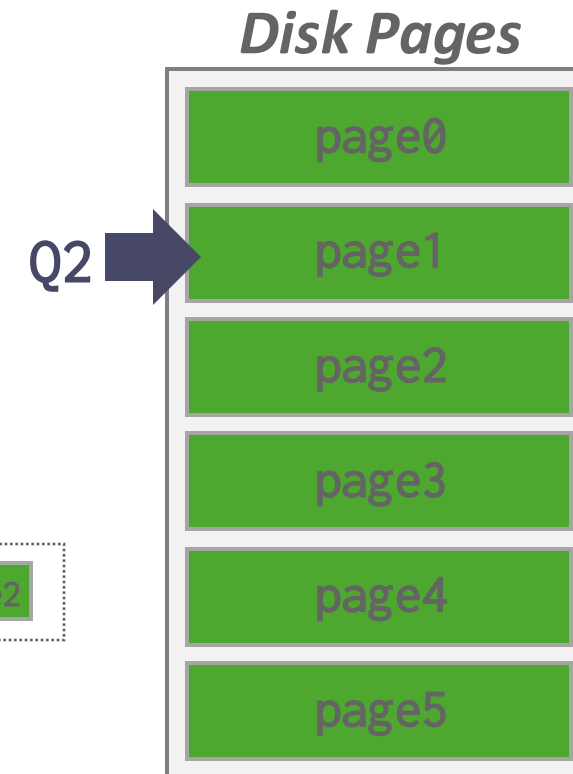
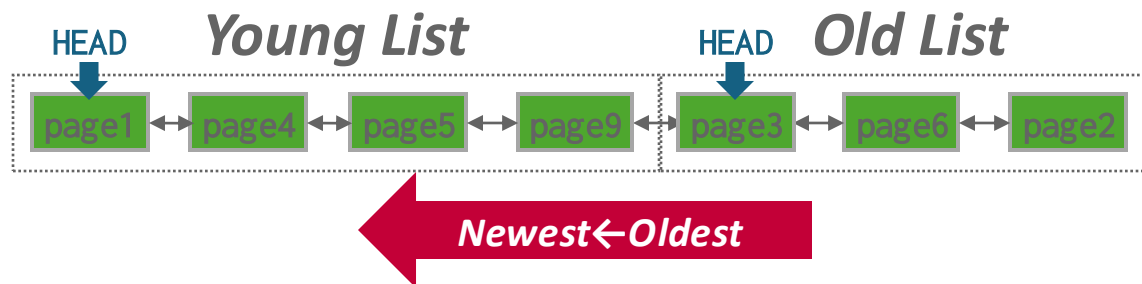
MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



MySQL Approximate LRU- k

- Single LRU linked list but with two entry points (“old” vs “young”).
 - New pages are always inserted to the head of the old list.
 - If pages in the old list is accessed again, then insert into the head of the young list.



Better Policies: Localization

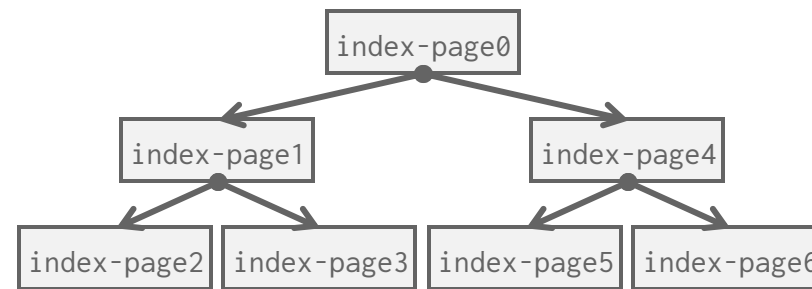
- The DBMS chooses which pages to evict on a per query basis. This minimizes the pollution of the buffer pool from each query.
 - Keep track of the pages that a query has accessed.
- Example: Postgres maintains a small ring buffer that is private to the query.

Better Policies: Priority Hints

- The DBMS knows about 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

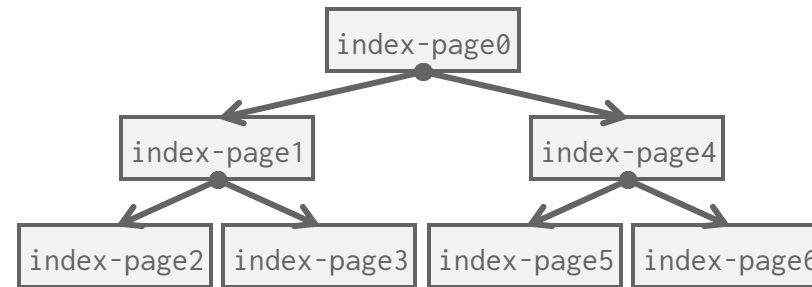
- The DBMS knows about the context of each page during query execution.
- It can provide hints to the buffer pool on whether a page is important or not.



Better Policies: Priority Hints

- The DBMS knows about the context of each page during query execution.
- It can provide hints to the buffer pool on whether a page is important or not.

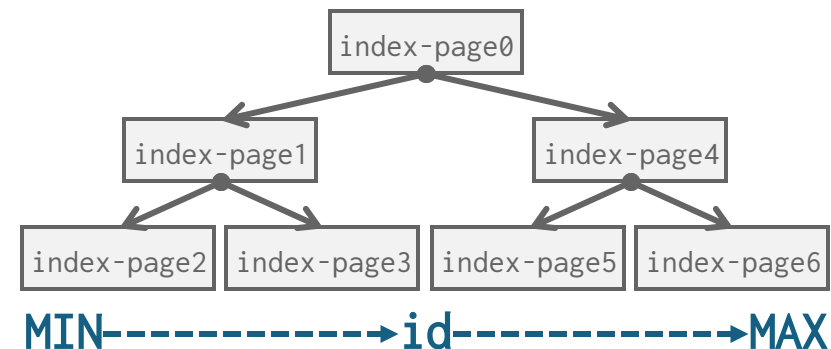
Q1 INSERT INTO A VALUES (*id++*)



Better Policies: Priority Hints

- The DBMS knows about the context of each page during query execution.
- It can provide hints to the buffer pool on whether a page is important or not.

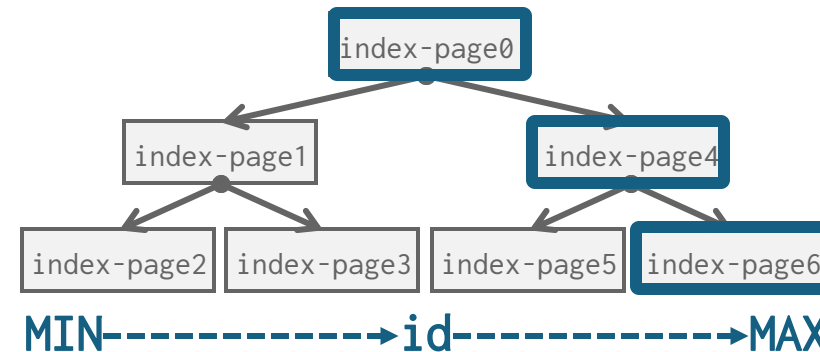
Q1 INSERT INTO A VALUES (*id++*)



Better Policies: Priority Hints

- The DBMS knows about the context of each page during query execution.
- It can provide hints to the buffer pool on whether a page is important or not.

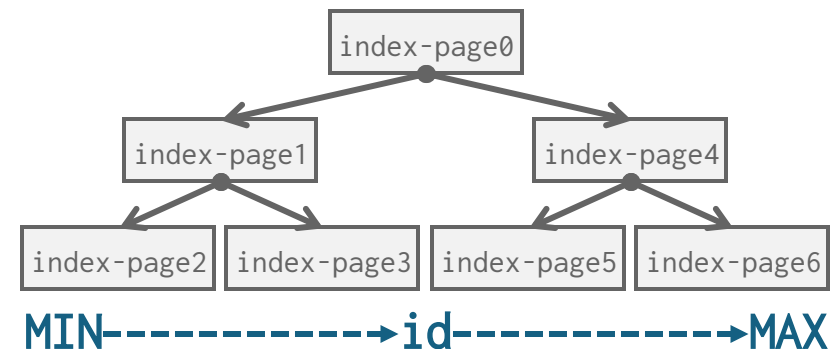
Q1 INSERT INTO A VALUES (*id++*)



Better Policies: Priority Hints

- The DBMS knows about 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++*)

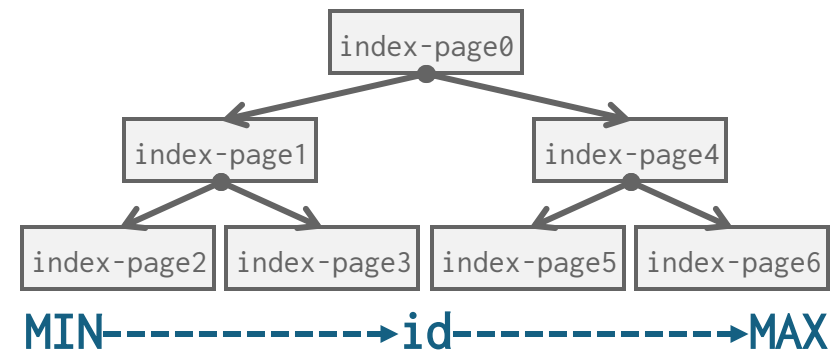


Better Policies: Priority Hints

- The DBMS knows about 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 = ?

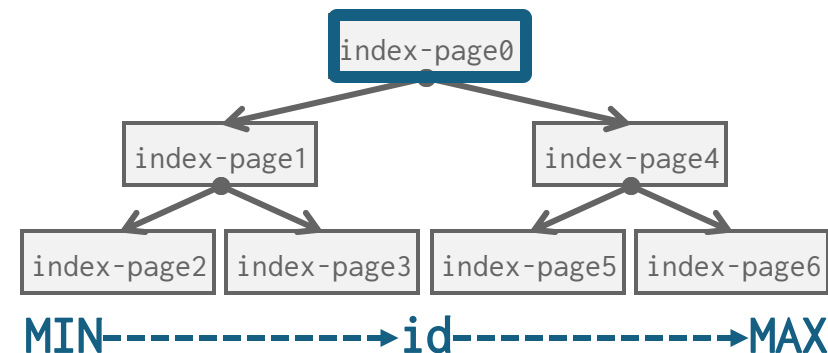


Better Policies: Priority Hints

- The DBMS knows about 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 Path:** If a page in the buffer pool is not dirty, then the DBMS can simply “drop” it.
- **Slow Path:** 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 the system doesn't write dirty pages before their log records are written...

Disk I/O Scheduling

Observation

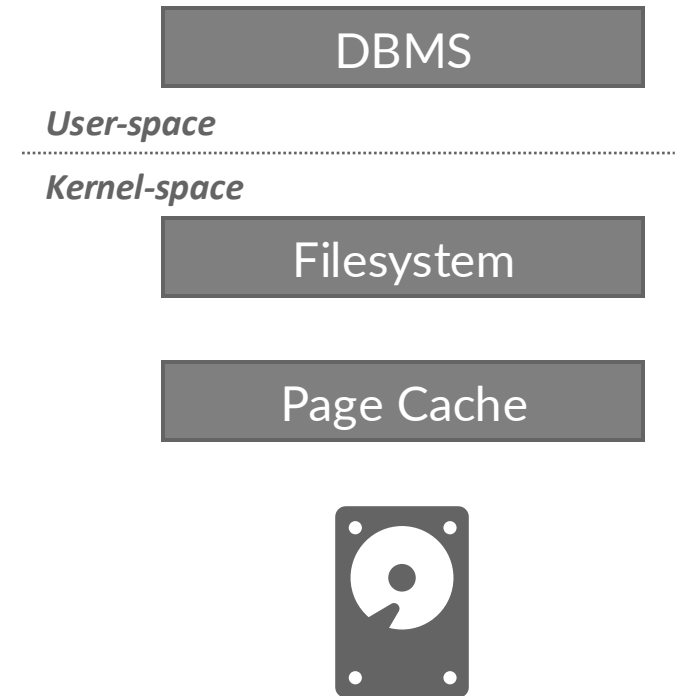
- OS/hardware tries to maximize disk bandwidth by reordering and batching I/O requests.
- But they do not know which I/O requests are more important than others.
- Many DBMSs tell you to switch Linux to use the deadline or noop (FIFO) scheduler.
 - Example: [Oracle](#), [Vertica](#), [MySQL](#)

Disk I/O Scheduling

- The DBMS maintain internal queue(s) to track page read/write requests from the entire system.
- Compute priorities based on several factors:
 - Sequential vs. Random I/O
 - Critical Path Task vs. Background Task
 - Table vs. Index vs. Log vs. Ephemeral Data
 - Transaction Information
 - User-based SLAs
- The OS doesn't know these things and is going to get into the way...

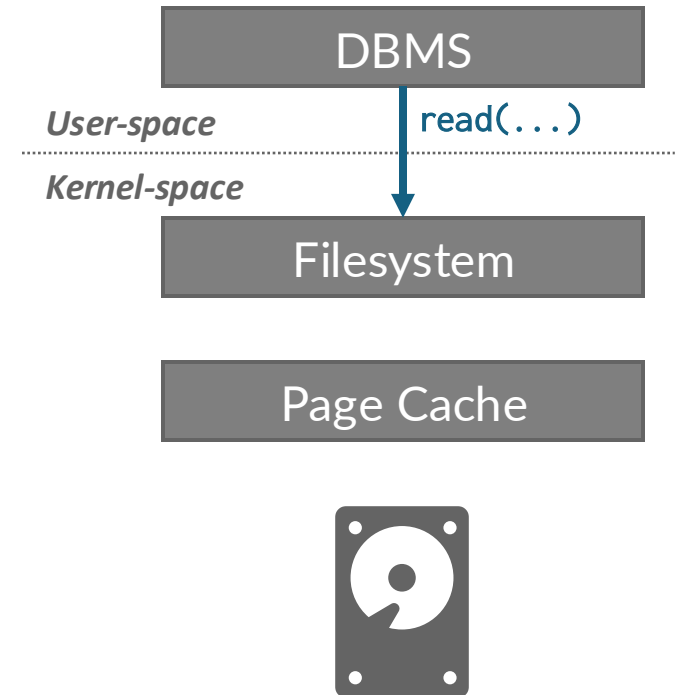
OS Page Cache

- Most disk operations go through the OS API. Unless the DBMS tells it not to, the OS maintains its own filesystem cache (aka page cache, buffer cache).
- Most DBMSs use direct I/O (**O_DIRECT**) to bypass the OS's cache.
 - Redundant copies of pages.
 - Different eviction policies.
 - Loss of control over file I/O.



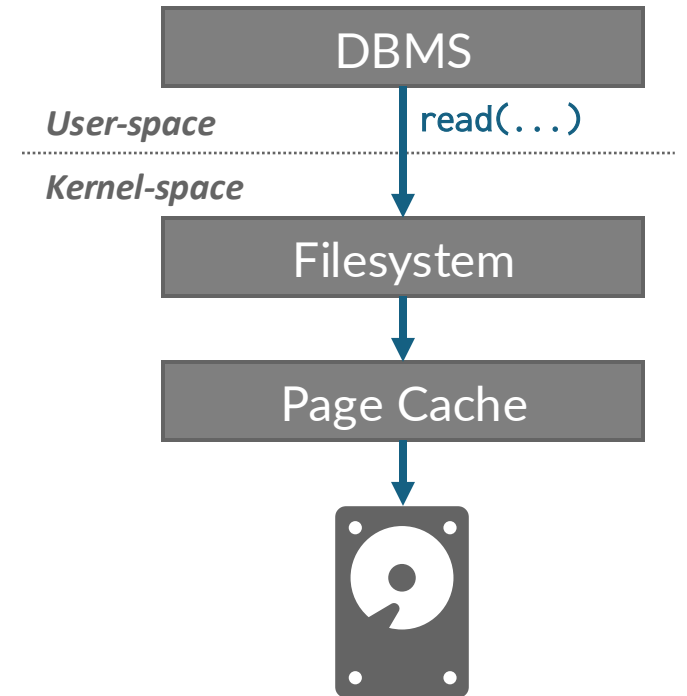
OS Page Cache

- Most disk operations go through the OS API. Unless the DBMS tells it not to, the OS maintains its own filesystem cache (aka page cache, buffer cache).
- Most DBMSs use direct I/O (**O_DIRECT**) to bypass the OS's cache.
 - Redundant copies of pages.
 - Different eviction policies.
 - Loss of control over file I/O.



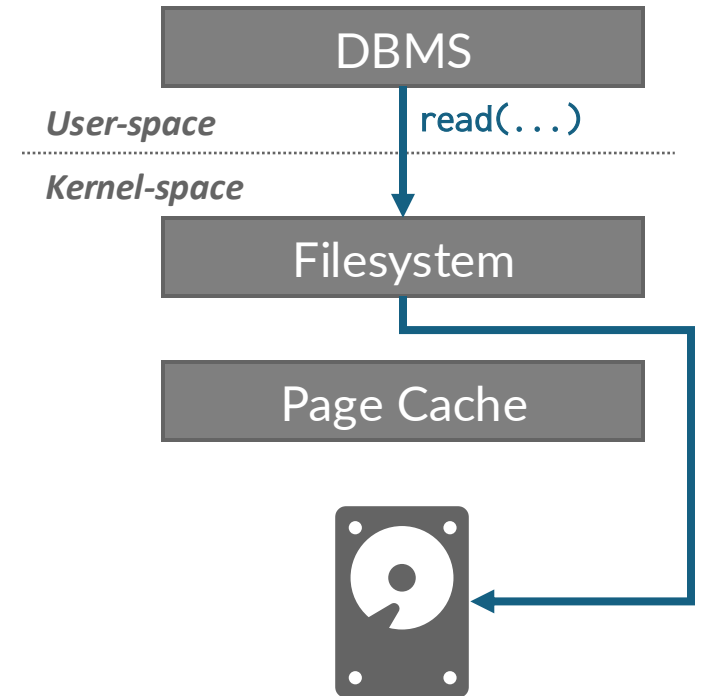
OS Page Cache

- Most disk operations go through the OS API. Unless the DBMS tells it not to, the OS maintains its own filesystem cache (aka page cache, buffer cache).
- Most DBMSs use direct I/O (`O_DIRECT`) to bypass the OS's cache.
 - Redundant copies of pages.
 - Different eviction policies.
 - Loss of control over file I/O.



OS Page Cache

- Most disk operations go through the OS API. Unless the DBMS tells it not to, the OS maintains its own filesystem cache (aka page cache, buffer cache).
- Most DBMSs use direct I/O (**O_DIRECT**) to bypass the OS's cache.
 - Redundant copies of pages.
 - Different eviction policies.
 - Loss of control over file I/O.

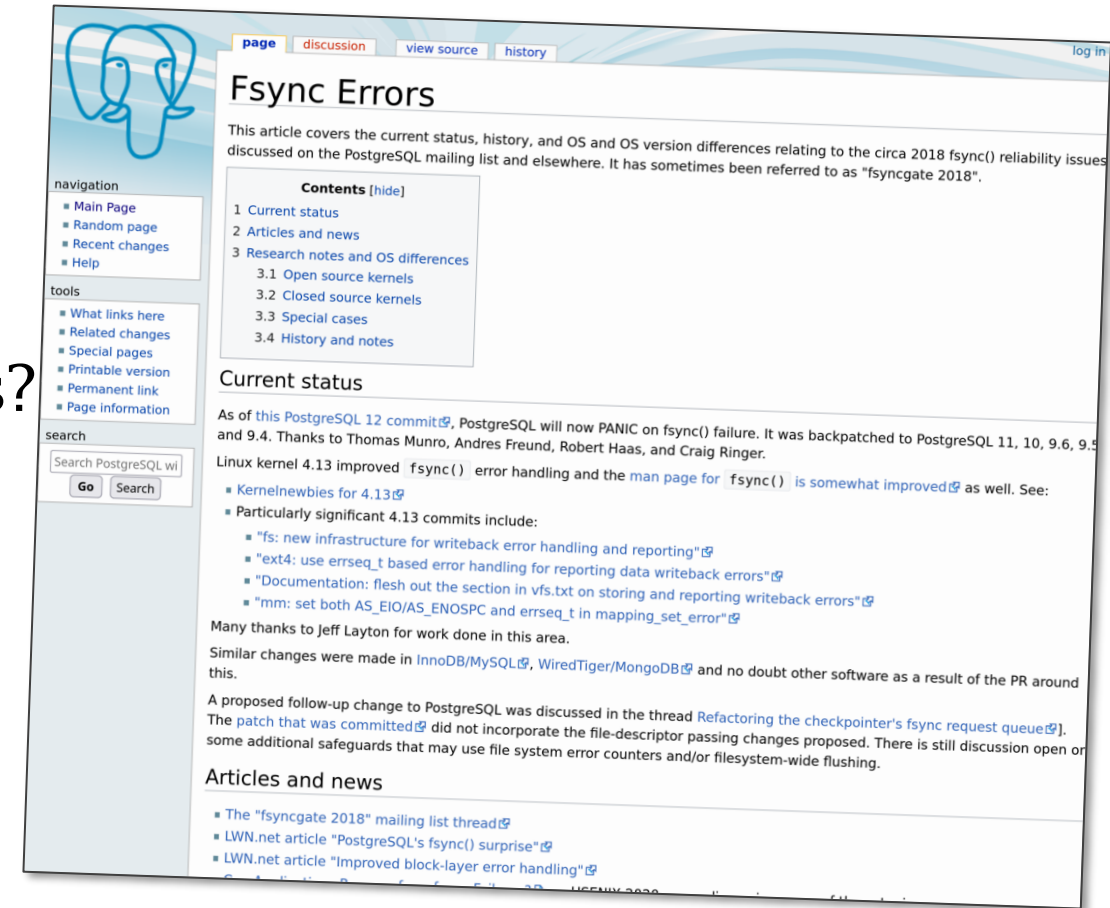


FSYNC Problems

- If the DBMS calls `fwrite`, what happens?
- If the DBMS calls `fsync`, what happens?
- If `fsync` fails (EIO), what happens?
 - Linux marks the dirty pages as clean.
 - If the DBMS calls `fsync` again, then Linux tells you that the flush was successful.

FSYNC Problems

- If the DBMS calls **fwrite**, what happens?
- If the DBMS calls **fsync**, what happens?
- If **fsync** fails (EIO), what happens?
 - Linux marks the dirty pages as clean.
 - If the DBMS calls **fsync** again, then Linux tells you that the flush was successful.



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 almost always manage memory better than the OS.
- Leverage the semantics about the query plan to make better decisions:
 - Evictions
 - Allocations
 - Pre-fetching
- Next: Hash Table