

Carnegie Mellon University

Database Systems

Memory & Disk Management



15-445/645 FALL 2024 » PROF. ANDY PAVLO

CMU-DB IAP VISIT DAY (TUE SEPT 17)

Info Session #1 (9:30-10:30am)

- DataStax: GHC 7101
- dbtLabs: GHC 7501
- Firebolt: GHC 8115

Info Session #2 (10:30-11:30am)

- ClickHouse: GHC 7101
- RelationalAI: GHC 7501
- StarTree: GHC 8115

Info Sessions #3 (11:30-12:30pm)

- Neon: GHC 7101
- PingCAP TiDB: GHC 7501
- Weaviate: GHC 8115

**Carnegie
Mellon
University**

Database Group
Industry Affiliates

<https://db.cs.cmu.edu/affiliates/visit2024>

LAST CLASS

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

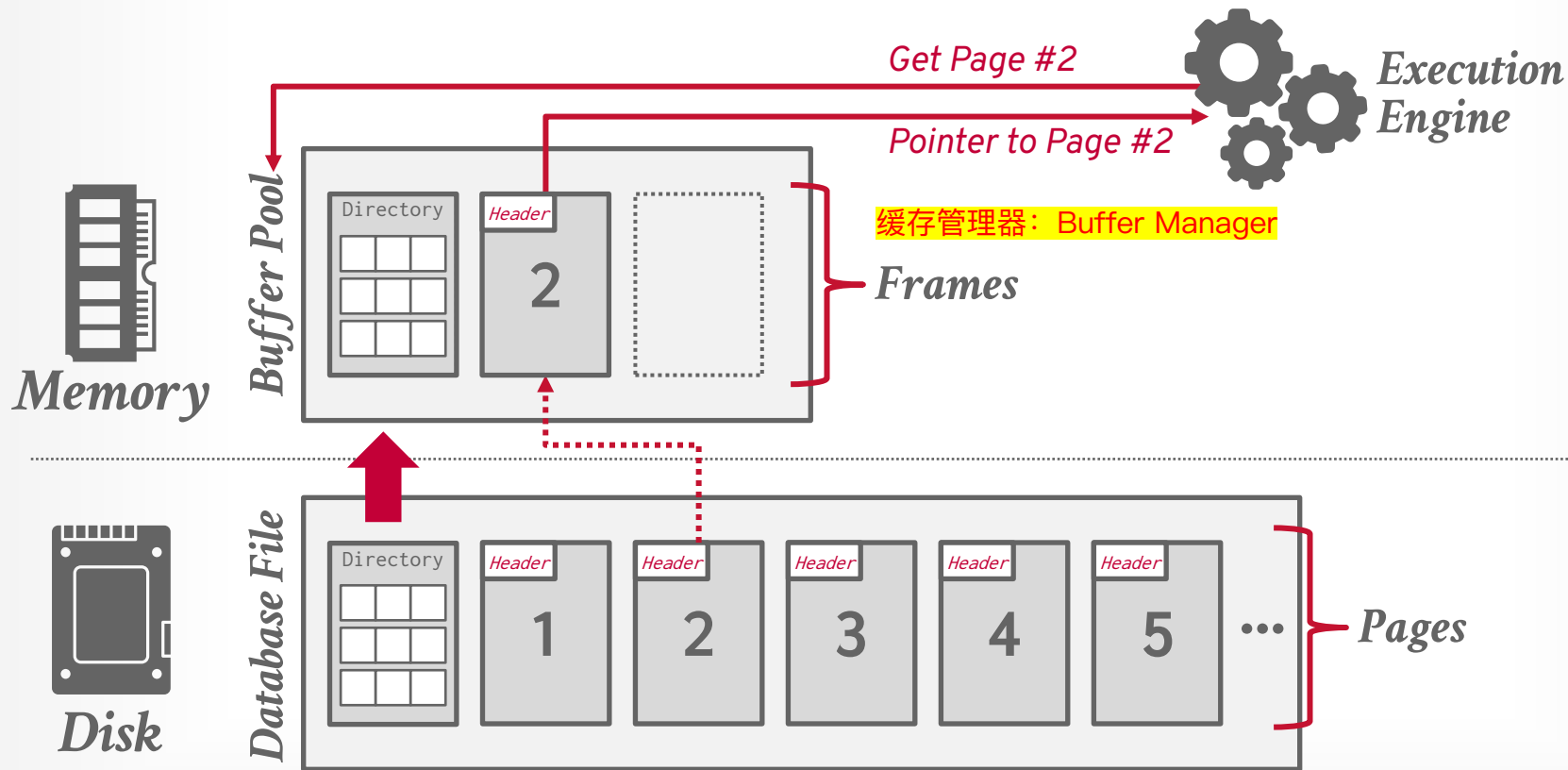
Spatial Control:

- 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.

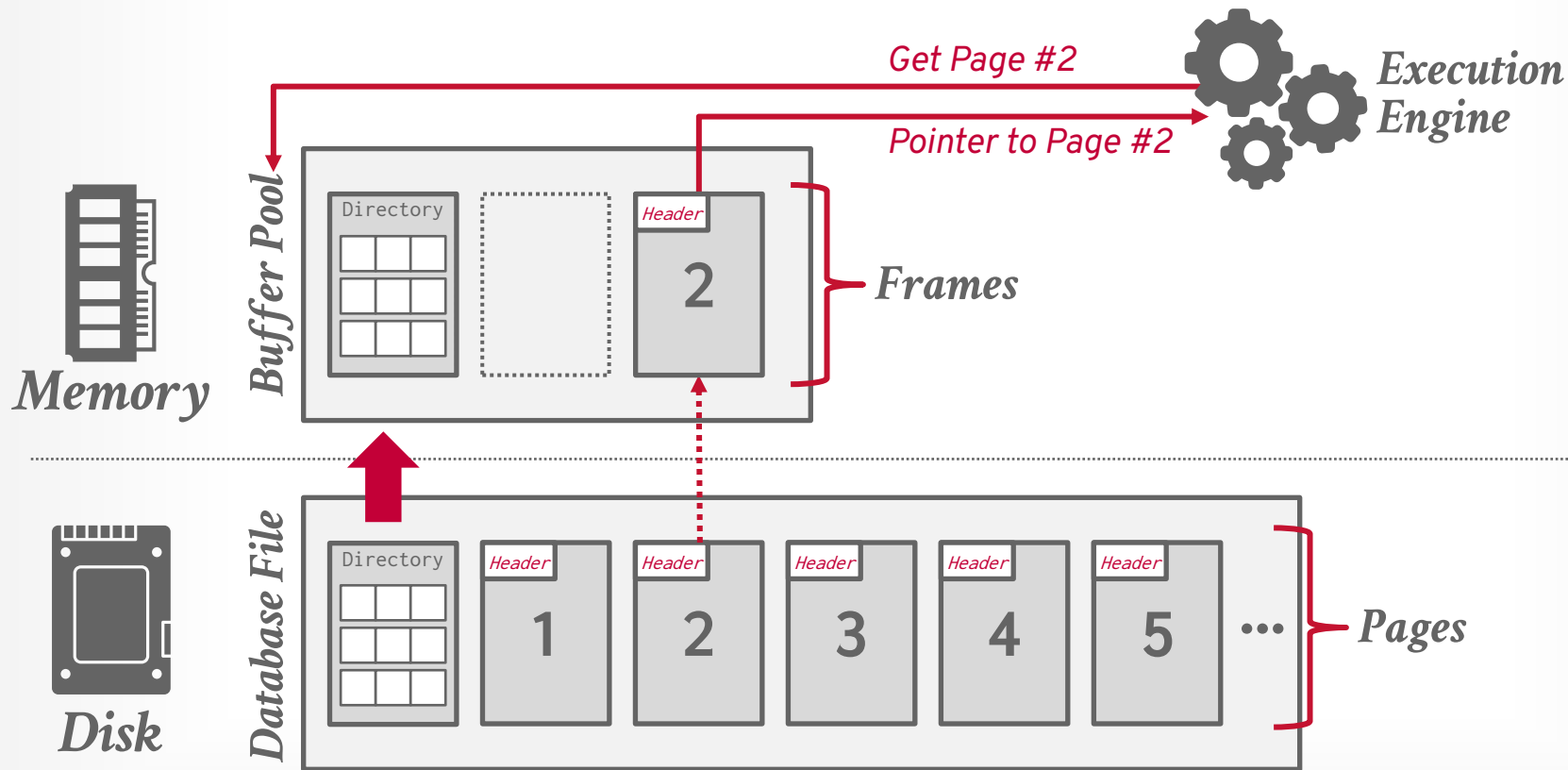
Temporal Control:

- 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



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. 有些缓存区的数据不需要持久化落盘.

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

TODAY'S AGENDA

Buffer Pool Manager

Why MMAP Will Murder Your DBMS

Disk I/O Scheduling

Replacement Policies

Other Memory Pools

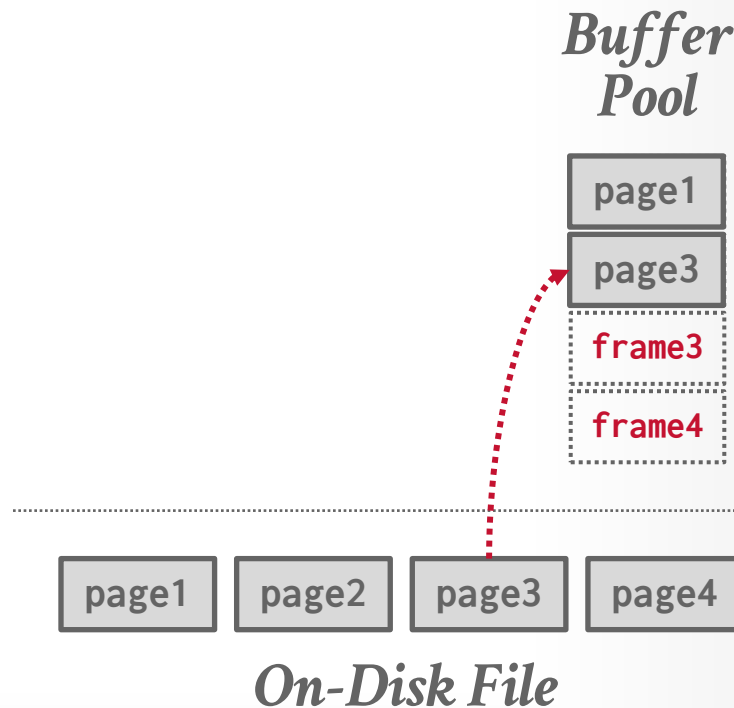
BUFFER POOL ORGANIZATION

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.

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



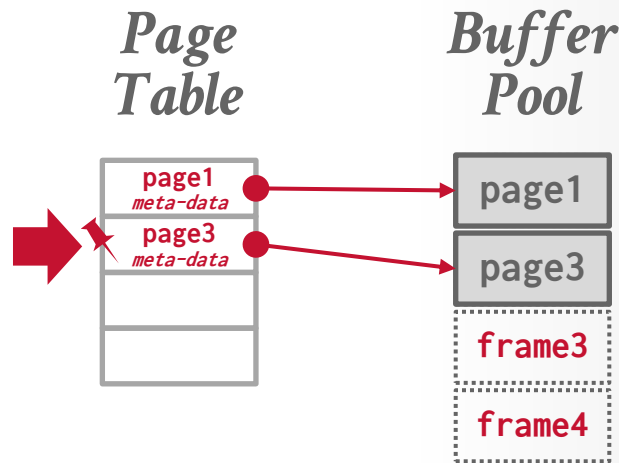
BUFFER POOL META-DATA

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**



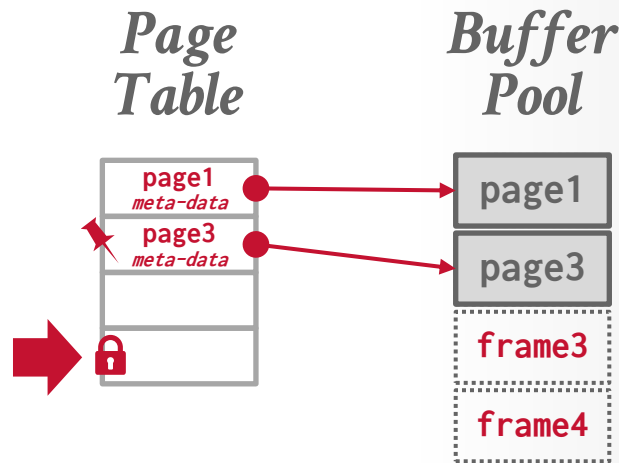
BUFFER POOL META-DATA

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**



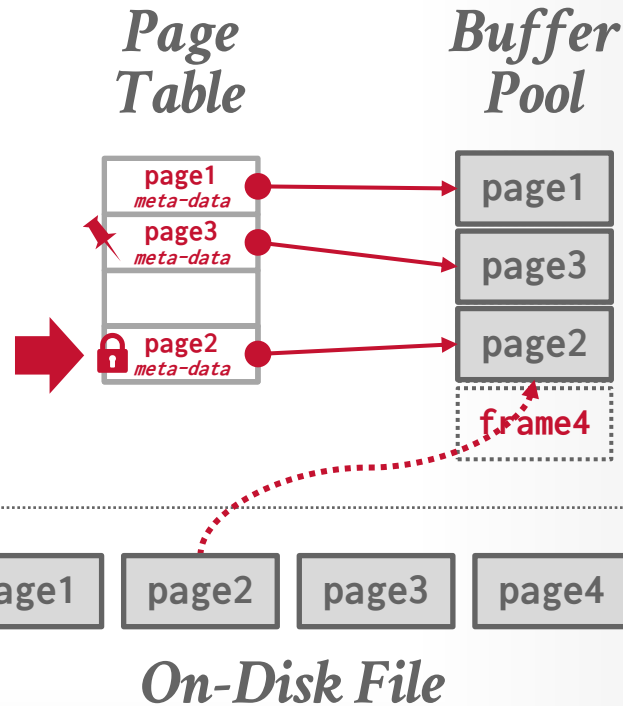
BUFFER POOL META-DATA

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**



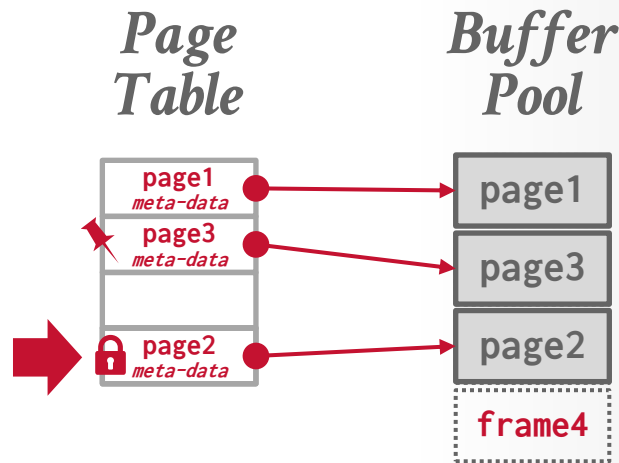
BUFFER POOL META-DATA

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**



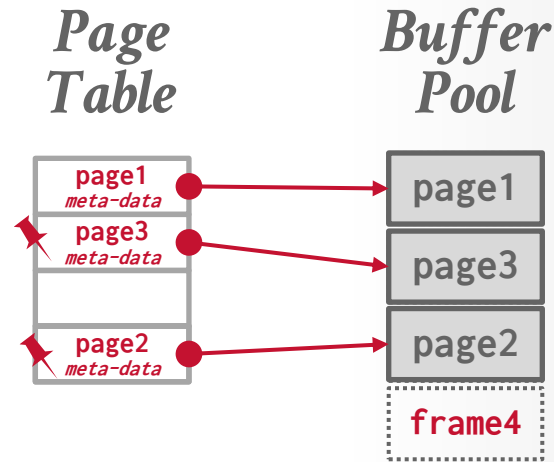
BUFFER POOL META-DATA

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**



LOCKS VS. LATCHES

Locks:

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

Latches:

保护关键区域的互斥锁.

- 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

The **page directory** is the mapping from page ids to page locations in the database files. 存储在磁盘的数据结构.

→ 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.

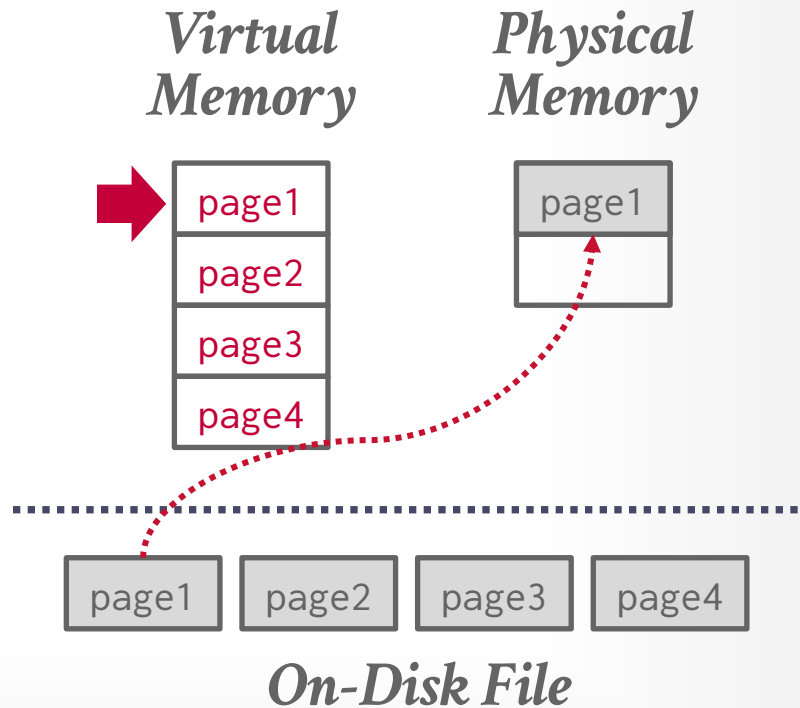
→ This is an in-memory data structure that does not need to be stored on disk.

WHY NOT USE THE OS?

Use OS memory mapping (**mmap**) to store the contents of a file into the address space of a program.

OS is responsible for moving file pages in and out of memory, so the DBMS doesn't need to worry about it.

What if DBMS allows multiple threads to access **mmap** files to hide page fault stalls?

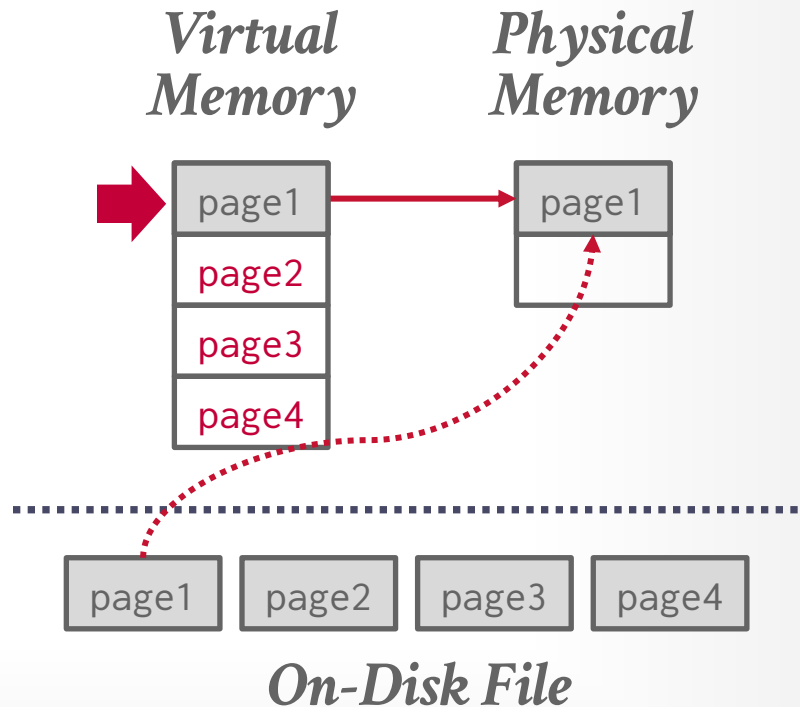


WHY NOT USE THE OS?

Use OS memory mapping (**mmap**) to store the contents of a file into the address space of a program.

OS is responsible for moving file pages in and out of memory, so the DBMS doesn't need to worry about it.

What if DBMS allows multiple threads to access **mmap** files to hide page fault stalls?

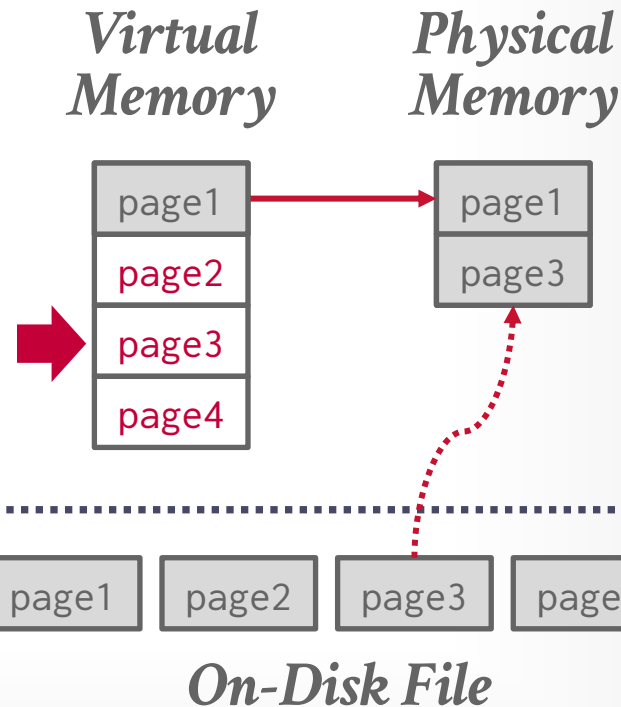


WHY NOT USE THE OS?

Use OS memory mapping (**mmap**) to store the contents of a file into the address space of a program.

OS is responsible for moving file pages in and out of memory, so the DBMS doesn't need to worry about it.

What if DBMS allows multiple threads to access **mmap** files to hide page fault stalls?

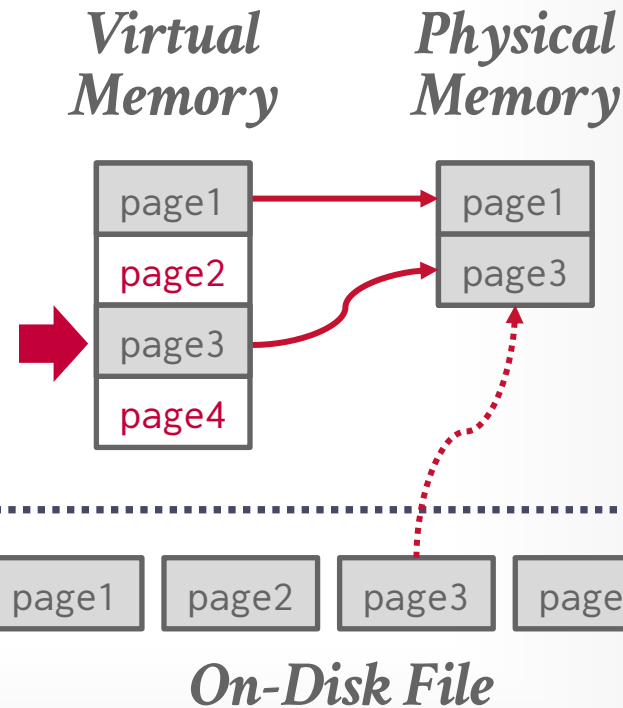


WHY NOT USE THE OS?

Use OS memory mapping (**mmap**) to store the contents of a file into the address space of a program.

OS is responsible for moving file pages in and out of memory, so the DBMS doesn't need to worry about it.

What if DBMS allows multiple threads to access **mmap** files to hide page fault stalls?

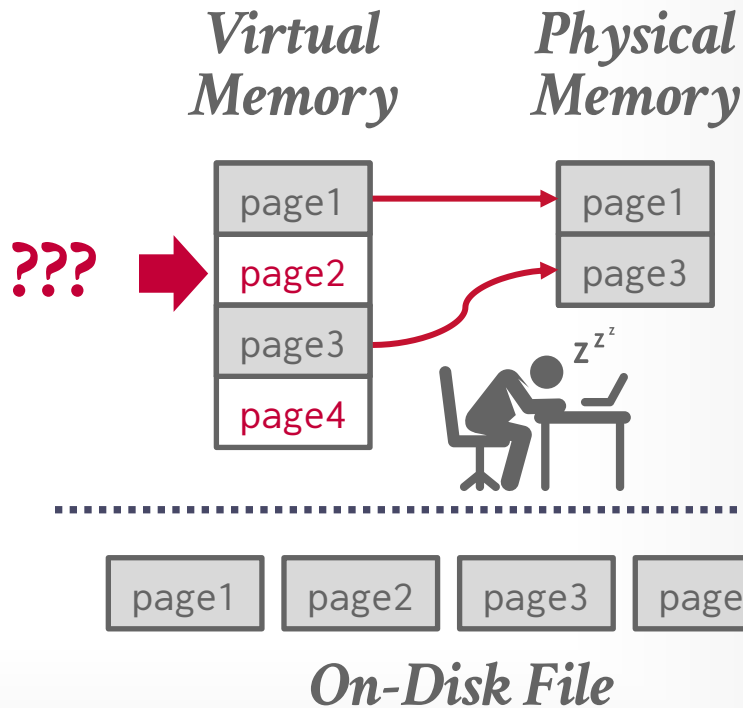


WHY NOT USE THE OS?

Use OS memory mapping (**mmap**) to store the contents of a file into the address space of a program.

OS is responsible for moving file pages in and out of memory, so the DBMS doesn't need to worry about it.

What if DBMS allows multiple threads to access **mmap** files to hide page fault stalls?



MEMORY MAPPED I/O PROBLEMS

Problem #1: Transaction Safety

→ OS can flush dirty pages at any time.

Problem #2: I/O Stalls

→ DBMS doesn't know which pages are in memory. The OS will stall a thread on page fault.

Problem #3: Error Handling

→ Difficult to validate pages. Any access can cause a **SIGBUS** that the DBMS must handle.

Problem #4: Performance Issues

→ OS data structure contention. TLB shootdowns.

WHY NOT USE THE OS?

There are some solutions to some of these problems:

- **advise**: Tell the OS how you expect to read certain pages.
- **lock**: Tell the OS that memory ranges cannot be paged out.
- **sync**: Tell the OS to flush memory ranges out to disk.

Using these syscalls to get the OS to behave correctly is just as onerous as managing memory yourself.

Full Usage



Weaviate

Partial Usage



MongoDB®



SingleStore



SQLite



influxdb

WHY NOT USE THE OS?

There are some solutions to some of these problems:

- **madvise**: Tell the OS how you expect to read certain pages.
- **mlock**: Tell the OS that memory ranges cannot be paged out.
- **msync**: Tell the OS to flush memory ranges out to disk.

Using these syscalls to get the OS to behave correctly is just as onerous as managing memory yourself.

Full Usage



Partial Usage



WHY NOT USE THE OS?

DBMS (almost) always wants to control things itself and can do a better job than the OS.

- Flushing dirty pages to disk in the correct order.
- Specialized prefetching.
- Buffer replacement policy.
- Thread/process scheduling.

The OS is **not** your friend.

WHY NOT USE

DBMS (almost) always wants to
and can do a better job than the

- Flushing dirty pages to disk in the
- Specialized prefetching.
- Buffer replacement policy.
- Thread/process scheduling.

The OS is **not** your friend.

Are You Sure You Want to Use MMAP in Your Database Management System?

Andrew Crotty
Carnegie Mellon University
andrewcr@cs.cmu.edu

Viktor Leis
University of Erlangen-Nuremberg
viktor.leis@fau.de

Andrew Pavlo
Carnegie Mellon University
pavlo@cs.cmu.edu

ABSTRACT

Memory-mapped (mmap) file I/O is an OS-provided feature that maps the contents of a file on secondary storage into a program's address space. The program then accesses pages via pointers as if the file resided entirely in memory. The OS transparently loads pages only when the program references them and automatically evicts pages if memory fills up.

mmap's perceived ease of use has seduced database management system (DBMS) developers for decades as a viable alternative to implementing a buffer pool. There are, however, severe correctness and performance issues with mmap that are not immediately apparent. Such problems make it difficult, if not impossible, to use mmap correctly and efficiently in a modern DBMS. In fact, several popular DBMSs initially used mmap to support larger-than-memory databases but soon encountered these hidden perils, forcing them to switch to managing file I/O themselves after significant engineering costs. In this way, mmap and DBMSs are like coffee and spicy food: an unfortunate combination that becomes obvious after the fact.

Since developers keep trying to use mmap in new DBMSs, we wrote this paper to provide a warning to others that mmap is not a suitable replacement for a traditional buffer pool. We discuss the main shortcomings of mmap in detail, and our experimental analysis demonstrates clear performance limitations. Based on these findings, we conclude with a prescription for when DBMS developers might consider using mmap for file I/O.

1 INTRODUCTION

An important feature of disk-based DBMSs is their ability to support databases that are larger than the available physical memory. This functionality allows a user to query a database as if it resided entirely in memory, even if it does not fit all at once. DBMSs achieve this illusion by reading pages of data from secondary storage (e.g., HDD, SSD) into memory on demand. If there is not enough memory for a new page, the DBMS will evict an existing page that is no longer needed in order to make room.

Traditionally, DBMSs implement the movement of pages between secondary storage and memory in a buffer pool, which interacts with secondary storage using system calls like read and write. These file I/O mechanisms copy data to and from a buffer in user space, with the DBMS maintaining complete control over how and when it transfers pages.

Alternatively, the DBMS can relinquish the responsibility of data movement to the OS, which maintains its own file mapping and

page cache. The POSIX mmap system call maps a file on secondary storage into the virtual address space of the caller (i.e., the DBMS), and the OS will then load pages lazily when the DBMS accesses them. To the DBMS, the database appears to reside fully in memory, but the OS handles all necessary paging behind the scenes rather than the DBMS's buffer pool.

On the surface, mmap seems like an attractive implementation option for managing file I/O in a DBMS. The most notable benefits are ease of use and low engineering cost. The DBMS no longer needs to track which pages are in memory, nor does it need to track how often pages are accessed or which pages are dirty. Instead, the DBMS can simply access disk-resident data via pointers as if it were accessing data in memory while leaving all low-level page management to the OS. If the available memory fills up, then the OS will free space for new pages by transparently evicting (ideally unneeded) pages from the page cache.

From a performance perspective, mmap should also have much lower overhead than a traditional buffer pool. Specifically, mmap does not incur the cost of explicit system calls (i.e., read/write) and avoids redundant copying to a buffer in user space because the DBMS can access pages directly from the OS page cache.

Since the early 1980s, these supposed benefits have enticed DBMS developers to forgo implementing a buffer pool and instead rely on the OS to manage file I/O [36]. In fact, the developers of several well-known DBMSs (see Section 2.3) have gone down this path, with some even touting mmap as a key factor in achieving good performance [20].

Unfortunately, mmap has a hidden dark side with many subtle problems that make it undesirable for file I/O in a DBMS. As we describe in this paper, these problems involve both data safety and system performance concerns. We contend that the engineering steps required to overcome them negate the purported simplicity of working with mmap. For these reasons, we believe that mmap adds too much complexity with no commensurate performance benefit and strongly urge DBMS developers to avoid using mmap as a replacement for a traditional buffer pool.

The remainder of this paper is organized as follows. We begin with a short background on mmap (Section 2), followed by a discussion of its main problems (Section 3) and our experimental analysis (Section 4). We then discuss related work (Section 5) and conclude with a summary of our guidance for when you might consider using mmap in your DBMS (Section 6).

2 BACKGROUND

This section provides the relevant background on mmap. We begin with a high-level overview of memory-mapped file I/O and the POSIX mmap API. Then, we discuss real-world implementations of mmap-based systems.

<https://db.cs.cmu.edu/mmap-cidr2022>

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.

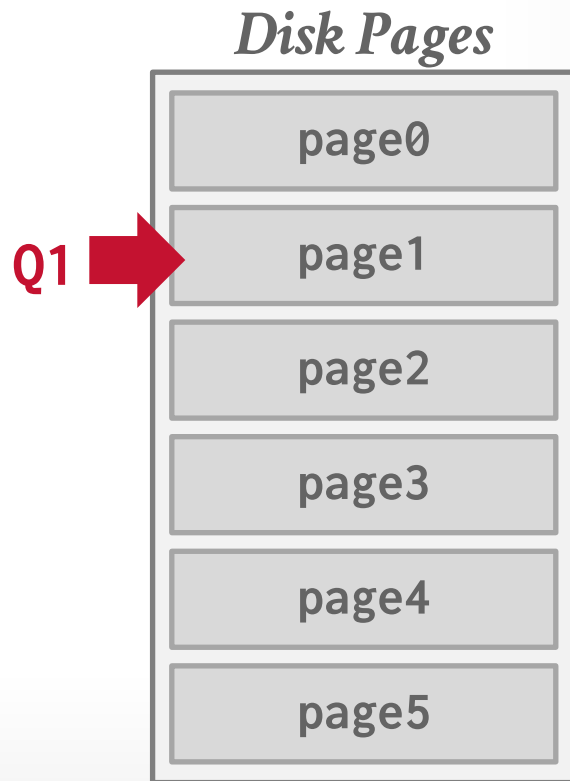
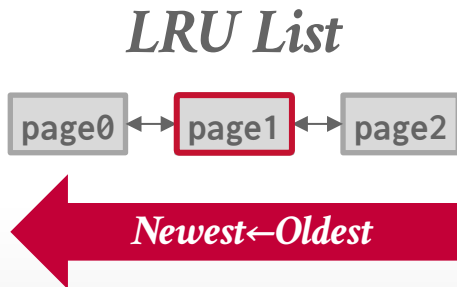
Goals:

- Correctness
- Accuracy
- Speed
- Meta-data overhead

LEAST-RECENTLY USED

Maintain a single timestamp of when each page was last accessed. When the DBMS needs to evict a page, select the one with the oldest timestamp.

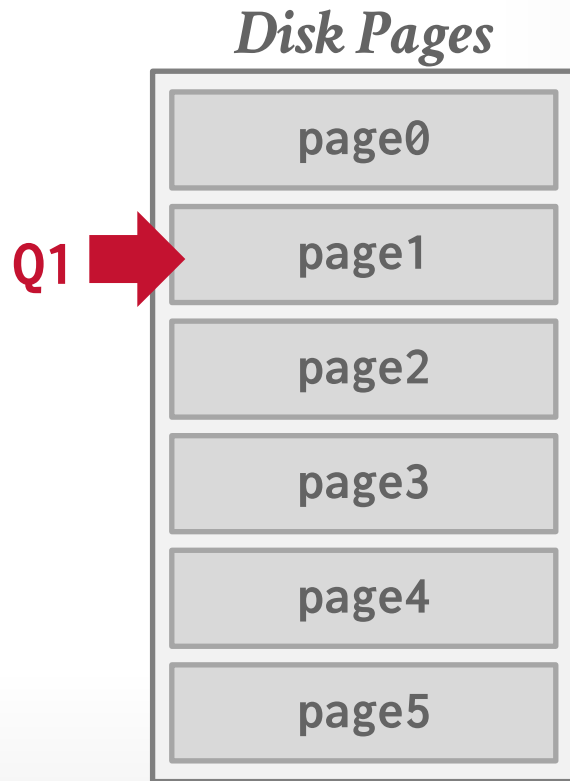
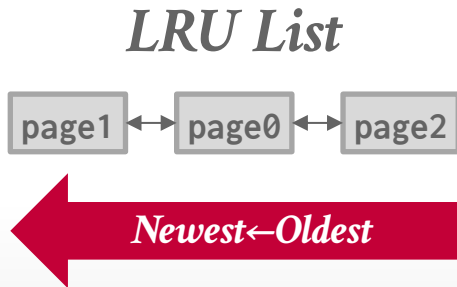
→ Keep the pages in sorted order to reduce the search time on eviction.



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.

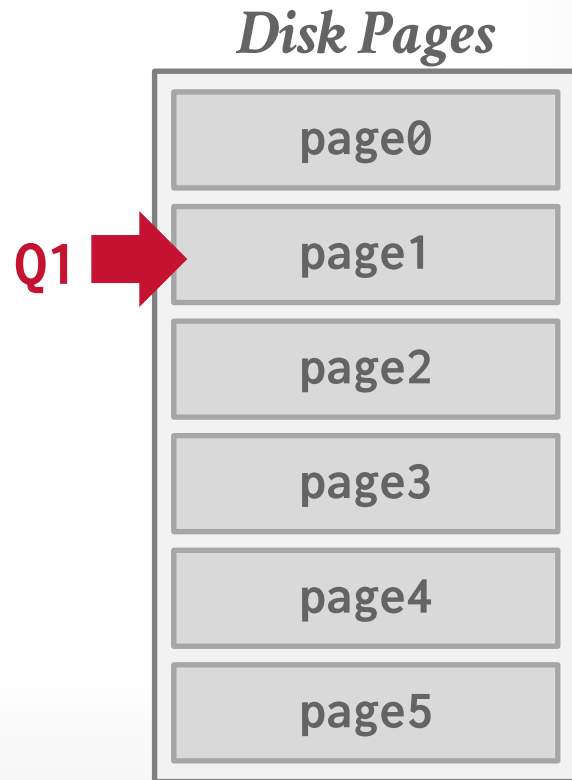
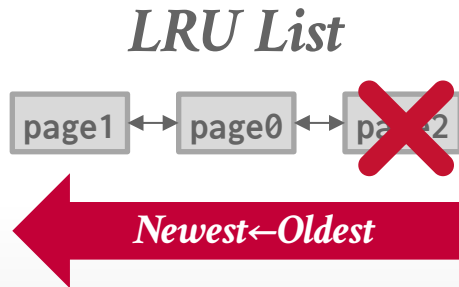
→ Keep the pages in sorted order to reduce the search time on eviction.



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

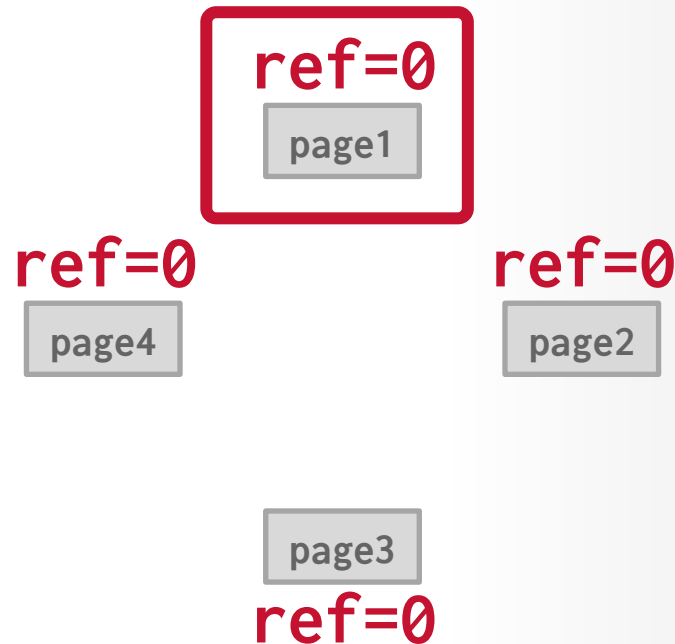
Approximation of LRU that does not need a separate timestamp per page.

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

表示上次检查后某个页是否被修改.

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



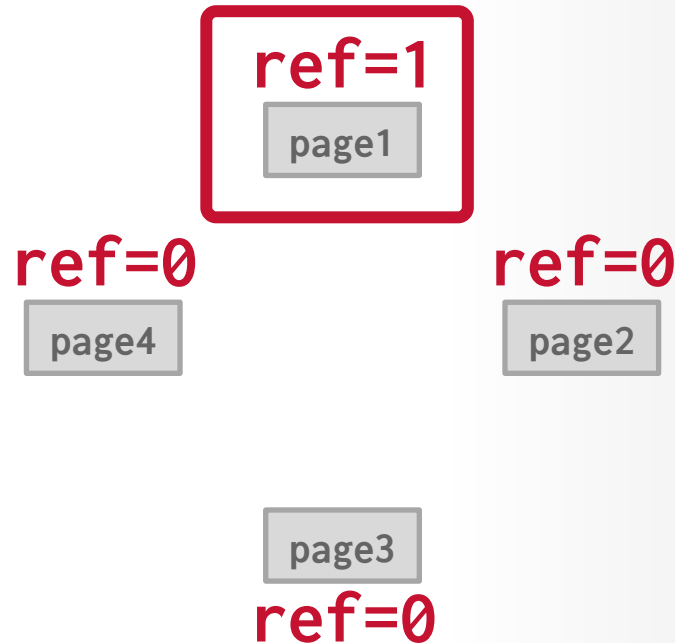
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



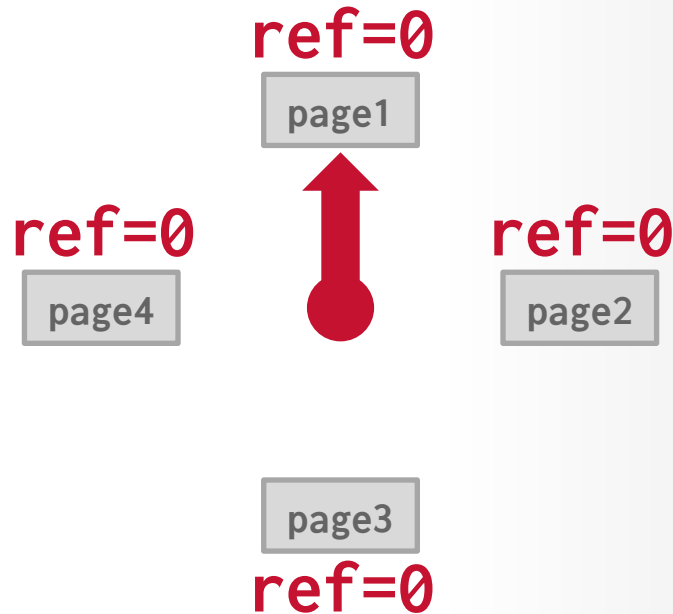
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



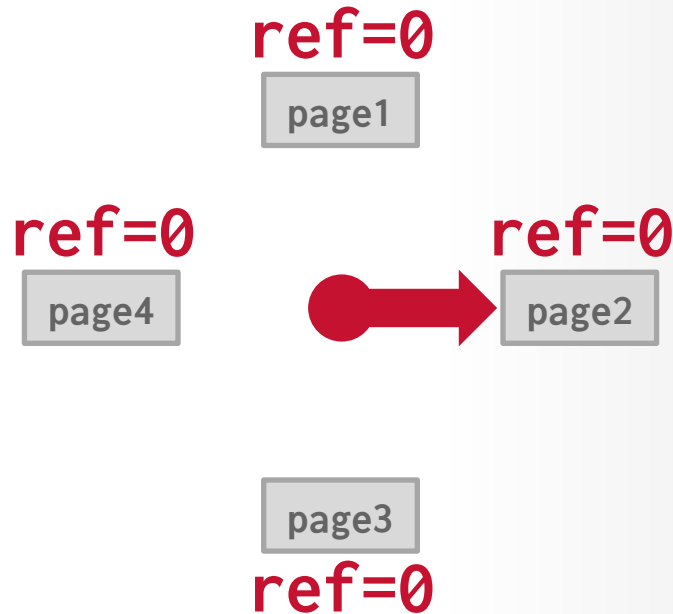
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



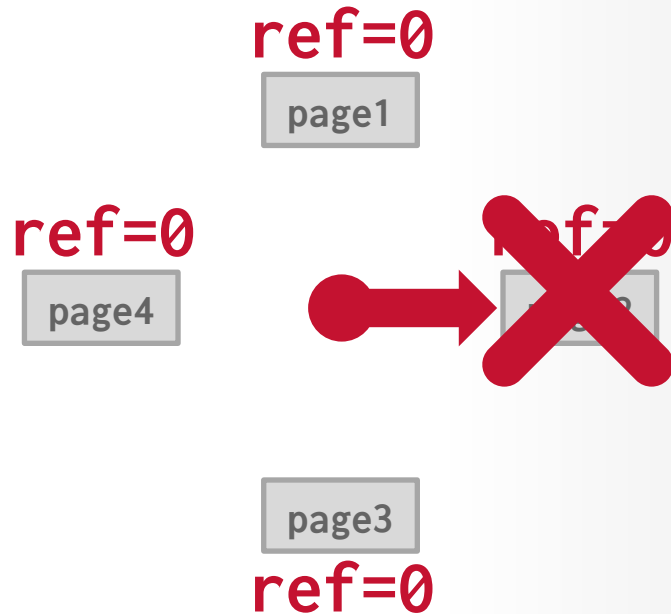
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



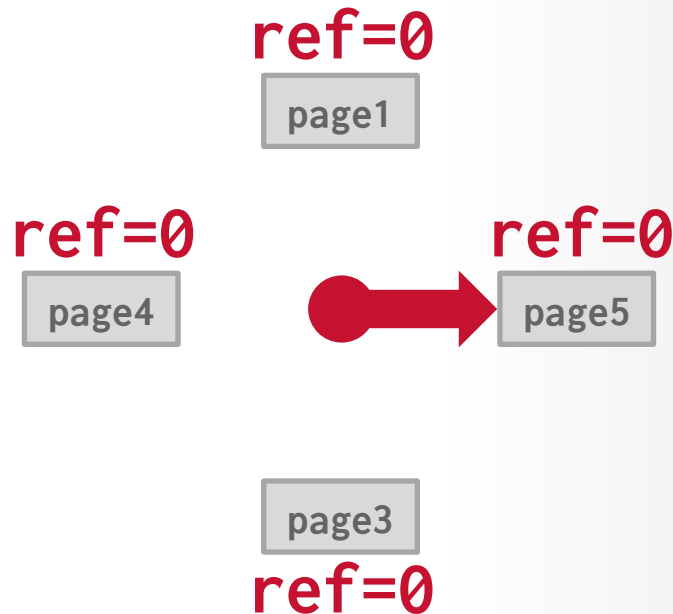
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



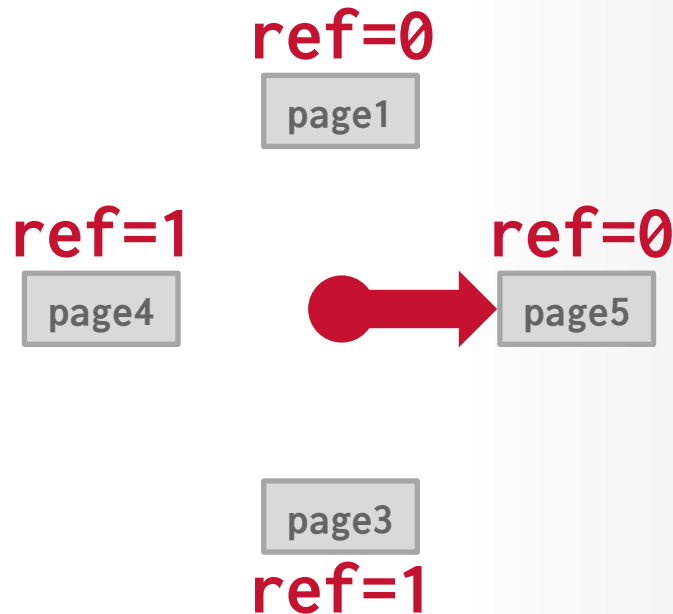
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



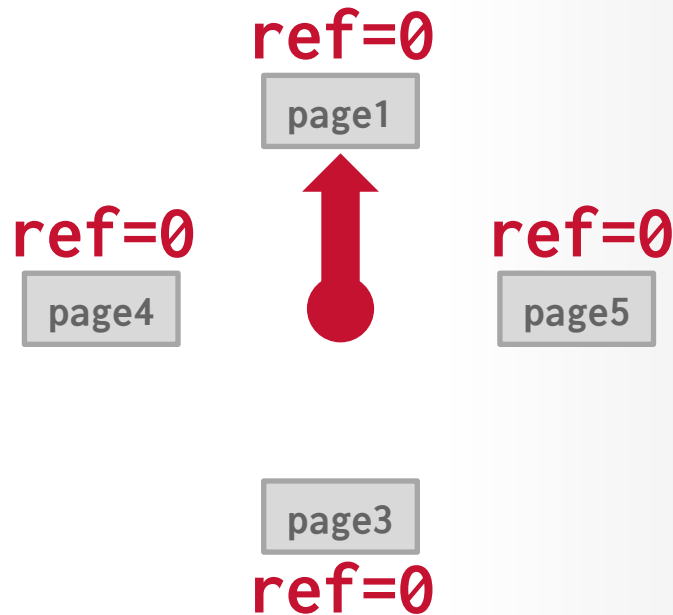
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



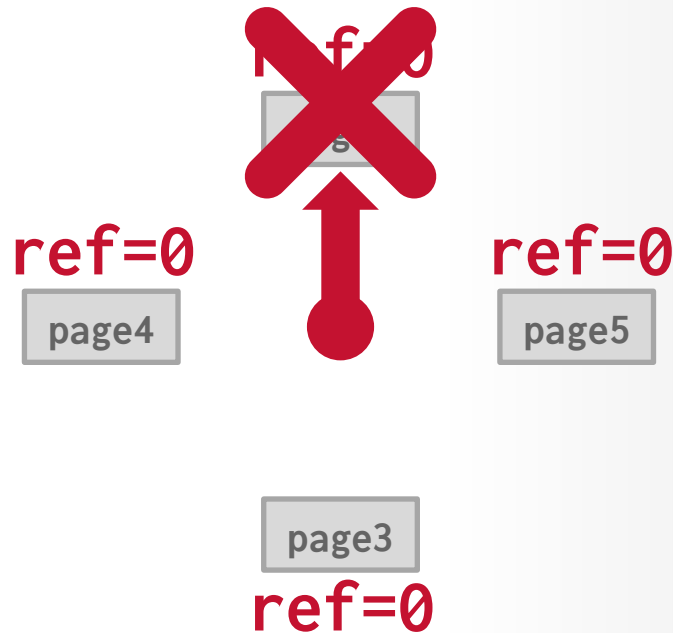
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



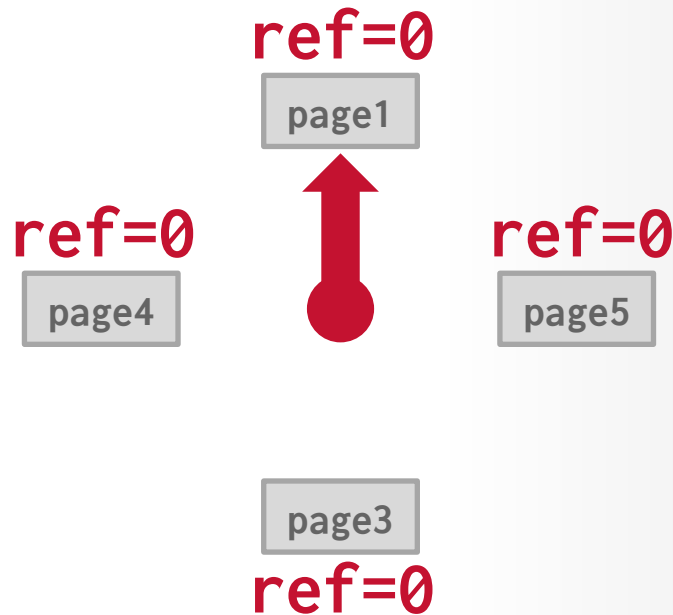
CLOCK

Approximation of LRU that does not need a separate timestamp per page.

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

Organize pages in a circular buffer with a "clock hand" that sweeps over pages in order:

- As the hand visits each page, check if its bit is set to 1.
- If yes, set to zero. If no, then evict.



OBSERVATION

LRU + CLOCK replacement policies are susceptible to **sequential flooding**. 顺序扫描导致的缓存池泛滥.

- A query performs a sequential scan that reads every page in a table one or more times (e.g., blocked nested-loop joins).
- 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 id = 1`

Buffer Pool



Disk Pages



SEQUENTIAL FLOODING

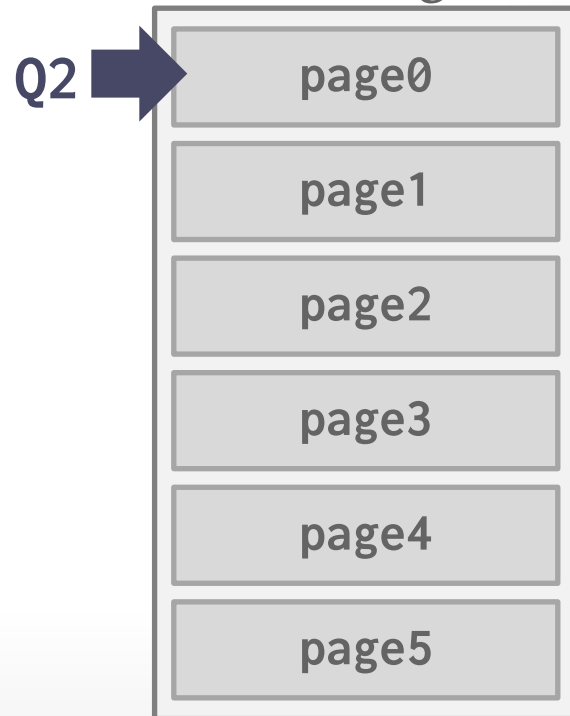
Q1 `SELECT * FROM A WHERE id = 1`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



Disk Pages



SEQUENTIAL FLOODING

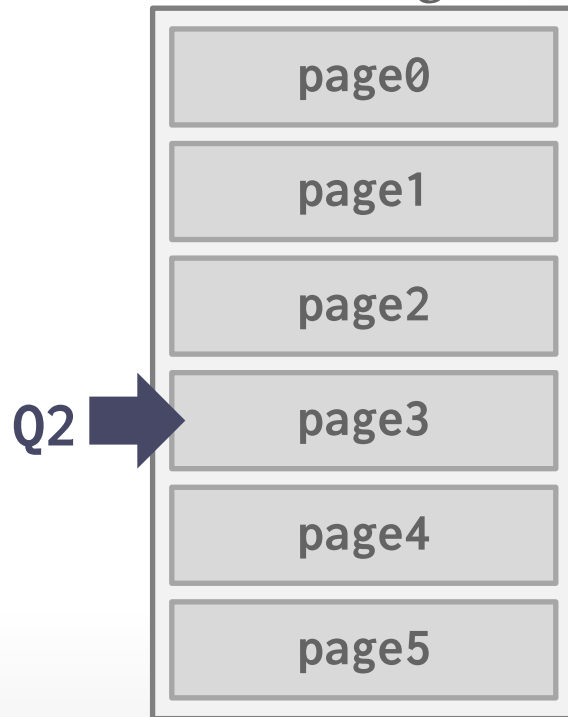
Q1 `SELECT * FROM A WHERE id = 1`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



Disk Pages



SEQUENTIAL FLOODING

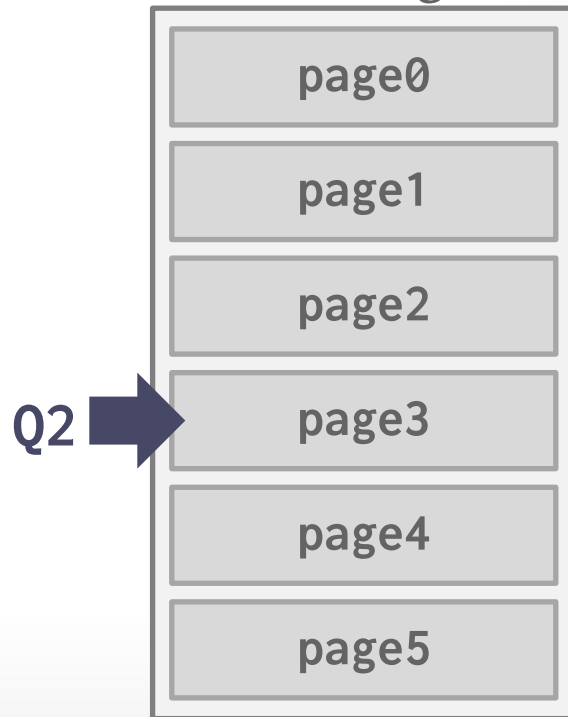
Q1 `SELECT * FROM A WHERE id = 1`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



Disk Pages



SEQUENTIAL FLOODING

Q1 SELECT * FROM A WHERE id = 1

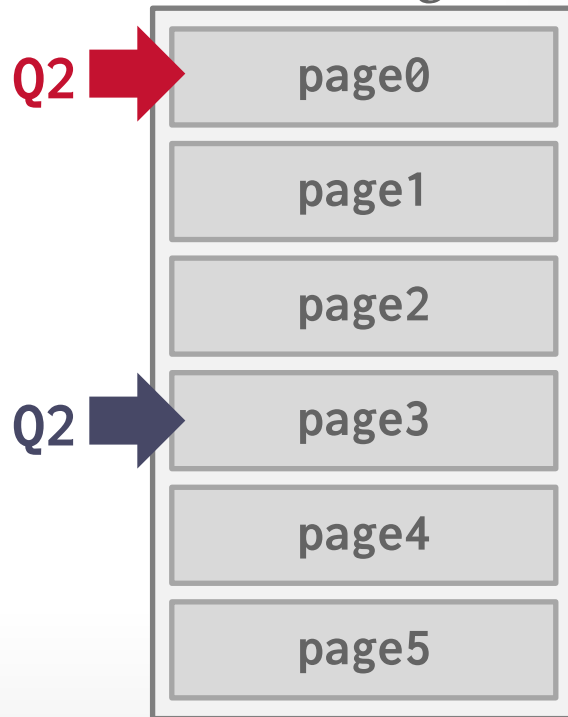
Q2 SELECT AVG(val) FROM A

Q3 SELECT * FROM A WHERE id = 1

Buffer Pool



Disk Pages



SEQUENTIAL FLOODING

Q1 SELECT * FROM A WHERE id = 1

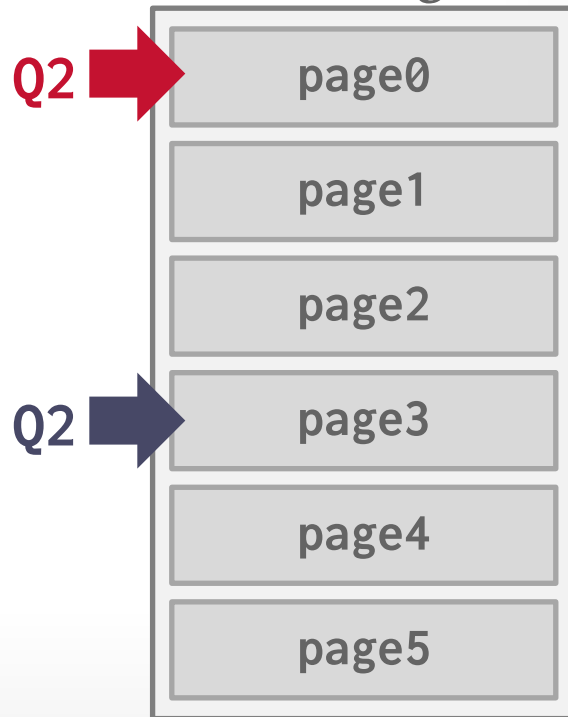
Q2 SELECT AVG(val) FROM A

Q3 SELECT * FROM A WHERE id = 1

Buffer Pool



Disk Pages



BETTER POLICIES: LRU-K

保留额外的元数据信息：跟踪每页最后被访问的 K 个时间戳。

Track the history of last K references to each page as timestamps and compute the interval between subsequent accesses.

→ Can distinguish between reference types

Use this history to estimate the next time that page is going to be accessed.

→ Replace the page with the oldest "K-th" access.

→ Balances recency vs. frequency of access.

→ Maintain an ephemeral in-memory cache for recently evicted pages to prevent them from always being evicted.

Weaving Relations for Cache Performance

Anastasia Ailamaki[‡]
Carnegie Mellon University
anastasi@cs.cmu.edu

David J. DeWitt
Univ. of Wisconsin-Madison
dewitt@cs.wisc.edu

Mark D. Hill
Univ. of Wisconsin-Madison
markhill@cs.wisc.edu

Marios Skounakis
Univ. of Wisconsin-Madison
skounakis@cs.wisc.edu

Abstract

Relational database systems have traditionally optimized for IO performance and organized records sequentially on disk pages using the N-ary Storage Model (NSM) (i.e., shared pages). Recent research, however, indicates that cache utilization and performance is becoming increasingly important on modern platforms. In this paper, we first demonstrate that in-page data placement is the key to high cache performance and that NSM exhibits low cache utilization on modern platforms. Next, we propose a new data organization model (called PAX (Partition Attributes Across)), that significantly improves cache performance by grouping together all values of each attribute within each page. Because PAX only affects layout inside the pages, it incurs no storage penalty and does not affect IO behavior. According to our experimental results, when compared to NSM (or PAX exhibits superior cache and memory bandwidth utilization, saving at least 75% of NSM's stall time due to data cache misses, (b) major selection queries and updates on memory-resident relations execute 17-25% faster, and (c) TPC-H queries involving IO execute 11-48% faster.

1 Introduction

The communication between the CPU and the secondary storage (IO) has been traditionally recognized as the major database performance bottleneck. To optimize data transfer to and from mass storage, relational DBMSs have long organized records in striped disk pages using the N-ary Storage Model (NSM). NSM stores records contiguously starting from the beginning of each disk page, and uses an offset (slot) table at the end of the page to locate the beginning of each record [27].

Unfortunately, most queries use only a fraction of each record. To minimize unnecessary IO, the Decomposition Storage Model (DSM) was proposed in 1985 [16]. DSM partitions an n -attribute relation vertically into n sub-relations, each of which is accessed only when the corresponding attribute is needed. Queries that involve multiple attributes from a relation, however, must spend

tremendous additional time to join the participating sub-relations together. Except for Sybase QJ [33], today's relational DBMSs use NSM for general-purpose data placement [20][29][32].

Recent research has demonstrated that modern database workloads, such as decision support systems and spatial applications, are often bound by delays related to the processor and the memory subsystem rather than IO [20][31][26]. When running commercial database systems on a modern processor, data requests that miss in the cache hierarchy (i.e., requests for data that are not found in any of the caches and are transferred from main memory) are a key memory bottleneck [1]. In addition, only a fraction of the data transferred to the cache is useful to the query; the item that the query processing algorithm requests and the transfer unit between the memory and the processor are typically not the same size. Loading the cache with useless data (a) wastes bandwidth, (b) pollutes the cache, and (c) possibly forces replacement of information that may be needed in the future, incurring even more delays. The challenge is to repair NSM's cache behavior without compromising its advantages over DSM.

This paper introduces and evaluates **Partition Attributes Across (PAX)**, a new layout for data records that combines the best of the two worlds and exhibits performance superior to both placement schemes by eliminating unnecessary accesses to main memory. For a given relation, PAX stores the same data on each page as NSM. Within each page, however, PAX groups all the values of a particular attribute together as a minipage. During a sequential scan (e.g., to apply a predicate on a fraction of the records), PAX fully utilizes the cache resources, because on each miss a number of a single attribute's values are loaded into the cache together. At the same time, all parts of the record are on the same page. To reconstruct a record one needs to perform a mini-join among minipages, which incurs minimal cost because it does not have to look beyond the page.

[‡] Work done while author was at the University of Wisconsin-Madison. Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.

Proceedings of the 27th VLDB Conference, Roma, Italy, 2001



Microsoft®

SQL Server®

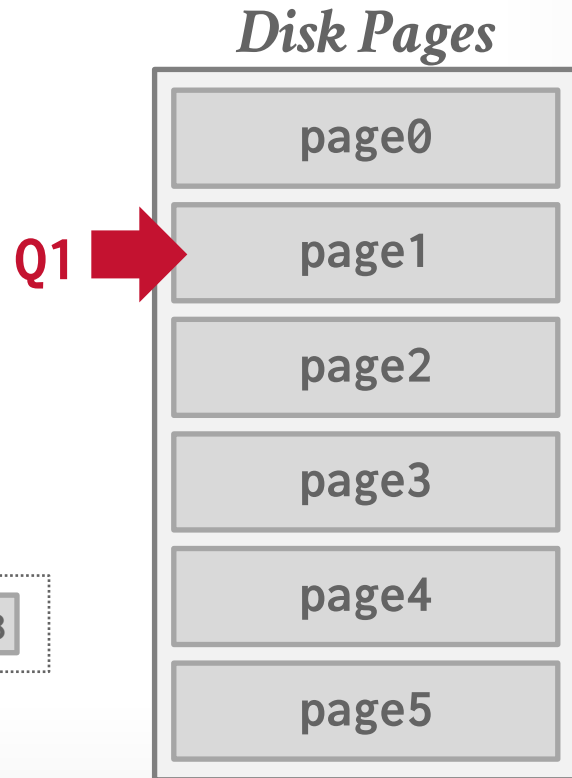
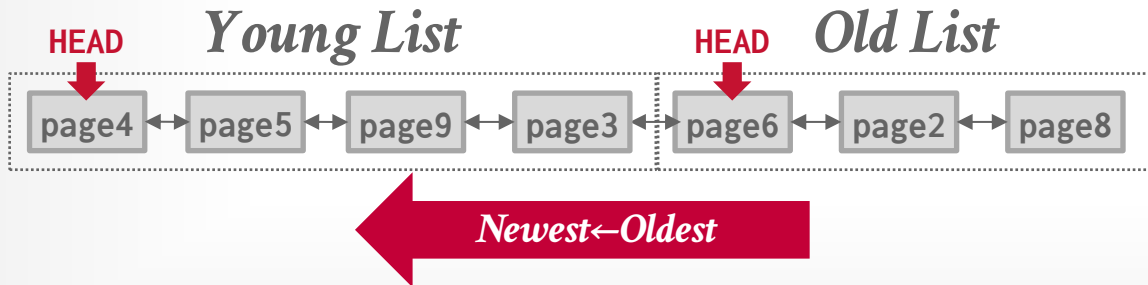


PostgreSQL

MYSQL APPROXIMATE LRU-K

Single LRU linked list but with two entry points ("old" vs "young").

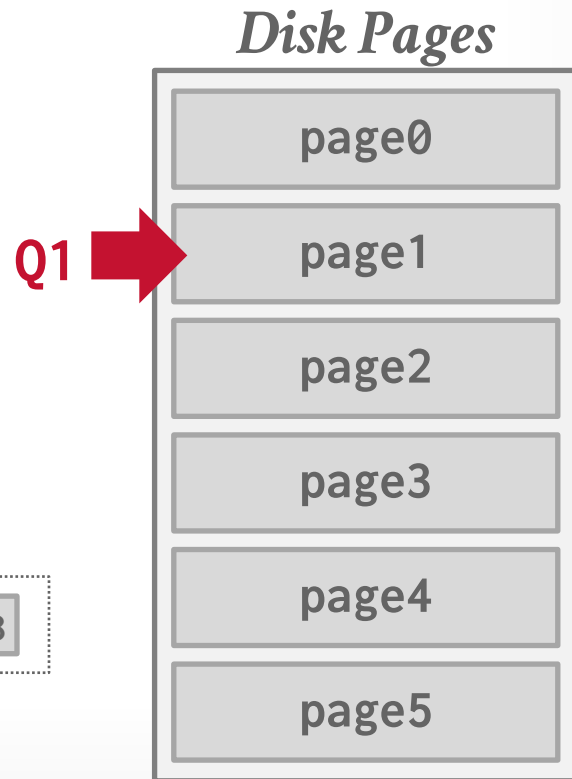
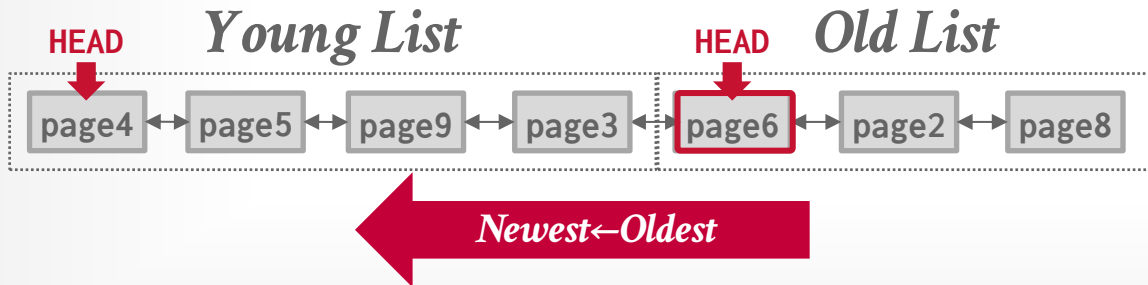
- New pages are always inserted to the head of the old list.
- If pages in the old list is accessed again, then insert into the head of the young list.



MYSQL APPROXIMATE LRU-K

Single LRU linked list but with two entry points ("old" vs "young").

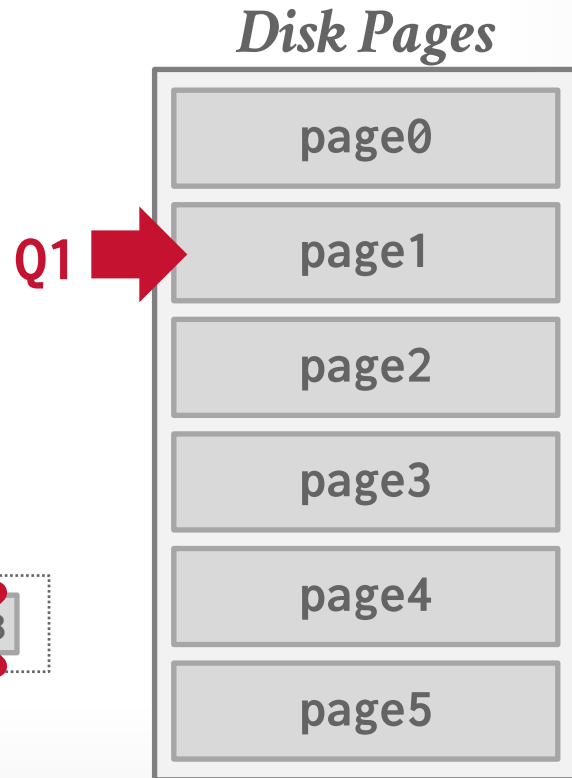
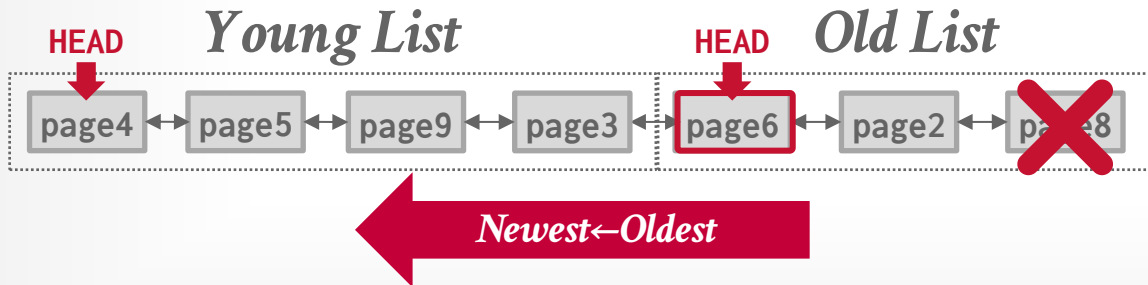
- New pages are always inserted to the head of the old list.
- If pages in the old list is accessed again, then insert into the head of the young list.



MYSQL APPROXIMATE LRU-K

Single LRU linked list but with two entry points ("old" vs "young").

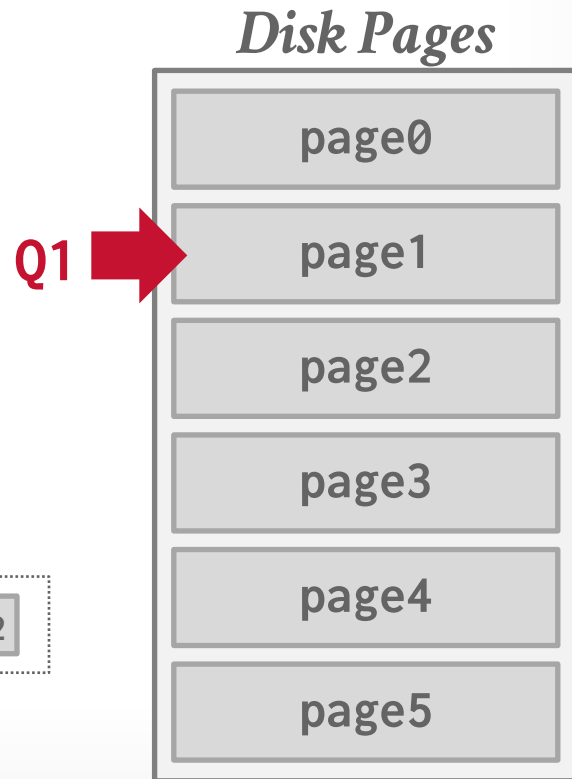
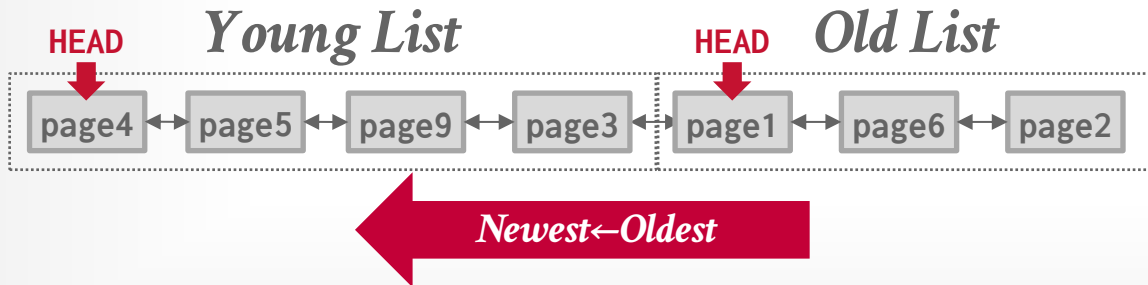
- New pages are always inserted to the head of the old list.
- If pages in the old list is accessed again, then insert into the head of the young list.



MYSQL APPROXIMATE LRU-K

Single LRU linked list but with two entry points ("old" vs "young").

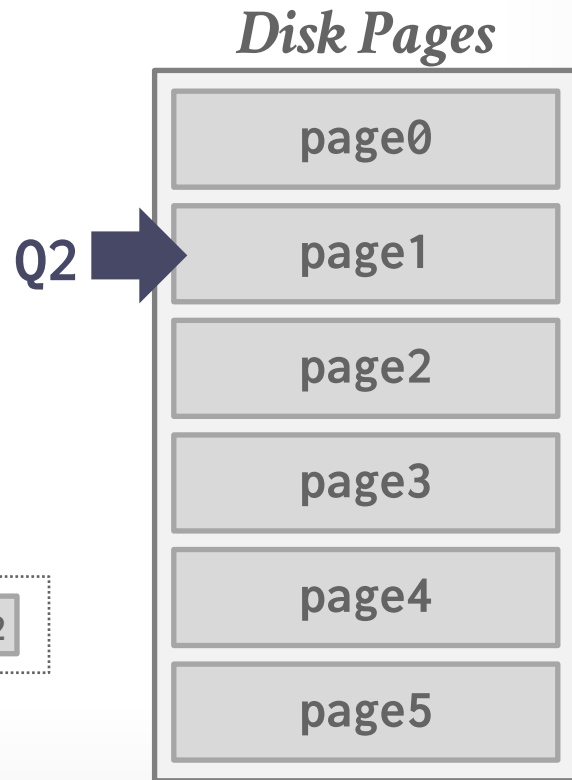
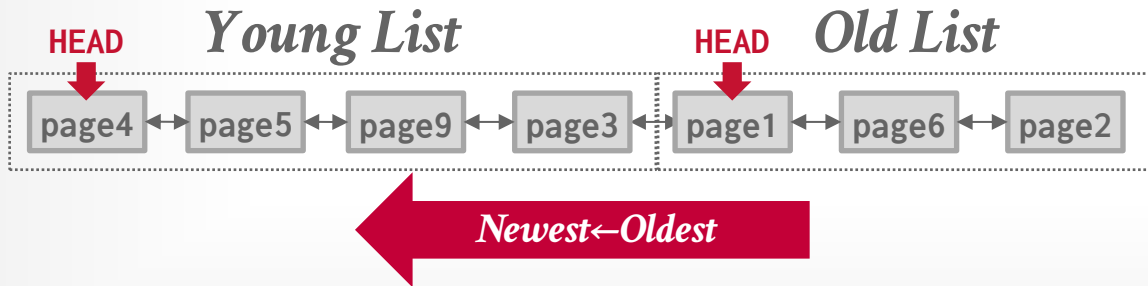
- New pages are always inserted to the head of the old list.
- If pages in the old list is accessed again, then insert into the head of the young list.



MYSQL APPROXIMATE LRU-K

Single LRU linked list but with two entry points ("old" vs "young").

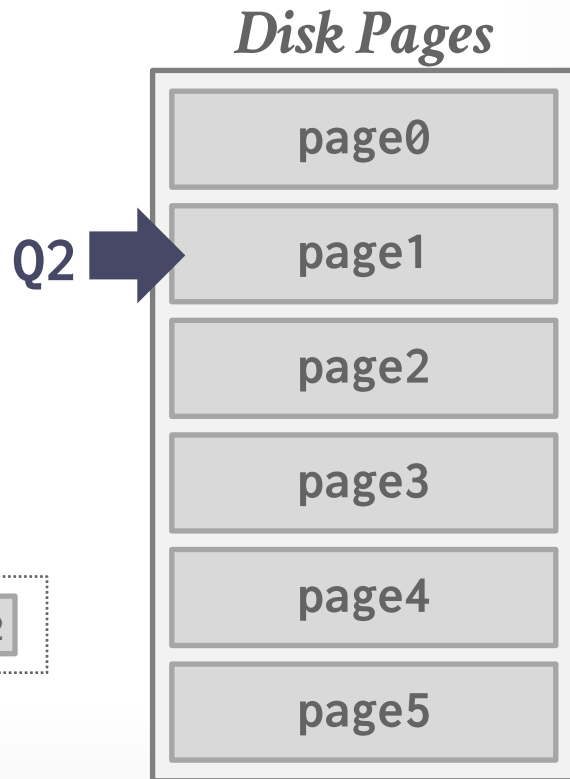
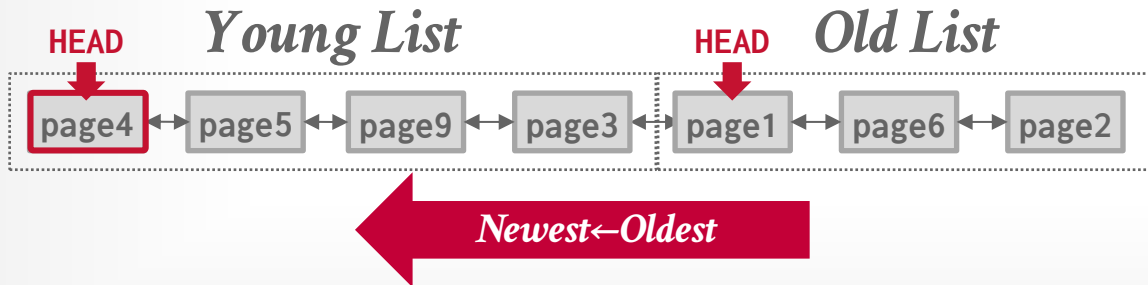
- New pages are always inserted to the head of the old list.
- If pages in the old list is accessed again, then insert into the head of the young list.



MYSQL APPROXIMATE LRU-K

Single LRU linked list but with two entry points ("old" vs "young").

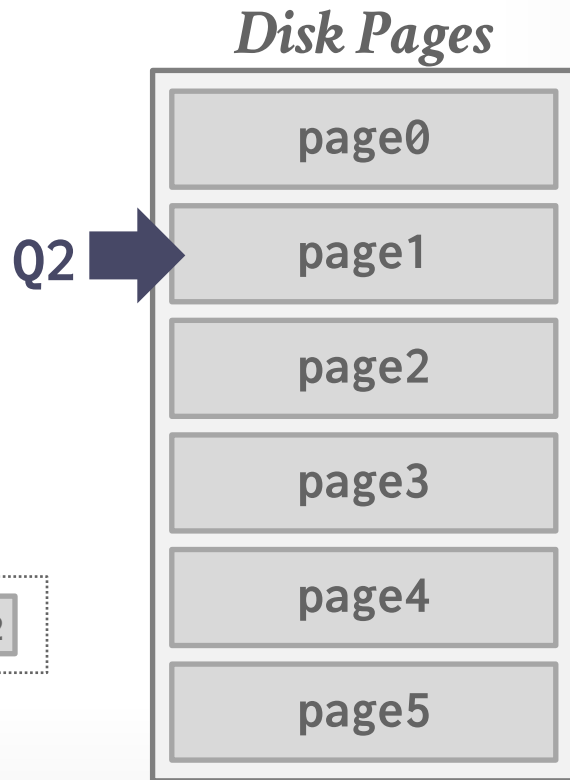
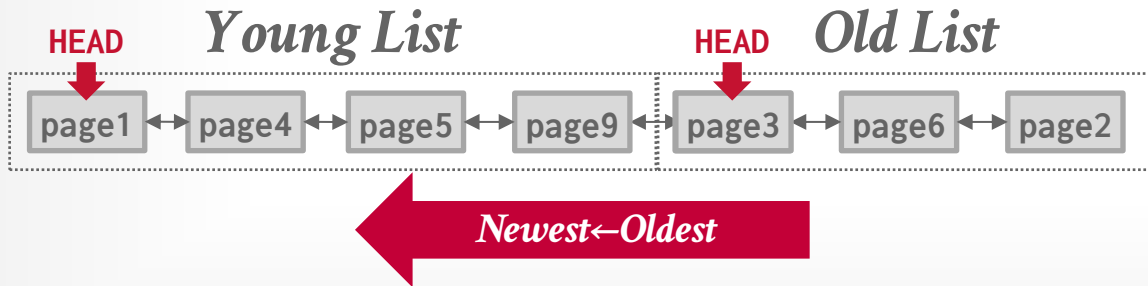
- New pages are always inserted to the head of the old list.
- If pages in the old list is accessed again, then insert into the head of the young list.



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.

使用缓存池的部分空间用于存放查询访问的数据页.

Example: Postgres assigns a limited number of buffer of buffer pool pages to a query and uses it as a circular ring buffer.

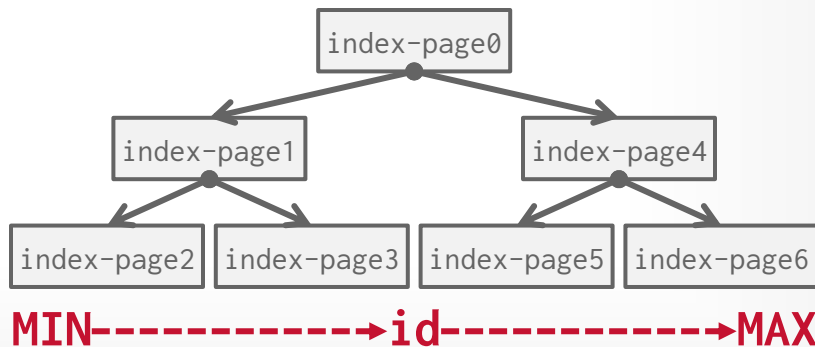
BETTER POLICIES: PRIORITY HINTS

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.

给数据库系统一个提示：某些页比较重要，尽量不要驱逐它们。

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

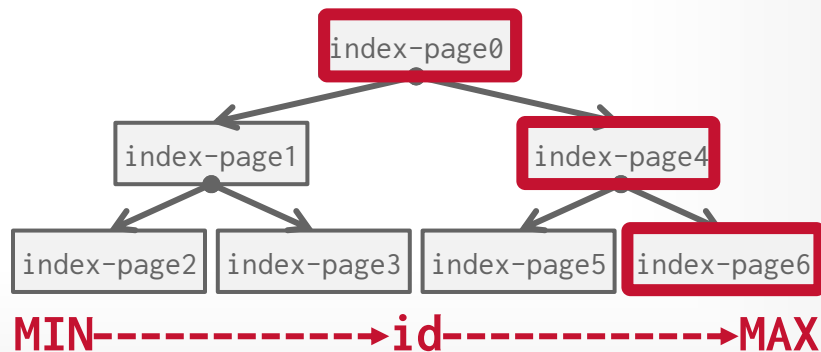


BETTER POLICIES: PRIORITY HINTS

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++*)



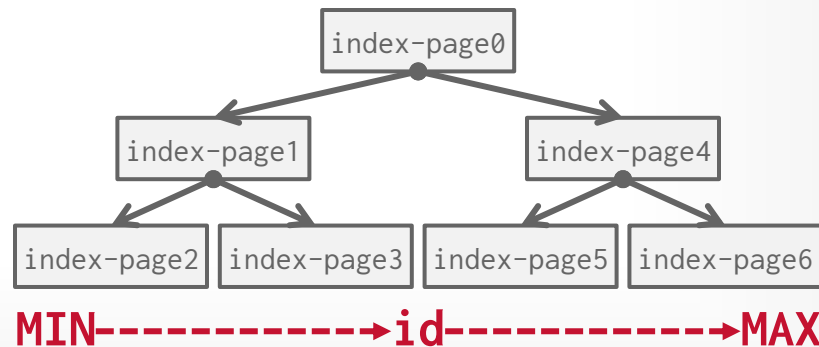
BETTER POLICIES: PRIORITY HINTS

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

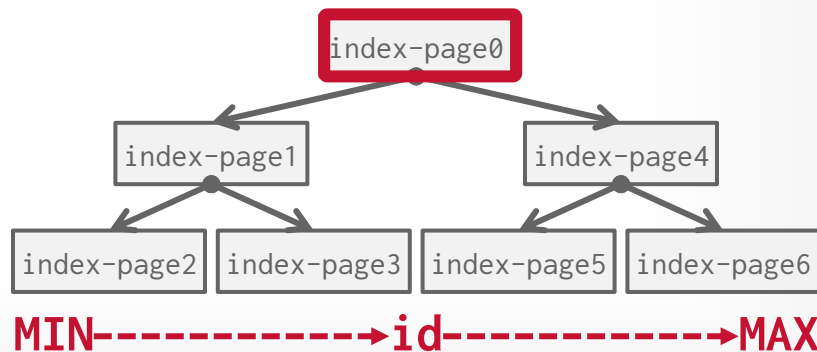
The DBMS knows about the context of each page during query execution.

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

例如：索引结构的根节点会被经常访问，应该留在缓存池中。

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

后台写入器定期将缓存池中的脏页写入磁盘, 但该页可继续驻留在内存中, 避免该缓存页淘汰时查询线程阻塞.

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...

在决定写出哪些脏页时, 利用缓存池驱逐策略找出哪些脏页更有可能被替换, 并优先刷新它们.

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. 关闭 Linux 提供的其他调度.

→ 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

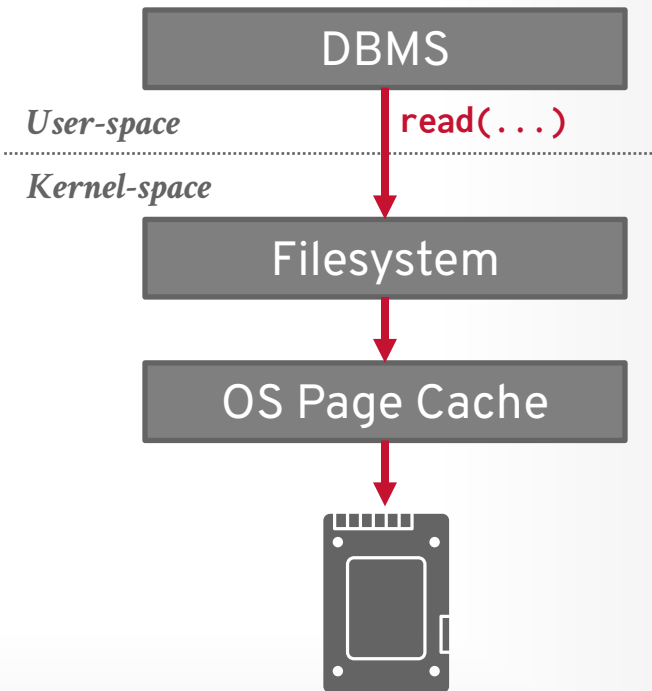
The OS doesn't know these things and is going to get into the way...

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.



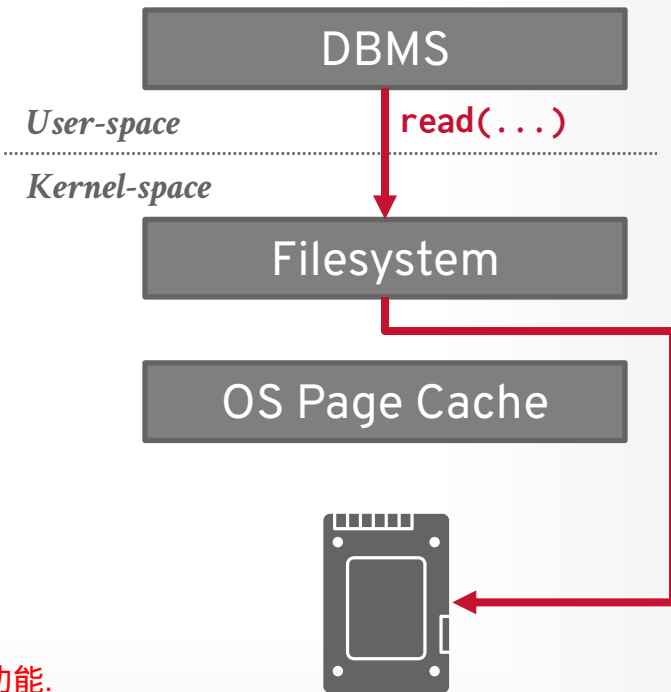
OS PAGE CACHE

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.

Postgres 依赖于 OS's page cache 机制，后续版本在尝试使用 O_DIRECT 功能。



OS PAGE C

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.

Krishnakumar R • 3rd+
Group Engineering Manager, PostgreSQL engine @ MicroS...
4mo • [Follow](#) ...

Direct IO in PostgreSQL and double buffering

The following was an experiment I had shown in my talk on PostgreSQL and Kernel interactions at PGDay Chicago last week :-)

The left side shows the default setting. When contents from a table are read, it will get cached both in the postgres buffer pool and kernel page cache. The third command shows the page details from the pg buffer pool, and the last command (uses fallocate utility) shows info on how much the file corresponding to the table (refresh note: PostgreSQL uses files for its data storage) is cached in the kernel. Note that PG has 8K block size while Kernel has 4K pages (x64 in this case).

On the right you can see developer debug setting which is present from PG16 onwards for enabling direct io is switched on for 'data'. This results in the pages no longer cached in kernel page cache and only cached in buffer pool of pg. As resultant you can see from the output from fallocate not pages are cached in page cache.

#postgres #PostgreSQL #Kernel #PageCache #Linux #LinuxKernel

```

postgres=# show debug_io_direct;
debug_io_direct
(1 row)

postgres=# select * from map limit 30;
 i | t
---+---
 1 | 
200 | Two Hundred
300 | Three Hundred
400 | Four Hundred
500 | Five Hundred
600 | Six Hundred
700 | Seven Hundred
800 | Eight Hundred
900 | Nine Hundred
1000 | Thousand
(10 rows)

postgres=# select bufferid,relfilnode from pg_buffercache where relfilnode=16384;
 bufferid | relfilnode
(1 row)
 0        | 16384

postgres=# \! fallocate -i fallocate -i pg/data/bsu/s/16384
BS PAGE SIZE FILE
0 0K Install pg/data/bsu/s/16384

```

```

postgres=# show debug_io_direct;
debug_io_direct
(1 row)

postgres=# select * from map limit 30;
 i | t
---+---
 100 | Hundred
200 | Two Hundred
300 | Three Hundred
400 | Four Hundred
500 | Five Hundred
600 | Six Hundred
700 | Seven Hundred
800 | Eight Hundred
900 | Nine Hundred
1000 | Thousand
(10 rows)

postgres=# select bufferid,relfilnode from pg_buffercache where relfilnode=16384;
 bufferid | relfilnode
(1 row)
 0        | 16384

postgres=# \! fallocate -i fallocate -i pg/data/bsu/s/16384
BS PAGE SIZE FILE
0 0K Install pg/data/bsu/s/16384

```

FSYNC PROBLEMS

If the DBMS calls **fwrite**, what happens?

If the DBMS calls **fsync**, what happens?

If **fsync** fails (EIO), what happens? 阻塞写入，然后返回。

→ Linux marks the dirty pages as clean.

→ If the DBMS calls **fsync** again, then Linux tells you that the flush was successful. Since the DBMS thought the OS was its friend, it assumed the write was successful...



*Don't
Do This!*

If the DBMS ca

If the DBMS ca

If **fsync** fails (

→ Linux marks t

→ If the DBMS

the flush was

was its friend



*Don't
Do This!*

Navigation: Main Page, Random page, Recent changes, Help

Tools: What links here, Related changes, Special pages, Printable version, Permanent link, Page information

Search: Search PostgreSQL wiki

Fsync Errors

This article covers the current status, history, and OS and OS version differences relating to the circa 2018 fsync() reliability issue discussed on the PostgreSQL mailing list and elsewhere. It has sometimes been referred to as "fsyncgate 2018".

Contents [hide]

- 1 Current status
- 2 Articles and news
- 3 Research notes and OS differences
 - 3.1 Open source kernels
 - 3.2 Closed source kernels
 - 3.3 Special cases
 - 3.4 History and notes

Current status

As of [this PostgreSQL 12 commit](#), PostgreSQL will now PANIC on fsync() failure. It was backpatched to PostgreSQL 11, 10, 9.6, 9.5 and 9.4. Thanks to Thomas Munro, Andres Freund, Robert Haas, and Craig Ringer.

Linux kernel 4.13 improved fsync() error handling and the [man page for fsync\(\)](#) is somewhat improved as well. See:

- Kernelnewbies for 4.13
- Particularly significant 4.13 commits include:
 - "fs: new infrastructure for writeback error handling and reporting"
 - "ext4: use erseq_t based error handling for reporting data writeback errors"
 - "Documentation: flesh out the section in vfs.txt on storing and reporting writeback errors"
 - "mm: set both AS_EIO/AS_ENOSPC and erseq_t in mapping_set_error"

Many thanks to Jeff Layton for work done in this area.

Similar changes were made in [InnoDB/MySQL](#), [WiredTiger/MongoDB](#) and no doubt other software as a result of the PR around this.

A proposed follow-up change to PostgreSQL was discussed in the thread [Refactoring the checkpoint's fsync request queue](#). The [patch that was committed](#) did not incorporate the file-descriptor passing changes proposed. There is still discussion open on some additional safeguards that may use file system error counters and/or filesystem-wide flushing.

Articles and news

- The "fsyncgate 2018" mailing list thread
- LWN.net article "PostgreSQL's fsync() surprise"
- LWN.net article "Improved block-layer error handling"

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

每个缓存池实例独立管理，减少页表竞争。

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

- Avoids contention on LRU tracking meta-data.



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

Buffer Pool #1



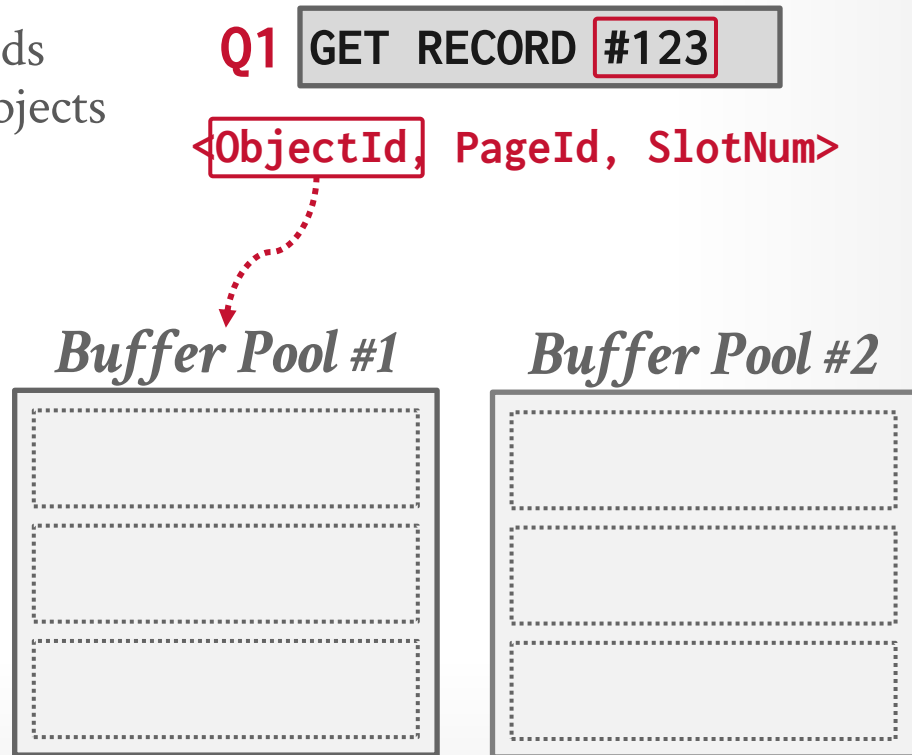
Buffer Pool #2



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.



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**

$\text{HASH}(123) \% n$

Approach #2: Hashing 哈希和轮询

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

MySQL 采用这种方法.

Buffer Pool #1



Buffer Pool #2



PRE-FETCHING

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

→ Examples: Sequential vs. Index Scans

Some DBMS prefetch to fill in empty frames upon start-up.

Buffer Pool



Disk Pages

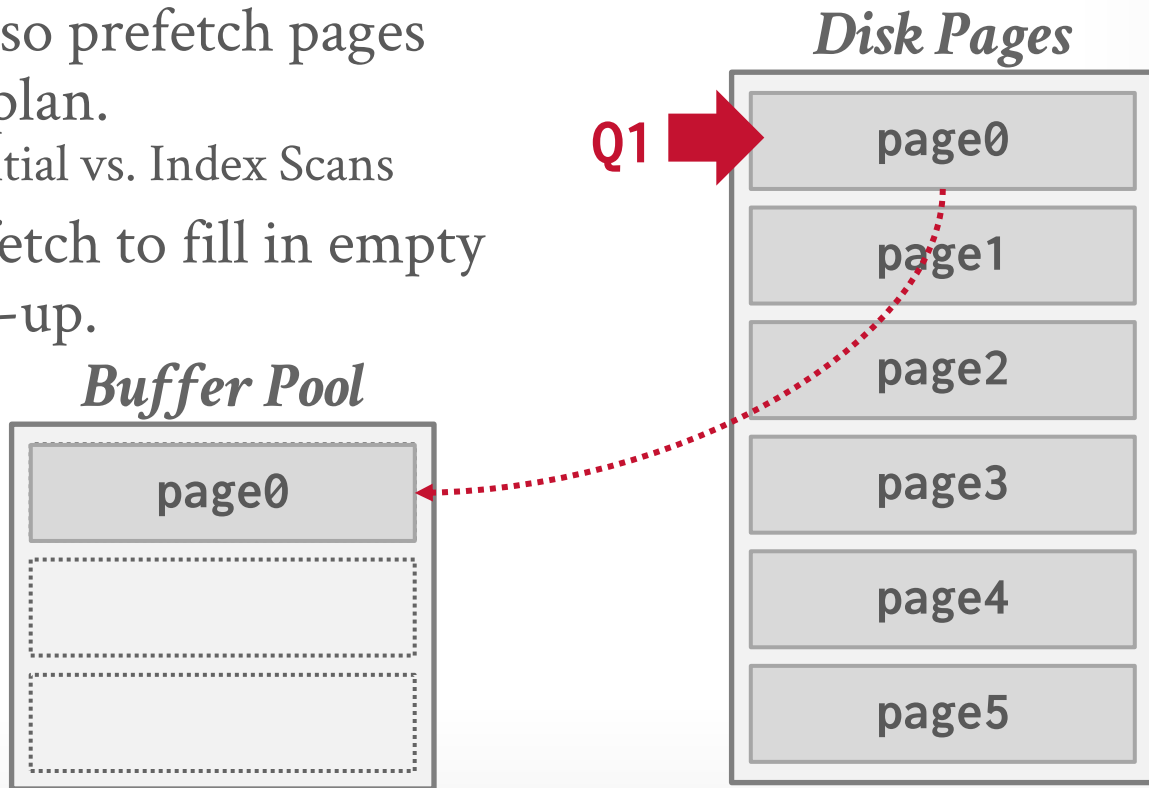


PRE-FETCHING

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

→ Examples: Sequential vs. Index Scans

Some DBMS prefetch to fill in empty frames upon start-up.

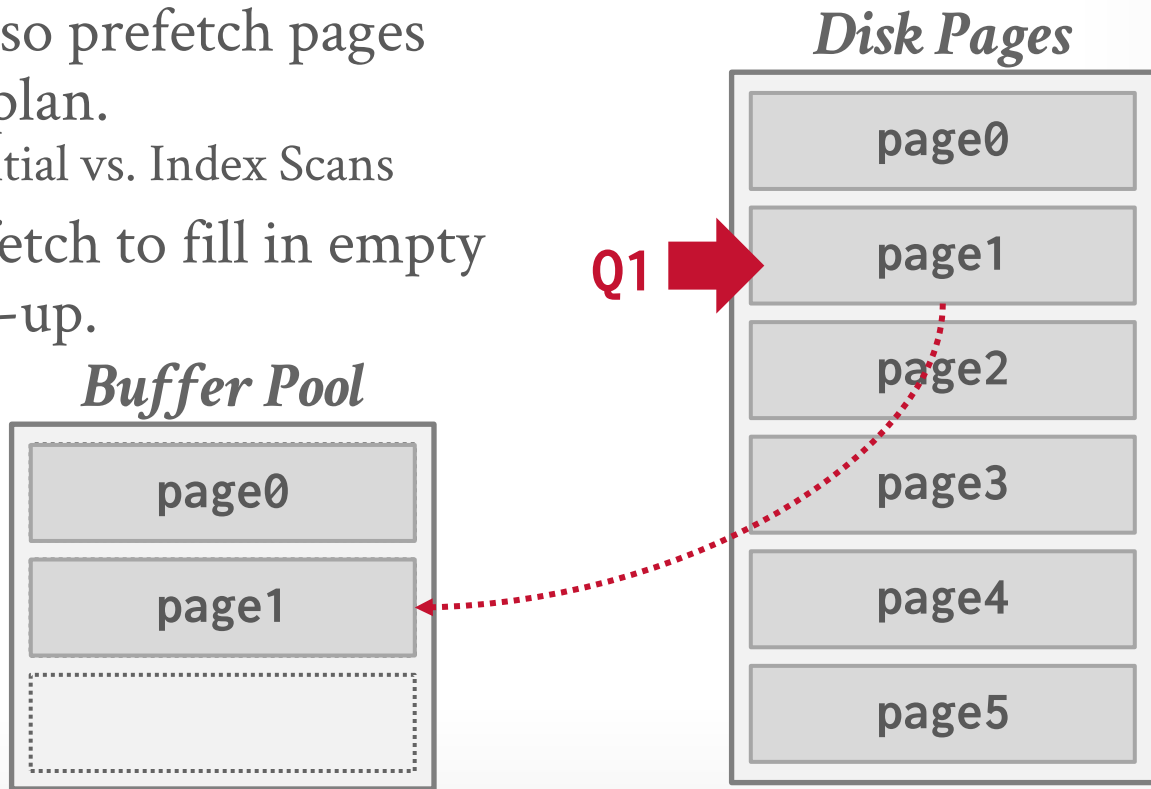


PRE-FETCHING

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

→ Examples: Sequential vs. Index Scans

Some DBMS prefetch to fill in empty frames upon start-up.



PRE-FETCHING

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

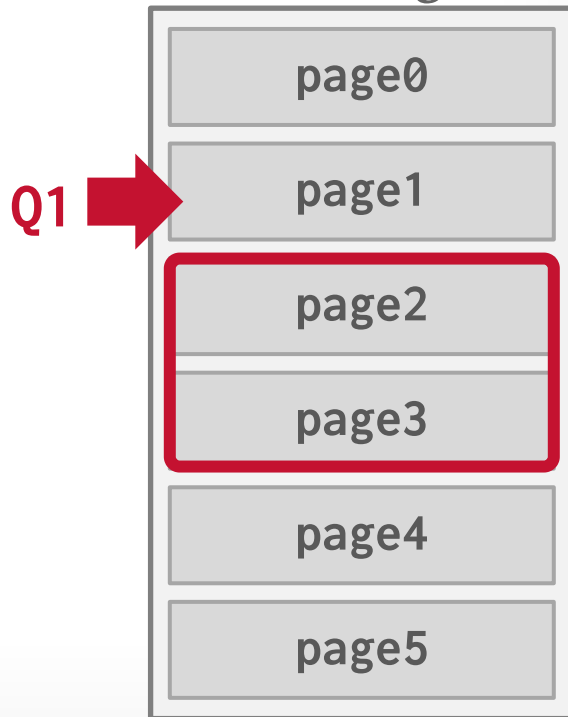
→ Examples: Sequential vs. Index Scans

Some DBMS prefetch to fill in empty frames upon start-up.

Buffer Pool



Disk Pages

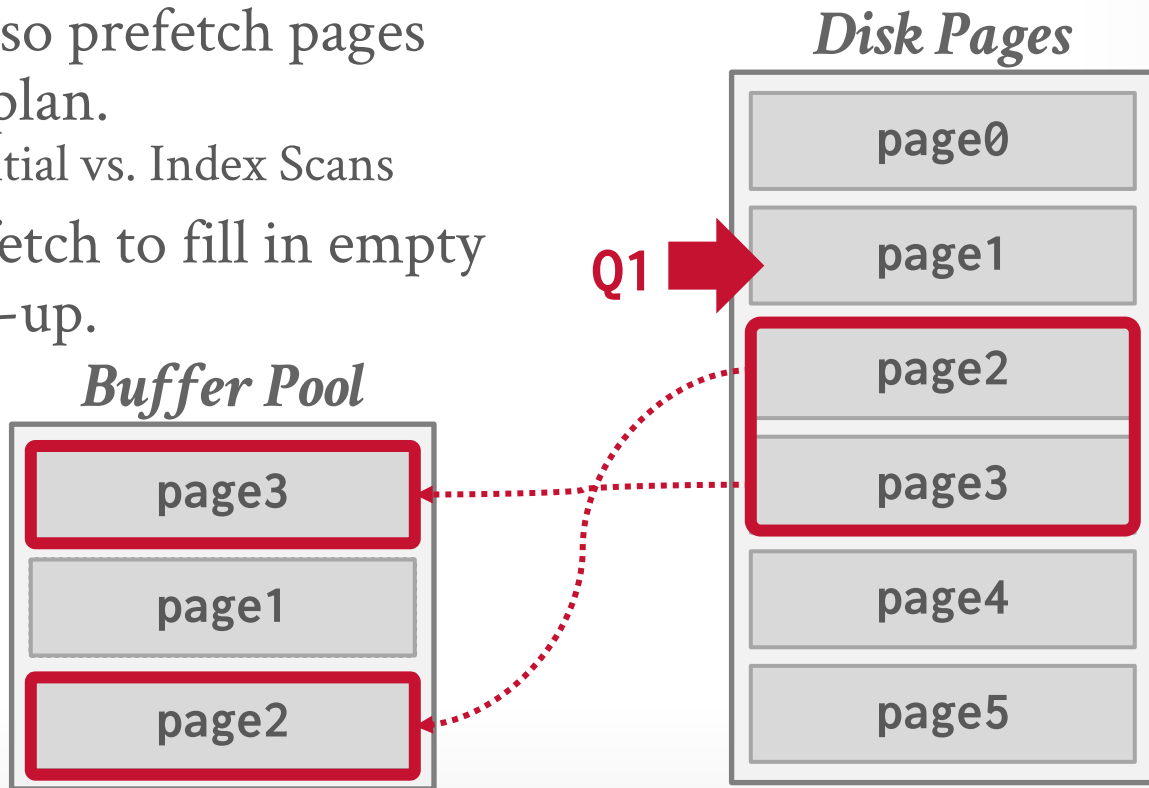


PRE-FETCHING

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

→ Examples: Sequential vs. Index Scans

Some DBMS prefetch to fill in empty frames upon start-up.



PRE-FETCHING

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

→ Examples: Sequential vs. Index Scans

Some DBMS prefetch to fill in empty frames upon start-up.

Buffer Pool



Disk Pages



PRE-FETCHING

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

→ Examples: Sequential vs. Index Scans

Some DBMS prefetch to fill in empty frames upon start-up.

Buffer Pool



Disk Pages



PRE-FETCHING

Q1

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

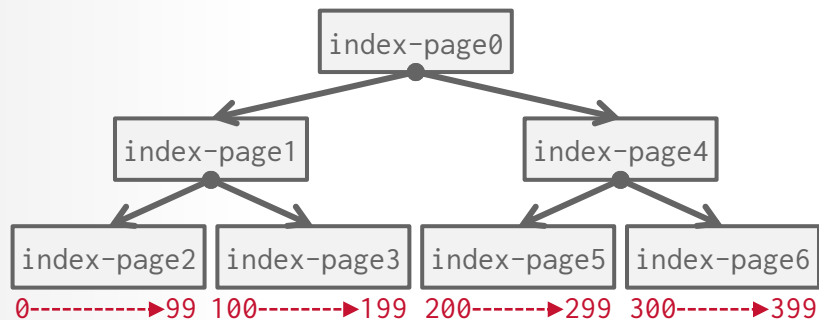
Buffer Pool



Disk Pages



PRE-FETCHING



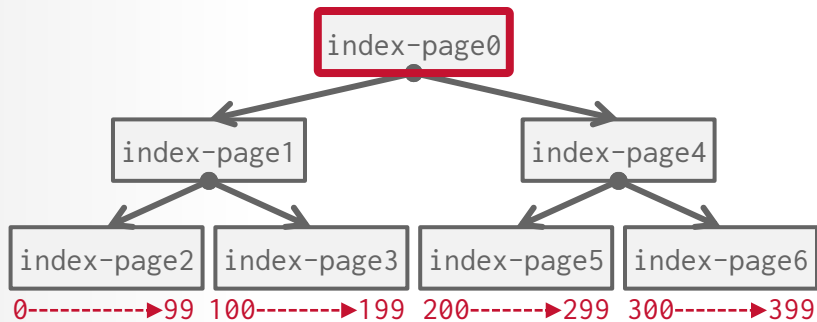
Buffer Pool



Disk Pages



PRE-FETCHING



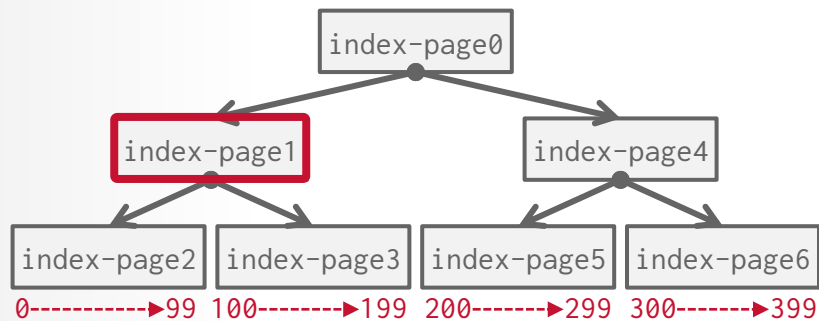
Buffer Pool



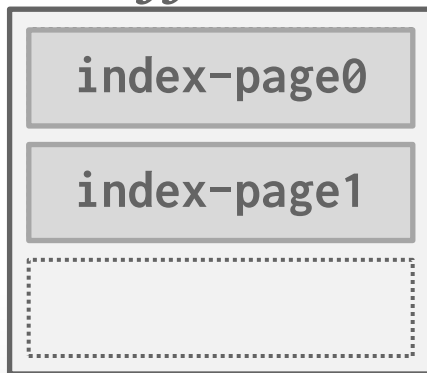
Disk Pages



PRE-FETCHING



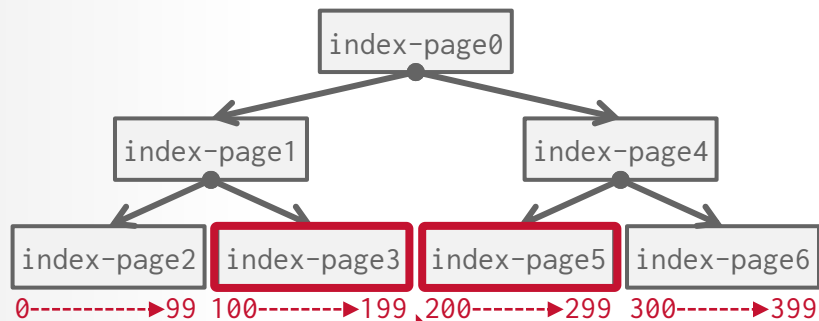
Buffer Pool



Disk Pages



PRE-FETCHING



Buffer Pool



Disk Pages

Q1 →



SCAN SHARING

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

- Also called *synchronized scans*.
- This is different from result caching.

Examples:

- Fully supported in DB2, MSSQL, Teradata, and Postgres.
- Oracle only supports cursor sharing for identical queries.



SCAN SHARING

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

- Also called *synchronized scans*.
- This is different from result caching.

Examples:

- Fully supported in DB2, MSSQL, Teradata, and Postgres.
- Oracle only supports cursor sharing for identical queries.




SCAN SHARING

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

- Also called *synchronized scans*.
- This is different from result caching.

For a textual match to occur, the text of the SQL statements or PL/SQL blocks must be character-for-character identical, including spaces, case, and comments. For example, the following statements cannot use the same shared SQL area:

```
SELECT * FROM employees;  
SELECT * FROM Employees;  
SELECT *   FROM employees;
```

 Copy

ORACLE® reSQL

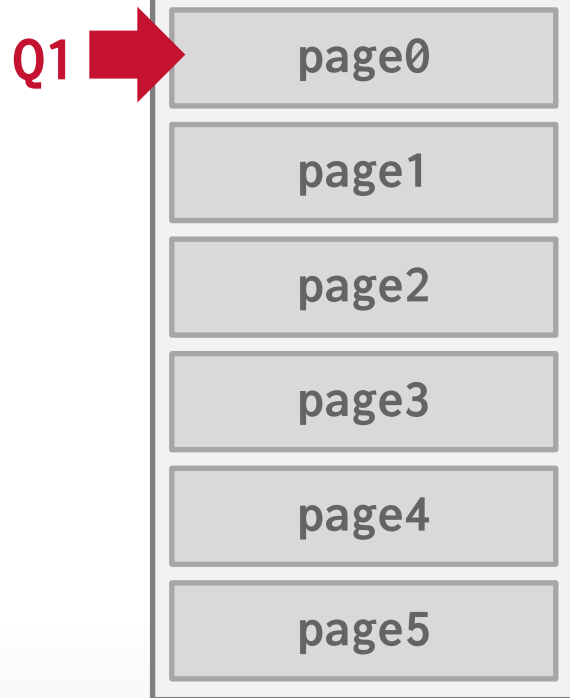
SCAN SHARING

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



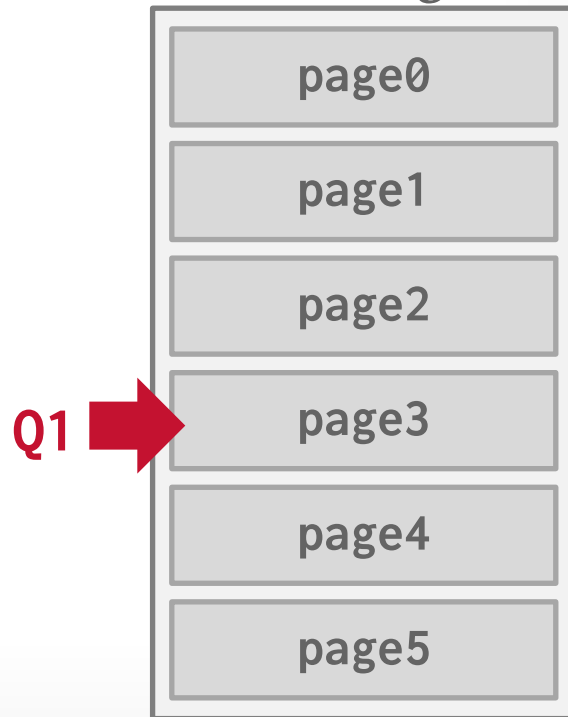
SCAN SHARING

Q1 `SELECT SUM(val) FROM A`

Buffer Pool



Disk Pages



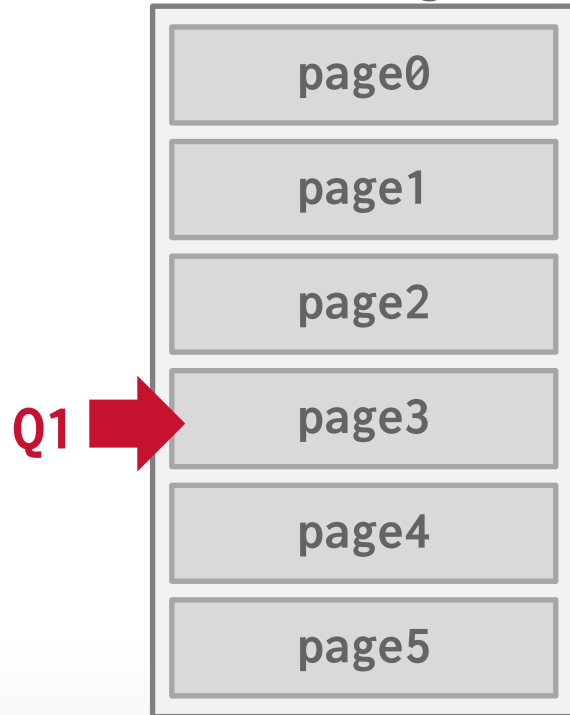
SCAN SHARING

Q1 `SELECT SUM(val) FROM A`

Buffer Pool



Disk Pages



SCAN SHARING

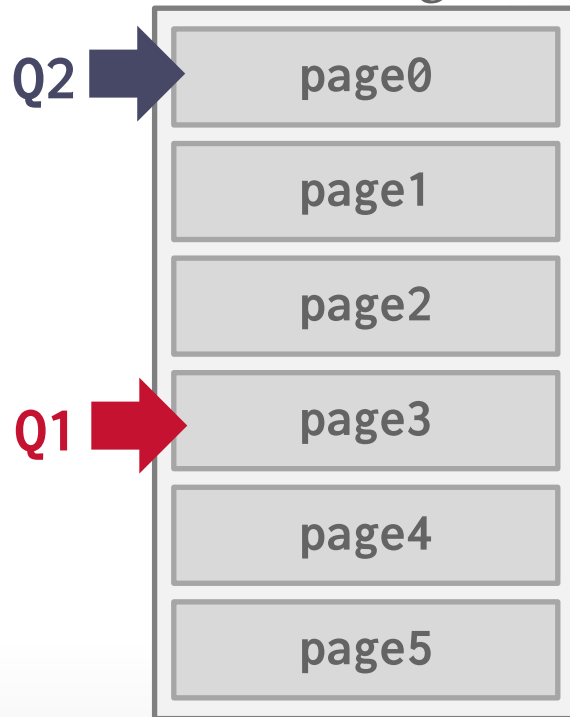
Q1 `SELECT SUM(val) FROM A`

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



Q2 Q1 →

Disk Pages



SCAN SHARING

Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



Disk Pages



SCAN SHARING

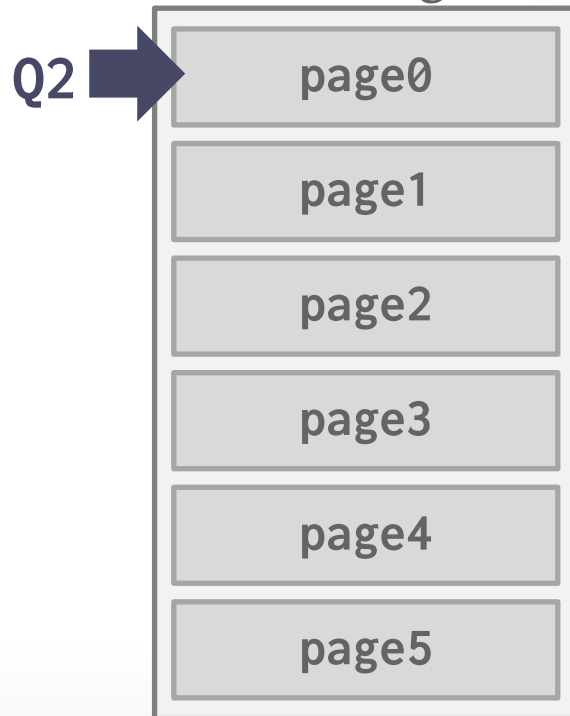
Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



Disk Pages



SCAN SHARING

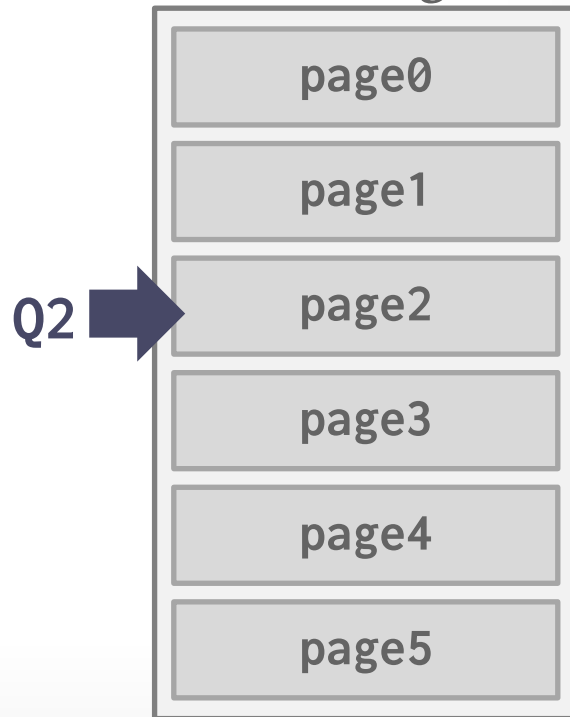
Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A`

Buffer Pool



Disk Pages



SCAN SHARING

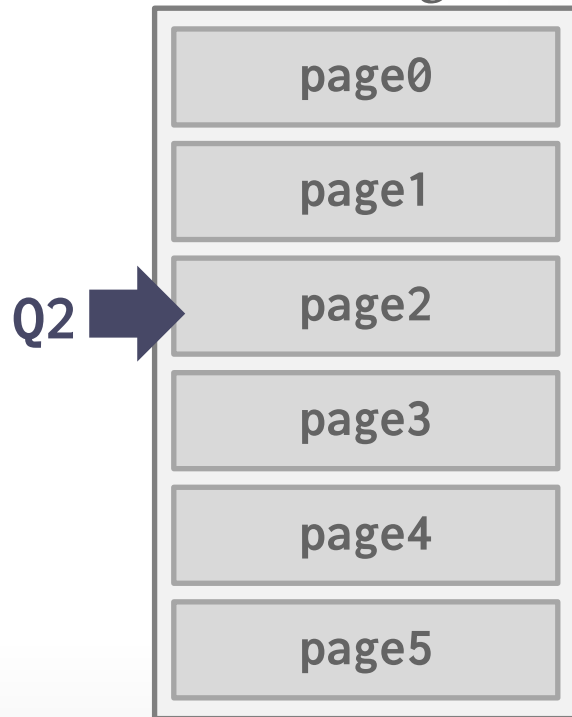
Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A LIMIT 100`

Buffer Pool



Disk Pages



SCAN SHARING

Q1 `SELECT SUM(val) FROM A`

Q2 `SELECT AVG(val) FROM A LIMIT 100`

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.
- Can also be used for temporary data (sorting, joins).

Called "Light Scans" in Informix.

ORACLE®



Microsoft®
SQL Server®

Informix®

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 CLASS

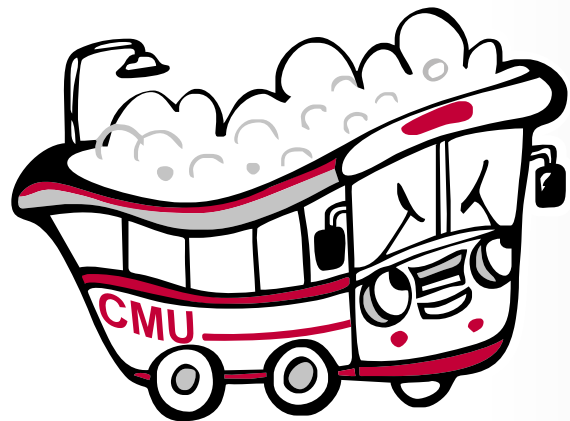
Hash Tables

PROJECT #1

You will build the first component of your storage manager.

- LRU-K Replacement Policy
- Disk Scheduler
- Buffer Pool Manager Instance

We will provide you with the basic APIs for these components.



BusTub

Due Date:
Sunday Sept 29th @ 11:59pm

TASK #1 – LRU-K REPLACEMENT POLICY

Build a data structure that tracks the usage of pages using the LRU-K policy.

General Hints:

- Your **LRUKReplacer** needs to check the "pinned" status of a **Page**.
- If there are no pages touched since last sweep, then return the lowest page id.

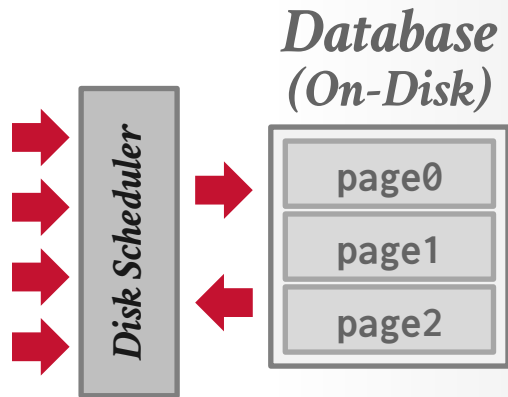
TASK #2 – DISK SCHEDULER

Create a background worker to read/write pages from disk.

- Single request queue.
- Simulates asynchronous IO using **`std::promise`** for callbacks.

It's up to you to decide how you want to batch, reorder, and issue read/write requests to the local disk.

Make sure it is thread-safe!

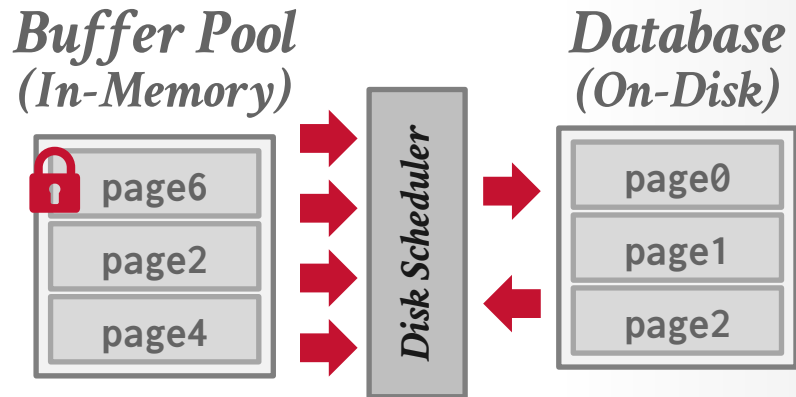


TASK #3 – BUFFER POOL MANAGER

Use your LRU-K replacer to manage the allocation of pages.

- Need to maintain internal data structures to track allocated + free pages.
- Implement page guards.
- Use whatever data structure you want for the page table.

Make sure you get the order of operations correct when pinning!



THINGS TO NOTE

Do **not** change any file other than the six that you must hand in. Other changes will not be graded.

The projects are cumulative.

We will **not** be providing solutions.

Post any questions on Piazza or come to office hours, but we will **not** help you debug.

CODE QUALITY

We will automatically check whether you are writing good code.

- [Google C++ Style Guide](#)
- [Doxygen Javadoc Style](#)

You need to run these targets before you submit your implementation to Gradescope.

- **make format**
- **make check-clang-tidy-p1**

EXTRA CREDIT

Gradescope Leaderboard runs your code with a specialized in-memory version of BusTub.

The top 20 fastest implementations in the class will receive extra credit for this assignment.

- **#1:** 50% bonus points
- **#2–10:** 25% bonus points
- **#11–20:** 10% bonus points

Student with the most bonus points at the end of the semester will receive a BusTub schwag!



PLAGIARISM WARNING



The homework and projects must be your own original work. They are **not** group assignments.

You may **not** copy source code from other people or the web.

Plagiarism is **not** tolerated. You will get lit up.
→ Please ask me if you are unsure.

See [CMU's Policy on Academic Integrity](#) for additional information.