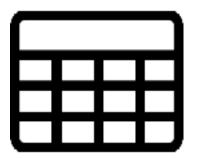
Tables Management



Database System Technology
Niv Dayan



We've decided to allow groups of both 2 and 3. Load is same. We recommend 3.

Please start forming groups and start step 1

Undergrads and grads can form groups

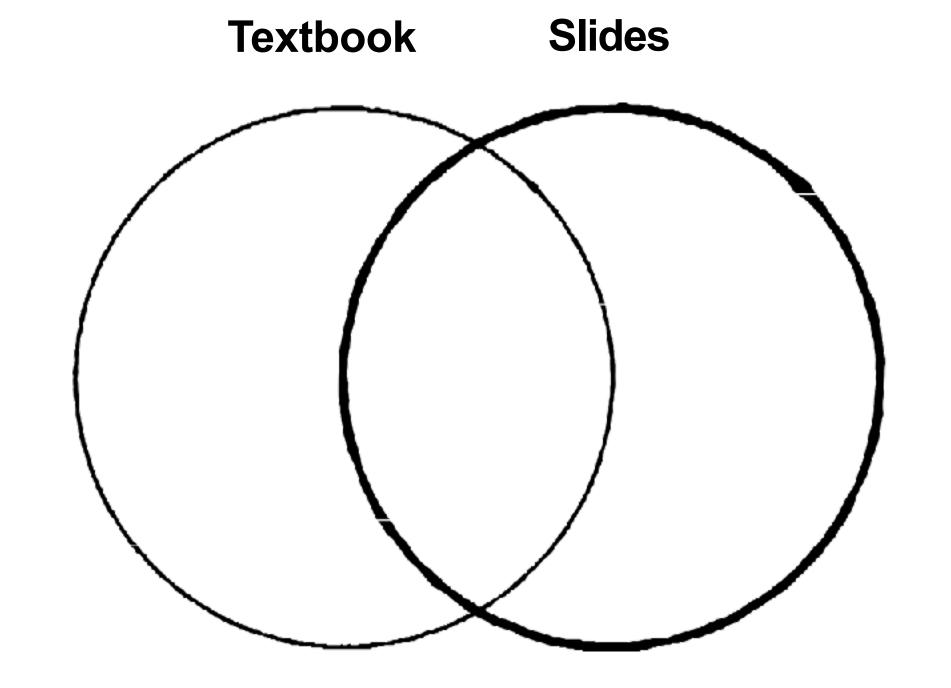
We enabled "Search for Teammates" on Piazza

In person classes from this Thursday:)



This Thursday: Navid

There to solidify your understanding and get a historical perspective



Only material in the slides will appear in the midterm/exam

Database Tables

A database consists of multiple tables

Customers

ID Name email Addr

Orders

ID	Customer ID	Product ID	Date

Database Tables

A database consists of multiple tables

How do we store them in storage efficiently?

Customers

ID	Name	email	Addr

Orders

ID	Customer ID	Product ID	Date



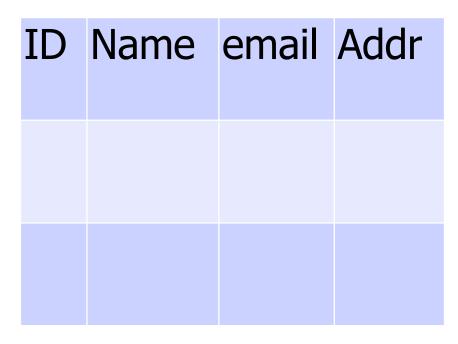


Operations to Efficiently Support

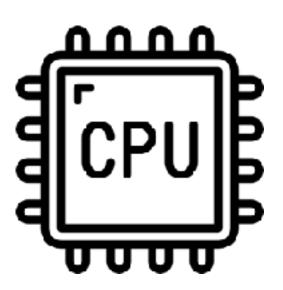
- 1. Scans
- 2. Deletes
- 3. Updates
- 4. Insertions

- e.g., select * from Customers
- e.g., delete from Customers where name = "..."
- e.g., update Customers set email = "..." where name = ""
- e.g., Insert into Customers (,,,)

Customers

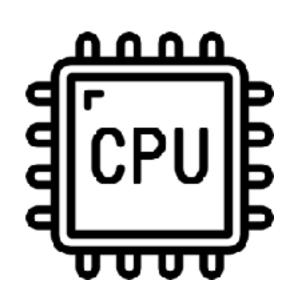


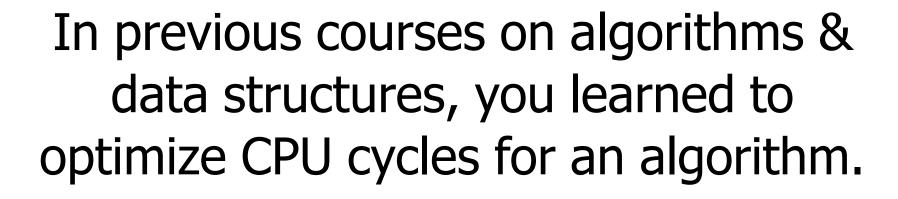
Optimizing for Data Movement



In previous courses on algorithms & data structures, you learned to optimize CPU cycles for an algorithm.

Optimizing for Data Movement







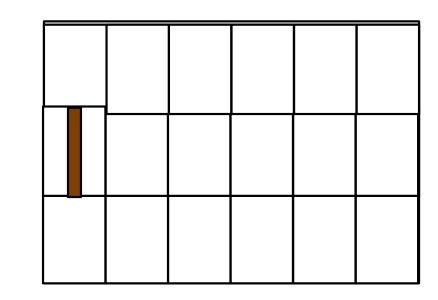


As storage devices are far slower, in this course we focus on optimizing data movement.

Reading/writing from storage at units of less than ≈4KB does not pay off.







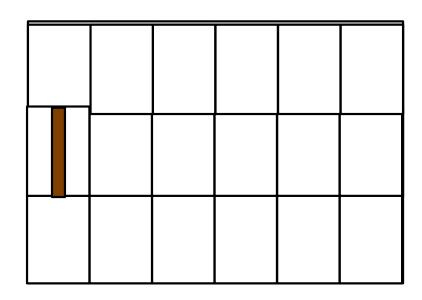
Storage





Reading/writing from storage at units of less than ≈4KB does not pay off.

Why? (Different reasons for disk and SSDs)



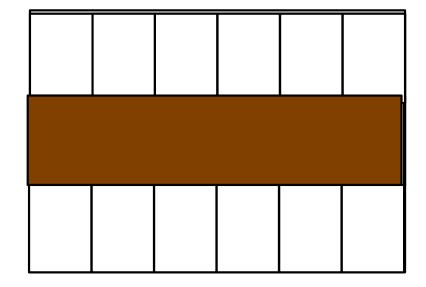
Storage

Reading/writing from storage at units of less than ≈4KB does not pay off.

Reading/writing at very large units consumes memory and is less flexible for applications







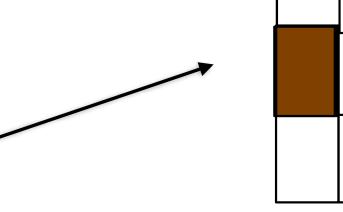
Storage

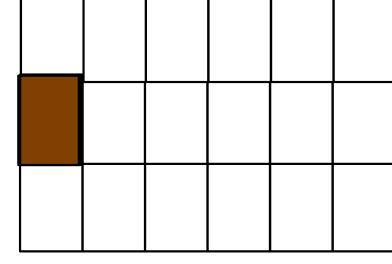




To balance, DBs use ≈4KB as the read/write unit. This is known as a database page.

An I/O (input/output) is one read or write request of one database page.





Storage

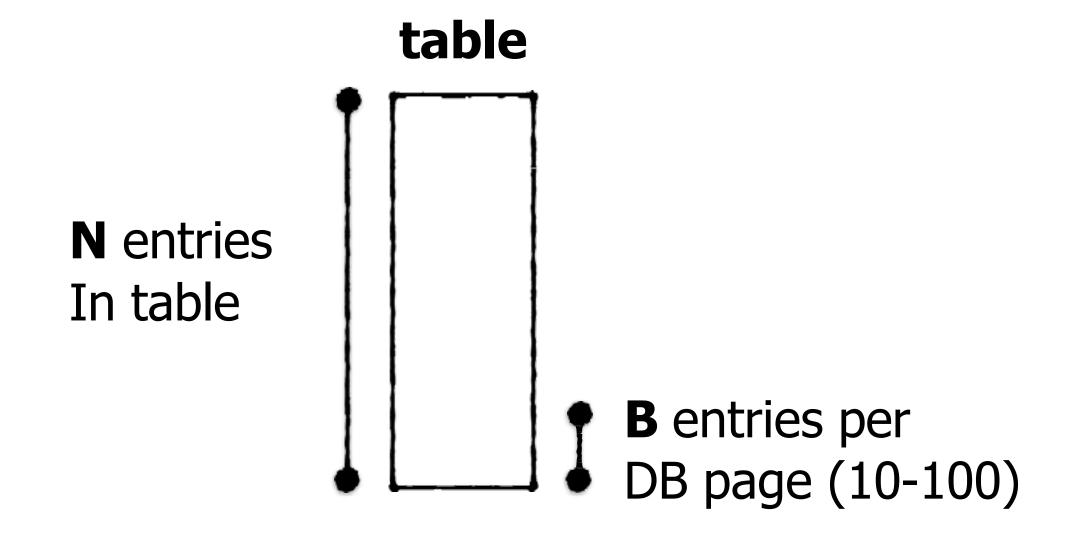
We will shortly propose algorithms to support scans/delete/updates/inserts

To reason about such algorithms, we need a cost model



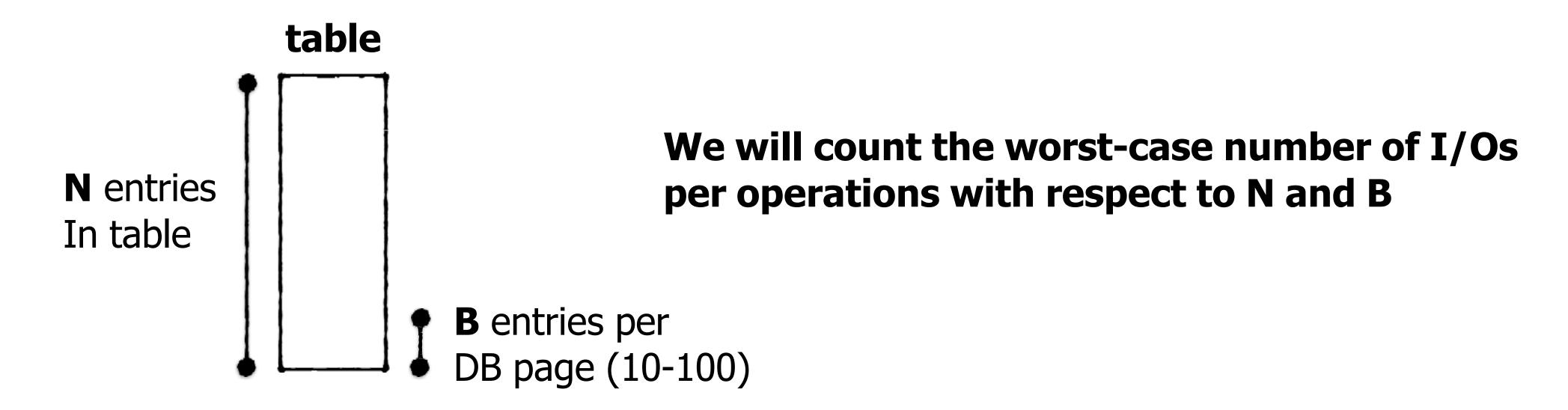
We will shortly propose algorithms to support scans/delete/updates/inserts

To reason about such algorithms, we need a cost model



We will shortly propose algorithms to support scans/delete/updates/inserts

To reason about such algorithms, we need a cost model



This model is imperfect. It ignores many characteristics of storage.

This model is imperfect. It ignores many characteristics of storage.

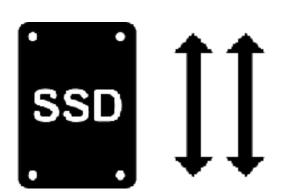


Ignores that sequential access is faster than random on disk

This model is imperfect. It ignores many characteristics of storage.



Ignores that sequential access is faster than random on disk

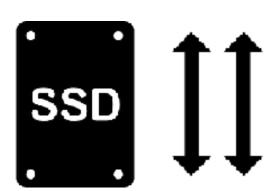


Ignores that SSD asynchronous I/O are faster

This model is imperfect. It ignores many characteristics of storage.



Ignores that sequential access is faster than random on disk



Ignores that SSD asynchronous I/O are faster

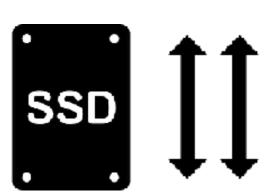


Ignores SSD garbagecollection due to random writes

This model is imperfect. It ignores many characteristics of storage.



Ignores that sequential access is faster than random on disk



Ignores that SSD asynchronous I/O are faster



Ignores SSD garbagecollection due to random writes

However, it's useful due to its simplicity.

Operations

1. Scans e.g., select * from Customers

2. Deletes e.g., delete from Customers where name = "..."

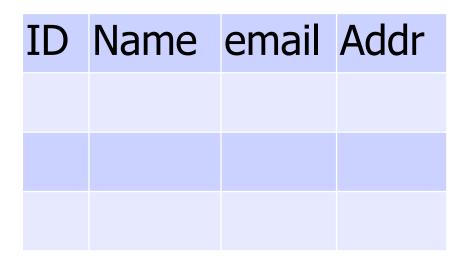
3. Updates e.g., update Customers set email = "..." where name = ""

4. Insertions e.g., Insert into Customers (,,,)

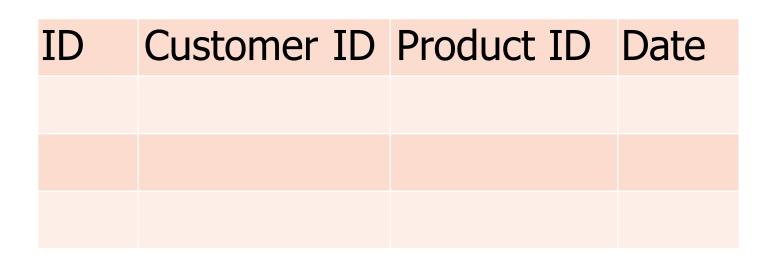
Scans - How not to Support Them

Customers

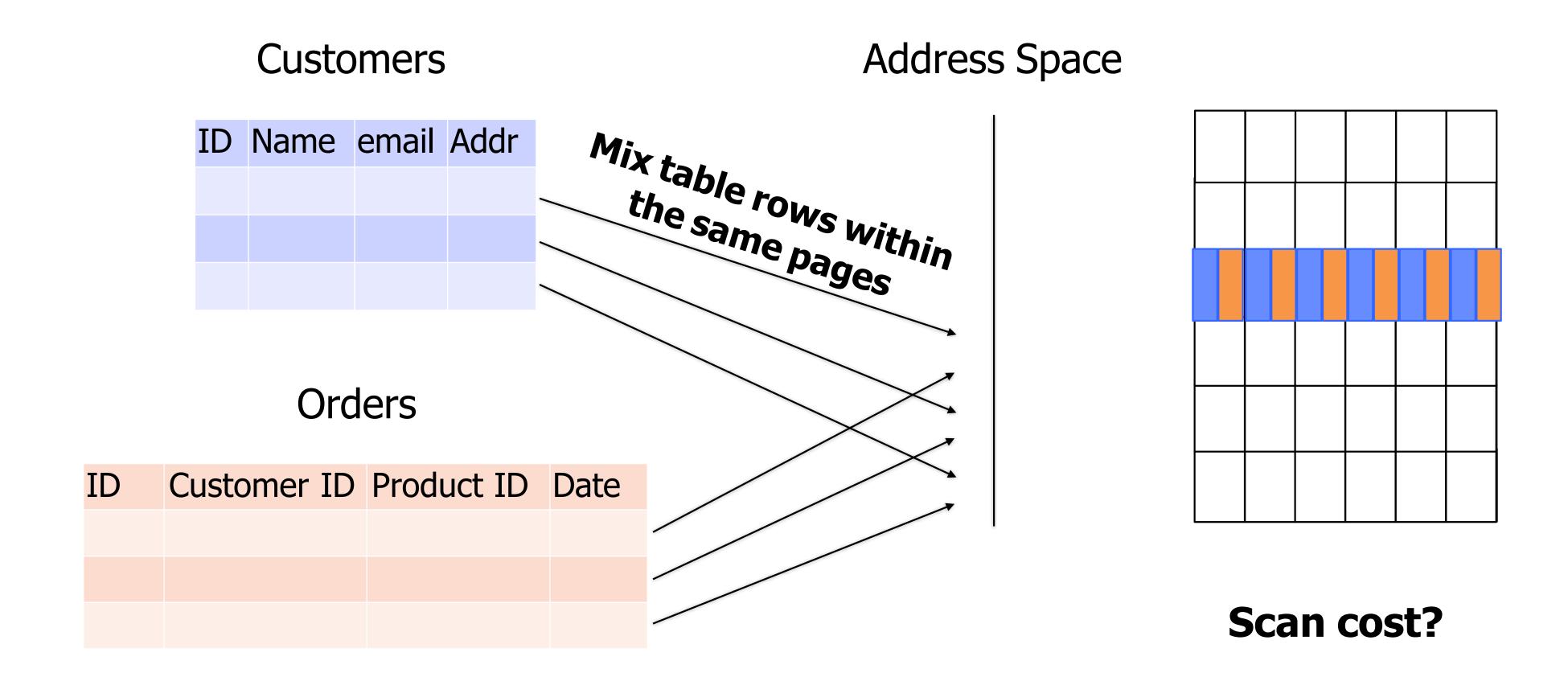
Address Space



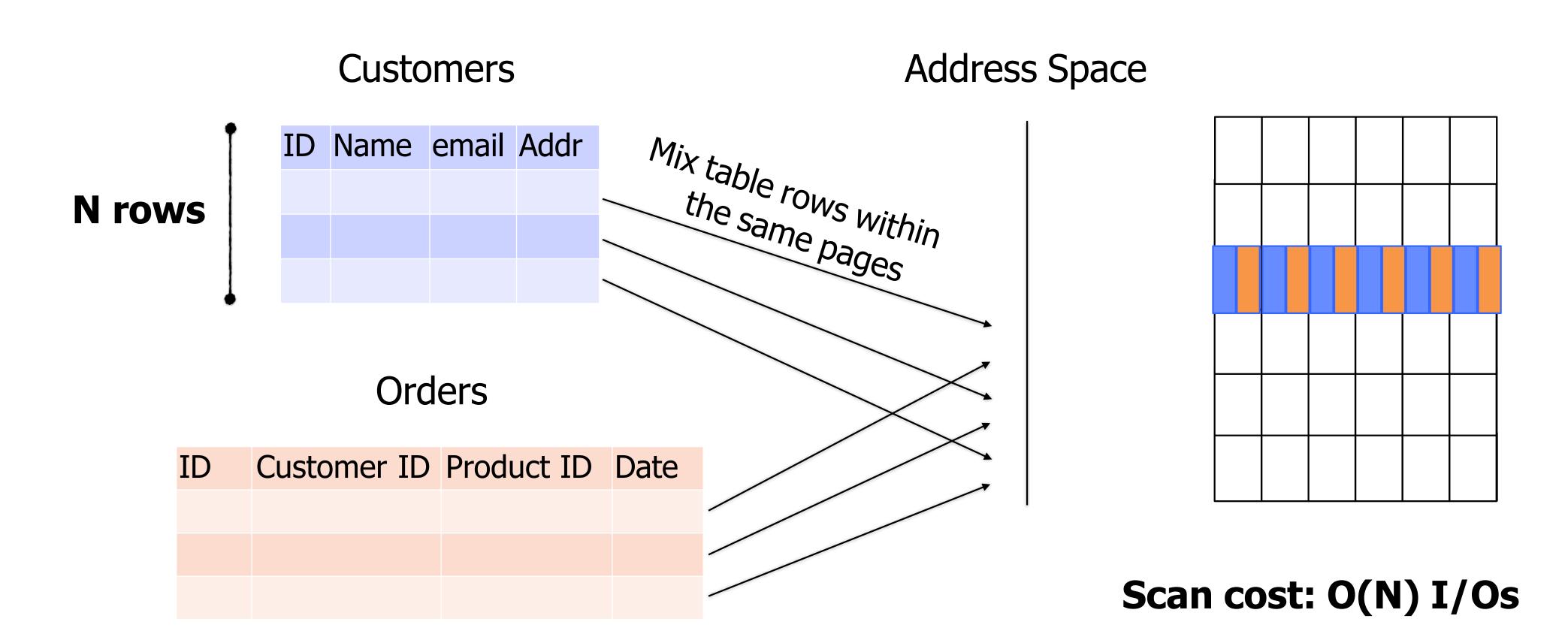
Orders

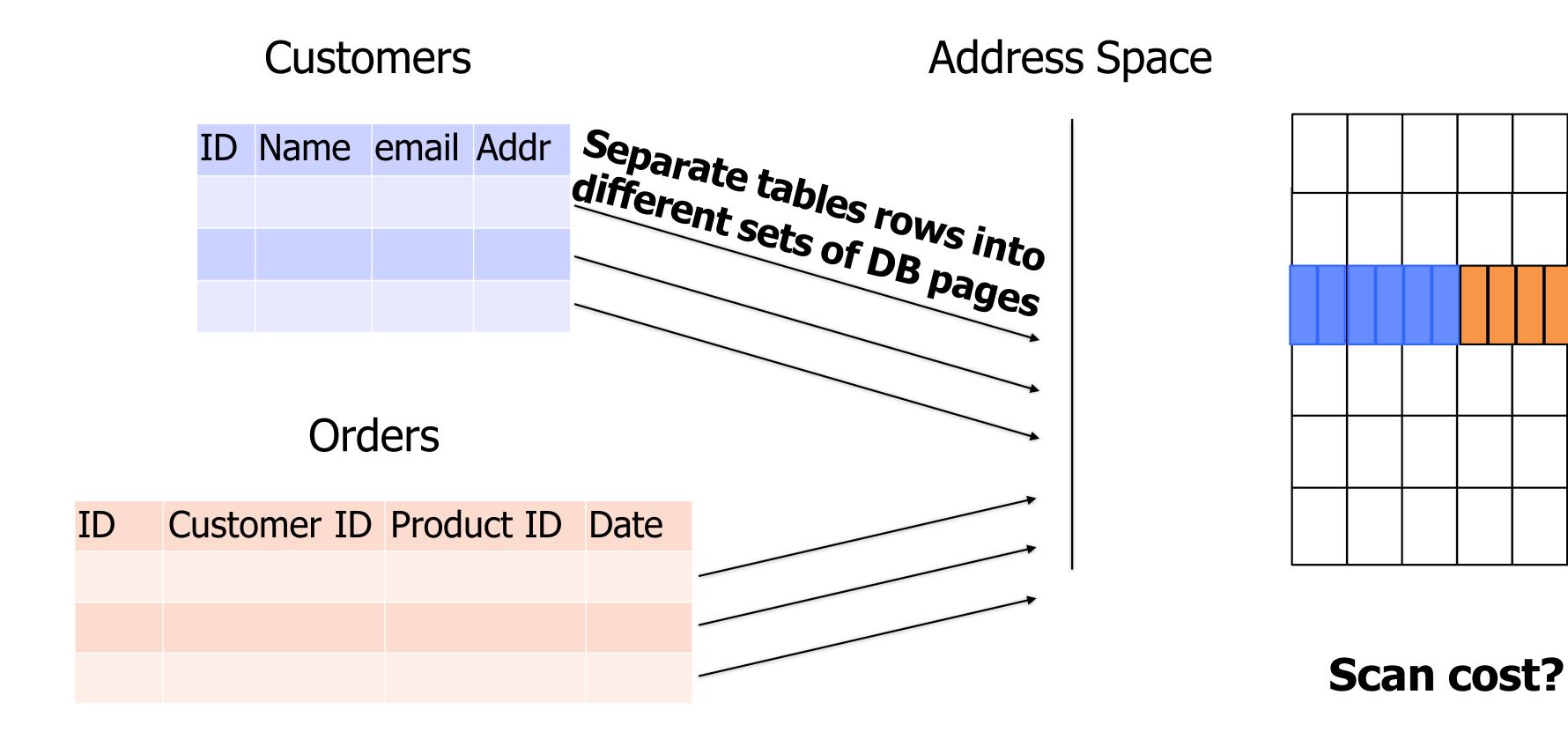


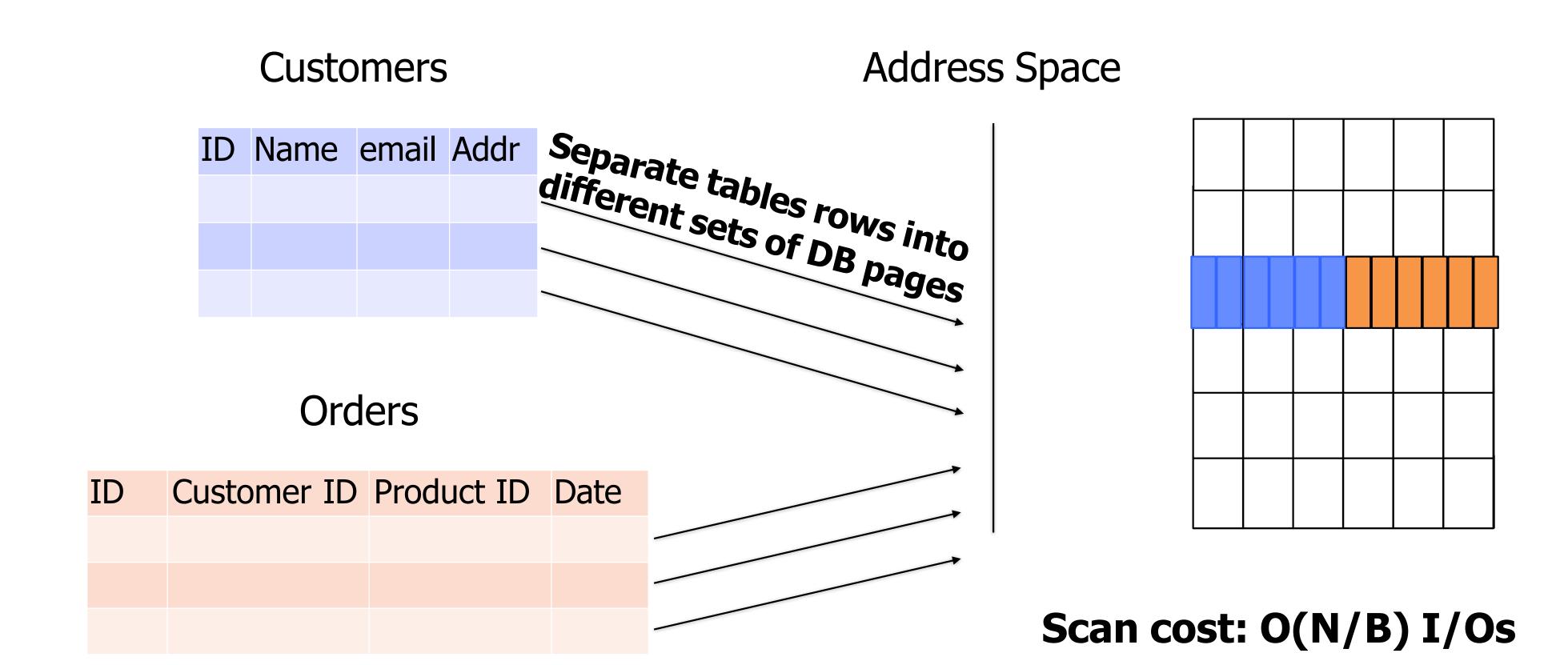
Scans - How not to Support Them



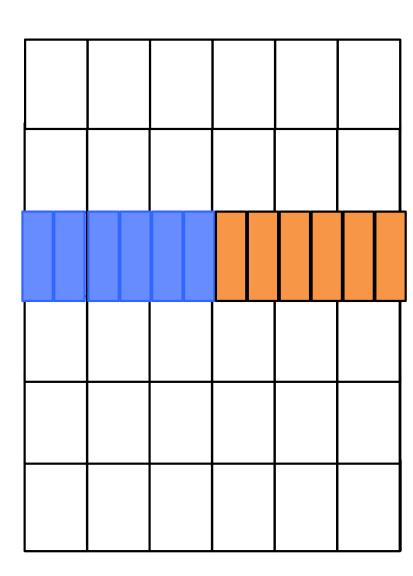
Scans - How not to Support Them





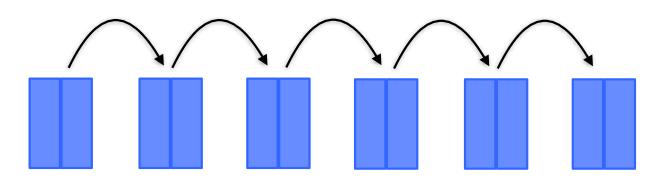


Which pages belong to which table?



Which pages belong to which table?

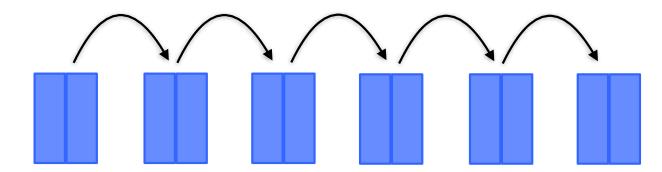
Simplest Solution: Linked List



Which pages belong to which table?

Simplest Solution: Linked List

Problem?

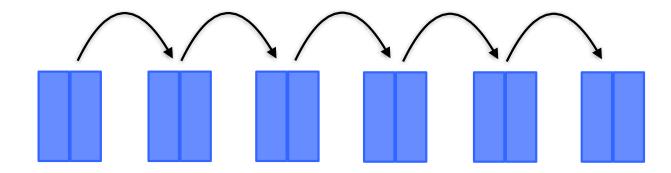


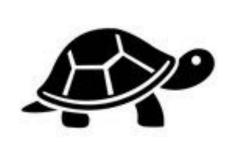
Which pages belong to which table?

Simplest Solution: Linked List

Problem: entails synchronous I/Os, which do not exploit SSD parallelism

Solution:





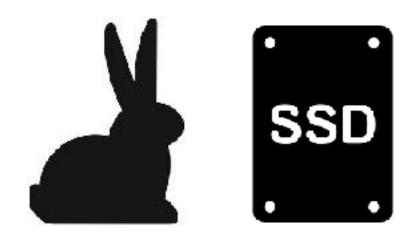


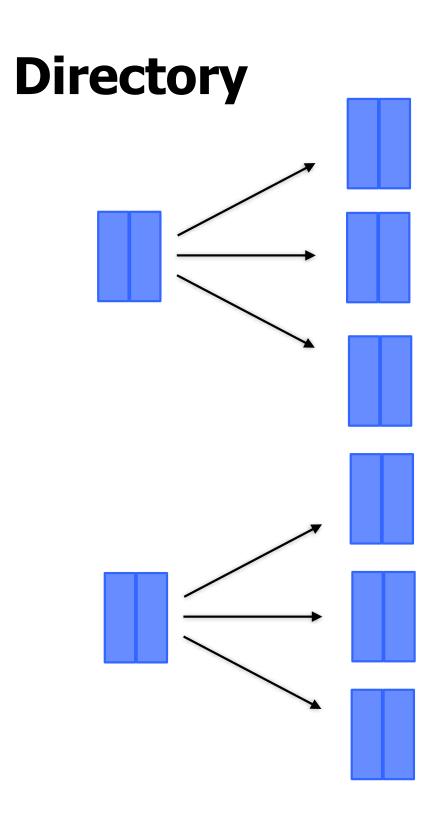
Which pages belong to which table?

Simplest Solution: Linked List

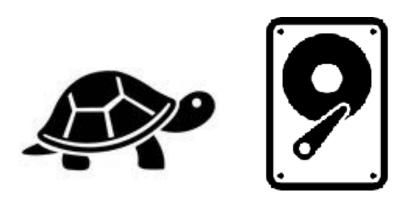
Problem: entails synchronous I/Os, which do not exploit SSD parallelism

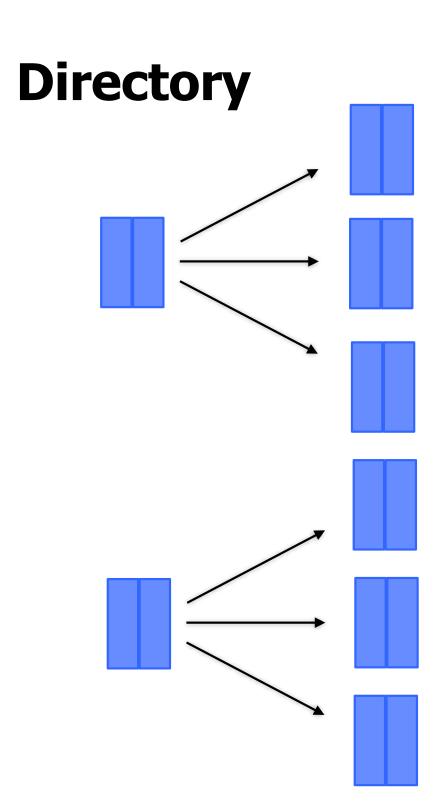
Solution: **Employ directory to allow reading many pages asynchronously**



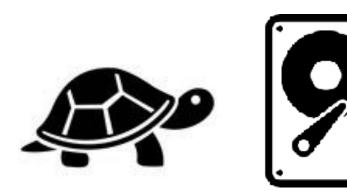


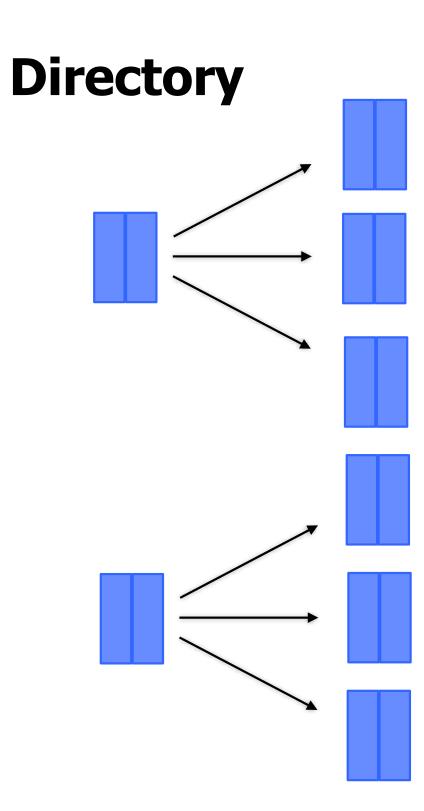
Problem for disk?





Problem: small I/Os, which do not saturate a disks's sequential bandwidth

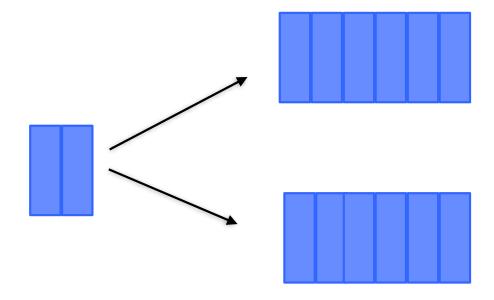




Problem: small I/Os, which do not saturate a disks's sequential bandwidth

Solution: Store multiple database pages contiguously along "extents" (8-64 pages)

Directory





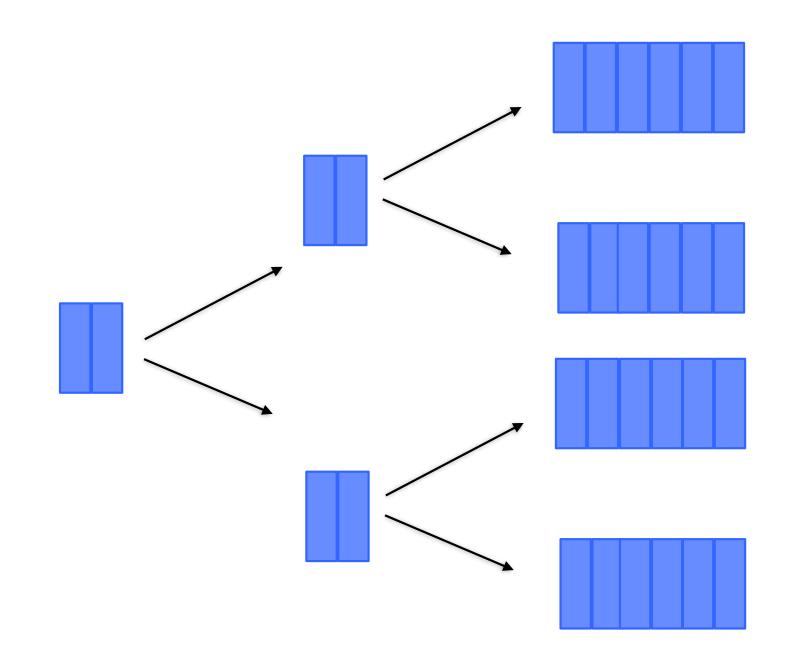


Problem: small I/Os, which do not saturate a disks's sequential bandwidth

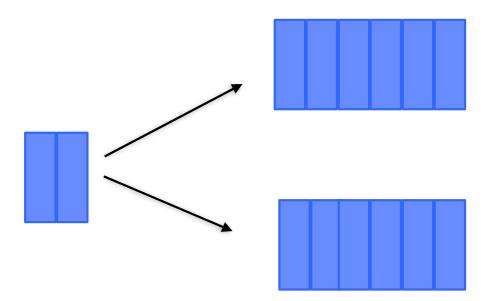
Solution: Store multiple database pages contiguously along "extents" (8-64 pages)

Bonus: Saves some metadata

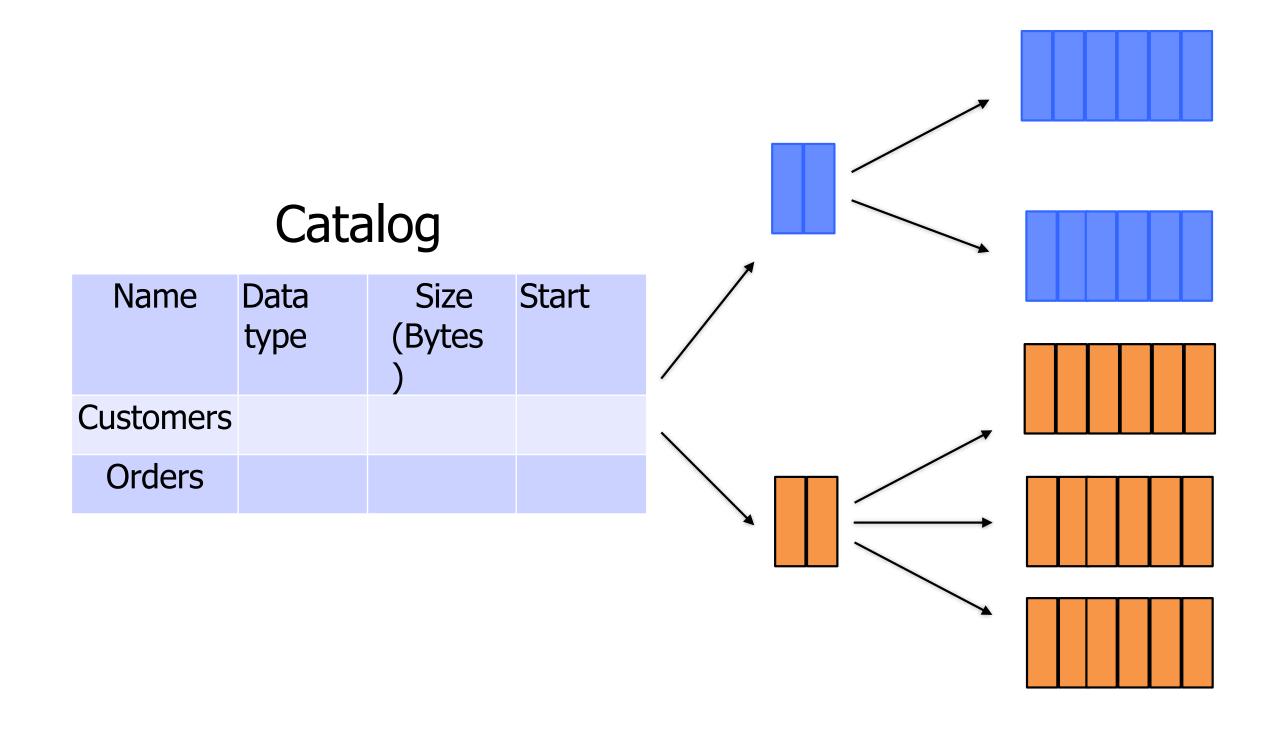
File can grow as a tree if it gets large



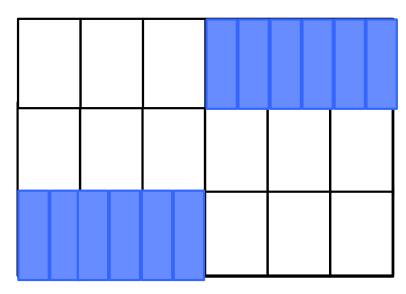
How to keep track of directories of all files?



How to keep track of directories of all files?

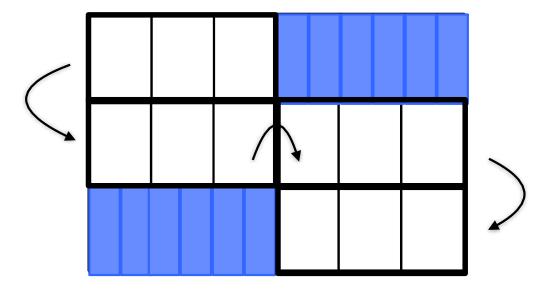


How to keep track of free pages/extents?



How to keep track of free pages/extents?

Solution 1: linked list (slower)



How to keep track of free pages/extents?

Solution 1: linked list (slower)

Solution 2: bitmap (takes space)

0 1
0 0
1 0

Operations

1. Scans

e.g., select * from Customers

2. Deletes

e.g., delete from Customers where name = "..."

3. Updates

e.g., update Customers set email = "..." where name = ""

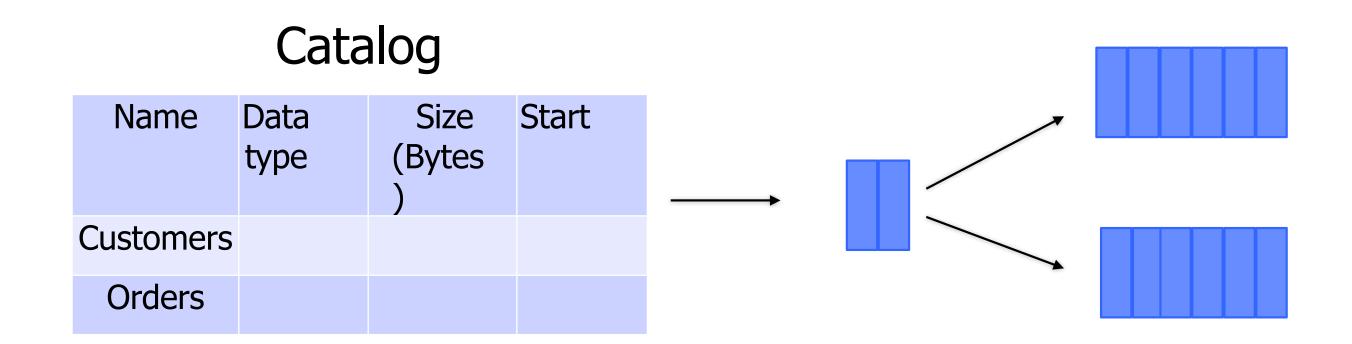
4. Insertions

e.g., Insert into Customers (,,,)

Supporting Deletes

e.g., delete from Customers where name = "..."

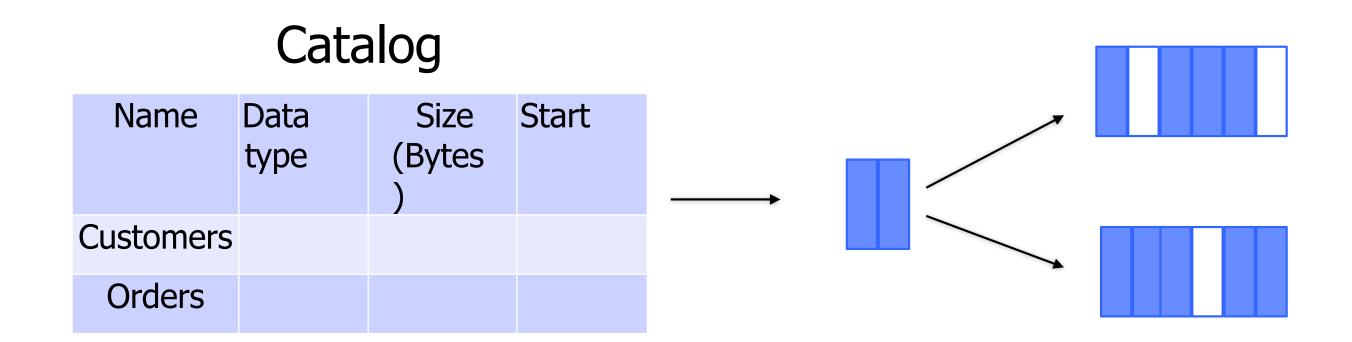
Simplest solution?



Supporting Deletes

e.g., delete from Customers where name = "..."

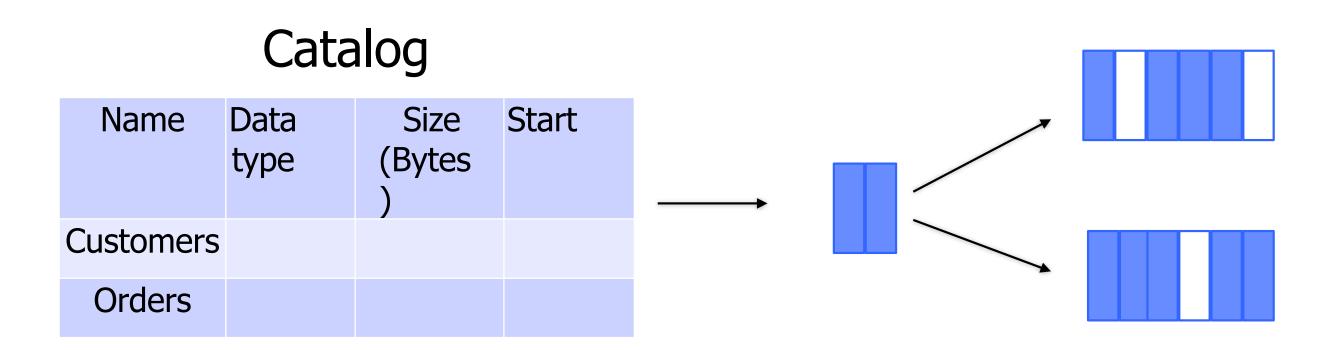
Simplest solution? Scan the table. Create "holes".



Supporting Deletes

e.g., delete from Customers where name = "..."

Simplest solution? Scan of the table. Creates "holes".



Cost: O(1) write and O(N/B) reads.

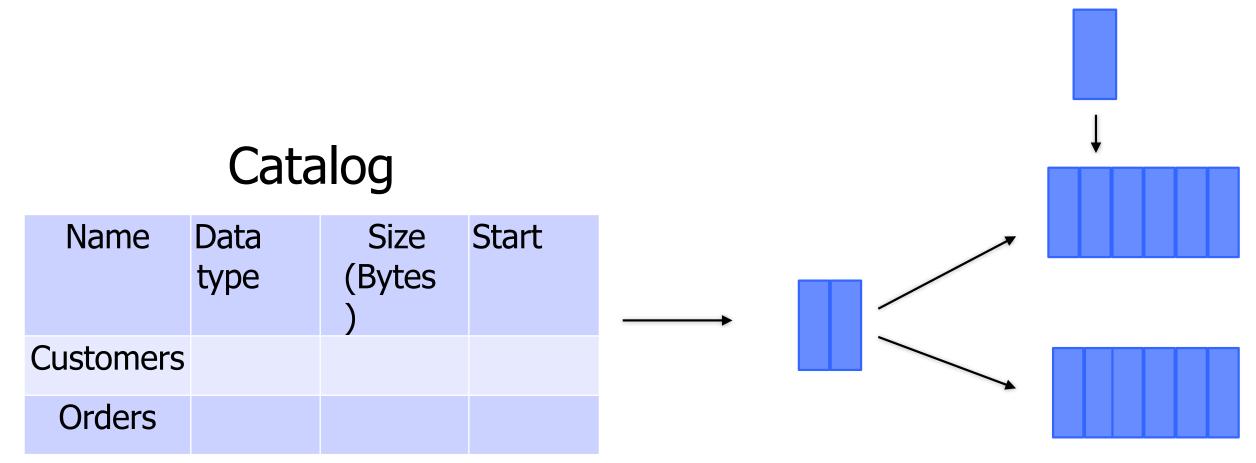
Operations

- 1. Scans
- 2. Deletes
- 3. Updates
- 4. Insertions

Supporting Updates

e.g., update Customers set email = "..." where name = ""

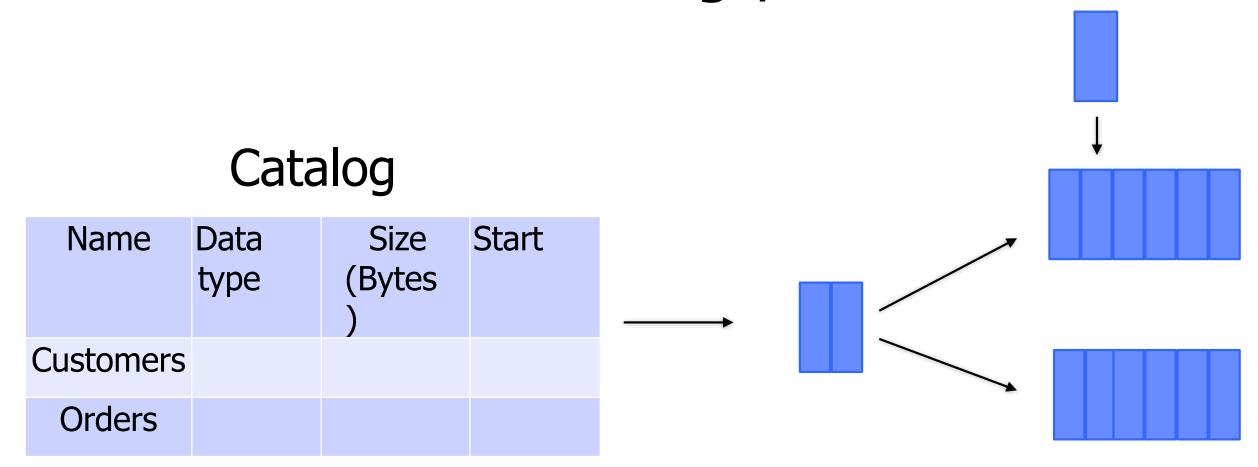
Scan and update.



Supporting Updates

e.g., update Customers set email = "..." where name = ""

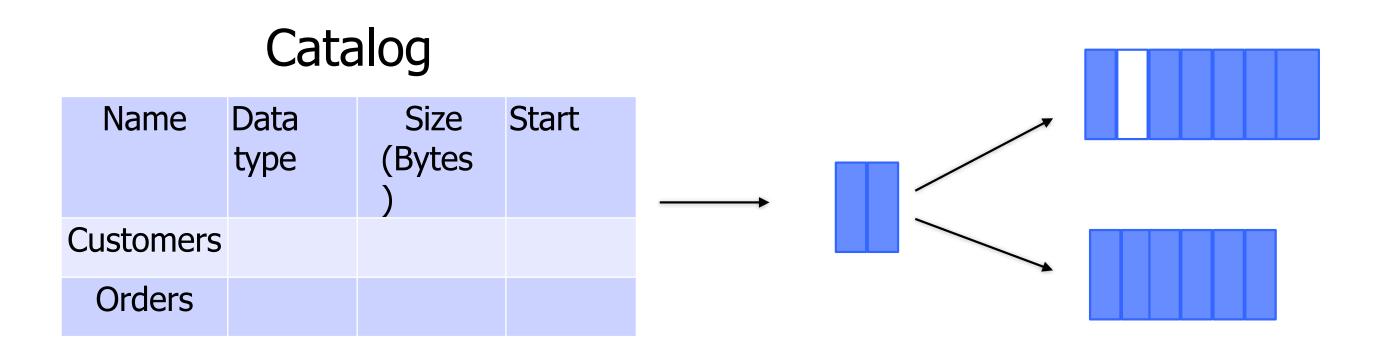
Scan and update. If newer version is too large, delete & reinsert



Supporting Updates

e.g., update Customers set email = "..." where name = ""

Scan and update. If newer version is too large, delete & reinsert



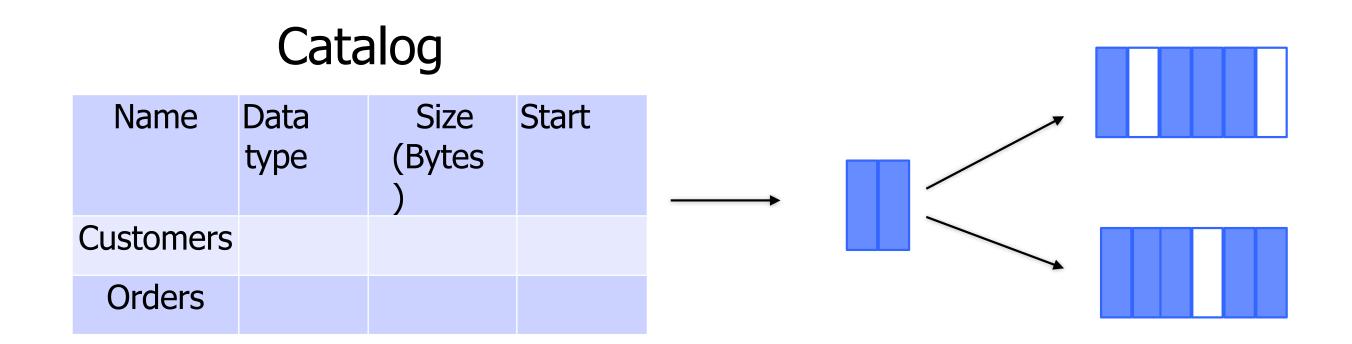
Cost: O(1) write and O(N/B) reads

Operations

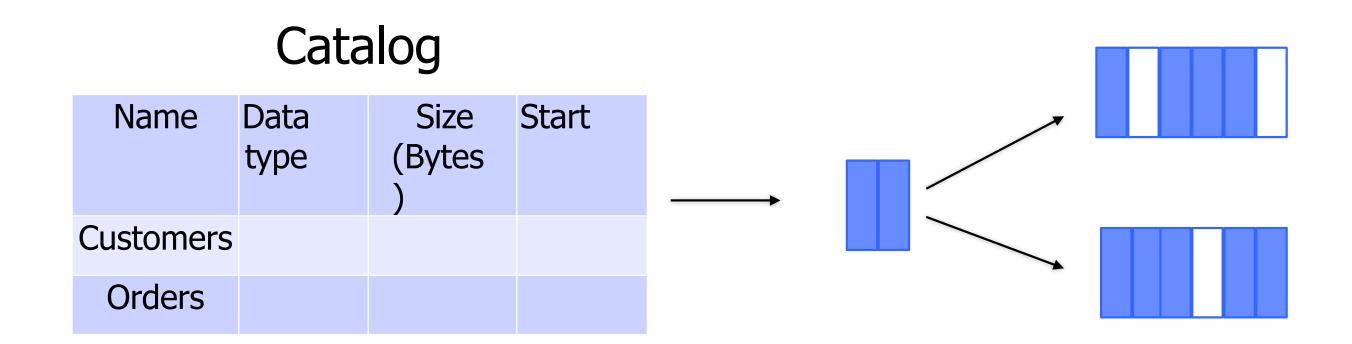
- 1. Scans
- 2. Deletes
- 3. Updates
- 4. Insertions

e.g., Insert into Customers (,,,)

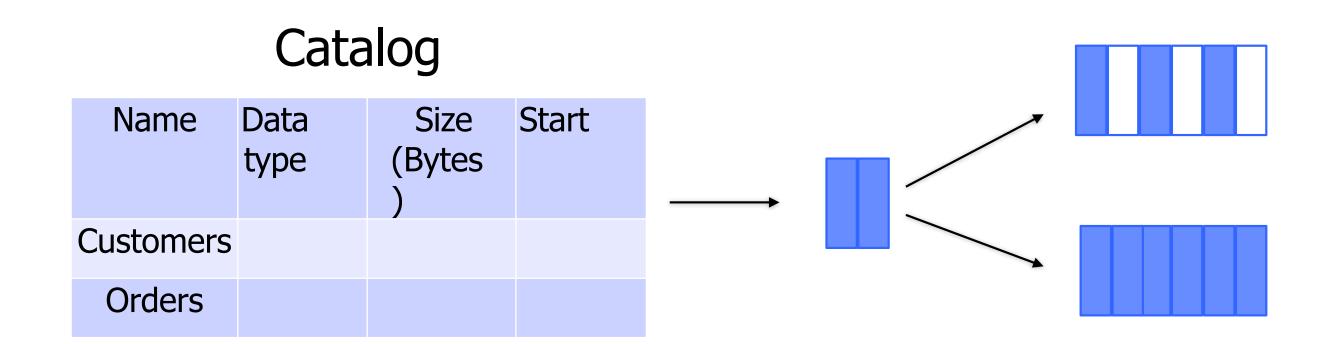
Solutions?



(1) Scan & find space. Cost: O(N/B) reads and O(1) write.

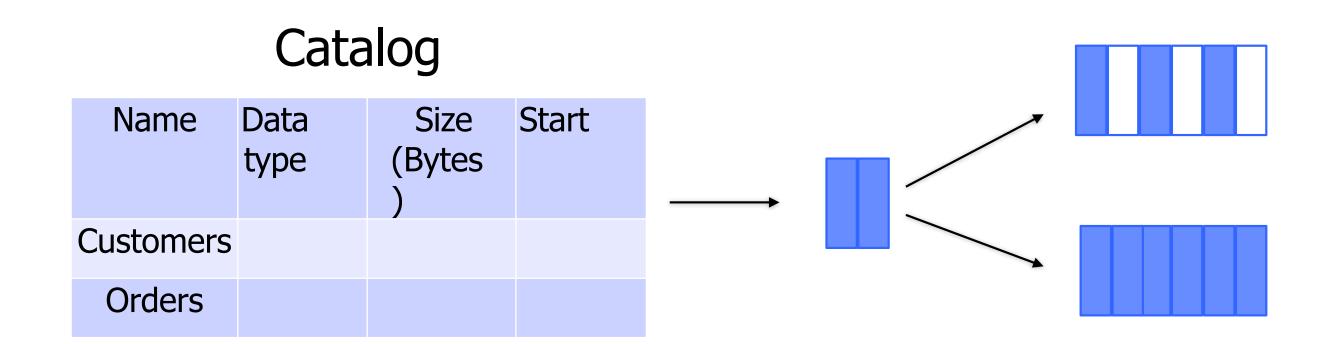


- (1) Scan & find space. Cost: O(N/B) reads and O(1) write.
- (2) Separate Linked list of pages with free space.



- (1) Scan & find space. Cost: O(N/B) reads and O(1) write.
- (2) Separate Linked list of pages with free space.

 Cost: O(1) reads & O(1) write for fixed-sized entries

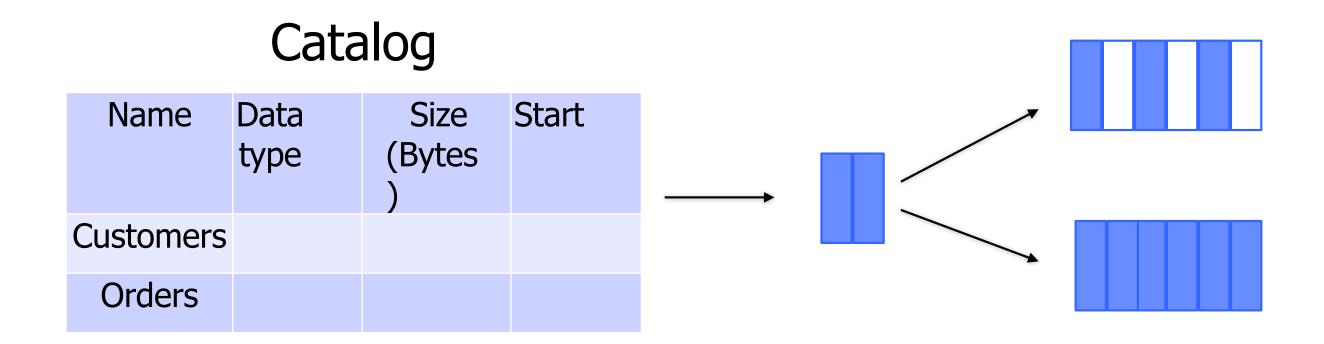


(1) Scan & find space. Cost: O(N/B) reads and O(1) write.

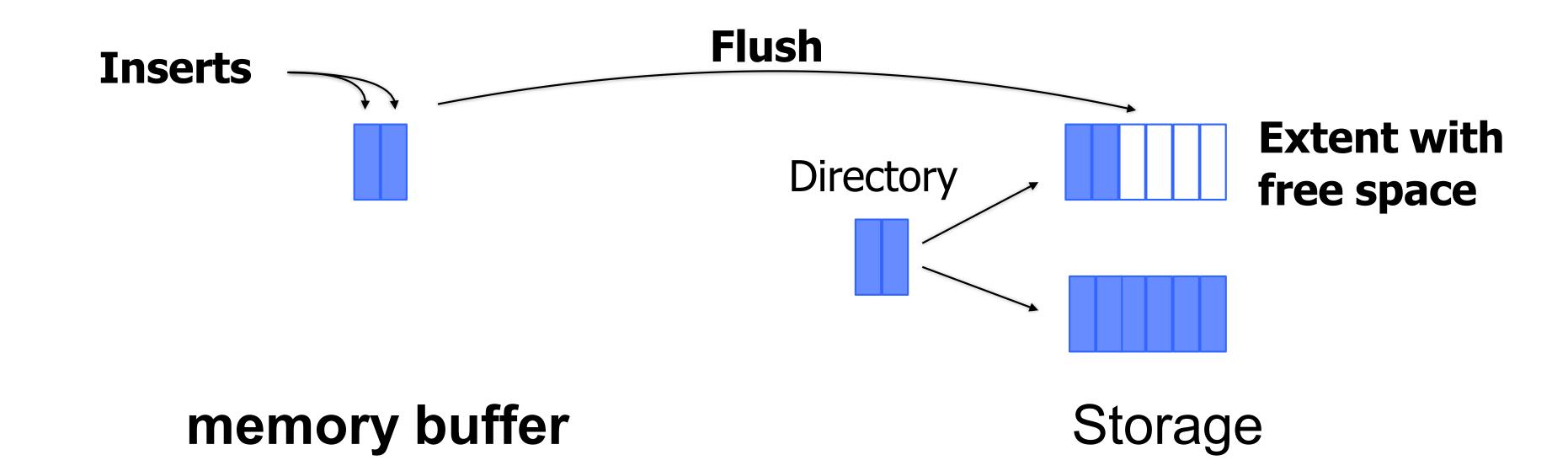
(2) Separate Linked list of pages with free space.

Cost: O(1) reads & O(1) write for fixed-sized entries

Cost: O(N/B) reads & O(1) write for variable-sized entries

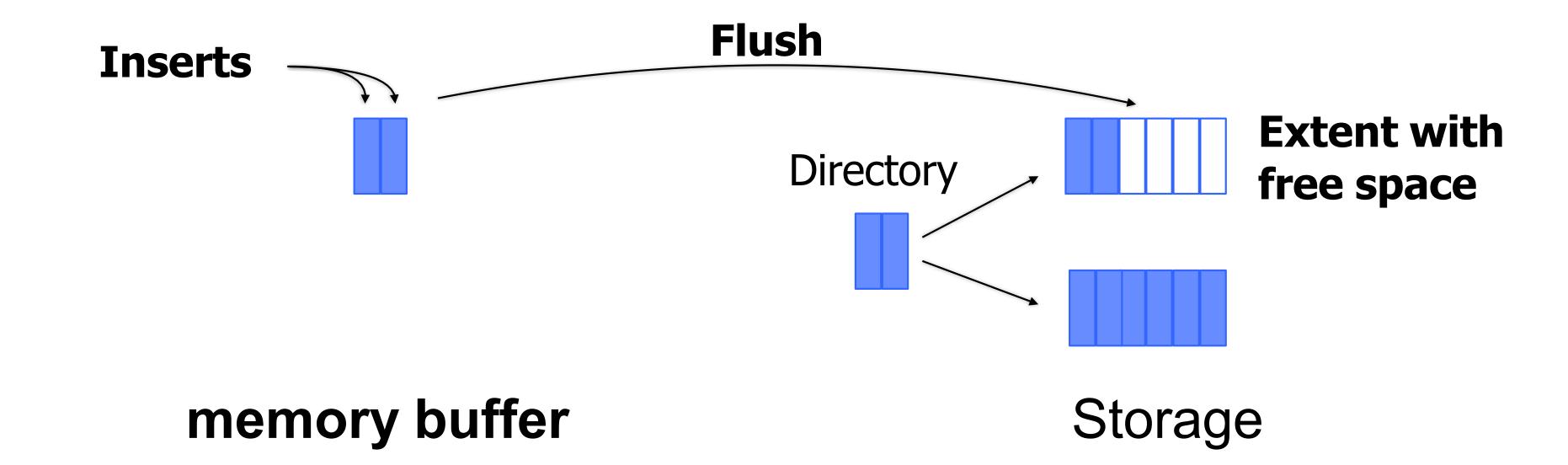


(3) buffer insertions in memory until a page fills up & append to extent



(3) buffer insertions in memory until a page fills up & append to extent

Cost: No reads and O(1/B) of a write



- (1) Scan & find space. Cost: O(N/B) reads and O(1) write.
- (2) Separate Linked list of pages with free space.Cost: O(1) reads & O(1) write for fixed-sized entries Cost:O(N/B) reads & O(1) write for variable-sized entries
- (3) buffer insertions in memory until a page fills up & append to extent Cost: No reads and O(1/B) of a write

Recall each page is 4 KB

Suppose rows are fixed-sized

How to organize rows within a slot?

Recall each page is 4 KB

Suppose rows are fixed-sized

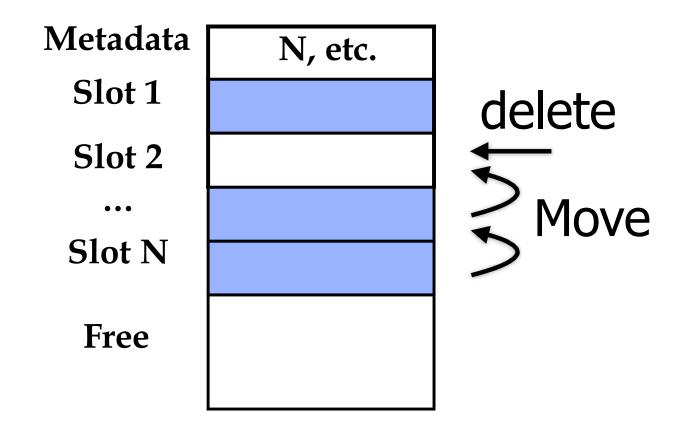
How to organize rows within a slot?

Metadata	N, etc.
Slot 1	
Slot 2	
• • •	
Slot N	
Free	

Recall each page is 4 KB

Suppose rows are fixed-sized

How to organize rows within a slot?

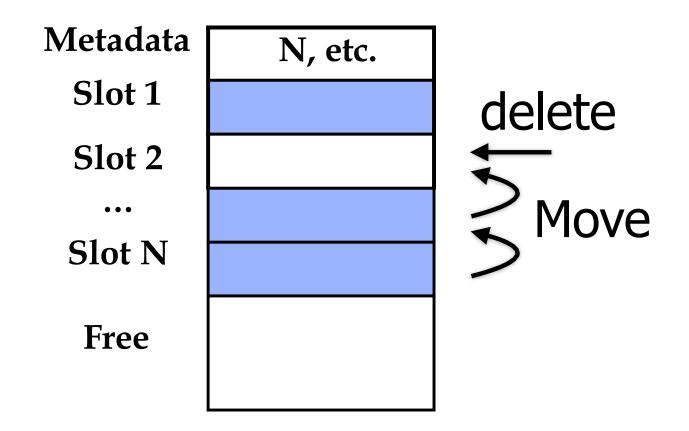


Need to reorganize due to deletes

Recall each page is 4 KB

Suppose rows are fixed-sized

How to organize rows within a slot?



Metadata
Free Bitmap
Slot 1
Slot 2
Slot 3
...
Slot N

Need to reorganize due to deletes

No reorganization, requires more space

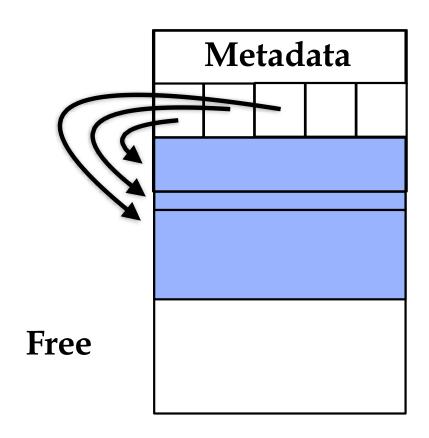
Recall each page is 4 KB

Suppose rows are variable-length

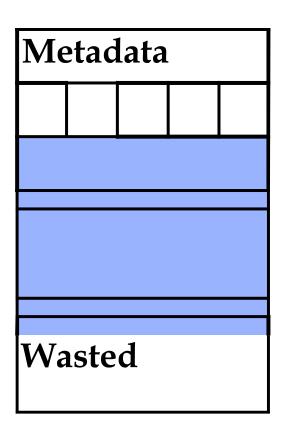
Solutions?

Recall each page is 4 KB

Suppose rows are **variable-length**

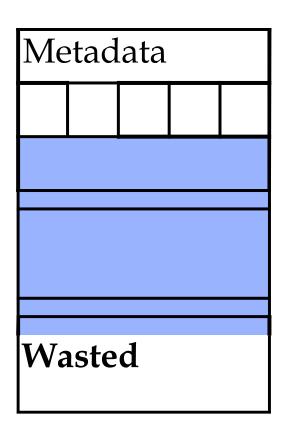


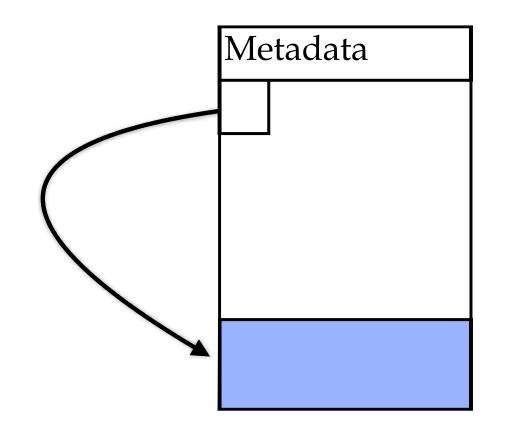
Recall each page is 4 KB
Suppose rows are **variable-length**



If entries are small, we waste space at the end, or we must push all content up to clear space

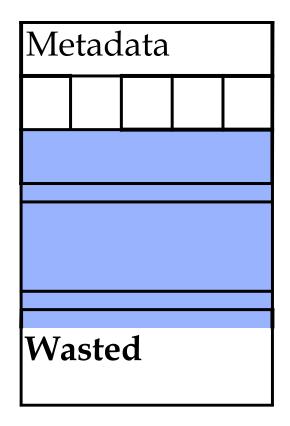
Recall each page is 4 KB
Suppose rows are **variable-length**

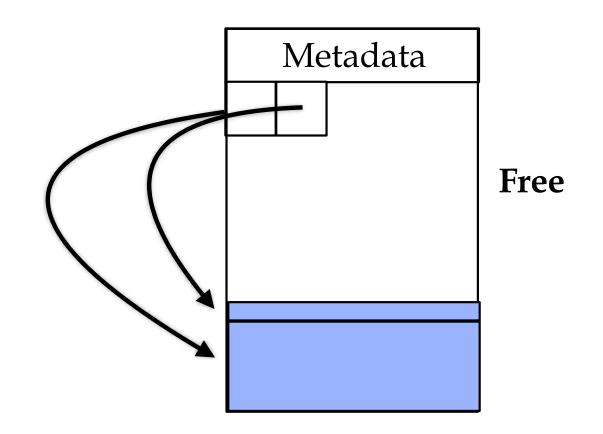




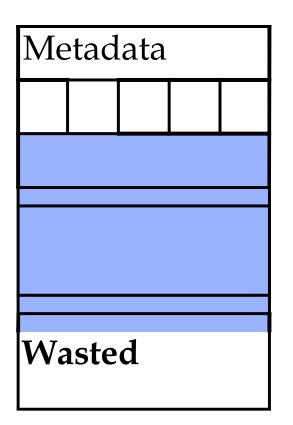
Store data from end of page

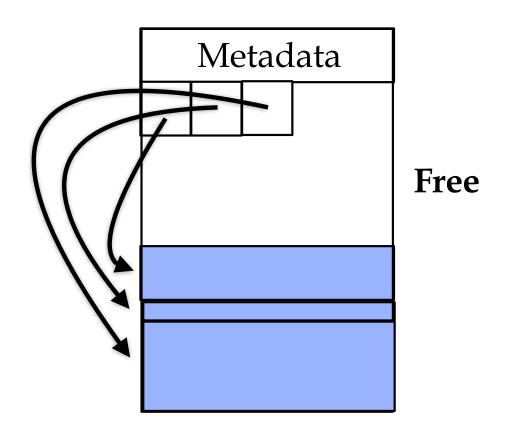
Recall each page is 4 KB
Suppose rows are variable-length





Recall each page is 4 KB
Suppose rows are **variable-length**





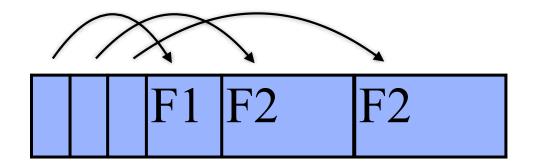
Minimal space wastage, and no need to move data

Variable-Sized Record Organization

Delimiters

Pointers



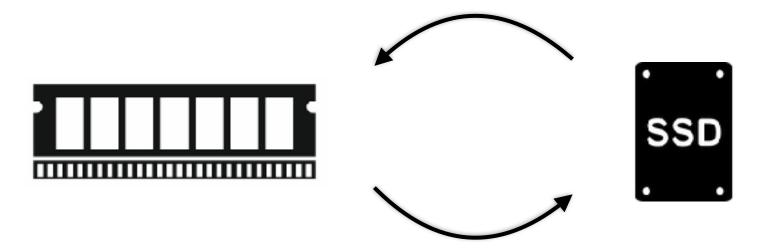


Smaller No random access More space Random access (faster)

Break

Then let's now move to buffer management

Buffer Management



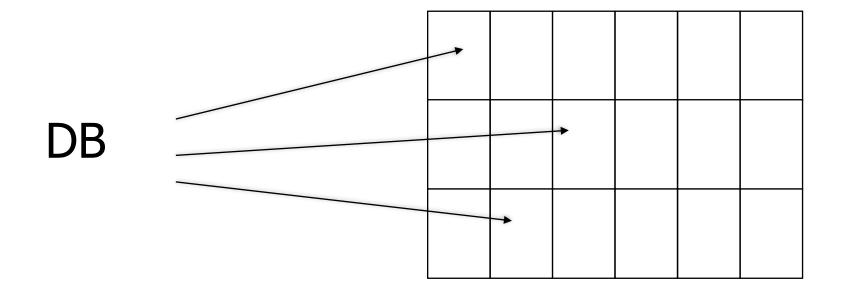
Database System Technology
Niv Dayan

Context

A DB is reading and writing aligned 4KB storage pages







Storage

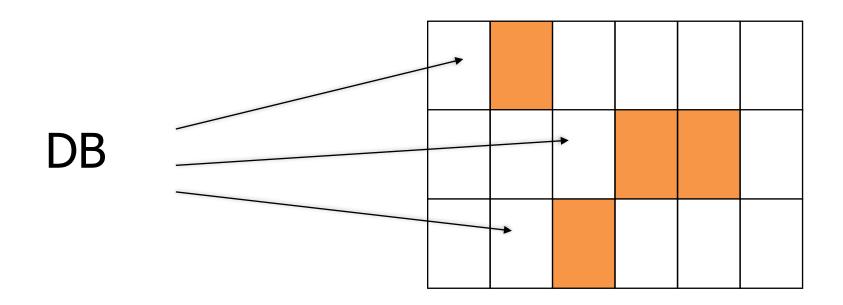
Context

A DB is reading and writing aligned 4KB storage pages

Suppose orange pages are frequently accessed ("hot")







Storage

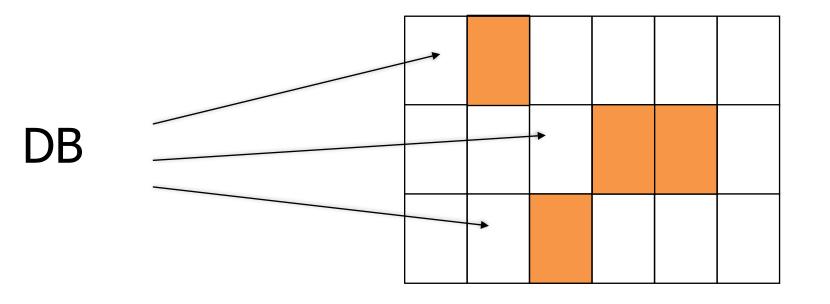
Context

A DB is reading and writing aligned 4KB storage pages

Suppose orange pages are frequently accessed ("hot")







Retrieving these pages over and over is expensive!

Storage

Keep copies of hot pages in memory







DB —

memory

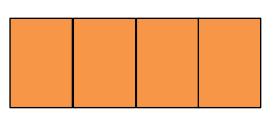
Storage

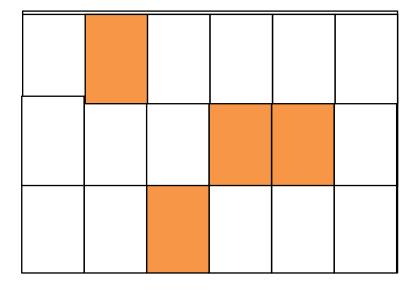
How to structure this buffer pool?











memory

Storage

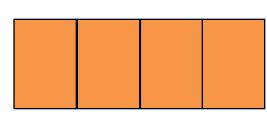
How to structure this buffer pool?

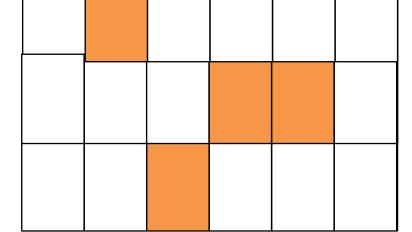
hash table (more details later)







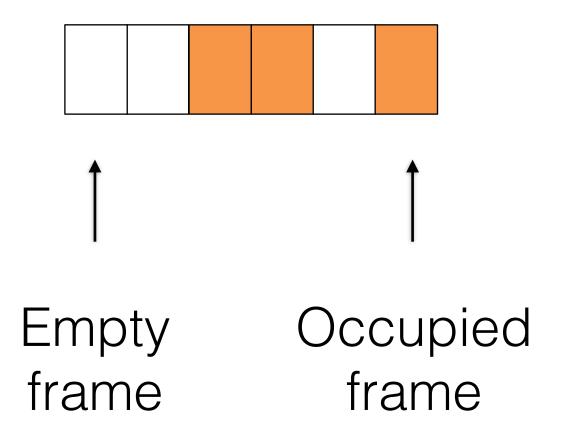




memory

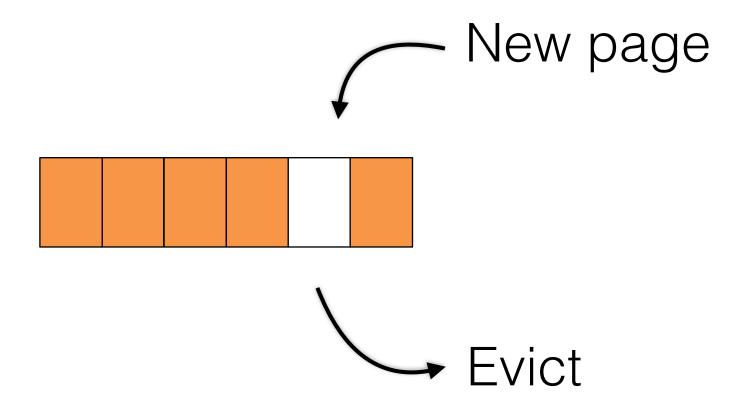
Storage

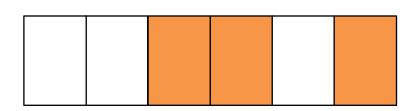
Consist of frames, each containing one page of data (e.g., 4 KB)



Consist of frames, each containing one page of data (e.g., 4 KB)

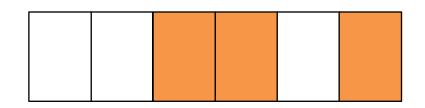
Eventually it fills up. Must evict pages to clear space.





Each frame must keep some metadata

- (1)?
- (2)?

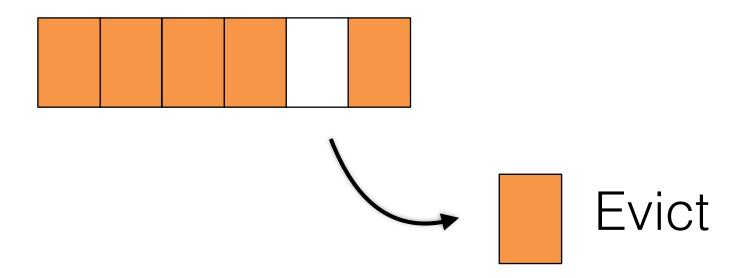


Each frame must keep some metadata

- (1) Pin count How many users are currently using this page
- (2) Dirty flag indicates whether the page has been updated

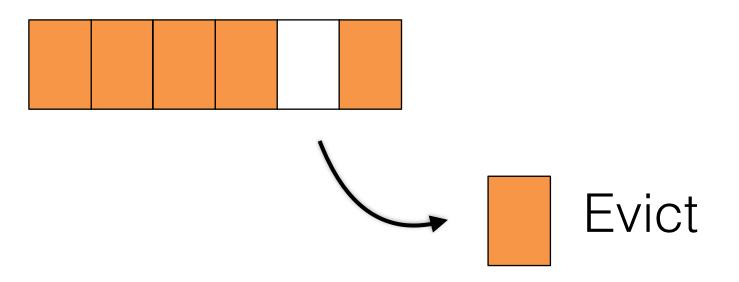
Eviction Policy

Which page to evict when we run out of space?



Eviction Policy

Which page to evict when we run out of space?

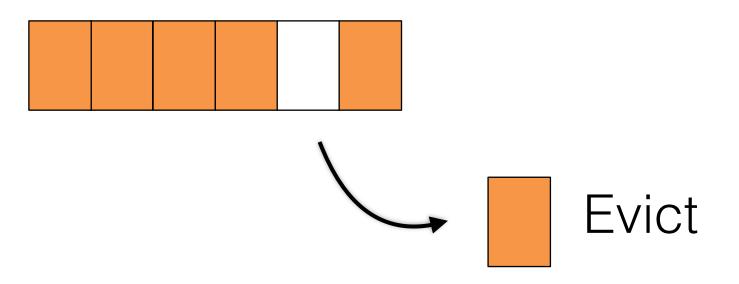


Considerations:

- (1) Avoid evicting a page that is likely to be used again
- (2) Avoid excessive metadata or CPU overheads to make decision

Eviction Policy

Which page to evict when we run out of space?



Big impact on number of I/Os and CPU efficiency

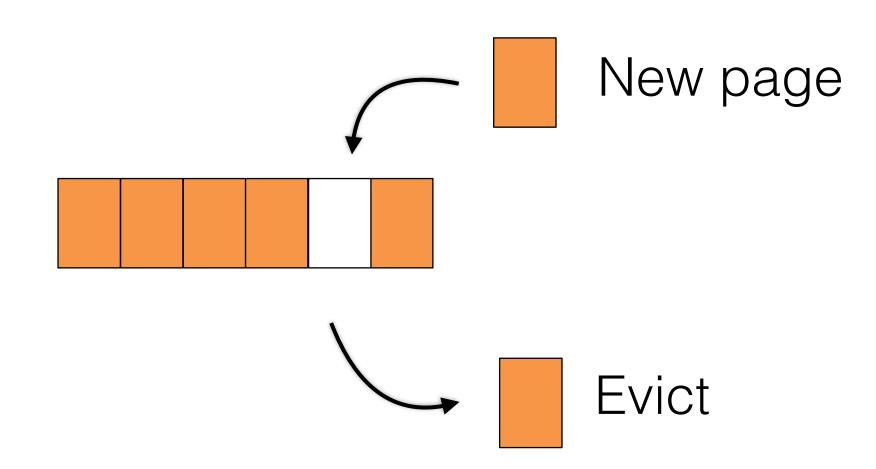
Depends on the access pattern

We'll cover 5 eviction policies

Evict whichever page collides in the hash table with a new page

Pros: ?

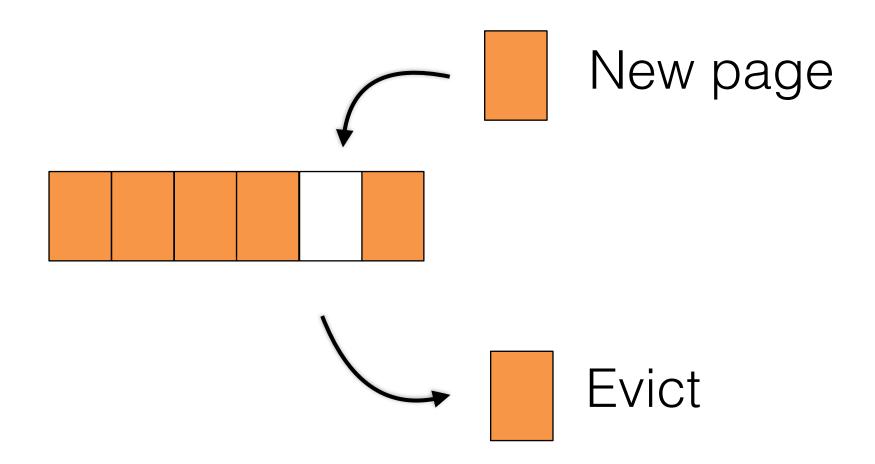
Con: ?



Evict whichever page collides in the hash table with a new page

Pros: Simple, CPU-efficient, no extra metadata

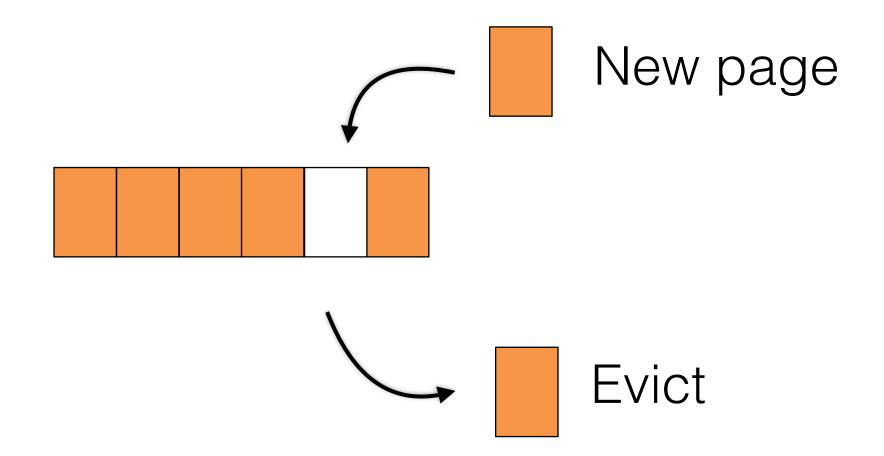
Con: ?



Evict whichever page collides in the hash table with a new page

Pros: Simple, CPU-efficient, no extra metadata

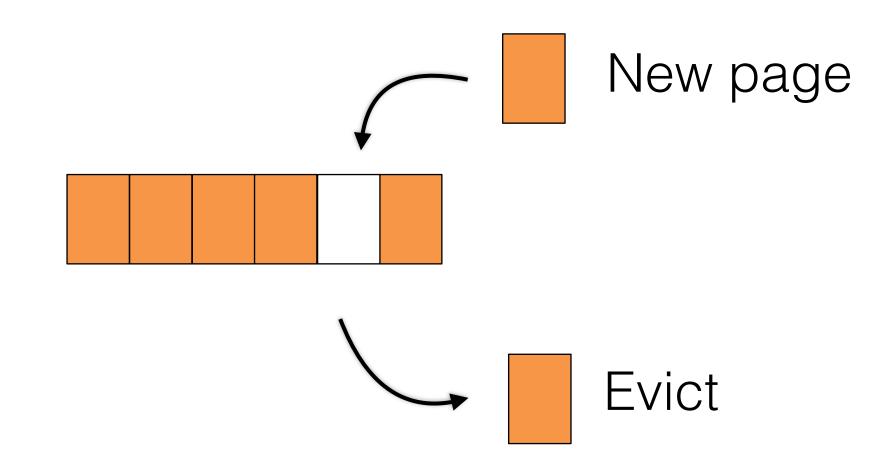
Con: May evict a frequently used page



Evict whichever page collides in the hash table with a new page

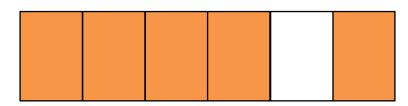
Pros: Simple, CPU-efficient, no extra metadata

Con: May evict a frequently used page
Can we improve this?



Evict Page that was inserted the longest time ago

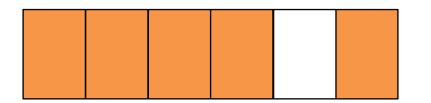
Rationale?



Evict Page that was inserted the longest time ago

Rationale? Less likely to be used again

Implementation?

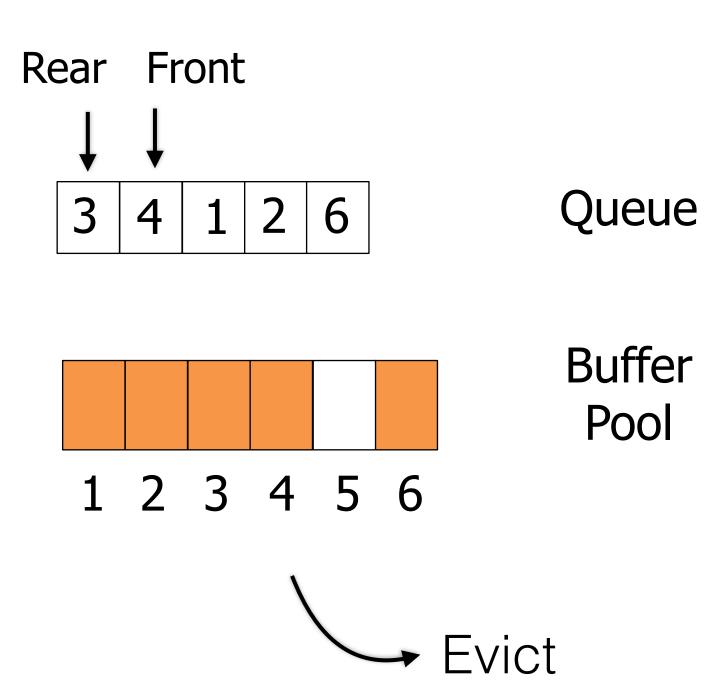


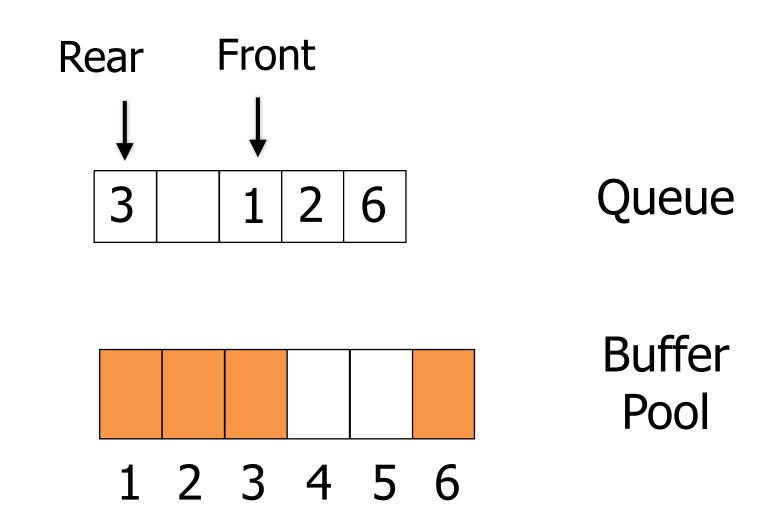
Evict Page that was inserted the longest time ago

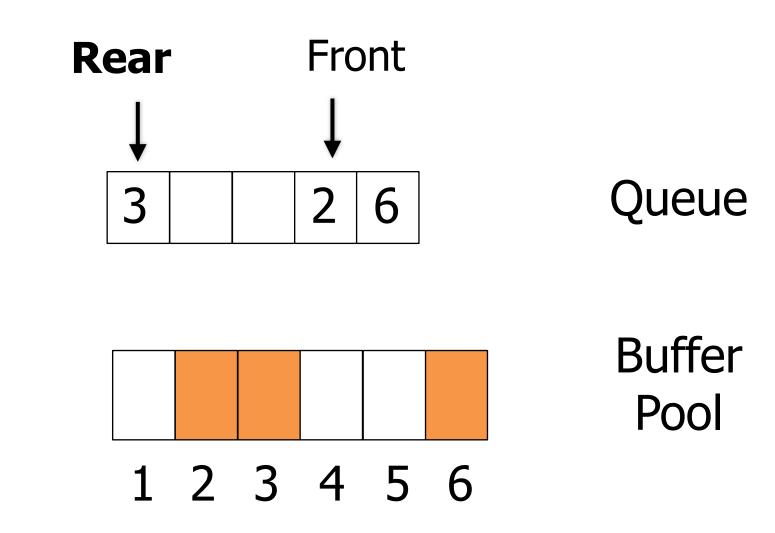
Rationale? Less likely to be used again

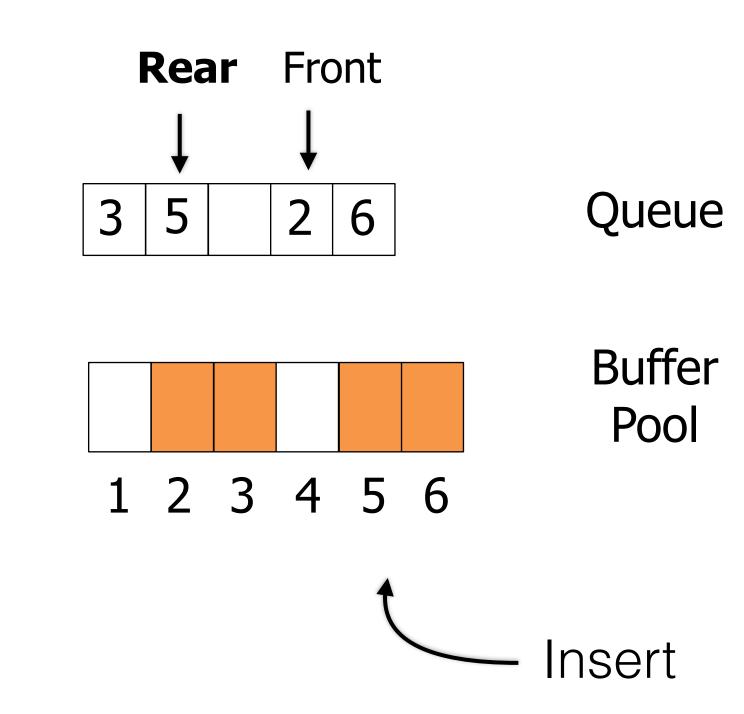
Implementation? Using a queue

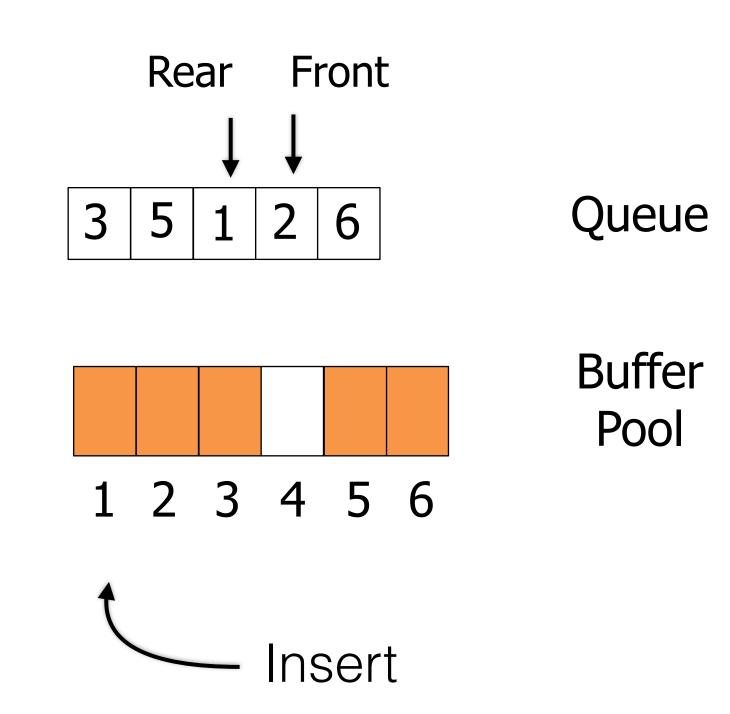
(i.e., array with front/rear pointers)



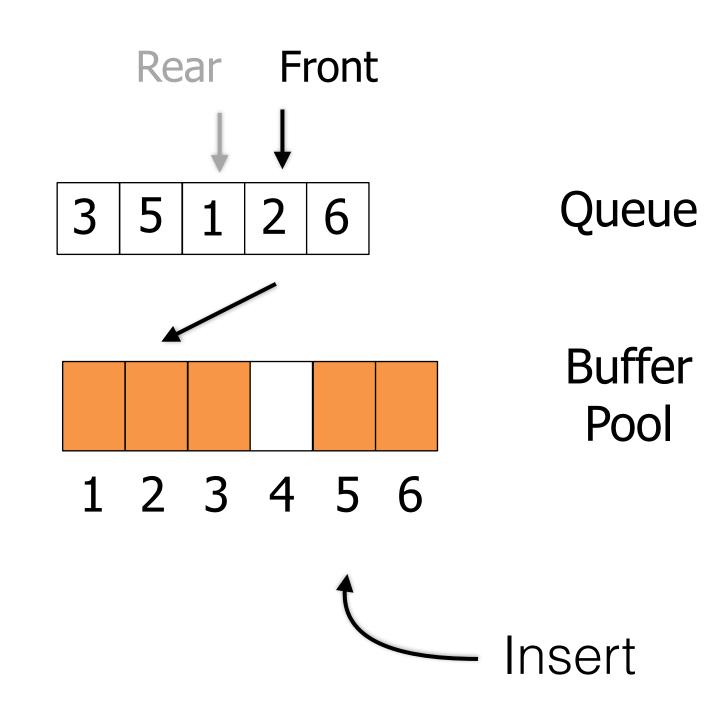






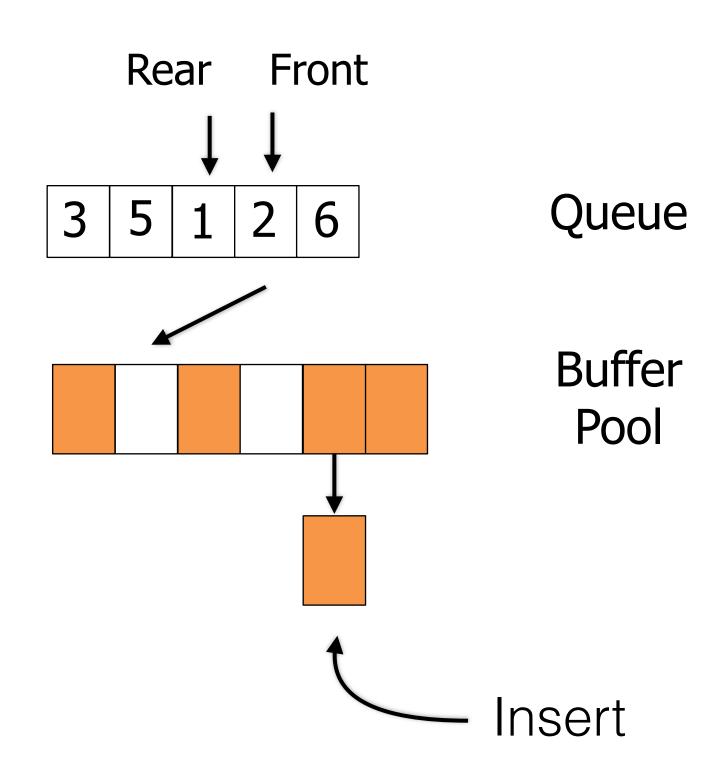


Unlike the random policy, pages we evict now have different frames than pages we insert.



Unlike the random policy, pages we evict now have different frames than pages we insert.

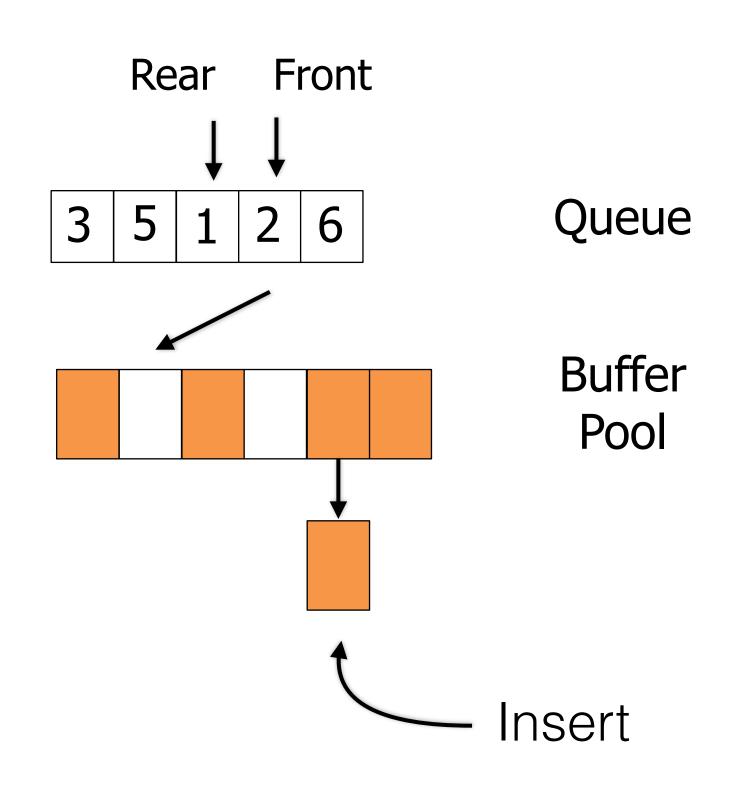
We need a hash collision resolution algo. (e.g., linear probing, chaining)



Unlike the random policy, pages we evict now have different frames than pages we insert.

We need a hash collision resolution algo. (e.g., linear probing, chaining)

Need 10-20% extra capacity in the table to reduce likelihood of collisions

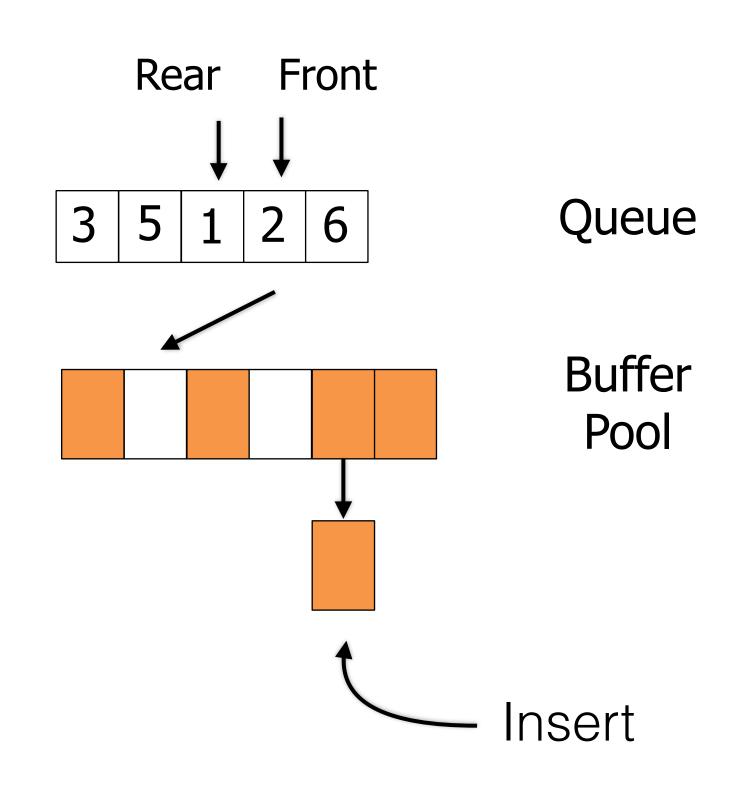


Unlike the random policy, pages we evict now have different frames than pages we insert.

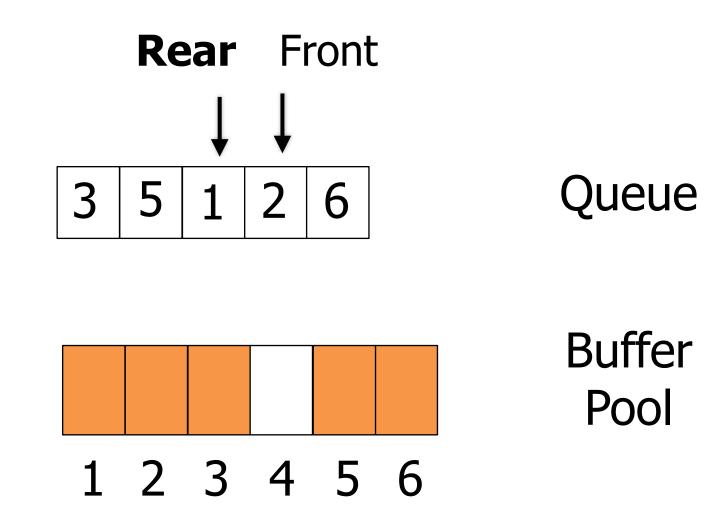
We need a hash collision resolution algo. (e.g., linear probing, chaining)

Need 10-20% extra capacity in the table to reduce likelihood of collisions

Extra space and CPU cost, but necessary for all policies but the random one

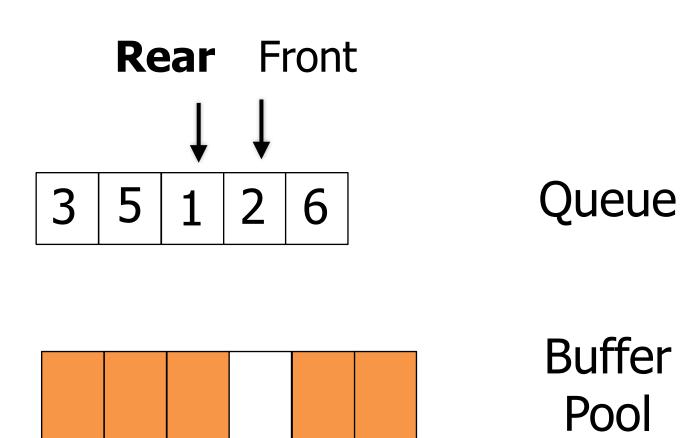


Problem?



Problem?

Oldest page may still be frequently used.

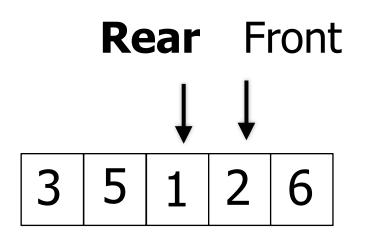


1 2 3 4 5 6

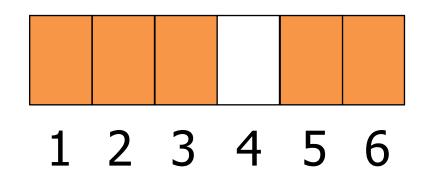
Problem?

Oldest page may still be frequently used.

Can we address this?



Queue



Buffer Pool

Evict page that was used last the longest time ago

Evict page that was used last the longest time ago

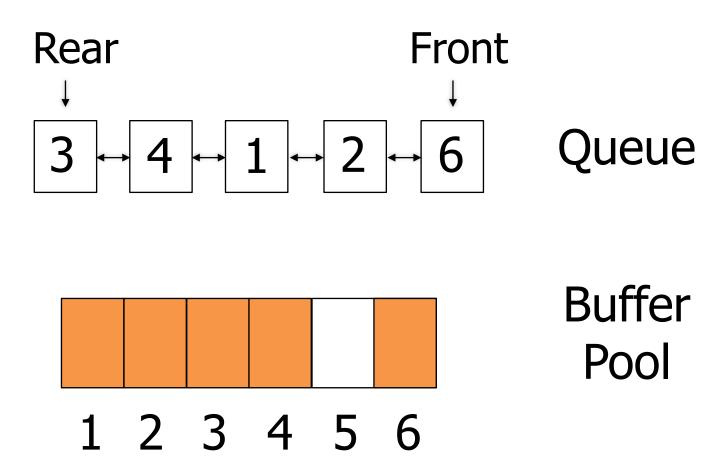
Implementation?

Evict page that was used last the longest time ago

Implementation? A random-access queue

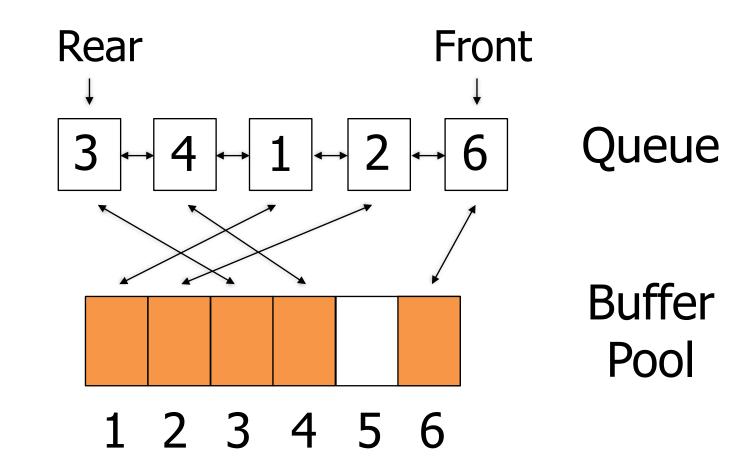
i.e., We must be able to move an element at a random location back to the front

Use a doubly-linked list as the queue



Use a doubly-linked list as the queue

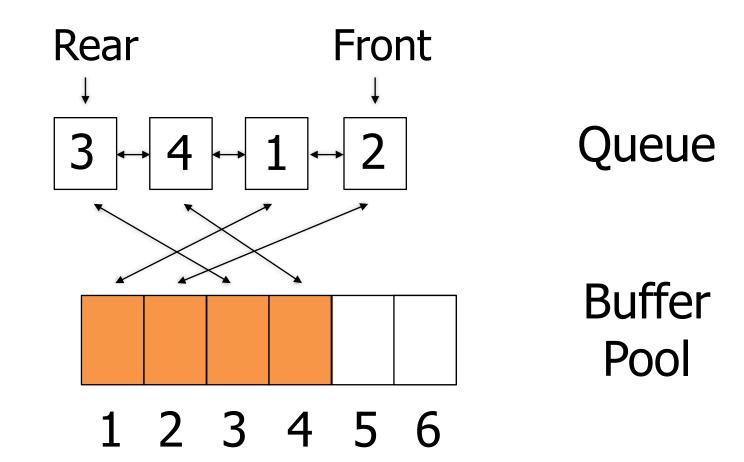
With pointers linking nodes and buckets



Use a doubly-linked list as the queue

With pointers linking nodes and buckets

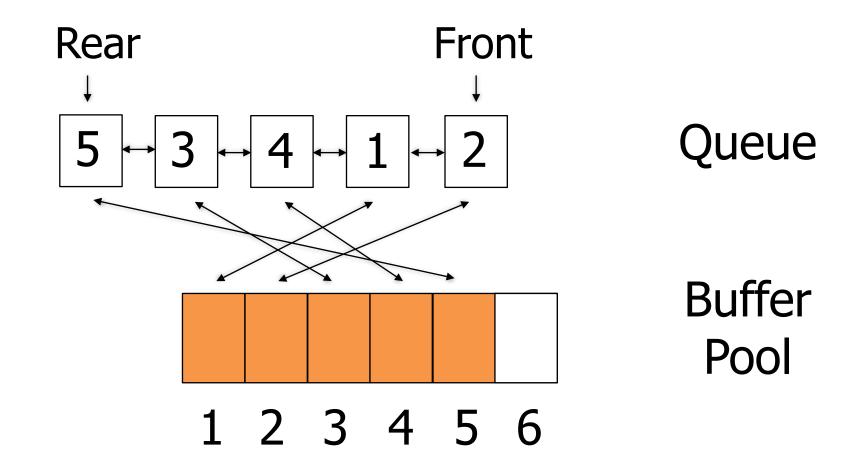
Load to rear and evict front (as before)



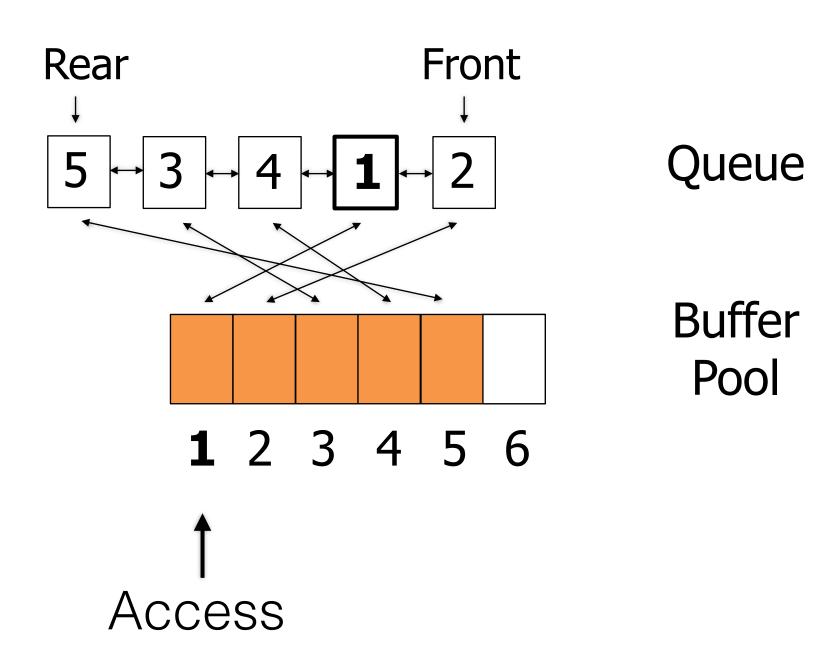
Use a doubly-linked list as the queue

With pointers linking nodes and buckets

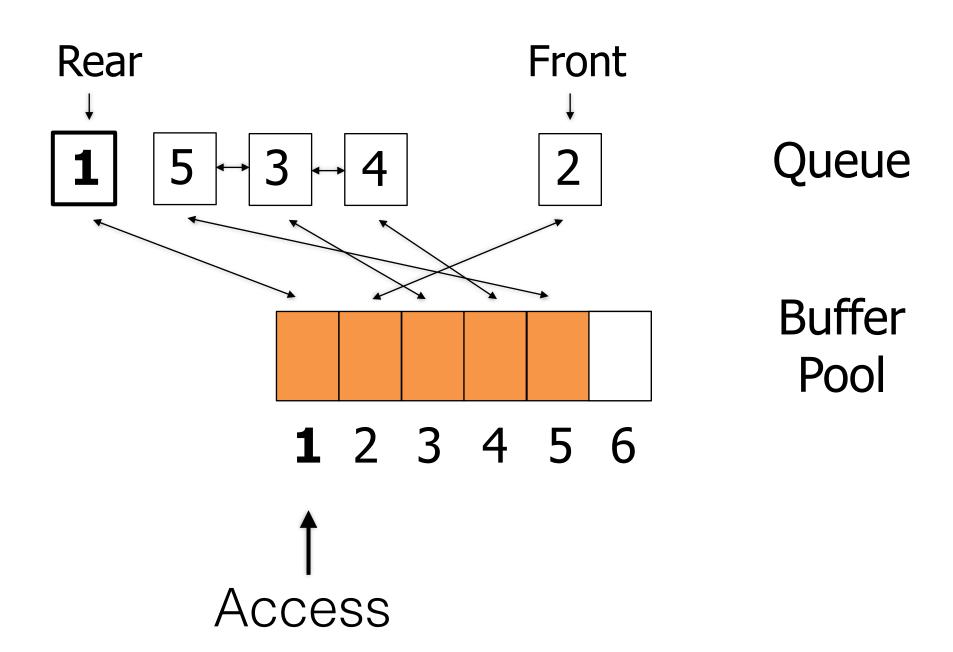
Load to rear and evict front (as before)



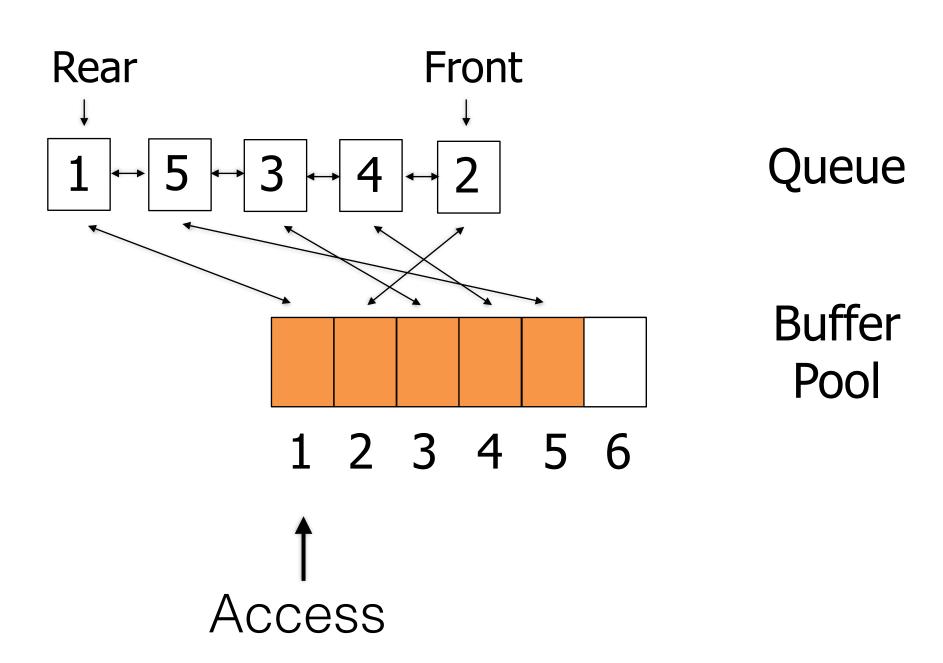
During access, return entry to rear



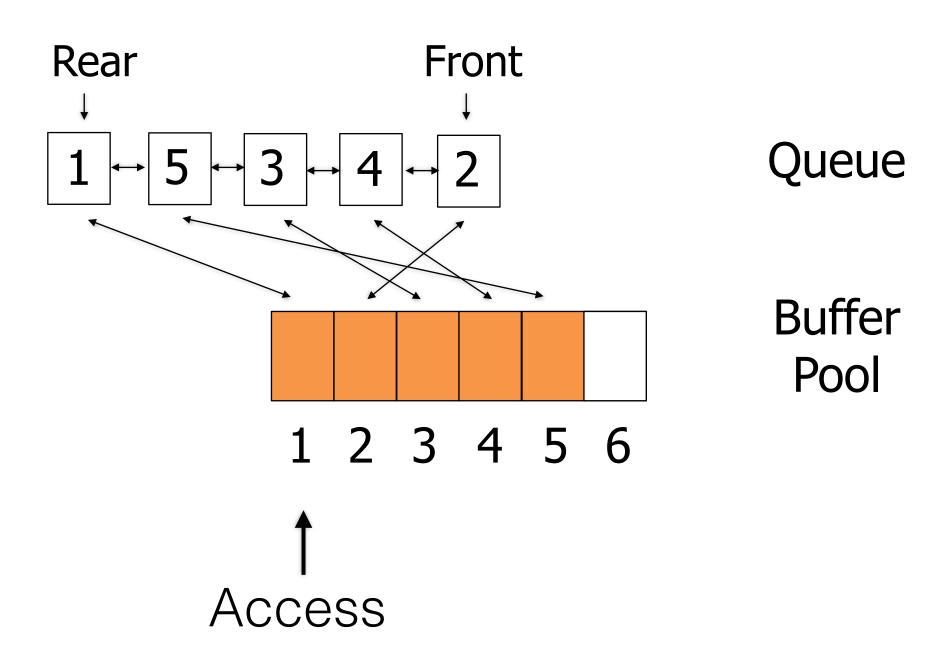
Implementation? Doubly-linked list



Implementation? Doubly-linked list

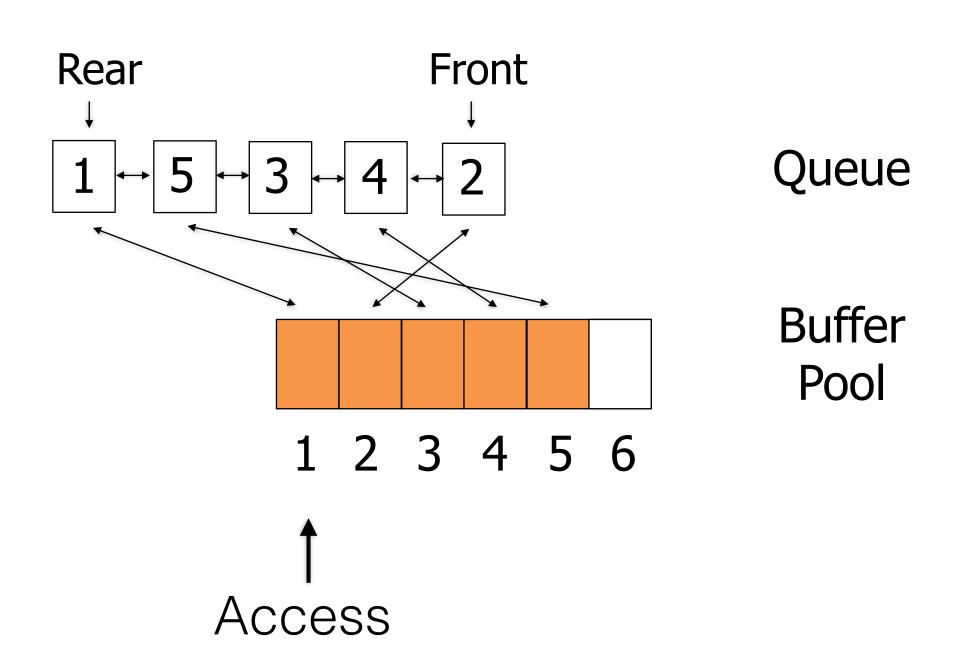


Problems?



Problems?

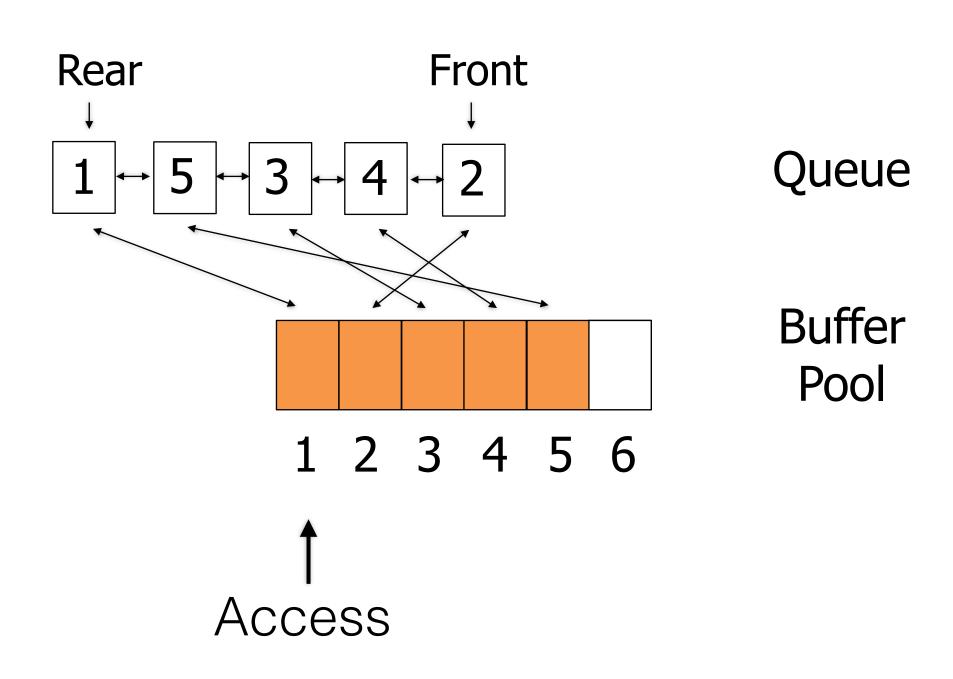
- (1) CPU overhead to update queue for each access
- (2) Linked lists are less efficient than arrays due to pointer chasing
- (3) Metadata overhead for pointers



Problems?

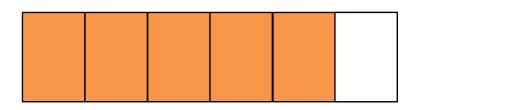
- (1) CPU overhead to update queue for each access
- (2) Linked lists are less efficient than arrays due to pointer chasing
- (3) Metadata overhead for pointers

Better ideas?:)



Traverse hash table circularly as a clock. Evict any entry not used since last traversal.

Implementation?



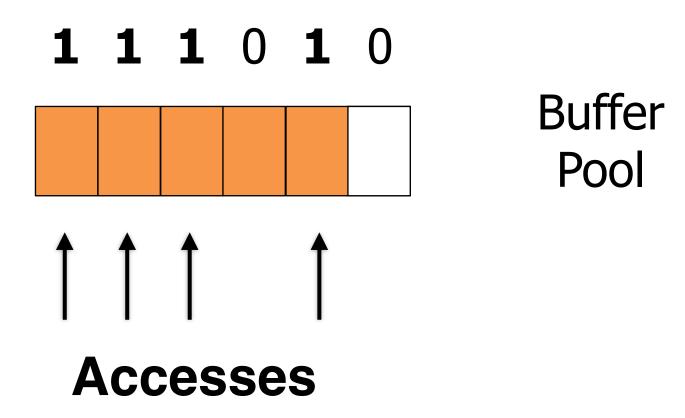
Buffer Pool

Traverse hash table circularly as a clock. Evict any entry not used since last traversal.

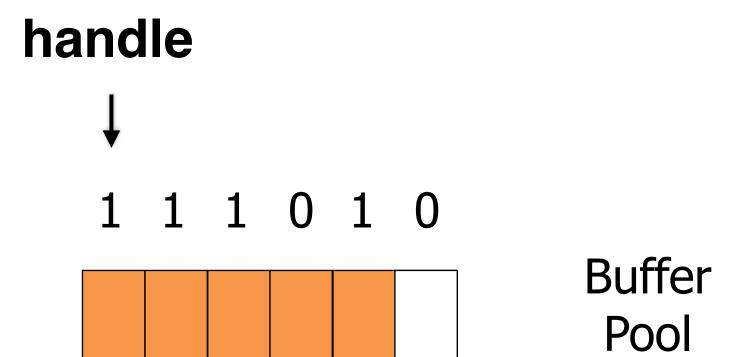
Implementation? Bitmap

0 0 0 0 0 0 Buffer Pool

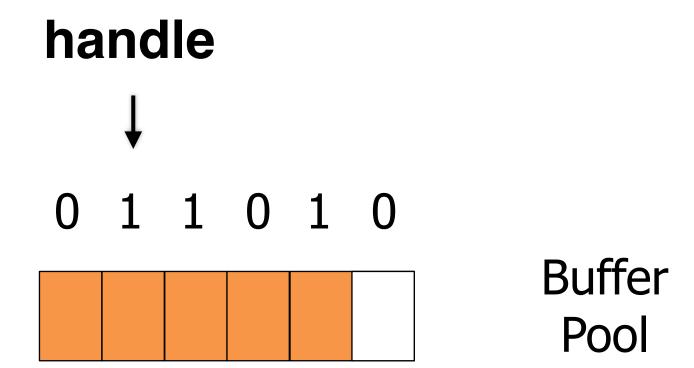
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



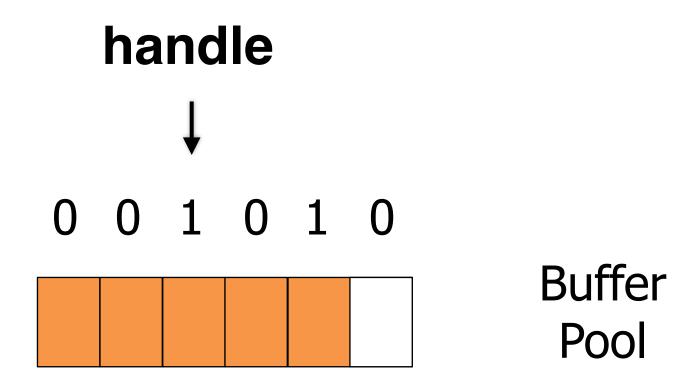
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



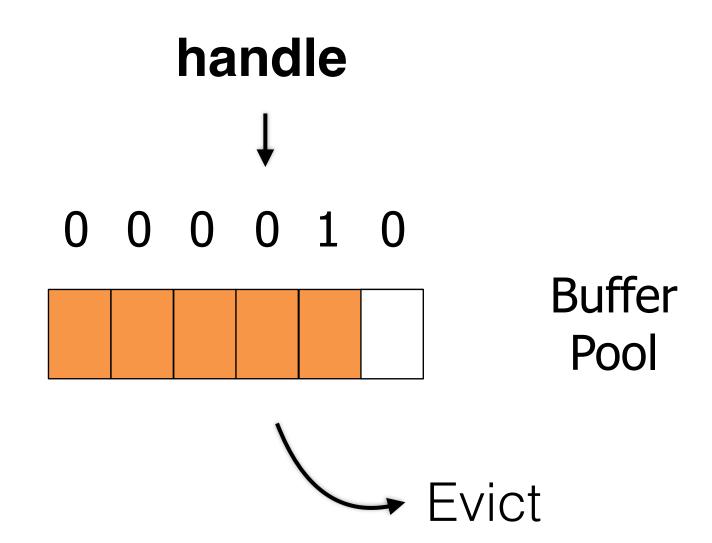
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



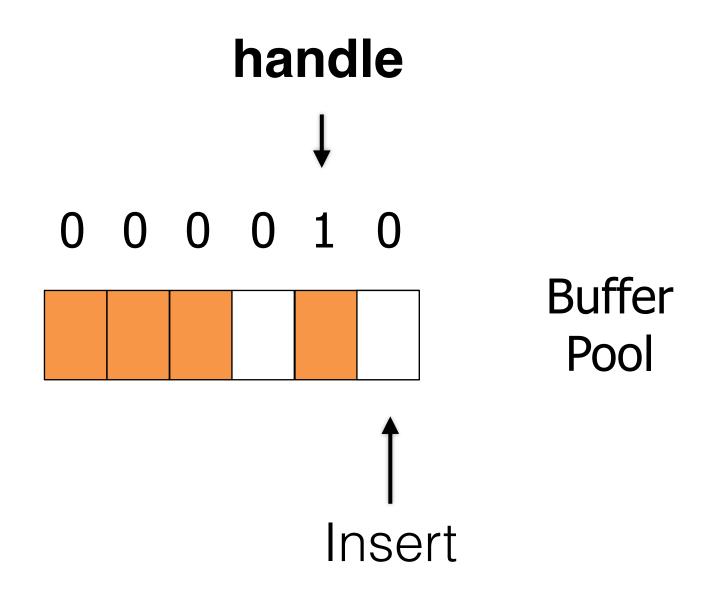
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



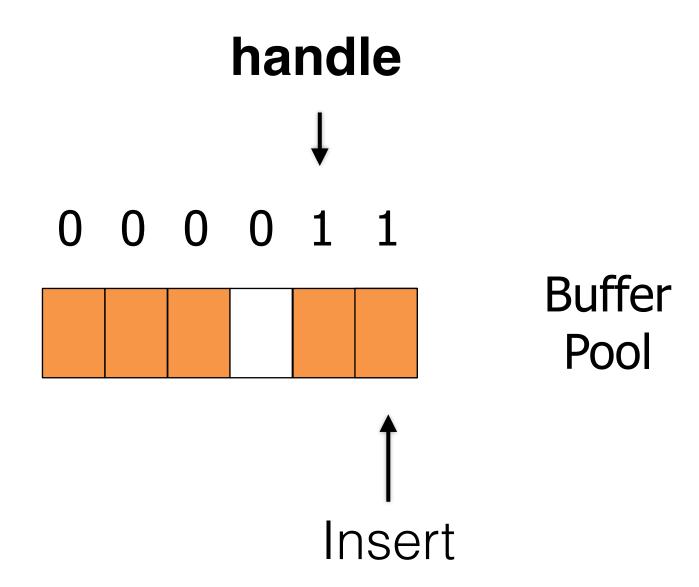
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



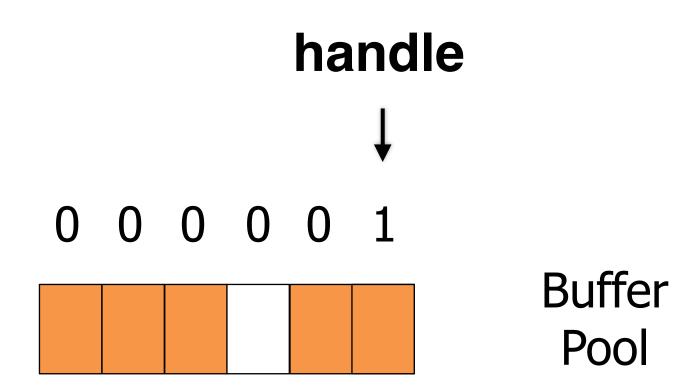
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



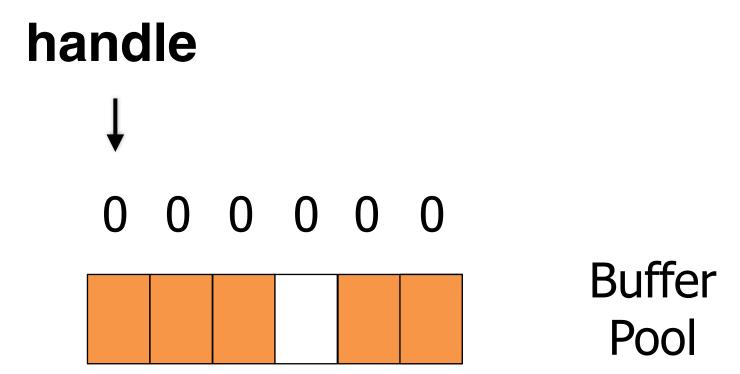
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



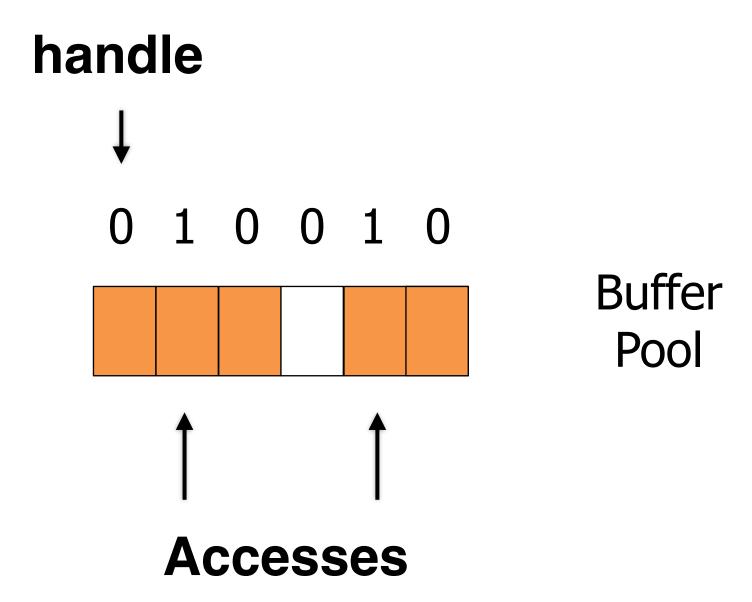
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



Traverse hash table circularly as a clock. Evict any entry not used since last traversal.



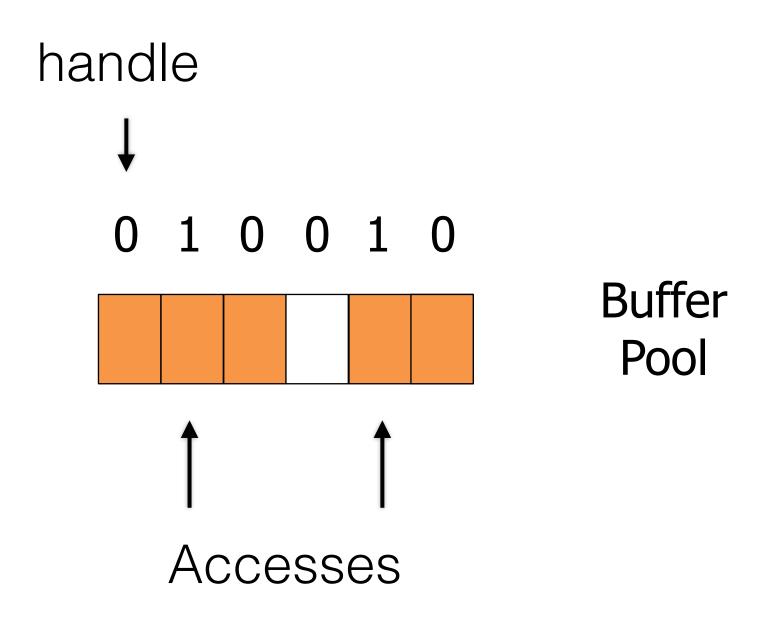
Traverse hash table circularly as a clock. Evict any entry not used since last traversal.

Advantages

- (1) lower overheads as there is no queue
- (2) bitmap takes little extra space

Disadvantages

(1) can evict "hotter" pages than LRU, But still better than FIFO



Summary

	Eviction Effectiveness	CPU
Random	Worst	Best
FIFO	Moderate	Good
LRU	Best	Worst
Clock	Close to Best	Good

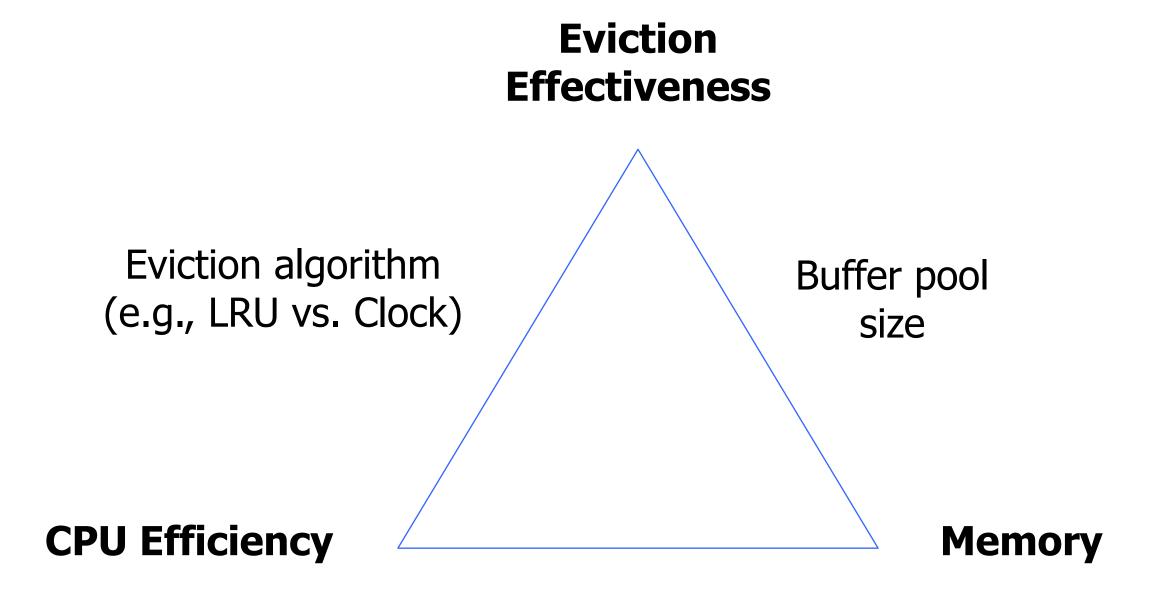
Two trade-offs

Eviction Effectiveness

Eviction algorithm (e.g., LRU vs. Clock)

CPU Efficiency

Two trade-offs



LRU and Clock are good for small random reads. How do they respond to large sequential reads?