

## CSE 444: Database Internals

### Lecture 4 Data storage and (more) buffer management

CSE 444 - Spring 2015

1

## Homework Logistics

- Homework instructions are in a pdf file
- Submit a single pdf or word file with your solution, or
- Submit a hard copy in class

CSE 444 - Spring 2015

2

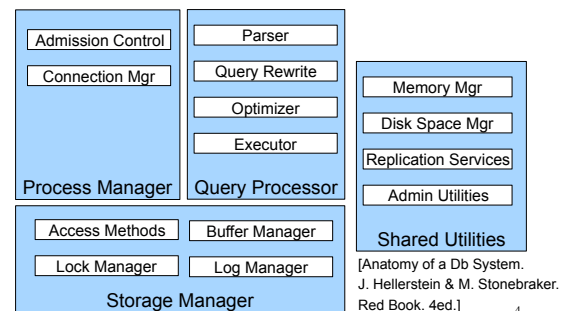
## Important Note

- Lectures show principles
- You need to think through what you will actually implement in SimpleDB!
  - Try to implement the simplest solutions
- If you are confused, tell us!

CSE 444 - Spring 2015

3

## DBMS Architecture



4

## Today: Starting at the Bottom

Consider a relation storing tweets:  
`Tweets(tid, user, time, content)`

How should we store it on disk?

CSE 444 - Spring 2015

5

## Design Exercise

- Design choice: **One OS file for each relation**
  - This does not always have to be the case! (e.g., SQLite uses one file for whole database)
  - DBMSs can also use disk drives directly
- An OS file provides an API of the form
  - Seek to some position (or "skip" over B bytes)
  - Read/Write B bytes

CSE 444 - Spring 2015

6

## First Principle: Work with Pages

- Reading/writing to/from disk
  - Seeking takes a long time!
  - Reading sequentially is fast
- To simplify buffer manager, want to cache a collection of same-sized objects
- Solution: Read/write **pages** of data
  - A page should correspond to a disk block

CSE 444 - Spring 2015

7

## Continuing our Design

Key questions:

- How do we organize pages into a file?
- How do we organize data within a page?

First, **how could we store some tuples on a page?**

Let's first assume all tuples are of the same size

```
Tweets(tid int, user char(10),
       time int, content char(140))
```

CSE 444 - Spring 2015

8

## Page Formats

Issues to consider

- 1 page = 1 disk block = fixed size (e.g. 8KB)
- Records:
  - Fixed length
  - Variable length
- **Record id = RID**
  - Typically RID = (PageID, SlotNumber)

**Why do we need RID's in a relational DBMS ?**  
See future discussion on indexes and transactions

CSE 444 - Spring 2015

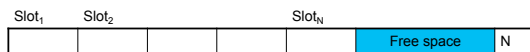
9

## Page Format Approach 1

Fixed-length records: packed representation

Divide page into slots. Each slot can hold one tuple

Record ID (RID) for each tuple is (PageID, SlotNb)



How do we insert a new record?

Number of records

CSE 444 - Spring 2015

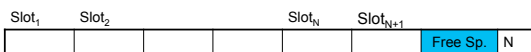
10

## Page Format Approach 1

Fixed-length records: packed representation

Divide page into slots. Each slot can hold one tuple

Record ID (RID) for each tuple is (PageID, SlotNb)



How do we insert a new record?

Number of records

CSE 444 - Spring 2015

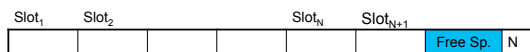
11

## Page Format Approach 1

Fixed-length records: packed representation

Divide page into slots. Each slot can hold one tuple

Record ID (RID) for each tuple is (PageID, SlotNb)



How do we insert a new record?

Number of records

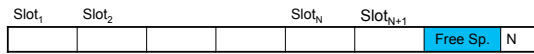
How do we delete a record?

CSE 444 - Spring 2015

12

## Page Format Approach 1

Fixed-length records: packed representation  
Divide page into slots. Each slot can hold one tuple  
Record ID (RID) for each tuple is (PageID, SlotNb)



How do we insert a new record? Number of records

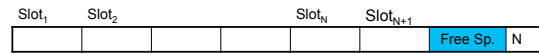
How do we delete a record? What is the problem?

CSE 444 - Spring 2015

13

## Page Format Approach 1

Fixed-length records: packed representation  
Divide page into slots. Each slot can hold one tuple  
Record ID (RID) for each tuple is (PageID, SlotNb)



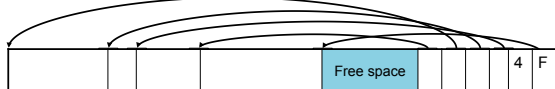
How do we insert a new record? Number of records

How do we delete a record? Cannot move records! (Why?)

How do we handle variable-length records?

14

## Page Format Approach 2



Header contains slot directory

+ Need to keep track of nb of slots

+ Also need to keep track of free space (F)

Slot directory  
Each slot contains  
<record offset, record length>

Can handle variable-length records

Can move tuples inside a page without changing RIDs

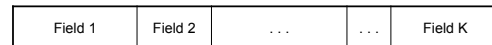
RID is (PageID, SlotID) combination

CSE 444 - Spring 2015

15

## Record Formats

Fixed-length records → Each field has a fixed length  
(i.e., it has the same length in all the records)



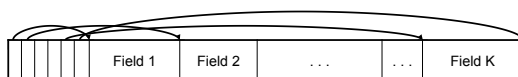
Information about field lengths and types is in the catalog

CSE 444 - Spring 2015

16

## Record Formats

Variable length records



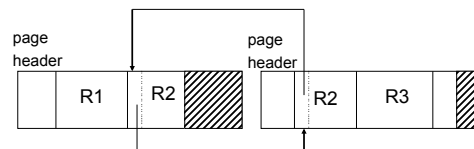
Record header

Remark: NULLS require no space at all (why ?)

CSE 444 - Spring 2015

17

## Long Records Across Pages



- When records are very large
- Or even medium size: saves space in blocks
- Commercial RDBMSs avoid this

CSE 444 - Spring 2015

18

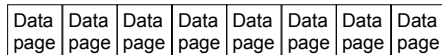
- Large objects
  - Binary large object: BLOB
  - Character large object: CLOB
- Supported by modern database systems
- E.g. images, sounds, texts, etc.
- Storage: attempt to cluster blocks together

Our key questions:

- How do we organize pages into a file?
- How do we organize data within a page?

Now, how should we group pages into files?

A sequence of pages (implementation in SimpleDB)

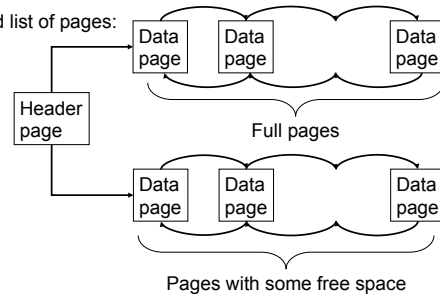


Some pages have space and other pages are full  
Add pages at the end when need more space

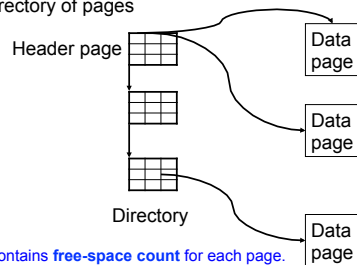
Works well for small files

But finding free space requires scanning the file

Linked list of pages:



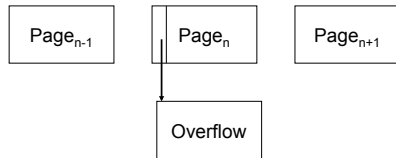
Better: directory of pages



Directory contains **free-space count** for each page.  
Faster inserts for variable-length records

- File is unsorted (= **heap file**)
  - add it wherever there is space (easy ☺)
  - add more pages if out of space
- File is sorted
  - Is there space on the right page ?
    - Yes: we are lucky, store it there
  - Is there space in a neighboring page ?
    - Look 1-2 pages to the left/right, shift records
  - If anything else fails, create **overflow page**

## Overflow Pages



- After a while the file starts being dominated by overflow pages: time to reorganize

CSE 444 - Spring 2015

25

## Modifications: Deletions

- Free space in page, shift records
  - Be careful with slots
  - RIDs for remaining tuples must NOT change
- May be able to eliminate an overflow page

CSE 444 - Spring 2015

26

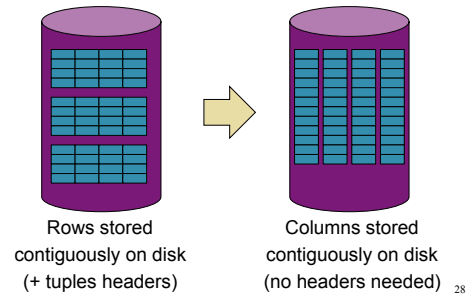
## Modifications: Updates

- If new record is shorter than previous, easy 😊
- If it is longer, need to shift records
  - May have to create overflow pages

CSE 444 - Spring 2015

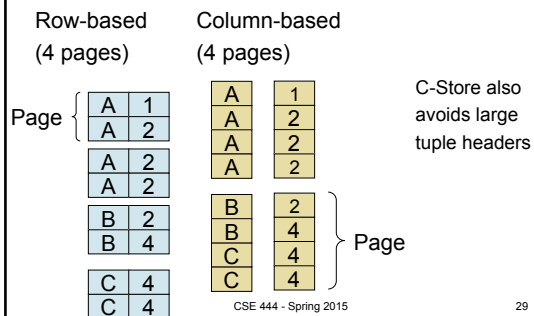
27

## Alternate Storage Manager Design: Column Store



28

## More Detailed Example



CSE 444 - Spring 2015

29

## Continuing our Design

We know how to store tuples on disk in a heap file

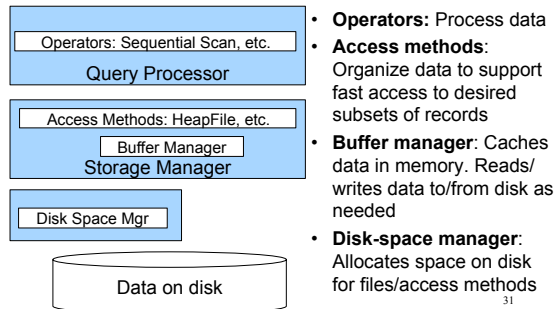
How do these files interact with rest of engine?

- Also see lecture 3

CSE 444 - Spring 2015

30

## How Components Fit Together



31

## Access Methods

- Operators view relations as collections of records
- The access methods worry about how to organize these collections

CSE 444 - Spring 2015

32

## Heap File Access Method API

- **Create** or **destroy** a file
- **Insert** a record
- **Delete** a record with a given rid (rid)
  - rid: unique tuple identifier (more later)
- **Get** a record with a given rid
  - Not necessary for sequential scan operator
  - But used with indexes (more next lecture)
- **Scan** all records in the file

CSE 444 - Spring 2015

33

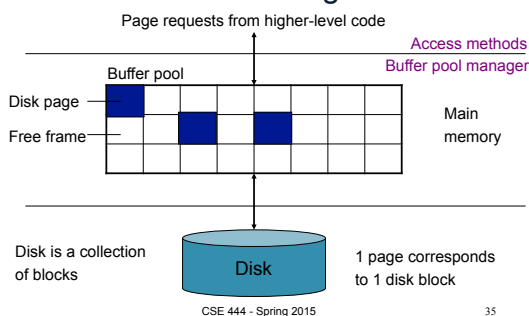
## Buffer Manager

- Brings pages in from memory and caches them
- Eviction policies
  - Random page (ok for SimpleDB)
  - Least-recently used
  - The “clock” algorithm (see whiteboard or book)
- Keeps track of which **pages are dirty**
  - A dirty page has changes not reflected on disk
  - Implementation: Each page includes a dirty bit

CSE 444 - Spring 2015

34

## Buffer Manager



CSE 444 - Spring 2015

35

## Pushing Updates to Disk

- When inserting a tuple, HeapFile inserts it on a page but does not write the page to disk
- When deleting a tuple, HeapFile deletes tuple from a page but does not write the page to disk
- The buffer manager worries when to write pages to disk (and when to read them from disk)
- When need to add a new page to the file, HeapFile adds page to the file on disk and then gets it again through the buffer manager

CSE 444 - Spring 2015

36

## Conclusion

- Row-store storage managers are most commonly used today
- They offer high-performance for transactions
- But column-stores win for analytical workloads
- They are gaining traction in that area
- Final discussion: OS vs DBMS
  - OS files vs DBMS files
  - OS buffer manager vs DBMS buffer manager

CSE 444 - Spring 2015

37