



A note on (confusing) terminology

Block = Unit of transfer for disk read/write

- 64KB - 128KB is a good number today
- Book says 4KB

Page = Fixed size contiguous chunk of memory

- Assume same size as block
- Refer to corresponding blocks on disk

For simplicity we use Block and Page interchangeably.



Disk Space Management

Lowest layer of DBMS, manages space on disk

- Mapping pages to locations on disk
- Loading pages from disk to memory
- Saving pages back to disk & ensuring writes

Higher levels call upon this layer to:

- read/write a pages
- allocate/de-allocate logical pages

Request for a sequence of pages best satisfied by pages stored sequentially on disk

- Physical details hidden from higher levels of system
- Higher levels may assume **Next Page** is fast!



Disk Space Management Implementation

Proposal 1: Talk to the device directly

- Could be very fast if you knew the device well
- What happens when devices change?

Proposal 2: Run over filesystem (FS)

- Allocate single large "contiguous" file and assume sequential / nearby byte access are fast
- Most FS optimize for sequential access and temporal locality (buffer cache on hot items)
 - Sometimes disable FS buffering
- May span multiple files on multiple disks / machines



Typically sits on top of local file system

Get Page 4



Get Page 5



Disk Space Management

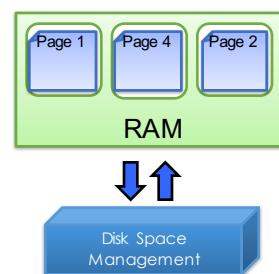


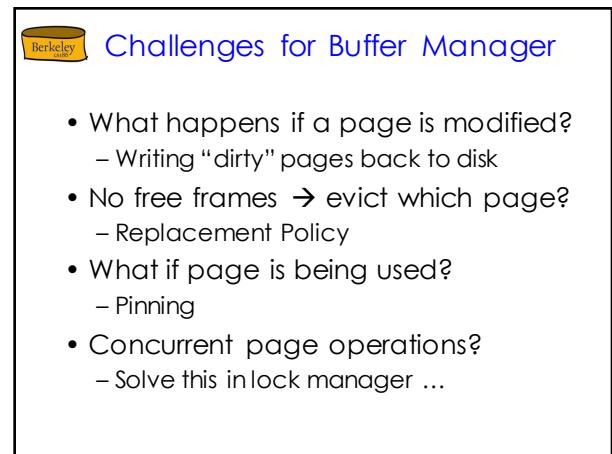
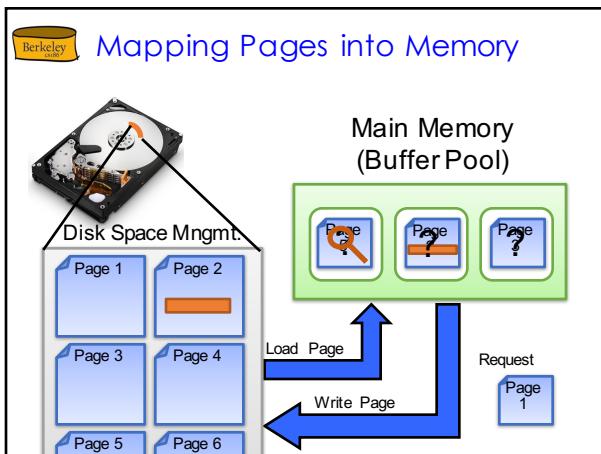
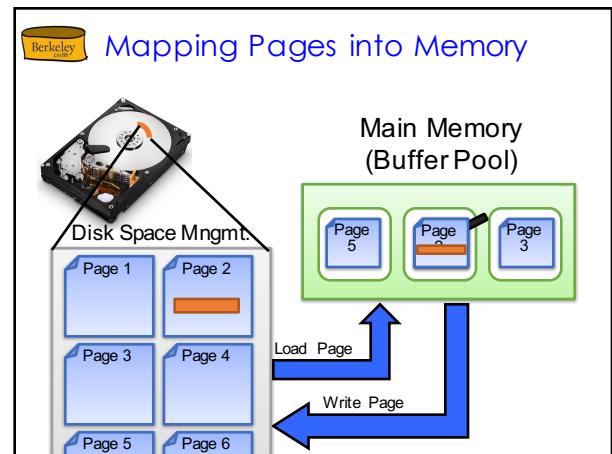
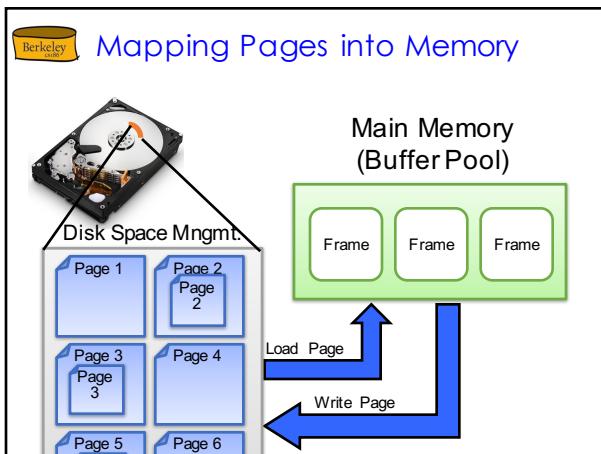
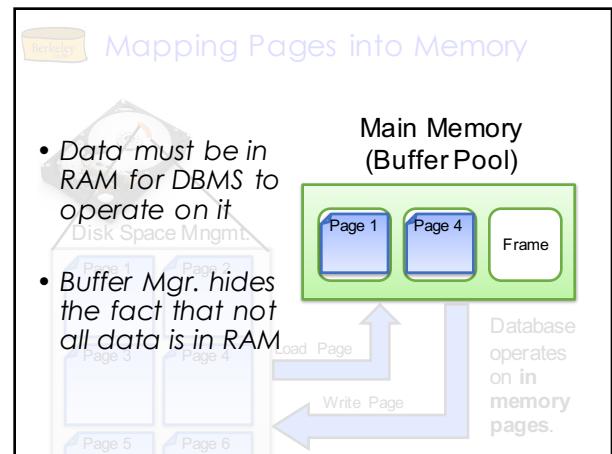
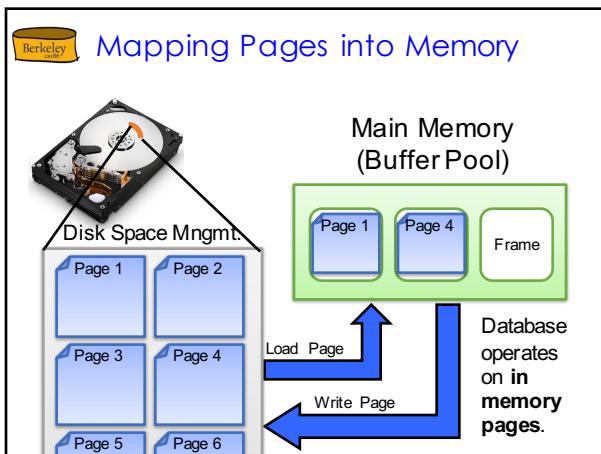
- Provide API to read and write pages to device
- Pages: block level organization of bytes on disk
- Ensures next locality and abstracts FS/Device details



Architecture of a DBMS

Illusion of operating in memory



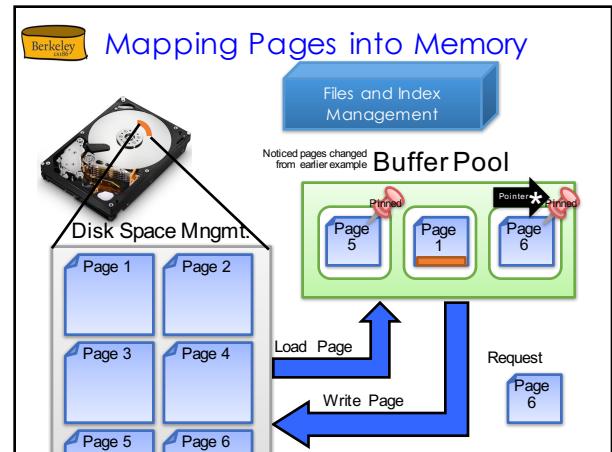


 When a Page is Requested ...

- Buffer pool information "table" contains: $\langle \text{frame\#}, \text{pageid}, \text{pin_count}, \text{dirty} \rangle$

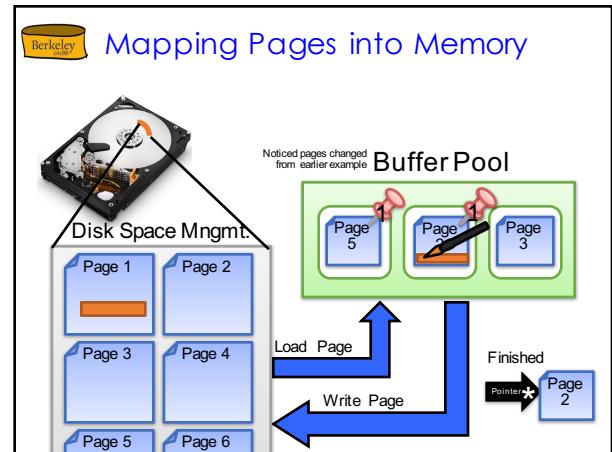
1. If requested page is not in pool:
 - a. Choose a frame for *replacement*. **Only "un-pinned" pages are candidates!**
 - b. If frame "dirty", write current page to disk
 - c. Read requested page into frame
2. Pin the page and return its address.

If requests can be predicted (e.g., sequential scans) pages can be pre-fetched several pages at a time!



 After Requestor Finishes

- Requestor of page must:
 1. indicate whether page was modified via *dirty* bit.
 2. unpin it (soon preferably!) why?
- Page in pool may be requested many times,
 - a *pin count* is used.
 - To pin a page: *pin_count++*
 - A page is a candidate for replacement iff *pin count == 0* ("unpinned")
- CC & recovery may do additional I/Os upon replacement.
 - Write-Ahead Log protocol; more later!

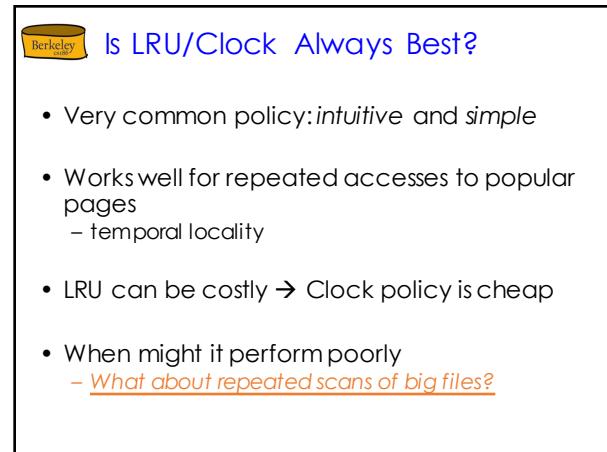
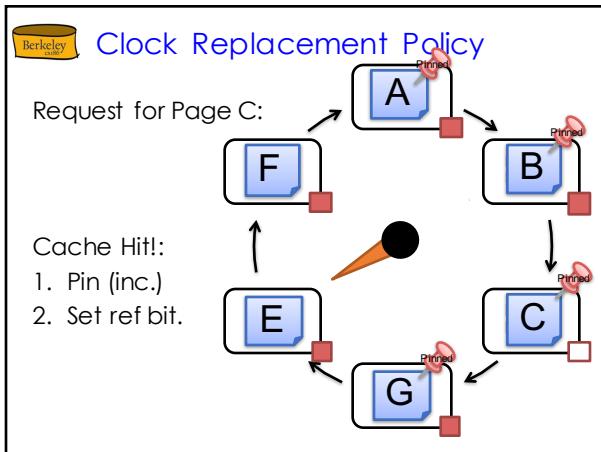
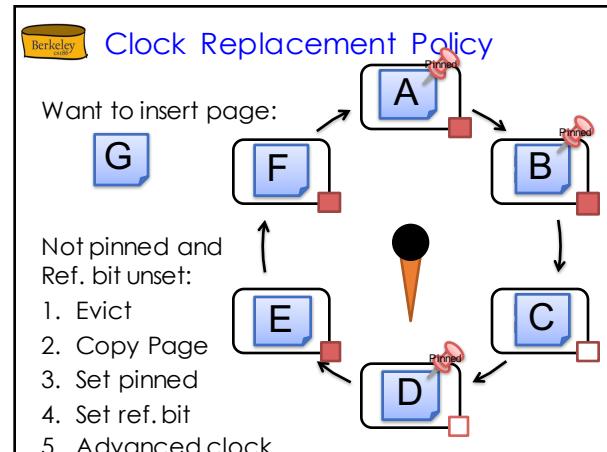
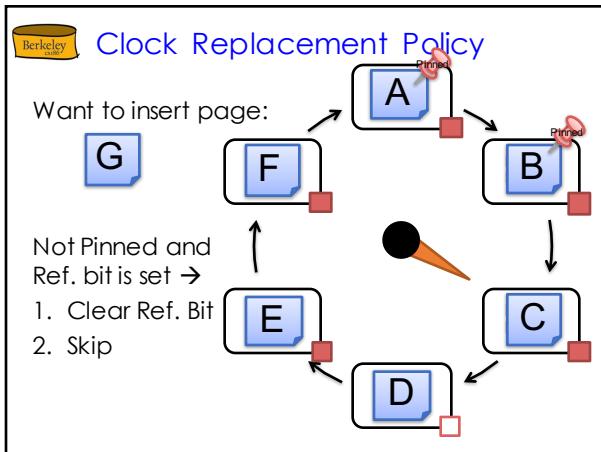
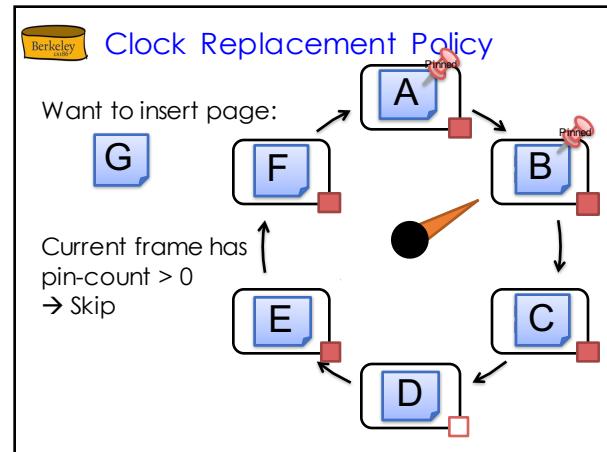
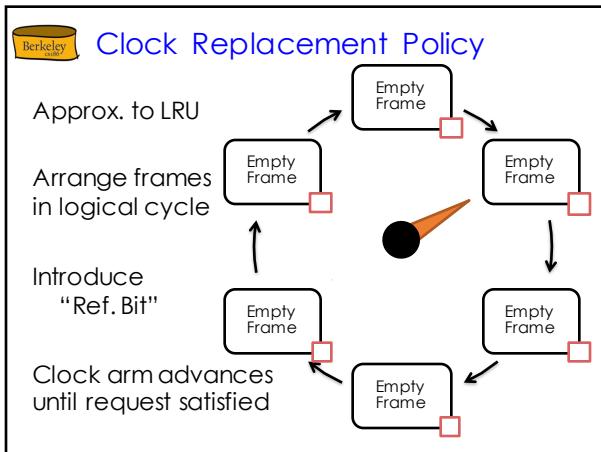


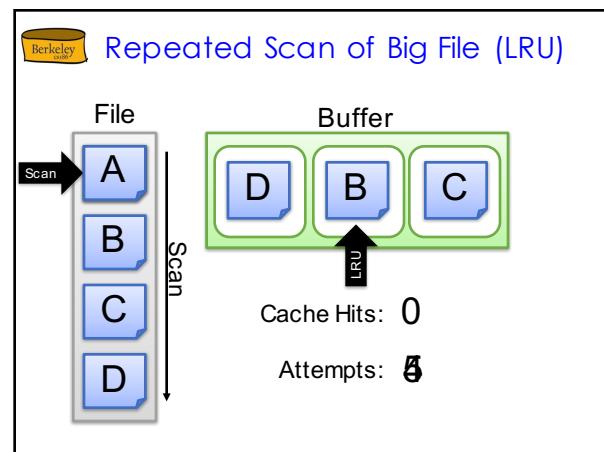
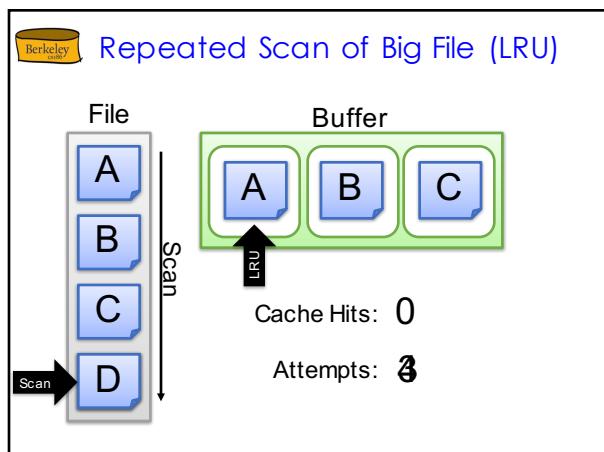
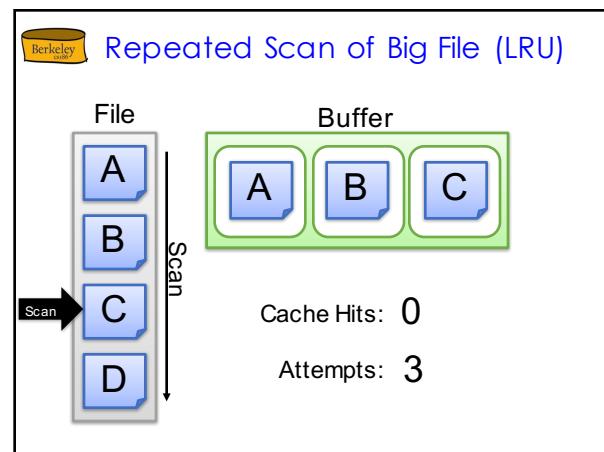
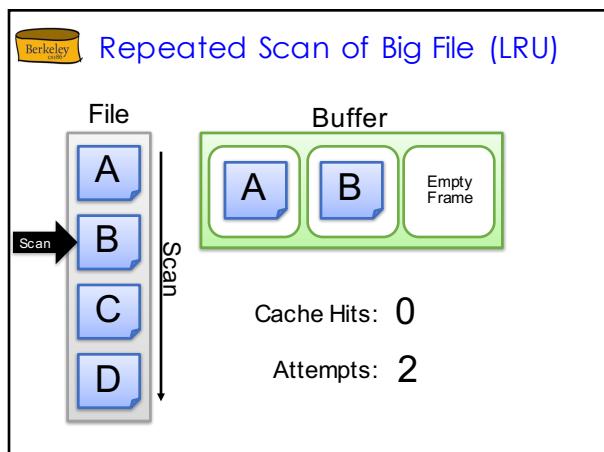
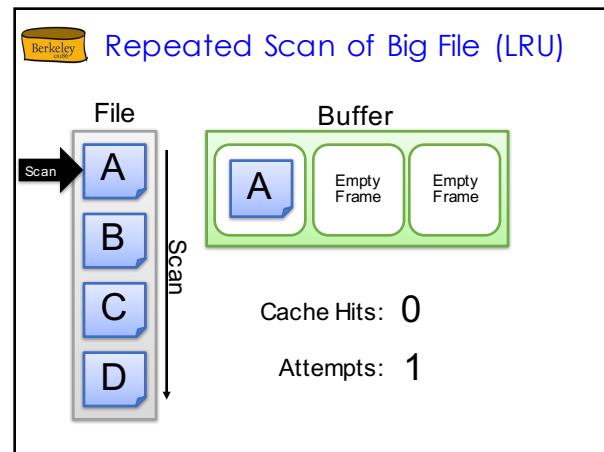
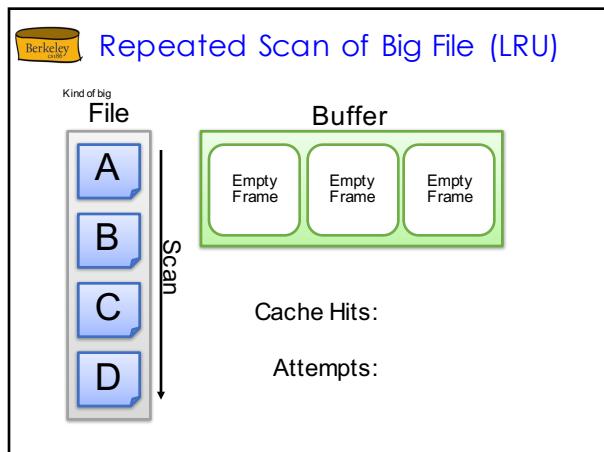
 Page Replacement Policy

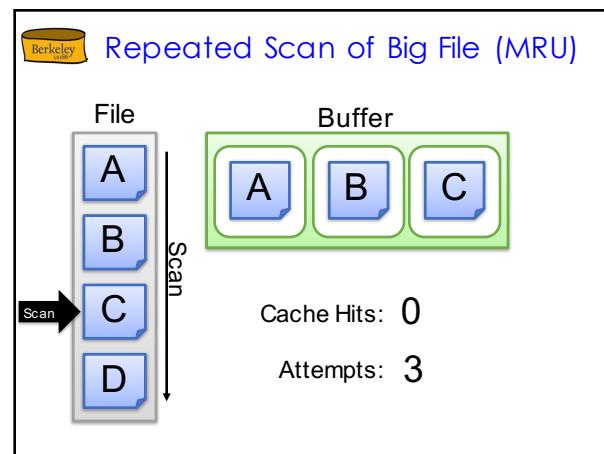
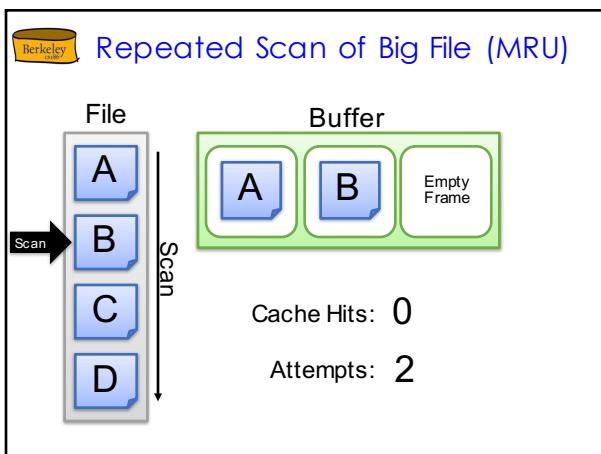
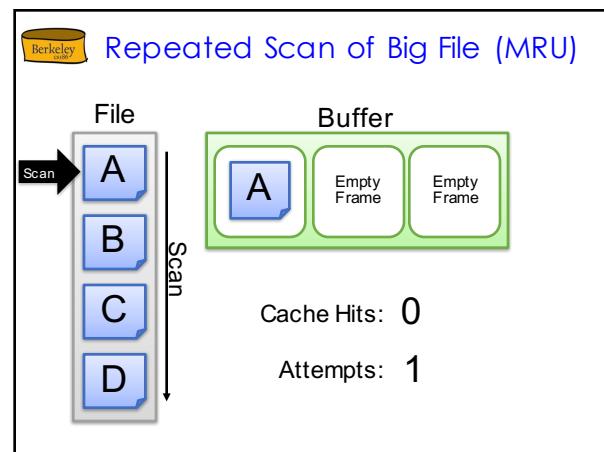
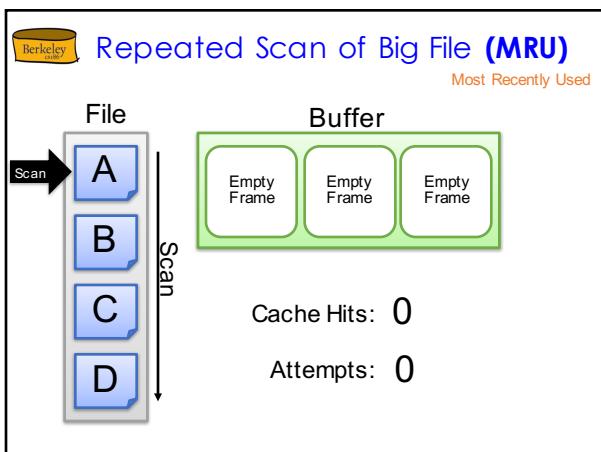
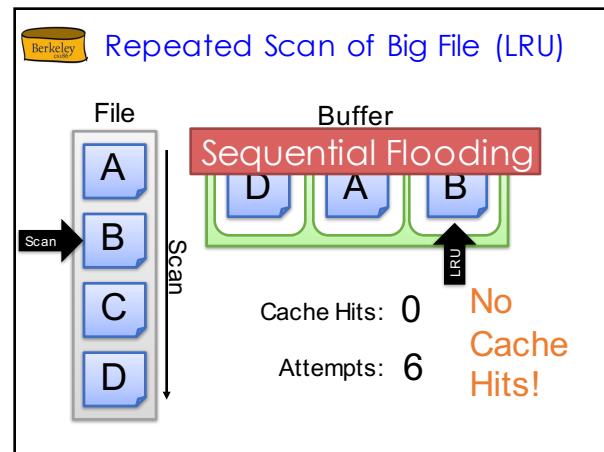
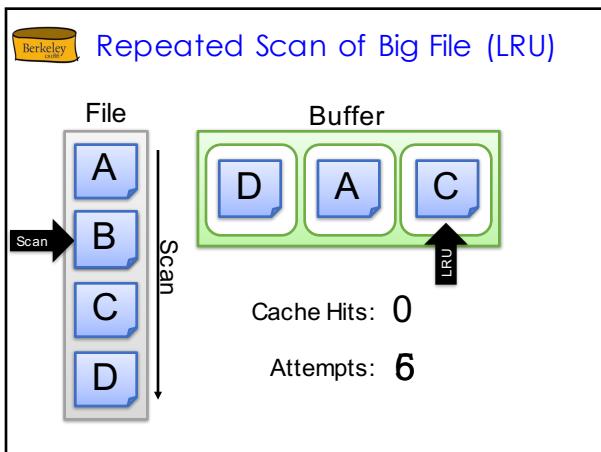
- Page is chosen for replacement by a *replacement policy*:
 - Least-recently-used (LRU), Clock
 - Most-recently-used (MRU)
- Policy can have big impact on #I/O's;
 - Depends on the *access pattern*.

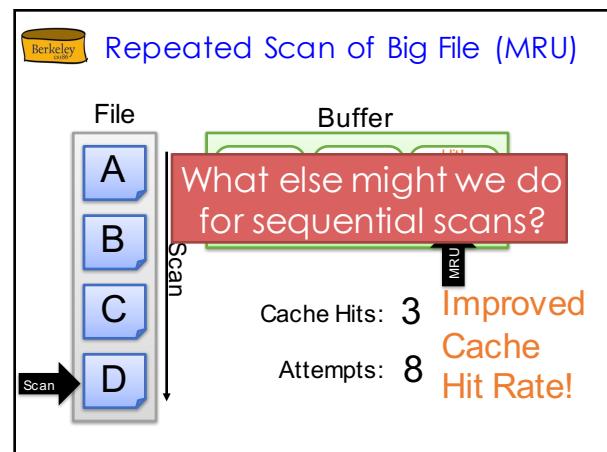
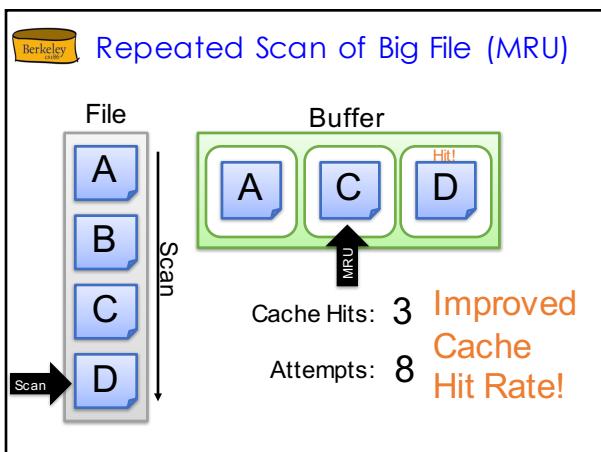
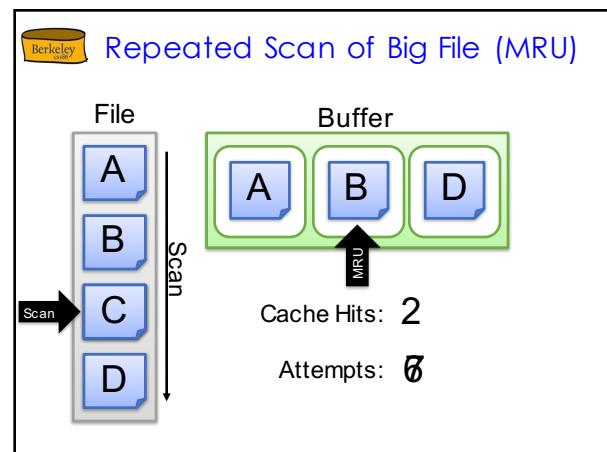
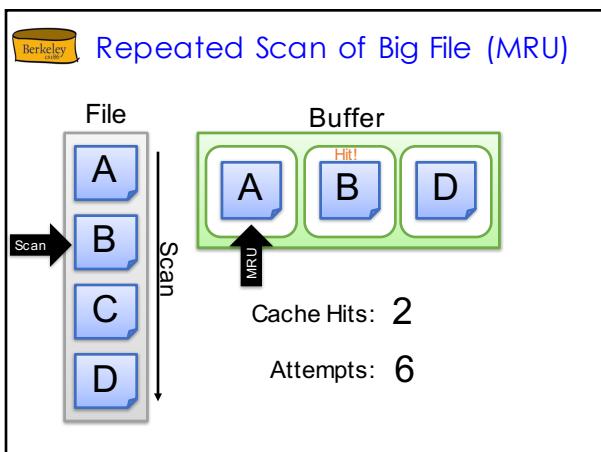
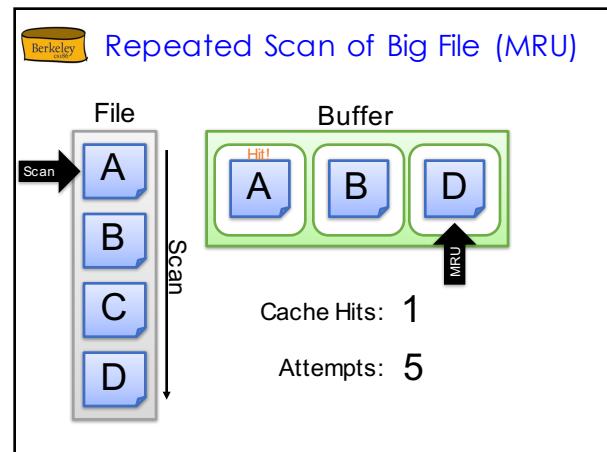
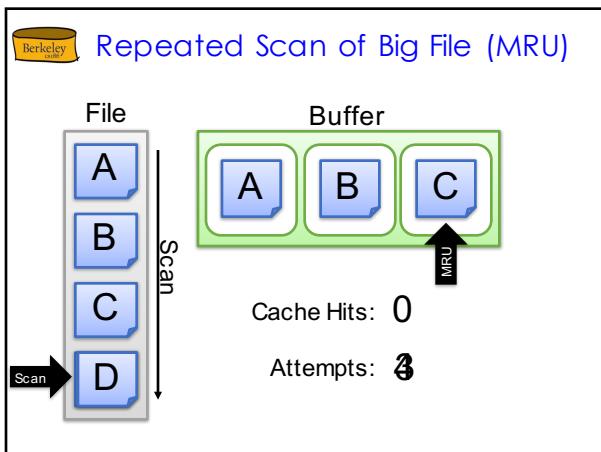
 LRU Replacement Policy

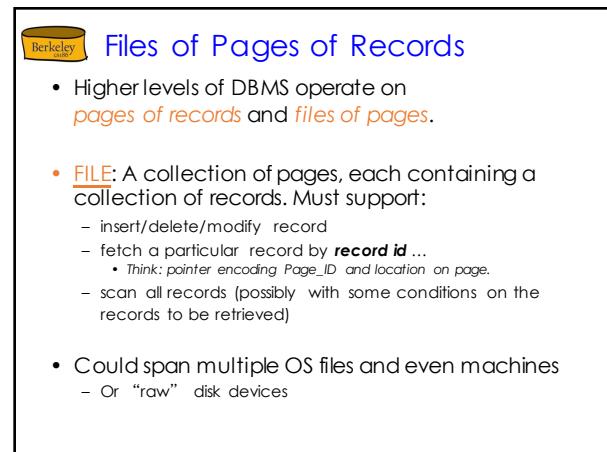
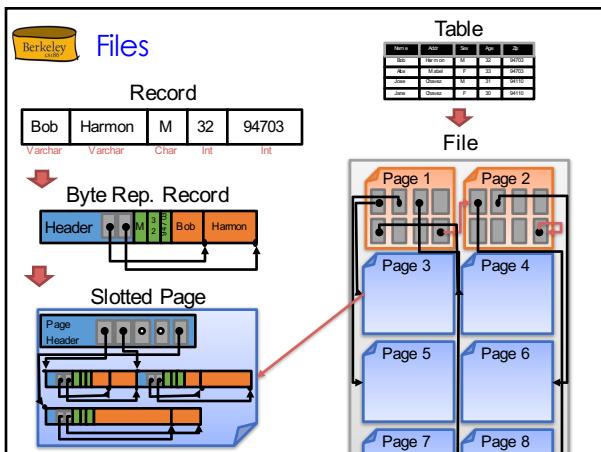
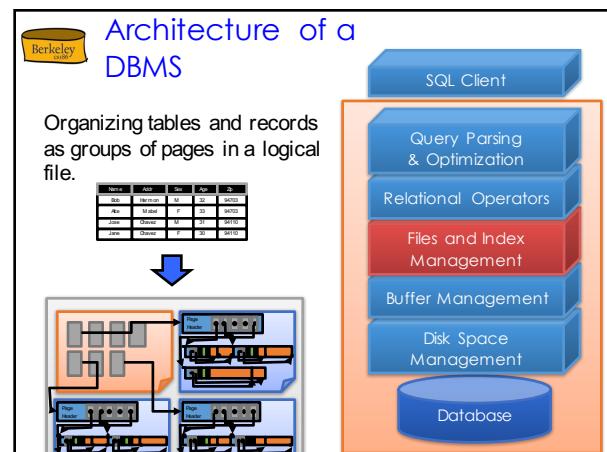
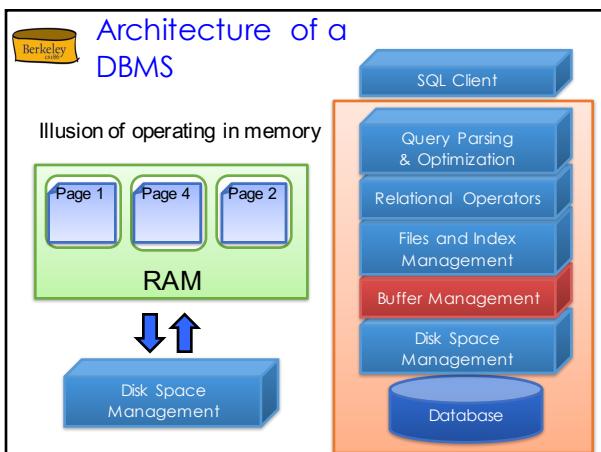
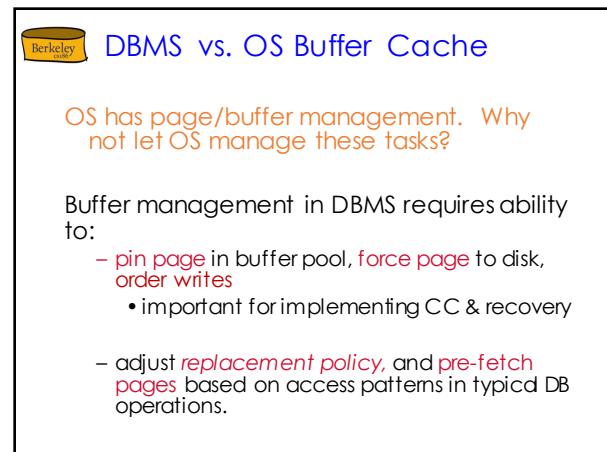
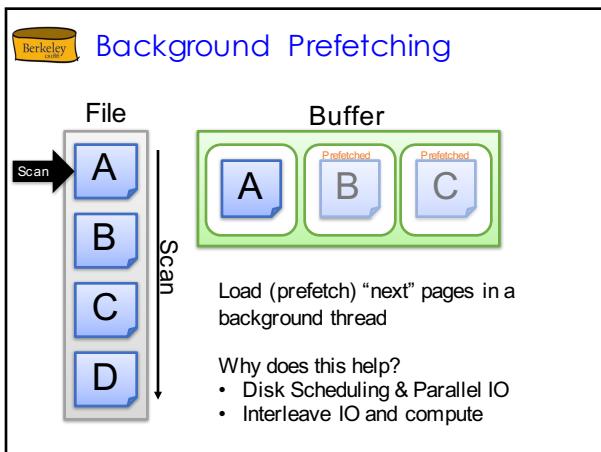
- Least Recently Used (LRU)
 - Pinned Frame: not available to replace
 - track time each frame last *unpinned* (end of use)
 - replace the frame which least recently unpinned
- Very common policy: intuitive and simple
 - Works well for repeated accesses to popular pages (temporal locality)
 - Can be costly. Why?
 - Need to maintain heap data-structure
 - Solution?











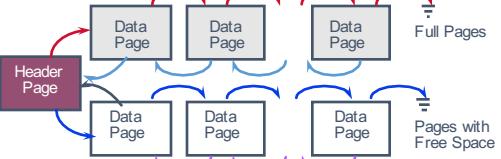
 **Many Kinds of Database Files**

- Unordered Heap Files
 - Records placed arbitrarily across pages
- Clustered Heap File and Hash Files
 - Records and pages are grouped
- Sorted Files
 - Page and records are in sorted order
- Index Files
 - B+ Trees, Hash Tables, ...
 - May contain records or point to records in other files

 **Basic Unordered Heap Files**

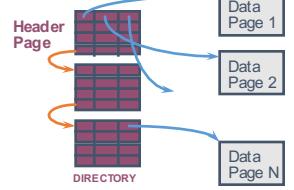
- Collection of records in no particular order
 - Not to be confused with "heap data-structure"
- As file shrinks/grows, pages (de)allocated
- To support record level operations, we must:
 - keep track of the **pages** in a file
 - keep track of **free space** on pages
 - keep track of the **records** on a page
- There are many alternatives for keeping track of this, we'll consider two in this class

 **Heap File Implemented as a List**



- Header page ID and Heap file name stored elsewhere
 - Database "catalog"
- Each page contains 2 "pointers" plus **free-space** and **data**.
- What is wrong with this?
 - How do I find a page with enough space for a 20 byte record?

 **Better: Use a Page Directory**



- Directory entries include **#free bytes** on the page.
- Header pages accessed often → likely in cache
 - What eviction policy is best here?
- Finding a page to fit a record required far fewer page loads than linked list (Why?)
 - One header page load reveals free space of many pages.

 **Indexes (sneak preview)**

- A Heap file allows us to retrieve records:
 - by specifying the **record id** (page id + slot)
 - by scanning all records sequentially
- Would like to fetch records by value, e.g.,
 - Find all students in the "CS" department
 - Find all students with a gpa > 3 AND blue hair
- **Indexes:** file structures for efficient **value-based queries**

 **Overview**

Table

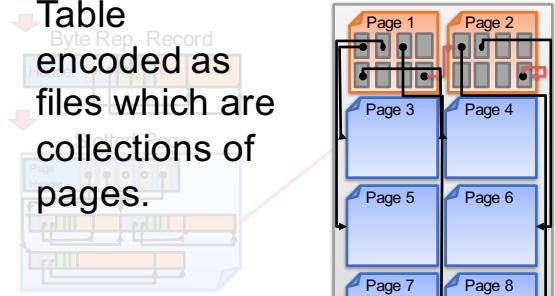
ID	Name	Age	Major	GPA
Alice	Melanie	21	30	3.5
Bob	Steve	22	31	3.4
Cathy	David	23	32	3.6
Dave	Charles	24	33	3.7

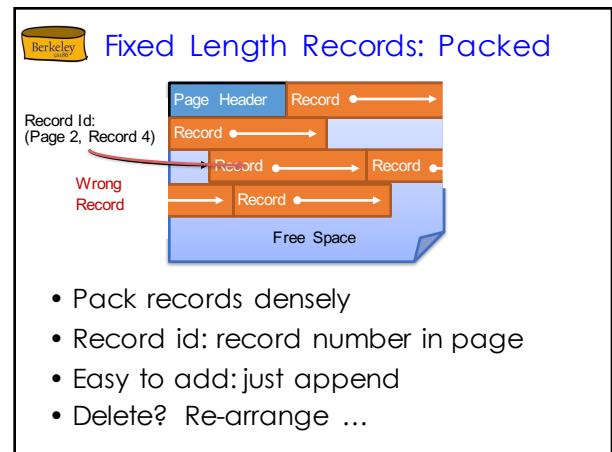
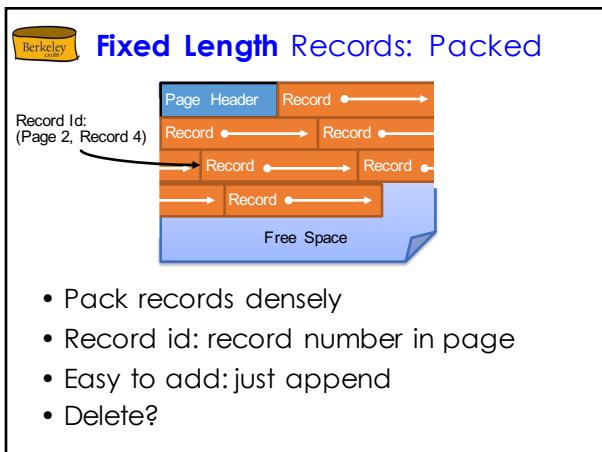
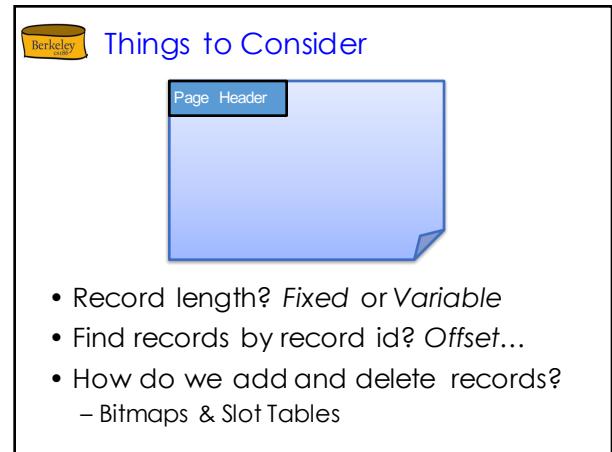
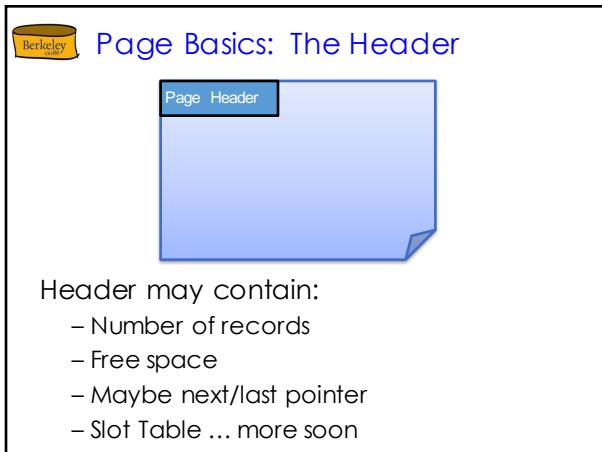
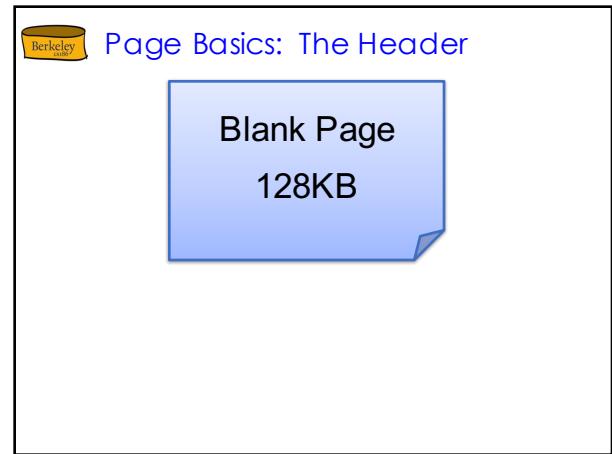
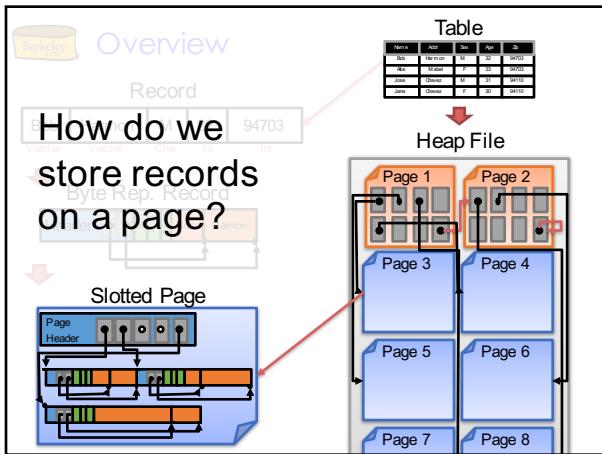
Record

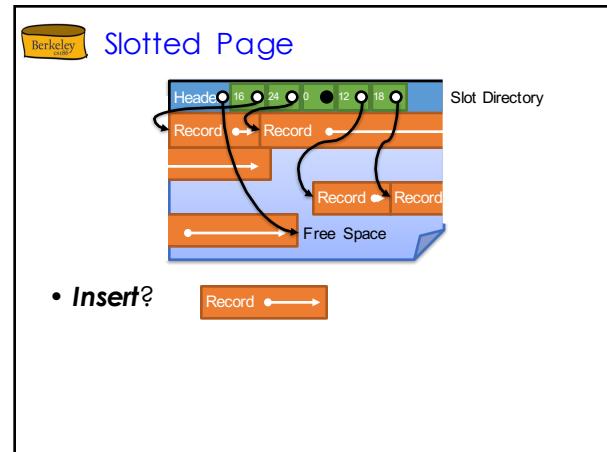
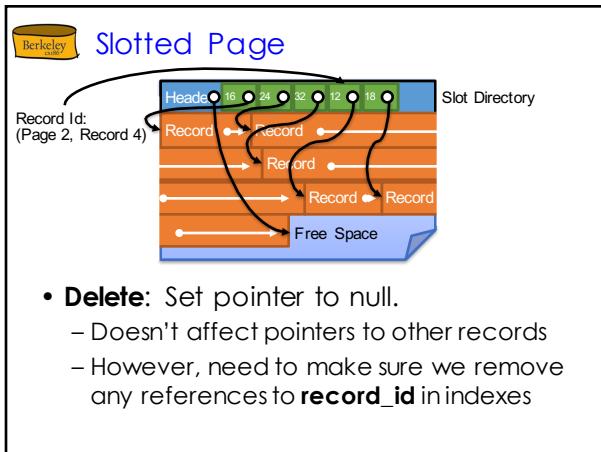
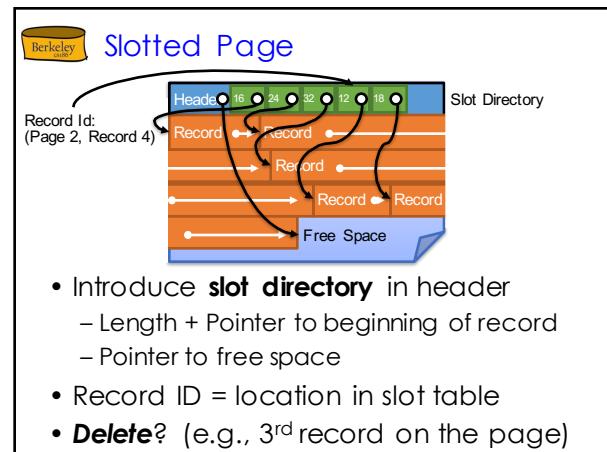
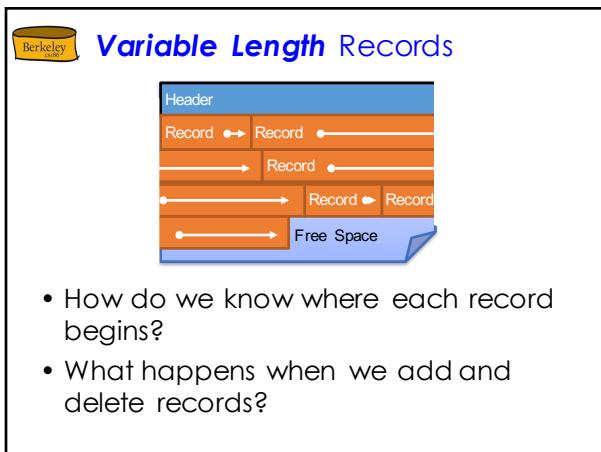
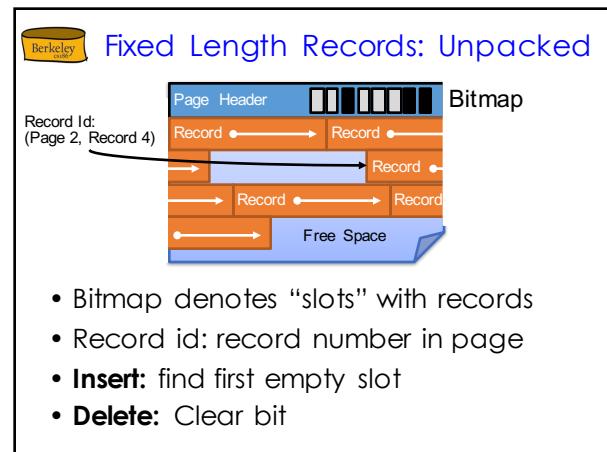
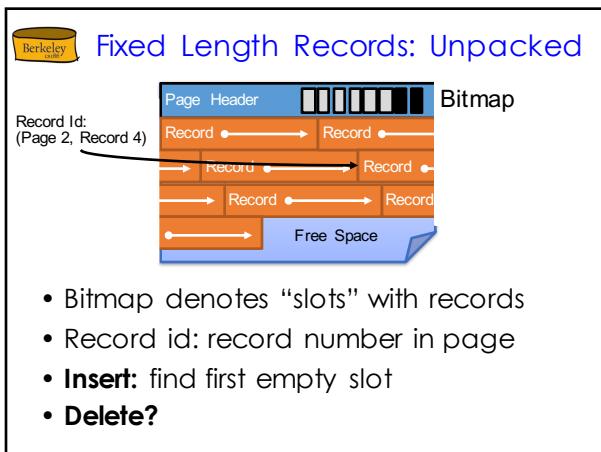
Summary M 32 94703

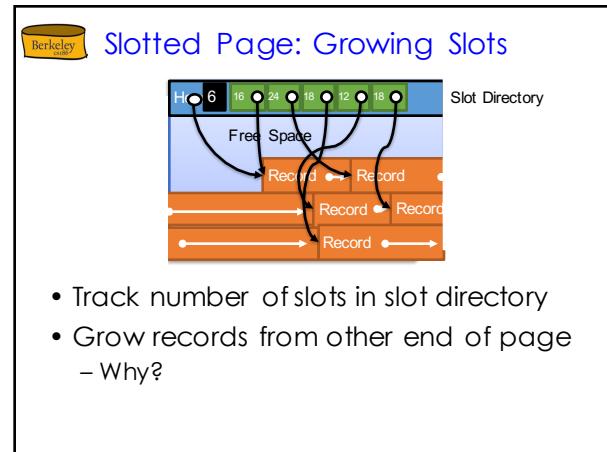
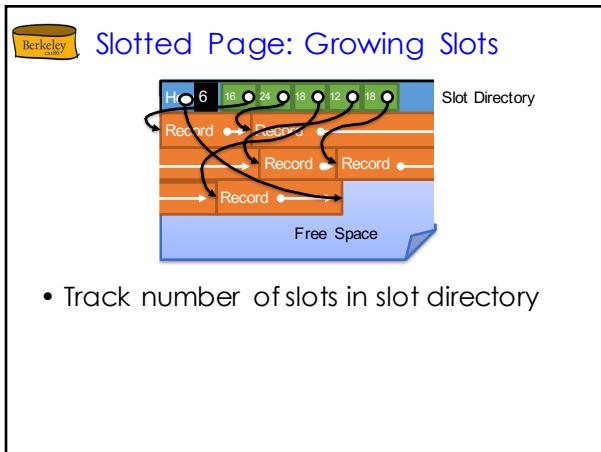
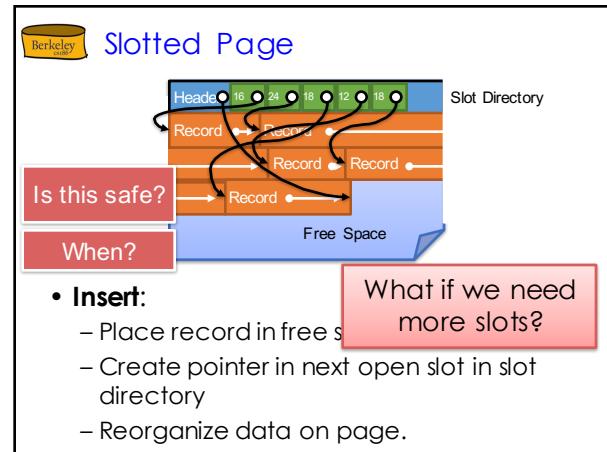
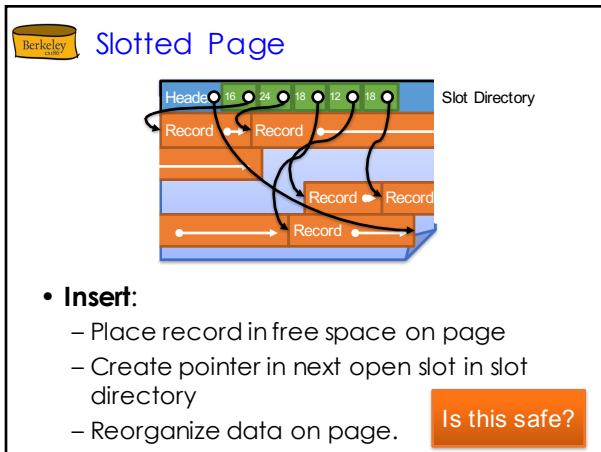
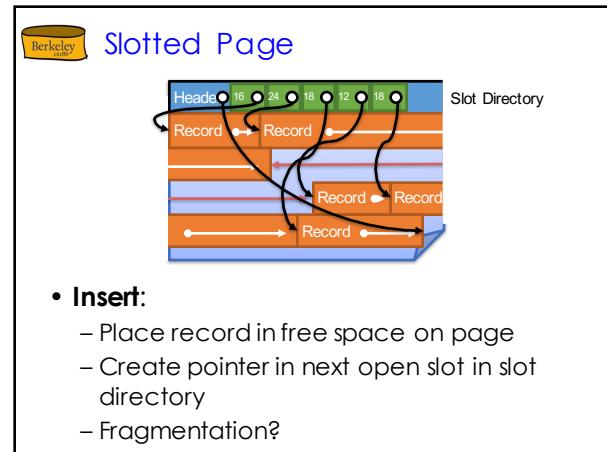
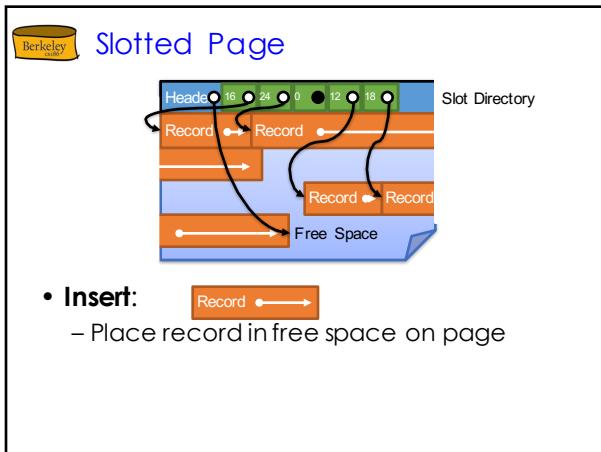
Table
Byte Rep. Record

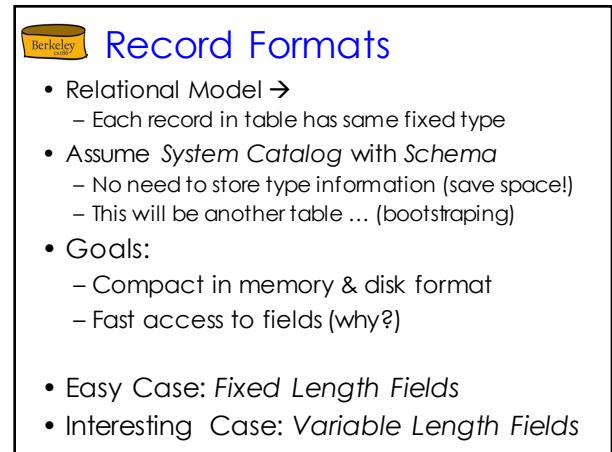
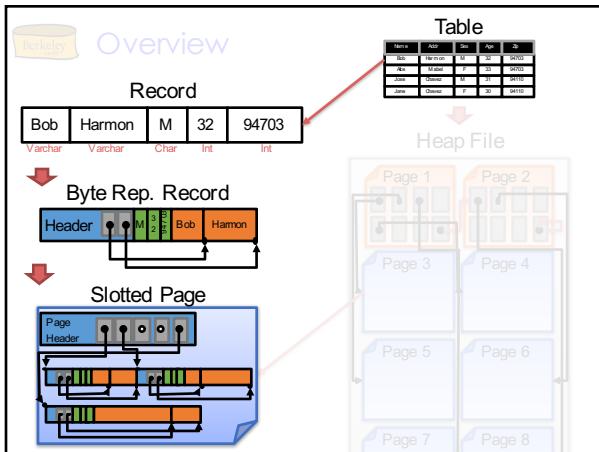
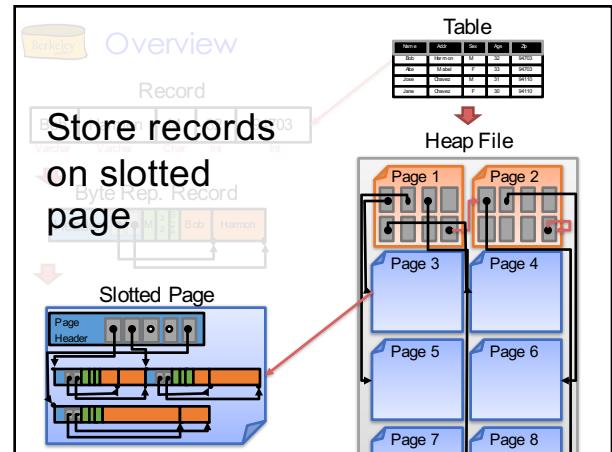
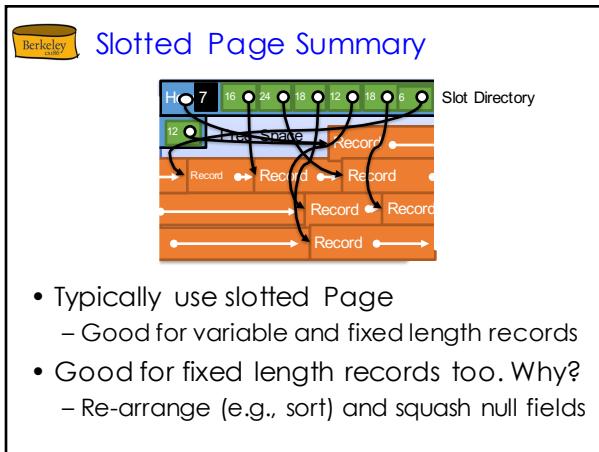
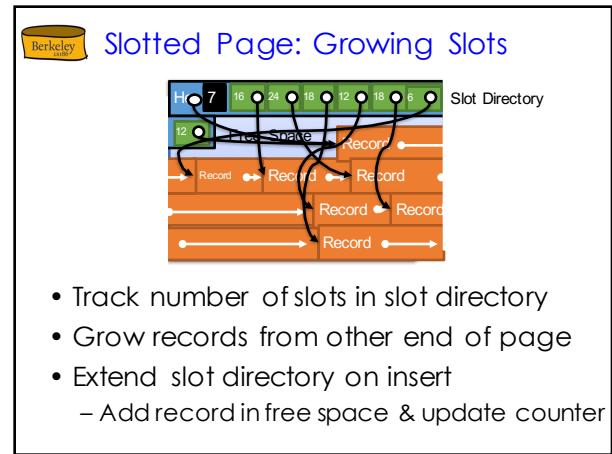
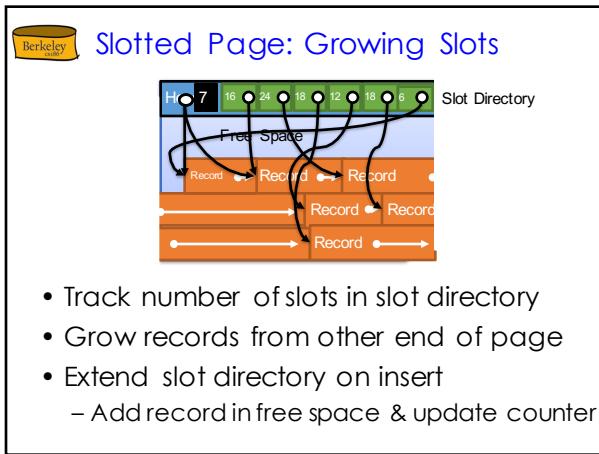
Table encoded as files which are collections of pages.



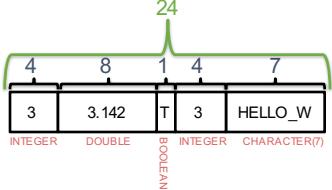








 Record Formats: Fixed Length



- Field types same for all records in a file.
 - Type info stored separately in *system catalog*
- On disk byte representation same as in memory
- Finding i 'th field?
 - done via arithmetic (fast)
- Compact? (Nulls?)

 Record Formats: Variable Length

What happens if fields are variable length?

Record

Bob	Big, St.	M	32	94703
VARCHAR	VARCHAR	CHAR	INT	INT

Could store with padding? (Fixed Length)

Wasted Space

Bob	Big, St.	M	32	94703
CHAR(20)	CHAR(18)	CHAR	INT	INT

Field Not Big Enough

Alice	Boulevard of the Allies	32	94703
CHAR(20)	CHAR(18)	CHAR	INT

 Record Formats: Variable Length

What happens if fields are variable length?

Record

Bob	Big, St.	M	32	94703
VARCHAR	VARCHAR	CHAR	INT	INT

Could use delimiters (i.e., CSV):

Comma Separated Values (CSV)

Bob	,	Big, St.	,	M	,	32	,	94703
VARCHAR		VARCHAR		CHAR		INT		INT

- Issues?

 Record Formats: Variable Length

What happens if fields are variable length?

Record

Bob	Big, St.	M	32	94703
VARCHAR	VARCHAR	CHAR	INT	INT

Could use delimiters (i.e., CSV):

Comma Separated Values (CSV)

Bob	,	Big, St.	,	M	,	32	,	94703
VARCHAR		VARCHAR		CHAR		INT		INT

- Requires scan to access field
- What if text contains commas?

 Record Formats: Variable Length

What happens if fields are variable length?

Record

Bob	Big, St.	M	32	94703
VARCHAR	VARCHAR	CHAR	INT	INT

Store length information before fields:

Variable Length Fields with Offsets

M	32	94703	3	Bob	8	Big, St.
CHAR	INT	INT	VARCHAR	VARCHAR		

Move all variable length fields to end
→ enable fast access

- Requires scan to access field
- What if text contains commas?

 Record Formats: Variable Length

What happens if fields are variable length?

Record

Bob	Big, St.	M	32	94703
VARCHAR	VARCHAR	CHAR	INT	INT

Introduce a record header:

Header	Bob	Big, St.	M	32	94703
CHAR	VARCHAR	CHAR	INT	INT	INT

- Requires scan to access field. Why?
- What if text contains commas?

Record Formats: Variable Length

What happens if fields are variable length?

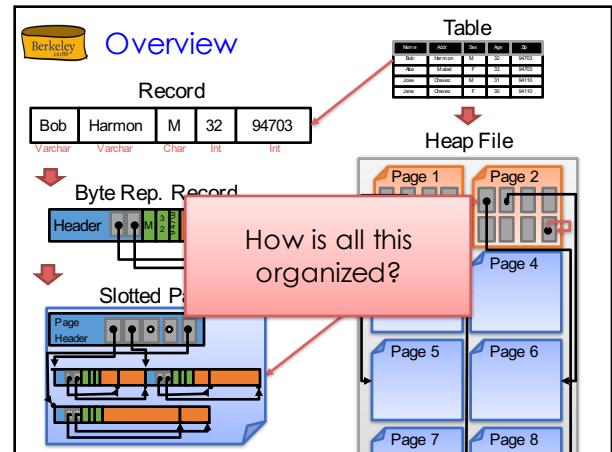
Record

Bob	Big, St.	M	32	94703
VARCHAR	VARCHAR	CHAR	INT	INT

Introduce a record header:

Header	M	32	94703	Bob	Big, St.
CHAR	INT	INT	VARCHAR	VARCHAR	

- Direct access & no “escaping”, other adv.?
 - Handle null fields → useful for fixed length



System Catalogs

- For each relation:
 - name, file location, file structure (e.g., Heap file)
 - attribute name and type, for each attribute
 - index name, for each index
 - integrity constraints
- For each index:
 - structure (e.g., B+ tree) and search key fields
- For each view:
 - view name and definition
- Plus statistics, authorization, buffer pool size, etc.

☞ Catalogs are themselves stored as relations!

pg_attribute

```
postgres=# \d pg_catalog.pg_attribute
              Table "pg_catalog.pg_attribute"
   Column    | Type | Modifiers
-----+-----+-----+
 attrelid   | oid  | not null
 attname    | name | not null
 atttypid   | oid  | not null
 attstattarget | integer | not null
 attlen     | smallint | not null
 attnum     | smallint | not null
 attndims   | integer | not null
 attstorageoff | integer | not null
 atttypmod  | integer | not null
 attbyval   | boolean | not null
 attstorage | "char"  | not null
 attalign   | "char"  | not null
 attnotnull | boolean | not null
 attinhdef  | boolean | not null
 attisdropped | boolean | not null
 attislocal  | boolean | not null
 attinhcount | integer | not null
 attcollation | oid  | not null
 attaci     | aciitem[] |
 attoptions  | text[]  |
 attfdwoptions | text[]  |
Indexes:
  "pg_attribute_relid_attname_index" UNIQUE, btree (attrelid, attname)
  "pg_attribute_relid_attnum_index" UNIQUE, btree (attrelid, attnum)
postgres@precise64:~$
```

Summary

- Disk manager loads and stores pages
 - Block level reasoning
 - Abstracts device and file system; provides fast next
- Buffer manager brings pages into RAM
 - page pinned while reading/writing
 - dirty pages written to disk
 - good replacement policy essential for performance
- DBMS “File” tracks collection of pages, records within each.
 - Heap-files: unordered records organized with directories

Summary (Contd.)

- Slotted page format
 - Variable length records and intra-page reorg
- Variable length record format
 - Direct access to i'th field and null values.
- Catalog relations store information about relations, indexes and views.

