# *Storage and Indexing*

## Computer Science Department
## Columbia University

# Why Not Store Everything in Main Memory?

❖ *Main memory is volatile.* We lose data between runs.

❖ *Cost.* 1GB of RAM will cost you ~ $30
  - At the end of 2011, 1GB of disk storage fell to $1.68 from $18.95 in 2005 (91% decline), according to IDC; it is expected to cost just pennies in a few years.
  - Solid state (flash) memory is a bit more expensive.

❖ Typical storage hierarchy:
  - Main memory (RAM) for currently used data.
  - (Solid state, not widely used yet but it will soon)
  - Disk for the main database (secondary storage).
  - Tapes for archiving older versions of the data (tertiary storage); not used for the operational part of the DB.

# *Disks (storage device of choice)*

❖ Data on disk is stored and retrieved in units called *disk blocks* or *pages*.

- READ: transfer data (a number of disk blocks) from disk to main memory (RAM).
- WRITE: transfer data (a number of disk blocks) from RAM to disk.
- Both are high-cost operations, relative to in-memory operations, so must be planned carefully!

❖ Unlike RAM, time to retrieve a disk page varies depending upon location on disk.

- Therefore, relative placement of pages on disk has major impact on DBMS performance!
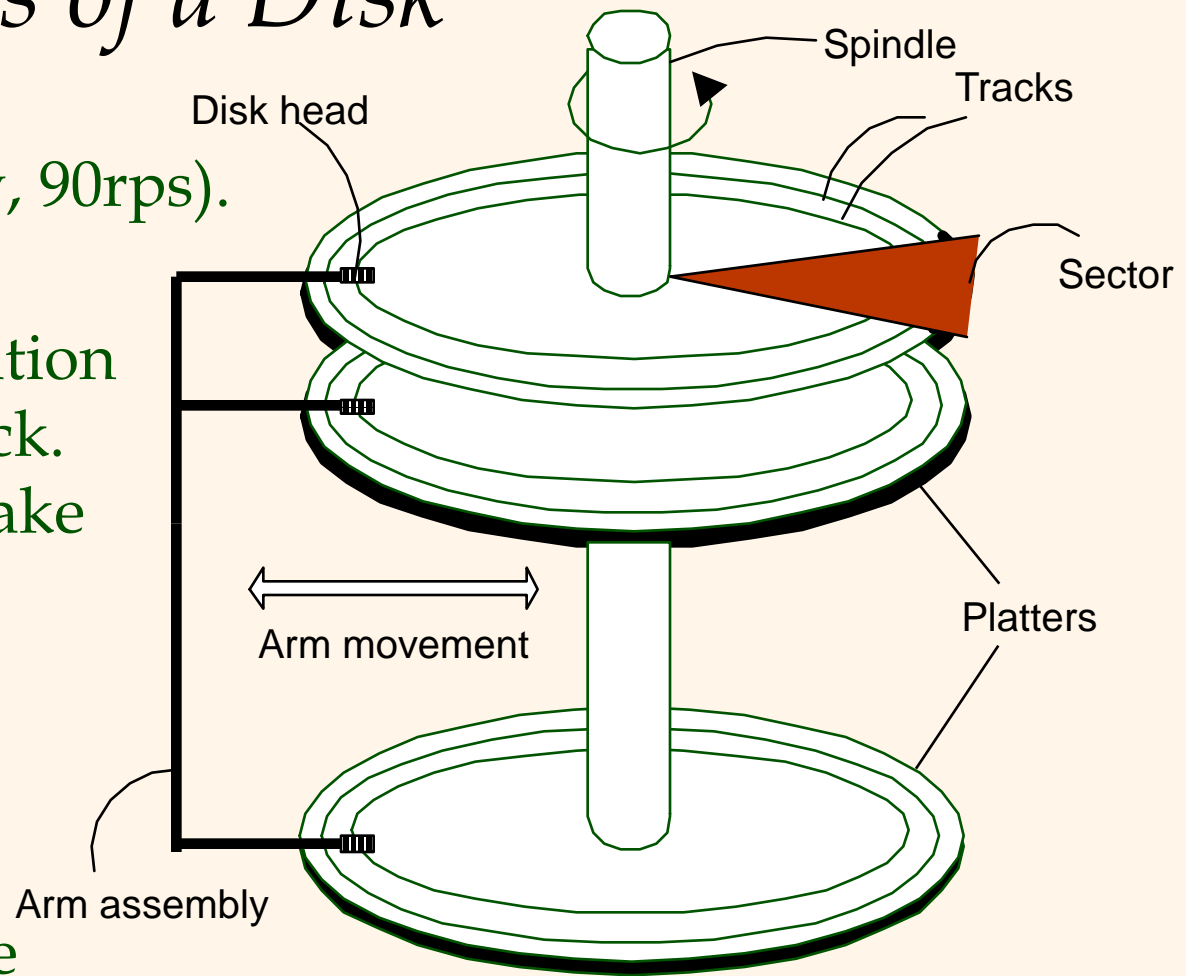
# *Components of a Disk*

❖ The platters spin (say, 90rps).

❖ The arm assembly is moved in or out to position a head on a desired track. Tracks under heads make a *cylinder* (imaginary!).

❖ Only one head reads/writes at any one time.

❖ *Block size* is a multiple of *sector size* (which is fixed).

Spindle

Tracks

Disk head

Sector

Arm movement

Platters

Arm assembly

# *Accessing a Disk Page*

❖ Time to access (read/write) a disk block:
  ▪ *seek time* (moving arms to position disk head on track) **+**
  ▪ *rotational delay* (waiting for block to rotate under head) **+**
  ▪ *transfer time* (actually moving data to/from disk surface)
❖ Seek time and rotational delay dominate.
  ▪ Seek time varies from about 1 to 20msec
  ▪ Rotational delay varies from 0 to 10msec
  ▪ Transfer rate is about 1msec per 4KB page
❖ Key to lower I/O cost: reduce seek/rotation delays!
❖ Compare disk access time of ~10msec with memory access time of ~60nsecs

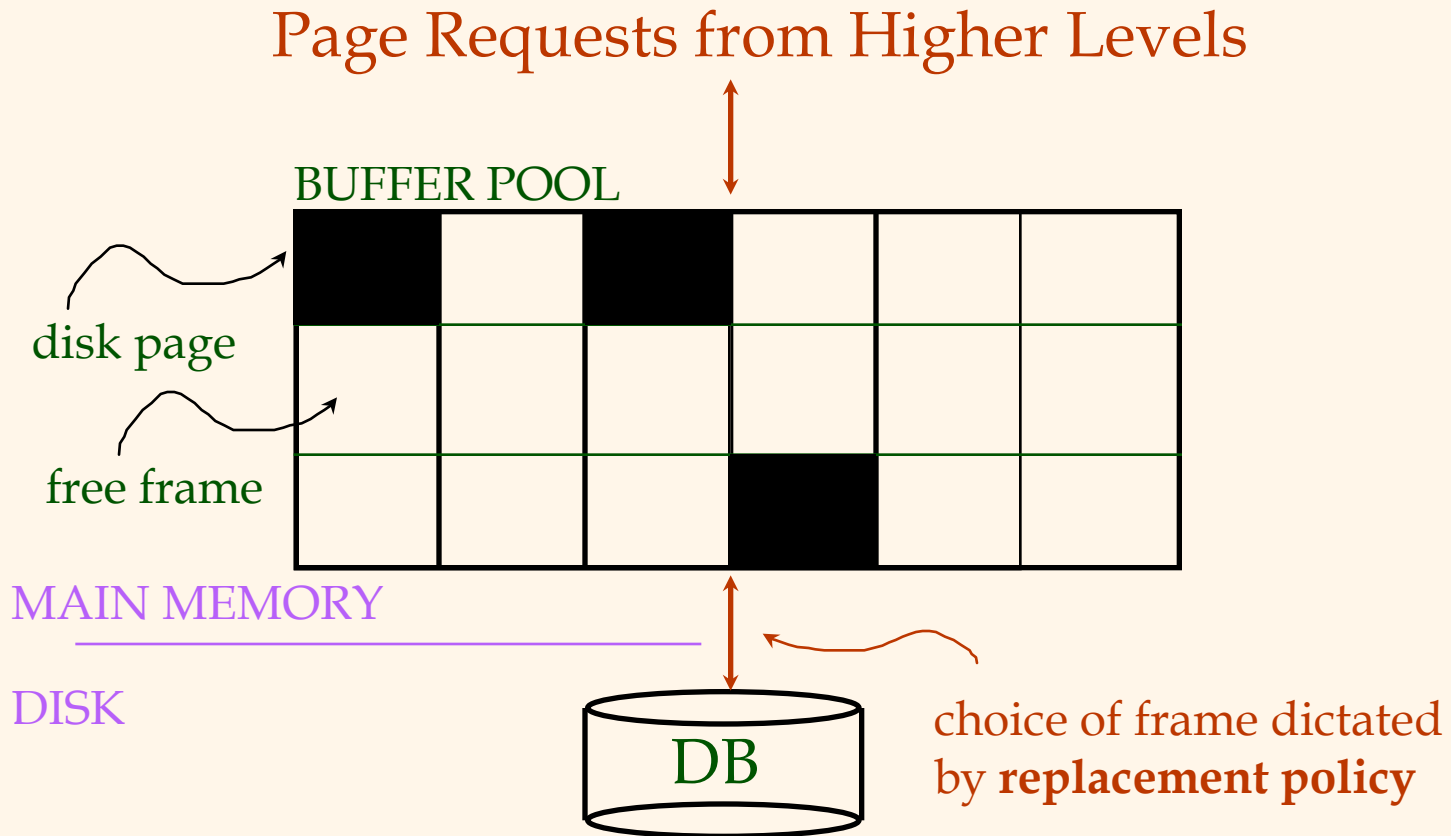# *Arranging Pages on Disk*

❖ `*Next*' block concept:
  ▪ blocks on same track, followed by
  ▪ blocks on same cylinder, followed by
  ▪ blocks on adjacent cylinder

❖ Blocks in a file should be arranged sequentially on disk (by `next'), to minimize seek and rotational delay.

❖ For a sequential scan, *pre-fetching* several pages at a time is a big win!

# *Disk Space Management*

❖ Lowest layer of DBMS software manages space on disk.

❖ Higher levels call upon this layer to:

  ▪ allocate(int n): allocate n pages, return the page id of the first page of the sequence of n pages

  ▪ deallocate(int pid, int n): free n consecutive pages starting at page id pid

❖ Request for a *sequence* of pages must be satisfied by allocating the pages sequentially on disk! Higher levels don't need to know how this is done, or how free space is managed.

# *Buffer Management in a DBMS*

Page Requests from Higher Levels

BUFFER POOL

disk page

free frame

MAIN MEMORY

DISK

DB

choice of frame dictated by **replacement policy**

- ❖ *Data must be in RAM for DBMS to operate on it!*
- ❖ *Table of <frame#, pageid> pairs is maintained.*
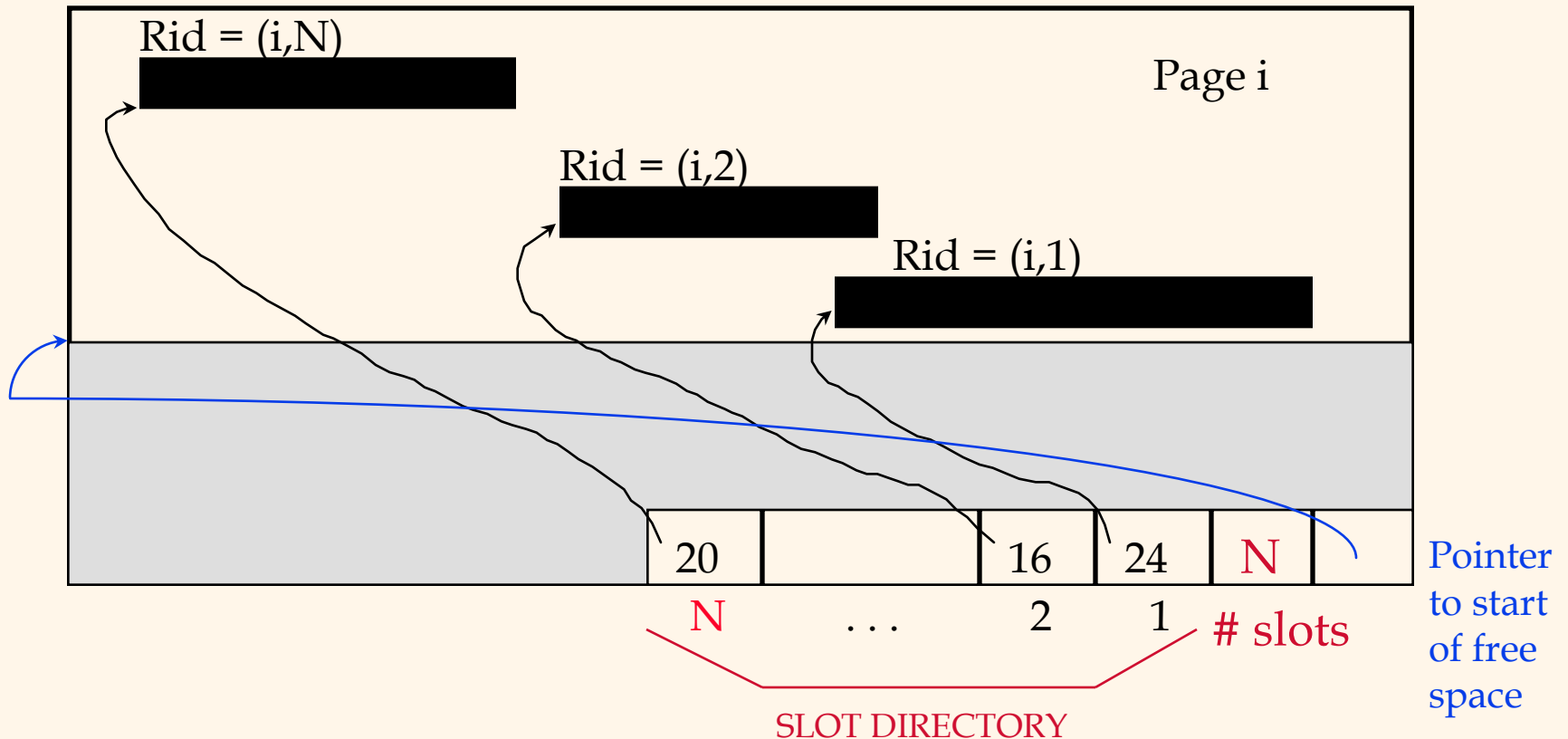
# *When a Page is Requested ...*

❖ If requested page is not in pool:
- Choose a frame for *replacement*
- If  frame is dirty, write it to disk
- Read requested page into chosen frame

❖ *Pin* the page and return its address.

❖ *If requests can be predicted (e.g., sequential scans) pages can be <u>pre-fetched</u> - several pages at a time!*

# DBMS *vs. OS File System*

OS does disk space & buffer mgmt: why not let OS manage these tasks?

❖ Differences in OS support: portability issues

❖ Some limitations, e.g., files can't span disks.

❖ Buffer management in DBMS requires ability to:
  - pin a page in buffer pool, force a page to disk (important for implementing CC & recovery),
  - adjust *replacement policy*, and pre-fetch pages based on access patterns in typical DB operations.

# Page Formats: Variable Length Records

Rid = (i,N)

Page i

Rid = (i,2)

Rid = (i,1)

| 20 | | 16 | 24 | N | |
| N | . . . | 2 | 1 | # slots | |

Pointer to start of free space

SLOT DIRECTORY

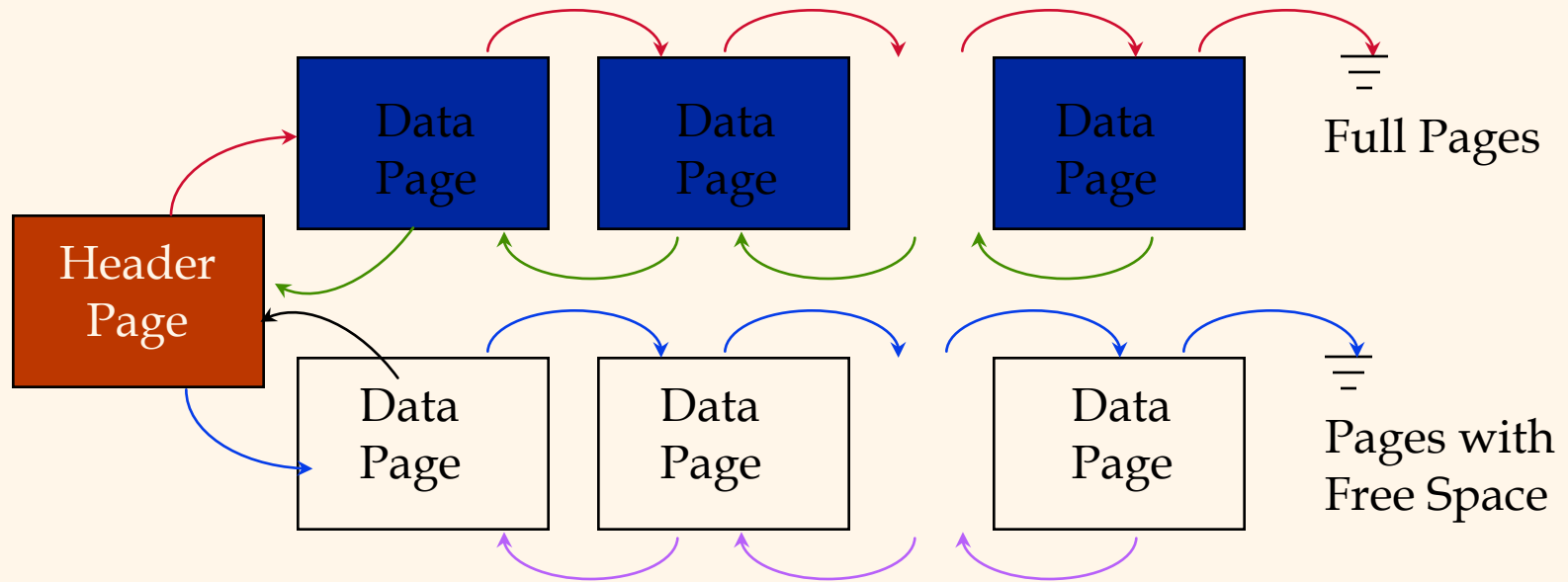- *Can move records on page without changing rid; so, attractive for fixed-length records too.*

# *Files of Records*

❖ Page or block is OK when doing I/O, but higher levels of DBMS operate on *records*, and *files of records*.

❖ <u>FILE</u>: A collection of pages, each containing a collection of records. Must support:

- insert/delete/modify record
- read a particular record (specified using *record id*)
- scan all records (possibly with some conditions on the records to be retrieved)
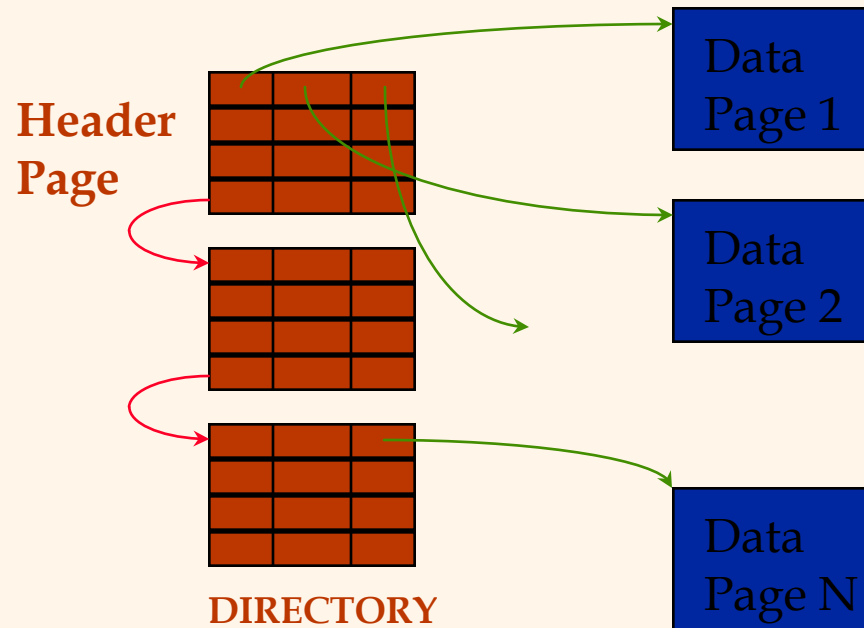
# *Unordered (Heap) Files*

❖ Simplest file structure contains records in no particular order.

❖ As file grows and shrinks, disk pages are allocated and 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.

# *Heap File Implemented as a List*



- ❖ The header page id and Heap file name must be stored someplace.
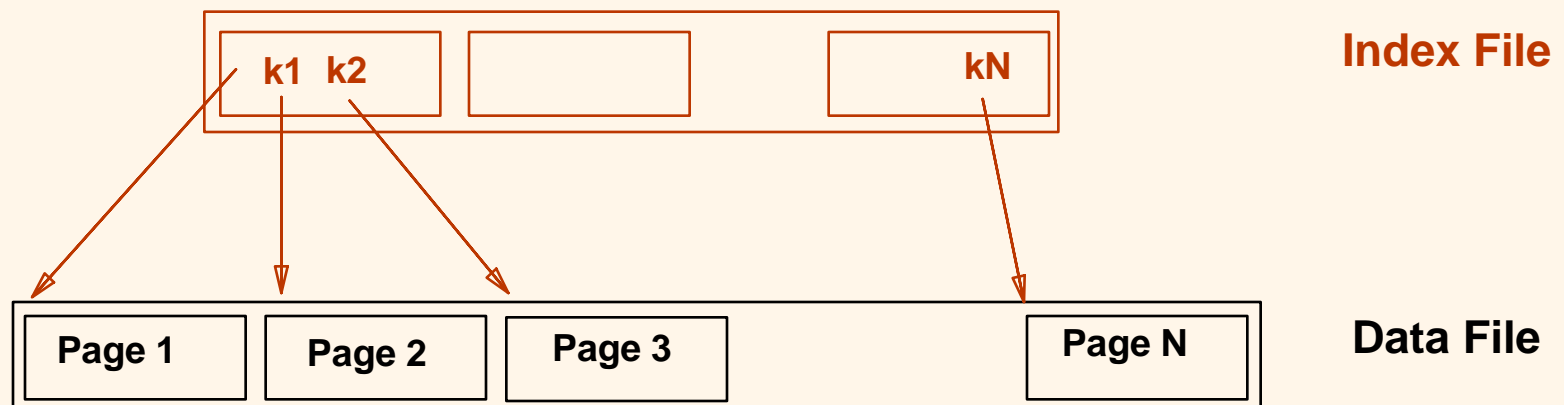- ❖ Each page contains 2 `pointers' plus data.

# *Heap File Using a Page Directory*

**Header Page**

**DIRECTORY**

Data Page 1

Data Page 2

Data Page N

❖ The entry for a page can include the number of free bytes on the page.

❖ The directory is a collection of pages; linked list implementation is just one alternative.

- ▪ *Much smaller than linked list of all HF pages*!

# *Index File*

❖ Example: *Find all students with gpa > 3.0*
- Assume sorted data, do binary search to find first such student, then scan to find others.
- Cost of binary search can be quite high.

❖ Simple idea:  Create an index file and *do binary search on (smaller) index file!*

| k1  k2 | | kN | **Index File** |
|---|---|---|---|

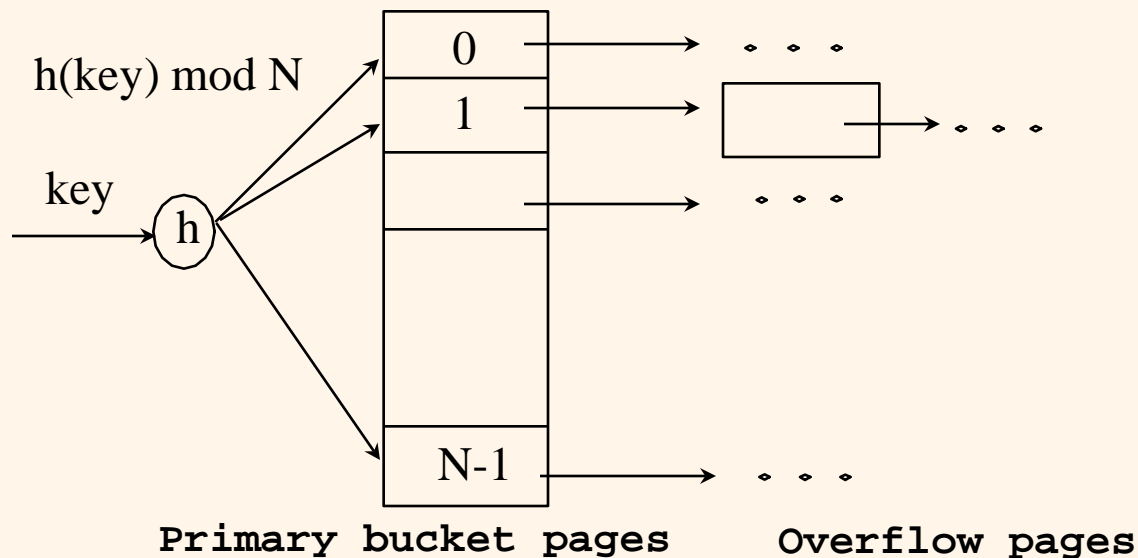| Page 1 | Page 2 | Page 3 | | Page N | **Data File** |
|---|---|---|---|---|---|

# *Indexes*

❖ An *index* on a file speeds up selections on the *search key fields* for the index.

- Any subset of the fields of a relation can be the search key for an index on the relation.

- *Search key* is not the same as *key* (minimal set of fields that uniquely identify a record in a relation).

❖ An index contains a collection of *data entries*, and supports efficient retrieval of all data entries **k\*** with a given key value **k**.

# *Hash-Based Indexes*

❖ Good for equality selections.

- Index is a collection of *buckets*. Bucket = *primary* page plus zero or more *overflow* pages.

- *Hashing function* **h**:  **h**(*r*) = bucket in which record *r* belongs. **h** looks at the *search key* fields of *r*.

# *Hash-Based Indexing*

❖ # primary pages fixed, allocated sequentially, never de-allocated; overflow pages if needed.

❖ **h**(*k*) mod N = bucket to which data entry with key *k* belongs. (N = # of buckets)



h(key) mod N

key

0
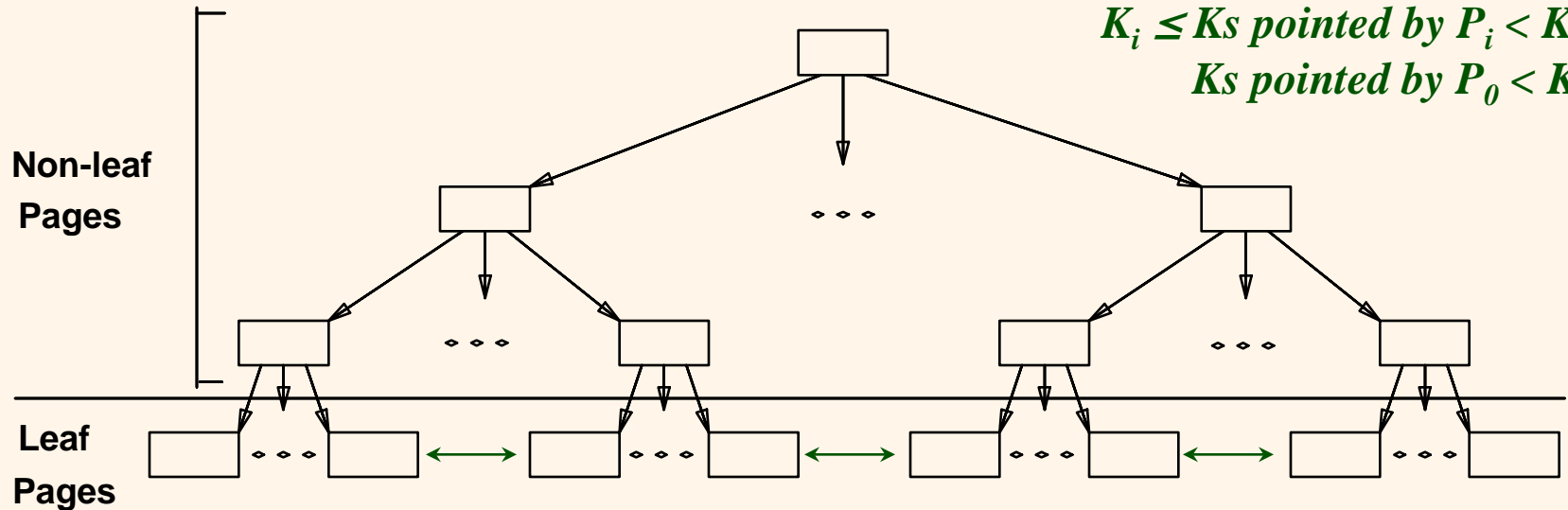1

N-1

**Primary bucket pages**          **Overflow pages**

# A Tree-Structured Index: B+ Tree Indexes

**Rule:**
$K_i \leq$ *Ks pointed by* $P_i < K_{i+1}$
*Ks pointed by* $P_0 < K_1$

**Non-leaf Pages**

**Leaf Pages**
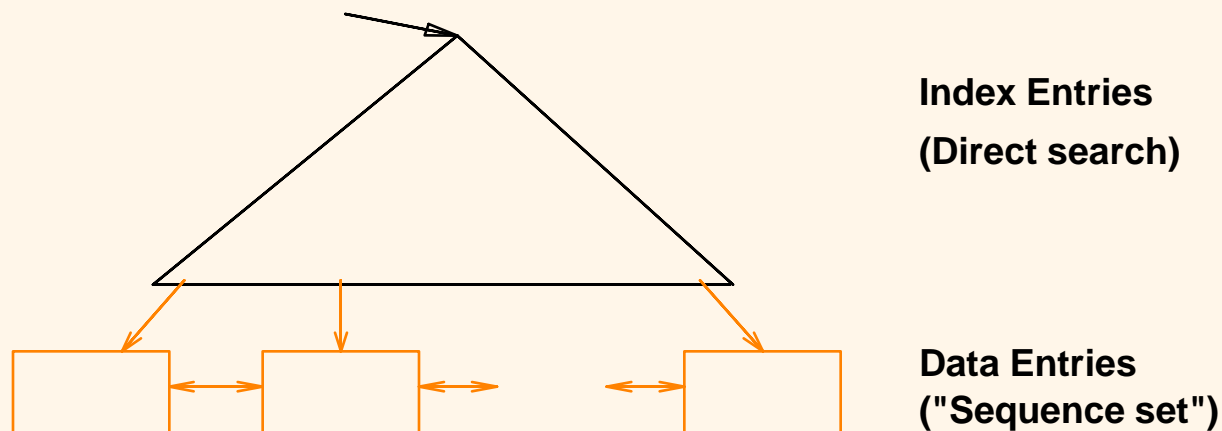


❖ Leaf pages contain *data entries*, and are chained (prev & next)
❖ Non-leaf pages contain *index entries* and direct searches:

**index entry**

| $P_0$ | $K_1$ | $P_1$ | $K_2$ | $P_2$ | ◇   ◇   ◇ | $K_m$ | $P_m$ |
|---|---|---|---|---|---|---|---|

# B+ Tree: Most Widely Used Index

❖ Keeps tree <u>always</u> *height-balanced*.

❖ Search starts at the root and follows pointers to leaf nodes and to data records.

❖ Nodes are at least 50% full (except for root).

❖ Supports equality and range-searches efficiently.

❖ For all practical purposes, what is the performance of a B-tree search?

**Index Entries**
**(Direct search)**

**Data Entries**
**("Sequence set")**

# *Example B+ Tree*

❖ Search begins at root, and key comparisons direct it to a leaf.

❖ Search for 5*, 15*, all data entries >= 24* ...

**Root**

| | 13 | 17 | 24 | 30 | |

| 2* | 3* | 5* | 7* | | 14* | 16* | | | | 19* | 20* | 22* | | | 24* | 27* | 29* | | | 33* | 34* | 38* | 39* |

*Based on the search for 15*, we <u>know</u> it is not in the tree!*

# *Alternatives for Data Entry k\* in Index*

❖ Three alternatives:

1. Data record with key value **k**
2. <**k**, rid of data record with search key value **k**>
3. <**k**, list of rids of data records with search key **k**>

❖ Choice of alternative for data entries is orthogonal to the indexing technique used to locate data entries with a given key value **k**.

- Examples of indexing techniques: B+ trees, hash-based structures

- Typically, index contains auxiliary information that directs searches to the desired data entries

# *Alternatives for Data Entries (Contd.)*

❖ Alternative 1:
  ▪ If this is used, index structure is a file organization for data records (instead of a Heap file or sorted file).

  ▪ At most one index on a given collection of data records can use Alternative 1.  (Otherwise, data records are duplicated, leading to redundant storage and potential inconsistency.)

  ▪ If data records are very large,  # of pages containing data entries is high.  Implies size of auxiliary information in the index is also large, typically.

# *Alternatives for Data Entries (Contd.)*

❖ Alternatives 2 and 3:

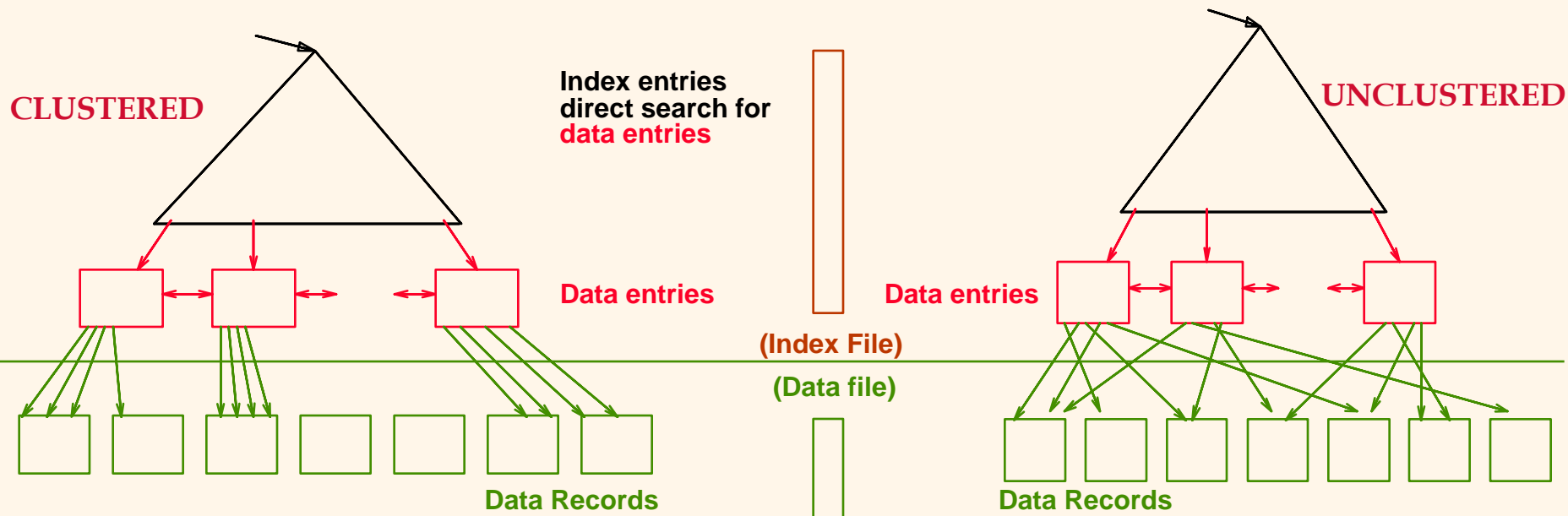- Data entries typically much smaller than data records.  So, better than Alternative 1 with large data records, especially if search keys are small. (Portion of index structure used to direct search, which depends on size of data entries, is much smaller than with Alternative 1.)

- Alternative 3 more compact than Alternative 2, but leads to variable sized data entries even if search keys are of fixed length.

# *Index Classification*

❖ *Clustered* vs. *unclustered*:  If order of data records is the same as, or `close to', order of data entries, then called clustered index.

  ▪ Alternative 1 implies clustered; in practice, clustered also implies Alternative 1 (since sorted files are rare).

  ▪ A file can be clustered on at most one search key.

  ▪ Cost of retrieving data records through index varies *greatly* based on whether index is clustered or not!

# *Clustered vs. Unclustered Index*

❖ Example with Alternative (2), ie., data entries of the form **<k**, rid of data record with search key value **k>**.



CLUSTERED

Index entries
direct search for
data entries

Data entries

(Index File)

(Data file)

Data Records

UNCLUSTERED
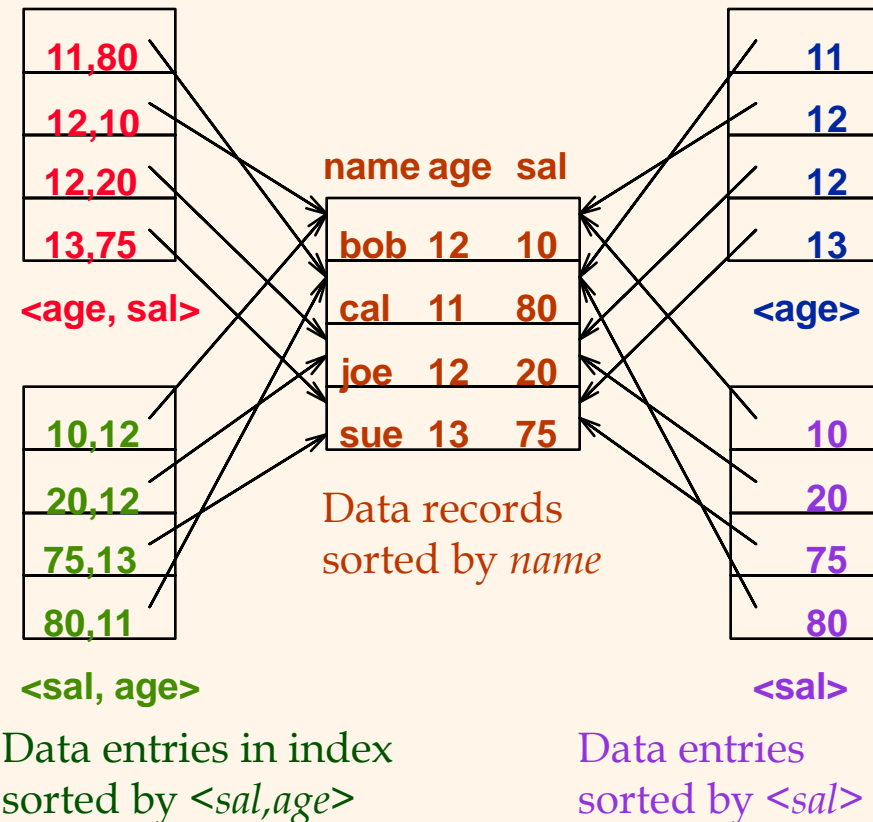
Data entries

Data Records

# *Index Classification (Contd.)*

❖ *Composite Search Keys*: Search on a combination of fields.

- Equality query: Every field value is equal to a constant value. E.g. wrt <sal,age> index:
  - age=20 and sal =75
- Range query: Some field value is not a constant. E.g.:
  - age =20; or age=20 and sal > 10

❖ Data entries in index sorted by search key to support range queries.
- Lexicographic order, or
- Spatial order.

Examples of composite key indexes using lexicographic order.

| 11,80 |
| 12,10 |
| 12,20 |
| 13,75 |

**<age, sal>**

| 10,12 |
| 20,12 |
| 75,13 |
| 80,11 |

**<sal, age>**

Data entries in index sorted by <*sal,age*>

| name | age | sal |
|------|-----|-----|
| bob | 12 | 10 |
| cal | 11 | 80 |
| joe | 12 | 20 |
| sue | 13 | 75 |

Data records sorted by *name*

| 11 |
| 12 |
| 12 |
| 13 |

**<age>**

| 10 |
| 20 |
| 75 |
| 80 |

**<sal>**

Data entries sorted by <*sal*>

# *Understanding the Workload*

❖ For each query in the workload:

  ▪ Which relations does it access?

  ▪ Which attributes are retrieved?

  ▪ Which attributes are involved in selection/join conditions? How selective are these conditions likely to be?

❖ For each update in the workload:

  ▪ Which attributes are involved in selection/join conditions? How selective are these conditions likely to be?

  ▪ The type of update (INSERT/DELETE/UPDATE), and the attributes that are affected.

# *Choice of Indexes*

❖ What indexes should we create?

  ▪ Which relations should have indexes?  What field(s) should be the search key?  Should we build several indexes?

❖ For each index, what kind of an index should it be?

  ▪ Clustered?  Hash/tree?

# *Choice of Indexes (Contd.)*

❖ One approach: Consider the most important queries in turn. Consider the best plan using the current indexes, and see if a better plan is possible with an additional index. If so, create it.

   ▪ Obviously, this implies that we must understand how a DBMS evaluates queries and creates query evaluation plans!

   ▪ For now, we discuss simple 1-table queries.

❖ Before creating an index, must also consider the impact on updates in the workload!

   ▪ Trade-off: Indexes can make queries go faster, updates slower. Require disk space, too.

# *Index Selection Guidelines*

❖ Attributes in WHERE clause are candidates for index keys.

- Exact match condition suggests hash index.
- Range query suggests tree index.
  - Clustering is especially useful for range queries; can also help on equality queries if there are many duplicates.

❖ Multi-attribute search keys should be considered when a WHERE clause contains several conditions.

- Order of attributes is important for range queries.
- Such indexes can sometimes enable index-only strategies for important queries.
  - For index-only strategies, clustering is not important!

❖ Try to choose indexes that benefit as many queries as possible.  Since only one index can be clustered per relation, choose it based on important queries that would benefit the most from clustering.

# *Index Specification in SQL*

❖ No specification in the standard!!

❖ In practice all commercial DBMSs support indexing (no exception).


CREATE INDEX Ind1 ON Students
WITH STRUCTURE = BTREE
KEY = (age, gpa)

# *Summary*

❖ Disks provide cheap, non-volatile storage.

  ▪ Random access, but cost depends on location of page on disk; important to arrange data sequentially to minimize *seek* and *rotation* delays.

❖ Buffer manager brings pages into RAM.

  ▪ Page stays in RAM until released by requestor.

  ▪ Written to disk when frame chosen for replacement (which is sometime after requestor releases the page).

  ▪ Choice of frame to replace based on *replacement policy*.

  ▪ Tries to *pre-fetch* several pages at a time.

# *Summary (Contd.)*

❖ DBMS vs. OS File Support
  - DBMS needs features not found in many OS's, e.g., forcing a page to disk, controlling the order of page writes to disk, files spanning disks, ability to control pre-fetching and page replacement policy based on predictable access patterns, etc.

❖ Variable length record format with field offset directory offers support for direct access to i'th field and null values.

❖ Slotted page format supports variable length records and allows records to move on page.

# *Summary (Contd.)*

❖ File layer keeps track of pages in a file, and supports abstraction of a collection of records.

  ▪ Pages with free space identified using linked list or directory structure (similar to how pages in file are kept track of).

❖ Indexes support efficient retrieval of records based on the values in some fields.

# *Summary*

❖ **If selection queries are frequent, sorting the file or building an *index* is important.**
  - Hash-based indexes only good for equality search.
  - Sorted files and tree-based indexes best for range search; also good for equality search.  (Files rarely kept sorted in practice; B+ tree index is better.)

❖ **Index is a collection of data entries plus a way to quickly find entries with given key values.**

❖ **Data entries can be actual data records, <key, rid> pairs, or <key, rid-list> pairs.**
  - Choice orthogonal to the *indexing technique* used.

# *Summary (Contd.)*

❖ Can have several indexes on a given file of data records, each with a different search key.

❖ Indexes can be classified as clustered vs. unclustered, primary vs. secondary, and dense vs. sparse.

❖ Understanding the nature of the *workload* for the application, and the performance goals, is essential to developing a good design.

❖ Indexes must be chosen to speed up important queries (and perhaps some updates!).
  ▪ Choose indexes that can help many queries, if possible.
  ▪ Build indexes to support index-only strategies.
  ▪ Clustering: only one index on a given relation can be clustered!