# CMSC 424 – Database design
## Lecture 13
## Storage: Files

Mihai Pop

# Recap

- Databases are stored on disk
  - cheaper than memory
  - non-volatile (survive power loss)
  - large capacity
- Operating systems are designed for "general" use – do not perform optimally when used to manage database storage
- Most DBMSs replace the OS and manage disk storage directly
  - Specialized buffer management (MRU policy might be better than LRU, pinned records, etc.)
  - Specialized storage of files (today)

# File Organization

- The database is stored as a collection of *files*.  Each file is a sequence of *records*.  A record is a sequence of fields.
- One approach:
  - assume record size is fixed
  - each file has records of one particular type only
  - different files are used for different relations
  This case is easiest to implement; will consider variable length records later.
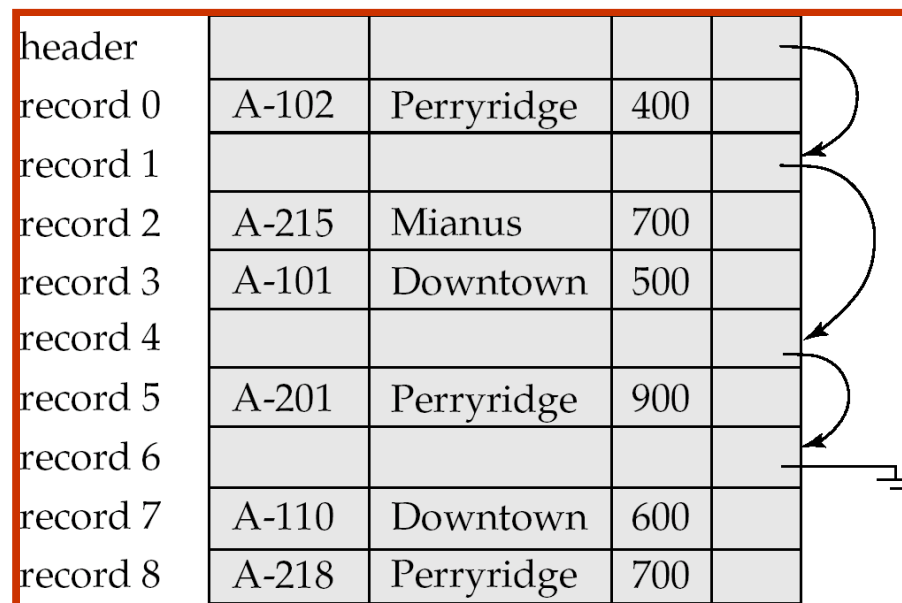
# Fixed-Length Records

- Simple approach:
  - Store record $i$ starting from byte $n * (i - 1)$, where $n$ is the size of each record.
  - Record access is simple but records may cross blocks
    - Modification: do not allow records to cross block boundaries

- Deletion of record $i$: alternatives:
  - move records $i + 1, \ldots, n$ to $i, \ldots, n - 1$
  - move record $n$ to $i$
  - do not move records, but link all free records on a *free list*

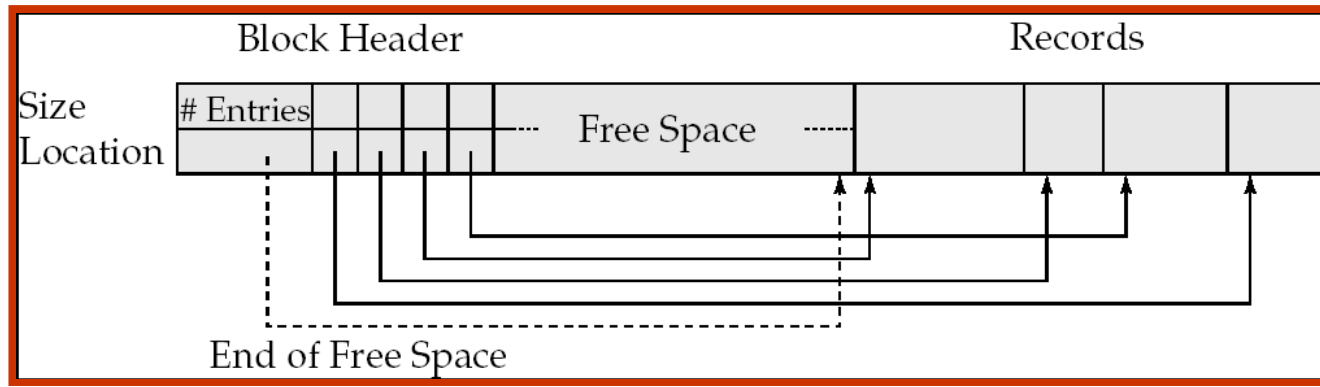| record 0 | A-102 | Perryridge | 400 |
|----------|-------|------------|-----|
| record 1 | A-305 | Round Hill | 350 |
| record 2 | A-215 | Mianus     | 700 |
| record 3 | A-101 | Downtown   | 500 |
| record 4 | A-222 | Redwood    | 700 |
| record 5 | A-201 | Perryridge | 900 |
| record 6 | A-217 | Brighton   | 750 |
| record 7 | A-110 | Downtown   | 600 |
| record 8 | A-218 | Perryridge | 700 |

# Free Lists

- Store the address of the first deleted record in the file header.
- Use this first record to store the address of the second deleted record, and so on
- Can think of these stored addresses as pointers since they "point" to the location of a record.
- More space efficient representation:  reuse space for normal attributes of free records to store pointers.  (No pointers stored in in-use records.)

| | | | | |
|---|---|---|---|---|
| header | | | | |
| record 0 | A-102 | Perryridge | 400 | |
| record 1 | | | | |
| record 2 | A-215 | Mianus | 700 | |
| record 3 | A-101 | Downtown | 500 | |
| record 4 | | | | |
| record 5 | A-201 | Perryridge | 900 | |
| record 6 | | | | |
| record 7 | A-110 | Downtown | 600 | |
| record 8 | A-218 | Perryridge | 700 | |

# Variable-Length Records

- Variable-length records arise in database systems in several ways:
  - Storage of multiple record types in a file.
  - Record types that allow variable lengths for one or more fields.
  - Record types that allow repeating fields (used in some older data models).

# Variable-Length Records: Slotted Page Structure



- Slotted page header contains:
  - number of record entries
  - end of free space in the block
  - location and size of each record
- Records can be moved around within a page to keep them contiguous with no empty space between them; entry in the header must be updated.
- Pointers should not point directly to record — instead they should point to the entry for the record in header.

# Organization of Records in Files

- **Heap** – a record can be placed anywhere in the file where there is space

- **Sequential** – store records in sequential order, based on the value of the search key of each record

- **Hashing** – a hash function computed on some attribute of each record; the result specifies in which block of the file the record should be placed

- Records of each relation may be stored in a separate file. In a **multitable clustering file organization** records of several different relations can be stored in the same file
  - Motivation: store related records on the same block to minimize I/O

# Sequential File Organization

- Suitable for applications that require sequential processing of the entire file
- The records in the file are ordered by a search-key

| A-217 | Brighton | 750 | |
|-------|----------|-----|---|
| A-101 | Downtown | 500 | |
| A-110 | Downtown | 600 | |
| A-215 | Mianus | 700 | |
| A-102 | Perryridge | 400 | |
| A-201 | Perryridge | 900 | |
| A-218 | Perryridge | 700 | |
| A-222 | Redwood | 700 | |
| A-305 | Round Hill | 350 | |

# Sequential File Organization (Cont.)

- Deletion – use pointer chains
- Insertion –locate the position where the record is to be inserted
  - if there is free space insert there
  - if no free space, insert the record in an overflow block
  - In either case, pointer chain must be updated
- Need to reorganize the file from time to time to restore sequential order

| A-217 | Brighton | 750 | |
|-------|----------|-----|--|
| A-101 | Downtown | 500 | |
| A-110 | Downtown | 600 | |
| A-215 | Mianus | 700 | |
| A-102 | Perryridge | 400 | |
| A-201 | Perryridge | 900 | |
| A-218 | Perryridge | 700 | |
| A-222 | Redwood | 700 | |
| A-305 | Round Hill | 350 | |
| | | | |
| A-888 | North Town | 800 | |

# Multitable Clustering File Organization

Store several relations in one file using a **multitable clustering** file organization

| customer_name | account_number |
|:---:|:---:|
| Hayes | A-102 |
| Hayes | A-220 |
| Hayes | A-503 |
| Turner | A-305 |

| customer_name | customer_street | customer_city |
|:---:|:---:|:---:|
| Hayes | Main | Brooklyn |
| Turner | Putnam | Stamford |

Multitable clustering organization of *customer* and *depositor:*

| Hayes | Main | Brooklyn |
|-------|-------|-----------|
| Hayes | A-102 | |
| Hayes | A-220 | |
| Hayes | A-503 | |
| Turner | Putnam | Stamford |
| Turner | A-305 | |

- good for queries involving *depositor* ⋈ *customer*, and for queries involving one single customer and his accounts
- bad for queries involving only customer
- results in variable size records
- Can add pointer chains to link records of a particular relation

# Data Dictionary Storage

- Data dictionary (also called system catalog) stores metadata; that is, data about data, such as:

- Information about relations
    - names of relations
    - names and types of attributes of each relation
    - names and definitions of views
    - integrity constraints

- User and accounting information, including passwords

- Statistical and descriptive data
    - number of tuples in each relation

- Physical file organization information
    - How relation is stored (sequential/hash/…)
    - Physical location of relation

- Information about indices (Chapter 12)

# Data Dictionary Storage (Cont.)

- Catalog structure
  - Relational representation on disk
  - specialized data structures designed for efficient access, in memory

- A possible catalog representation:

*Relation_metadata =* (*relation_name, number_of_attributes,*
                      *storage_organization, location*)

*Attribute_metadata =* (*attribute_name, relation_name, domain_type,*
                      *position, length*)

*User_metadata =* (*user_name, encrypted_password, group*)

*Index_metadata =* (*index_name, relation_name, index_type,*
                      *index_attributes*)

*View_metadata =* (*view_name, definition*)

# Indexing...rationale

- Remember the "join" function
  - assume tables R1($\underline{A1}$, B), R2($\underline{A2}$, C)

```
for t1 in R1
  for t2 in R2
    if (t1[A1] == t2[A2])
      output (t1[A1], t1[B], t2[C])
    end
  end
end
```

    running time - #tuples in R1 * # tuples in R2
- Can we do better?
  - what if the tables were written in sorted files

# Indexing...rationale

- Better algorithm

```
while not end of R1 or R2
  while (t1[A1] < t2[A2])
    t1 = next
  end
  while (t1[A1] > t2[A2])
    t2 = next
  end
  foreach t1 & t2 st. t1[A1] == t2[A2]
    output (t1[A1], t1[B], t2[C])
  end
  advance t1 and t2 to next difference
end
```

running time min (# tuples in R1, # tuples in R2) + "size of largest cluster of equal keys"

# Indexing...rationale

- Sorting makes things faster
- What if we have more than one key on which we join?
- Store a separate index for each key
  - file of pointers to the records
  - order of pointers in the index corresponds to ordering of key values
- Multiple indices – we can sort the same file in different ways

# Basic Concepts

- Indexing mechanisms used to speed up access to desired data.
  - E.g., author catalog in library
- **Search Key** - attribute to set of attributes used to look up records in a file.
- An **index file** consists of records (called **index entries**) of the form

| search-key | pointer |
|------------|---------|

- Index files are typically much smaller than the original file
- Two basic kinds of indices:
  - **Ordered indices:** search keys are stored in sorted order
  - **Hash indices:** search keys are distributed uniformly across "buckets" using a "hash function".

# Index Evaluation Metrics

- Access types supported efficiently.  E.g.,
  - records with a specified value in the attribute
  - or records with an attribute value falling in a specified range of values.
- Access time
- Insertion time
- Deletion time
- Space overhead

# Ordered Indices

- In an **ordered index**, index entries are stored sorted on the search key value.  E.g., author catalog in library.
- **Primary index**: in a sequentially ordered file, the index whose search key specifies the sequential order of the file.
  - Also called **clustering index**
  - The search key of a primary index is usually but not necessarily the primary key.
- **Secondary index**: an index whose search key specifies an order different from the sequential order of the file.  Also called non-clustering index.
- Index-sequential file: ordered sequential file with a primary index.

# Dense Index Files

- Dense index — Index record appears for every search-key value in the file.

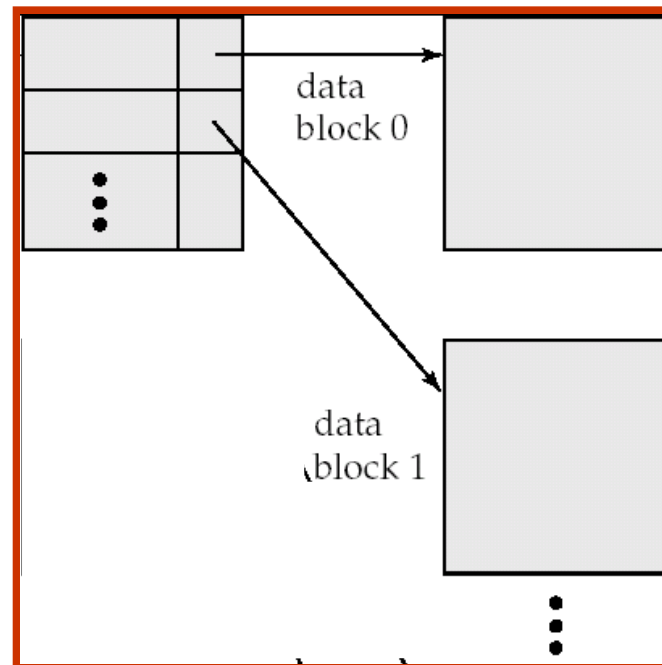| | | | | |
|---|---|---|---|---|
| Brighton | | A-217 | Brighton | 750 |
| Downtown | | A-101 | Downtown | 500 |
| Mianus | | A-110 | Downtown | 600 |
| Perryridge | | A-215 | Mianus | 700 |
| Redwood | | A-102 | Perryridge | 400 |
| Round Hill | | A-201 | Perryridge | 900 |
| | | A-218 | Perryridge | 700 |
| | | A-222 | Redwood | 700 |
| | | A-305 | Round Hill | 350 |

# Sparse Index Files

- Sparse Index: contains index records for only some search-key values.
  - Applicable when records are sequentially ordered on search-key
- To locate a record with search-key value *K* we:
  - Find index record with largest search-key value < *K*
  - Search file sequentially starting at the record to which the index record points

| | | A-217 | Brighton | 750 | |
|---|---|---|---|---|---|
| Brighton | | A-101 | Downtown | 500 | |
| Mianus | | A-110 | Downtown | 600 | |
| Redwood | | A-215 | Mianus | 700 | |
| | | A-102 | Perryridge | 400 | |
| | | A-201 | Perryridge | 900 | |
| | | A-218 | Perryridge | 700 | |
| | | A-222 | Redwood | 700 | |
| | | A-305 | Round Hill | 350 | |

# Sparse Index Files (Cont.)

- Compared to dense indices:
  - Less space and less maintenance overhead for insertions and deletions.
  - Generally slower than dense index for locating records.
- **Good tradeoff**: sparse index with an index entry for every block in file, corresponding to least search-key value in the block.
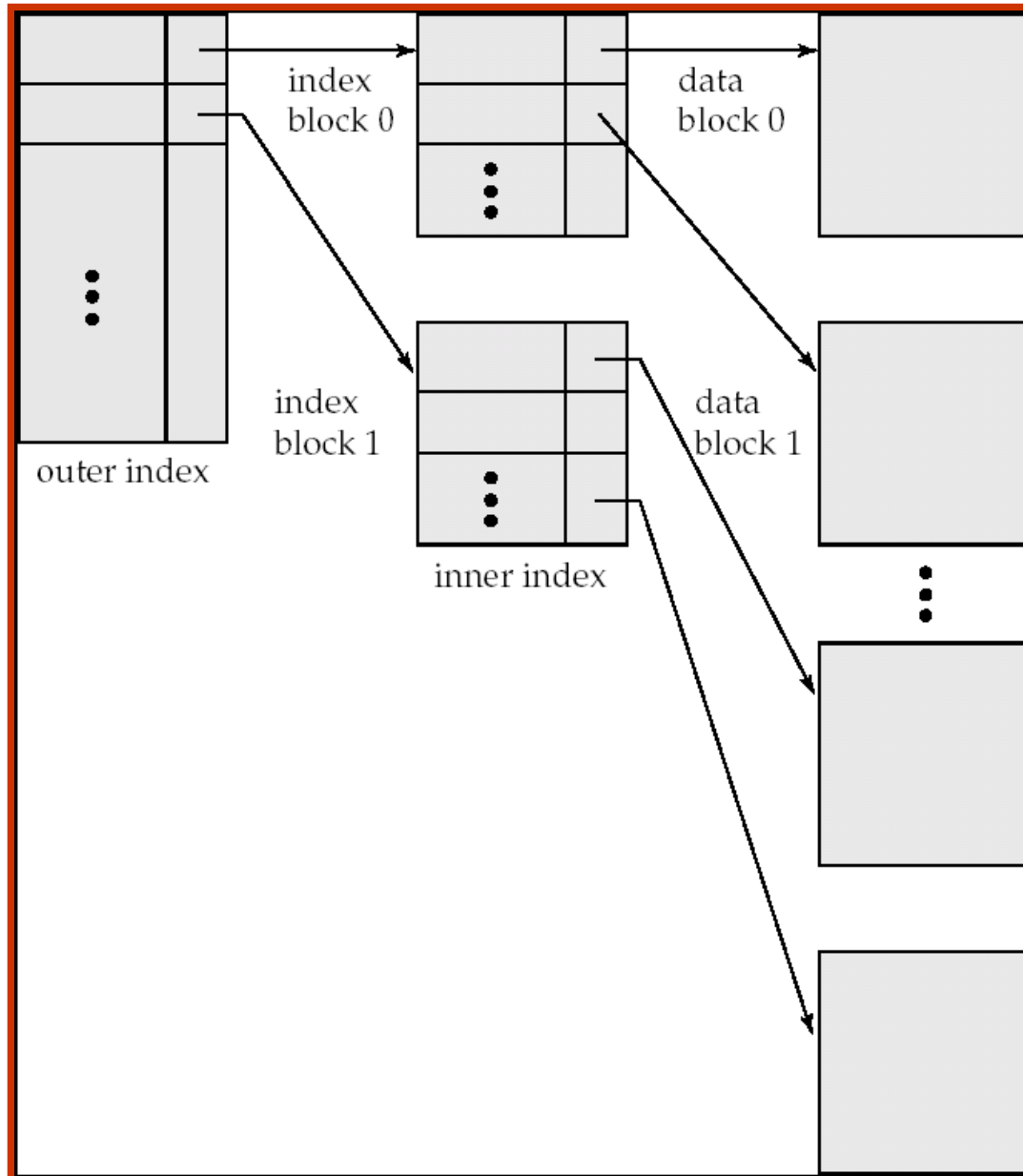


data block 0

data block 1

# Multilevel Index

- If primary index does not fit in memory, access becomes expensive.

- Solution: treat primary index kept on disk as a sequential file and construct a sparse index on it.

  - outer index – a sparse index of primary index

  - inner index – the primary index file

- If even outer index is too large to fit in main memory, yet another level of index can be created, and so on.

- Indices at all levels must be updated on insertion or deletion from the file.

# Multilevel Index (Cont.)

# Index Update:  Deletion

- If deleted record was the only record in the file with its particular search-key value, the search-key is deleted from the index also.

- Single-level index deletion:
  - **Dense indices** – deletion of search-key:similar to file record deletion.
  - **Sparse indices** –
    - if an entry for the search key exists in the index, it is deleted by replacing the entry in the index with the next search-key value in the file (in search-key order).
    - If the next search-key value already has an index entry, the entry is deleted instead of being replaced.

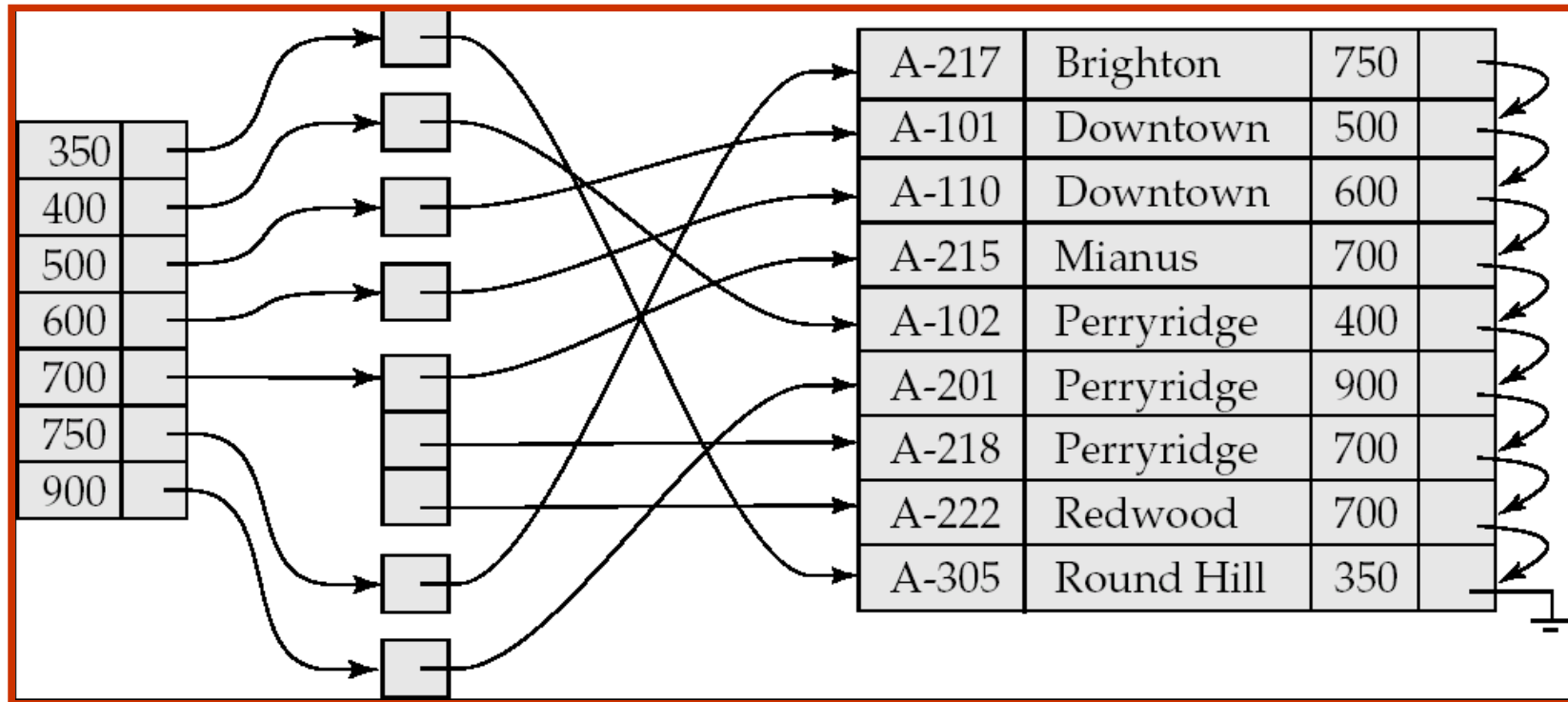| Brighton | | A-217 | Brighton | 750 | |
|----------|--|-------|----------|-----|--|
| Mianus | | A-101 | Downtown | 500 | |
| Redwood | | A-110 | Downtown | 600 | |
| | | A-215 | Mianus | 700 | |
| | | A-102 | Perryridge | 400 | |
| | | A-201 | Perryridge | 900 | |
| | | A-218 | Perryridge | 700 | |
| | | A-222 | Redwood | 700 | |
| | | A-305 | Round Hill | 350 | |

# Index Update:  Insertion

- Single-level index insertion:
  - Perform a lookup using the search-key value appearing in the record to be inserted.

  - **Dense indices** – if the search-key value does not appear in the index, insert it.

  - **Sparse indices** – if index stores an entry for each block of the file, no change needs to be made to the index unless a new block is created.
    - If a new block is created, the first search-key value appearing in the new block is inserted into the index.

- Multilevel insertion (as well as deletion) algorithms are simple extensions of the single-level algorithms

# Secondary Indices

- Frequently, one wants to find all the records whose values in a certain field (which is not the search-key of the primary index) satisfy some condition.
    - Example 1: In the *account* relation stored sequentially by account number, we may want to find all accounts in a particular branch
    - Example 2: as above, but where we want to find all accounts with a specified balance or range of balances
- We can have a secondary index with an index record for each search-key value

# Secondary Indices Example



**Secondary index on *balance* field of *account***

- Index record points to a bucket that contains pointers to all the actual records with that particular search-key value.

- Secondary indices have to be dense

# Primary and Secondary Indices

- Indices offer substantial benefits when searching for records.
- BUT: Updating indices imposes overhead on database modification --when a file is modified, every index on the file must be updated,
- Sequential scan using primary index is efficient, but a sequential scan using a secondary index is expensive
  - Each record access may fetch a new block from disk
  - Block fetch requires about 5 to 10 milliseconds
    - versus about 100 nanoseconds for memory access

# Next...

- B+-trees
- Hashing


- Have a good break!