# CS 443
# Database Management Systems

Sayyed Nezhadi

Pashootan Vaezipoor

https://piazza.com/utoronto.ca/winter2018/csc443/home

# What Is a DBMS?

- A *Database* is a very large, integrated collection of data.

- Models real-world *enterprise.*
    - Entities (e.g., students, courses)
    - Relationships (e.g., Madonna is taking CS564)

- A *Database Management System (DBMS)* is a software package designed to store and manage databases.

# Files vs. DBMS

- Application must stage large datasets between main memory & secondary storage (e.g., buffering, page-oriented access, etc.)
- Special code for different queries
- Must protect data from inconsistency due to multiple concurrent users
- Crash recovery
- Security and access control

# Why Use a DBMS?

- Data independence and efficient access
- Reduced application development time
- Data integrity and security
- Uniform data administration
- Concurrent access, recovery from crashes.

# Why Study Databases??

- Big Data: massive datasets
  - Streaming Data
  - Parallelism (load balancing)
  - Data Analysis, Data Mining: wide-scale distribution
  - Distributed databases
- Complex: complex datatypes and their associated lookups
  - complex base types: geographic data, multimedia, scientific data, CAD data
  - complex objects
  - extensible query processing engines
  - indexing new data types
- Old problems: the data integration problem
  - schema integration: trying to figure out how different schemas fit together. Hard!!! (data cleaning)
  - DBMS integration: trying to semi-transparently glue different kinds of database systems together
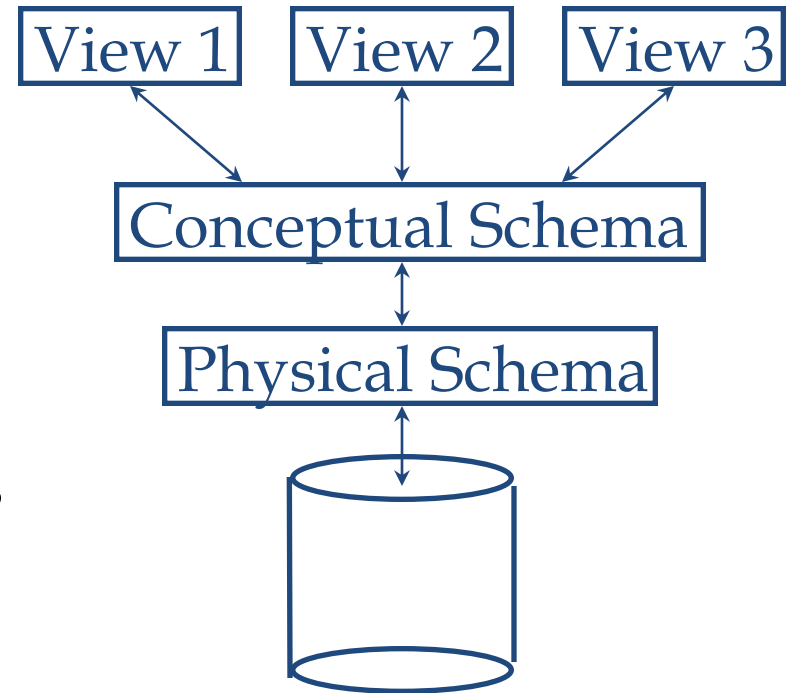- Major Conferences: SIGMOD , VLDB and KDD

# Data Models

- A *data model* is a collection of abstract concepts for:
  - describing the elements of data, their relation and properties *(schema)*
  - how data can be organized, stored and manipulated *(logical data model)*.

- Logical data models:
  - Hierarchical database model
  - Network model
  - Relational model
  - Object model
  - Document model
  - Star schema
  - …

- Physical data models:
  - Inverted Index
  - Flat Files

# Data Models

- A *schema* is a description of a particular collection of data.

- The *relational model of data* is the most widely used model today.

  - Main concept: *relation*, basically a table with rows and columns.

  - Every relation has a *schema*, which describes the columns, or fields.

# Levels of Abstraction

- Many *views*, single *conceptual (logical) schema* and *physical schema*.
  - Views describe how users see the data.
  - Conceptual schema defines logical structure
  - Physical schema describes the files and indexes used.

☞ *Schemas are defined using DDL; data is modified/queried using DML.*

# Example: University Database

- Example Conceptual schema:
  - *Students(sid: string, name: string, login: string,*
    *age: integer, gpa:real)*
  - *Courses(cid: string, cname:string, credits:integer)*
  - *Enrolled(sid:string, cid:string, grade:string)*
- Example Physical schema:
  - Relations stored as unordered files.
  - Index on first column of Students.
- Example External Schema (View):
  - *Course_info(cid:string,enrollment:integer)*

# Data Independence

- Applications insulated from how data is structured and stored.

- *Logical data independence*:  Protection from changes in *logical* structure of data.

- *Physical data independence*:   Protection from changes in *physical* structure of data.

☛ *One of the most important benefits of using a DBMS!*

# Concurrency Control

- Concurrent execution of user programs is essential for good DBMS performance.
  - Because disk accesses are frequent, and relatively slow, it is important to keep the cpu humming by working on several user programs concurrently.
- Interleaving actions of different user programs can lead to inconsistency
  - check is cleared while account balance is being computed
- DBMS ensures such problems don't arise:  users can pretend they are using a single-user system.

# Transaction: An Execution of a DB Program

- Key concept is *transaction,* which is an *atomic* sequence of database actions (reads/writes).
- Each transaction, executed completely, must leave the DB in a *consistent state* if DB is consistent when the transaction begins.
  - Users can specify *integrity constraints* on the data, and the DBMS will enforce these constraints.
  - Beyond this, the DBMS does not really understand the semantics of the data. (e.g., it does not understand how the interest on a bank account is computed).
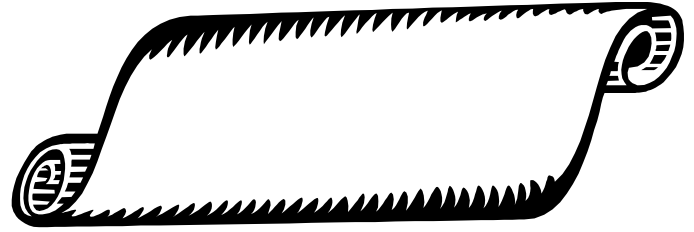
# Scheduling Concurrent Transactions

- DBMS ensures that execution of $\{T_1, \ldots, T_n\}$ is equivalent to some _serial_ execution $T_1' \ldots T_n'$.
  - Before reading/writing an object, a transaction requests a lock on the object, and waits till the DBMS gives it the lock. All locks are released at the end of the transaction. (Strict 2PL locking protocol.)
  - Idea: If an action of $T_i$ (say, writing X) affects $T_j$ (which perhaps reads X), one of them, say $T_i$, will obtain the lock on X first and $T_j$ is forced to wait until $T_i$ completes; this effectively orders the transactions.
  - What if $T_j$ already has a lock on Y and $T_i$ later requests a lock on Y? (Deadlock!) $T_i$ or $T_j$ is aborted and restarted!

# Ensuring Atomicity

- DBMS ensures *atomicity* (all-or-nothing property) even if system crashes in the middle of a Xact.
- Idea: Keep a *log* (history) of all actions carried out by the DBMS while executing a set of Xacts:
  - Before a change is made to the database, the corresponding log entry is forced to a safe location.  (*WAL protocol*; OS support for this is often inadequate.)
  - Write Ahead Log (WAL), if log entry wasn't saved before the crash, corresponding change was not applied to database!
- After a crash, the effects of partially executed transactions are *undone* using the log.
  - Write Ahead Log (WAL), if log entry wasn't saved before the crash, corresponding change was not applied to database!
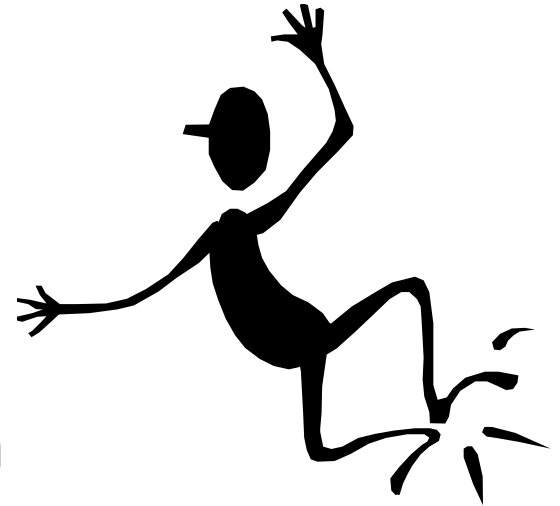
# The Log

- The following actions are recorded in the log:
  - *Ti writes an object*:  The old value and the new value.
    - Log record must go to disk *before* the changed page!
  - *Ti commits/aborts*:  A log record indicating this action.
- Log records chained together by Xact id, so it's easy to undo a specific Xact (e.g., to resolve a deadlock).
- Log is often *duplexed* and *archived* on "stable" storage.
- All log related activities (and in fact, all activities such as lock/unlock, dealing with deadlocks etc.) are handled transparently by the DBMS.

# Databases make these folks happy ...

- End users and DBMS vendors
- DB application programmers
- *Database administrator (DBA)*
  - Designs logical /physical schemas
  - Handles security and authorization
  - Data availability, crash recovery
  - Database tuning as needs evolve

  *Must understand how a DBMS works!*

# Structure of a DBMS

**These layers must consider concurrency control and recovery**

- A typical DBMS has a layered architecture.

- The figure does not show the concurrency control and recovery components.

- This is one of several possible architectures; each system has its own variations.

| Query Optimization and Execution |
|---|
| Relational Operators |
| Files and Access Methods |
| Buffer Management |
| Disk Space Management |

DB

# Summary

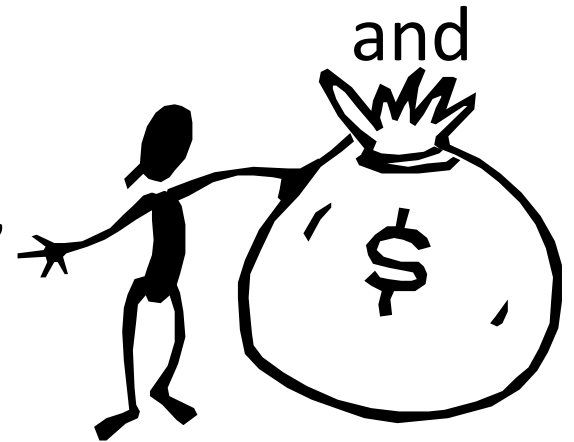- DBMS used to maintain, query large datasets.
- Benefits include recovery from system crashes, concurrent access, quick application development, data integrity and security.
- Levels of abstraction give data independence.
- A DBMS typically has a layered architecture.
- DBAs hold responsible jobs                    and are well-paid! ☺
- DBMS R&D is one of the broadest, most exciting areas in CS.

# Overview of Storage and Indexing

Chapter 8

# Data on External Storage

- Disks: Can retrieve random page at fixed cost
  - But reading several consecutive pages is much cheaper than reading them in random order
- Tapes: Can only read pages in sequence
  - Cheaper than disks; used for archival storage
- File organization: Method of arranging a file of records on external storage.
  - Record id (rid) is sufficient to physically locate record
  - Indexes are data structures that allow us to find the record ids of records with given values in index search key fields
- Architecture: Buffer manager stages pages from external storage to main memory buffer pool. File and index layers make calls to the buffer manager.

# Alternative File Organizations

Many alternatives exist, *each ideal for some situations, and not so good in others:*

- Heap (random order) files:  Suitable when typical access is a file scan retrieving all records.

- Sorted Files:  Best if records must be retrieved in some order, or only a `range' of records is needed.

- Indexes: Data structures to organize records via trees or hashing.

  - Like sorted files, they speed up searches for a subset of records, based on values in certain ("search key") fields
  - Updates are much faster than in sorted files.

# 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**.
  - Given data entry k*, we can find record with key k in at most one disk I/O. (Details soon …)

# B+ Tree Indexes



**Non-leaf Pages**

**Leaf Pages (Sorted by search key)**

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

**index entry**

| $P_0$ | $K_1$ | $P_1$ | $K_2$ | $P_2$ | $\diamond$  $\diamond$  $\diamond$ | $K_m$ | $P_m$ |
|---|---|---|---|---|---|---|---|

# Example B+ Tree



**Root**

**17**

Note how data entries in leaf level are sorted

Entries <= 17          Entries > 17

| 5 | 13 | | | |

| 27 | 30 | | | |

| 2* | 3* | | | | 5* | 7* | 8* | | | 14* | 16* | | | | 22* | 24* | | | 27* | 29* | | | | 33* | 34* | 38* | 39* |

- Find 28*? 29*? All > 15* and < 30*
- Insert/delete:  Find data entry in leaf, then change it. Need to adjust parent sometimes.
  – And change sometimes bubbles up the tree

# Hash-Based Indexes

- Good for equality selections.
- Index is a collection of *buckets.*
  - Bucket = *primary* page plus zero or more *overflow* pages.
  - Buckets contain data entries.
- *Hashing function* **h**: **h**(*r*) = bucket in which (data entry for) record *r* belongs. **h** looks at the *search key* fields of *r.*
  - *No need for "index entries" in this scheme.*

# B+Tree VS Hashing



Index entries

(Direct search for data entries)

Index file

Data entries

Data records

Data

Data file

h(age)=00

Smith, 44, 3000
Jones, 40, 6003
Tracy, 44, 5004

age

h1

h(age) = 01

Ashby, 25, 3000
Basu, 33, 4003
Bristow,29,2007

h(age)=10

Cass, 50, 5004
Daniels, 22, 6003

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

- In a data entry k\* we can store:
  - Data record with key value **k,** or
  - <**k**, rid of data record with search key value **k**>, or
  - <**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

- *Primary* vs. *secondary*:  If search key contains primary key, then called primary index.
  - *Unique* index:  Search key contains a candidate key.
- *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!

# Cost Model for Our Analysis

We ignore CPU costs, for simplicity:

- **B:** The number of data pages

- **R:** Number of records per page

- **D:** (Average) time to read or write disk page

- Measuring number of page I/O's ignores gains of pre-fetching a sequence of pages; thus, even I/O cost is only approximated.

- Average-case analysis; based on several simplistic assumptions.

☛ *Good enough to show the overall trends!*

# Comparing File Organizations

- Heap files (random order; insert at eof)

- Sorted files, sorted on *<age, sal>*

- Clustered B+ tree file, Alternative (1), search key *<age, sal>*

- Heap file with unclustered B + tree index on search key *<age, sal>*

- Heap file with unclustered hash index on search key *<age, sal>*

# Operations to Compare

- Scan: Fetch all records from disk
- Equality search
- Range selection
- Insert a record
- Delete a record

# Assumptions in Our Analysis

- Heap Files:
  - Equality selection on key; exactly one match.
- Indexes:
  - Alt (2), (3): data entry size = 10% size of record
  - Hash: No overflow buckets.
    - 80% page occupancy => File size = 1.25 data size
  - Tree: 67% occupancy (this is typical).
    - Implies file size =  1.5 data size

# Assumptions (contd.)

- Scans:
  - Leaf levels of a tree-index are chained.
  - Index data-entries plus actual file scanned for unclustered indexes.

- Range searches:
  - We use tree indexes to restrict the set of data records fetched, but ignore hash indexes.

# Cost of Operations

|  | (a) Scan | (b) Equality | (c) Range | (d) Insert | (e) Delete |
|---|---|---|---|---|---|
| (1) Heap |  |  |  |  |  |
| (2) Sorted |  |  |  |  |  |
| (3) Clustered |  |  |  |  |  |
| (4) Unclustered Tree index |  |  |  |  |  |
| (5) Unclustered Hash index |  |  |  |  |  |

☛ *Several assumptions underlie these (rough) estimates!*

# Cost of Operations

| | (a) Scan | (b) Equality | (c) Range | (d) Insert | (e) Delete |
|---|---|---|---|---|---|
| (1) Heap | BD | 0.5BD | BD | 2D | Search +D |
| (2) Sorted | BD | $D\log_2 B$ | $D(\log_2 B +$ # pgs with match recs) | Search + BD | Search +BD |
| (3) Clustered | 1.5BD | $D\log_F 1.5B$ | $D(\log_F 1.5B$ + # pgs w. match recs) | Search + D | Search +D |
| (4) Unclust. Tree index | BD(R+0.15) | $D(1 + \log_F 0.15B)$ | $D(\log_F 0.15B$ + # pgs w. match recs) | Search + 2D | Search + 2D |
| (5) Unclust. Hash index | BD(R+0.125) | 2D | BD | Search + 2D | Search + 2D |

☛ *Several assumptions underlie these (rough) estimates!*

# 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.

# Examples of Clustered Indexes

- B+ tree index on E.age can be used to get qualifying tuples.
  - How selective is the condition?
  - Is the index clustered?
- Consider the GROUP BY query.
  - If many tuples have *E.age* > 10, using *E.age* index and sorting the retrieved tuples may be costly.
  - Clustered *E.dno* index may be better!
- Equality queries and duplicates:
  - Clustering on *E.hobby* helps!

```
SELECT  E.dno
FROM  Emp E
WHERE  E.age>40
```

```
SELECT  E.dno,  COUNT (*)
FROM  Emp E
WHERE  E.age>10
GROUP BY E.dno
```
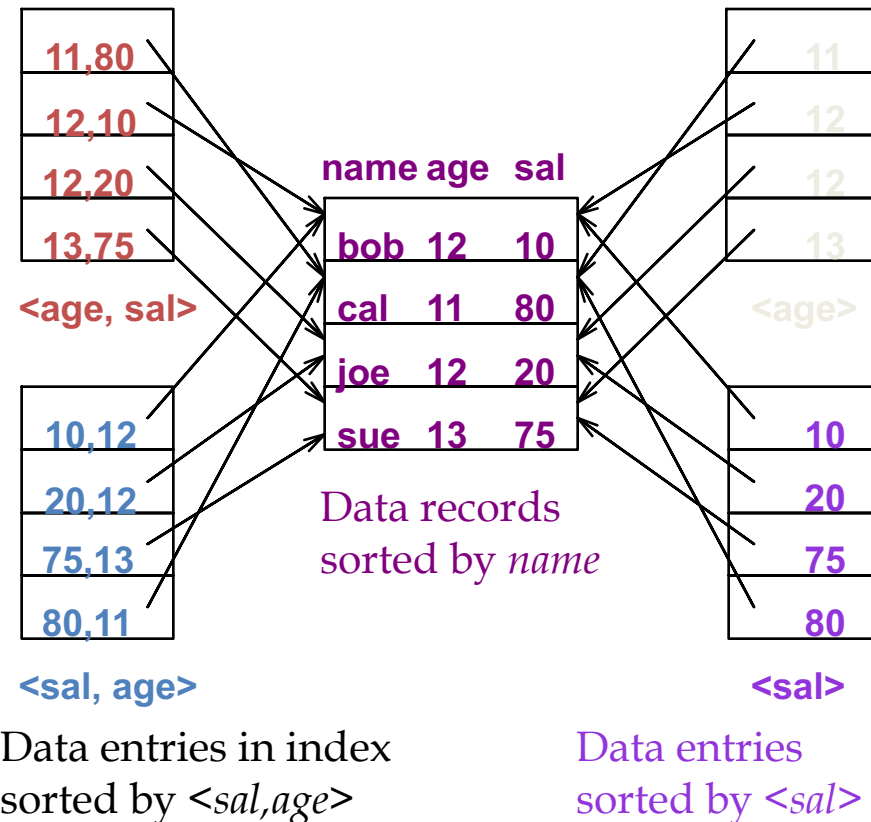
```
SELECT  E.dno
FROM  Emp E
WHERE  E.hobby=Stamps
```

# Indexes with Composite Search Keys

- *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 | | 11 |
| 12,10 | | 12 |
| 12,20 | | 12 |
| 13,75 | | 13 |

**<age, sal>**    **<age>**

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

Data records sorted by *name*

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

**<sal, age>**    **<sal>**

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

Data entries sorted by *<sal>*

# Composite Search Keys

- To retrieve Emp records with *age*=30 AND *sal*=4000, an index on *<age,sal>* would be better than an index on *age* or an index on *sal*.
  - Choice of index key orthogonal to clustering etc.
- If condition is: 20<*age*<30 AND 3000<*sal*<5000:
  - Clustered tree index on *<age,sal>* or *<sal,age>* is best.
- If condition is: *age*=30 AND 3000<*sal*<5000:
  - Clustered *<age,sal>* index much better than *<sal,age>* index!
- Composite indexes are larger, updated more often.

# Index-Only Plans

- A number of queries can be answered without retrieving any tuples from one or more of the relations involved if a suitable index is available.

*<E.dno>*

```
SELECT  E.dno, COUNT(*)
FROM  Emp E
GROUP BY  E.dno
```

*<E.dno,E.sal>*

*Tree index!*

```
SELECT  E.dno, MIN(E.sal)
FROM  Emp E
GROUP BY  E.dno
```

*<E. age,E.sal>*
or
*<E.sal, E.age>*

*Tree index!*

```
SELECT AVG(E.sal)
FROM  Emp E
WHERE  E.age=25 AND
  E.sal BETWEEN 3000 AND 5000
```

# Index-Only Plans (Contd.)

- Index-only plans are possible if the key is <dno,age> or we have a tree index with key <age,dno>
  - Which is better?
  - What if we consider the second query?

```
SELECT  E.dno,  COUNT (*)
FROM  Emp E
WHERE  E.age=30
GROUP BY E.dno
```

```
SELECT  E.dno,  COUNT (*)
FROM  Emp E
WHERE  E.age>30
GROUP BY E.dno
```

# Index-Only Plans (Contd.)

- Index-only plans can also be found for queries involving more than one table; more on this later.

# Summary

- Many alternative file organizations exist, each appropriate in some situation.
- 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.

# Summary (Contd.)

- Data entries can be actual data records, <key, rid> pairs, or <key, rid-list> pairs.
  - Choice orthogonal to *indexing technique* used to locate data entries with a given key value.
- 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. Differences have important consequences for utility/performance.

# Summary (Contd.)

- Understanding the nature of the *workload* for the application, and the performance goals, is essential to developing a good design.
  - What are the important queries and updates?  What attributes/relations are involved?
- Indexes must be chosen to speed up important queries (and perhaps some updates!).
  - Index maintenance overhead on updates to key fields.
  - Choose indexes that can help many queries, if possible.
  - Build indexes to support index-only strategies.
  - Clustering is an important decision; only one index on a given relation can be clustered!
  - Order of fields in composite index key can be important.