# Data Wrangling and Data Analysis
# Advanced SQL
# DB Design & Indexing for Improved Performance

**Hakim Qahtan**

Slides prepared by

**Yannis Velegrakis**

Department of Information and Computing Sciences

Utrecht University

Utrecht University

# Relational Database (Revisted)

relation = table

❖ *Relational database:* a set of *relations*

❖ *Relation:* made up of 2 parts:
- *Instance* : a *table,* with rows and columns.
  #Rows = *cardinality,* #fields = *degree / arity.*
- *Schema* : specifies name of relation, plus name and type of each column.
  - E.G. Students(*sid*: string, *name*: string, *login*: string,                 *age*: integer, *gpa*: real).

❖ Can think of a relation as a *set* of rows or *tuples* (i.e., all rows are distinct).

every tuple should be referred to as unique - SELECT still finds duplicated tho; for instance if we were looking for a name and then there's two students with same name. DISTINCT for unique values

Utrecht University

# Example Instance of Students Relation

| s i d | n a m e | l o g i n | a g e | g p a |
|---|---|---|---|---|
| 5 3 6 6 6 | J o n e s | j o n e s @ c s | 1 8 | 3 . 4 |
| 5 3 6 8 8 | S m i t h | s m i t h @ e e c s | 1 8 | 3 . 2 |
| 5 3 6 5 0 | S m i t h | s m i t h @ m a t h | 1 9 | 3 . 8 |

❖ Cardinality = 3, degree = 5, all rows distinct

❖ Do all columns in a relation instance have to be distinct?   YES

# Creating Relations in SQL

❖ Creates the Students relation.
  ❖ Observe that the type (domain) of each field is specified, and enforced by the DBMS whenever tuples are added or modified.

between the () we put columns, data type,

❖ As another example, the Enrolled table holds information about courses that students take.

relation between two tables studentID and course

CREATE TABLE Students
(sid: CHAR(20),
name: CHAR(20),
login: CHAR(10),
age: INTEGER,
gpa: REAL)

CREATE TABLE Enrolled
(sid: CHAR(20),
cid: CHAR(20),
grade: CHAR(2))

Utrecht University

# Destroying and Altering Relations

DROP TABLE  Students

> won't actually be deleted but name will be changed

❖ Destroys the relation Students.  The schema information *and* the tuples are deleted.

ALTER TABLE  Students
        ADD COLUMN firstYear: integer

> adding column with no values means that all values will be NULL

❖ The schema of Students is altered by adding a new field; every tuple in the current instance is extended with a *null* value in the new field.

Utrecht University

# Adding and Deleting Tuples

❖Can insert a single tuple using:

INSERT INTO Students

VALUES  (53688, 'Smith' , 'smith@ee' , 18, 3.2)

what is the difference?

INSERT INTO Students (name, sid, login, age, gpa)

VALUES  ('Smith' , 53688, 'smith@ee' , 18, 3.2)

❖ Can delete all tuples satisfying some condition (e.g., name = Smith):

DELETE

FROM Students S

AS

WHERE S.name =  'Smith'

delete all records where the name is Smith

Utrecht University

# Integrity Constraints (ICs)

❖ **IC:** condition that must be true for *any* instance of the database; e.g., *domain constraints.*
  ▪ ICs are specified when schema is defined. constraints are defined with the schema
  ▪ ICs are checked when relations are modified.

❖ A *legal* instance of a relation is one that satisfies all specified ICs.
  ▪ DBMS should not allow illegal instances.

❖ If the DBMS checks ICs, stored data is more faithful to real-world meaning.
  ▪ Avoids data entry errors, too!

Utrecht University

# Primary Key Constraints

first constraint that is used;
all the value(s) for the primary key
have to be unique. It can only exist
in one record like that

- A set of fields is a *key* for a relation if :
    1. No two distinct tuples can have same values in all key fields, and
    2. This is not true for any subset of the key.
        - Part 2 false? A *superkey*.
        - If there's >1 key for a relation, one of the keys is chosen (by DBA) to be the *primary key*.

- E.g., *sid* is a key for Students. (What about *name*?) The set {*sid, gpa*} is a superkey.

name is not a primary key

can identify uniquely

values have to be distinct
AND there's no smaller
amount of attributes can be
used to identify. Primary key
is minimum number of
attributes that can identify
each record

Utrecht University

# Primary and Candidate Keys in SQL

❖ Possibly many *candidate keys*  (specified using UNIQUE),
  one of which is chosen as the *primary key*.

❖ "For a given student and course, there is a single grade." vs. "Students can take only one course, and receive a single grade for that course; further, no two students in a course receive the same grade."

❖ Used carelessly, an IC can prevent the storage of database instances that arise in practice!

CREATE TABLE Enrolled
  (sid CHAR(20)
    cid  CHAR(20),
    grade CHAR(2),
    PRIMARY KEY  (sid,cid) )

one has to be very careful with choosing a primary key (constrain) otherwise one would keep values from being added that are correct.

CREATE TABLE Enrolled
  (sid CHAR(20)
    cid  CHAR(20),
    grade CHAR(2),
    PRIMARY KEY  (sid),
    UNIQUE (cid, grade) )

that wouldn't allow two students to have the same grade

Utrecht University

# Foreign Keys, Referential Integrity

- *Foreign key* : Set of fields in one relation that is used to `refer' to a tuple in another relation.  (Must correspond to primary key of the second relation.)  Like a `logical pointer'.

- E.g. *sid* is a foreign key referring to Students:
  - Enrolled(*sid*: string, *cid*: string, *grade*: string)
  - If all foreign key constraints are enforced,  *referential integrity* is achieved, i.e., no dangling references.

    you cannot find a student in the enrolled table that isn't in the student table

  - Can you name a data model w/o referential integrity?
    - Links in HTML!

foreign key will identify a record in another table - e.g. student ID refers to student and can be used to identify student in the course table

Utrecht University

# Foreign Keys in SQL

❖ Only students listed in the Students relation should be allowed to enroll for courses.

CREATE TABLE Enrolled
(esid CHAR(20), cid CHAR(20), grade CHAR(2),
PRIMARY KEY (esid,cid),
FOREIGN KEY (esid) REFERENCES Students )

esid = sid
if one found a
value in enrolled
for esid that isn't
in students table
= dangling

Enrolled

| esid | cid | grade |
|------|------------|-------|
| 53666 | Carnatic101 | C |
| 53666 | Reggae203 | B |
| 53650 | Topology112 | A |
| 53666 | History105 | B |

Students

| sid | name | login | age | gpa |
|-----|------|-------|-----|-----|

Utrecht University

# Enforcing Referential Integrity

- Consider Students and Enrolled; *sid* in Enrolled is a foreign key that references Students.

- What should be done if an Enrolled tuple with a non-existent student id is inserted? *(Reject it!)*

- What should be done if a Students tuple is deleted?
  - Also delete all Enrolled tuples that refer to it.
  - Disallow deletion of a Students tuple that is referred to.
  - Set sid in Enrolled tuples that refer to it to a *default sid*.
  - (In SQL, also: Set sid in Enrolled tuples that refer to it to a special value *null,* denoting `unknown´ or `inapplicable´.)

- Similar if primary key of Students tuple is updated.

Utrecht University

# Referential Integrity in SQL

❖SQL/92 and SQL:1999 support all 4 options on deletes and updates.

- Default is NO ACTION (*delete/update is rejected*)

- CASCADE  (also delete all tuples that refer to deleted tuple)

- SET NULL / SET DEFAULT  (sets foreign key value of referencing tuple)

if you try to add a record where the reference doesn't exist it will delete it by default

CREATE TABLE Enrolled
  (sid CHAR(20),
   cid CHAR(20),
   grade CHAR(2),
   PRIMARY KEY  (sid,cid),
   FOREIGN KEY (sid)
     REFERENCES Students
       ON DELETE CASCADE
       ON UPDATE SET DEFAULT )

ON DELETE CASCADE
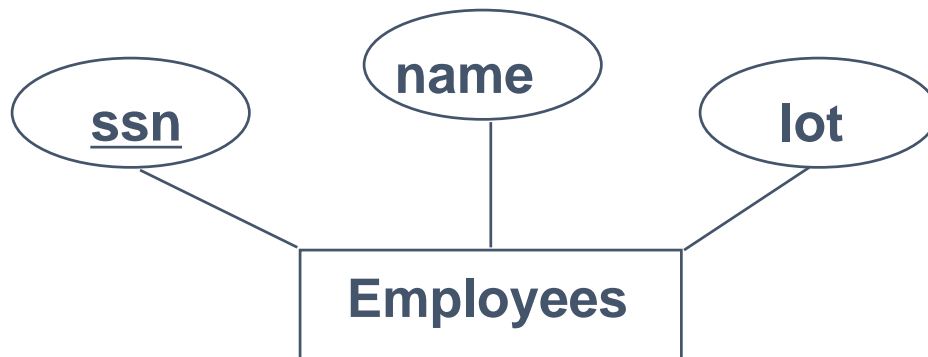= when we delete a student don't allow the id in enrolled - so when update don't use it

# Where do ICs Come From?

❖ ICs are based upon the semantics of the real-world enterprise that is being described in the database relations.

❖ We can check a database instance to see if an IC is violated, but we can NEVER infer that an IC is true by looking at an instance.

- An IC is a statement about *all possible* instances!
- From example, we know *name* is not a key, but the assertion that *sid* is a key is given to us.

I cannot say if the rule is correct

❖ Key and foreign key ICs are the most common; more general ICs supported too.

# Logical DB Design: Entities as Tables



CREATE TABLE Employees
    (ssn CHAR(11),
    name CHAR(20),
    lot  INTEGER,
    PRIMARY KEY  (ssn))

character means text length - VAR character means that it is up to 20 usually. if name is less than 20 it will till accept it

Utrecht University

# Logical DB Design: Relationships as Tables

❖ A relationship between two entities is expressed as a table (stand alone) that is linked to the entities through foreign keys.

❖ That table has a Foreign Key to each key of the tables of the entities to which it is linked, plus some additional descriptive attributes.

CREATE TABLE Works_In(
  ssn  CHAR(11),
  did  INTEGER,
  since  DATE,
  PRIMARY KEY (ssn, did),
  FOREIGN KEY (ssn)
      REFERENCES Employees,
  FOREIGN KEY (did)
      REFERENCES Departments)

employee working at specific department - so relation represents relationship between employee and department primary key for new table (ssn,did) Foreign key from employee (ssn) and Foreign key from department (did)

# Examples of Relationships as Tables ... sort of.

```
CREATE TABLE Manages(
    ssn  CHAR(11),
    did  INTEGER,
    since  DATE,
    PRIMARY KEY  (did),
    FOREIGN KEY (ssn) REFERENCES Employees,
    FOREIGN KEY (did) REFERENCES Departments)
```

❖ Since each department has a unique manager, we could instead combine Manages and Departments.

```
CREATE TABLE Dept_Mgr(
    did  INTEGER,
    dname  CHAR(20),
    budget  REAL,
    ssn  CHAR(11),
    since  DATE,
    PRIMARY KEY  (did),
    FOREIGN KEY (ssn) REFERENCES Employees)
```
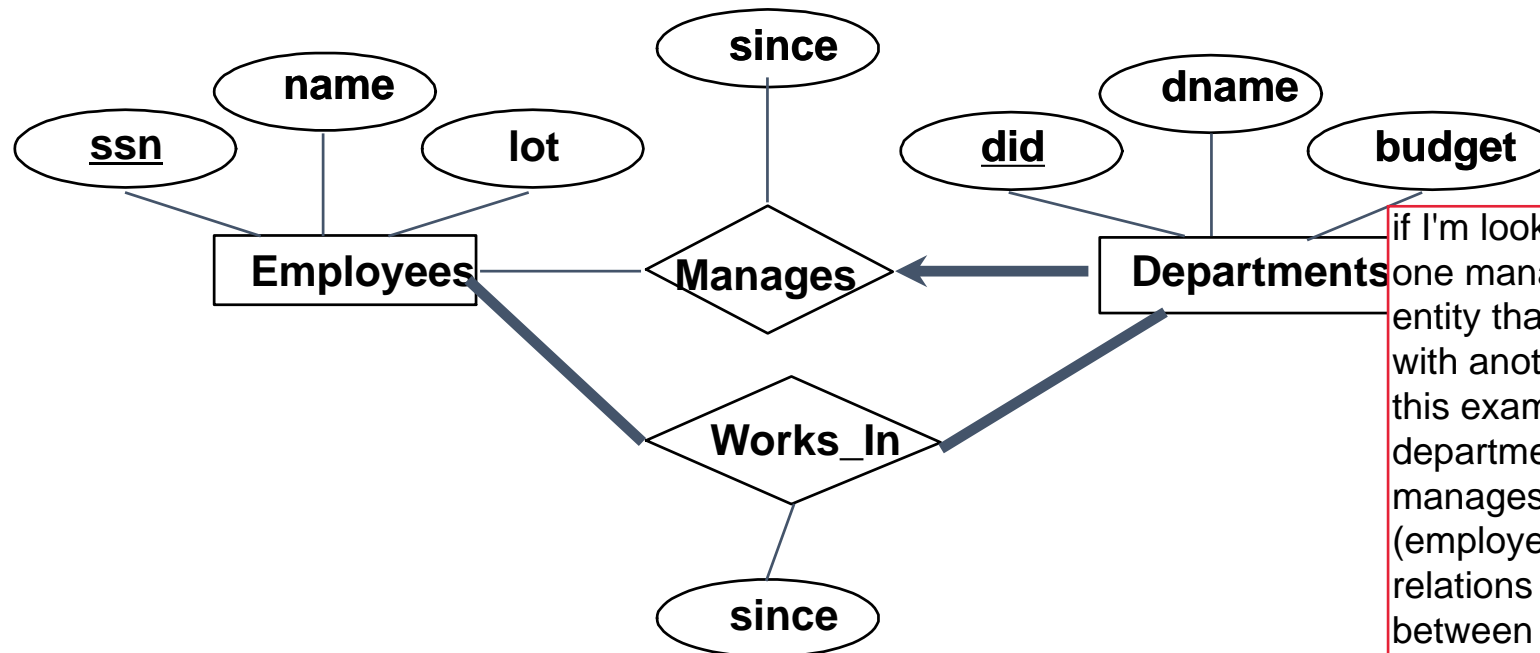
combining the two tables instead of two

Utrecht University

# Review: Participation Constraints

❖Does every department have a manager?

- If so, this is a *participation constraint*:  the participation of Departments in Manages is said to be *total* (vs. *partial*).

  - Every *did* value in Departments table must appear in a row of the Manages table (with a non-null *ssn* value!)



if I'm looking at this department: - one manager but maybe there's a entity that doesn't have a relation with another in the other relation. In this example every did from department must appear in manages table. Two entity relations (employee and department) and two relations that define relationships between the tables (entities).

Utrecht University

# Participation Constraints in SQL

❖ We can capture participation constraints involving one entity set in a binary relationship, but little else (without resorting to CHECK constraints).

```
CREATE TABLE  Dept_Mgr(
   did  INTEGER,
   dname  CHAR(20),
   budget  REAL,
   ssn  CHAR(11) NOT NULL,       every department needs a manager - no department without manager
   since  DATE,                   (ssn)
   PRIMARY KEY  (did),
   FOREIGN KEY  (ssn) REFERENCES Employees,
     ON DELETE NO ACTION)        if I delete a manger I don't want the whole thing deleted
```

Utrecht University

# Another ISA Hierarchy example



Utrecht University

# Translating ISA Hierarchies to Relations

❖ *General approach:*

- 3 relations: Employees, Hourly_Emps and Contract_Emps.

  - *Hourly_Emps*:  Every employee is recorded in Employees.  For hourly emps, extra info recorded in Hourly_Emps (*hourly_wages*, *hours_worked*, <u>*ssn*</u>); must delete Hourly_Emps tuple if referenced Employees tuple is deleted).

  - Queries involving all employees easy, those involving just Hourly_Emps require a join to get some attributes.

    need more attributes from hourly employees table

❖ Alternative:  Just Hourly_Emps and Contract_Emps.

- *Hourly_Emps*:  <u>*ssn*</u>, *name, lot, hourly_wages, hours_worked.*

- Each employee must be in one of these two subclasses.

Utrecht University

# Views

❖A *view* is just a relation, but we store a *definition*, rather than a set of tuples.

CREATE VIEW YoungActiveStudents (name, grade)
  AS SELECT S.name, E.grade
  FROM Students S, Enrolled E
  WHERE S.sid = E.sid and S.age<21

returns a set of record. not information will be stored but stored as definition

❖ Views can be dropped using the DROP VIEW command.

▪ How to handle DROP TABLE if there's a view on the table?

• DROP TABLE command has options to let the user specify this.

Utrecht University

# Views and Security

- Views can be used to present necessary information (or a summary), while hiding details in underlying relation(s).
  - Given YoungStudents, but not Students or Enrolled, we can find students s who have are enrolled, but not the *cid's* of the courses they are enrolled in.

Utrecht University

# Relational Model: Summary

❖ A tabular representation of data.

❖ Simple and intuitive, currently the most widely used.

❖ Integrity constraints can be specified by the DBA, based on application semantics. DBMS checks for violations.   two main constraints:

- Two important ICs: primary and foreign keys
- In addition, we *always* have domain constraints.

# The NULL

- Null is a special value

- Comparison to any other value **always false**

- which is why we have the special operator ISNULL

# Designing Good Databases

**FullInfo**

| id | name | carId | brand | color |
|---|---|---|---|---|
| A | John | 2 | VW | black |
| B | Nick | 3 | VW | Red |
| C | Mary | null | null | null |
| A | John | 3 | VW | Red |

**Person**

| id | name | carId |
|----|------|-------|
| A | John | 2 |
| B | Nick | 3 |
| C | Mary | null |
| A | John | 3 |

**Car**

| carId | brand | color |
|-------|-------|-------|
| 2 | VW | black |
| 3 | VW | Red |

FULL OUTER JOIN

**FullInfo**

| id | name | carId | brand | color |
|----|------|-------|-------|-------|
| A | John | 2 | VW | black |
| B | Nick | 3 | VW | Red |
| C | Mary | null | null | null |
| A | John | 3 | VW | Red |

Which one do you like?

**Person**

| id | name | carId |
|---|---|---|
| A | John | 2 |
| B | Nick | 3 |
| C | Mary | 2 |
| A | John | 3 |

**Car**

| carId | brand | color |
|---|---|---|
| 2 | VW | black |
| 3 | VW | Red |

PERSON JOIN CAR

**FullInfo**

| id | name | carId | brand | color |
|---|---|---|---|---|
| A | John | 2 | VW | black |
| B | Nick | 3 | VW | Red |
| C | Mary | 2 | VW | black |
| A | John | 3 | VW | Red |

projection to get a new table with less data (specific)

Projection on id, name, carId

**Person**

| id | name | carId |
|---|---|---|
| A | John | 2 |
| B | Nick | 3 |
| C | Mary | 2 |
| A | John | 3 |

Projection on carId, brand, color

**Car**

| carId | brand | color |
|---|---|---|
| 2 | VW | black |
| 3 | VW | Red |

Utrecht University

**FullInfo**

| id | name | carId | brand | color |
|---|---|---|---|---|
| A | John | 2 | VW | black |
| B | Nick | 3 | VW | Red |
| C | Mary | 2 | VW | black |
| A | John | 3 | VW | Red |

Projection on id, name, carId

**Person**

| id | name | carId |
|---|---|---|
| A | John | 2 |
| B | Nick | 3 |
| C | Mary | 2 |
| A | John | 3 |

Projection on carId, brand, color

**Car**

| carId | brand | color |
|---|---|---|
| 2 | VW | black |
| 3 | VW | Red |

PERSON JOIN CAR

**FullInfo**

| id | name | carId | brand | color |
|---|---|---|---|---|
| A | John | 2 | VW | black |
| B | Nick | 3 | VW | Red |
| C | Mary | 2 | VW | black |
| A | John | 3 | VW | Red |

**FullInfo**

| id | name | carId | brand | color |
|---|---|---|---|---|
| A | John | 2 | VW | black |
| B | Nick | 3 | VW | Red |
| C | Mary | 2 | VW | black |
| A | John | 3 | VW | Red |

Projection on id, name, carId, brand

Projection on brand, color

**Person**

| id | name | carId | brand |
|---|---|---|---|
| A | John | 2 | VW |
| B | Nick | 3 | VW |
| C | Mary | 2 | VW |
| A | John | 3 | VW |

**Car**

| brand | color |
|---|---|
| VW | black |
| VW | Red |

missing another identifying column
so that every VW can be black or red -
carID

PERSON JOIN CAR

**FullInfo**

| id | name | carId | brand | color |
|---|---|---|---|---|
| A | John | 2 | VW | black |
| B | Nick | 3 | VW | black |
| C | Mary | 2 | VW | black |
| A | John | 3 | VW | black |
| A | John | 2 | VW | Red |
| B | Nick | 3 | VW | Red |
| C | Mary | 2 | VW | Red |
| A | John | 3 | VW | Red |

**There is clearly
a problem here !!**

# Dependencies

- How do we know when a design is good? Follow some basic rules
  - 1 table for each entity
  - 1 table for each relationship

  > like the entities person and car - one table each and one table for each relationship

- Decompose Big Relations to Normal Forms
  - A dependency is an expression of the form:
    - Office → Department, or
    - Passport → Nationality

    > if I have two columns in a table there can be dependencies between the columns- I can't find a passport with two nationalities etc.

  - BCNF
    - For every functional dependency in a relation
      - Either it is trivial
      - Or the left part is a key

  - Dependency preserving
    - For each dependency, all the attributes are in the same relation
      - *This is not the only case but it would suffice for this course*

# Overview of Storage and Indexing

# Data on External Storage

- <u>Disks:</u> Can retrieve random page at fixed cost
  - But reading several consecutive pages is much cheaper than reading them in random order
- <u>Tapes (Rare today):</u> Can only read pages in sequence
  - Cheaper than disks; used for archival storage
- <u>File organization:</u> 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
- <u>Architecture:</u> Buffer manager stages pages from external storage to main memory buffer pool. File and index layers make calls to the buffer manager.

Utrecht University

# 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   *if I always have the samw query like "looking for student ID' it works fast*
  - Updates are much faster than in sorted files.
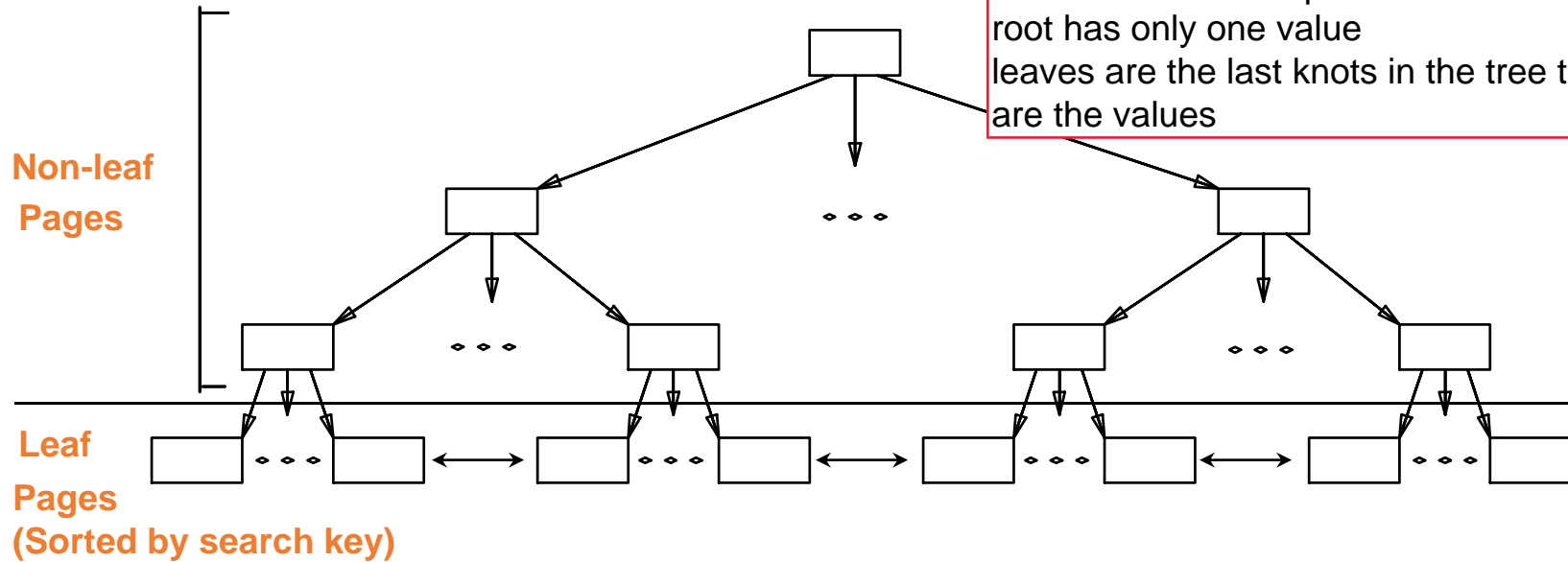
Utrecht University

# 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.**

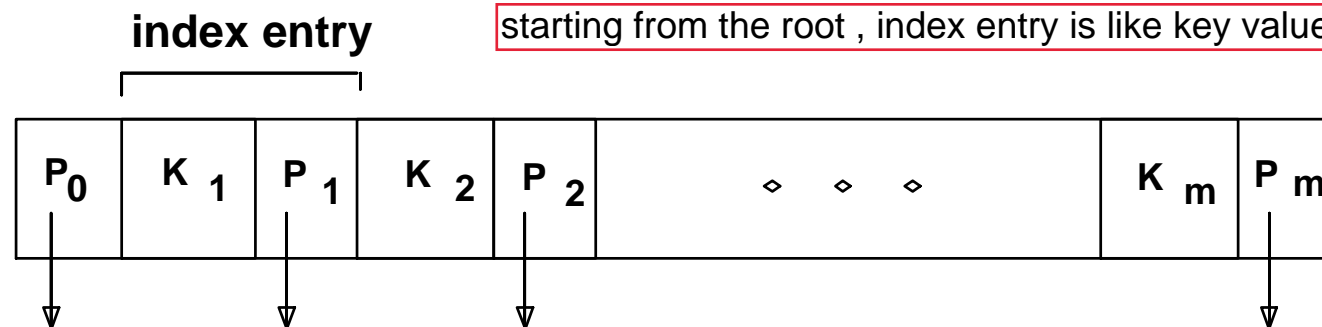  - Given data entry k*, we can find record with key k in at most one disk I/O. (Details soon …)

# B+ Tree Indexes

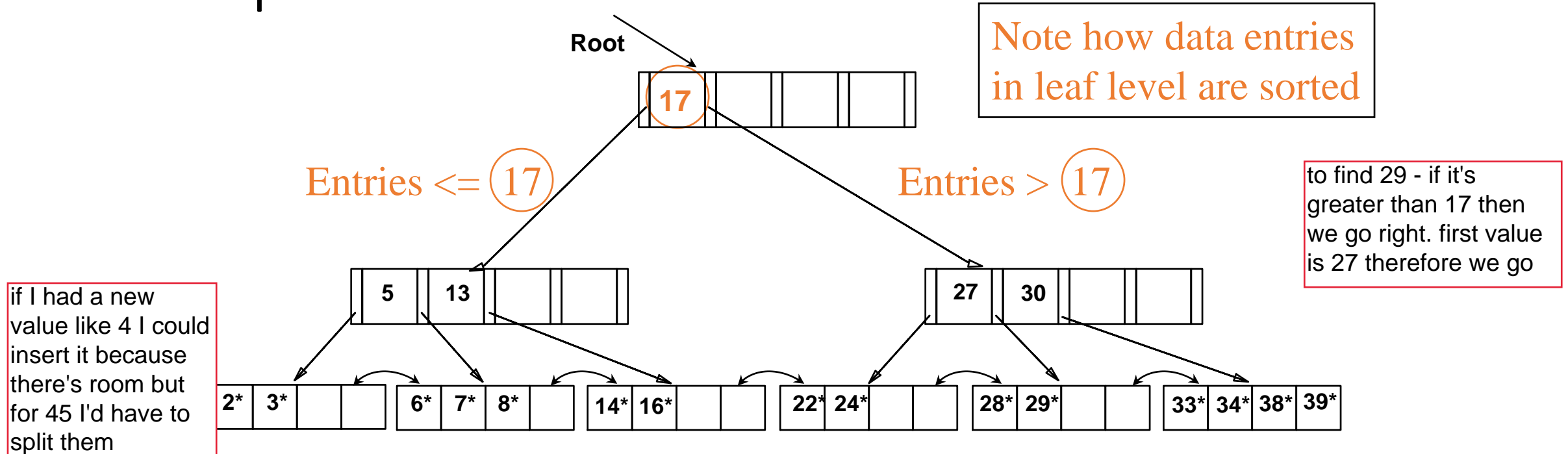**Non-leaf Pages**

**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**

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

Utrecht University

# Example B+ Tree

**Root**

17

Entries <= (17)

Entries > (17)

to find 29 - if it's greater than 17 then we go right. first value is 27 therefore we go

5   13

27   30

if I had a new value like 4 I could insert it because there's room but for 45 I'd have to split them

2*  3*

6*  7*  8*

14*  16*

22*  24*

28*  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

Utrecht University

# Hash-Based Indexes

data is stored in a binary representation - hash functions converts from decimal representation into binary representation
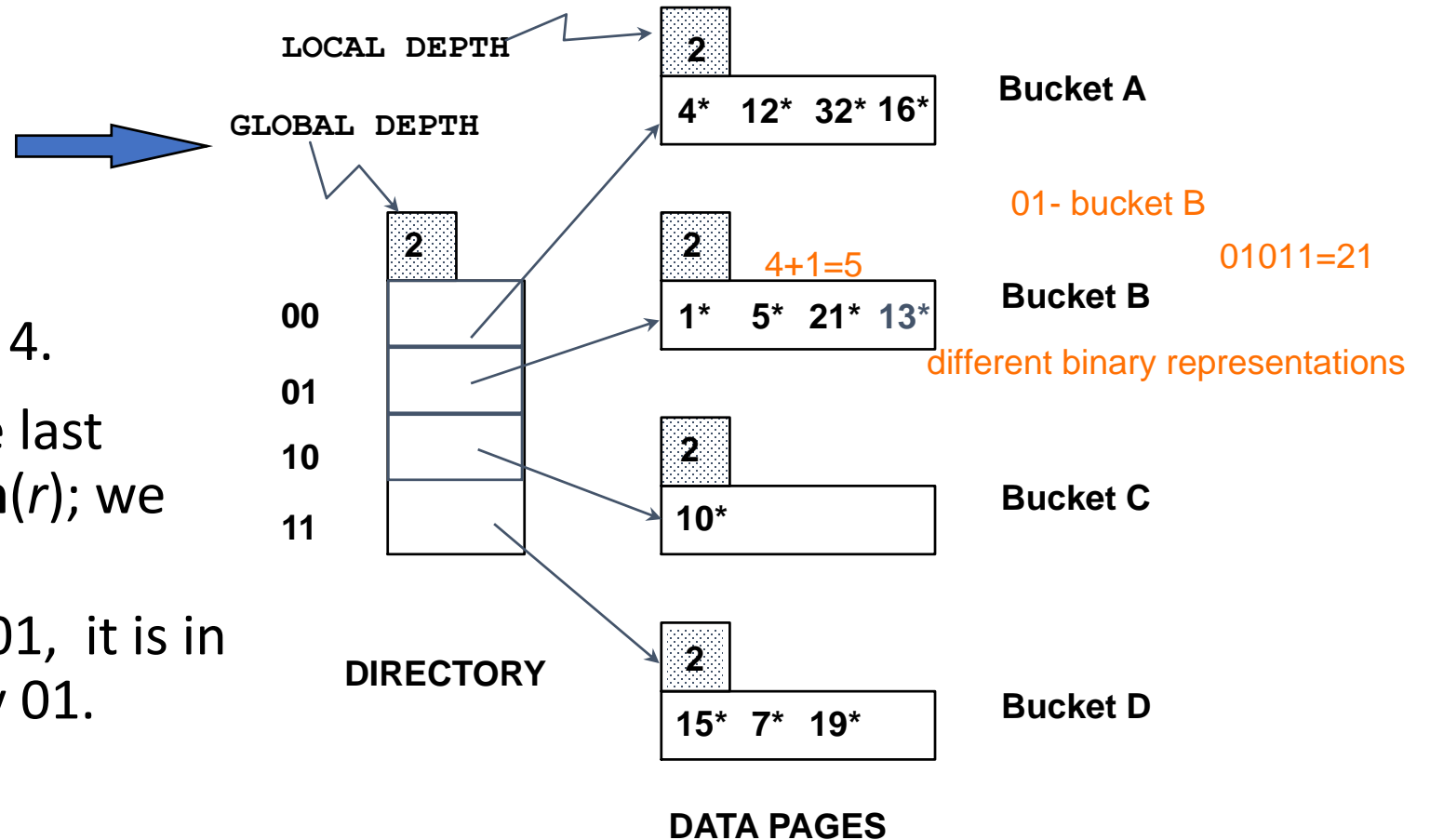
- Good for equality selections.

- 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.*

Utrecht University

binary representation of the values
e.g. Bucket B

# Example

- Directory is array of size 4.
- To find bucket for *r*, take last `*global depth*` # bits of **h**(*r*); we denote *r* by **h**(*r*).
  - If **h**(*r*) = 5 = binary 101, it is in bucket pointed to by 01.

LOCAL DEPTH

GLOBAL DEPTH

**2**

**2**

| 4* | 12* | 32* | 16* |

**Bucket A**

01- bucket B

01011=21

**2**    4+1=5

| 1* | 5* | 21* | 13* |

**Bucket B**

different binary representations

00

01

10

11

**2**

| 10* |

**Bucket C**

**2**

| 15* | 7* | 19* |

**Bucket D**

DIRECTORY

DATA PAGES

❖ **<u>Insert</u>**:  If bucket is full, *<u>split</u>* it (*allocate new page, re-distribute*).

❖ *If necessary,* double the directory.  (As we will see, splitting a bucket does not always require doubling; we can tell by comparing *global depth* with *local depth* for the split bucket.)

if I need to split the buckets I also need to increase the number of buckets in the entry
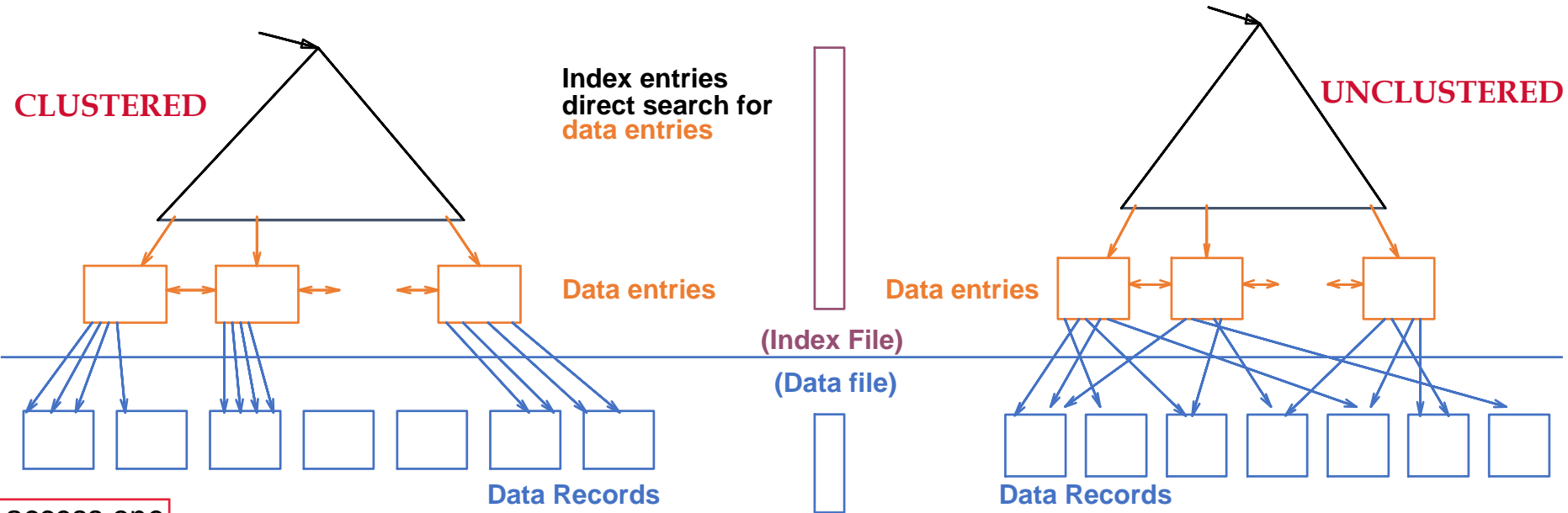
# Index Classification

- *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.
  - 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!

Utrecht University

# Clustered vs. Unclustered Index

means putting similar values together
one query will access a set of records



CLUSTERED

Index entries
direct search for
data entries

Data entries

(Index File)

(Data file)

Data Records

UNCLUSTERED

Data entries

Data Records

if we access one
knot we can go
to a specific set
of values/knots

Utrecht University

# 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

expecting almost all values from the table - then grouped afterwards

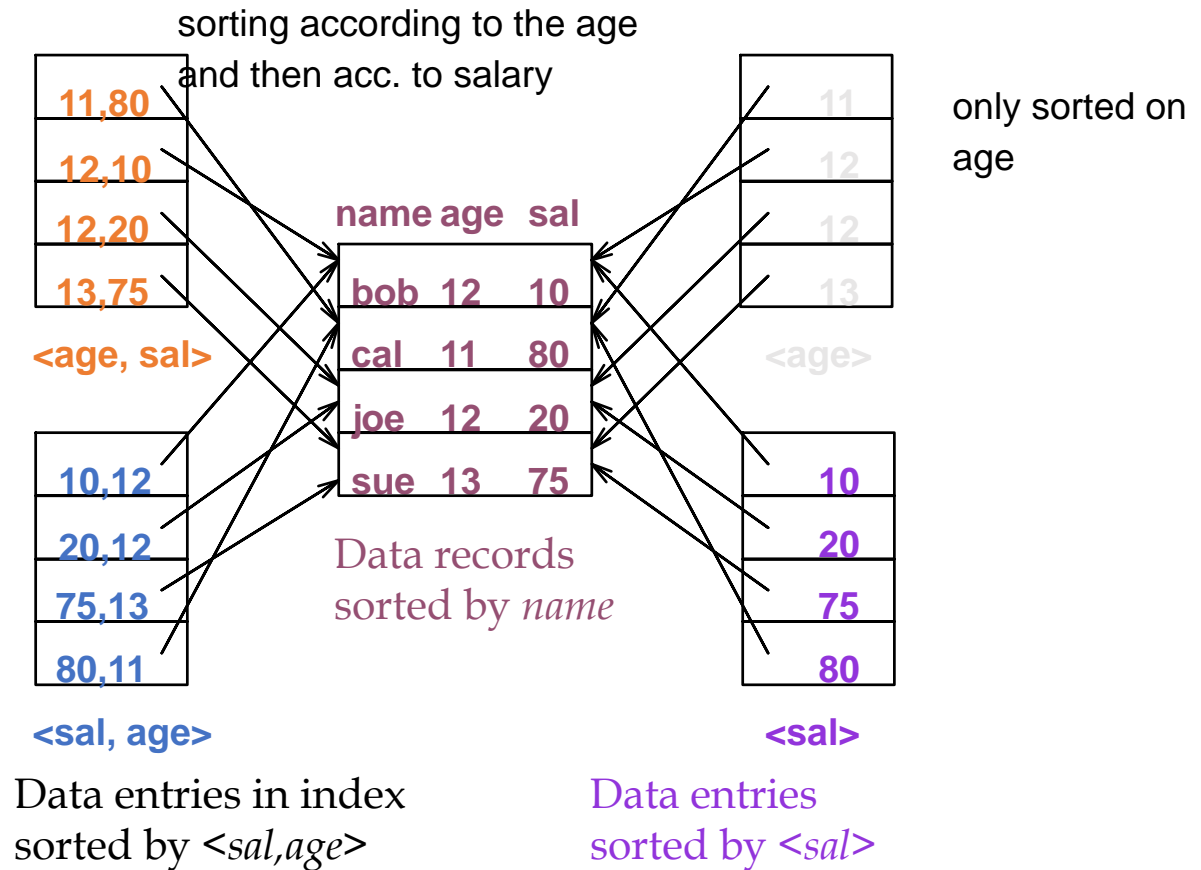SELECT  E.dno
FROM  Emp E
WHERE  E.hobby=Stamps

Utrecht University

# 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.

  search keys from multiple attributes
  more attributes will be used in the index

Examples of composite key indexes using lexicographic order.

sorting according to the age and then acc. to salary



only sorted on age

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

**<age, sal>**

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

11
12
12
13

<age>

Data records sorted by *name*

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

<sal, age>

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

10
20
75
80

**<sal>**

Data entries sorted by <*sal*>

# Composite Indexes for Single Attribute Search

composite because it has TWO attributes

- An index on *<age,sal>* can be used to retrieve records with *age*=30  (or age>10) but is NOT of any help for retrieving records with sal=20 (or sal>30)

Utrecht University

# 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 wouldn't be helpful
clustering is the way to go when
range on both queries

so age first!
age is fixed and then
salary between 3k and 5k
otherwise you'd get a lot
more records

Utrecht University

# 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

looking in employee relation and trying to find the number

index on employee department number and employee salary

minimum salary

can be answered by only accessing the index for that relation

*<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

Utrecht University

# 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  equality query
GROUP BY E.dno

index starting with equality range

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

Utrecht University

# 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?

can be decide

Utrecht University

# 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.

# Operations to Compare

- Scan: Fetch all records from disk

- Equality search

- Range selection

- Insert a record

- Delete a record

Utrecht University

# 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.

Utrecht University

# 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.

# A trick for finding the max

- To find the max of a set (e.g., best student in class)
  - It is enough to find all those that are not best and remove them from the set that contains everything.
  - Those that remain are by definition … the best

always try to eliminate the options that are impossible

Utrecht University

# Suggested References for Further Reading

You can read the related sections from:

Database Management Systems, 3rd Edition

by Raghu Ramakrishnan, Johannes Gehrke

Utrecht University