

# INTRODUCTION

A quick disclaimer before reading the notes:  
they were taken by me during the lectures,  
they do not replace the professor's work and  
are not sufficient for passing the exams.

Moreover they might contain mistakes, so  
please double check all that you read. The  
notes are freely readable and can be shared  
(always remembering to credit me and to not  
obscure this page), but **can't** be modified.

Thank you and hope these notes are useful!

-Francesca Cinelli



DMAV1

## Introduction To databases and database management Systems

↳ collection of related data items within a specific business process or problem setting

↳ software package used to define, create, use, and maintain a database

↳ database system

→ database approach: all data stored, managed centrally by DBMS

→ raw data + metadata

↳ data definitions that are stored in the catalog of DBMS

### Database model / schema

- description of DB data at different levels of detail and specifies the various data items, relationships etc.
- specified during design and doesn't change frequently
- stored in catalog (heart of DBMS)

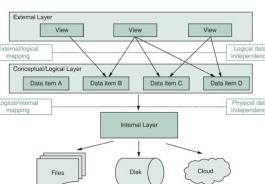
### Database state / instance

- data in DB at particular moment
- changes

### Models

- conceptual model : high-level description of data items with their characteristics and relationships
- logical model: translation / mapping of the conceptual model toward a specific implementation environment
  - ↳ relational
- internal model: mapping of ↑ that represents the data's physical storage details
- external model: subsets of data items in logical model, views, tailored toward the needs of specific applications or groups of users

### 3-layer architecture



### Catalog

- contains data definitions / metadata of DB app.
- stores definitions of views, logical and internal data models and synchronizes these three data models to ensure their consistency

### Database languages

- data definition language (DDL) express external, logical and internal schemas } SQL
- data manipulation language (DML) retrieve, insert, delete, modify data

### Information

- structured: objects represented by short strings of symbols and numbers
- unstructured: texts written in a natural language
- semi-structured: have a structure but irregular/relative

### Information system

a set of data physically organized in secondary memory and managed in such a way as to allow its creation, updating and interrogation in homogeneous aggregates

## Data independence

- physical - from logical and external
- logical - from external

## Advantages of DB

- no redundancy
- integrity - correctness of data, constraints have to be satisfied
- security - protect data from unauthorized access
- concurrency control

L transaction = sequence of operations constituting a single logical transaction

L ACID

- atomicity: all-or-nothing property
- consistency: DB from one consistent state to another
- isolation: effect of concurrent transactions should be the same as if they had been executed in i.
- durability: the DB changes made by a transaction declared successful can be made permanent in all circumstances

→ backup and recovery - transaction log or dump

## Performance indicators

- response time: time between a request and its termination
- throughput rate: transactions processed per unit of time
- space utilization: space used to store raw data + metadata

## Relational model

- data + relationships = values
- no explicit pointers

object = record / tuple  
attribute = info of interest

table = set of records of homogeneous type  
 $r = \text{instances} + R = \text{schemes}$  (To write qries)

## Definitions

- domain: a set of values (columns of table)
- cartesian product:  $D_1 \times D_2 \times \dots \times D_k$  the set  $\{(v_1, v_2, \dots, v_k) \mid v_1 \in D_1, v_2 \in D_2, \dots, v_k \in D_k\}$
- mathematical relation: subset of  $\uparrow$  of  $k$  domains (relation of degree  $k$ )
  - elements = tuples /ennuples
  - # tuples = cardinality
  - $t[i]$  =  $i^{\text{th}}$  component of tuple  $t$
- attribute = name + domain (set of attributes = schema)
- NULL values = lack of info or info not applicable, polymorphic value (doesn't belong to any domain but can replace values in any domain), 2 null values are different, null ≠ 0
- intra-relational constraints = defined on single attribute values or between values of same tuple or between tuples of the same relation
- inter-relational constraints = defined between multiple relations
- key = attribute / set of attributes that uniquely identifies a tuple → if  $t_1[x] \neq t_2[x]$
- primary key = most used key / one with smaller # of attributes → can never be NULL
- foreign key = it is p.key in other relation
- functional dependencies = semantic link between 2 sets of attributes  $x$  and  $y$  belonging to  $R$  ( $x \rightarrow y$ ) →  $t_1[x] = t_2[x] \rightarrow t_1[y] = t_2[y]$

## Relational algebra 1

- > notation for specifying queries about the contents of relations
- > formal language to interrogate a relational DB, consisting of a set of unary and binary operators that if applied to one or two relation instances, generate a new relation instance
- > procedural language - describes exact order in which operators have to be applied to relational instance to obtain the desired result
- > 4 types of operators (single, set, combine, renaming) that can also be combined with each other

## Operators

### 1. Projection (single)

- vertical cut on a relation, selects a subset of attributes
- $\pi_{A_1, A_2, \dots, A_k}(R)$
- select the columns of R corresponding to attributes  $A_1, A_2, \dots, A_k$
- ex.  $\pi_{\text{name}}(\text{Customer}) \rightarrow \begin{array}{|c|} \hline \text{Rossi} \\ \text{Bianchi} \\ \text{Verdi} \\ \hline \end{array}$
- no duplicates in result (might be loss of tuples)

### 2. Selection (single)

- horizontal cut on the relation, selects all rows that meet some constraints
- $\sigma_C(R)$
- select the tuples of R that meet condition C
  - ↳ Boolean expression whose terms are in the form  $A \theta B$  or  $A \theta a'$
  - ↳ item in domain of A
  - ↳ comparison operator
  - ↳ attributes same domain
- no loss of tuples

### 3. Union (set)

- creates a new relation instance containing all tuples belonging to at least one of the operand instances
- $R_1 \cup R_2$
- performed on union compatible operands only
  - ↳ have same # attributes, corresponding attributes have same domain
- can use projection first to make them compatible to union
- is commutative

### 4. Difference (set)

- creates a relation containing the tuples in the first operand relation but not in the second
- $R_1 - R_2$
- applied only to union compatible operands
- not commutative

### 5. Intersection (set)

- creates a relation containing the items that are present in both the operand relations
- $R_1 \cap R_2 = (R_1 - (R_1 - R_2))$
- applied only to union compatible operands
- is commutative

## Relational algebra 2

### 6. Cartesian product (combine)

- creates a relation with tuples obtained by combining all the tuples in the first relation with all the tuples in the second relation
- $r_1 \times r_2$
- used when info needed is contained in multiple relations
- associative and commutative

### 7. Natural join (combine)

- selects the tuples in the result of the Cartesian product that satisfy the condition:

$$* R_1 A_1 = R_2 A_1 \wedge R_1 A_2 = R_2 A_2 \wedge \dots \wedge R_1 A_k = R_2 A_k$$

where  $A_1, A_2, \dots, A_k \rightarrow$  attributes in common, so duplicate attributes are dropped

$$\cdot r_1 \bowtie r_2 = \pi_{xy} (\sigma_c (r_1 \times r_2))$$

attribute  $r_1$       attribute  $r_2$  not in  $r_1$

- special cases:

→ 2 relations have some attributes in common but the attributes have no values in common

result = empty set

→ 2 relations do not have attributes with the same name so condition \* cannot be evaluated

result = becomes a Cartesian product

- associative and commutative

### 8. $\theta$ -join (combine)

- selects the tuples resulting from the Cartesian product and satisfying the following condition:

$$A \theta B$$

where  $\theta$  is a comparison operator ( $\theta \in \{ \leq, =, \geq, \leq, \geq \}$ )

$$\cdot r_1 \bowtie_{\theta B} r_2 = \theta_{A \theta B} (r_1 \times r_2)$$

### 9. Renaming

- creates a copy of the relation in which an attribute is renamed

$$\cdot R_2 = \rho_{A \rightarrow B} (R_1)$$

## Relational algebra 3

→ join of a table with itself can be done

→ examples on slides 05

## Designing a relational database - problems and constraints

- redundancy

- waste of memory space } happens when we have a database with a single schema and ≠ concepts  
- anomalies

→ update anomaly → if some information changes we have to change that info everywhere

→ insertion anomaly → cannot insert some data if it's not related to other data in Table

→ deletion anomaly → if delete some data other important correlated data will be cancelled

## Solution:

A database schema is good if it has no redundancies and no anomalies.

To achieve this we need to correlate in the same schema only things that represent the same concept, so each concept has to be represented separately in a distinct relation.

We have to make a database in which a single "table" can be achieved by join operations between multiple tables.

How to identify the concepts represented in a relation? By the key

→ an attribute/group of attributes that determine a particular functional dependence which is a particular kind of constraint

## Constraints

↳ the representation in the database schema of a condition that is valid in the reality of interest

↳ an instance of the database is legal if it satisfies all constraints ("honest" representation of reality)

↳ DBMS allows us to define the constraints together with the DB schema, to verify if an instance is legal, to prevent insertion of tuples that would violate constraints

↳ functional dependencies express constraints between subsets of attributes of the schema itself

## Functional dependencies

### notation

- $R = A_1, A_2, \dots, A_n$
- $(A, B, C, \dots)$  single attributes
- $(X, Y, \dots)$  sets of attributes
- $X \cup Y = X \cup Y$
- tuple = function that associates to each attribute  $A_i$  in  $R$  a value  $t_i[A_i]$  in the domain
- $t_1[X] = t_2[X]$  if  $\forall A \in X \quad t_1[A] = t_2[A]$  ( $t_1$  and  $t_2$  coincide)
- $t_1[X] \neq t_2[X]$  if  $\exists A \in X \quad t_1[A] \neq t_2[A]$
- instance of  $R$  = set of tuples on  $R$

functional dependence on  $R$  = ordered pair of non-empty subsets  $X, Y \subseteq R$

-  $X \rightarrow Y$  ( $X$  functionally determines  $Y$  or  $Y$  is functionally dependent on  $X$ )

-  $X$  = determinant,  $Y$  = dependent

- an instance  $r$  of  $R$  satisfies the f.d.  $X \rightarrow Y$  if:  $\forall t_1, t_2 \in r: t_1[X] = t_2[X] \Rightarrow t_1[Y] = t_2[Y]$

↳ implies

note: if  $t_1[X] \neq t_2[X]$  the dependence is

$$\begin{aligned} &\checkmark \text{ if } (t_1[X] = t_2[X]) \text{ then } (t_1[Y] = t_2[Y]) \\ &\times \text{ if } (t_1[X] \neq t_2[X]) \text{ then } (t_1[Y] \neq t_2[Y]) \end{aligned}$$

satisfied whatever the values of  $t_1[Y]$  and  $t_2[Y]$

remark:  $A \rightarrow B \equiv \bar{A} \cup B$

if  $F = \{A \rightarrow B, B \rightarrow C\}$  is satisfied then  $A \rightarrow C$  is also satisfied

if  $A \rightarrow BC$  then  $A \rightarrow B$  and  $A \rightarrow C$  BUT if  $AB \rightarrow C$  then not necessarily  $A \rightarrow C$  and  $B \rightarrow C$

• closure of  $F$  ( $F^+$ ) = set of f.d. that are satisfied by each legal instance of  $R$  ( $F \subseteq F^+$ )

• a subset  $k$  of  $R$  is a key of  $R$  if  $k \rightarrow R \in F^+$  and there is no proper subset  $k' \subseteq k$  such that  $k' \rightarrow R$

•  $X \rightarrow Y \in F^+ \Leftrightarrow \forall A \in Y, X \rightarrow A \in F^+$  (Property of f.d.)

How to calculate  $F^+$ ?

## Closure of F

$F^A$  is a set of functional dependencies on R so that:

- if  $X \rightarrow Y \in F$  then  $X \rightarrow Y \in F^A$
- if  $Y \subseteq X \in R$  then  $X \rightarrow Y \in F^A$  (reflexivity)
- if  $X \rightarrow Y \in F^A$  then  $XZ \rightarrow YZ \in F^A$ , for each  $Z \in R$  (augmentation)
- if  $X \rightarrow Y \in F^A$  and  $Y \rightarrow Z \in F^A$  then  $X \rightarrow Z \in F^A$  (transitivity)
- if  $X \rightarrow Y \in F^A$  and  $X \rightarrow Z \in F^A$  then  $X \rightarrow YZ \in F^A$  (union rule)
- if  $X \rightarrow Y \in F^A$  and  $Z \subseteq Y$  then  $X \rightarrow Z \in F^A$  (decomposition rule)
- if  $X \rightarrow Y \in F^A$  and  $WY \rightarrow Z \in F^A$  then  $WX \rightarrow Z \in F^A$  (pseudotransitivity rule)

Armstrong  
axioms

a

b

c

Proofs:

- a if  $X \rightarrow Y \in F^A$  by augmentation  $X \rightarrow XY \in F^A$ , if  $XY \rightarrow Z \in F^A$  by augmentation  $X \rightarrow YZ \in F^A$ , since  $X \rightarrow XY \in F^A$  and  $XY \rightarrow YZ \in F^A$  by transitivity  $X \rightarrow YZ \in F^A$
- b if  $Z \subseteq Y$  by reflexivity  $Y \rightarrow Z \in F^A$ , since  $X \rightarrow Y \in F^A$  and  $Y \rightarrow Z \in F^A$  by transitivity  $X \rightarrow Z \in F^A$
- c if  $X \rightarrow Y \in F^A$  by augmentation  $WX \rightarrow WY \in F^A$ , since  $WX \rightarrow WY \in F^A$  and  $WY \rightarrow Z \in F^A$  by transitivity  $WX \rightarrow Z \in F^A$

Let  $X \subseteq R$ , the closure of  $X$  w.r.t.  $F$  ( $X^+$ ) is defined as  $X_F^+ = \{A \mid X \rightarrow A \in F^A\}$

$$X \subseteq X_F^+$$

Lemma: Let R be a schema and F a set of functional dependencies on R:  $X \rightarrow Y \in F^A \Leftrightarrow Y \subseteq X^+$

Proof: let  $Y = A_1, A_2, \dots, A_n$

if ( $\Leftarrow$ ): since  $Y \subseteq X^+$ , for each  $i=1, \dots, n$  we have that  $X \rightarrow A_i \in F^A$  by the union rule,  $X \rightarrow Y \in F^A$

only if ( $\Rightarrow$ ): since  $X \rightarrow Y \in F^A$ , for each  $i=1, \dots, n$ ,  $X \rightarrow A_i \in F^A$  ( $A_i \in X^+$ ), for each  $i=1, \dots, n$   $Y \subseteq X^+$

## Theorem

$$F^+ = F^A$$

Proof:  $F^A \subseteq F^+$

$\bar{F}_i^A$  = f.d. with at most i applications of Armstrong axioms

base case:  $i=0 \rightarrow \bar{F}_0^A = F \subseteq F^+$

induction step:  $i > 0 \rightarrow \bar{F}_i^A \subseteq F^+ \Rightarrow \bar{F}_{i+1}^A \subseteq F^+ \quad X \rightarrow Y \in \bar{F}_{i+1}^A$

3 cases

- reflexivity  $Y \subseteq X \quad t_1, t_2 \in r \quad r = \text{instance } R$   
 $t_1[x] = t_2[x] \Rightarrow t_1[Y] = t_2[Y] \Rightarrow X \rightarrow Y \in F^+$

- augmentation  $V \rightarrow W \in F_i^A \quad Z \subseteq R \quad X = VZ \text{ and } Y = WZ$   
 $t_1[X] = t_2[X] \rightarrow t_1[VZ] = t_2[VZ] \rightarrow t_1[V] = t_2[V] \wedge t_1[Z] = t_2[Z] \rightarrow$   
 $\rightarrow t_1[W] = t_2[W] \wedge t_1[Z] = t_2[Z] \rightarrow t_1[WZ] = t_2[WZ] \rightarrow t_1[Y] = t_2[Y]$

- transitivity  $X \rightarrow Z, Z \rightarrow Y \in F_i^A \quad X \rightarrow Z, Z \rightarrow Y \in F^+$   
 $t_1[X] = t_2[X] \rightarrow t_1[Z] = t_2[Z] \rightarrow t_1[Y] = t_2[Y]$

Proof:  $F^+ \subseteq F^A$

$$r \quad X^+ \quad R - X^+$$

$$X \rightarrow Y \in F^+$$

$t_1$	1.....1	1.....1
$t_2$	1.....1	0.....0

• prove r is legal

$$V \rightarrow W \in F(\subseteq F^A) \quad t_1[V] = t_2[V] \Rightarrow V \subseteq X^+ \Rightarrow X \rightarrow V \in F^A \Rightarrow X \rightarrow W \in F^A \Rightarrow W \subseteq X^+ \Rightarrow t_1[W] = t_2[W]$$

$$X \subseteq X^+ \Rightarrow t_1[X] = t_2[X] \Rightarrow t_1[Y] = t_2[Y] \Rightarrow Y \subseteq X^+ \Rightarrow X \rightarrow Y \in F^A$$

$$\text{So } F^+ = F^A \rightarrow |F^+| >> 2^{|R|}$$

↳ lemma

## Third Normal Form 1

→ a relational schema is in 3NF if the only non-trivial f.d. that must be satisfied by any legal instance are those of type  $k \rightarrow X$  where  $k$  contains a key or  $X$  is contained in a key  
 → formally

$X \rightarrow Y$  non-trivial if  $Y \notin X$

$X \rightarrow Y \in F^+$

→  $R$  is in 3NF if  $\forall X \rightarrow A \in F^+, A \notin X$ ,  $A$  belongs to a key (is prime) or/and  $X$  contains a key (it is a superkey)

→ to check 3NF we only need to check the dependencies  $X \rightarrow A \in F^+$  obtained by decomposing the dependents of the dependencies in  $F$

ex.  $F = \{A \rightarrow B, B \rightarrow CD\}$  only check  $A \rightarrow B, B \rightarrow C, B \rightarrow D$  but not  $A \rightarrow C, A \rightarrow D$  that are obtained by transitivity  
 if  $A$  is the only key then the schema  $R$  for which  $F$  holds is not in 3NF

→ if the keys contain all dependents then for sure  $R$  will be in 3NF (since all dependents are prime)

## Third Normal Form 2

→ partial dependencies: when a dependency is a consequence of another dependency

$X \rightarrow A \in F^+ | A \notin X$  is a partial dependence on  $R$  if  $A$  is not prime and  $X$  is properly contained in a key of  $R$  ( $X - k = \emptyset$ )

→ transitive dependencies: when a dependency is a consequence of two functional dependencies

$X \rightarrow A \in F^+ | A \notin X$  is a transitive dependence on  $R$  if  $A$  is not prime and for each key  $k$  of  $R$  we have that  $X$  is not properly contained in  $k$  ( $X - k \neq \emptyset$ ) and  $k - X \neq \emptyset$

→ schema is in 3NF iff there are neither partial nor transitive dependencies in  $F$

→ goal: produce a schema in which every relation is in 3NF

→ if we produce schema that is not in 3NF we can decompose it so they are in 3NF

→ if we join decompositions we must produce a legal instance, we must preserve all dependencies in  $F^+$  and not lose any information

## Boyce-Codd Normal Form

a relation is in BCNF when every determinant in it is a superkey (a key is also a superkey)

→ if a relation respects BCNF it is also in 3NF, opposite not true

→ not always possible to decompose a schema in BCNF but always possible to do so with 3NF

## Closure of $X$

· when obtaining schemas in 3NF we have to preserve all dependencies in  $F^+$ , but calculating it takes exponential time. But by the lemma  $X \rightarrow Y \in F^+ \Leftrightarrow Y \subseteq X^+$  we only need to calculate  $X^+$  to know that  $X \rightarrow Y \in F^+$

### Algorithm for calculating $X^+$

input:  $R, F, X \subseteq R$  ( $X$  may be a single attribute)

output:  $X_F^+$

begin

$Z = X$

$S = A | Y \rightarrow V \in F, A \in V \wedge Y \subseteq Z$

while  $S \subsetneq Z$  do:

begin

$Z = Z \cup S$

$S = A | Y \rightarrow V \in F, A \in V \wedge Y \subseteq Z$

end

end

insert into  $S$  the attributes that make up the right-hand parts of dep. in  $F$  whose left-hand part is contained in  $Z$ . At first  $Z$  is just  $X$ , so we insert attributes that are functionally determined by  $X$ ; once these have entered  $Z$ , from these we add more (by transitivity)

Proof Let us denote by  $Z^{(0)}$  the initial value of  $Z$  ( $Z^{(0)} = X$ ) and by  $Z^{(i)}$  and  $S^{(i)}$  the values of  $Z$  and  $S$  after the  $i^{\text{th}}$  execution of the while loop; it is easy to see that  $Z^{(i)} \subseteq Z^{(i+1)}$ , for each  $i$ . Let  $j$  be such that  $S(j) \subseteq Z(j)$  ( $Z(j)$  = value of  $Z$  when algo terminates) we will prove that  $A \in Z_j \Leftrightarrow A \in X^+$

part one:  $Z_j \subseteq X^+$  by induction

base case:  $i=0$

since  $Z^{(0)} = X$  and  $X \subseteq X^+$   $\rightarrow Z^{(0)} \subseteq X^+$

induction step:  $i > 0 \quad Z^{(i+1)} \subseteq X^+$

$A \in Z_{i+1} = Z_i \cup S_i \rightarrow A \in Z_i \subseteq X^+$

$\rightarrow A \in S_i \Rightarrow \exists \gamma \rightarrow V \in F \text{ s.t. } A \in V \wedge \gamma \in Z_i$

$\gamma \subseteq Z_i \subseteq X^+ \Rightarrow X \rightarrow \gamma \in F^A \quad \left. \begin{array}{l} X \rightarrow V \in F \\ V \subseteq X^+ \end{array} \right\} \Rightarrow V \subseteq X^+ \quad \left. \begin{array}{l} A \in V \\ A \in Z_i \end{array} \right\} A \in X^+$

$\Rightarrow Z_j \subseteq X^+$  (since we proved for any  $i$ )

part two:  $X^+ \subseteq Z_j \quad \frac{F^+}{A \in X^+ \Rightarrow X \rightarrow A \in F^A}$

consider

$n$	$Z_i$	$R - Z_i$
$t_1$	1...1	1...1
$t_2$	1...1	0...0

legal?  $\downarrow \quad \gamma \rightarrow V \in F \quad t_1[\gamma] = t_2[\gamma] \Rightarrow \gamma \subseteq Z_j \quad \left. \begin{array}{l} V \subseteq S_j \subseteq Z_j \\ \gamma \rightarrow V \in F \end{array} \right\}$

Yes!

$\Rightarrow t_1[V] = t_2[V]$

$X = Z_0 \subseteq Z_j \Rightarrow t_1[X] = t_2[X] \Rightarrow t_1[A] = t_2[A] \Rightarrow A \in Z_j$  (contradiction)

algorithm is polynomial time

$\Rightarrow Z_j = X^+$

## Finding the keys of a schema

$\rightarrow$  a subset  $K$  of a relation schema  $R$  is a key of  $R$  if:

1.  $K \rightarrow R \in F^+$

2. there is no proper subset  $K'$  of  $K$  such that  $K' \rightarrow R \in F^+$

remark 1  $\rightarrow$  start from those with higher cardinality, if their closure does not contain  $R \rightarrow$  don't check  $C$

remark 2  $\rightarrow$  attributes that never appear on the right of the functional dependencies of  $F$  must be in all the keys

remark 3  $\rightarrow$  brute force approach not best

remark 4  $\rightarrow$  every shortcut in calculations must be justified

$\rightarrow$  if the closure of  $X$  contains  $R$  we have identified a superkey but we always have to check its minimality, meaning that we have to check that no subset of that closure is a key

## Preserving $F$

Decomposition of  $R$  is a set of subsets  $R_1, R_2, \dots, R_k$  of  $R$  covering  $R$  that is  $R = \bigcup_{i=1}^k R_i$

Equivalence, let  $F$  and  $G$  be two sets of f.d.;  $F$  and  $G$  are equivalent ( $F = G$ ) if  $F^+ = G^+$ , how to check?

Lemma: if  $F \subseteq G^+$  then  $F^+ \subseteq G^+$

Proof:  $f \in F^+ \Rightarrow f \in G^+$

1.  $f \in F \subseteq G^+$   $\{f_1, \dots, f_n\} \xrightarrow{G^+} f \in F^+ \quad f_i \in F \subseteq G^+ = G^A$

2.  $f \notin F \rightarrow f \in F^+ - F \quad \{g_1, \dots, g_n\}_{\notin F} \xrightarrow{G^+} f \rightarrow \text{for every } i \xrightarrow{G^+} f \in G^+$

## Preserving

$\rho = R_1, R_2, \dots, R_k$  a decomposition of  $R$

$\rightarrow \rho$  preserves  $F$  if  $F \in G = \bigcup_{i=1}^k \pi_{R_i}(F)$  where  $\pi_{R_i}(F) = \{x \rightarrow y \in F^+ \mid xy \subseteq R_i\}$

check if  $F$  is preserved by checking  $F \subseteq G^+$ :

$$\forall x \rightarrow y \in F \Rightarrow y \subseteq X_G^+ \Rightarrow x \rightarrow y \in G^A \Rightarrow x \rightarrow y \in G^+$$

How to compute  $X_G^+$ ?

**Input:** a schema  $R$ , a set  $F$  of functional dependencies on  $R$ ,

**Output:** the closure of  $X$  with respect to  $G = \bigcup_{i=1}^k \pi_{R_i}(F)$ , in variable  $Z$

```

begin
Z = X
S = ∅
for i = 1 to k
    S = S ∪ (Z ∩ R_i)^+ ∩ R_i
while S ⊈ Z
    begin
        Z = Z ∪ S
        for i = 1 to k
            S = S ∪ (Z ∩ R_i)^+ ∩ R_i
    end
end

```

Theorem:  $Z_f = X_G^+$

Proof:  $Z_f \subseteq X_G^+$  by induction

base case:  $i = 0$

$$Z_0 = X \subseteq X_G^+$$

induction case:  $i > 0 \quad Z_i \subseteq X_G^+ \Rightarrow Z_{i+1} \subseteq X_G^+$

$$A \in Z_{i+1} = Z_i \cup S_i$$

$\cdot A \in Z_i \subseteq X_G^+$  by induction hypothesis

$$\begin{aligned} \cdot A \in S_i = \bigcup_{j=1}^k (Z_i \cap R_j)^+ \cap R_j &\Rightarrow \exists j \text{ s.t. } A \in (Z_i \cap R_j)^+ \cap R_j \Rightarrow A \in R_j \wedge A \in (Z_i \cap R_j)^+ \\ &\Rightarrow (Z_i \cap R_j) \rightarrow A \in F^A = F^+ \Rightarrow \underbrace{Z_i \cap R_j}_{\subseteq Z_i \subseteq X_G^+} \rightarrow A \in \pi_{R_j}(F) \subseteq G \subseteq G^+ \Rightarrow x \rightarrow A \in G^+ \Leftrightarrow A \in X_G^+ \Rightarrow Z_{i+1} \subseteq X_G^+ \end{aligned}$$

Proof:  $X_G^+ \subseteq Z_f \quad |y \in V \Rightarrow y \in X_G^+ \subseteq Z_f| \quad$  (not part of program)

$$X = Z_0 \subseteq Z_f \Rightarrow X_G^+ \subseteq (Z_f)_G^+$$

$$S' = \{A \mid \exists y \rightarrow y \in G \mid A \in V \wedge y \in Z_f\}$$

$$\text{let } A \in S' \Rightarrow \exists j \in \{1 \dots k\} \text{ s.t. } y \rightarrow y \in \pi_{R_j}(F) \Rightarrow S' \subseteq Z_f \subseteq Z_f$$

$$X_G^+ \subseteq Z_f = (Z_f)_G^+$$

## Lossless join

Def: let  $R$  be a relation schema; A decomposition  $\rho = R_1, R_2, \dots, R_k$  has a lossless join if for each legal instance  $r$  of  $R$  we have  $r = \pi_{R_1}(r) \bowtie \pi_{R_2}(r) \bowtie \dots \bowtie \pi_{R_k}(r)$

Thm: let  $R$  be a relation schema and let  $\rho = R_1, R_2, \dots, R_k$  be a decomposition of  $R$ ; for each legal instance  $r$  of  $R$ , denoted by  $m_\rho(\rho) := \pi_{R_1}(r) \bowtie \pi_{R_2}(r) \bowtie \dots \bowtie \pi_{R_k}(r)$  we have:

1.  $r \subseteq m_\rho(r)$
2.  $\pi_{R_i}(m_\rho(r)) = \pi_{R_i}(r)$
3.  $m_\rho(m_\rho(r)) = m_\rho(r)$

Proof:

1.  $t \in r$ . consider  $t[R_i] \in \pi_{R_i}(r) \quad t \in \{t[R_i]\} \bowtie \dots \bowtie \{t[R_k]\} \quad \rightarrow t \in m_\rho(r) \rightarrow r \subseteq m_\rho(r)$   
 $\pi_{R_1}(r) \bowtie \dots \bowtie \pi_{R_k}(r) = m_\rho(r)$

2.  $\rightarrow \pi_{R_i}(r) \subseteq \pi_{R_i}(m_\rho(r))$

$t_i \in \pi_{R_i}(r) \quad \exists t \in r \text{ s.t. } t_i = t[R_i] \Rightarrow t \in m_\rho(r) \Rightarrow t_i \in \pi_{R_i}(m_\rho(r))$

$\rightarrow \pi_{R_i}(m_\rho(r)) \subseteq \pi_{R_i}(r)$

$|t_i \in \pi_{R_i}(m_\rho(r))| \Rightarrow \exists t' \in m_\rho(r) \text{ s.t. } t_i = t'[R_i] \Rightarrow \exists t_1, \dots, t_k \in r \text{ s.t. } t'[R_i] = t_i[R_i] \in \pi_{R_i}(r)$

$t_i = t'[R_i] = t_i[R_i] \in \pi_{R_i}(r)$  contradiction

3.  $m_\rho(m_\rho(r)) = \pi_{R_1}(m_\rho(r)) \bowtie \dots \bowtie \pi_{R_k}(m_\rho(r)) = \pi_{R_1}(r) \bowtie \dots \bowtie \pi_{R_k}(r) = m_\rho(r)$

## Algorithm

begin  
construct a table  $r$  as follows:

$r$  has  $R$  columns and  $k$  rows

at the intersection of the  $i$ -th row and  $j$ -th column put:

the symbol  $a_j$  if the attribute  $A \in R_{i,j}$

the symbol  $b_{ij}$  otherwise

repeat

for every  $X \rightarrow Y \in F$

if there are two tuples  $t_1$  and  $t_2$  in  $r$  such that  $t_1[X] = t_2[X]$  and  $t_1[Y] \neq t_2[Y]$

then for every attribute  $A_j$  in  $Y$ : if  $t_1[A_j] = a_j$  then  $t_2[A_j] = t_1[A_j]$

else  $t_1[A_j] = t_2[A_j]$

until  $r$  has a line with all "a" or  $r$  has not changed;

if  $r$  has a row with all "a" then: it has a lossless join

else: it does not have a lossless join  $\rightarrow$  has a lossy join

end

ex.  $R = ABCD$   $P = \{ABC, ABD\}$   $F = \{AB \rightarrow C; AB \rightarrow D\}$

r	$\overbrace{\quad \quad \quad \quad}^j$			
	A	B	C	D
$R_1 = ABC$	$a_1$	$a_2$	$a_3$	$b_{4,1}^{a_4}$
$R_2 = ABD$	$a_1$	$a_2$	$b_{2,3}^{a_3}$	$a_4$

$\Rightarrow$  has a lossless join.

Thm: let  $R$  be a relation scheme,  $F$  a set of functional dependencies on  $R$  and let  $P = R_1, R_2, \dots, R_k$  be a decomposition of  $R$ ; the algorithm correctly decides whether  $P$  has a lossless join.

$S$  has a lossless join iff algorithm returns an instance  $r$  containing the tuple  $t^a = (a \dots a)$

Proof:

$\Rightarrow$  let  $r^o$  be the initial instance

let  $r^f$  be the final instance

$t_i^o[R_i] = (a, \dots, a) = t_i^f[R_i] \in \Pi_{R_i}(r^f)$

$t^a \in \Pi_{R_1}(r^f) \times \dots \times \Pi_{R_k}(r^f) = P_P(r^f)$  but  $r^f$  is legal  $\Rightarrow t^a \in r^f$   
 $"r^f"$  since assume has lossless join

## Minimal covers

How to obtain a "good" decomposition?

Is it always possible to get it? Yes, sometimes even several

$\rightarrow$  is the input to the decomposition algorithm

$\rightarrow$  given  $F$  there can exist several equivalent minimal covers

Def: a minimal cover of  $F$  is a set of functional dependencies  $G$ , equivalent to  $F$ , such that:

- for each functional dependency in  $G$ , the dependent is a singleton (each dependent is non-redundant)  $\rightarrow$  always possible thanks to decomposition rule
- for each dependency  $X \rightarrow A$  in  $G$ , there not exists  $X' \subset X$  such that  $G \equiv (G - \{X \rightarrow A\}) \cup \{X' \rightarrow A\}$  (each determinant is non-redundant)  $\rightarrow$  not possible to functionally determine  $A$  by a subset of  $X$
- there not exists any  $X \rightarrow A$  in  $G$ , such that  $G \equiv G - \{X \rightarrow A\}$  (each dependency is non-redundant)  
 $\rightarrow$  not possible to functionally determine  $A$  through other dependencies

How calculated?

in polynomial time in 3 steps:

1. using the decomposition rule, the dependents are reduced to singletons (right side of dependency)
2. every functional dependency  $[A_1, A_2, \dots, A_i, A_i, A_{i+1}, \dots, A_n] \rightarrow [A_j]$  in  $F$  such that  $F \equiv F - \{ \cdot \cup \{A_i\} \rightarrow A_j \}$  is replaced by "subset of  $X$ "; if the latter already belongs to  $F$  the original dependency is simply deleted; the process is repeated recursively on "subset of  $X$ "; the process ends when no f.d. can be further reduced

3. every functional dependency  $X \rightarrow A$  in  $F$  such that  $F \equiv F - \{X \rightarrow A\}$  is removed from  $F$  (because redundant)

2 and 3 are equivalence checks

2:  $\cdot \in G^+$  and "subset of  $X$ "  $\in F^+$

$\hookrightarrow$  not necessary

3:  $X \rightarrow A \in G^+$  or  $A \in X_G^+$

## Decomposition algorithm

- given  $R$  and  $F$  which is a minimal cover, compute a decomposition  $\rho = R_1, R_2, \dots, R_k$  such that
  - for each  $i=1, \dots, k$   $R_i$  is in 3NF
  - $\rho$  preserves  $F$
  - $\rho$  has no lossless join

## Algorithm

```

begin
S = ∅
ρ = ∅
for each A ∈ R, such that A is not involved in any functional dependency in F do
    S = S ∪ A
if S ≠ ∅ then:
    R = R - S
    ρ = ρ ∪ {S}
if there is a functional dependency in F that involves all the attributes in R
    then: ρ = ρ ∪ {R}
else:
    for each X → A ∈ F do
        ρ = ρ ∪ {XA}
end

```

Theorem: the decomposition algorithm computes a decomposition  $\rho$  of  $R$  such that: each relational schema in  $\rho$  is in 3NF and  $\rho$  preserves  $F$

## Proof:

- $\rho$  preserves  $F$ 
  - Let  $G = \bigcup_{i=1}^k T_{R_i}(F)$ ; since for each f.d.  $X \rightarrow A \in F$  we have that  $XA \in \rho$ , we have that this dependency of  $F$  will be in  $G$ , hence  $F \subseteq G$  and  $F^+ \subseteq G^+$
- each schema is in 3NF
  1. if  $S \in \rho$ , each attribute in  $S$  is part of the key and so  $S$  is in 3NF
  2. if  $R \in \rho$ , there is a functional dependency in  $F$  involving all attributes of  $R$ ; as  $F$  is a minimal cover, this d. will have the form  $R \cdot A \rightarrow A$ ; as  $F$  is a minimal cover, it cannot exist a f.d.  $X \rightarrow A$  in  $F$ , such that  $X \subseteq R \cdot A$  and, then,  $R \cdot A$  is a key of  $R$ ; let  $Y \rightarrow B$  be a d. in  $F$ ; if  $B = A$  and  $F$  is a minimal cover, then  $Y = R \cdot A$  ( $Y$  is a superkey); if  $B \neq A$  then  $B \in R \cdot A$  and so  $B$  is prime
  3. if  $XA \in \rho$ , as  $F$  is a minimal cover, it cannot exist a f.d.  $X' \rightarrow A$  in  $F$  such that  $X' \subset X$  and hence  $X$  is a key of  $XA$ ; let  $Y \rightarrow B$  be any dependency in  $F$ , such that  $YB \subseteq XA$ ; if  $B = A$  then, as  $F$  is a minimal cover,  $Y = X$  ( $Y$  is a superkey); if  $B \neq A$  then  $B \in X$ , so,  $B$  is prime

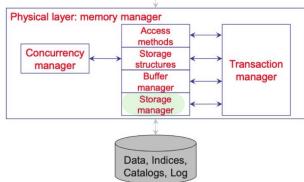
it can be either 1+2 or 1+3 or 3 only

## What about the lossless join?

it can be demonstrated that, to have a lossless join, we just need to add to  $\rho$  a subschema containing one of the keys of  $R$

# Physical Organization

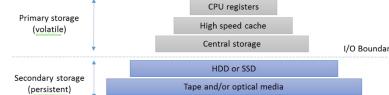
→ physical layer : memory management



## storage hierarchy

- top, primary storage:

- CPU
- Cache
- central storage, RAM
- contains DB buffer and runtime code of applications and DBMS



- bottom, secondary storage

- persistent storage media
- HDD, SSD
- contains physical DB files

## Hard Disk Drives

→ stores data on circular platters, covered with magnetic particles

→ contains a hard disk controller

→ they are directly accessible storage devices (DASD)

→ platters are secured on a spindle, that rotates at a constant speed

→ read/write heads can be positioned on arms, that are fixed to an actuator reading from a block, or writing to a block implies

└ positioning the actuator → seek time : time to position the head on the correct track

└ wait until the desired sector has rotated under the read/write head → rotation delay

└ transfer time : depends on block size, density of magnetic particles and rotation speed of disks

└ response time: service time + queuing time

    └ seek time + rotational delay + transfer time

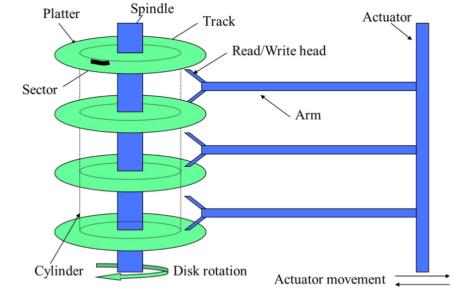
└  $t_{rba}$  (random block access): expected time to retrieve/write disk block independently of previous r/w

$$t_{rba} = \text{seek} + \text{ROT}/2 + \text{BS}/\text{TR}$$

└  $t_{sba}$  (sequential block access): expected time to sequentially retrieve disk block with r/w head already in correct position

$$t_{sba} = \text{ROT}/2 + \text{BS}/\text{TR}$$

→ ROT: rotation time, BS: block size, TR: transfer rate



## From logical To physical

Logical data model (general terminology)	Logical data model (relational setting)	Internal data model
Attribute type and attribute	Column name and (cell) value	Data item or field (bits or characters representing a specific value)
(Entity) record	Row or tuple	Stored record (collection of data items; represents all attributes of an entity)
(Entity) record type	Table or relation	Physical file or data set (similar entities, e.g., students)
Set of (entity) record types	Set of tables or relations	Physical database (collection of stored files)
Logical data structures	Foreign keys	Physical storage structures (logical interrelations like foreign keys)

## The physical database

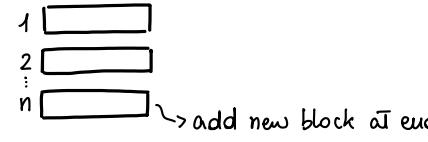
- a database consists of a set of files
- each file can be seen as a collection of pages/blocks, of fixed size
- each page stores several records (logical tuples)
- a record consists of several fields, with fixed/variable size (tuple's attributes)

## Primary file organization methods

- linear search: each record in file is retrieved and assessed against search key ] search and manage data
- hashing and indexing: direct access to the location containing the records corresponding to the search key

### 1. Heap file

- simplest
- only linear search
- on average check  $\frac{n}{2}$  blocks



$$\text{Search time: } \frac{n}{2} \cdot t_{\text{sba}}$$

### 2. Sequential file

- records in ascending/descending order of search key
- linear and binary search



$$\text{Search time: } \log_2(n) \cdot t_{\text{rba}}$$

```
- Selection criterion: record with search key value  $K_\mu$ 
- Set  $l = 1$ ;  $h = \text{number of blocks in file}$  (suppose records are in ascending order of search key  $K$ )
- Repeat until  $h \geq 1$ 
  -  $i = (l + h) / 2$ , rounded to nearest integer
  - Retrieve block  $i$  and examine key values  $K_j$  of records in block  $i$ 
  - if any  $K_j = K_\mu$  record is found!
  - else if  $K_\mu > \text{all } K_j$  continue with  $l = i+1$ 
  - else if  $K_\mu < \text{all } K_j$  continue with  $h = i-1$ 
  - else record is not in file
```

- updating done in batch, since have to shift all
- often combined with indexes

### 3. Random file (hashing)

- assumes direct relationship between value of search key and physical location

- record retrieved with few block accesses

- hashing algorithm defines key-to-address transformation

→ record physical address calculated from its key

- most effective when using primary key

- doesn't guarantee that all keys are mapped to different hash values

| collision occurs when records assigned to same bucket (synonyms)

| if more synonyms than slots for a bucket → overflow → more block accesses to retrieve overflow records

| should distribute keys evenly

- popular hashing technique is division:  $\text{address}(\text{key}_i) = \text{key}_i \bmod M$  where  $M$  is a prime number (close to, but larger, than # of available addresses), remainder of modulo division is the record address

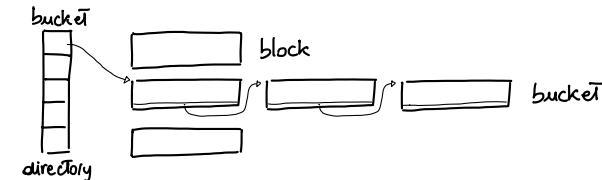
- retrieving non-overflow record:  $1 \cdot t_{\text{rba}}$  to first block of bucket, possibly followed by  $m \cdot t_{\text{sba}}$

- additional block accesses for overflow records depending on percentage of overflow and handling technique

- required number of buckets  $b$ :  $b = \lceil r / (bs \cdot lf) \rceil$   $r = \# \text{records}$ ,  $bs = \text{bucket size}$ ,  $lf = \text{loading factor}$  (how much we want the block to be filled)

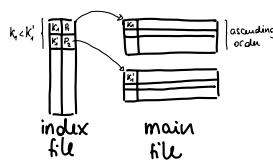
$$\text{Search time} = \frac{n}{2b}$$

- open addressing: overflow records stored in next free slot after full bucket where a record would normally have been stored



## 4. Indexed sequential file indexed sequential access method - ISAM

- file divided into intervals or partitions
- each interval is represented by index entry (in index file)



↴ search key value of first record in interval  
 ↴ pointer to physical position of first record in interval  
 ↴ block pointer  
 ↴ record pointer

### - dense or sparse indexing

↓  
 index entry for all possible values of search key

↓  
 index entry for only some of search key values

↳ each entry = group of records

] dense is faster but more storage  
 ↳ still less storage than data/main file

- data file is ordered on unique key, index is defined over this unique search key

- Search time:  $\log_2(nbi) + t_{rba}$  with  $nbi \ll n$  ( $nbi = \# \text{ of blocks in index}$ ) (index based search)  
 ↳ access block in main file

## Tree-structured indexes

### Properties:

- all children of a node are the descendants
- a node without children is a leaf node
- a tree where all leaf nodes are at the same level is called balanced
- nodes are stored in top-down-left-right sequence
  - ↳ letters are record keys, numbers are level in tree
  - ↳ can only navigate in sequential way

### trees can be implemented by linked lists

- ↳ each node has pointer to next sibling
- ↳ single bit often added if node is leaf node (0)
- ↳ parent-child and sibling-sibling navigation is supported

### binary search trees

- ↳ left subtree only contain key values that are lower than key value in original node, right subtree higher

### multiway search trees

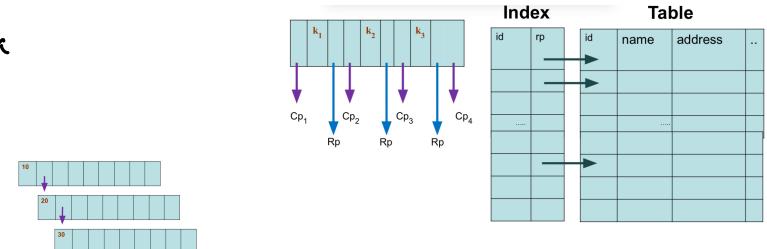
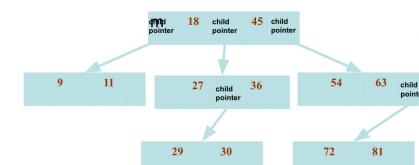
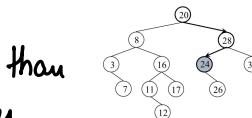
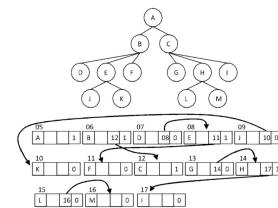
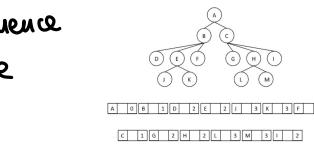
- ↳ node has  $n$  key values and  $n+1$  pointers,  $n+1 \leq m$
- ↳ all key values in subtree of  $P_i$  are less than  $k_i$
- ↳ all key values in subtree of  $P_{i+1}$  are bigger than  $k_i$
- ↳ a node can have from 1 to  $m-1$  values and the number of sub-trees is from 0 to  $1+i$ , where  $i$  is the number of key values in the node
- ↳  $m$  is a fixed upper limit that defines how many key values can be stored in the node

↳ an  $m$ -way search tree has  $m-1$  keys

↳ each node corresponds to a disk block

↳ besides child pointers ( $Cp_i$ ) we consider record pointers ( $Rp$ ) after every key

problem: no control in the way keys are inserted → should fill up first node



## → B-tree

- | a self-balanced m-way search tree
- | every node has at most m children
- | every node except the root node and leaf nodes have at least  $m/2$  children
- | root node has at least two children if it is not a leaf node
- | all leaf nodes are at same level (balanced)
- | the creation process is bottom-up
- | each node corresponds to a disk block and nodes are kept between half full and full
- | every node contains a set of search key values, a set of tree pointers that refer to child nodes, a set of data pointers that refer to data records, or blocks with data records, that correspond to search key values and that are stored separately and are not part of B-tree search
- | recursively starting from root
- | cost  $\leq$  tree height - 1 + 1
  - ↓  
root node in RAM
  - accessing the  
data file
- |  $h_{\min} = \lceil \log_m (N+1) - 1 \rceil$        $N = \# \text{ nodes}$        $m = \max \# \text{ children for node}$
- |  $h_{\max} = \lfloor \log_d ((N+1)/2) \rfloor$        $d = \min \# \text{ children of node}$

