

#### Motivation

- Can we be sure that the translation of an EERdiagram into a relational model results in a good database design ?
- Confronted with a deployed database, how can we be sure that it is well-designed?
- What is a good database design?
  - □ Four informal measures.
  - □ One formal measure: Normalization.

# Good design: Informal measures

Easy-to-explain semantics of the relational schema.

Minimal redundant information in tuples.

□ Why ? Redundancy causes waste of space and update anomalies:

- Insertion anomalies.
- Deletion anomalies.
- Modification anomalies.

EMP(	<u>EMPID</u> ,	EMPNAME,	DEPTNAME,	DEPTMGR)
	123	Smith	Research	999
	333	Wong	Research	999
	888	Borg	Administration	null
				3

Good design: Informal measures

#### Minimal number of NULL values in tuples.

- □ Why ?
  - Efficient use of space.
  - Avoid costly outer joins.
  - Ambiguous interpretation (unknown vs. doesn't apply).
- Disallow the possibility of generating spurious tuples.

□ How ? By joining only on foreign key/primary key attributes.

# Definitions

Relational schema: The header of the table.

×	<u>EmplD</u>	<u>Dept</u>	Work%	EmpName
	<u>100</u>	<u>Dev</u>	50	Baker
~	<u>100</u>	<u>Support</u>	50	Baker
	<u>200</u>	<u>Dev</u>	80	Miller

- Relation: The data in the table.
- Relation is a set, i.e. no duplicate tuples exist.

# **Functional dependencies**

- Let R be a relational schema with the attributes A<sub>1</sub>,...,A<sub>n</sub> and let X and Y be subsets of {A<sub>1</sub>,...,A<sub>n</sub>}.
- Let r(R) denote a relation in the relational schema R.

We say that X functionally determines Y, i.e.  $X \rightarrow Y$ , if for each pair of tuples  $t_1, t_2 \in r(R)$  and for all relations r(R) we have that if  $t_1[X] = t_2[X]$  then  $t_1[Y] = t_2[Y]$ .

 Despite the mathematical definition, a functional dependency cannot be discovered from the relation. It is a property of the semantics of attributes in the relational schema.

R= {ID, NAME, BIRTHDATE, TEL, CITY, ZIP}

```
 \begin{array}{c} \text{ID} \rightarrow \text{NAME} \\ \text{ID} \rightarrow \text{BIRTHDATE} \\ \text{ID} \rightarrow \text{TEL} \\ \text{ID} \rightarrow \text{CITY} \\ \text{ID} \rightarrow \text{ZIP} \end{array} \right) \text{ID} \rightarrow \text{NAME}, \text{BIRTHDATE}, \text{TEL}, \text{CITY}, \text{ZIP}
```



#### Inference rules

- 1. If  $X \supseteq Y$  then  $X \rightarrow Y$ , or  $X \rightarrow X$  (reflexive rule)
- 2.  $X \rightarrow Y \models XZ \rightarrow YZ$  (augmentation rule)
- 3.  $X \rightarrow Y, Y \rightarrow Z \models X \rightarrow Z$  (transitive rule)
- 4.  $X \rightarrow YZ \models X \rightarrow Y$  (decomposition rule)
- 5.  $X \rightarrow Y, X \rightarrow Z \models X \rightarrow YZ$  (union or additive rule)
- 6.  $X \rightarrow Y, WY \rightarrow Z \models WX \rightarrow Z$  (pseudotransitive rule)

#### Inference rules

Textbook, page 341:

"... X  $\rightarrow$  A, and Y  $\rightarrow$  B does *not* imply that XY  $\rightarrow$  AB." Prove that this statement is wrong.

Prove inference rules 4, 5 and 6 by using only inference rules 1, 2 and 3.

# Definitions

- True for any relation. So, it depends on the semantics of the attributes
- Superkey: A set of attributes that uniquely (but not necessarily minimally) identifies a tuple of a relation.
- Key: A set of attributes that uniquely and **minimally** identifies a tuple of a relation.
- Candidate key: If there is more than one key in a relation, every key is called a candidate key.
- In other words, X is a candidate key if  $X \rightarrow A_1, \dots, A_n \setminus X$ .
- Primary key: A particular candidate key is chosen to be the primary key.
- Prime attribute: Every attribute that is part of a candidate key (vs. nonprime attribute).

# Good design: Formal measure, normalization

- 1NF, 2NF, 3NF, BCNF (4NF, 5NF).
- Minimize redundancy.
- Minimize update anomalies.
- The higher the normal form, the less the redundancy. Moreover, relations become more but smaller.
- Join operations are needed to recover the original relations.
- This may harm running time. So, normalization is not mandatory. In some cases, even denormalization may be preferred. In these cases, you may want to deal with redundancy via coding (e.g. procedures, triggers, views, etc.).

### First normal form (1NF)

The relational model does not have non-atomic values.

#### $\mathbf{R}_{non1NF}$

<u>ID</u> Name		LivesIn		
<u>100</u>	Pettersson	{Stockholm, Linköping}		
<u>101</u>	Andersson	{Linköping}		
<u>102</u>	Svensson	{Ystad, Hjo, Berlin}		

Nori	nalization

R1 <sub>1NF</sub>				
<u>ID</u>	Name			
<u>100</u>	Pettersson			
<u>101</u>	Andersson			
<u>102</u>	Svensson			

This is how we dealt with multivalued attributes in the translation.

 $R2_{1NF}$ 

<u>ID</u>	<u>LivesIn</u>
<u>100</u>	<u>Stockholm</u>
<u>100</u>	<u>Linköping</u>
<u>101</u>	<u>Linköping</u>
<u>102</u>	<u>Ystad</u>
<u>102</u>	<u>Hjo</u>
<u>102</u>	<u>Berlin</u>

11

#### Second normal form (2NF)

The relational model does not have any set of nonprime attributes that is functionally dependent on part of a candidate key.

R <sub>non2NF</sub>						
<u>EmpID</u>	<u>Dept</u>	Work%	EmpName			
<u>100</u>	<u>Dev</u>	50	Baker			
<u>100</u>	<u>Support</u>	50	Baker			
<u>200</u>	<u>Dev</u>	80	Miller			

Norr	nalization

R1 <sub>2NF</sub>		
<u>EmpID</u>	EmpName	
<u>100</u>	Baker	
<u>200</u>	Miller	



R2 <sub>2NF</sub>					
<u>EmplD</u>	<u>Dept</u>	Work%			
<u>100</u>	<u>Dev</u>	50			
<u>100</u>	<u>Support</u>	50			
<u>200</u>	<u>Dev</u>	80			

#### Second normal form (2NF)

- The relational model does not have any set of nonprime attribute that is functionally dependent on part of a candidate key.
  - Why not ? Because, a part of a candidate key can have repeated values in the relation and, thus, so can have the set of non-prime attributes, which implies redundancy and thus waste of space and update anomalies.
  - □ Solution:

- Assume that  $X \rightarrow Y$  violates 2NF in a relational schema R.
- Create a new relational schema R2(X,Y).
- Remove Y from R.
- The relational schema consists now of R and R2.

#### Third normal form (3NF)

102

The relational model does not have any set of nonprime attributes that is functionally dependent on a set of attributes that is not a candidate key.

non31	Rnon3NF						
<u>ID</u>	Name	Zip	City		ID	Name	Zip City
<u>100</u>	Andersson	58214	Linköping			$\uparrow$	<u>^</u>
<u>101</u>	Björk	10223	Stockholm				
<u>102</u>	Carlsson	58214	Linköping				
R1 <sub>3NF</sub> R2 <sub>3NF</sub>							
Normalization		<u>ID</u>	Name	Zip		<u>Zip</u>	City
		<u>100</u>	Andersson	58214		<u>58214</u>	Linköping
		<u>101</u>	Björk	10223		<u>10223</u>	Stockholm

58214

Carlsson

#### Third normal form (3NF)

- The relational model does not have any set of non-prime attributes that is functionally dependent on a set of attributes that is not a candidate key.
  - Why not ? Because, a set of attributes that is not a candidate key can have repeated values in the relation and, thus, so can have the set of non-prime attributes, which implies redundancy and thus waste of space and update anomalies.
  - □ Note that 3NF implies 2NF.
  - □ Solution:
    - Assume that  $X \rightarrow Y$  violates 3NF in a relational schema R.
    - Create a new relational schema R2(X,Y).
    - Remove Y from R.
    - The relational schema consists now of R and R2.

#### Little summary

- $\mathsf{X} \to \mathsf{Y}$
- 2NF and 3NF do nothing if Y is prime.
- Assume Y is non-prime.
- 2NF = decompose if X is part of a candidate key.
- 3NF = decompose if X is not a candidate key.
- 3NF = X is a candidate key or Y is prime.

If Y is prime but X is not a candidate key, then X can have repeated values in the relation and, thus, so can have Y. Is not this a problem just because Y is prime ?

#### Boyce-Codd normal form (BCNF)

- In every functional dependency in the relational model, the determinant is a candidate key.
- Example: Let R(<u>A,B</u>,C,D) denote a relational schema with functional dependencies AB→CD, C→B. Then R is in 3NF but not in BCNF.
  - □ C is a determinant but not a candidate key.
  - □ Decompose R into R1(<u>A,C</u>,D) with AC  $\rightarrow$  D and R2(<u>C</u>,B) with C  $\rightarrow$  B.
- In general:
  - $\Box$  Assume that X  $\rightarrow$  Y violates BCNF in a relational schema R.
  - $\Box$  Create a new relational schema R2(X,Y).
  - □ Remove Y from R.
  - □ The relational schema consists now of R and R2.
  - $\Box$  You may have to find a new primary key for R.

17

#### Little summary

- $\mathsf{X} \to \mathsf{Y}$
- 2NF and 3NF do nothing if Y is prime.
- Assume Y is non-prime.
- 2NF = decompose if X is part of a candidate key.
- 3NF = decompose if X is not a candidate key.
- 3NF = X is a candidate key or Y is prime.
- Assume Y is prime.
- BCNF = decompose if X is not a candidate key.



#### Normalization: Example

Functional dependencies:

PID  $\rightarrow$  PersonName PID, Country  $\rightarrow$  NumberVisitsCountry Country  $\rightarrow$  Continent Continent  $\rightarrow$  ContinentArea

• What are the candidate keys for R? Use the inference rules to show that  $X \rightarrow A_1, ..., A_n \setminus X$ .

#### Normalization: Exmple

Country → Continent, Continent → ContinentArea then Country → ContinentArea (transitive rule) Country → Continent, ContinentArea (additive rule) PID, Country → PID, Continent, ContinentArea (augmentation rule) PID, Country → Continent, ContinentArea (decomposition rule) PID → PersonName then PID, Country → PersonName (augmentation + decomposition rules) PID, Country → NumberVisitsCountry then PID, Country → Continent, ContinentArea, PersonName, NumberVisitsCountry (additive rule)

<u>PID, Country</u> is the only candidate key for R and, thus, its primary key.

21

### Normalization: Example

ls

R (<u>PID,Country</u>,Continent,ContinentArea,PersonName,NumberVisitsCountry) in 2NF ?

No, PersonName depends on a part of a candidate key (PID), then R1(PID, PersonName)

R2(PID, Country, Continent, ContinentArea, NumberVisitsCountry)

Is R1 in 2NF ? Yes.

Is R2 in 2NF ? No, Continent and ContinentArea depend on a part of a candidate key (Country), then

R1(<u>PID</u>, PersonName)

R21(Country, Continent, ContinentArea)

R22(<u>PID, Country</u>, NumberVisitsCountry) Now, R1, R21, R22 are in 2NF.

2NF: No non-prime attribute should be functionally dependent on a part of a candidate key.



Are R1, R21, R22 in 3NF?

3NF: No nonprime attribute should be functionally dependent on a set of attributes that is not a candidate key.

R22(<u>PID, Country</u>, NumberVisitsCountry) R1(<u>PID</u>, PersonName)

Yes, because they have only one non-prime attribute.

R21(<u>Country</u>, Continent, ContinentArea) No, Continent determines ContinentArea, then R211(<u>Country</u>, Continent) R212(<u>Continent</u>, ContinentArea)

Now, R1, R22, R211, R212 are in 3NF.

Are R1, R22, R211, R212 in BCNF?

**BCNF:** Every determinant is a candidate key.

R22(<u>PID, Country</u>, NumberVisitsCountry) R1(<u>PID</u>, PersonName) R211(<u>Country</u>, Continent) R212(<u>Continent</u>, ContinentArea)

Yes, they are in BCNF.

Can the original relation R be recovered by joining R1, R22, R211 and R212 without generating spurious tuples? Yes. Mind the foreign keys created during normalization !

# Desirable properties of normalization

- Keep all the attributes from the original relational model (true in our method).
- Preserve all the functional dependencies from the original relational model (false in our method).
- Lossless join (true in our method).
  - It must be possible to join the smaller relations produced by normalization and recover the original relation without generating spurious tuples.