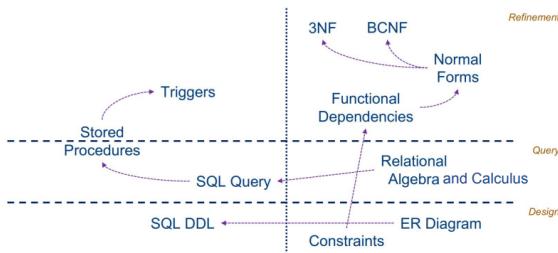


# CS2102 Database Sys Summary

AY23/24 Sem 1, github.com/gerteck

## Topics & Objectives

- **Design:** Entity-Relationship (ER) Model, Functional Dependencies, Normal Forms
- **Implementation:** SQL (Data definition language, Queries, Stored procedures, Triggers)
- **Theory:** Relational Calculus and algebra
- Module covers fundamental concepts and techniques for:
  - Understanding and practice of design & implementation of database applications and management of data with relational db management systems.
  - Design of ER data models to capture data requirements, translate to relational database schema, refine using schema decompositions to avoid anomalies.
  - Use SQL to define relational schemas, write queries.
  - Reason about correctness using concepts of formal query lang (relational calculus & algebra) and apply knowledge to develop database applications.



## 1. Database Management Sys DBMS

### Challenges for Data-Intensive applications

- **Efficiency:** Fast access to information in volumes of data
- **Transactions:** "All or nothing" changes to data
- **Data Integrity:** Parallel access and changes to data
- **Recovery:** Fast and reliable handling of failures (e.g. HD-D/Sys crash, power outage, network disruption)
- **Security:** Fine-grained data access rights

### File-based data management to DBMS

- Complex, low level code, Often similar requirements across different programs

- **Problems:** High development effort, Long development times, Higher risk of (critical) errors
- **DBMS:** Set off universal and powerful functionalities for data management, with faster application development, higher stability, less errors.

## Core concepts of DBMS

- **ACID Transaction:** Finite sequence of database operations (reads and/or writes), smallest logical unit of work
- **Atomicity:** either all effects of T are reflected in the database or none ("all or nothing")
- **Consistency:** the execution of T guarantees to yield a correct state of the database
- **Isolation:** execution of T is isolated from the effects of concurrent transactions
- **Durability:** after commit of T, its effects are permanent even in case of failures

## Concurrent Execution

### Concurrent Execution — Common Problems

T <sub>1</sub> (B, 500)	T <sub>2</sub> (B, 100)
begin	
read(B) <i>1000</i>	
B = B + 500 <i>1500</i>	
begin	
read(B) <i>1000</i>	B = B + 100 <i>1100</i>
write(B) <i>1500</i>	
commit	
	write(B) <i>1100</i>
	commit

Final balance B = 1,100  
(effect of T<sub>1</sub> overwritten)

→ Lost Update

T <sub>1</sub> (B, 500)	T <sub>2</sub> (B, 100)
begin	
read(B) <i>1000</i>	
B = B + 500 <i>1500</i>	
write(B) <i>1500</i>	write(B) <i>1500</i>
begin	
read(B) <i>1000</i>	B = B + 100 <i>1100</i>
write(B) <i>1500</i>	write(B) <i>1100</i>
abort	commit

Final balance B = 1,600  
(when it should be 1,100)

→ Dirty Read

T <sub>1</sub> (B, 500)	T <sub>2</sub> (B, 100)
begin	
read(B) <i>1000</i>	
B = B + 500 <i>1500</i>	
begin	
read(B) <i>1000</i>	B = B + 100 <i>1100</i>
write(B) <i>1500</i>	write(B) <i>1100</i>
commit	commit
	read(B) <i>1100</i>
	...

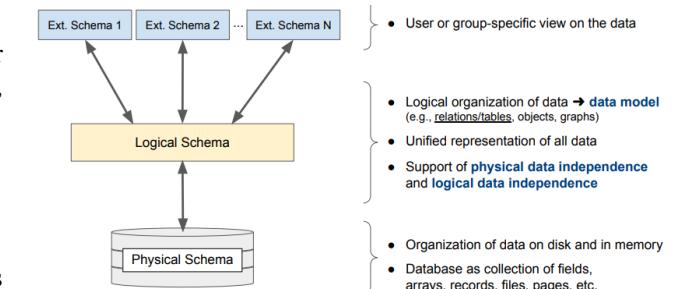
Balance B is retrieved twice  
but the values differ

→ Unrepeatable Read

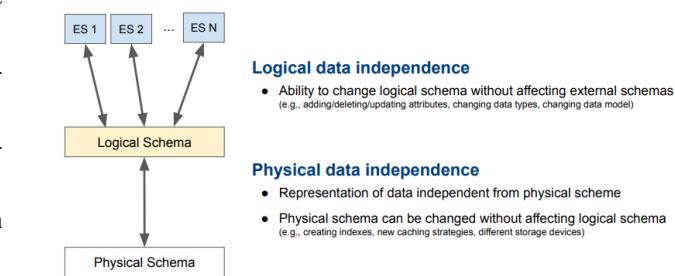
Require Serializable transaction execution:

- A concurrent execution of a set of transactions is serializable if this execution is equivalent to some serial execution of the same set of transactions
- Two executions are equivalent if they have the same effect on the data
- **DBMS:** Support concurrent executions of transactions to optimize performance, Enforce serializability of concurrent executions to ensure integrity of data

## Data Abstraction



## Data Independence



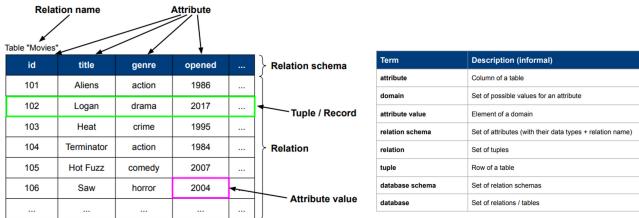
## Terminology / Definitions

- **Data Model:** Collection of concepts for describing data
- **Schema:** Description of structure of DB using data model
- **Schema Instance:** Content of a DB at a particular time

## Relational Data Model

Data is modelled by relations, and each relation has a definition called a relation schema. This schema specifies attributes (columns) and data constraints (e.g. domain constraints)

- **Relation:** Can be seen as Tables with rows and columns:
  - No. of cols = Degree/Arity, No. of rows = Cardinality
  - Each row is called a tuple/record. It has a component for each attribute of the relation.
  - A relation is thus a set of tuples and an instance of the relation schema, i.e. of a single table.
- **Domain:** Set of atomic values, e.g. integers. All values for an attribute is either in this domain or null.
- **Relational database schema:** Set of relation schemas and their data constraints, i.e. of multiple tables
- **Relational database:** Instance of the schema and is a collection of tables.



## Integrity Constraints

Condition that restricts the data that can be stored in a database instance. A legal relation instance is a relation that satisfies all specified ICs.

- **Domain Constraints:** Restrict the attribute values of relations, e.g. only integers allowed
- **Key Constraints:**
  - **Superkey:** A superkey is a subset of attributes in a relation that uniquely identifies its tuples.
  - **Key:** A key is a superkey which is minimal, i.e. no proper subset of itself is a superkey.
  - **Candidate keys:** Set of all possible keys for a relation. One of these keys is selected as the primary key.
  - **Primary key:** Chosen candidate key for a relation. Cannot be null (entity integrity constraint). Underlined in relation schema. Prime attribute: Attribute of a primary key (cannot be null)
- **Foreign Key Constraints:**
  - **Foreign key:** A foreign key refers to the primary key of a second relation (which can be itself)
  - Each foreign key value must be the primary key value in the referenced relation or be null (foreign key constraint)
  - Also known as referential integrity constraints.

Term	Description (informal)
(candidate) key	Minimal set of attributes that uniquely identify a tuple in a relation
primary key	Selected key (in case of multiple candidate keys)
foreign key	Set of attributes that is a key in referenced relation
prime attribute	Attribute of a (candidate) key

- Terminology: DB. vs DBS vs. DBMS

$$\text{DBS} = \text{DBMS} + n^* \text{DB} \quad (n > 0)$$

## 2. SQL: Structured Query Language

- Declarative language: focus on what to compute, not on how to compute
- Contains two parts: Data Definition Language and Data Manipulation Language

### Datatypes

type	description
boolean	logical Boolean (true/false)
integer	signed 4-byte integer
float8	double precision floating-point number (8 bytes)
numeric(p, s)	number with p significant digits and s decimal places
char(n)	fixed-length character string
varchar(n)	variable-length character string
text	variable-length character string
date	calendar date (year, month, day)
timestamp	date and time

### Integrity Constraints 2

- A **consistent state** of the database is a state which complies with the business rules as defined by the structural constraints and the integrity constraints in the schema.
- If an integrity constraint is violated by an operation or a transaction, the operation or the transaction is aborted and rolled back and its changes are undone, otherwise, it is committed and its changes are effective for all users.
- Five main kinds of integrity constraints in SQL: **NOT NULL, PRIMARY KEY, UNIQUE, FOREIGN KEY, CHECK.**

### Primary Key

A primary key is a set of columns that uniquely identifies a record in the table. Each table has at most one primary key. The primary key can be one column or a combination of columns.

```
-- Declare primary key as column constraint
CREATE TABLE customers (
    firstname VARCHAR(64) NOT NULL ,
    lastname VARCHAR(64) NOT NULL ,
    email VARCHAR(64) UNIQUE NOT NULL ,
    id VARCHAR(16) PRIMARY KEY,
    UNIQUE (firstname, lastname));
```

### Composite Primary Key & NOT NULL

A not null constraint guarantees that no value of the column can be set to null. A not null constraint is always declared as a row constraint. When it is explicit, it is declared with the keyword NOT NULL.

```
CREATE TABLE games (
    name VARCHAR(32) ,
    version CHAR(3) ,
    price NUMERIC NOT NULL ,
    PRIMARY KEY (name, version) );
```

### Data Insertion Populating Tables

```
INSERT INTO customers VALUES(
    'Carole', 'Yoga', 'cyoga@email.org',
    'Carole89');
```

### Deleting Tables

```
DELETE FROM customers
-- DROP deletes content \& table definition
DROP TABLE customers
DROP TABLE IF EXISTS downloads
```

### Unique

A unique constraint on a column or a combination of columns guarantees the table cannot contain two records with the same value in the corresponding column or combination of columns.

```
CREATE TABLE customers (
    firstname VARCHAR(64) NOT NULL ,
    lastname VARCHAR(64) NOT NULL ,
    email VARCHAR(64) UNIQUE NOT NULL ,
    id VARCHAR(16) PRIMARY KEY,
    UNIQUE (firstname, lastname));
```

### Foreign Key

- A foreign key constraint enforces **referential integrity**. The values in the columns for which the constraint is declared must exist in the corresponding columns of the referenced table.
- Referenced columns are usually required to be the primary key of the referenced table. Some systems relax this.

- A foreign key is declared using the keyword REFERENCES as a row constraint and the keywords FOREIGN KEY and REFERENCES as a table constraint.

```
CREATE TABLE downloads (
customerid VARCHAR (16) REFERENCES customers
(id),
name VARCHAR (32)
version CHAR (3),
FOREIGN KEY ( name, version) REFERENCES games
(name, version) );
```

## Check

- Check constraint enforces any other condition that can be expressed in SQL. Declared as row or table constraint.

```
CREATE TABLE games (
name VARCHAR(32),
version CHAR(3),
PRIMARY KEY (name, version)
price NUMERIC NOT NULL CHECK (price > 0)
-- or as table constraint:
CHECK (price > 0) );
```

## Update and Delete Propagation

- The annotations ON UPDATE/DELETE with the option CASCADE propagate the update or deletion, when there are chains of foreign key dependencies.

```
CREATE TABLE downloads(
id VARCHAR (16) REFERENCES customers (id)
ON UPDATE CASCADE
ON DELETE CASCADE,
name VARCHAR(32),
version CHAR(3),
PRIMARY KEY (id, name, version),
FOREIGN KEY (name, version) REFERENCES
games(name, version)
ON UPDATE CASCADE
ON DELETE CASCADE);
```

- Generally a good idea to constraint all columns not to be null unless there is a good design or tuning reason for not doing so.
- Think carefully about which foreign keys should be subject to cascade.

- Good idea to defer all the constraints that can be deferred. These are checked at the end of a transaction and not immediately after each operation.

## Querying Tables

### Print Table

- Wildcard '\*' to include all attributes

```
SELECT *
FROM customers;
```

### View

- We can give a name to a query, called a view. Once created, a view can be queried like a table.
- Creating a view is generally a better option than creating and populating a table, temporary or not.

```
CREATE VIEW sg\_customers AS
SELECT c.firstname, c.lastname, c.email, c.id
FROM customers c
WHERE country = 'Singapore';
SELECT * FROM sg\_customers;
```

## 3. SQL Queries

### Printing one Table

```
SELECT firstname, lastname
FROM customers
WHERE country = 'Singapore';
```

### DISTINCT and ORDER BY

- Selecting a subset of columns may result in duplicate row even if original table has a primary key.
- DISTINCT keyword eliminates eventual duplicates, requests results contain distinct rows.
- Both DISTINCT and ORDER BY involve sorting and conceptually ORDER BY is applied before SELECT DISTINCT.

```
SELECT DISTINCT name, version
FROM downloads
ORDER BY name ASC, version DESC;
```

## WHERE

- Returns rows that evaluate to true, filter rows on a Boolean condition
- Uses Boolean operators such as AND, OR and NOT, and various comparison operators such as >, <, >=, <=, <>, IN, LIKE and BETWEEN AND
- Does not return rows that evaluate to unknown/null!
- '-' matches single char
- '%' matches any sequence of zero or more chars

```
SELECT firstname, lastname
FROM customers
WHERE country IN ('Singapore', 'Indonesia')
AND (dob BETWEEN '2001-01-01' AND
      '2000-12-01' OR since >= '2016-12-01')
AND lastname LIKE 'B%'
```

- PostgreSQL use "—" for concatenation

```
-- Not to collect GST below 30 cents:
SELECT name || ' ' || version AS game, price
      * 1.07
FROM games
WHERE price * 0.07 < 0.3
```

## De Morgan's Laws

```
SELECT name
FROM games
-- all 3 are the same:
WHERE (version = '1.0' or version = '1.1')
WHERE version IN ('1.0', '1.1')
WHERE NOT (version <> '1.0' AND version <> '1.1');
```

## NULL value

- Every domain has additional value, null. Ambiguous, could be "unknown", "does not exists", or both. In SQL it is generally (but not always) "unknown".
- With null values, the logic of SQL is a three valued logic with unknown.
- Use IS (NOT) NULL for comparison with null
- COALESCE() returns the first non-null of its argument.
- COUNT(\*) counts NULL values.
- COUNT(att)AVG(att)MAX(att)MIN(att) eliminate null values

## Cross Join

- Cross join & Cartesian Product & cross product, represented by comma **CROSS JOIN**,
- Cross join with **WHERE** clause: add condition that FK columns equal to the corresponding PK columns.
- Systematically define table variables (e.g. **games AS g**)

```
SELECT *
FROM customers c, downloads d, games g
WHERE d.id = c.id
AND d.name = g.name
AND d.version = g.version
```

## 4. Algebraic SQL Queries

### Inner Join

- **JOIN** interpreted as **INNER JOIN** • Inner joins combine records from two tables whenever there are matching values in a field common to both tables.

```
SELECT *
FROM customers c JOIN downloads d ON d.id =
  c.id
JOIN games g ON d.name = g.name AND d.version =
  g.version;
```

### Natural Join

- If we give the same name to columns that are the same, can use natural join. Joins the rows that have the same values for their columns that have the same names. It also prints one of the two equated columns

```
SELECT *
FROM customers c NATURAL JOIN downloads d
  NATURAL JOIN games g;
```

### Outer Join

- outer join keeps the columns of the rows in the left (left outer join), right (right outer join) or in both (full outer join) tables that do not match anything in the other table according to the join condition and pad the remaining columns with null values.
- Better to avoid outer joins whenever possible as they introduce null values.

• **RIGHT (OUTER) JOIN, LEFT (OUTER) JOIN**  
• **FULL (OUTER) JOIN**

-- finds customers, never downloaded a game

```
SELECT c.id FROM customers c
LEFT JOIN downloads d ON c.id = d.id
WHERE d.id IS NULL;
```

### Set Operations

- Union, intersect and non-symmetric difference.
- Eliminate duplicates: **UNION, INTERSECT, EXCEPT**
- Keep duplicates: **UNION ALL, INTERSECT ALL**  
**EXCEPT ALL**

## 4. Aggregate SQL Queries

### Aggregate Functions

- The values of a column can be aggregated aggregation functions such as **COUNT()**, **SUM()**, **MAX()**, **MIN()**  
**AVG()**, **STDDEV()** etc.

```
SELECT COUNT(*)
FROM customers c;
```

-- ALL is default and omitted

-- DISTINCT needed

```
SELECT COUNT(ALL DISTINCT c.country)
FROM customers c;
```

-- Finds min, max, avg and stddev, TRUNC()
 displays 2 dp

```
SELECT MAX(g.price),
MIN(g.price),
TRUNC(AVG(g.price), 2) AS ave,
TRUNC(STDDEV(g.price), 2) AS std
FROM games g;
```

### Group By

- The GROUP BY clause creates groups of records that have the same values for the specified fields before computing the aggregate functions.
- Groups are formed after the rows have been filtered by the WHERE clause
- Recommended (and required by SQL standard) to include attributes projected in the SELECT clause in the GROUP BY clause.
- The order of columns in the GROUP BY clause does not change the meaning of the query.

```
SELECT c.country, COUNT(*)
FROM customers c
WHERE c.dob >= '2000-01-01'
GROUP BY c.country
```

```
SELECT c.country, EXTRACT( YEAR FROM c.since)
  AS regyear, COUNT(*) AS total
FROM customers c, downloads d
WHERE c.id = d.id
GROUP BY c.country, regyear,
ORDERBY regyear, c.country;
```

### Having

- Aggregate functions can be used in conditions. However, agg. functions not allowed in WHERE.
- **HAVING** clause to add conditions to be checked after the evaluation of the GROUP BY.
- **HAVING** can only involve aggregate functions, columns listed in the GROUP BY clause and subqueries.

```
SELECT c.country,
FROM customers c
GROUP BY c.country
HAVING COUNT(*) >= 100;
```

## 5. Nested SQL Queries

### Subqueries

- In FROM clause: Must be enclosed in parenthesis, Table alias mandatory, Column aliases optional
- Not recommended, can be written as simple query

```
SELECT cs.lastname, d.name
FROM (SELECT *
      FROM customers c
      WHERE c.country = 'Singapore') AS cs,
      downloads d
WHERE cs.id = d.id;
```

- In WHERE clause, also can be written as simple query.
- Never use a comparison to a subquery without specifying the quantifier ALL or ANY

```
SELECT g1.name, g1.version, g1.price
FROM games g1
WHERE g1.price >= ALL (
    SELECT g2.price
    FROM games g2);
-- or do:
WHERE g1.price = ALL (
    SELECT MAX(g2.price)
    FROM games g2);
-- Note HAVING g.price=MAX(g.price) will not
  work
```

### Exists

- EXISTS** evaluates to true if the subquery has some results.
- Generally correlated. If uncorrelated, likely unnecessary.

```
SELECT c.id FROM customers c
WHERE NOT EXISTS (
    SELECT d.id
    FROM downloads d
    WHERE c.id = d.id);
-- same as
WHERE c.id NOT IN (
    SELECT d.id
    FROM downloads d);
-- same as
WHERE c.id <> ALL (
    SELECT d.id
    FROM downloads d);
```

```
-- Find countries with most customers
SELECT c1.country
FROM customers c1
GROUP BY c1.country
HAVING COUNT(*) >= ALL (
    SELECT COUNT(*)
    FROM customers c2
    GROUP BY c2.country);
```

### SQL Conditionals

**CASE expression** • Generic conditional, like if/else statement, Used in SELECT, ORDER BY, etc

```
-- Regular if else
CASE
    WHEN cond1 then res1
    WHEN cond2 then res2
    ...
    WHEN condN then resN
    ELSE res0
END
```

```
-- Switch like statement
CASE expression
    WHEN val1 then res1
    WHEN val2 then res2
    ...
    WHEN valN then resN
    ELSE res0
END
```

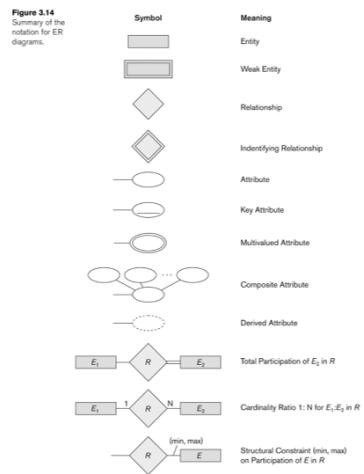
### Conceptual evaluation of queries

FROM → WHERE → GROUP BY → HAVING → SELECT  
→ ORDER BY → LIMIT/OFFSET

## 6. ER Model

Data Modeling Using the Entity-Relationship (ER) Model

### NOTATION for ER diagrams

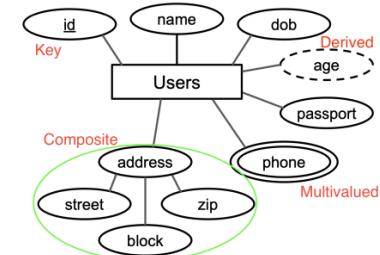


### Entity

- Objects that are distinguishable from other objects
- Entity set:** Collection of entities of the same type
- In ER diagrams, an entity type is displayed in a rectangular box

### Attributes

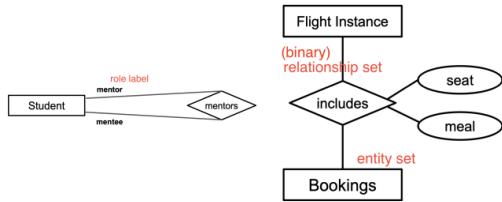
- Attributes** are properties used to describe an entity
- Each attribute has a value set (or data type) associated with it – e.g. integer
- Key attribute(s)** uniquely identifies each entity
- Composite attribute** composed of multiple other attributes
- Multivalued attribute** may consist of more than one value for a given entity
- Derived attribute** derived from other attributes



- Attributes are displayed in ovals, multivalued attributes double ovals

## Relationship

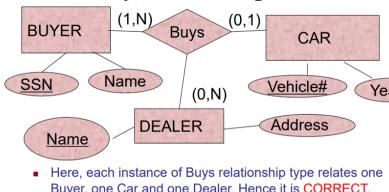
- A **relationship** relates two or more distinct entities with a specific meaning.
- **Degree** of a relationship type is the number of participating entity types. A n-ary relationship set involves n entity roles. Typically binary or ternary.



- **Relationship Set:** Collection of relationships of the same type, can have their own attributes that further describe the relationship
- Most **relationship attributes** are used with M:N relationships: (In 1:N relationships, they can be transferred to the entity type on the N-side of the relationship)
- We represent the relationship type as **Diamond-shaped box**, connected to the participating entity types.
- Relationship type typically readable from left to right and top to bottom.

## N-ary relationships

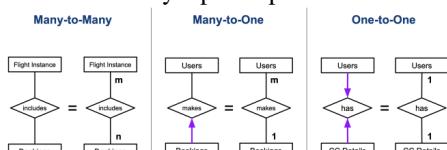
- **Implicit Constraint** In ternary relationship, every instance of relationship must have one instance of each entity.
- Rule extends to n-ary relationships



- **Avoid Ternaries for Easy Modeling:** e.g. “objectifying” the relationship type “Interview” into an entity type ‘Interview’.

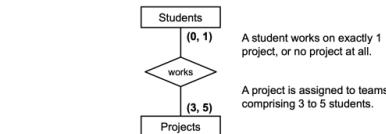
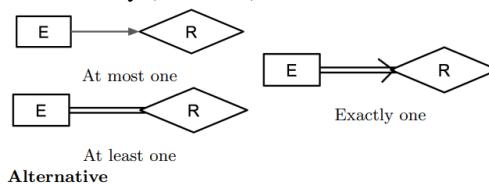
## Cardinality (Ratio) constraints

- **Upper bound** for entity's participation



## Participation / Existence Dependency constraints

- **Lower bound** for entity's participation
- Partial (default): participation not mandatory
- Total: mandatory (at least 1)



## Recursive Relationship Type

- A relationship type between the same participating entity type in **distinct roles**.
- In ER diagram, need to display role names to distinguish participations.

## Dependency constraints

**Weak Entity Types:** No key attribute and is identification dependent on another entity.

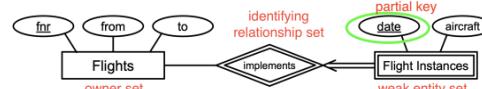
- Participate in an identifying relationship type with an owner or identifying entity type
- Two kinds of dependencies:

**Existence (no weak entity) dependency**

**Identification (weak entity) dependency**

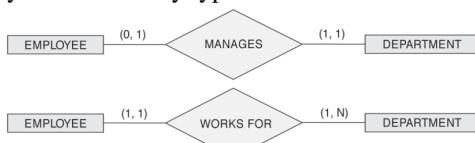
### Partial Key

- Set of attributes of weak entity set that uniquely identifies a weak entity, for a given owner entity.



## (min,max) notation for relationship constraints

- Read the min,max numbers next to the entity type and looking away from the entity type



## 7. EER Model

### Enhanced ER or Extended ER

#### IS-A Hierarchies

- **“Is a” relationship** - used to model generalization/specialization of entity sets.
- **Hierarchy:** constraint that every subclass has only one superclass (single inheritance); basically a tree structure.
- **Lattice:** a subclass can be subclass of more than one superclass (multiple inheritance).

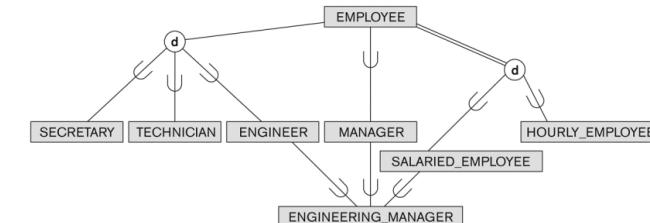


Figure 4.6

A specialization lattice with shared subclass ENGINEERING\_MANAGER.

- **Subclass:** subclass member is entity in a **distinct specific role**. Entity inherits all attributes and relationships of superclass.
- **Specialization:** process of defining a set of subclasses of a superclass, based on **distinguishing characteristics**
- **Type of Specialization:** Can be Predicated defined, Attribute defined (Written beside joining line) or User defined.
- **Generalization:** reverse of specialization, based upon some **distinguishing characteristics**.

## Constraints on Specialization / Gen

- **Disjointness Constraint:** Subclasses must be disjoint, entity can be member of at most one. (specified by **d** in EER)
- If not disjoint, specialization is **overlapping**, (specified by **o** in EER)
- **Completeness Constraint:** Total specifies every entity must be a member of some subclass. (specified by **double line** in EER)
- **Partial** allows entity not to belong to any subclass. (specified by **single line** in EER)
- Hence, we have four types: Disjoint/Overlapping x Total-/Partial. Note generalization usually is total.

## 8. Relational Mapping

- Preserve info, maintain constraints, minimize null.

### 1. Map Regular Entity Types

- Create new table (relation R), include all simple attributes.

### 2. Map Weak Entity Types

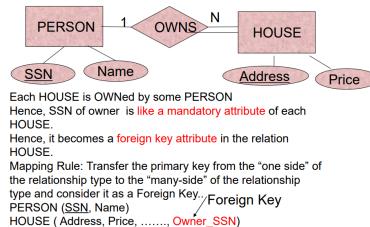
- Weak entity type, create table with corresponding FK.

### 3. Map Binary 1:1 Relationship Type

- Option 1: 2 Foreign Key (2 relations tables) option.
- 2: Merged relation option: Combine relationship set and either entity set into **one** table
- 3: Cross-reference (3 relations, one lookup table).

### 4. Map Binary 1:N Relationship Type

- For 1:N create table of participating entity (N-side).
- Include FK in table & attributes of the 1:N relation type.



### 5. Map Binary N:M Relationship Type

- Represent relationship set with a new table. Include as FK the PK of relations, combination will form PK of table.

### 6. Map Multi-valued Attributes

- Create new table / relation for each multivalued attribute.

### 7. Map N-ary relationship type

- For each n-ary relationship type R, where n > 2, create new table. Include FK of relations, and attributes.

### Summary:

Table 9.1 Correspondence between ER and Relational Models

ER MODEL	RELATIONAL MODEL
Entity type	Entity relation
1:1 or 1:N relationship type	Foreign key (or <i>relationship</i> relation)
M:N relationship type	<i>Relationship</i> relation and two foreign keys
<i>n</i> -ary relationship type	<i>Relationship</i> relation and <i>n</i> foreign keys
Simple attribute	Attribute
Composite attribute	Set of simple component attributes
Multivalued attribute	Relation and foreign key
Value set	Domain
Key attribute	Primary (or secondary) key

### 8. Map Specialization or Generalization

- Option 1: Multiple tables (relations) - Superclass and subclasses
- 2: Table each for each subclass relations, inherit superclass attributes
- 3: Single table with one type attribute
- 4: Single table with multiple type attributes
- For specialisation hierarchies with one superclass and n subclasses, possibility of mapping from 1 relation to (n+1) relations. Design highly subjective, try to maintain appropriate attributes to determine subclass identity.

## 9. Relational Calculus & Algebra

- Both are formal query languages. A query is composed of a collection of operators called relational operators.
- Relational Calculus** (declarative language: what is to be done rather than how to do it). Order not specified, concerned with result we have to obtain.
- Relational Algebra:** (procedural language) Order specified in which operations have to be performed.
- Relational algebra basis of implementation of relational DBMS.
- SQL queries translated into execution plans, orchestrate physical implementation of relational algebra operators.

### Predicate Logic

- Predicate / first order logic:** formulae built from:
- predicates** (lowercase), **operators** ( $=$ , etc.) **constants** (lower case), **variables** (upper case, quantified or free), **connectives** ( $\rightarrow$ ), and **quantifiers** ( $\forall$  and  $\exists$ ).
- Existential quantifier**,  $\exists$ : disjunction:  
 $\exists X, F[X] \equiv F[a] \vee F[b] \vee \dots$
- The universal quantifier**, *forall*: conjunction.  
 $\forall X, F[X] \equiv F[a] \wedge F[b] \wedge \dots$
- De Morgan laws** applies to quantifiers.  
 $\neg(\exists X, F[X]) \equiv \forall X, \neg F[X]$   
 $\neg(\forall X, F[X]) \equiv \exists X, \neg F[X]$

### Relational Calculus (Tuple)

- Concerned with **result we have to obtain**.
- In **Tuple Relation Calculus (TRC)** variables values of rows of tables (tuples.), is the theoretical basis of SQL.
- Material Implication:**  $P \rightarrow Q \equiv \neg P \vee Q$
- Relational Calculus denoted as:**

$$\{t | P(t)\}$$

- t:** set of tuples, **P:** condition which is true for the given set of tuples.
- E.g.  $\{T | \exists T_1 (T_1 \in \text{department} \wedge T_1.\text{faculty} = \text{'School of Computing'} \wedge T.\text{department} = T_1.\text{department})\}$

## Relational Algebra

- Order is specified in which operations have to be performed.

### Unary Operators

#### Selection, $\sigma_c$

- For each tuple  $T \in R, T \in \sigma_c(R)$ , means selection condition  $c$  evaluates to true for tuple  $t$ .
- E.g. Find the customers from Singapore.

$$\sigma_{c.\text{country} = \text{'Singapore'}}(\rho(\text{customers}, c))$$

- Equivalent to:

```
SELECT * FROM customers c
WHERE c.country = 'Singapore'
```

- Condition** is boolean expression of form:

expression	example
attribute <b>op</b> constant	$\sigma_{\text{start}=2020}(\text{Projects})$
<i>attr<sub>1</sub></i> <b>op</b> <i>attr<sub>2</sub></i>	$\sigma_{\text{start}=\text{end}}(\text{Projects})$
<i>expr<sub>1</sub></i> $\wedge$ <i>expr<sub>2</sub></i>	$\sigma_{\text{start}=2020 \wedge \text{end}=2021}(\text{Projects})$
<i>expr<sub>1</sub></i> $\vee$ <i>expr<sub>2</sub></i>	$\sigma_{\text{start}=2020 \vee \text{end}=2021}(\text{Projects})$
$\neg \text{expr}$	$\sigma_{\neg(\text{start}=2020)}(\text{Projects})$
( <i>expr</i> )	-

- op**  $\in \{=, <, <, \leq, \geq, >\}$
- Precedence: () , **op**,  $\neg$ ,  $\wedge$ ,  $\vee$
- null** comparison is **unknown**, arithmetic with **null** is **null**

### Projection $\pi_l$

- Projects columns of a table specified in **list l**.
- (**SELECT xx, yy FROM games**)
- Order of attribute in l matters.**
- Duplicates removed.
- Examples:**

$$\pi_{g.name, g.version, g.price}(\rho(games, g))$$

Teams		
en	pn	hours
Sarah	BigAI	10
Sam	BigAI	5
Sam	BigAI	3

$\pi_{pn, en}(\text{Teams})$	
pn	en
BigAI	Sarah
BigAI	Sam

### Renaming, $\rho_l$

- Can change name of relation, of attributes, or both.
- Change name of relation:  $\rho(R_1, R_2)$
- Change attribute names:  $\rho(R_1, R_1(a_1 \rightarrow b_1, a_2 \rightarrow b_2))$

### Set Operations

- Set operations include  $\cup, \cap, \times$ , set difference ( $\setminus$ )
- Intersection able to express with union and set difference:  
 $R \cap S = (R \cup S) - ((R - S) \cup (S - R))$
- Union Compatability:** two relations must be union compatible. Have same number of attributes, corresponding attributes have same or compatible domains.
- (i.e. relations must have same columns).
- Cross Product:** (Cartesian Product) Forms all possible pairs of tuples from two relations.

$$R_1 \times R_2$$

### Join Operations

- Combines  $\times, \sigma_c, \pi_l$  into a single op.
- Simple relational algebra expressions

### Inner Joins

- Eliminates tuples that do not satisfy matching criteria (i.e. selection)
- Is a selection from cross product  
 $R \bowtie_C S = \sigma_C (R \times S)$
- Example:  
 $\rho(\text{customers}, c) \bowtie_{d.id=c.id} \rho(\text{downloads}, d)$

### Relational Calculus & Algebra

- 4 Steps to construct calculus and algebra queries:**
  - Construct SQL query you are familiar with (difficult)
  - From query, map the tables that you need (yellow)
  - From query, map the conditional statements (blue)
  - From query, map the columns you need to print (green)

Q1 (A)	Q2 (A)
Q: Find the different departments in School of Computing $(T_1   T_1 \in \text{department} \wedge T_1.\text{faculty} = \text{'School of Computing'}) \wedge T_1.\text{department} = T_2.\text{department}$	Q: Find the different departments in School of Computing. $\rho_{d.department}(\sigma_{d.faculty = \text{'School of Computing'}}(\rho_{d.department}))$
Equivalent Query $\text{SELECT DISTINCT } d.department$ $\text{FROM department } d$ $\text{WHERE } d.faculty = \text{'School of Computing'}$	Equivalent Query $\text{SELECT DISTINCT } d.department$ $\text{FROM department } d$ $\text{WHERE } d.faculty = \text{'School of Computing'}$

# 10. Programming with SQL

## Writing Database Applications

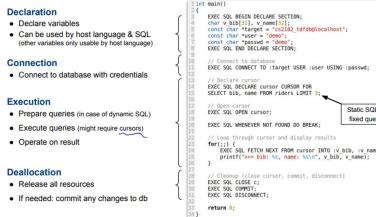
- Interactive SQL:** Directly writing SQL statements to an interface. (e.g. PostgreSQL's psql cli, pgAdmin).
- Non-interactive SQL:** SQL statements included in application written in host language.
- 2 Main alternatives: **Statement Level Interface (SLI), and Call Level Interface (CLI).**
- Crudely, SLI = CLI in disguise, as SLI preprocessor generates CLI code.

## Statement Level Interface (SLI)

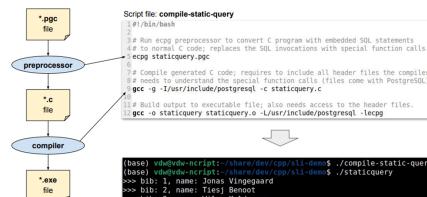
- Code is mix of host language statements and SQL statements (e.g. embedded SQL, dynamic SQL).
- Basic process for SLI:** Write code that mixes host language with SQL, preprocess code using a preprocessor, compile code into exe program.

### Statement Level Interface (SLI)

#### SLI — Common Steps



#### SLI — Preprocessing, Compiling, Running Code



#### SLI — Dynamic SQL

- Dynamic SQL:**
  - SQL query is generated at runtime
  - Example on the right: number of riders are specified as command line parameter



## Call Level Interface (CLI)

- Application completely written in host language, while **SQL statements are strings** passed as **arguments** to host language procedures or libraries
- E.g. ODBC (Open DataBase Connectivity), JDBC (Java DB Connectivity), psycopg library for Python - PostgreSQL.

### CLI — Static SQL Example

```

import psycopg # Host language library (here psycopg for Python)

# Connect to database
connection = psycopg.connect("host=localhost dbname=cs2102_tdfdb user=demo password=demo")

# Create cursor
cursor = connection.cursor()

# Open cursor by executing query (string parameter passed to execute() method)
cursor.execute("SELECT bid, name FROM riders LIMIT 3")

# Loop over all results until no next tuple is returned
while True:
    row = cursor.fetchone()
    if row is None:
        break
    else:
        print(f">>>> bid: {row[0]}, name: {row[1]}")

# Clean up
cursor.close()
connection.commit()
connection.close()

>>> bid: 1, name: Jonas Vinggaard
>>> bid: 2, name: Tiesj Benoot
>>> bid: 3, name: Wilco Kelderman
  
```

### CLI — Dynamic SQL Example

```

import psycopg # Host language library (here psycopg for Python)

# Set a user-defined value (here: maximum number of riders returned)
limit = 5

# Connect to database
connection = psycopg.connect("host=localhost dbname=cs2102_tdfdb user=demo password=demo")

# Create cursor
cursor = connection.cursor()

# Open cursor by executing query (string parameter passed to execute() method)
cursor.execute("SELECT * FROM riders LIMIT %s", (limit,))

# Loop over all results until no next tuple is returned
while True:
    row = cursor.fetchone()
    if row is None:
        break
    else:
        print(f">>>> bid: {row[0]}, name: {row[1]}")

# Clean up
cursor.close()
connection.commit()
connection.close()

>>> bid: 1, name: Jonas Vinggaard
>>> bid: 2, name: Tiesj Benoot
>>> bid: 3, name: Wilco Kelderman
>>> bid: 4, name: Stepp Kost
>>> bid: 5, name: Christophe Laporte
  
```

## SQL Injection Attack

- Class of cyber attacks on dynamic SQL, goal is to execute unintended (malicious) SQL statements.
- Typical cause:** dynamic queries are generated by merging / concatenating strings.
- Common attack point:** Omnipresent form fields in web interfaces. Entered values define some SQL statement.
- Key Points:** Don't manually merge values to a query, don't use % or + operator to merge values, use provided methods.

# 11. SQL Functions and Procedures

- Tasks **requiring multiple DB operations** common, involve any combination of reads and writes.
- E.g. update user password: check user exists → check new password differs from old → if ok, update password (3 separate requests/accesses to DB)
- Problems:** Application and DB may run on different machines, poor performance or DB becomes bottleneck.
- Different DB operations only loosely connected, difficult to ensure “all or nothing” behavior.
- Approach:** Move (some) application logic into DB, group DB operations that form task together, treat task as single DB operation.

## 11. Stored Functions and Procedures

- Collection of SQL statements and procedural logic,** precompiled and reusable code, allows execute multiple database operations as a single unit.
- Procedural Logic:** Relevant for application logic that requires assignments, conditionals or loops, and queries that cannot be expressed using basic SQL.
- ISO standard:** SQL/PSM (Persistent Stored Modules). Different DBMS have their own flavor.
- Advantages:** better performance, code reuse, ease of maintenance, added security.
- Disadvantages:** testing & debugging more challenging, limited portability / vendor lock in, no simple versioning of code, not the most intuitive language.

# Stored Functions, Procedures

## Syntax: Stored Functions

```

CREATE [OR REPLACE] FUNCTION <name> (<arg_1>, <arg_2>, ...)
RETURNS <type> AS
$$
DECLARE
    -- Declaration of variables
BEGIN
    -- Sequence of SQL statement
    -- and/or procedural logic
    [RETURN ...]
END;
$$
LANGUAGE <language>;

```

**Name of function**: The name of the function.

**Function arguments (mode, name, type)**: The arguments of the function, including their mode (IN, OUT, INOUT), name, and type.

**Return type**: The return type of the function.

**Body of function**: The body of the function, enclosed within dollar quotes and treated as a string.

**Language used in body (mainly: plpgsql or sql)**: The language used in the body of the function.

- **CREATE OR REPLACE** helps to re-declare function/procedure if already previously defined
- Code is enclosed within `$$ <> $$`
- Calling a function: (USE SELECT, e.g.)  
`SELECT * FROM swap(2, 3);`
- Call a procedure: (USE CALL, e.g.)  
`CALL transfer('Alice', 'Bob', 100);`

## Syntax: Stored Procedures

```

CREATE PROCEDURE add_bonus_proc(sid INT, amount INT)
AS
$$
    UPDATE students
    SET points = points + amount
    WHERE id = sid;
$$
LANGUAGE sql;

CALL add_bonus(3, 5);

```

No output / result, but table gets updated

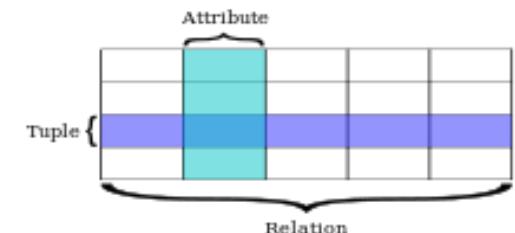
- Syntax essentially same for procedures and functions, but procedures invoked using CALL command.
- **Obvious Difference:** Procedures no RETURNS clauses.
- **Functions** must return something (but can be VOID).
- **Procedures** do not have to return anything, but can (using INOUT and OUT params).

## Function Arguments for Functions

- **Each argument described by 3 values**
- **Mode:** of argument (mainly IN, OUT, INOUT)
- **Name:** of argument (optional)
- **Type:** datatype of argument. (e.g. INT, VARCHAR)

IN	OUT	INOUT
Default	Explicitly specified	Explicitly specified
Value is passed to a function	Value is returned by a function	Value is passed to the function which returns another updated value
Behaves like constants	Behaves like an uninitialized variable	Behaves like an initialized variable
Value <u>cannot</u> be assigned	Value <u>must</u> be assigned	Value <u>can/should</u> be assigned

## Return and Type



Return	Type
One existing tuple from table	<table_name>
Set of tuples from table	SETOF <table_name>
Single new tuple	RECORD
Set of new tuples	SETOF RECORD or TABLE(attributes...)
No return value	VOID, or use PROCEDURE instead of FUNCTION
Trigger	TRIGGER

## sql VS plpgsql

- **sql:** Use where body consists of only SQL statements, often a wrapper of single / few SQL statements. Simpler syntax, no {BEGIN ... END}
- **PL/pgSQL:** Procedural Lang/ PostgreSQL, allows writing of procedural code providing control flows, variables, error handling. Statements generated at runtime, used for trigger functions.

### Function

```

CREATE OR REPLACE FUNCTION swap
    (INOUT val1 INT, INOUT val2 INT)
RETURNS RECORD AS $$

DECLARE
    temp INT;
BEGIN
    temp := val1; val1 := val2; val2 := temp;
END;
$$ LANGUAGE plpgsql;

```

### Procedure

```

CREATE OR REPLACE PROCEDURE transfer
    (src TEXT, dst TEXT, amt NUMERIC)
AS $$

    UPDATE Accounts
    SET bal = bal - amt WHERE name = src;
    UPDATE Accounts
    SET bal = bal + amt WHERE name = dst;
$$ LANGUAGE sql;

```

Important: If we use RECORD, we must have at least two OUT parameters. But if we use TABLE construct, we can just have one attribute.

## Stored Functions vs. Procedures

- Procedures can **commit or roll back transactions** during execution, cannot be involved in DML commands (select, insert, update, delete).
- Procedures invoked in isolation using `CALL`, functions invoked in `SELECT` statements.
- **Best practice:** return value(s): create function, no return value: create procedure.

## Assignments (of values of variables)

- **Basic Assignment** with `:=`, e.g. `age := 29;`
  - **Assignment of query result** to declared variable(s):  
`SELECT ... INTO ...`
- ```

SELECT points INTO mark
FROM students WHERE id = sid;

```

## Control Structures:

### • Conditionals:

- 4 types of **IF** expressions
  - **IF ... THEN ... END IF**
  - **IF ... THEN ... ELSE ... END IF**
  - **IF ... THEN ... ELSIF ... THEN ... ELSE ... END IF**
- 2 types of **CASE** expressions
  - **CASE ... WHEN ... THEN ... ELSE ... END CASE**
  - **CASE WHEN ... THEN ... ELSE ... END CASE**

### • Simple Loops

- **LOOP ... END LOOP** (typically requires **EXIT...WHEN...** to jump out of loop)
- **WHILE ... LOOP ... END LOOP**
- **FOR ... IN ... LOOP ... END LOOP**

- No curly braces or colons, hence additionaly keywords to indicate where loop begins and ends.
- **END IF** for conditionals, **END LOOP** for loops.
- **Simple example:** Compute sum of first n integers, if n is negative, return 0.

```
CREATE FUNCTION sum_n(IN n INT)
RETURNS INT AS $$$
DECLARE sum INT;
BEGIN
    sum := 0;
    IF n <= 0 THEN
        RETURN sum;
    END IF;
    FOR val IN 1..n
    LOOP
        sum := sum + val;
    END LOOP;
    RETURN sum;
END; $$$
LANGUAGE plpgsql;
SELECT * FROM sum_n(5);
```

We can also raise an exception if n is negative:

```
IF n <= 0 THEN
    RAISE EXCEPTION 'n<0 error';
END IF;
```

## Errors & Messages

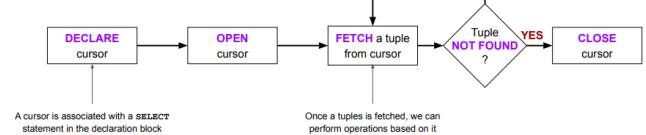
- **RAISE** keyword. 6 raise levels in PostgreSQL.

|                 |                                                                                                                                                                                                                                                                                                                                           |
|-----------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| RAISE DEBUG     |                                                                                                                                                                                                                                                                                                                                           |
| RAISE LOG       |                                                                                                                                                                                                                                                                                                                                           |
| RAISE INFO      |                                                                                                                                                                                                                                                                                                                                           |
| RAISE NOTICE    |                                                                                                                                                                                                                                                                                                                                           |
| RAISE WARNING   |                                                                                                                                                                                                                                                                                                                                           |
| RAISE EXCEPTION | <ul style="list-style-type: none"><li>• Generate messages of different priority levels</li><li>• Whether messages of a particular priority are reported to the client, depends on the PostgreSQL configuration</li></ul> <ul style="list-style-type: none"><li>• Raises an error</li><li>• Typically aborts current transaction</li></ul> |

## Cursors

- **Purpose:** Declare on a query, access each indiv row.
- Helps avoids memory overrun when the query result is large (don't access whole query at once).

### • General workflow



```
CREATE FUNCTION compute_points_gaps()
```

```
RETURNS TABLE(
    name TEXT, points INT, gap INT) AS $$$
DECLARE
    c CURSOR FOR (SELECT * FROM students ORDER
        BY points DESC);
    s RECORD; prev INT;
BEGIN
    prev := -1;
    OPEN c;
    LOOP
        FETCH c INTO s;
        EXIT WHEN NOT FOUND;
        name := s.name;
        points := s.points;
        IF prev >= 0 THEN
            gap := prev - s.points;
        ELSE
            gap := 0;
        END IF;
        RETURN NEXT; -- (OUTPUT) TABLE TUPLE
        prev := s.points;
    END LOOP;
    CLOSE c;
END; $$$
LANGUAGE plpgsql;
-- ** To check NULL: IF high IS NULL THEN,
-- not IF high = NULL THEN
```

## Advantage of Cursor (over for loops etc.):

- Flexible 'navigation' through query results in **different directions**.
- **FETCH** to move row and read data.
- **MOVE** only to move to row (no read).

## Cursor Directions

|                   |                                                                                                                                                                                                                                                                                                                              |
|-------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| NEXT              | Fetch the next row (default)                                                                                                                                                                                                                                                                                                 |
| PRIOR             | Fetch the prior row                                                                                                                                                                                                                                                                                                          |
| FIRST             | Fetch the first row of the query (same as ABSOLUTE 1)                                                                                                                                                                                                                                                                        |
| LAST              | Fetch the last row of the query (same as ABSOLUTE -1)                                                                                                                                                                                                                                                                        |
| ABSOLUTE <i>n</i> | <ul style="list-style-type: none"> <li>Fetch the <i>n</i>-th row of the query, if <i>n</i> &gt;= 0</li> <li>Fetch abs(<i>n</i>)-th row from the end, if <i>n</i> &lt; 0.</li> <li><b>ABSOLUTE 0</b> positions before the first row</li> </ul>                                                                                |
| RELATIVE <i>n</i> | <ul style="list-style-type: none"> <li>Fetch the <i>n</i>-th succeeding row, if <i>n</i> &gt;= 0</li> <li>Fetch the abs(<i>n</i>)-th prior row, if <i>n</i> &lt; 0</li> <li>Position before first row or after last row if <i>n</i> is out of range</li> <li><b>RELATIVE 0</b> re-fetches the current row, if any</li> </ul> |
| FORWARD           | Fetch the next row (same as NEXT)                                                                                                                                                                                                                                                                                            |
| BACKWARD          | Fetch the prior row (same as PRIOR).                                                                                                                                                                                                                                                                                         |

## Examples

- Using FETCH ABSOLUTE, FETCH NEXT, FETCH RELATIVE to calculate median points().
- Dynamic cursors:** Cursors can also have inputs, which are taken from function inputs, that affect the query results.
- SELECT median\_points(TRUE);  
vs SELECT median\_points(FALSE);

## Cursors — Example (beyond NEXT)

```
CREATE OR REPLACE FUNCTION median_points()
RETURNS NUMERIC AS
$$
DECLARE
    c CURSOR FOR (SELECT * FROM students ORDER BY points DESC);
    s1 RECORD; s2 RECORD; num_students INT;
BEGIN
    OPEN c;
    SELECT COUNT(*) INTO num_students FROM students;
    IF num_students%2 = 1 THEN
        FETCH ABSOLUTE (num_students+1)/2 FROM c INTO s1;
        RETURN s1.points;
    ELSE
        FETCH ABSOLUTE num_students/2 FROM c INTO s1;
        FETCH NEXT FROM c INTO s2;
        RETURN (s1.points+s2.points)/2;
    END IF;
    CLOSE c;
END;
$$
LANGUAGE plpgsql;
```

## Dynamic Cursors — Example

```
CREATE OR REPLACE FUNCTION median_points(IN has_graduated BOOLEAN)
RETURNS NUMERIC AS
$$
DECLARE
    c CURSOR (grad BOOLEAN) FOR (SELECT * FROM students
                                    WHERE graduated = grad
                                    ORDER BY points DESC);
    s1 RECORD; s2 RECORD; num_students INT;
BEGIN
    OPEN c(has_graduated);
    SELECT COUNT(*) INTO num_students
    FROM students WHERE graduated = has_graduated;
    IF num_students%2 = 1 THEN
        FETCH ABSOLUTE (num_students+1)/2 FROM c INTO s1;
        RETURN s1.points;
    ELSE
        FETCH ABSOLUTE num_students/2 FROM c INTO s1;
        FETCH NEXT FROM c INTO s2;
        RETURN (s1.points+s2.points)/2;
    END IF;
    CLOSE c;
END;
$$
LANGUAGE plpgsql;
```

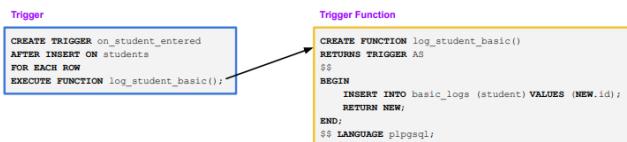
## 12. Triggers

### Motivation

- Model constraints**, where user is not forced to call stored procedure or function.
- Example constraints: restricted change in data, derived values, set cardinality constraints.
- Application Requirement:** E.g. automatic logging of changes.
- We could create a stored procedure that combines insertion and logging, but users can circumvent this by directly running **INSERT INTO** instead of calling procedure.

### Trigger Concept (Trigger fires Trigger Function)

- Triggers has **Event-Condition-Action** rule.
- When event occurs, test condition, and if satisfied, perform action.
- ECA rule** is split into 2 parts:  
Event and Condition under **Trigger**, and  
Action under **Trigger Function**.



## Triggers

- Triggers can listen on multiple event types.

```
CREATE TRIGGER
    on_student_modified_advanced
AFTER INSERT OR DELETE OR UPDATE ON students
FOR EACH ROW
EXECUTE FUNCTION log_student_advanced();
```

- Triggers options: **Event**, **Timing**, **Granularity**.
- Event:** Situation that fires trigger.
- Timing:** When trigger is fired (BEFORE or AFTER) for tables, (INSTEAD OF) for views.
- Granularity:** Specifies if triggered for each affect row or only once. (FOR EACH ROW / F.E. STATEMENT)

### Trigger Event types:

|                             |           |          |
|-----------------------------|-----------|----------|
| INSERT ON table             | → TG_OP = | 'INSERT' |
| DELETE ON table             |           | 'DELETE' |
| UPDATE [OF column] ON table |           | 'UPDATE' |

### Access to transition variables:

|        | NEW | OLD |
|--------|-----|-----|
| INSERT | ✓   | ✗   |
| UPDATE | ✓   | ✓   |
| DELETE | ✗   | ✓   |

### Trigger Timing:

|               | RETURN value                                         |                  |
|---------------|------------------------------------------------------|------------------|
| AFTER         | NULL tuple                                           | non-NULL tuple t |
| BEFORE        | Trigger fires before the operation is attempted      |                  |
| INSTEAD OF    | Trigger fires if an operation on a view is attempted |                  |
| AFTER INSERT  | No tuple inserted                                    | Tuple t inserted |
| BEFORE UPDATE | No tuple updated                                     | Tuple t updated  |
| BEFORE DELETE | No tuple deleted                                     | Tuple t deleted  |
| AFTER UPDATE  |                                                      | No effects!      |
| AFTER DELETE  |                                                      |                  |

## Views (Recap) and Triggers

- Views:** virtual table, a permanently named query.
- Query results not permanently stored, executed each time query used, hides complexity, heavily used.

```
CREATE VIEW <name> AS
    SELECT ... FROM ... ;
```

- Updateable Views:** Must have only one entry in FROM clause, no GROUP BY / LIMIT etc, no UNION, INTERSECT, no aggregate functions. Otherwise direct modification not possible.

- Triggers useful for (non-updateable) views.  
**(Trigger Timing: (INSTEAD OF) for views.)**

## Trigger Granularity:

- Row-level triggers
  - Trigger function is executed for each affected row
  - Keyword: **FOR EACH ROW**
- Statement-level triggers
  - Trigger function is executed once for each transaction (no matter how many rows are affected)
  - Keyword: **FOR EACH STATEMENT**
  - Ignored return value of trigger function (Enforcing a rollback requires **RAISE EXCEPTION!**)

Example for a statement-level trigger

- Prohibit the deletion of rows from the logs
- Show warning to user only once no matter how many rows the user attempted to delete

```
CREATE TRIGGER on_delete_from_log
BEFORE DELETE ON advanced_log
FOR EACH STATEMENT
EXECUTE FUNCTION show_warning();
```

```
CREATE FUNCTION show_warning()
RETURNS TRIGGER AS
$$
BEGIN
  RAISE EXCEPTION 'Do not DELETE the logs!!!!';
  RETURN NULL;
END;
$$ LANGUAGE plpgsql;
```

## Trigger Conditions

- Execute trigger function only if condition is true.
- Instead of having condition in trigger function, we move to trigger, and only execute if condition is true.
- **Condition:** generally can formulate any boolean expression, but no SELECT, OLD for INSERT, NEW for DELETE, in the WHEN() clause.

```
CREATE TRIGGER on_student_updated_advanced
AFTER UPDATE ON students
FOR EACH ROW WHEN (NEW.points <> OLD.points)
EXECUTE FUNCTION log_student_advanced();
```

## Deferrable Triggers

- Triggers run immediately for every statement that fire them.
- Operations of multiple statements yielding intermediate inconsistent states
- **Deferred triggers:** Run trigger only at end of transactions.

```
CREATE CONSTRAINT TRIGGER on_account_modified
AFTER INSERT OR DELETE OR UPDATE ON accounts
DEFERRABLE INITIALLY IMMEDIATE
FOR EACH ROW
EXECUTE FUNCTION check_balance();
```

- Only work for AFTER and FOR EACH ROW triggers.
- Both CONSTRAINT and DEFERRABLE must be specified.
- INITIALLY DEFERRED: trigger deferred by default.
- INITIALLY IMMEDIATE: trigger not deferred by default (but can be deferred on demand).

```
BEGIN TRANSACTION;
SET CONSTRAINTS on_account_modified DEFERRED;
UPDATE accounts SET balance = balance - 50
  WHERE id = 10;
UPDATE accounts SET balance = balance + 50
  WHERE id = 11;
END TRANSACTION;
```

## Other Trigger Notes

- **Trigger Order:** Triggers for same event on same table:
- Order of activation: BEFORE statement-level, BEFORE row-level, AFTER row-level, AFTER statement-level.
- Within each, fired in alphabetic order. If BEFORE row-level trigger returns NULL, subsequent triggers on same row omitted.

## Trigger Functions

- Trigger functions do not take in (ordinary) arguments.
- Must have return type TRIGGER
- Must be defined before trigger itself.
- “**Input**”: Special internal data structure from trigger.

## Useful Data available in trigger function

- When function is called as a trigger, several special variables created automatically in top level block.

|                  |                                                                         |
|------------------|-------------------------------------------------------------------------|
| <b>TG_NAME</b>   | Name of the trigger that fired                                          |
| <b>TG_OP</b>     | Operation that fired the trigger ( <b>INSERT, UPDATE, DELETE</b> )      |
| <b>TG_WHEN</b>   | Time when the trigger was fired ( <b>BEFORE, AFTER, or INSTEAD OF</b> ) |
| <b>NEW</b>       | Record holding the <u>new</u> row for <b>INSERT/UPDATE</b> operations   |
| <b>OLD</b>       | Record holding the <u>old</u> row for <b>UPDATE/DELETE</b> operations   |
| ...              | ...                                                                     |
| <b>TG_ARGV[]</b> | Array of arguments from the <b>CREATE TRIGGER</b> statement.            |

## Trigger Function Example

```
CREATE OR REPLACE FUNCTION
log_student_advanced()
RETURNS TRIGGER AS $$ BEGIN
  IF TG_OP = 'INSERT' THEN
    INSERT INTO points_log_advanced VALUES
    (NEW.id, TG_OP, NULL, NEW.points,
     DEFAULT);
    RETURN NEW;
  ELSIF (TG_OP = 'DELETE') THEN
    INSERT INTO points_log_advanced VALUES
    (OLD.id, TG_OP, OLD.points, NULL,
     DEFAULT);
    RETURN OLD;
  ELSIF (TG_OP = 'UPDATE') THEN
    INSERT INTO points_log_advanced VALUES
    (OLD.id, TG_OP, OLD.points,
     NEW.points, DEFAULT);
    RETURN NEW;
  END IF;
END; $$ LANGUAGE plpgsql;
```

# 13. Basics of Functional Dependencies

## Chapter Outline

1. Informal Design Guidelines for Relational Databases
2. Functional Dependencies (FDs)
3. 3 Normal Forms based on Primary Keys
4. 4 General Normal Form Definitions for 2NF, 3NF (Multiple Candidate Keys)
5. BCNF (Boyce-Codd Normal Form)

## Goals

- Relational Database Design as a practical activity followed in large organizations worldwide. Relational model dominates the commercial market.
- Have some informal guidelines that can point out problems with relational design.
- Understand theoretical basis for analyzing designs called **functional dependencies**.
- Understand and utilize process of **normalization** to “improve” / “purify” poor designs.
- Understand formal basis for synthesizing good relations strictly based on knowledge of dependencies among attributes.

## Informal Design Guidelines for Relational Databases

- **Relational Database Design:** Grouping of attributes to form “good” relation schemas.
- Two levels of relation schemas’: Logical “user view” level & storage “base relation” level.

## Criteria for “Good” Design

1. **minimality:** should express information with minimum number of distinct relations.
2. **lack of redundancy:** should minimize amount of redundancy among relations.
3. **Information preservation:** should preserve all information captured by the conceptual design (in terms of entity types, relationship types, attributes, and constraints).
4. **consistency:** among the relations
5. **efficiency:** (beyond course scope). Typically addressed in the physical design.

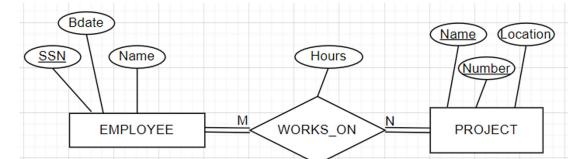
## Guidelines Summary

1. **Guideline 1: Informally, each tuple in a relation should represent one entity or relationship instance. (Applies to individual relations and their attributes)**
  - Bottom line: Design a schema that can be explained easily relation by relation. The semantics of attributes should be easy to interpret.
2. **Guideline 2: Design a schema that does not suffer from update anomalies: (Insertion, Deletion, Modification anomalies).**
  - If there are any anomalies present, then note them so that applications can be made to take them into account. (e.g. introduced for reasons as attributes needed for reporting or accounting purposes).
  - Introduction of anomalies referred to as De-Normalization.
3. **Guideline 3: Relations should be designed such that their tuples will have as few NULL values as possible.**
  - Attributes that are NULL frequently could be placed in separate relations (with the primary key).
  - Reasons for nulls: Attribute not applicable or invalid, Attribute value unknown (may exist), Value known to exist, but unavailable.
4. **Guideline 4: Avoid generation of “spurious data” when tables are joined – an absolute “MUST”.**
  - Bad designs for a relational database may result in erroneous results for certain JOIN operations. Generating bad data cannot be accepted at any cost.
  - The “lossless join” (non-additive) property is used to guarantee that join operation will not create bad data.
  - The relations should be designed to satisfy the lossless join condition. No spurious tuples should be generated by doing a natural-join of any relations.
5. Tool for analysis: **functional dependency**.  
Methodology for “fixing” (bad) designs is specified by process of **“Normalization.”**

## Guideline 1: Don’t mix Relations

- **Guideline 1: Informally, each tuple in a relation should represent one entity or relationship instance. (Applies to individual relations and their attributes)**
  - Do not mix attributes of different entities in same relation.
  - Only use foreign keys to refer to other entities.
  - Entity and relationship attributes be kept apart as much as possible.
- E.g. (EMPLOYEES, DEPARTMENTS, PROJECTS) attributes should not be mixed in the same relation (table).

ER Diagram: (Portion)



Bad Relational Design:

EMP\_PROJ(Emp#, Proj#, Ename, Pname, No\_hours)

Update / Modification Anomaly:

- Changing name of project number P1 from “X” to “Y” may cause update to be made for all 100 employees working on project P1.

Insert Anomaly:

- Cannot insert a project unless an employee is assigned to it. + Converse on inserting employee.

Delete Anomaly:

- When project deleted, result in deleting all employees who work on that project.
- Conversely, if employee sole employee on project, deleting employee -> delete correspond. project

Correct Relational Design:

EMPLOYEE( Ssn, Fname, .., Lname,.., Dno<sub>(FK)</sub>)

PROJECT ( Pnumber, Pname, Plocation, Dnum<sub>(FK)</sub>)

WORKS\_ON (Ssn<sub>(FK)</sub>, Pnumber<sub>(FK)</sub>, No-hours)

Optimization:

- Top 2 relations strictly stand for an entity type
- Third relation strictly stands for a relationship type
- No mixing of entity and relationship information as done in previous example
- **NET RESULT:** Design DOES NOT suffer from any of the anomalies.

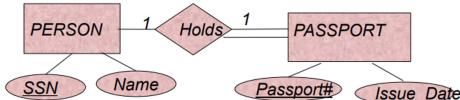
## Guideline 2: Avoid designs with anomalies

- Design a schema that does not suffer from the insertion, deletion and update anomalies.
- If there are any anomalies present, then note them so that applications can be made to take them into account.
- Corollary of guideline 1:** Basically states that when you are forced to mix descriptor attributes of an entity type into the table for another entity type or a relationship type (for performance reasons or reporting requirements etc.), you should document them and take care of the consistency preservation via the application.

## Guideline 3: Minimize Null values in Tuples

- Relations should be designed such that their tuples will have as few NULL values as possible
- Attributes that are NULL frequently could be placed in separate relations (with the primary key)
- Reasons for nulls:**
  - Attribute not applicable or invalid
  - Attribute value unknown (may exist)
  - Value known to exist, but unavailable.

### ER Diagram: (Portion)



Suppose only 60% of persons in the database of 2 million persons hold a Passport

### Bad Relational Design:

PERSON ( SSN, Name, Passport#, Issue\_date )

### Null Values:

- Will contain 40% of tuples with NULL values for Passport# and Issue\_date.

### Correct Relational Design:

PERSON ( SSN, Name ) – 2 million tuples

PASSPORT ( Passport#, Issue\_date, SSN<sub>(FK)</sub> ) – 1.2 million tuples

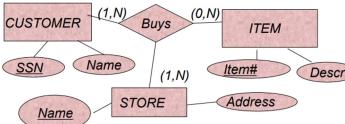
### Optimization:

- No null values

## Guideline 4: Relations satisfy lossless (Non-Additive) join condition to avoid generating spurious tuples.

- Bad designs for relational database: **erroneous results for certain JOIN operations**. The relations should be designed to satisfy the lossless join condition, guarantee meaningful results for join operations.
- No spurious tuples should be generated by doing a natural-join of any relations. Applies to a natural join among any pairs or collections of relations.
- How to know whether a decomposition will be lossless: Use functional dependencies and algorithm (algo 15.3, chapter 15 test losslessness property of an n-way decomposition).

### ER Diagram: (Portion)



Each instance of BUYS relationship type relates one customer, one item, one store. Customer buys at least one to N items, an item may not be sold, but can be sold in any number of BUYS transaction (0, N).

### Correct Relational Design:

BUYS( SSN<sub>(FK)</sub>, Item#<sub>(FK)</sub>, Store\_name<sub>(FK)</sub>, ..... )

Right way to map this ternary relationship is to create one relation for "BUYS".

### Bad Relational Design:

Cust\_Item ( Ssn, Item# )

Cust\_Store( Ssn, Store\_name )

The two resulting tables:

| Cust_Item | Cust_Store |
|-----------|------------|
| s1, x1    | t1         |
| s1, x2    | t2         |
| s2, x1    | t1         |
| s2, x2    | t2         |

Select \* From Cust\_Item x Inner Join Cust\_Store t on x.ssn = t.ssn

| Ssn | Item# | Store_name |
|-----|-------|------------|
| s1  | x1    | t1         |
| s1  | x1    | t2         |
| s1  | x2    | t1         |
| s1  | x2    | t2         |
| s2  | x1    | t2         |
| s2  | x1    | t1         |
| s2  | x2    | t2         |
| s2  | x2    | t1         |

- The **BAD tuples are spurious** (incorrect/invalid) data that resulted from the BAD design
- The ternary relation BUYS (Cust\_ssn, Item#, Store#) was wrongly decomposed into the 2 relations.
- There are two important properties of decompositions:**
  - a) Non-additive or losslessness of the corresponding join
  - b) Preservation of the functional dependencies.

## Functional Dependencies (FDs)

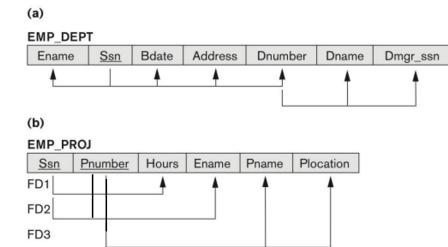
- Used to specify **formal measures** of the "goodness" of relational designs.
- And keys (derived from FDs) are used to define **normal forms** for relations.
- FDs are **constraints** that are derived from the meaning of and interrelationships among the data attributes.
- A set of attributes X **functionally determines** a set of attributes Y if **value of X determines a unique value for Y**.

## Defining Functional Dependencies

- $X \rightarrow Y$  holds if whenever two tuples have the same value for X, they must have the same value for Y.
- For any two tuples  $t_1$  and  $t_2$  in any relation instance:  $r(R)$ : If  $t_1[X] = t_2[X]$ , then  $t_1[Y] = t_2[Y]$
- $X \rightarrow Y$  in R specifies constraint on all relat<sup>n</sup> instances  $r(R)$ .
- Written as  $X \rightarrow Y$ ; can be displayed graphically on a relation schema as in Figures. ( denoted by the arrow).
- FDs are derived from the real-world constraints on the attribute.

### Bad Relational Design:

Figure 14.3  
Two relation schemas subject to update anomalies: (a) EMP\_DEPT and (b) EMP\_PROJ.



### Examples of FD Constraints:

- $SSN \rightarrow Ename$
- $Pnumber \rightarrow \{Pname, Plocation\}$
- $\{SSN, Pnumber\} \rightarrow Hours$

- An FD is a (semantic, logical) property of the attributes in the schema R.
- The constraint must hold on every relation instance  $r(R)$ .
- If K is a key of R, then K functionally determines all attributes in R.
- We never have two distinct tuples with  $t_1[K] = t_2[K]$ , i.e., the projection of each tuple on the K column(s) must yield distinct values.

## Defining FDs from Instances

- To **define** FDs, need to understand meaning of attributes involved and relationship between them.
- Given the instance (population) of a relation, all we can conclude is that an FD **may exist** between certain attributes.
- What we can definitely conclude is – that certain FDs **do not exist** because there are tuples that show a violation of those dependencies.

### Defining Functional Dependencies (from Instance):

- A relation  $R(A, B, C, D)$  with its extension.
- Which FDs **may exist** in this relation?

| A  | B  | C  | D  |
|----|----|----|----|
| a1 | b1 | c1 | d1 |
| a1 | b2 | c2 | d2 |
| a2 | b2 | c2 | d3 |
| a3 | b3 | c4 | d3 |

#### FD Constraints that may exist:

- $\{A, B\} \rightarrow C, \{A, C\} \rightarrow D,$
- $\{A, B\} \rightarrow \{C, D\}$

#### FD that cannot exist (ruled out): $\times$

- $A \rightarrow B, B \rightarrow A, C \rightarrow D, D \rightarrow C,$
- $A \rightarrow \{B, C\}, C \rightarrow \{B, D\}, \{B, C\} \rightarrow A.$

| TEACH | Teacher         | Course   | Text |
|-------|-----------------|----------|------|
| Smith | Data Structures | Bartram  |      |
| Smith | Data Management | Martin   |      |
| Hall  | Compilers       | Hoffman  |      |
| Brown | Data Structures | Horowitz |      |

#### FD Constraints that may exist:

- $\text{Text} \rightarrow \text{Course}$

#### FD that cannot exist (ruled out): $\times$

- $\text{Teacher} \rightarrow \text{Course}, \text{Teacher} \rightarrow \text{Text},$
- $\text{Course} \rightarrow \text{Text}$

## Armstrong Axioms

### Armstrong's Axioms:

#### Armstrong's Axioms

- Reflexivity**  $AB \rightarrow A$
- Augmentation**  $A \rightarrow B \Rightarrow AC \rightarrow BC$
- Transitivity**  $A \rightarrow B \ \& \ B \rightarrow C \Rightarrow A \rightarrow C$

#### 1. Axiom of Reflexivity

- A set of attributes  $\rightarrow$  A subset of the attributes

#### 2. Axiom of Augmentation

- If  $A \rightarrow B$
- Then  $AC \rightarrow BC$  (*for any C*)
- Also,  $AA \rightarrow BA$  which means  $A \rightarrow AB$

#### 3. Axiom of Transitivity

- If  $A \rightarrow B$  and  $B \rightarrow C$
- Then  $A \rightarrow C$

## Further Topics in Functional Dependencies

### Inference Rules for FDs

- Definition:** An FD  $(X \rightarrow Y)$  is inferred from or implied by a set of dependencies F specified on R if  $(X \rightarrow Y)$  **holds in every legal relation state r of R**:

– Whenever r satisfies all the dependencies in F,  $(X \rightarrow Y)$  also holds in r.

- Inference:** Given a set of FDs F, we can infer additional FDs that hold whenever the FDs in F hold.

### Armstrong's Inference Rules

- IR1. (Reflexive):** If Y subset-of X, then  $X \rightarrow Y$
  - IR2. (Augmentation):** If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$ .  
• (Notation:  $XZ$  stands for  $X \cup Z$ )
  - IR3. (Transitive):** If  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z$
- IR1, 2, 3 form a **sound and complete** set of inference rules
- Sound:** Given a set F that holds in R, every dependency that can be inferred using the rules will hold in every state of R.
- Complete:** These rules and all other extended rules that hold can be applied to a set F of dependencies in R until no more dependencies can be inferred.

### Extended Axioms:

#### Extended Axioms

- Decomposition:**  $X \rightarrow YZ \Rightarrow (X \rightarrow Y), (X \rightarrow Z)$
- Union:**  $(X \rightarrow Y), (X \rightarrow Z) \Rightarrow X \rightarrow YZ$
- Pseudotransitivity:**  $(X \rightarrow Y), (WY \rightarrow Z) \Rightarrow WX \rightarrow Z$

#### A. Rule of Decomposition

- If  $A \rightarrow BC$
- Then  $A \rightarrow B$  and  $A \rightarrow C$

#### B. Proof:

- $A \rightarrow BC$  Given
- $BC \rightarrow B$  Reflexivity  $B \subseteq BC$
- Hence,  $A \rightarrow B$  Transitivity (1) and (2)
- $BC \rightarrow C$  Reflexivity  $C \subseteq BC$
- Now,  $A \rightarrow BC$  and  $BC \rightarrow C$ ; Hence,  $A \rightarrow C$ , by Transitivity

#### B. Rule of Union

- If  $A \rightarrow B$  and  $A \rightarrow C$
- Then  $A \rightarrow BC$

#### C. Proof:

- $A \rightarrow B$  Given
- $A \rightarrow C$  Given
- Hence,  $A \rightarrow AB$  Augmentation of (1) with A
- $AB \rightarrow BC$  Augmentation of (2) with B
- Now,  $A \rightarrow AB$  and  $AB \rightarrow BC$ . Thus,  $A \rightarrow BC$

## Closure

### Two types of Closure:

- Closure of a set F of FDs:** is the set  $F^+$  (which is called “Closure of F” or “F closure”) of all FDs that can be inferred from F.
  - Closure of a set of attributes X with respect to F** is the set  $X^+$  of all attributes ( called Closure of X) that are functionally determined by X.
- $X^+$  can be calculated by repeatedly applying IR1, IR2, IR3 using the FDs in F.

### Set of Attributes in closure of X w.r.t F algo[15.1]

#### Algorithm to determine Closure:

- Algorithm 15.1.** Determining  $X^+$ , the Closure of X under F
- Input:** A set F of FDs on a relation schema R, and a set of attributes X, which is a subset of R.

```
X* := X;
repeat
    oldX* := X*;
    for each functional dependency Y  $\rightarrow Z$  in F do
        if  $X^* \supseteq Y$  then  $X^* := X^* \cup Z$ ;
    until  $(X^* = \text{old}X^*)$ ;
```

#### Example:

- Same instructor may offer same course# in an assigned classroom on different days – different sectionids.
- Different instructors may choose different texts for the same course.

```
CLASS (Classid, Course#, Instr_name, Credit_hrs, Text, Publisher, Classroom, Capacity);
```

Let F, (set of functional dependencies for above relation) include the following f.d.s:

- FD1: Classid  $\rightarrow$  Course#, Instr\_name, Credit\_hrs, Text, Publisher, Classroom, Capacity;
- FD2: Course#  $\rightarrow$  Credit\_hrs;
- FD3: (Course#, Instr\_name)  $\rightarrow$  Text, Classroom;
- FD4: Text  $\rightarrow$  Publisher;
- FD5: Classroom  $\rightarrow$  Capacity

These f.d.s above represent meaning of the individual attributes and relationships among them, defines certain rules about the classes.

#### Closure of (sets of) attributes for some example sets:

```
{Classid}+ = {Classid, Course#, Instr_name, Credit_hrs, Text, Publisher, Classroom, Capacity}
{Course#}+ = {Course#, Credit_hrs}
{Course#, Instr_name}+ = {Course#, Instr_name, Credit_hrs, Text, Publisher, Classroom, Capacity}
```

## Equivalence of Sets of FDs

- Two sets of FDs F and G are equivalent if:
  - Every FD in F can be inferred from G, and
  - Every FD in G can be inferred from F
  - Hence, F and G are equivalent if  $F^+ = G^+$ .
- Definition (Covers):**
  - F covers G if every FD in G can be inferred from F.
  - (i.e., if  $G^+$  subset-of  $F^+$ )
- F and G are equivalent if F covers G and G covers F.
- To prove equivalence of two sets of FDs:** e.g.  $F_1 = F_2$ , derive  $F_1$  from  $F_2$ , and derive  $F_2$  from  $F_1$ .

## Minimal Cover of F.D.s

- Apply inference rules to expand on a set F of FDs to arrive at  $F^+$ , its closure.
- Think in opposite direction to reduce set F to its **minimal form** so that the minimal set is still equivalent to the original set F.
- **Definition:** An attribute in a functional dependency (on LHS) is considered extraneous attribute if we can remove it without changing the closure of the set of dependencies.
- Formally, given F, (set of functional dependencies) and a functional dependency  $X \rightarrow A$  in F:
  - Attribute set Y is extraneous in X if:
  - Y is a subset of X, and
  - F logically implies  $(F - (X \rightarrow A)) \cup \{ (X - Y) \rightarrow A \}$
  - (aka replace  $(X \rightarrow A)$  with  $(X - Y) \rightarrow A$  gives back F)

## Minimal Sets of F.D.s

A set of FDs is **minimal** if it satisfies the following conditions:

1. Every dependency in F has a single attribute for its RHS.
2. We cannot replace any dependency  $X \rightarrow A$  in F with a dependency  $Y \rightarrow A$ , where Y is a proper subset-of X and still have a set of dependencies that is equivalent to F.
3. We cannot remove any dependency from F and have a set of dependencies that is equivalent to F.

## Minimal Sets of FDs as Basis for design of relations

- Every set of FDs F has an equivalent minimal set.
- Can be several equivalent minimal sets for given set F of FDs.
- No simple algorithm for computing minimal set of FDs that is equivalent to a set F of FDs. The process of Algorithm 15.2 is used until no further reduction is possible.
- Synthesis approach to design a set of relations, (Lecture 12), starts with all possible F.D.s among a set of attributes that we wish to store, computes their minimal cover and then proceeds to design a set of relations in a specific Normal Form. (Algorithms of Ch.15.)

## Computing Minimal Sets of FDs Algo [15.2]

- **Algorithm 15.2. Finding a Minimal Cover F for a Set of Functional Dependencies E**
  - **Input:** A set of functional dependencies E.
  - 1. Se  $f := E$ .
  - 2. Replace each functional dependency  $X \rightarrow \{A_1, A_2, \dots, A_n\}$  in F by the n functional dependencies  $X \rightarrow A_1, X \rightarrow A_2, \dots, X \rightarrow A_n$ .
  - 3. For each functional dependency  $X \rightarrow A$  in F
    - for each attribute B that is an element of X
      - if  $\{F - \{X \rightarrow A\} \cup \{(X - \{B\}) \rightarrow A\}\}$  is equivalent to F, then replace  $X \rightarrow A$  with  $(X - \{B\}) \rightarrow A$  in F.  
*(\* The above constitutes a removal of the extraneous attribute B from X \*)*
  - 4. For each remaining functional dependency  $X \rightarrow A$  in F if  $\{F - \{X \rightarrow A\}\}$  is equivalent to F, then remove  $X \rightarrow A$  from F.  
*(\* The above constitutes a removal of the redundant dependency  $X \rightarrow A$  from F \*)*

We illustrate algorithm 15.2 with the following:  
Let the given set of FDs be  $E : \{B \rightarrow A, D \rightarrow A, AB \rightarrow D\}$ . We have to find the minimum cover of E.

- All above dependencies are in canonical form; i.e., they have only one attribute on the RHS. So we have completed step 1
- of Algorithm 15.2 and can proceed to step 2. In step 2 we need to determine if  $AB \rightarrow D$  has any redundant attribute on the left-hand side; that is, can it be replaced by  $B \rightarrow D$  or  $A \rightarrow D$ ?
  - Since  $B \rightarrow A$ , by augmenting with B on both sides (IR2), we have  $BB \rightarrow AB$ , or  $B \rightarrow AB$  (i). However,  $AB \rightarrow D$  as given (ii).
  - Hence by the transitive rule (IR3), we get from (i) and (ii),  $B \rightarrow D$ . Hence  $AB \rightarrow D$  may be replaced by  $B \rightarrow D$ .
  - We now have a set equivalent to original E, say  $E' : \{B \rightarrow A, D \rightarrow A, B \rightarrow D\}$ . No further reduction is possible in step 2 since all FDs have a single attribute on the left-hand side.
  - In step 3 we look for a redundant FD in  $E'$ . By using the transitive rule on  $B \rightarrow D$  and  $D \rightarrow A$ , we derive  $B \rightarrow A$ . Hence  $B \rightarrow A$  is redundant in  $E'$  and can be eliminated.
  - Hence the minimum cover of E is  $\{B \rightarrow D, D \rightarrow A\}$ .

### Minimal Basis (Cover)

#### Example

- $F = \{A \rightarrow B, B \rightarrow C, A \rightarrow C\}$        $F_b = \{A \rightarrow B, B \rightarrow C\}$
- Is  $F_b$  a minimal basis for F?
  1.  $A \rightarrow C$  in F can be derived from  $F_b$ 
    - $F_b$  is F by removal of  $A \rightarrow C$
  2. All FDs in  $F_b$  are non-trivial and decomposed
  3. For any FD in  $F_b$ , if we remove an attribute from left hand side, then the FD cannot be derived from  $F_b$   
*(in fact, they have no left hand side!)*
  4. If any FD in  $F_b$  is removed, then some FD in F cannot be derived
- ∴  $F_b$  is a minimal basis (cover) for F

- Conditions
1.  $F_b \equiv F$
  2. Non-trivial and decomposed
  3. No redundant attributes on LHS
  4. No redundant FD

### Minimal Basis (Cover)

#### Example

- $F = \{A \rightarrow B, B \rightarrow C, A \rightarrow C\}$        $F_b = \{A \rightarrow B, AB \rightarrow C\}$
- Is  $F_b$  a minimal basis for F?
  1.  $B \rightarrow C$  in F can NOT be derived from  $F_b$
- ∴  $F_b$  is NOT equivalent to F and hence cannot be a minimal basis for F

- Conditions
1.  $F_b \equiv F$
  2. Non-trivial and decomposed
  3. No redundant attributes on LHS
  4. No redundant FD

## Normalization of Relations (summary)

- **Normalization:** The process of decomposing unsatisfactory "bad" relations by breaking up their attributes into smaller relations by the process of decomposition.
- **Normal form:** Condition using keys and FDs of a relation to certify whether a relation schema is in a particular normal form.
- **2NF, 3NF, BCNF** based on keys and FDs of a relation schema
- **4NF:** based on keys, multi-valued dependencies: MVDs;
- **5NF:** based on keys, join dependencies : JDs;
- Additional properties may be needed to ensure a good relational design (lossless join, dependency preservation; see Chapter 15)

## Designing a Set of Relations

### The Approach of Relational Synthesis (Bottom-up Design):

- Assumes that all possible functional dependencies are known.
- First constructs a minimal set of FDs
- Then applies algorithms that construct a target set of 3NF or BCNF relations.
- Additional criteria may be needed to ensure the the set of relations in a relational database are satisfactory (see Algorithm 15.3).

## Normal Forms Summary

### Normal Forms Defined Informally

- **1<sup>st</sup> normal form**
  - All attributes depend on **the key**
- **2<sup>nd</sup> normal form**
  - All attributes depend on **the whole key**
- **3<sup>rd</sup> normal form**
  - All attributes depend on **nothing but the key**
- **BCNF:**
  - A relation schema R is in Boyce-Codd Normal Form (BCNF) if whenever an FD  $X \rightarrow A$  holds in R, then  $X$  is a **super-key** of R.

## 14. Normalization for Relational DB.

### Properties of Decompositions

There are two important properties of decompositions:

1. Non-additive or losslessness of the corresponding join.
2. Preservation of the functional dependencies.
  - Property (a) is very important, **cannot be sacrificed**.
  - Property (b) is less stringent and may be sacrificed.
  - If the losslessness (non-additivity) of joins is not guaranteed, there will be chaos in terms of lot of spurious data generated by database queries and transactions.

### Important Points

#### Relational Database Design Guidelines:

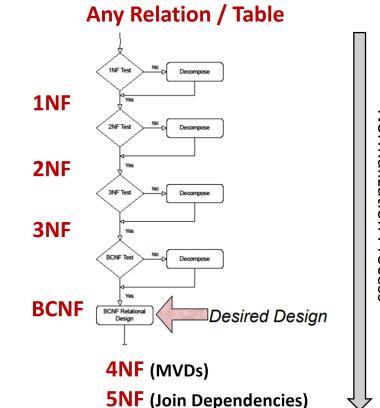
- **Guideline 1:** Informally, each tuple in a relation should represent one entity or relationship instance. (Applies to individual relations and their attributes)
- **Guideline 2:** Design a schema that does not suffer from update anomalies: (Insertion, Deletion, Modification anomalies).
- **Guideline 3:** Relations should be designed such that their tuples will have as few NULL values as possible.
- **Guideline 4:** Avoid generation of “spurious data” when tables are joined – an absolute “MUST”.

#### Functional Dependencies:

- **FDs:** used to specify formal measures of relational designs, are constraints derived from meaning and interrelationships of the data attributes, keys are used to define normal forms for relations.
- **Inference Rules:** Armstrong rules for inferring new FDs.
- **Closure:** Closure of FD as set of all FDs that can be derived from a given set.
- **Cover:** A set of FDs X covers another set Y if all FDs in set Y can be inferred from set X.
- **Equivalence:** Equivalence among two sets of FDs F and G based on F covering G and G covering F.
- **Minimum Cover:** Compute the minimum cover  $F_{min}$  of a set F as an equivalent set with no extraneous attributes on the LHS of any FD in  $F_{min}$  and no redundant FD in  $F_{min}$

### Normalization of Relations

- **Normalization:** Process of decomposing relations by breaking up their attributes into smaller relations. Done to eliminate anomalies and redundancy.
- **Normal form:** Condition using keys and FDs of a relation to certify whether a relation schema is optimal.
- **Practical use:** Normalization carried out in practice so resulting designs are of high quality, meet desirable properties.
- Practical utility of 4NF, 5NF normal forms when constraints are hard to understand or to detect are questionable. Normalize just up to BCNF.
- **Denormalization:** Process of storing join of higher normal form relations as a base relation—which is in a lower normal form.



### Definitions of Keys & Attributes Participating in Keys

- **Superkey** of a relation schema  $R = \{A_1, A_2, \dots, A_n\}$  is set of attributes S subset-of R with property:
  - no two tuples  $t_1$  and  $t_2$  in any legal relation state  $r$  of  $R$  will have  $t_1[S] = t_2[S]$ .
- **A key K** is a **superkey** with additional property that removal of any attribute from K will cause K not to be a superkey any more.
- **Candidate Key:** If relation schema  $>$  one key, each is called a candidate key.
  - One of candidate keys arbitrarily designated as primary key, others called candidate keys (alternate keys).
- **Prime attribute:** member of some candidate key.
- **Nonprime attribute:** not member of any candidate key.

### Example of Keys

#### Definitions of Keys and Attributes Participating in Keys:

Person (NRIC, Email, Name, Birthdate, Income, Postal\_code)

##### Superkey:

- (NRIC, Name)
- (NRIC, Age, Income) etc.

##### Candidate Key / Key (Minimal Superkey):

- NRIC
- Email (If every person requires unique email)

##### In SQL data definition UNIQUE specification can be used to declare candidate keys. E.G.:

- PRIMARY KEY (Ssn)
- UNIQUE (Email)

##### Informal Understanding:

- **Superkey:** Unique tuples
- **Key:** Superkey, where removal of any attribute cause it to not be unique. (minimal superkey)
- **Candidate Key:** When there is more than one key, all are candidates.
- **Prime Attribute:** Part of candidate key
- **Nonprime Attribute:**



### First Normal Form (1NF)

- **1NF:** Table said to be in 1NF if **every attribute possesses only atomic values**.
- **Disallows:** composite attributes, multivalued attributes, nested relations (attributes whose values for an individual tuple are non-atomic).
- **Functional Dependency Analysis:** For a relation to be 1NF, every attribute must functionally be dependent on primary key.
- As long as relation has a key, it is considered to be in 1NF.

#### Normalization into 1NF:

##### Relation Schema not in 1NF:

DEPARTMENT

| Dname          | Dnumber | Dmgr_ssn  | Dlocations                     |
|----------------|---------|-----------|--------------------------------|
| Research       | 5       | 333445555 | (Bellaire, Sugarland, Houston) |
| Administration | 4       | 987654321 | (Stafford)                     |
| Headquarters   | 1       | 888665555 | (Houston)                      |

| Dname          | Dnumber | Dmgr_ssn  | Dlocations                     |
|----------------|---------|-----------|--------------------------------|
| Research       | 5       | 333445555 | (Bellaire, Sugarland, Houston) |
| Administration | 4       | 987654321 | (Stafford)                     |
| Headquarters   | 1       | 888665555 | (Houston)                      |

##### Remedy Decomposition:

DEPARTMENT1 (Dnumber, Dname, Dmgr\_ssn )

DEPARTMENT2 (Dnumber, Dlocation)

##### Functional Dependency argument:

- Dnumber is primary key.
- For a relation to be in 1NF: every attribute must functionally be dependent on the primary key.
- F: (Dnumber  $\rightarrow$  Dname, Dmgr\_ssn)
- But Dnumber  $\not\rightarrow$  Dlocations.

Hence the DEPARTMENT relation does not meet the 1NF requirement.

##### Decomposition:

- It preserves original FDs in F above.
- DEPARTMENT2 has NO FD.

#### Relation Schema not in 1NF:

| Ssn       | Ename                | Pnumber | Hours |
|-----------|----------------------|---------|-------|
| 123456789 | Smith, John B.       | 1       | 32.5  |
| 666884444 | Narayan, Ramesh K.   | 2       | 7.5   |
| 453453453 | English, Joyce A.    | 3       | 4.00  |
|           |                      | 1       | 20.00 |
| 333445555 | Wong, Franklin T.    | 2       | 20.00 |
|           |                      | 3       | 10.00 |
|           |                      | 10      | 10.00 |
| 999887777 | Zelaya, Alicia J.    | 20      | 10.00 |
|           |                      | 30      | 30.00 |
| 987987987 | Jabbar, Ahmad V.     | 10      | 10.00 |
|           |                      | 30      | 35.00 |
| 987654321 | Wallace, Jennifer S. | 30      | 20.00 |
|           |                      | 20      | 15.00 |
| 888665555 | Borg, James E.       | 20      | NULL  |

| Ssn       | Ename                | Pnumber | Hours |
|-----------|----------------------|---------|-------|
| 123456789 | Smith, John B.       | 1       | 32.5  |
| 666884444 | Narayan, Ramesh K.   | 2       | 7.5   |
| 453453453 | English, Joyce A.    | 3       | 4.00  |
|           |                      | 1       | 20.00 |
| 333445555 | Wong, Franklin T.    | 2       | 20.00 |
|           |                      | 3       | 10.00 |
|           |                      | 10      | 10.00 |
| 999887777 | Zelaya, Alicia J.    | 20      | 10.00 |
|           |                      | 30      | 30.00 |
| 987987987 | Jabbar, Ahmad V.     | 10      | 10.00 |
|           |                      | 30      | 35.00 |
| 987654321 | Wallace, Jennifer S. | 30      | 20.00 |
|           |                      | 20      | 15.00 |
| 888665555 | Borg, James E.       | 20      | NULL  |

##### Functional Dependency argument:

- Ssn is primary key.
- F: {Ssn  $\rightarrow$  Ename}
- Ssn  $\not\rightarrow$  PROJ.
- But,
- $(Ssn, Pno) \rightarrow Hours$

##### Remedy Decomposition:

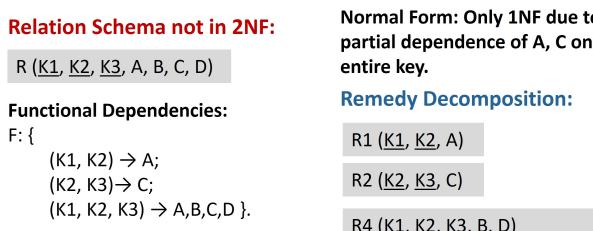
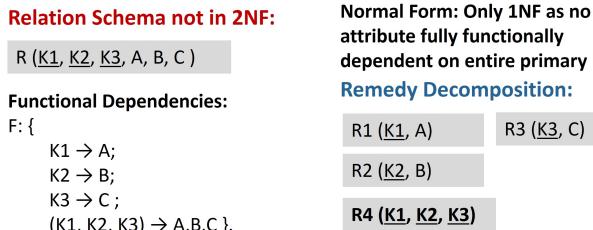
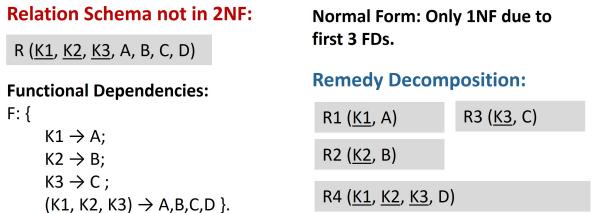
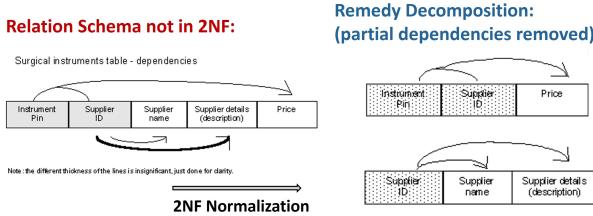
EMP\_PROJ

| Ssn       | Ename                | Proj | Hours |
|-----------|----------------------|------|-------|
| 123456789 | Smith, John B.       | 1    | 32.5  |
| 666884444 | Narayan, Ramesh K.   | 2    | 7.5   |
| 453453453 | English, Joyce A.    | 3    | 4.00  |
|           |                      | 1    | 20.00 |
| 333445555 | Wong, Franklin T.    | 2    | 20.00 |
|           |                      | 3    | 10.00 |
|           |                      | 10   | 10.00 |
| 999887777 | Zelaya, Alicia J.    | 20   | 10.00 |
|           |                      | 30   | 30.00 |
| 987987987 | Jabbar, Ahmad V.     | 10   | 10.00 |
|           |                      | 30   | 35.00 |
| 987654321 | Wallace, Jennifer S. | 30   | 20.00 |
|           |                      | 20   | 15.00 |
| 888665555 | Borg, James E.       | 20   | NULL  |

EMP\_PROJ2

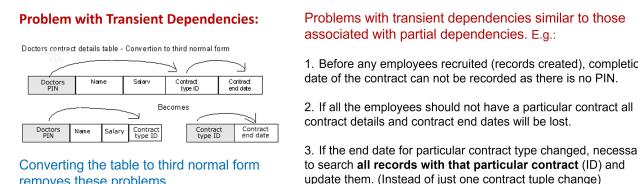
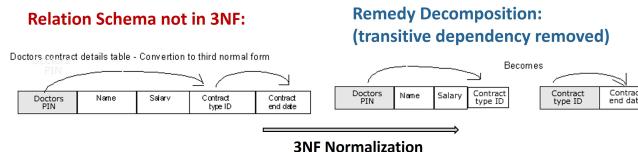
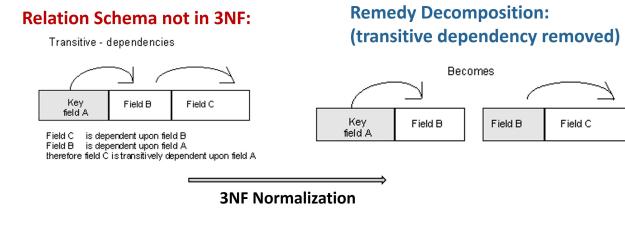
## Second Normal Form (2NF)

- **2NF:** A relation schema R in 2NF if every **non-prime attribute A** in R is **fully functionally dependent** on (every) **ENTIRE primary key**.
- **Understanding:** (don't mix relations).
- R decomposed into 2NF relations via **2NF normalization**: decompose so as to achieve full functional dependence on entire primary key in each (new) relation.
- **Functional Dependency Analysis:** For 2NF, non-prime attribute must be **fully functionally dependent** on primary key.
- By converting to 2NF, **update anomalies are eliminated**.



## Third Normal Form (3NF)

- **3NF:** A relation schema R in 3NF if in 2NF and no non-prime attribute A in R is transitively dependent on the primary key.
  - If whenever a FD (X → A) holds in R, then either:
    - (a) X is a superkey of R, or
    - (b) A is a prime attribute of R. (member of candidate key)
  - **(Every non key field** is non-transitively dependent on the primary key.)
- a) catches 2NF violations due to non-full functional dep, and 3NF violations due to transitive dependencies.
- b) allows certain dependencies up to 3NF (which are disallowed in BCNF)



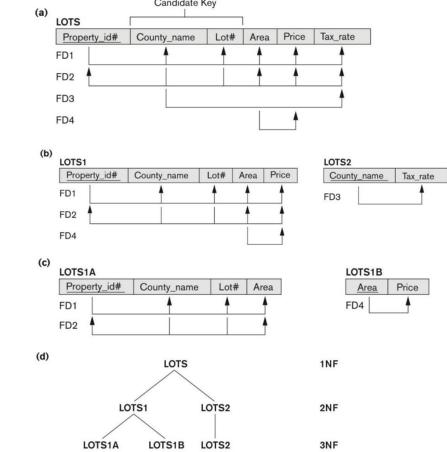
### Note on 3NF: When dependent on Candidate Key, not problematic.

#### NOTE:

- In X → Y and Y → Z, with X as the primary key, we consider this a problematic case only if Y is not a candidate key.
- When Y is a candidate key, there is no problem with the transitive dependency, and it does not cause 3NF violation.
- E.g., Consider EMP (SSN, Emp#, Salary).
  - Here, SSN → Emp# → Salary; and Emp# → SSN. i.e., it is a candidate key.
  - Hence, there is no problematic transitive dependency; So EMP is in 3NF and Salary is an attribute that is directly dependent on the key.

## Successive Normalization into 2NF and 3NF

**Figure 14.12**  
Normalization into 2NF and 3NF. (a) The LOTS relation with its functional dependencies FD1 through FD4. (b) Decomposing into the 2NF relations LOTS1 and LOTS2. (c) Decomposing LOTS1 into the 3NF relations LOTS1A and LOTS1B. (d) Progressive normalization of LOTS into a 3NF design.



## Boyce-Codd Normal Form (BCNF, 3NF+)

- **BCNF:** Relation schema R in BCNF if:
  - If whenever a FD (X → A) holds in R, then:
    - (a) X is a superkey of R.
- Considered strong form of 3NF. Relation NOT in BCNF should be decomposed to achieve BCNF, meet property that every FD has **LHS which must be superkey**.
- However, decomposition possibly forgoes preservation of all functional dependencies in the decomposed relations.
- **“More but not All” violation** method for BCNF:
  1. Derive closure for each attribute subset.
  2. If any closures has “more” attributes on the RHS compared to LHS, but “not all” attributes compared to the table, then it violates BCNF.
  3. Can stop the algorithm and determine that it is not in BCNF

### To check BCNF - More but not all!

- R(A, B, C, D) with FDs AB → C, C → D, and D → A
- 1. Compute the closure of each attribute subset
  - $A^+$  = {A},  $B^+$  = {B},  $C^+$  = {ACD},  $D^+$  = {AD}, ...
- Stop right there...
- Take a look at  $C^+ = \{ACD\}$ 
  - $C^+$  contains more attributes than {C} does
  - $C^+$  does not contain all attributes in R
- “More but not all”, which is a violation of BCNF
- So R is not in BCNF

Possible attribute subsets

$\{A\}^+, \{B\}^+, \{C\}^+, \{D\}^+, \{A,B\}^+, \{A,C\}^+, \{A,D\}^+, \{B,C\}^+, \{B,D\}^+, \{A,B,C\}^+, \{A,B,D\}^+, \{B,C,D\}^+, \{A,B,C,D\}^+$

## BCNF Normalization

relation TEACH in 3NF but not BCNF:

| TEACH   | Student           | Course     | Instructor |
|---------|-------------------|------------|------------|
| Narayan | Database          | Mark       |            |
| Smith   | Database          | Navathe    |            |
| Smith   | Operating Systems | Ammar      |            |
| Smith   | Theory            | Schulman   |            |
| Wallace | Database          | Mark       |            |
| Wallace | Operating Systems | Ahamad     |            |
| Wong    | Database          | Omiecinski |            |
| Zeloya  | Database          | Navathe    |            |
| Narayan | Operating Systems | Ammar      |            |

- Three possible decompositions for relation TEACH
- D1: {student, instructor} and {student, course} ✗
  - D2: {course, instructor} and {course, student} ✗
  - D3: {instructor, course} and {instructor, student} ✓

All three decompositions will lose fd1.

We have to settle for sacrificing the functional dependency preservation. But we cannot sacrifice the losslessness (nonadditivity) property after decomposition

Out of the above three, only the 3rd decomposition D3 will not generate spurious tuples after join.(and hence has the non-additivity property).

TEACH1 (Instructor, Course)

TEACH2 (Instructor, Student)

Functional Dependencies:

- fd1: {student, course} → instructor
- fd2: instructor → course

Not in BCNF due to fd2.  
Redundancy present in relation!

## Testing Binary Decomposition for Lossless Join

- **Binary Decomposition:** Decomposition of a relation R into two relations.
- **NJB Property:** (Non-additive Join test for Binary decompositions).
- A Decomposition  $D = \{R_1, R_2\}$  of R has the **lossless join property** with respect to a set of functional dependencies F on R if and only if either:
  - The f.d.  $((R_1 \cap R_2) \rightarrow (R_1 - R_2))$  is in  $F^+$ , or
  - The f.d.  $((R_1 \cap R_2) \rightarrow (R_2 - R_1))$  is in  $F^+$
- aka  $(R_1 \text{ INTERSECT } R_2) \rightarrow (R_1 - R_2) \text{ or } (R_2 - R_1)$ .

Applying NJB test:

| TEACH   | Student           | Course     | Instructor |
|---------|-------------------|------------|------------|
| Narayan | Database          | Mark       |            |
| Smith   | Database          | Navathe    |            |
| Smith   | Operating Systems | Ammar      |            |
| Smith   | Theory            | Schulman   |            |
| Wallace | Database          | Mark       |            |
| Wallace | Operating Systems | Ahamad     |            |
| Wong    | Database          | Omiecinski |            |
| Zeloya  | Database          | Navathe    |            |
| Narayan | Operating Systems | Ammar      |            |

- Three possible decompositions for relation TEACH
- D1: {student, instructor} and {student, course} ✗
  - D2: {course, instructor} and {course, student} ✗
  - D3: {instructor, course} and {instructor, student} ✓

Applying NJB test:

If you apply the NJB test to the 3 decompositions of the TEACH relation:

- D1 gives  $(\text{Student} \rightarrow \text{Instructor})$  or  $(\text{Student} \rightarrow \text{Course})$ , • none of which is true.
- D2 gives  $(\text{Course} \rightarrow \text{Instructor})$  or  $(\text{Course} \rightarrow \text{Student})$ , • none of which is true.
- However, in D3, we get:  $(\text{Instructor} \rightarrow \text{Course})$  or  $(\text{Instructor} \rightarrow \text{Student})$ .

Since  $(\text{Instructor} \rightarrow \text{Course})$  is indeed true, the NJB property is satisfied and D3 is determined as a non-additive (good) decomposition

## General Procedure for achieving BCNF when a relation fails BCNF

- Let R be the relation not in BCNF.
- Let X be a subset-of R, and let  $X \rightarrow Y$  be the FD that causes BCNF violation.
- Then R may be decomposed into two relations:
  - (i)  $(R - Y)$  and (ii)  $(X \cup Y)$ .
  - If either  $(R - Y)$  or  $(X \cup Y)$  not in BCNF, repeat the process.

## 15. Relational DB. Design Algorithms

Focus: Relational Database Design from bottom up, relational synthesis of relations instead of top down decomposition.

### Design by Analysis (“Top-down”)

- **Normalization** process, decomposing unsatisfactory relations by breaking up their attributes into smaller relations.
- “Improvement” of given design to make it better to eliminate anomalies, redundancy.
- **Normal form:** Metric for condition of relation, using keys and FDs of a relation.

### Design by Synthesis (“Bottom-Up”)

**Relational Synthesis Approach:** Synthesizing 3NF relations from functional dependencies (berNSTein, 1976).

- Assumes that all possible functional dependencies are known:
- Imagine all attributes thrown into a single relation: called **Universal relation**.
- Construct a **minimal set of FDs**, then apply algorithms that construct a target set of **3NF or BCNF relations**.
- Additional criteria may be needed to ensure set of relations in relational database are satisfactory (Algorithms 15.3, 15.4).

### Algorithm to determine key of Relation [15.2a]

**Algo 15.2a:** Finding a key K for R, given a set F of functional dependencies.

**Input:** Universal relation R and a set of functional dependencies F on attributes of R.

1. Set  $K := R$
2. For each attribute A in K:
  - Compute  $(K - A)^+$  with respect to F
  - If closure  $(K - A)^+$  contains all attributes in R, then set  $K := K - \{A\}$ , aka remove from key.
- **Idea:** Assume default key is entire table, remove an attribute at a time, compute cover and see if still equal original functional dependencies. We end up with the minimal set of attributes, and there may be multiple minimal sets (candidate keys).

Example:

- ORDER (order#, order\_date, customer\_id, amount, cust\_phone#)
- 1. Default key: entire relation
- 2. Start dropping attributes until you find combinations that uniquely determine each row in the table.
- Note: application semantics/rules govern what constitutes a key.
- customer\_id, order\_date ? : yes, but only if ....
- cust\_phone#, order\_date ? : yes, but only if ....

Eventually: Order# ?  
Other candidate keys may work if the rules for order processing support those keys

### Properties of Relational Decompositions

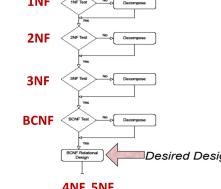
#### Relation Decomposition & Insufficiency of Normal Forms:

- **Universal Relation Schema:** A relation schema  $U = \{A_1, A_2, \dots, A_n\}$  that includes all the attributes of the database.
- **Universal relation assume:** All attribute names unique.
- **Decomposition:** Process of decomposing universal relation schema  $R$  into a set of relation schemas  $D = \{R_1, R_2, \dots, R_m\}$  that will become the relational database schema by using the functional dependencies.
- **Goal:** Each individual relation  $R_i$  in decomposition D be in BCNF/3NF, prevent generating spurious tuples.
- **Attribute preservation condition:** Each attribute in  $R$  will appear in at least one relation schema  $R_i$  in the decomposition so that no attributes are “lost”.
- **Dependency preservation property:** A decomposition  $D = \{R_1, R_2, \dots, R_m\}$  of R is dependency-preserving w.r.t. F if union of projections of F on each  $R_i$  in D is equivalent to F:

$$([\pi_{R1}(F)] \cup \dots \cup [\pi_{Rm}(F)])^+ = F^+$$

#### Top-Down Design Process: (Design by Analysis)

Any Relation / Table



Normalization Process

#### Bottom-Up Design Process: (Design by Synthesis)

Any Relation / Table



- **Dependency preservation property:** Given  $F$  on  $R$ , projection of  $F$  on  $R_i$ ,  $\pi_{R_i}(F)$ , where  $R_i$  is subset of  $R$ , is the set of dependencies  $X \rightarrow Y$  in  $F^+$  such that the attributes in  $X \cup Y$  are all contained in  $R_i$ .

- Hence, the projection of  $F$  on each relation schema  $R_i$  in decomposition  $D$  is set of FDs in  $F^+$  (closure of  $F$ ), where all left- and right-hand-side attributes are in  $R_i$ .

### • Non-additive (Lossless) Join Property of Decomposition:

Decomposition  $D = \{R_1, R_2, \dots, R_m\}$  of  $R$  has property w.r.t. set of dependencies  $F$  on  $R$  if, for every relation state  $r$  of  $R$  that satisfies  $F$ , the following holds, where  $*$  is natural join of all the relations in  $D$ :

$$*(\pi_{R_1}(r), \dots, (\pi_{R_m}(r)) = r$$

- Word loss in lossless refers to loss of information, not tuples. In fact, for “loss of information” a better term is “addition of spurious information”.

## Test Non-Additive Join Property for Binary D.

- See NJB test on previous page. Basically:

$$(R_1 \text{ INTERSECT } R_2) \rightarrow (R_1 - R_2) \text{ or } (R_2 - R_1)$$

## Algorithm[15.3] Test Non-Additive Join Property for N-ary Decomposition

### Lossless (Non-additive) Join Property of a Decomposition :

#### ■ Algorithm 15.3: Testing for Lossless Join Property

- **Input:** A universal relation  $R$ , a decomposition  $D = \{R_1, R_2, \dots, R_m\}$  of  $R$ , and a set  $F$  of functional dependencies.

1. Create an initial matrix  $S$  with one row  $i$  for each relation  $R_i$  in  $D$ , and one column  $j$  for each attribute  $A_j$  in  $R$ .
2. Set  $S(i,j) := b_{ij}$  for all matrix entries. (\* each  $b_{ij}$  is a distinct symbol associated with indices  $(i, j)$  \*).
3. For each row  $i$  representing relation schema  $R_i$  (for each column  $j$  representing attribute  $A_j$    
 {if (relation  $R_i$  includes attribute  $A_j$ ) then set  $S(i,j) := a_j$ ;});   
 (\* each  $a_j$  is a distinct symbol associated with index  $(j)$  \*)
4. Repeat the following loop until a complete loop execution results in no changes to  $S$  (for each functional dependency  $X \rightarrow Y$  in  $F$    
 {for all rows in  $S$  which have the same symbols in the columns corresponding to attributes in  $X$    
 {make the symbols in each column that correspond to an attribute in  $Y$  be the same in all these rows as follows:   
 If any of the rows has an “a” symbol for the column, set the other rows to that same “a” symbol in the column.   
 If no “a” symbol exists for the attribute in any of the rows, choose one of the “b” symbols that appear in one of the rows for the attribute and set the other rows to that same “b” symbol in the column };});   
 );
5. If a row is made up entirely of “a” symbols, then the decomposition has the lossless join property; otherwise it does not.

## N-ary Non-Additive Test Example

Figure 15.1 Nonadditive join test for n-ary decompositions.  
(a) Case 1: Decomposition of EMP\_PROJ into EMP\_PROJ1 and EMP\_LOCS fails test. Hence, **not a good decomposition**.  
(b) A decomposition of EMP\_PROJ that has the lossless join property.

(a)  $R = \{\text{Ssn, Ename, Pnumber, Pname, Plocation, Hours}\}$   
 $R_1 = \text{EMP\_LOCS} = \{\text{Ename, Plocation}\}$   
 $R_2 = \text{EMP\_PROJ1} = \{\text{Ssn, Pnumber, Hours, Pname, Plocation}\}$

$$D = \{R_1, R_2\}$$

$$F = \{\text{Ssn} \rightarrow \text{Ename}; \text{Pnumber} \rightarrow \{\text{Pname, Plocation}\}; \{\text{Ssn, Pnumber}\} \rightarrow \{\text{Hours}\}\}$$

| $R_1$ | $Ssn$    | $Ename$  | $Pnumber$ | $Pname$  | $Plocation$ | $Hours$  |
|-------|----------|----------|-----------|----------|-------------|----------|
| $R_1$ | $b_{11}$ | $a_2$    | $b_{13}$  | $b_{14}$ | $a_5$       | $b_{16}$ |
| $R_2$ | $a_1$    | $b_{22}$ | $a_3$     | $a_4$    | $a_5$       | $a_6$    |

(No changes to matrix after applying functional dependencies)

Now consider a “proper” decomposition:

(b)  $\text{EMP} = \{\text{Ssn, Ename}\}$     $\text{PROJECT} = \{\text{Pnumber, Pname, Plocation}\}$     $\text{WORKS\_ON} = \{\text{Ssn, Pnumber, Hours}\}$

$$D = \{R_1, R_2, R_3\}$$

(Case 2: Decomposition of EMP\_PROJ into EMP, PROJECT, and WORKS\_ON satisfies test.)

(c)  $R = \{\text{Ssn, Ename, Pnumber, Pname, Plocation, Hours}\}$   
 $R_1 = \text{EMP} = \{\text{Ssn, Ename}\}$   
 $R_2 = \text{PROJ} = \{\text{Pnumber, Pname, Plocation}\}$   
 $R_3 = \text{WORKS\_ON} = \{\text{Ssn, Pnumber, Hours}\}$

$$D = \{R_1, R_2, R_3\}$$

$$F = \{\text{Ssn} \rightarrow \text{Ename}; \text{Pnumber} \rightarrow \{\text{Pname, Plocation}\}; \{\text{Ssn, Pnumber}\} \rightarrow \{\text{Hours}\}\}$$

| $R_1$ | $Ssn$    | $Ename$  | $Pnumber$ | $Pname$  | $Plocation$ | $Hours$  |
|-------|----------|----------|-----------|----------|-------------|----------|
| $R_1$ | $a_1$    | $a_2$    | $b_{13}$  | $b_{14}$ | $b_{15}$    | $b_{16}$ |
| $R_2$ | $b_{21}$ | $b_{22}$ | $a_3$     | $a_4$    | $a_5$       | $b_{26}$ |
| $R_3$ | $a_1$    | $b_{32}$ | $a_3$     | $b_{34}$ | $b_{35}$    | $a_6$    |

(Original matrix S at start of algorithm)

| $R_1$ | $Ssn$    | $Ename$  | $Pnumber$ | $Pname$  | $Plocation$ | $Hours$  |
|-------|----------|----------|-----------|----------|-------------|----------|
| $R_1$ | $a_1$    | $a_2$    | $b_{13}$  | $b_{14}$ | $b_{15}$    | $b_{16}$ |
| $R_2$ | $b_{21}$ | $b_{22}$ | $a_3$     | $a_4$    | $a_5$       | $b_{26}$ |
| $R_3$ | $a_1$    | $b_{32}$ | $a_3$     | $b_{34}$ | $a_4$       | $b_{35}$ |

(Matrix S after applying the first two functional dependencies; last row is all “a” symbols so we stop)

(1) Because  $\text{Ssn} \rightarrow \text{Ename}$  and both these have ‘a’s in row1, the  $a_1$  in row3 causes  $\text{Ename}$  to be set to  $a_2$ .

(2) Because  $\text{Pnumber} \rightarrow \{\text{Pname, Plocation}\}$  and all these columns have an ‘a’ in row2, the  $a_3$  in row3 causes  $a_4, a_5$  to be set in row3.

## Relational Decomposition

See textbooks for proof.

- **Claim 1 (Dependency Preserving):** Always possible to find a dependency preserving decomposition  $D$  with respect to  $F$  such that each relation  $R_i$  in  $D$  is in 3NF.

- **Claim 2 (Preservation of non-additivity in successive decompositions):** If  $D = \{R_1, R_2, \dots, R_m\}$  of  $R$  has lossless join property w.r.t. set of FDs  $F$  on  $R$ , and if decomposition  $D_i = \{Q_1, Q_2, \dots, Q_k\}$  of  $R_i$  has lossless join property w.r.t. projection of  $F$  on  $R_i$ , then decomposition  $D_2 = \{R_1, R_2, \dots, R_{i-1}, Q_1, Q_2, \dots, Q_k, R_{i+1}, \dots, R_m\}$  of  $R$  has the non-additive join property w.r.t.  $F$ .

## Relational Synthesis Algorithm[15.4] (to 3NF)

### ■ Design of 3NF Schemas:

Algorithm 15.4 Relational Synthesis into 3NF with Dependency Preservation and Non-Additive (Lossless) Join Property

- Input: A universal relation  $R$  and a set of functional dependencies  $F$  on the attributes of  $R$ .

1. Find a minimal cover  $G$  for  $F$  (use Algorithm 15.2).

2. For each left-hand-side  $X$  of a functional dependency that appears in  $G$ ,

create a relation schema in  $D$  with attributes  $\{X \cup \{A_1\} \cup \{A_2\} \dots \cup \{A_k\}\}$ ,

where  $X \rightarrow A_1, X \rightarrow A_2, \dots, X \rightarrow A_k$  are the only dependencies in  $G$  with  $X$  as left-hand-side ( $X$  is the key of this relation).

3. If none of the relation schemas in  $D$  contains a key of  $R$ , then create one more relation schema in  $D$  that contains attributes that form a key of  $R$ . (Use Algorithm 15.4a to find the key of  $R$ )

4. Eliminate redundant relations from the resulting set of relations in the relational database schema. A relation  $R$  is considered redundant if  $R$  is a projection of another relation  $S$  in the schema; alternately,  $R$  is subsumed by  $S$ .

## Algo[15.4] Relational Synthesis Example

### ■ Example 1 of Algorithm 15.4.

Consider the following universal relation:  
 $U = \{\text{Emp\_ssn, Pno, Esal, Ephone, Dno, Pname, Plocation}\}$

- The following dependencies are present:

- FD1:  $\text{Emp\_ssn} \rightarrow \{\text{Esal, Ephone, Dno}\}$
- FD2:  $\text{Pno} \rightarrow \{\text{Pname, Plocation}\}$
- FD3:  $\text{Emp\_ssn, Pno} \rightarrow \{\text{Esal, Ephone, Dno, Pname, Plocation}\}$
- By virtue of FD3, the attribute set  $\{\text{Emp\_ssn, Pno}\}$  represents a key of the universal relation.
- Hence  $\mathcal{F}$ , the set of given FDs, includes  $\{\text{Emp\_ssn} \rightarrow \text{Esal, Ephone, Dno}; \text{Pno} \rightarrow \text{Pname, Plocation}; \text{Emp\_ssn, Pno} \rightarrow \text{Esal, Ephone, Dno, Pname, Plocation}\}$ .

- To determine min cover of  $\mathcal{F}$ , we note that FD3 can be expressed as (In Step2 of Algo 15.2)

$\text{Emp\_ssn, Pno} \rightarrow \text{Esal}$

$\text{Emp\_ssn, Pno} \rightarrow \text{Ephone}$

$\text{Emp\_ssn, Pno} \rightarrow \text{Dno}$

$\text{Emp\_ssn, Pno} \rightarrow \text{Pname}$

$\text{Emp\_ssn, Pno} \rightarrow \text{Plocation}$

In Step 3 of Algo 15.2 - the first three FDs above,  $\text{Pno}$  is extraneous and in the last two,  $\text{Emp\_ssn}$  is extraneous.

- Algorithm 15.4: second step produces relations  $R1$  and  $R2$  as:

$R1 = \{\text{Emp\_ssn, Esal, Ephone, Dno}\}$

$R2 = \{\text{Pno, Pname, Plocation}\}$

- Algorithm 15.4: third step - we generate a relation corresponding to the key  $\{\text{Emp\_ssn, Pno}\}$  of  $U$

- Hence, the resulting design contains:

$R1 = \{\text{Emp\_ssn, Esal, Ephone, Dno}\}$

$R2 = \{\text{Pno, Pname, Plocation}\}$

$R3 = \{\text{Emp\_ssn, Pno}\}$

We can easily see that the final design meets 3NF

## Algo[15.4] Relational Synthesis Example 2

### ■ Example 2 of Algorithm 15.4.

- Consider the relation schema LOTS1A
- Assume that this relation is given as a universal relation :  
 $U(\text{Property\_id}, \text{County}, \text{Lot\#}, \text{Area})$  with the following functional dependencies:

- FD1:  $\text{Property\_id} \rightarrow \text{Lot\#}, \text{County}, \text{Area}$
- FD2:  $\text{Lot\#}, \text{County} \rightarrow \text{Area}, \text{Property\_id}$
- FD3:  $\text{Area} \rightarrow \text{County}$

Represent the functional dependencies as the set in an abbreviated form  
 $F: \{P \rightarrow LCA, LC \rightarrow AP, A \rightarrow C\}$

If we apply the minimal cover Algorithm 15.2 to F, (in step 2) we first represent the set F as  
 $F: \{P \rightarrow L, P \rightarrow C, P \rightarrow A, LC \rightarrow A, LC \rightarrow P, A \rightarrow C\}$

In the set F: ,  $P \rightarrow A$  can be inferred from  $P \rightarrow LC$  and  $LC \rightarrow A$ ; hence  $P \rightarrow A$  by transitivity and is therefore redundant. Hence, Algorithm 15.2 (in step 4) removes this redundant FD from the set F:

- Thus, one possible minimal cover is
- Minimal cover FX:  $\{P \rightarrow LC, LC \rightarrow A, LC \rightarrow P, A \rightarrow C\}$
- Algorithm 15.4 - step 2 will produce design X using the above minimal cover FX as

**Design X:** R1 (P, L, C), R2 (L, C, A, P), and R3 (A, C)

Now step 4 of Algorithm 15.4 applies: To reiterate:

**Step 4 of Algorithm 15.4: Eliminate redundant relations from the resulting set of relations in the relational database schema. A relation R is considered redundant if R is a projection of another relation S in the schema; alternately, R is subsumed by S.**

- we find that R3 is subsumed by R2
- we find that R1 is also subsumed by R2
- Hence both of those relations R1 and R3 are redundant. Thus the 3NF schema that achieves both of the desirable properties is (after removing redundant relations), is:

**FINAL 3NF Design X:** R2 (L, C, A, P) which is same as the universal relation we started with; in other words it is identical to the relation

- LOTS1A (Property\_id, Lot#, County, Area) that we had determined to be in 3NF in Section 14.4.2.

**FINAL 3NF Design X:** R2 (L, C, A, P) which is same as the universal relation we started with.

- In other words it is identical to the relation

LOTS1A (Property\_id, Lot#, County, Area) that we had determined to be in 3NF in Section 14.4.2.

- Note: In the textbook we discuss an alternate min cover of F that leads to an alternate 3NF Design Y:

**Design Y:** S1 (P, A, L), S2 (L, C, P), and S3 (A, C)

(see pages 521-522 in the book for further details on how we arrived at this alternate 3NF design)

## Relational Synthesis Algorithm[15.5] (to BCNF)

### Cannot guarantee functional dependency preservation.

#### Design of BCNF Schemas

**Algorithm 15.5: Relational Decomposition into BCNF with Lossless (non-additive) join property**

- Input: A universal relation R and a set of functional dependencies F on the attributes of R.

1. Set D := {R};
2. While there is a relation schema Q in D that is not in BCNF do {

choose a relation schema Q in D that is not in BCNF;  
 find a functional dependency  $X \rightarrow Y$  in Q that violates BCNF;  
 replace Q in D by two relation schemas (Q - Y) and (X  $\cup$  Y);

}

**Assumption:** No null values are allowed for the join attributes.

## Algo[15.5] Relational Synthesis Example

### ■ Design of BCNF Schemas

Consider a simple relation for a High class restaurant:

CUST\_TABLE (Cust#, Table#, Date, Waiter#, Bill\_amount)

The attributes are self-explanatory

The FDs are:

Fd1:  $(\text{Cust\#}, \text{Table\#}, \text{Date}) \rightarrow \text{Waiter\#}, \text{Bill\_amount}$

Fd2 :  $\text{Waiter\#} \rightarrow \text{Table\#}$  (a fancy restaurant where a waiter waits on a single table!)

The relation CUST\_TABLE is in 3NF because if you apply the generalized definition of 3NF, Fd1 has LHS as superkey and Fd2 has RHS which is a prime attribute. However Fd 2 violates BCNF. Using Algorithm 15.5, ...  $X \rightarrow Y$  violates BCNF ..... where X is Waiter# and Y is Table#.

Hence the BCNF design is:

CUST\_TABLE 1 (Cust#, Date, Waiter#, Bill\_amount) and

WAITER (Waiter#, Table#)

Note that this meets non-additive decomposition property but loses Fd1.

## Notes on Normalization Algorithms

### Problems:

- Database designer must first specify all relevant functional dependencies among the database attributes.
- Algorithms are not deterministic in general.
- Not always possible to find a decomposition into relation schemas that preserves dependencies and allows each relation schema in the decomposition to be in BCNF (instead of 3NF as in Algorithm 15.5).

## Summary: Relational DB Schema Design Algos

Table 15.1 Summary of the Algorithms Discussed in This Chapter

| Algorithm | Input                                                         | Output                                                    | Properties/Purpose                                                      | Remarks                                                                                     |
|-----------|---------------------------------------------------------------|-----------------------------------------------------------|-------------------------------------------------------------------------|---------------------------------------------------------------------------------------------|
| 15.1      | An attribute or a set of attributes X, and a set of FDs F     | A set of attributes in the closure of X with respect to F | Determine all the attributes that can be functionally determined from X | The closure of a key is the entire relation                                                 |
| 15.2      | A set of functional dependencies F                            | The minimal cover of functional dependencies              | To determine the minimal cover of a set of dependencies F               | Multiple minimal covers may exist—depends on the order of selecting functional dependencies |
| 15.2a     | Relation schema R with a set of functional dependencies F     | Key K of R                                                | To find a key K (that is a subset of R)                                 | The entire relation R is always a default superkey                                          |
| 15.3      | A decomposition D of R and a set F of functional dependencies | Boolean result: yes or no for nonadditive join property   | Testing for nonadditive join decomposition                              | See a simpler test NJB in Section 14.5 for binary decompositions                            |
| 15.4      | A relation R and a set of functional dependencies F           | A set of relations in 3NF                                 | Nonadditive join and dependency-preserving decomposition                | May not achieve BCNF, but achieves all desirable properties and 3NF                         |
| 15.5      | A relation R and a set of functional dependencies F           | A set of relations in BCNF                                | Nonadditive join decomposition                                          | No guarantee of dependency preservation                                                     |

## Problems with Null Values and Dangling Tuples

• **Problems with NULL values:** When some tuples have NULL values for attributes that will be used to join individual relations in the decomposition that may lead to incomplete results.

• E.g. Joining on an attribute with NULL values may leave those tuples out in the joint result.

• In some cases, LEFT /RIGHT OUTER JOIN may be used to include these values.

• **Problems with Dangling Tuples:** When decomposition does not include NULL values (for attributes which allow NULL values), it may cause certain tuples to not appear in the result of a (natural) JOIN. These tuples are called dangling tuples.

## Some Other Dependencies in R. DBs

• Besides FDs, MVDs (multivalued dependencies) and JDs (join dependencies) used for defining normal forms 1NF upto 5NF, some other dependencies have been proposed.

• **Inclusion Dependencies:** Formalize two types of interrelational constraints which cannot be expressed using F.D.s or MVDs: **Referential integrity constraints, Class/subclass relationships.**

• Basically, represents dependencies between attributes, express the way foreign keys relates to the primary keys.

#### Definition:

■ An inclusion dependency  $R.X < S.Y$  between two sets of attributes : X of relation schema R, and Y of relation schema S, specifies the constraint that, at any specific time when r is a relation state of R and s a relation state of S, we must have

$$\pi_X(r(R)) \subseteq \pi_Y(s(S))$$

#### Note:

- The  $\subseteq$  (subset) relationship does not necessarily have to be a proper subset.
- The sets of attributes on which the inclusion dependency is specified—X of R and Y of S—must have the same number of attributes.
- In addition, the domains for each pair of corresponding attributes should be compatible.

• **Arithmetic Functions:** As long as a unique value of Y is associated with every X, we may consider that FD  $X \rightarrow Y$  exists.

• There may be some procedure that takes into account the arithmetic function functional dependency.

• The above dependencies are relevant during insertion/loading of data or query processing, but NOT relevant to normalization of the relation.

-end-