

## XXI. Database Design

Databases and DBMS  
 Data Models, Hierarchical, Network, Relational  
 Database Design  
 Restructuring an ER schema  
 Performance analysis  
 Analysis of Redundancies, Removing generalizations  
 Translation into a Relational Schema  
 The Training Company Example  
 Normal Forms and Normalization of Relational Schemata



## Databases

- A **database** is a collection of persistent data shared by a number of applications
- Databases have been founded on the concept of **data independence**: Applications should not have to know the organization of the data or the access strategy employed  
Need query processing facility, which generates automatically an access plan, given a query
- Databases also founded on the concept of **data sharing**: Applications should be able to work on the same data concurrently, without knowing of each others' existence.

Database procedures defined in terms of atomic operations called transactions

## Conventional Files vs Databases

### Files

Advantages -- many already exist; good for simple applications; very efficient

Disadvantages -- data duplication; hard to evolve; hard to build for complex applications

### Databases

Advantages -- Good for data integration; allow for more flexible formats (not just records)

Disadvantages -- high cost; drawbacks in a centralized facility

*The future is with databases!*

## Database Concepts

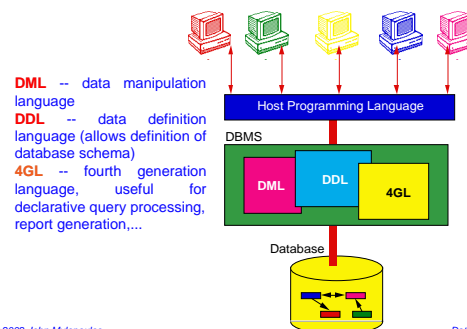
- **Data model** -- defines a set of data structures along with associated operations, for building and accessing a database  
 ✓ e.g., the relational model offers relations (tables) as data structure for building a database
- **Database management system (DBMS)** -- generic tool for building, accessing, updating and managing a database  
 ✓ E.g., Oracle, DB2, Access,... are all relational DBMSs
- **Database schema** -- describes the types and structure of the data stored in the database.  
 ✓ E.g., Employee(emp#, name, addr, sal, dept, mgr)
- **Transaction** -- an atomic operation on a database; looks like a procedure but has different semantics: when called, it either completes its execution, or aborts and undoes all changes it made to the database.  
 ✓ E.g., TransferFunds(fromAcct#, toAcct#, amount, date)

## Types of Databases

- **Conventional databases** -- (relational, network, hierarchical) consist of records of many different record types (database looks like a collection of files)
- **Object-Oriented databases** -- database consists of objects (and possibly associated programs); database schema consists of classes (which can be objects too).
- **Multimedia databases** -- database can store formatted data (i.e., records) but also text, pictures,...
- **Active databases** -- database includes event-condition-action rules
- **Deductive databases\*** -- like large Prolog programs
- **Hypertext databases** -- store and access efficiently HTML/XML documents; provide navigational facilities through a database, so that a user can retrieve and/or browse.

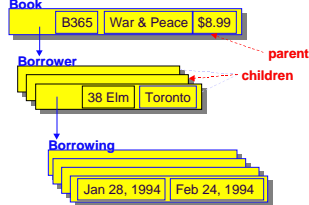
\* -- not available commercially

## Database Management Systems



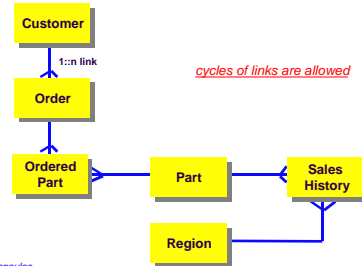
## The Hierarchical Data Model

- Database consists of **hierarchical record structures**; a field may have as value a list of records; every record has at most one parent



## The Network Data Model

- A database now consists of records with pointers (links) to other records. Offers a navigational view of a database.



## The Relational Data Model

- A database now consists of sets of records or (equivalently) sets of tuples (relations) or (equivalently) tables of tuples; no links allowed in the database.
- Every tuple is an element of exactly one relation and is identified uniquely by a primary key

Customer			Order			Ordered Part		
Cust#	Name	Address	Ord#	Date	Amount	Part#	Ord#	Quantity
1127	George	25 Mars St.	4237	11/2/93	\$65.87	2397	1997	980
1377	Maria	12 Love Ave	2908	6/5/93	\$126.88	2908	1997	100 doz.
1532	Manolis	1 Bloss St.	1552	12/1/93	\$284.21	6590	4237	40 doz.
...			...			...		

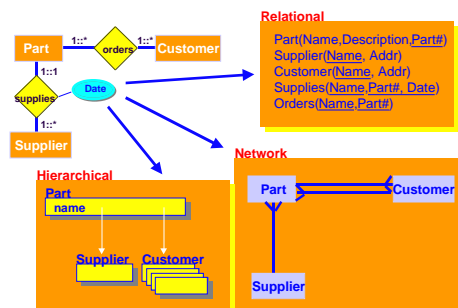
  

Part		
Part#	Desc	Quantity
2397	widget	12,960
2908	nut	16,000 doz.
6590	bolt	14,340 doz.
...		

## Comparing Data Models

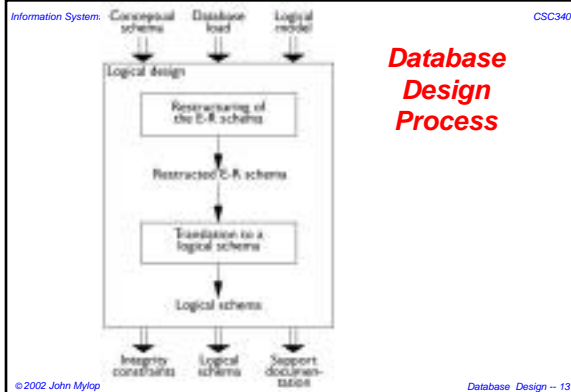
- The oldest DBMSs were hierarchical, dating back to the mid-60s. IMS (IBM product) is the most popular among them. Many old databases are hierarchical.
- The network data model came next (early '70s). At the time of its proposal, it was viewed as a breakthrough. It emphasized the role of the database programmer as "navigator", chasing links (pointers, actually) around a database.
- But, the network model was found to be in many respects too implementation-oriented, not insulating sufficiently the programmer from implementation features of network DBMSs.
- The relational model is the most recent arrival (early '80s) and it has taken over the database market. Relational databases are considered simpler than their hierarchical and network cousins because they don't allow any links/pointers (which are necessarily implementation-dependent).

## Designing a Database Schema



## (Relational) Database Design

- The aim of database design is to construct a relational schema that correctly and efficiently represents all of the information described by a class or Entity-Relationship diagram (or 'schema') produced during requirements analysis.
- From now, we'll only talk about transforming an E-R schema into a relational schema. Most of the transformation process applies for class diagrams as well.
- This is *not* just a simple translation from one model to another for two main reasons:
  - ✓ not all the constructs of the Entity-Relationship model can be translated naturally into the relational model;
  - ✓ the schema must be restructured in such a way as to make the execution of the projected operations as efficient as possible.



Information Systems Analysis and Design CSC340

**Logical Design Steps**

It is helpful to divide the logical design into two steps:

- **Restructuring of the Entity-Relationship schema**, based on criteria for the optimization of the schema and the simplification of the following step;
- **Translation into the logical model**, based on the features of the logical model (in our case, the relational model).

© 2002 John Mylopoulos Database Design -- 14

Information Systems Analysis and Design CSC340

**Performance Analysis**

- An ER schema can be restructured to optimize two parameters:
  - ✓ **Cost of an operation** (evaluated in terms of the number of occurrences of entities and relationships that are visited during the execution of an operation on the database);
  - ✓ **Storage requirements** (evaluated in terms of number of bytes necessary to store the data described by the schema).
- In order to study these parameters, we need to know:
  - ✓ Projected volume of data;
  - ✓ Projected operation characteristics.

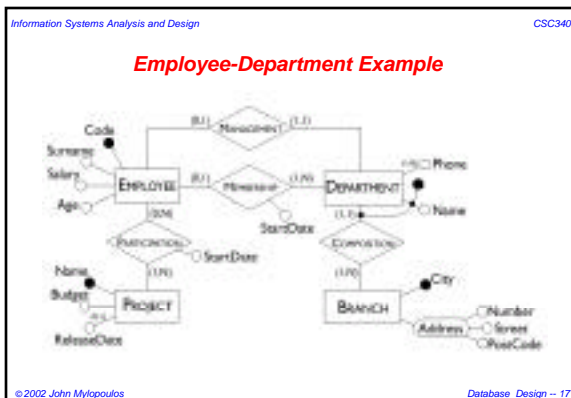
© 2002 John Mylopoulos Database Design -- 15

Information Systems Analysis and Design CSC340

**Cost Model**

- The cost of an operation is measured in terms of the number of disk accesses required. A disk access is, generally orders of magnitude more expensive than in-memory accesses, or CPU operations.
- For a coarse estimate of cost, we will assume that
  - ✓ a Read operation requires 1 disk access
  - ✓ A Write operation requires 2 disk accesses (read from disk, change, write back to disk)

© 2002 John Mylopoulos Database Design -- 16



Information Systems Analysis and Design CSC340

**Typical Operations**

- Operation 1: Assign an employee to a project.
- Operation 2: Find the record of an employee, including the department where she works, and the projects she works for.
- Operation 3: Find the records of all employees for a given department.
- Operation 4: For each branch, retrieve its departments, and for each department, retrieve the last names of their managers, and the list of their employees.
- **Note:** For class diagrams, these would be operations associated with database classes.

© 2002 John Mylopoulos Database Design -- 18

### Tables of Volumes and Operations

The volume of data and the general characteristics of the operations can be summed up using two special tables.

**Table of volumes**

Concept	Type	Volume
Branch	E	10
Department	E	80
Employee	E	2000
Project	E	500
Composition	R	80
Membership	R	1900
Management	R	80
Participation	R	6000

**Table of operations**

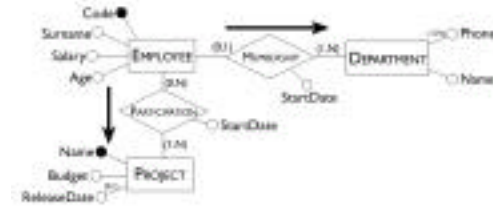
Operation	Type	Frequency
Operation 1	I	50 per day
Operation 2	I	100 per day
Operation 3	I	10 per day
Operation 4	B	2 per day

I - Interactive

B - Batch

### Navigation Schema for Operation 2

A navigation schema starts from the inputs to an operation and moves (via arrows) towards its outputs.



### Table of Accesses

This table evaluates the cost of an operation, using the table of volumes and the navigation schema.

Concept	Type	Accesses	Type
Employee	Entity	1	R
Membership	Relationship	1	R
Department	Entity	1	R
Participation	Relationship	3	R
Project	Entity	3	R

Average # of participations and projects per employee

Type: R - Read, W - Write, RW - Read&Write

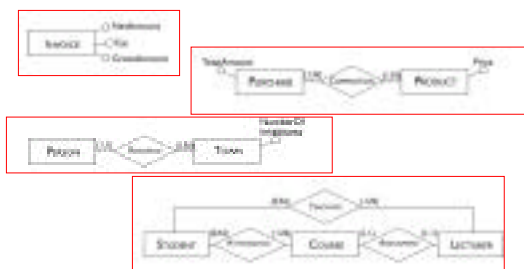
### Analysis Steps



### Analysis of Redundancies

- A redundancy in a conceptual schema corresponds to a piece of information that can be derived (that is, obtained through a series of retrieval operations) from other data in the database.
- An Entity-Relationship schema can contain various forms of redundancy.

### Examples of Redundancies



Information Systems Analysis and Design CSC340

### Deciding About Redundancies

- The presence of a redundancy in a database may be
  - ✓ **an advantage:** a reduction in the number of accesses necessary to obtain the derived information;
  - ✓ **a disadvantage:** because of larger storage requirements, (but, usually at negligible cost) and the necessity to carry out additional operations in order to keep the derived data consistent.
- The decision to maintain or delete a redundancy is made by comparing the cost of operations that involve the redundant information and the storage needed, in the case of presence or absence of redundancy.

© 2002 John Mylopoulos Database Design -- 25

Information Systems Analysis and Design CSC340

### Cost Comparison: An Example

In this schema the attribute NumberOfInhabitants is redundant.

© 2002 John Mylopoulos Database Design -- 26

Information Systems Analysis and Design CSC340

### Load and Operations for the Example

Concept	Type	Volume
Town	E	200
Person	E	1000000
Residence	R	1000000

Operation	Type	Frequency
Operation 1	I	500 per day
Operation 2	I	2 per day

- **Operation 1:** add a new person with the person's town of residence.
- **Operation 2:** print all the data of a town (including the number of inhabitants).

© 2002 John Mylopoulos Database Design -- 27

Information Systems Analysis and Design CSC340

### Table of Accesses, with Redundancy

Concept	Type	Accesses	Type
Person	Entity	1	W
Residence	Relationship	1	W
Town	Entity	1	W

Concept	Type	Accesses	Type
Town	Entity	1	R

© 2002 John Mylopoulos Database Design -- 28

Information Systems Analysis and Design CSC340

### Table of Accesses, without Redundancy

Concept	Type	Accesses	Type
Person	Entity	1	W
Residence	Relationship	1	W

Concept	Type	Accesses	Type
Town	Entity	1	R
Residence	Relationship	5000	R

© 2002 John Mylopoulos Database Design -- 29

Information Systems Analysis and Design CSC340

### Comparing the Cost of Operations

Presence of redundancy:

- ✓ Operation 1: 1,500 write accesses per day;
- ✓ The cost of operation 2 is almost negligible;
- ✓ Counting twice the write accesses, we have a total of 3,000 accesses a day.

Absence of redundancy.

- ✓ Operation 1: 1,000 write accesses per day;
- ✓ Operation 2 however requires a total of 10,000 read accesses per day;
- ✓ Counting twice the write accesses, we have a total of 12,000 accesses per day.

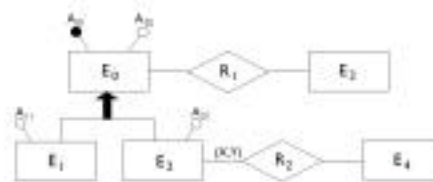
**Redundant data can improve performance, sometimes!**

© 2002 John Mylopoulos Database Design -- 30

### Removing Generalizations

- The relational model does not allow the direct representation of generalizations that may be present in an E-R diagram.
- We need, therefore, to transform these constructs into other constructs that are easier to translate: entities and relationships.

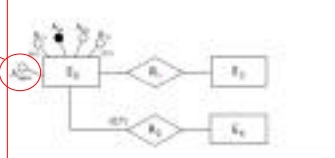
### A Schema with Generalizations



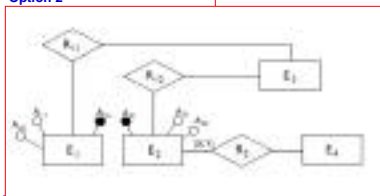
#### Option 1

#### Possible Restructurings

Note!

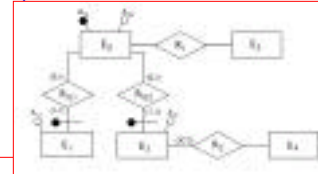


#### Option 2

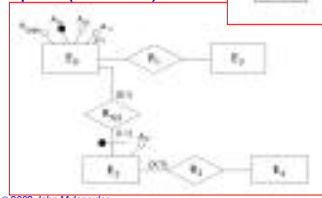


#### Option 3

#### ...Two More...



#### Option 4 (combination)



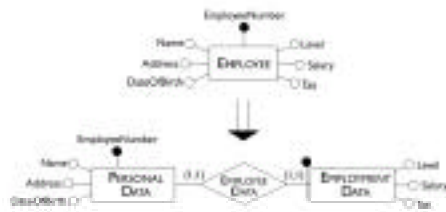
### General Rules For Removing Generalization

- Option 1 is convenient when the operations involve the occurrences and the attributes of  $E_0$ ,  $E_1$  and  $E_2$  more or less in the same way.
- Option 2 is possible only if the generalization is total and is useful when there are operations that apply only to occurrences of  $E_1$  or of  $E_2$ .
- Option 3 is useful when the generalization is not total and the operations refer to either occurrences and attributes of  $E_1$  ( $E_2$ ) or of  $E_0$ , and therefore make distinctions between child and parent entities.
- Available options can be combined (see option 4)

### Partitioning and Merging of Entities and Relationships

- Entities and relationships of an E-R schema can be partitioned or merged to improve the efficiency of operations, using the following principle: Accesses are reduced by separating attributes of the same concept that are accessed by different operations and by merging attributes of different concepts that are accessed by the same operations.
- The same criteria with those discussed for redundancies are valid in making a decision about this type of restructuring.

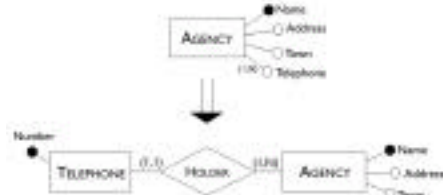
### Example of Partitioning



© 2002 John Mylopoulos

Database Design -- 37

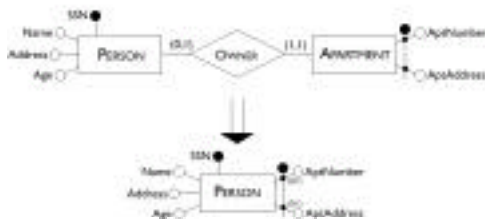
### Deletion of Multi-Valued Attribute



© 2002 John Mylopoulos

Database Design -- 38

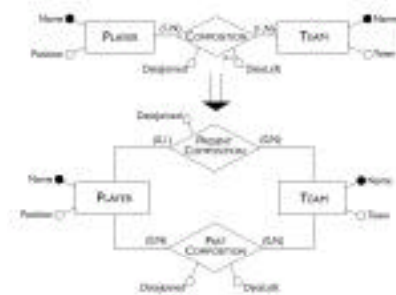
### Merging of Entities



© 2002 John Mylopoulos

Database Design -- 39

### Partitioning of a Relationship



© 2002 John Mylopoulos

Database Design -- 40

### Selection of Primary Identifiers

- The criteria for this decision are as follows.
  - ✓ Attributes with null values cannot form primary identifiers;
  - ✓ One or few attributes are preferable to many attributes;
  - ✓ An internal identifier with few attributes is preferable to an external one, possibly involving many entities;
  - ✓ An identifier that is used by many operations to access the occurrences of an entity is preferable to others.
- At this stage, if none of the candidate identifiers satisfies the above requirements, it is possible to introduce a further attribute to the entity. This attribute will hold special values (often called **codes**) generated solely for the purpose of identifying occurrences of the entity.

© 2002 John Mylopoulos

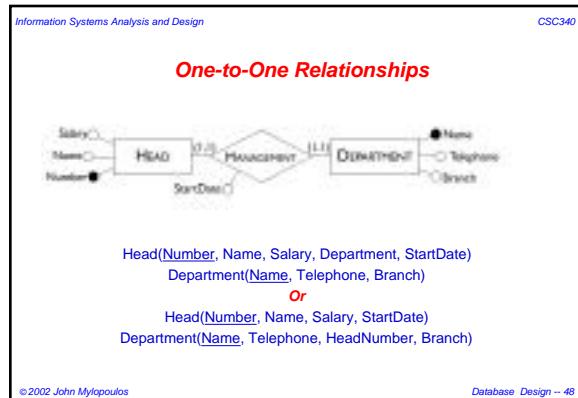
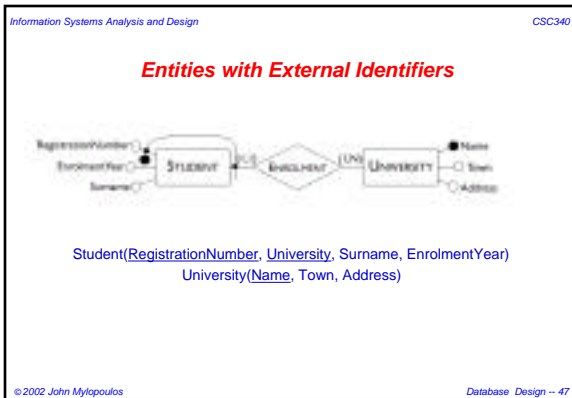
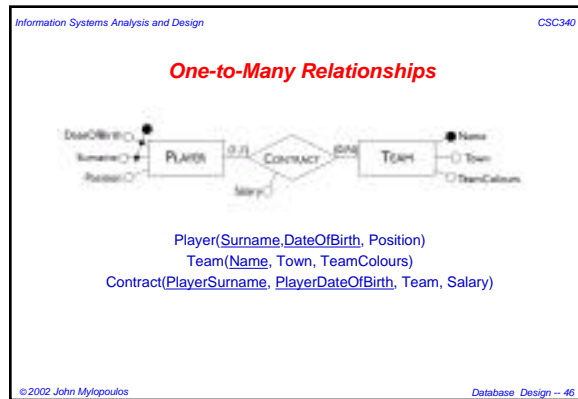
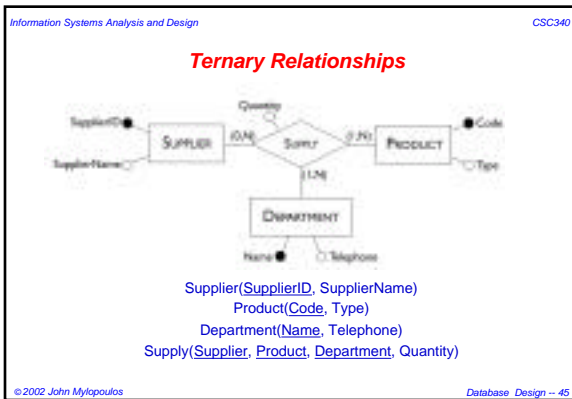
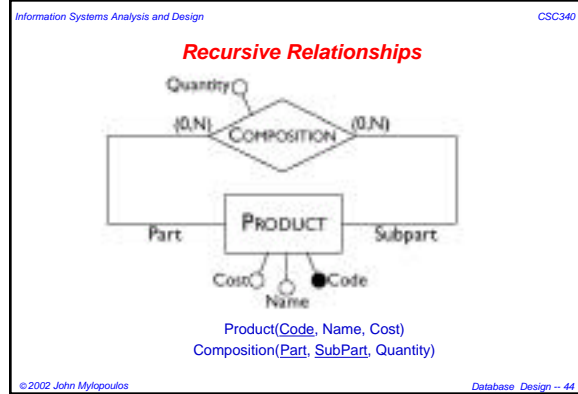
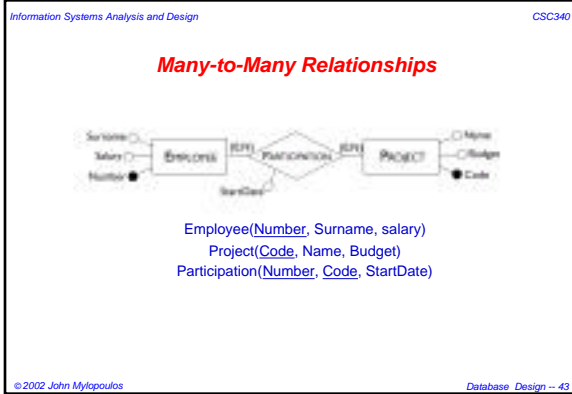
Database Design -- 41

### Translation into a Logical Schema

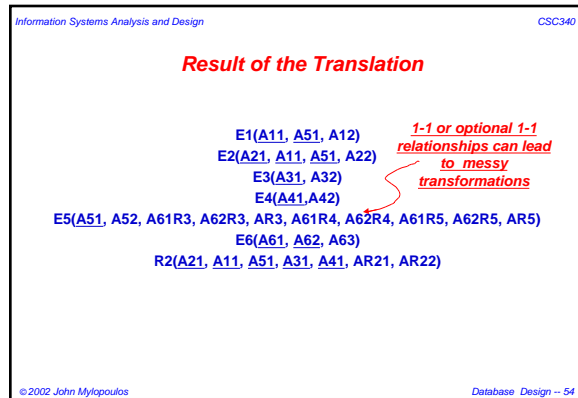
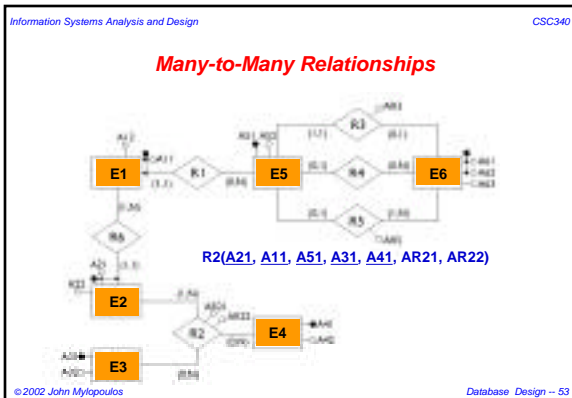
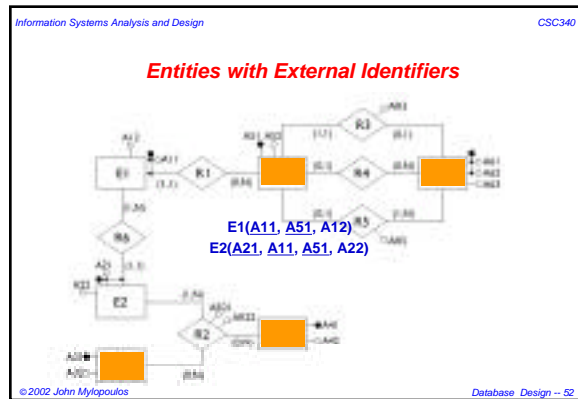
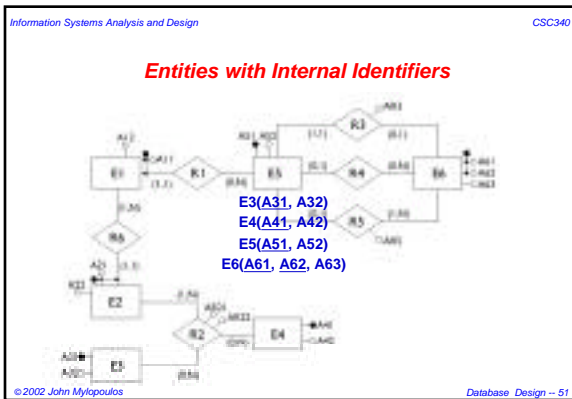
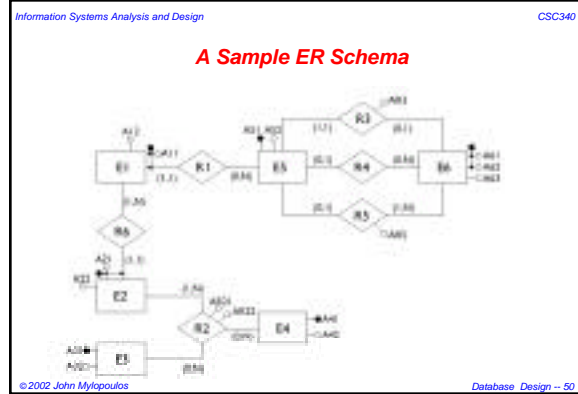
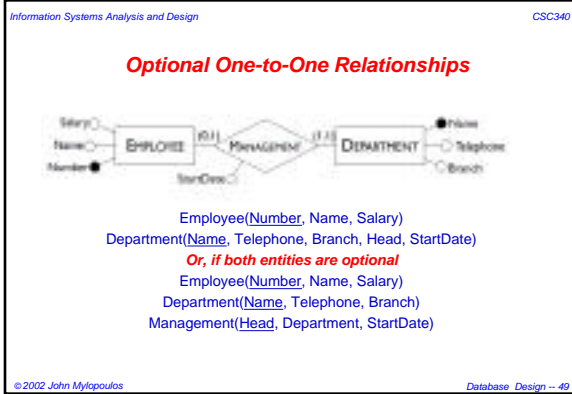
- The second step of logical design corresponds to a translation between different data models.
- Starting from an E-R schema, an equivalent relational schema is constructed. By "equivalent", we mean a schema capable of representing the same information.
- We will deal with the translation problem systematically, beginning with the fundamental case, that of entities linked by many-to-many relationships.

© 2002 John Mylopoulos

Database Design -- 42







### Summary of Transformation Rules

Type	Initial schema	Possible translation
Binary many-to-many relationship		$E_1(A_{1,1}, A_{1,2}, \dots, A_{1,n})$ $E_2(A_{2,1}, A_{2,2}, \dots, A_{2,m})$ $R(A_{1,1}, A_{1,2}, \dots, A_{1,n}, A_{2,1}, A_{2,2}, \dots, A_{2,m})$
Binary many-to-one relationship		$E_1(A_{1,1}, A_{1,2}, \dots, A_{1,n})$ $E_2(A_{2,1}, A_{2,2}, \dots, A_{2,m})$ $R(A_{1,1}, A_{1,2}, \dots, A_{1,n}, A_{2,1}, A_{2,2}, \dots, A_{2,m})$
One-to-one relationship with mandatory participation		$E_1(A_{1,1}, A_{1,2}, \dots, A_{1,n})$ $E_2(A_{2,1}, A_{2,2}, \dots, A_{2,m})$ $R(A_{1,1}, A_{1,2}, \dots, A_{1,n}, A_{2,1}, A_{2,2}, \dots, A_{2,m})$

© 2002 John Mylopoulos

Database Design -- 55

### ...More Rules...

Type	Initial schema	Possible translation
One-to-many relationship with optional participation		$E_1(A_{1,1}, A_{1,2}, \dots, A_{1,n})$ $E_2(A_{2,1}, A_{2,2}, \dots, A_{2,m})$ $R(A_{1,1}, A_{1,2}, \dots, A_{1,n}, A_{2,1}, A_{2,2}, \dots, A_{2,m})$
Relationship with optional participation		$E_1(A_{1,1}, A_{1,2}, \dots, A_{1,n})$ $E_2(A_{2,1}, A_{2,2}, \dots, A_{2,m})$ $R(A_{1,1}, A_{1,2}, \dots, A_{1,n}, A_{2,1}, A_{2,2}, \dots, A_{2,m})$

© 2002 John Mylopoulos

Database Design -- 56

### ...Even More Rules...

Type	Initial schema	Possible translation
One-to-one relationship with mandatory participation for both entities		$E_1(A_{1,1}, A_{1,2}, \dots, A_{1,n})$ $E_2(A_{2,1}, A_{2,2}, \dots, A_{2,m})$ $R(A_{1,1}, A_{1,2}, \dots, A_{1,n}, A_{2,1}, A_{2,2}, \dots, A_{2,m})$
One-to-one relationship with optional participation for one entity		$E_1(A_{1,1}, A_{1,2}, \dots, A_{1,n})$ $E_2(A_{2,1}, A_{2,2}, \dots, A_{2,m})$ $R(A_{1,1}, A_{1,2}, \dots, A_{1,n}, A_{2,1}, A_{2,2}, \dots, A_{2,m})$

© 2002 John Mylopoulos

Database Design -- 57

### ...and the Last One...

Type	Initial schema	Possible translation
One-to-one relationship with optional participation for both entities		$E_1(A_{1,1}, A_{1,2}, \dots, A_{1,n})$ $E_2(A_{2,1}, A_{2,2}, \dots, A_{2,m})$ $R(A_{1,1}, A_{1,2}, \dots, A_{1,n}, A_{2,1}, A_{2,2}, \dots, A_{2,m})$

© 2002 John Mylopoulos

Database Design -- 58

### The Training Company Revisited



© 2002 John Mylopoulos

Database Design -- 59

### Operational Requirements, Revisited

- operation 1:** insert a new trainee including all his or her data (to be carried out approximately 40 times a day);
- operation 2:** assign a trainee to an edition of a course (50 times a day);
- operation 3:** insert a new instructor, including all his or her data and the courses he or she is qualified to teach (twice a day);
- operation 4:** assign a qualified instructor to an edition of a course (15 times a day);
- operation 5:** display all the information on the past editions of a course with title, class timetables and number of trainees (10 times a day);
- operation 6:** display all the courses offered, with information on the instructors who are qualified to teach them (20 times a day);
- operation 7:** for each instructor, find the trainees all the courses he or she is teaching or has taught (5 times a week);
- operation 8:** carry out a statistical analysis of all the trainees with all the information about them, about the editions of courses they have attended and the marks obtained (10 times a month).

© 2002 John Mylopoulos

Database Design -- 60

## Database Load

Table of volumes

Concept	Type	Volume
Class	E	8000
CourseEdition	E	1000
Course	E	200
Instructor	E	300
Freelance	E	250
Permanent	E	50
Trainee	E	5000
Employee	E	4000
Professional	E	1000
Employer	E	8000
PastAttendance	R	10000
CurrentAttendance	R	500
Composition	R	8000
Type	R	1000
PastTeaching	R	900
CurrentTeaching	R	100
Qualification	R	500
CurrentEmployment	R	4000
PastEmployment	R	10000

Table of operations

Operation	Type	Frequency
Operation 1	I	40 per day
Operation 2	I	50 per day
Operation 3	I	2 per day
Operation 4	I	15 per day
Operation 5	I	10 per day
Operation 6	I	20 per day
Operation 7	I	5 per day
Operation 8	B	10 per month

## Access Tables

The attribute `NumberOfParticipants` in `CourseEdition` can be derived from the relationships `CurrentAttendance` and `PastAttendance`.

Operation 2 with redundancy

Concept	Type	Acc	Type
Trainee	E	1	R
CurrentAttendance	R	1	W
CourseEdition	E	1	R
CourseEdition	E	1	W

Operation 2 without redundancy

Concept	Type	Acc	Type
Trainee	E	1	R
CurrentAttendance	R	1	W

Operation 5 with redundancy

Concept	Type	Acc	Type
CourseEdition	E	5	R
Type	R	5	R
Course	E	1	R
Composition	R	40	R
Class	E	40	R

Operation 5 without redundancy

Concept	Type	Acc	Type
CourseEdition	E	5	R
Type	R	5	R
Course	E	1	R
Composition	R	40	R
Class	E	40	R
PastAttendance	E	50	R

## Analysis of Redundancy

- From the access tables we obtain (giving double weight to the write accesses):
  - presence of redundancy: for operation 2 we have 100 read accesses and 100 write accesses per day; for operation 5 we have 910 read accesses per day, for a total of 1,210 accesses per day;
  - without redundancy: for operation 2 we have 50 read accesses per day and 100 write accesses per day; for operation 5, we have 1,410 read accesses per day, for a total of 1,660 accesses per day.
- Thus, redundancy makes sense in this case, so we leave the attribute `NumberOfParticipants` as an attribute of the entity `CourseEdition`.

## Removing Generalizations

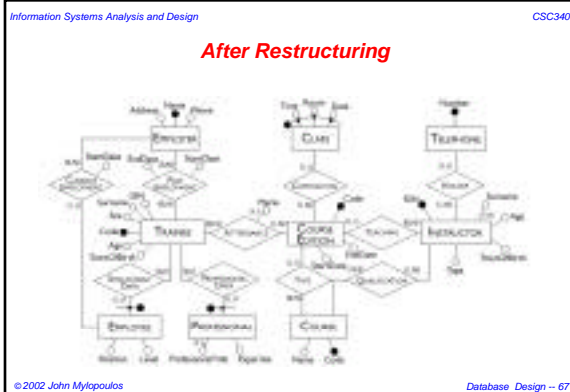
- For the generalization on instructors:
  - the relevant operations make no distinction between the child entities and these entities have no specific attributes;
  - we can therefore delete the child entities and add an attribute `Type` to the parent entity.
- For the generalization on trainees:
  - the relevant operations make no distinction between the child entities, but these entities have specific attributes;
  - we can therefore leave all the entities and add two relationships to link each child with the parent entity: in this way, we will have no attributes with possible null values on the parent entity and the dimension of the relations will be reduced.

## Partitioning and Merging of Concepts

- The relationships `PastTeaching` and `PresentTeaching` can be merged since they describe similar concepts between which the operations make no difference. A similar consideration applies to the relationships `PastAttendance` and `PresentAttendance`.
- The multi-valued attribute `Telephone` can be removed from the `Instructor` entity by introducing a new entity `Telephone` linked by a one-to-many relationship to the `Instructor` entity.

## Choice of Main Identifiers

- Trainee entity:
  - there are two identifiers: the social security number and the internal code;
  - it is far preferable to choose the latter: a social security number will require several bytes whereas an internal code, which serves to distinguish between 5000 occurrences, requires a few bytes.
- CourseEdition entity:
  - it is identified externally by the `StartDate` attribute and by the `Course` entity;
  - we can see however that we can easily generate for each edition a code from the course code: this code is simpler and can replace the external identifier.



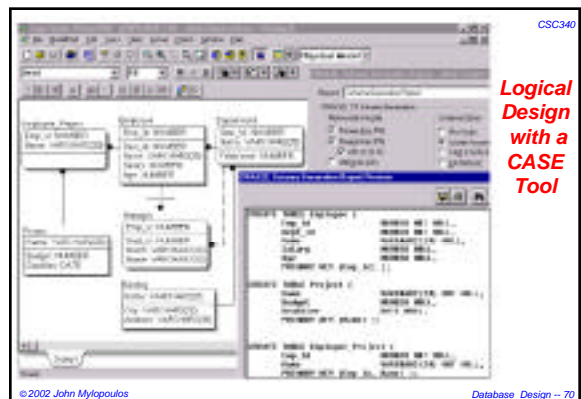
Information Systems Analysis and Design CSC340

### Translation into the Relational Model

CourseEdition(Code, StartDate, EndDate, Course, Instructor)  
 Class(Time, Room, Date, Edition)  
 Instructor(SSN, Surname, Age, TownOfBirth, Type)  
 Telephone(Number, Instructor)  
 Course(Code, Name)  
 Qualification(Course, Instructor)  
 Trainee(Code, SSN, Surname, Age, TownOfBirth, Sex)  
 Attendance(Trainee, Edition, Marks\*)  
 Employer(Name, Address, Telephone)  
 PastEmployment(Trainee, Employer, StartDate, EndDate)  
 Professional(Trainee, Expertise, ProfessionalTitle\*)  
 Employee(Trainee, Level, Position, Employer, StartDate)

© 2002 John Mylopoulos Database Design -- 68

- Information Systems Analysis and Design CSC340
- ### Logical Design Using CASE Tools
- The logical design phase is partially supported by database design tools:
    - ✓ the translation to the relational model is carried out by such tools semi-automatically;
    - ✓ the restructuring step is difficult to automate and CASE tools provide little or no support for it.
  - Most commercial CASE tools will generate automatically SQL code for the creation of the database.
  - Some tools allow direct connection with a DBMS and can construct the corresponding database automatically.
- © 2002 John Mylopoulos Database Design -- 69



Information Systems Analysis and Design CSC340

### What is a Good Relational Schema Like?

■ Some relational schemata are "better" representations than others. What are the criteria we can use to decide whether a diagram is better than another? Should we have more/fewer relations as opposed to attributes?

**Enter normal forms**

■ An attribute  $a$  (functionally) depends on a set of attributes  $a_1, a_2, \dots, a_n$  if these determine uniquely the value of  $a$  for every tuple of the relation where they appear together

$$a_1, a_2, \dots, a_n \rightarrow a$$

■ Example: For the relation  
 Course(name, title, instrName, rmName, address),  
 E.g. (csc340, "Analysis and Design", JM, RW117, "48 StG")  
 the title attribute depends on the name attribute. Likewise, the address attribute depends on the rmName attribute,  
 name  $\rightarrow$  title, also rmName  $\rightarrow$  address

© 2002 John Mylopoulos Database Design -- 71

Information Systems Analysis and Design CSC340

### Examples of Functional Dependencies

- Consider  
 Supplier(S#, SName, Status, Address)
- Here SName, Status, Address functionally depend on S# because S# uniquely determines the values of the other attributes of the Supplier relation  
 $S\# \rightarrow SName, Status, Address$
- Likewise, assuming that Lastname, Firstname uniquely identify people, we have  
 $Lastname, Firstname \rightarrow Salary, Address$

© 2002 John Mylopoulos Database Design -- 72

### What's Wrong with Un-Normalized Relations?

- Normalization helps produce a schema that is not redundant and does not suffer from **anomalies**.
- Consider `Emp(Emp#, Ename, Address, Dept, Mngr#)` with  
`Emp1# --> EName, Address, Dept, Mngr#, Dept --> Mngr#`
- **Insertion anomaly**: We can't add information about a new department and its manager until we have an employee in that department.
- **Deletion anomaly**: If we delete the only employee in a department, we lose information about the department (e.g., its manager)
- **Update anomaly**: If we update the `Mngr#` attribute of one tuple, we must do it for all, otherwise we have an inconsistent database.

*It's easy to get an inconsistent database when it's not normalized*

### How Do We Identify Functional Dependencies?

- 1) Think about the meaning of different attributes and try to think of situations where the value of `a` is not determined by the values of `a1, a2, ..., an`
- 2) Alternatively, if you are given sample values for the attributes of the relation (see below), check to ensure that every combination of values for `a1, a2, ..., an` has the same associated value for `a`

Name	Title	Instructor	Office	Tutors	Enrollment
csc148	"Intro."	Reiter	LP290G	4	133
csc228	"DP"	Clarke	SF285	3	124
csc238	"Logic"	Fich	SF254	3	85
csc324	"PLs"	Bonner	LP354	2	72
csc340	"SA"	Reiter	LP290G	2	121
csc408	"SE"	Clarke	SF285	3	88
csc434	"DM"	Fich	SF254	3	107

*What functional dependencies are appropriate here?*

### Normalizing Relational Schemas: 1NF

- 1) A relation is in **First Normal Form (1NF)** if it does not include any multi-valued attributes or any composite attributes.  
e.g., consider the relation  
`Course(name, title, instrName*, studentNo*, addr)`  
`Course` is not in 1NF because of two attribute groups that repeat (instructor and student groups)
- 2) To place a relation in 1NF, take each multi-valued attribute or composite attribute and promote it into a relation in its own right.

### An Example

- 1) For the `Course(name, title, instrName*, studentNo*, addr)`, example, assume that `addr` is a composite attribute consisting of a `streetNm`, `streetNo`, `city` and `postalCode`:  
`==> Course(name, title)`  
`CourseStud(name, studentNo)`  
`CourseInstr(name, instrName)`  
`CourseAddr(name, streetNm, streetNo, city, postalCode)`

### Normalizing Relational Schemas: 2NF

- 1) A relation is in **Second Normal Form (2NF)** if it is in 1NF and, moreover, all non-key attributes depend on all elements of its key, rather than a subset.
- 2) Consider  
`Room(street, number, bldgNm, room#, capacity, AVEquip)`
- `Room` is not in 2NF because its address attributes functionally depend on its `bldgNm` key element, rather than the combination (`room#`, `bldgNm`)
- 3) To place a 1NF relation into 2NF, take each non-key attribute that does not depend on the full key and move it to a relation identified by the partial key  
`==> Room(bldgNm, room#, capacity, AVEquip), Building(bldgNm, street, number)`

### Normalizing Relational Schemas: 3NF

- A relation is in **Third Normal Form (3NF)** if it is in 2NF and none of its non-key attributes depends on any other non-key attribute.
- Assuming that each course has only one instructor (why do we need this assumption?), `Course` is not in 3NF because `instrDept` depends on `instrNm`:  
`Course(name, year, sem, instrNm, instrDept, enrol#)`
- To place a 2NF relation into 3NF, move each non-key attribute that depends on other non-key attributes to another relation  
`==> Course(name, year, sem, instrNm, enrol#)`  
`Instructor(name, dept)`