



DMIF, University of Udine

---

# Data Management for Big Data

## *Entity-Relationship Diagrams*

Andrea Brunello

andrea.brunello@uniud.it

April 2022

The process of designing a database is typically articulated into different phases:

- ① Brainstorming meetings with IT personnel and all interested stakeholders
  - Collection of requirements
  - Design of a conceptual model (e.g., E-R model)
  - The E-R model has a specific notation, the *E-R diagram*, that helps all involved parts to discuss about the future database
- ② Translation of the conceptual model into a logical model
  - Typically, a set of translation rules is followed
  - At this stage, a specific DBMS technology must be chosen (e.g., relational DB)
- ③ Addition to the logical model of details regarding the physical, low-level aspects (e.g., usage of indexes)
  - The physical schema is thus obtained

- The initial phase of database design is to characterize fully the data needs of the prospective database users
- Such needs are encoded into a conceptual model, that acts as a connecting point between stakeholders and IT personnel
- We will consider the **Entity-Relationship** model, which makes use of *diagrams* to represent the overall logical structure of a database graphically
- Such diagrams are typically complemented by free text annotations, that describe the requirements that could not be encoded in the diagram, or present the kinds of operations that users are expected to perform on the data



# Entity-Relationship model

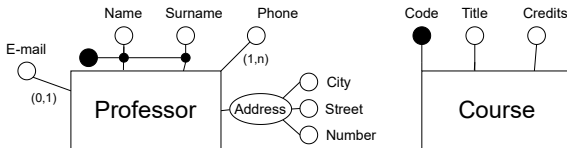
- An Entity-Relationship model describes a domain by means of:
  - Entity sets
  - Relationship sets
  - Attributes
- There are several ways in which entity sets, relationship sets and attributes can be represented in E-R diagrams
- We will consider the notation proposed in the book:

*Database Systems - Concepts, Languages and Architectures,*  
*P. Atzeni, S. Ceri, S. Paraboschi, and R. Torlone*



# Entity sets

- An **entity** is an object that exists in the considered domain and is distinguishable from other objects
  - Example in the University domain: a specific professor, student, or course
- An **entity set** is a set of entities of the same type that share the same properties
  - Example: the set of all professors, each characterized by a name, a surname, and a salary
- Such properties, possessed by all members of an entity set, are represented by **attributes**
- A subset of the attributes forms a **primary key** of the entity set, i.e, knowing the values of such attributes, it is possible to uniquely identify a member (entity) of the entity set

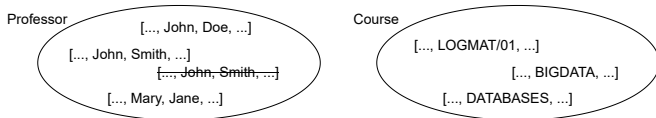


- Entity set *Professor* has the attributes:
  - E-mail*, optional
  - Name* and *Surname*, that together make the primary key
  - Phone*, multi-valued
  - Address*, composite, and made by *City*, *Street*, *Number*
- Entity set *Student* has the attributes:
  - Code*, that makes the primary key
  - Name*
  - Surname*



# Entity sets

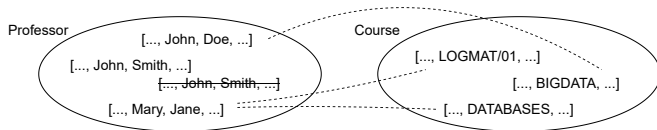
## Intuitive meaning



- An entity set cannot contain two entities with the same values on the primary key attributes
- Observe that, being a set, each entity set has a *trivial* primary key composed of all the attributes

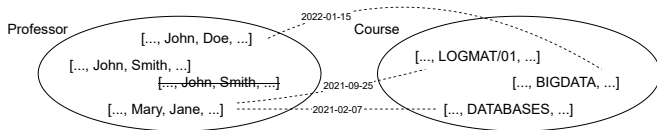
- A **relationship** is an association among several entities
  - E.g., *Professor:John Doe* may be linked with *Course:BIGDATA* by means of a relationship *Teaches*
  - The relationship is an association between actual entities
- A **relationship set**  $R \subseteq E_1 \times E_2 \times \dots \times E_n$  is a mathematical relation among  $n \geq 2$  entities ( $e_i$ ), each taken from its entity set ( $E_i$ )
  - $R = \{(e_1, e_2, \dots, e_n) \mid e_1 \in E_1, e_2 \in E_2, \dots, e_n \in E_n\}$
  - where  $(e_1, e_2, \dots, e_n)$  is a *relationship*
- For instance,  $(\text{John}, \text{Doe}, \text{BIGDATA}) \in \text{Teaches}$  is a relationship belonging to relationship set  $\text{Teaches} \subseteq \text{Professor} \times \text{Course}$





- The picture shows three relationships (dashed lines), which we assume belonging to relationship set *Teaches*:
  - John Doe* is associated to *BIGDATA*,  
(*John*, *Doe*, *BIGDATA*)
  - Mary Jane* is associated to *LOGMAT/01*,  
(*Mary*, *Jane*, *LOGMAT/01*)
  - Mary Jane* is associated to *DATABASES*,  
(*Mary*, *Jane*, *DATABASES*)

- Attributes may also be associated with relationships, e.g., to track when a professor started teaching a course
- Intuitively, a relationship is thus as follows:
  - (*John, Doe*, *BIGDATA*, 2022 – 01 – 15)





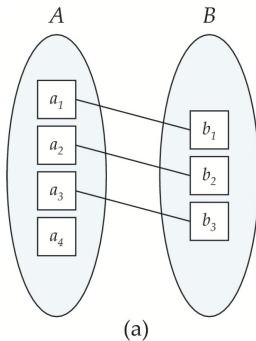
# Relationship sets

## Cardinality constraints

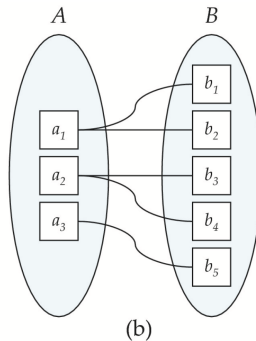
- Given entity sets  $E_1, E_2$ , and a relationship set  $R$  between them, entities may participate into relationships in different ways, for instance:
  - An entity of  $E_1$  may not participate to any relationship with an entity of  $E_2$  (a prof. may not be teaching any courses)
  - An entity of  $E_1$  may participate to  $> 1$  relationship with an entity of  $E_2$  (a prof. may be teaching more than one course)
- In general, for a binary relationship set we can identify the following **cardinalities**:
  - One to one
  - One to many
  - Many to one
  - Many to Many
- Moreover, entity sets may have a *partial* or *total* **participation** to a relationship set

# Relationship sets

## One to one and one to many

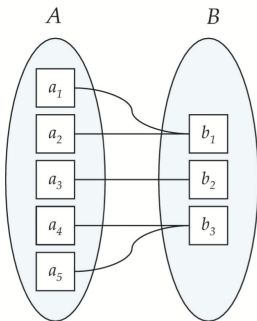


One to one



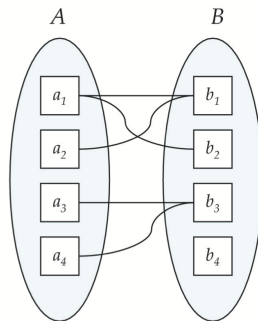
One to many

- Left:  $A$  has a partial participation, while  $B$  total
- Right: both  $A$  and  $B$  have a total participation



(a)

Many to  
one



(b)

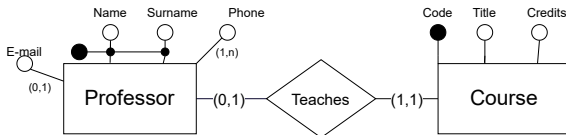
Many to many

- Left: both  $A$  and  $B$  have a total participation
- Right: both  $A$  and  $B$  have a total participation



# Relationship sets

## Notation for participation and cardinality constraints

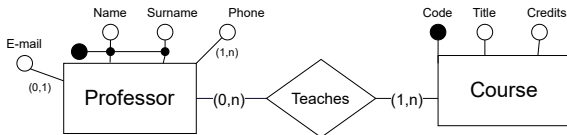


- A professor teaches in at most one course:  $(0,1)$
- A course is taught by one and only one professor:  $(1,1)$
- $(1,1)$  is assumed by default, and can be omitted
- Concerning the relationship set *Teaches*, this means that:
  - A given entity  $p_1 \in Professor$  can appear in at most one relationship  $(p_1, c) \in Teaches$  where  $c \in C$
  - A given entity  $c_1 \in Course$  appears in exactly one relationship  $(p, c_1) \in Teaches$  where  $p \in P$



# Relationship sets

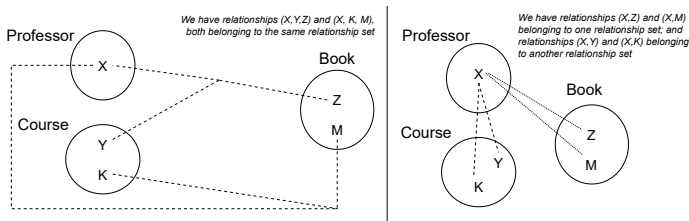
## Notation for participation and cardinality constraints



- A professor teaches in zero or more courses:  $(0,n)$
- A course is taught by one or more professors:  $(1,n)$
- Concerning the relationship set *Teaches*, this means that:
  - A given entity  $p_1 \in Professor$  can appear in any number of relationships  $(p_1, c) \in Teaches$  where  $c \in C$
  - A given entity  $c_1 \in Course$  appears in at least one relationship  $(p, c_1) \in Teaches$  where  $p \in P$
- Observe that the only values allowed in the parentheses are 0, 1,  $n$ . We cannot specify a number, e.g., 4

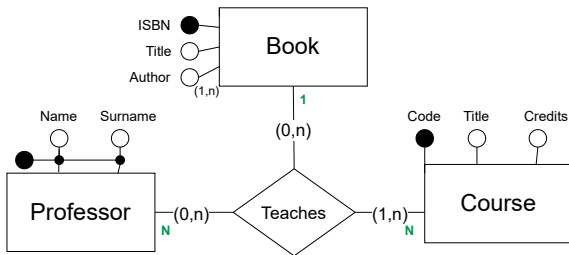
# Relationships involving $> 2$ entities

- Sometimes, it is necessary to involve more than 2 entities in a relationship (though in a DB most of them are binary, and nearly always at most ternary)
- For instance:
  - professor X teaches course Y using the book Z*
  - professor X teaches course K using the book M*



- Observe how the relationships on the right convey less information than those on the left





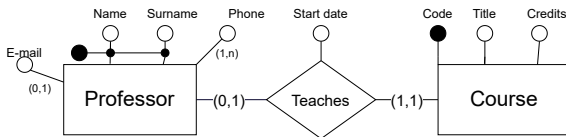
- Participation/cardinality constraints still apply to ternary relationships, although things complicate a little bit
- Here: a professor may teach zero or more courses, each course is taught by at least one professor, and a book is used within zero or more teaching activities, although a professor uses only a book within the same course



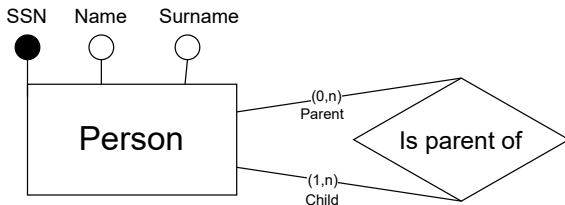
# Relationship sets

## Ternary relationship: setting participations/cardinalities

- If I have an entity of *Professor*, how many entities of *Book* and *Course* together can be associated to it?
  - Here, the answer is zero or more, so entity set *Professor* participates with  $(0, n)$  to *Teaches*
- As for the green numerosities, they allow us to solve some ambiguity. Without them, we wouldn't know if a professor may appear in more than one relationship because:
  - He may teach more than one course using the same book
  - He may teach only one course but using multiple books
  - He may teach multiple courses using multiple books
- Setting the green 1 close to *Book*, we say that, given a *Professor* and a *Course*, there can be only a *Book* associated to them both. On the contrary, given a *Professor* and a *Book*, we may have more courses (green  $N$  close to *Course*)



- Relationship sets can also have attributes
- This is quite natural when the attribute tells us something regarding the specific entities involved in the relationship
- For instance, here we track the date in which a given professor started teaching a given course



- Entity sets of a relationship need not be distinct
  - For instance, consider the case in which we want to track the prerequisites of a given course
- Each occurrence of an entity set plays a “role” in the relationship
- In that case, it is useful to write down the **roles** as labels on the arcs

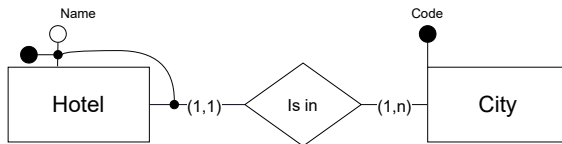


- Consider the following scenario:
  - We want to keep track of hotels, each characterized by a name
  - There may be multiple hotels with the same name in different cities, however, given a city, there cannot be two hotels with the same name
  - Each city is univocally identified by a code
- Intuitively, to univocally identify a hotel, its name is not enough; we also need to know the city where it is

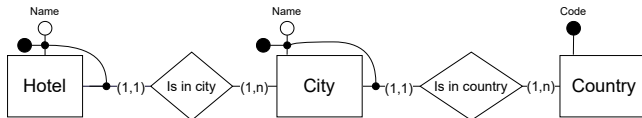


- The notion of **weak entity set** formalizes the above intuition. A weak entity set is one whose existence is dependent on another entity, called its **identifying entity**; instead of associating a primary key with a weak entity, to uniquely identify it we use the identifying entity, along with extra attributes called **partial keys**
- Each weak entity must be associated with a *a corresponding* identifying entity; that is, the weak entity set is said to be **existence dependent** on the identifying entity set. The identifying entity set is said to **own** the weak entity set that it identifies. The relationship associating the weak entity set with the identifying entity set is called the **identifying relationship**

- Observe the  $(1, 1)$  cardinality on the *Hotel* side

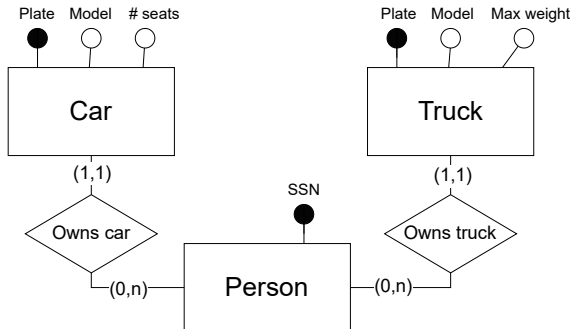


- We may also define chains of weak entity sets



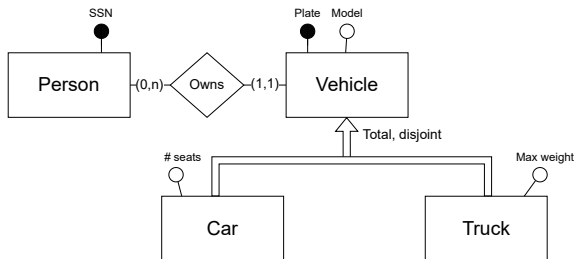
- A weak entity set may depend on more than one identifying entity

- It may happen that different entity sets share some commonalities regarding their attributes or relationships
  - E.g., a *Truck* and a *Car* may both be considered as a *Vehicle*





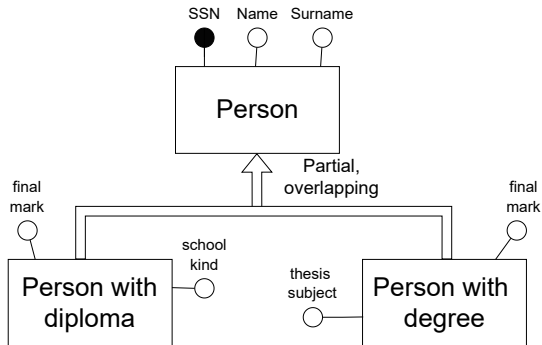
- The **specialization** construct allows us to handle such a scenario
- It is a similar concept as specialization/generalization in object oriented programming
- A “lower-level” entity set inherits all the attributes and relationship participations of the “higher-level” entity set which it is linked to





# Specialization

## Another example

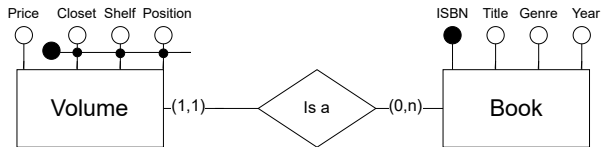




- A specialization may be:
  - **Total or Partial**
    - **Total:** each higher-level entity corresponds to at least one of the lower-level entities (e.g., a *Vehicle* is a *Car* or a *Truck*)
    - **Partial:** there may be higher-level entities that do not correspond to any lower-level entity (e.g., a *Person* may not hold a diploma nor a degree)
  - **Disjoint or Overlapping**
    - **Disjoint:** a higher-level entity corresponds to at most one of the lower-level entities (e.g., a *Vehicle* can be a *Car* or a *Truck*, but not both)
    - **Overlapping:** a higher-level entity may correspond to more lower-level entities (e.g., a *Person* may hold both a diploma as well as a degree)
- Of course, each lower-level entity corresponds to a higher-level entity
- Always write in the diagram the kind of specialization

# The *instance* of construct

- Consider the following scenario:
  - A library maintains a catalogue of published books. For each book, it records: ISBN, title, genre, and year of publication
  - Not all books are available at the library. Of the available ones, we record the selling price and the location of each copy: closet, shelf, position in the shelf
- We may notice the presence of two different entities:
  - Book*: a “virtual” book, with its ISBN
  - Copy*: the actual, “physical” volume, which can be considered as an instance of a book

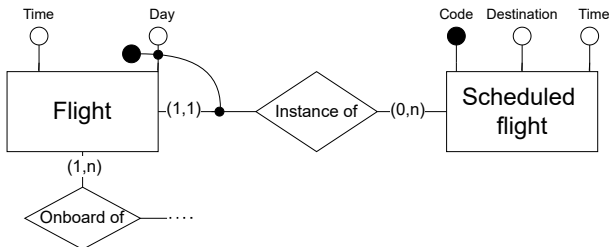




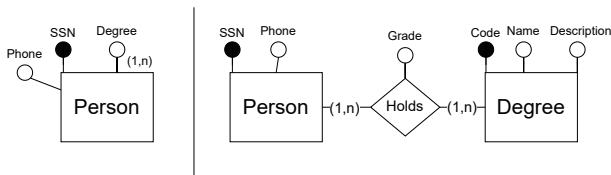
# The *instance of* construct

## Another example

- An airport wants to keep track of flights
- A flight is uniquely identified by a code, and characterized by a destination and a daily scheduled departure time (e.g., the AZ123 flight of 3:50 PM to Rome)
- For each occurrence of a flight, we also want to keep track of the actual departure time, its crew, [...]



- When should attributes be used, and when entity and relationship sets instead?
- It really depends on the domain, and on the amount of detail that we want to express





# Design example

- We want to design a database for the management of an University
- The database will have to deal with information regarding professors, students, and courses
- Each professor is uniquely identified by his/her social security number, and is characterized by a name, a surname, and a salary. Professors may be full or associate. For the former kind, we also want to keep track of the date when they became full professors
- Each professor belongs to a department, which is uniquely identified by its name, and characterized by an address (city, street, number) and a budget



## Design example (Cont.)

- Finally, full professors may be in charge of a commission, which is uniquely identified by a code and possibly characterized by a description
- The university gives several courses, each identified by a code and characterized by its credits and a possible prerequisite. A course may be held on different years, and each edition can be taught by a different professor
- Students participate to courses (a student may possibly follow the same course several times on different years). When a student passes the exam of a given course, we want to keep track of the exam mark. Each student is uniquely identified by his/her social security number, and characterized by a name and a surname
- Also, every students holds a student card, that has a unique code and has been released on a certain year