

# Data Management

## 02 Conceptual Design

**Matthias Boehm**

Graz University of Technology, Austria  
Computer Science and Biomedical Engineering  
Institute of Interactive Systems and Data Science  
BMK endowed chair for Data Management

# Announcements/Org

- **#1 Video Recording**

- Link in **TeachCenter** & **TUbe** (lectures will be public)



- **#2 Course Registrations SS20**

- Data Management (lectures/exercises): **143/144**
- Databases (combined lectures/exercises): **68**

Total:  
**211**

- **#3 Exercise 1**

- Task Description published last night (discussed in today's lecture)
- Deadline: **Nov 03** in TeachCenter

- **#4 Learning Analytics (LA)**

- **LA Research Study** by Carla Souta Barreiros
- <https://tc.tugraz.at/main/course/view.php?id=3124>  
(video and self-regulated learning questionnaire)

**“Participate and  
make a positive  
impact”**



# Agenda

- **DB Design Lifecycle**
- **ER Model and Diagrams**
- **Exercise 01 – Data Modeling**



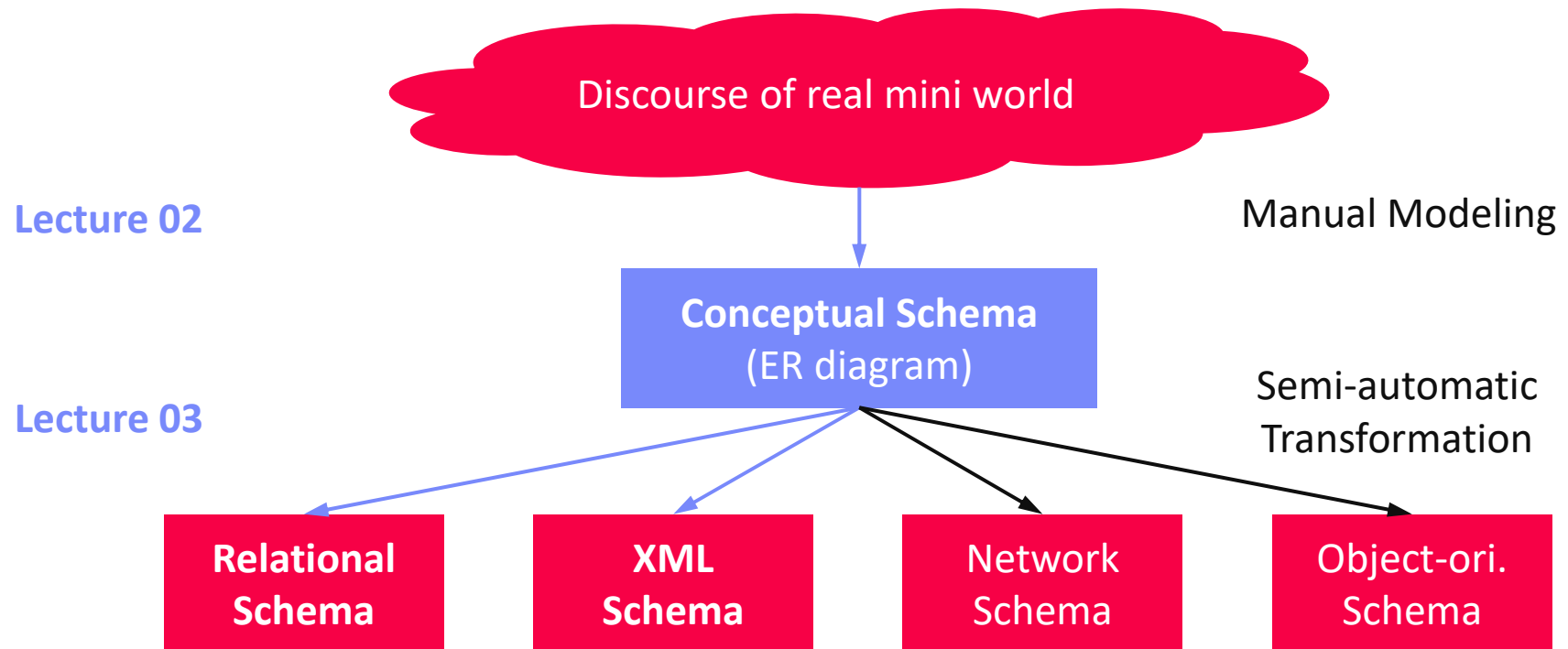
[**Credit:** Alfons Kemper, André Eickler: Datenbanksysteme - Eine Einführung, 10. Auflage. De Gruyter Studium, de Gruyter Oldenbourg 2015, ISBN 978-3-11-044375-2, pp. 1-879]

# DB Design Lifecycle

# Data Modeling

## ■ Data Model

- Concepts for describing data objects and their relationships (meta model)
- **Schema:** Description (structure, semantics) of specific data collection



# Data Models

## ■ Conceptual Data Models

- **Entity-Relationship Model (ERM)**, focus on data, ~1975
- Unified Modeling Language (UML), focus on data and behavior, ~1990

## ■ Logical Data Models

- **Relational** (Object/Relational)

- Key-Value
- Document (XML, JSON)
- Graph
- Time Series
- Matrix/Tensor

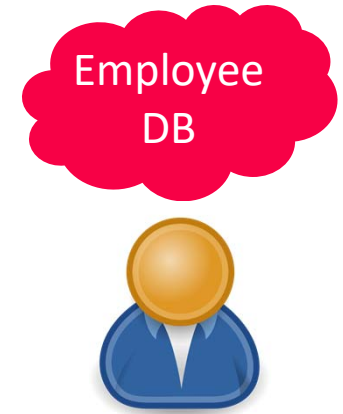
Partly covered  
in part B

- Object-oriented
- Network
- Hierarchical

Mostly obsolete

# DB Design Lifecycle Phases

- **#1 Requirements engineering**
  - Collect and analyze data and application requirements
  - ➔ Specification documents
- **#2 Conceptual Design** (lecture 02, exercise 1)
  - Model data semantics and structure, independent of logical data model
  - ➔ ER model / diagram
- **#3 Logical Design** (lecture 03, exercise 1)
  - Model data with implementation primitives of concrete data model
  - ➔ e.g., relational schema + integrity constraints, views, permissions, etc
- **#4 Physical Design** (lecture 07, exercise 3)
  - Model **user-level data organization** in a specific DBMS (and data model)
  - Account for deployment environment and performance requirements



# Relevance in Practice

## ■ Analogy ERM-UML

- **Model-driven development** (self-documenting, but quickly outdated)
- **But:** Once data is loaded, data model and schema harder to change

## ■ **Observation: Full-fledged ER modeling rarely used in practice**

- Often the logical schema (relational schema) is directly created, maintained and used for documentation
- **Reasons:** redundancy, indirection, single target (relational)
- Simplified ER modeling used for brainstorming and early ideas

## ■ Goals

- **Understanding of proper database design** from conceptual to physical schema
- ER modeling as a helpful **tool in database design**
- Schema transformation and normalization as blueprint for **good designs**



# Tool Support

## ■ #1 Visual Design Tools

- Draw ER diagrams in any presentation software (e.g., MS PowerPoint, LibreOffice)
- Many desktop or web-based tools support ER diagrams directly (e.g., MS Visio, creately.com)

## ■ #2 Design Tools w/ Code Generation

- Draw and validate ER diagrams
- Generate relational schemas as SQL DDL scripts
- **Examples:** SAP (Sybase) PowerDesigner, MS Visual Studio plugins (SQL server), etc.

➔ **Note:** For the exercises, please use basic drawing tools (existing tools use slightly diverging notations)

# Entity-Relationship (ER) Model and Diagrams



[Peter P. Chen: The Entity-Relationship Model - Toward a Unified View of Data. **ACM Trans. Database Syst.** **1(1)** 1976]

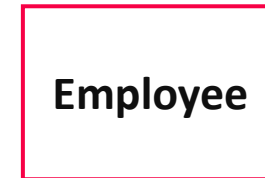
[Peter P. Chen: The Entity-Relationship Model: Toward a Unified View of Data. **VLDB** 1975]



# ER Diagram Components (Chen Notation)

## Entity Type (noun)

- Entities are objects of the real world
- An entity type (or **entity set**) represents a collection of entities



Weak  
entities



## Relationship Type (verb)

- Relationships are concrete associations of entities
- Relationship type (or **relationship set**) or relationship of entity types



$works \subseteq A \times B$

## Attribute

- Entities or relationships are characterized by attribute-value pairs
- Attribute types (or value sets) describe entity and relationship types
- Extended attributes: composite, multi-valued, derived

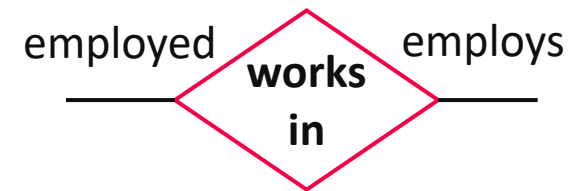


Multi-valued  
attributes



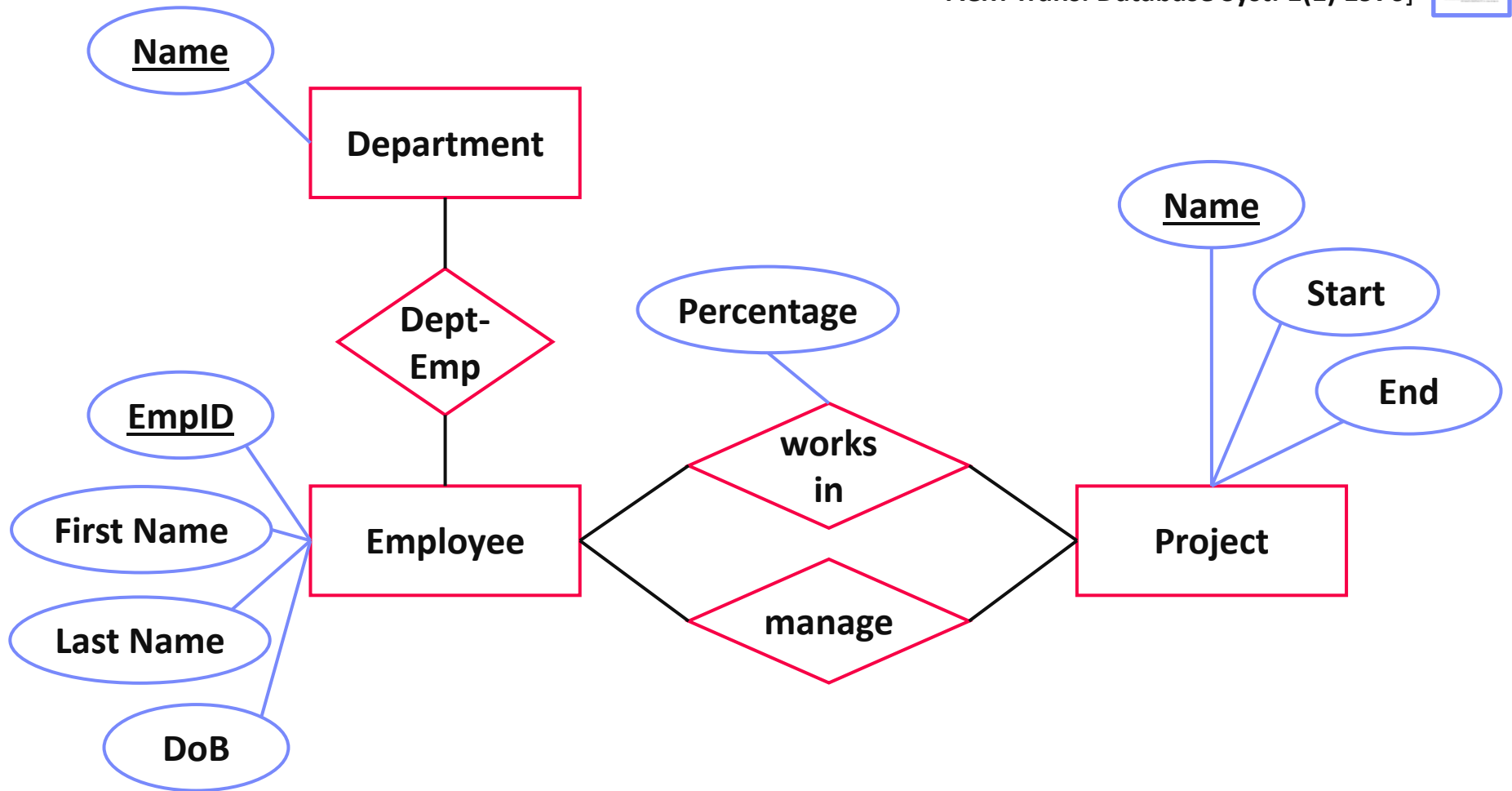
## ER Diagram Components (Chen Notation), cont.

- **Keys**
  - Attributes that uniquely identify an entity
  - Every entity type must have such a key
  - Natural or surrogate (artificial) keys
- **Role**
  - Optional description of relationship types
  - Useful for recursive relationships



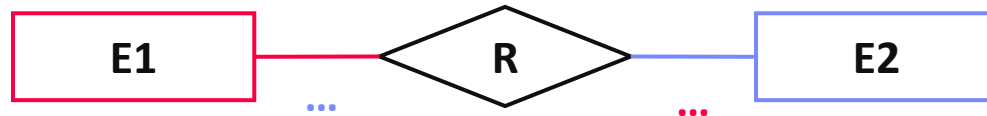
# An EmployeeDB Example

[Peter P. Chen: The Entity-Relationship Model - Toward a Unified View of Data. ACM Trans. Database Syst. 1(1) 1976]



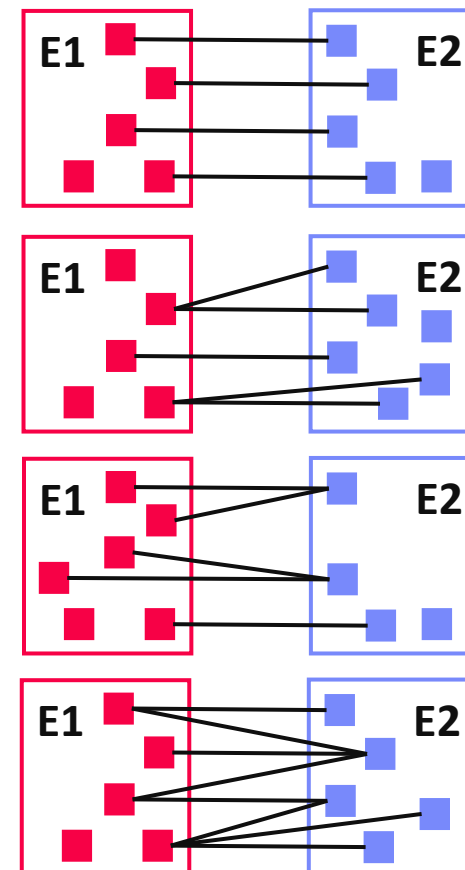
# Multiplicity/Cardinality in Chen Notation

1 .. [0,1]  
N ... [0,1,N]



$$R \subseteq E1 \times E2$$

- **1:1 (one-to-one)**  $\longleftrightarrow$ 
  - Each e1 relates to at most one e2
  - Each e2 relates to at most one e1
- **1:N (one-to-many)**  $\longleftarrow$ 
  - Each e1 relates to many e2 (0,1,...N)
  - Each e2 relates to at most one e1
- **N:1 (many-to-one)**  $\longrightarrow$ 
  - Symmetric to 1:N
- **N:M (many-to-many)**
  - Each e1 relates to many e2 (0,1,...M)
  - Each e2 related to many e1 (0,1,...N)

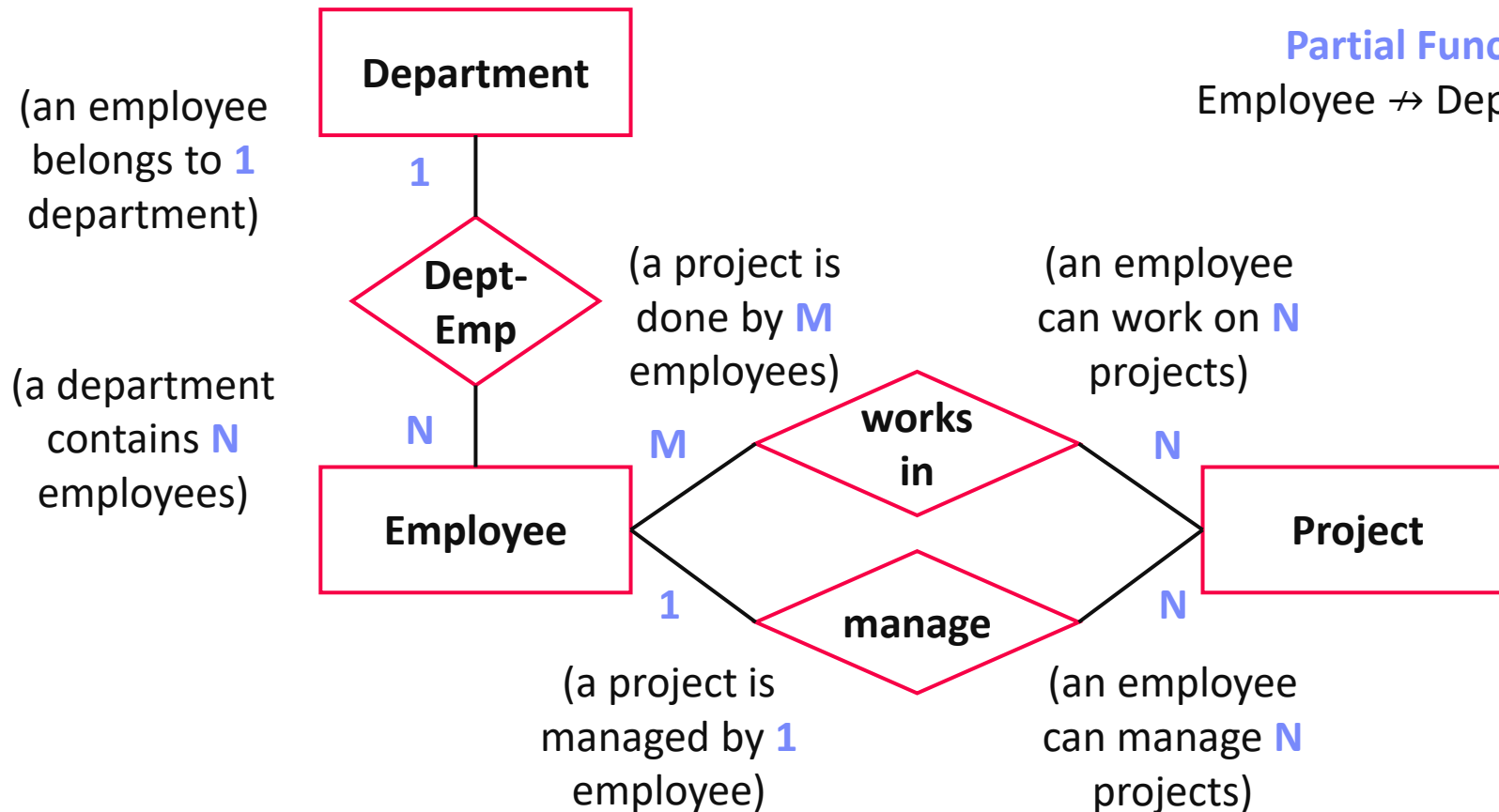


# An EmployeeDB Example, cont.

[Peter P. Chen: The Entity-Relationship Model - Toward a Unified View of Data. ACM Trans. Database Syst. 1(1) 1976]

Partial Function

Employee  $\rightarrow$  Department



## Multiplicity in Modified Chen Notation

- **Extension:** C (“choice”/“can”) to model 0 or 1, while 1 means exactly 1 and M means at least 1.

4 alternatives (1, C, M, MC)

→ 4\*4 = 16 combinations

(symmetric combinations omitted)

- **1:1** – [1] to [1]
- **1:C** – [1] to [0 or 1]
- **1:M** – [1] to [at least 1]
- **1:MC** – [1] to [arbitrary many]

1	1	1	1
0	1	1	1
0	0	1	1
0	0	0	1

$$\frac{n \cdot (n + 1)}{2}$$

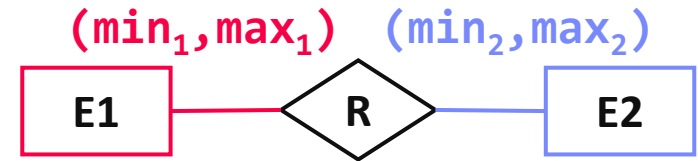
- **C:C** – [0 or 1] to [0 or 1] → see **1:1 in Chen**
- **C:M** – [0 or 1] to [at least 1]
- **C:MC** – [0 or 1] to [arbitrary many] → see **1:N in Chen**
- **M:M** – [at least 1] to [at least 1]
- **M:MC** – [at least 1] to [arbitrary many]
- **MC:MC** – [arbitrary many] to [arbitrary many] → see **M:N in Chen**



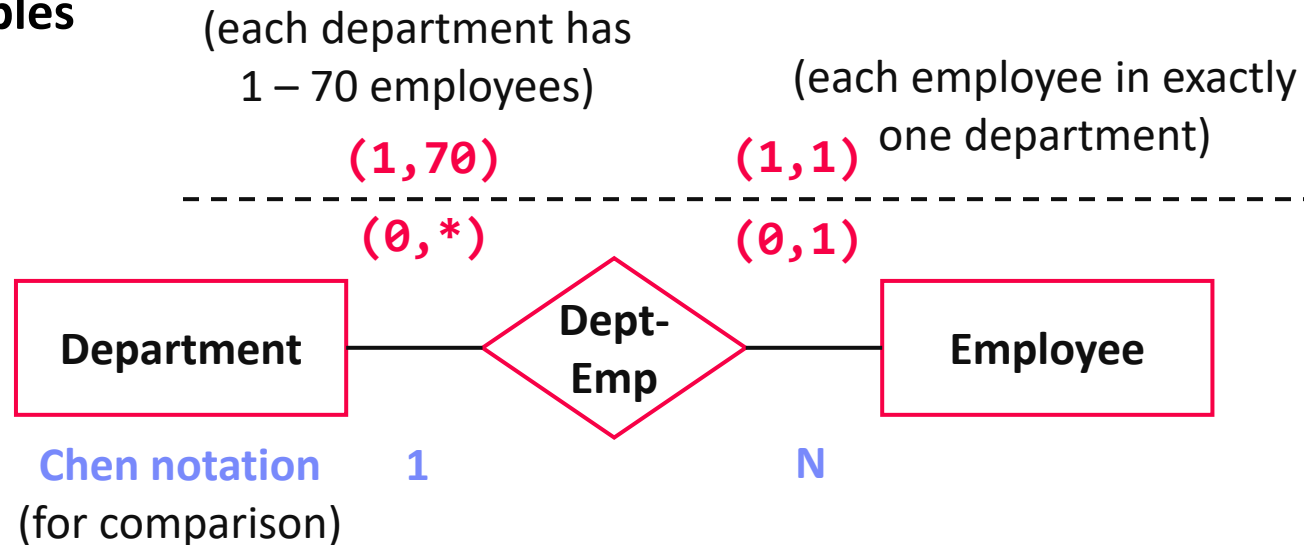
# (min,max)-Notation

- Alternative Cardinality Notation

- Indicate concrete min/max constraints (each entity is part of at least/at most x relationships)
- Chen and (min,max) notation generally incomparable
- Wildcard \* indicates arbitrary many (i.e., N)



- Examples

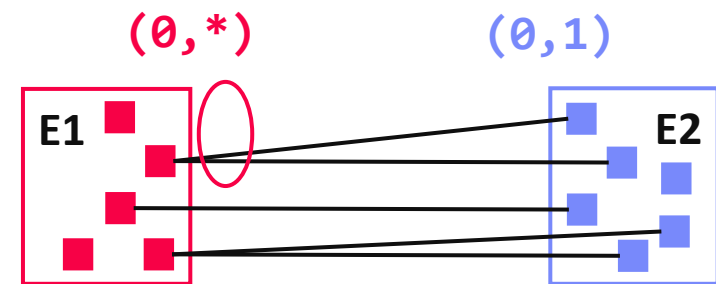


## (min,max)-Notation, cont.

- **Problem:** Where do these conflicting notations come from?

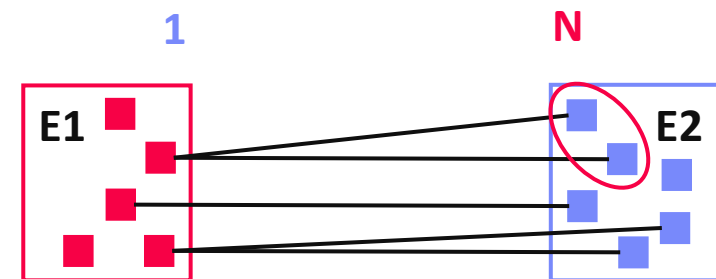
- **Understanding (min, max)-Notation**

- **Focus on relationships!**
- Describes number of outgoing relationships for each entity



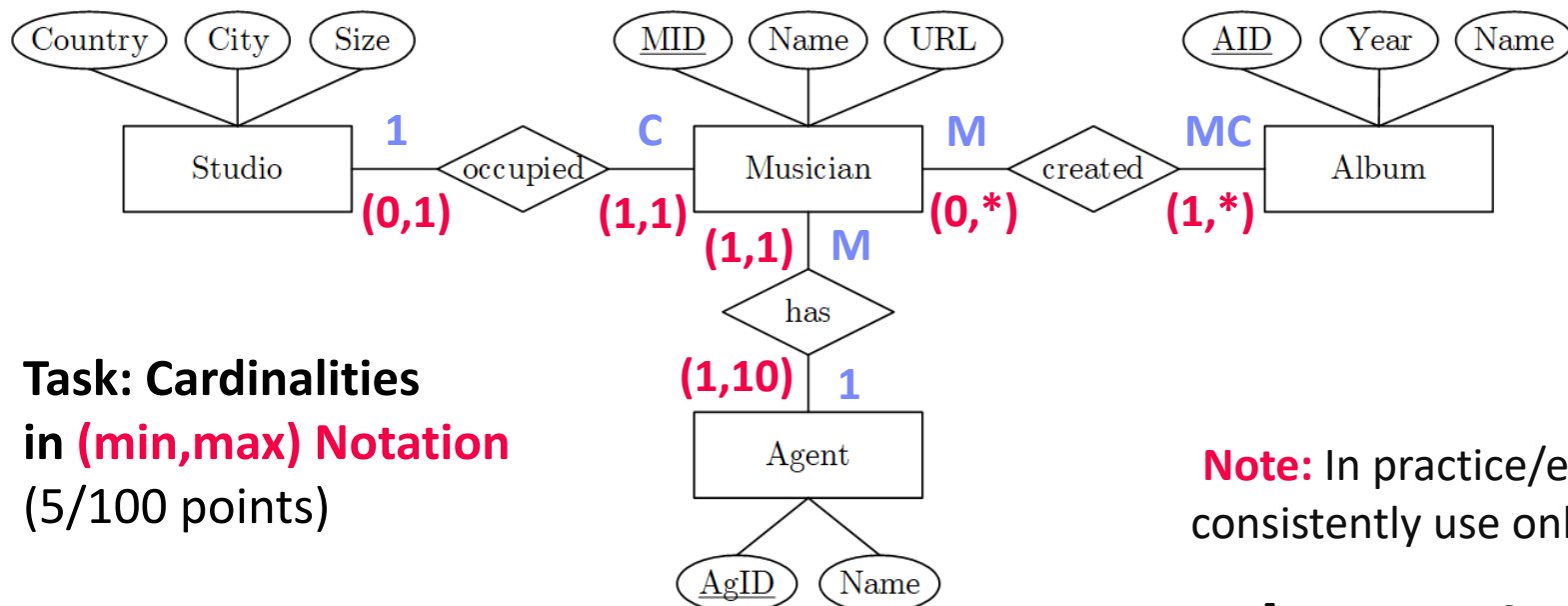
- **Understanding Chen- / Modified-Chen-Notation**

- **Focus on entities!**
- Describes number of target entities (over relationships) for each entity



# BREAK (and Test Yourself)

- Task: Cardinalities in Modified-Chen Notation** (prev. exam 6/100 points)
  - A musician might have created none or arbitrary many albums, and any album is created by at least one musician.
  - Every musician has exactly one agent, and an agent might be responsible for one to ten musicians.
  - Every musician occupies exactly one studio, and musicians never share a studio.



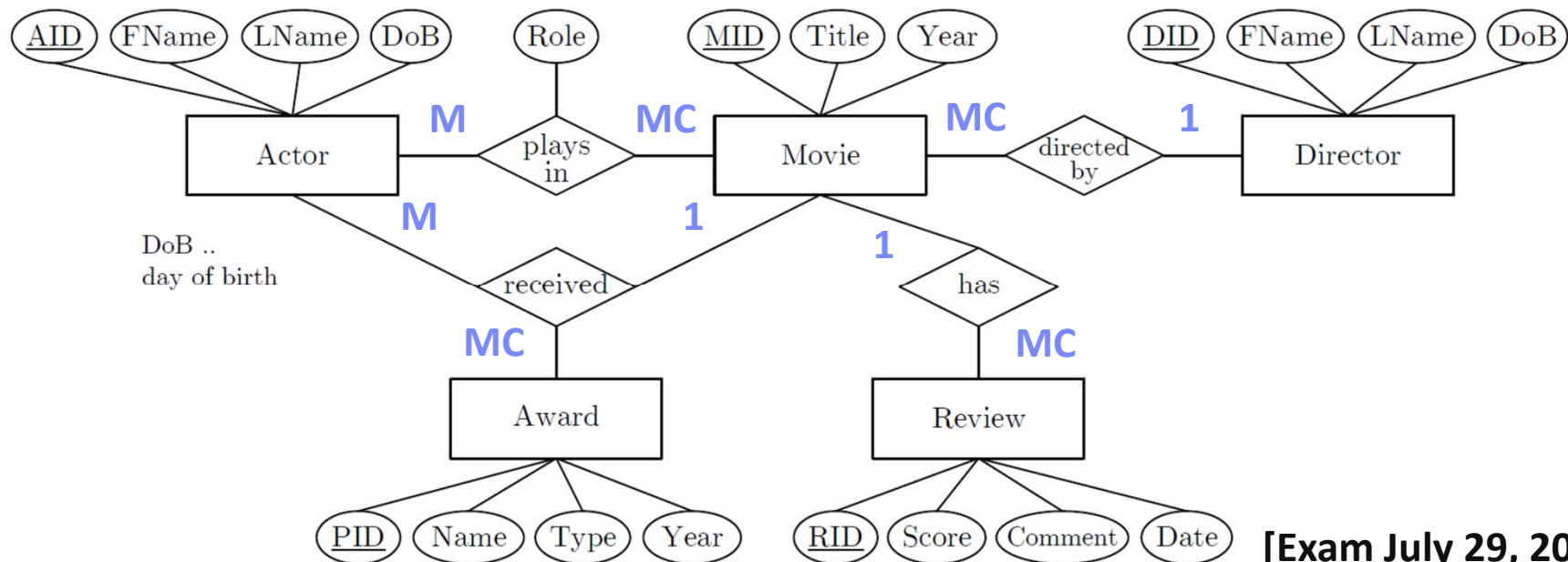
- Task: Cardinalities in (min,max) Notation** (5/100 points)

**Note:** In practice/exams, consistently use only one

[Exam June 24, 2019]

# BREAK (and Test Yourself), cont.

- Task: Cardinalities in Modified-Chen Notation** (prev. exam 9/100 points)
  - An actor may play roles in an arbitrary number of movies (incl. none), every movie has a cast of at least one but potentially many actors
  - A movie is directed by 1 director, directors produce arbitrary many movies
  - A movie review refers to 1 movie, but there can be 0-many reviews per movie
  - Actors (incl a single actor) may receive multiple awards for a single movie. An actor can receive only 1 per movie. Awards to 1-many actors are possible.



[Exam July 29, 2020]

# Weak Entity Types

## ■ Existence Dependencies

- Entities **E2** whose existence depends on the other entities **E1**
- Visualized as a special rectangle with double border
- Primary key is contains primary key of **E1**
- Relationship between strong and weak entity types **1:N** (sometimes **1:1**)

## ■ Examples

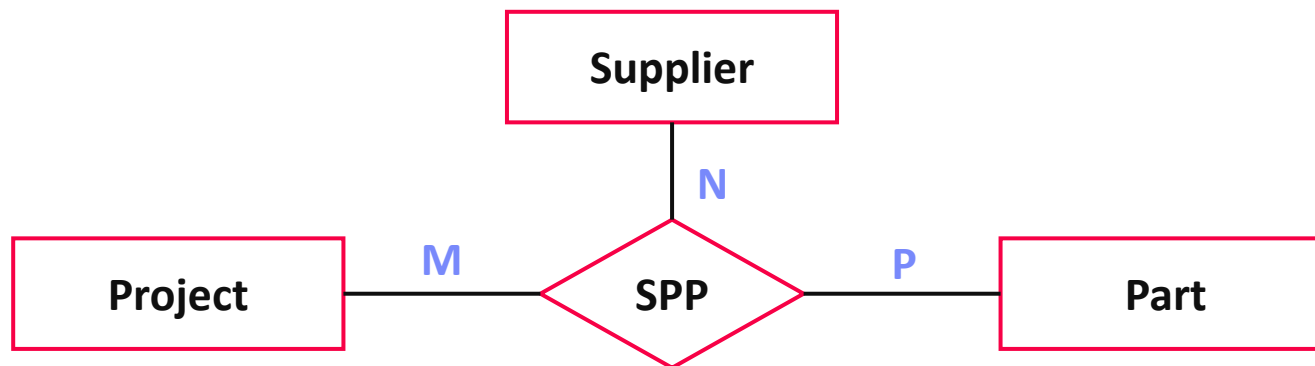
- Dependents of an employee (spouse, children)
- Rooms of a building



# N-ary Relationships

## ■ Use of n-ary relationships

- Relationship type among multiple entity types
- N-ary relationship can be converted to binary relationships
- Design choice: **simplicity** and **consistency constraints**



## ■ Multiplicity

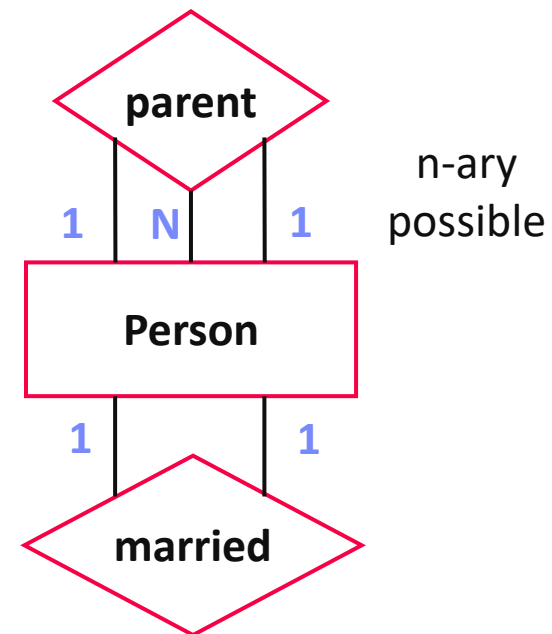
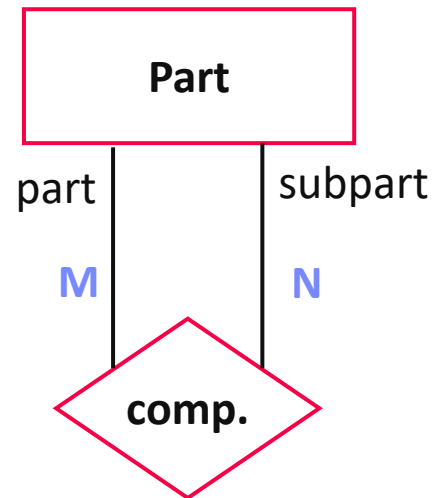
- 1 Project and 1 Supplier → supply **P** parts
- 1 Project and 1 Part → supplied by **N** suppliers (**1 instead of N?**)
- 1 Supplier and 1 Part → supply for **M** projects

# Recursive Relationships

## Definition

- Recursive relationships are relations between entities of the same type
- Use roles to differentiate cardinalities

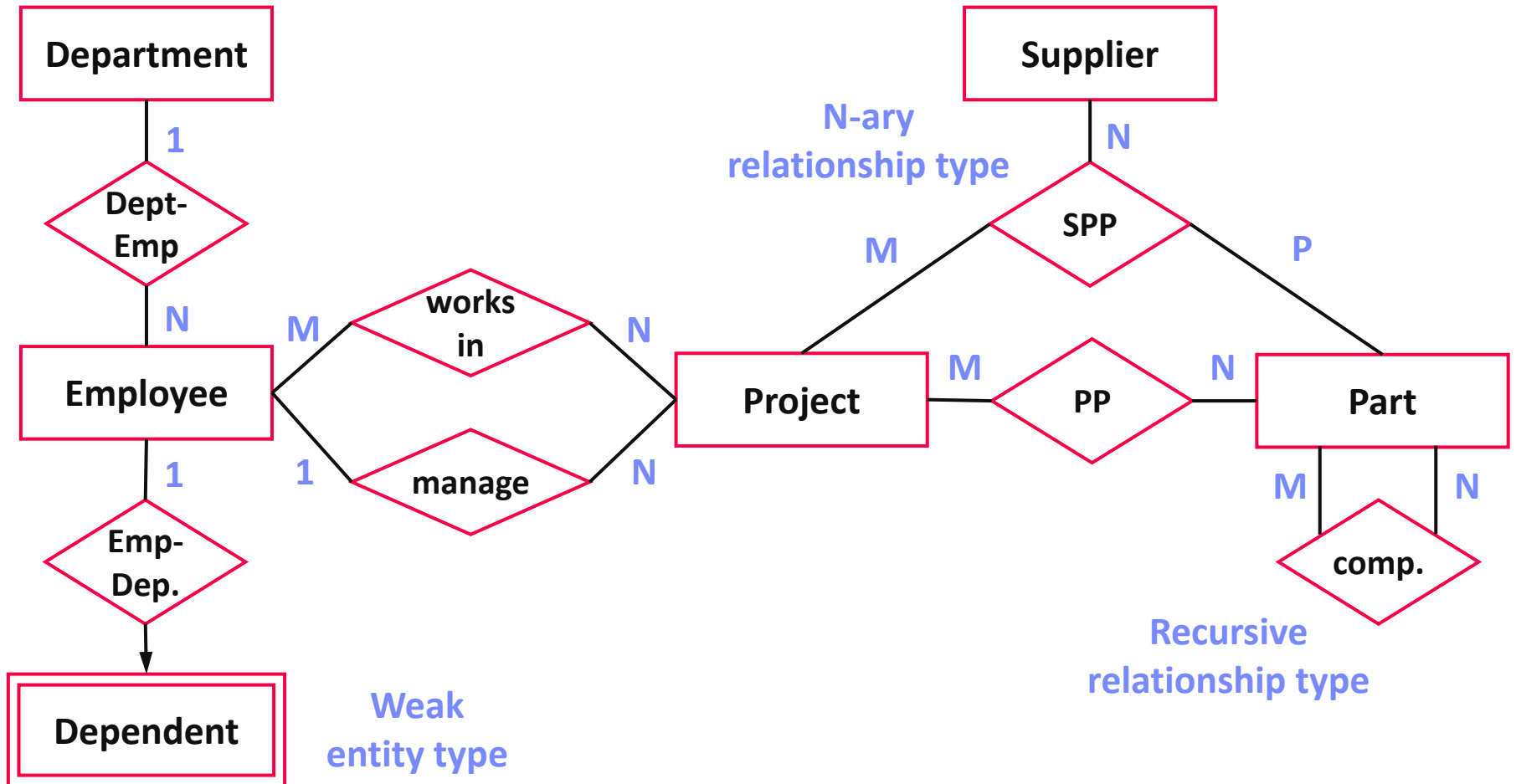
## Examples



- Beware of [at least 1] constraints in recursive relationships** (e.g., (min,max)-notation, or MC notation)

# An EmployeeDB Example, cont.

[Peter P. Chen: The Entity-Relationship Model - Toward a Unified View of Data. ACM Trans. Database Syst. 1(1) 1976]

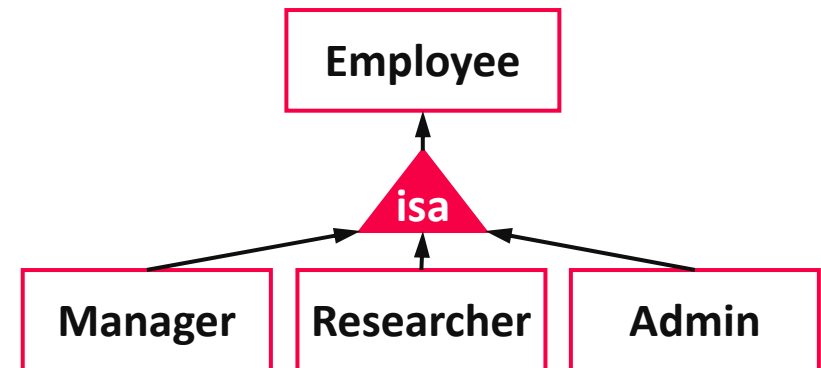




# Specialization and Aggregation

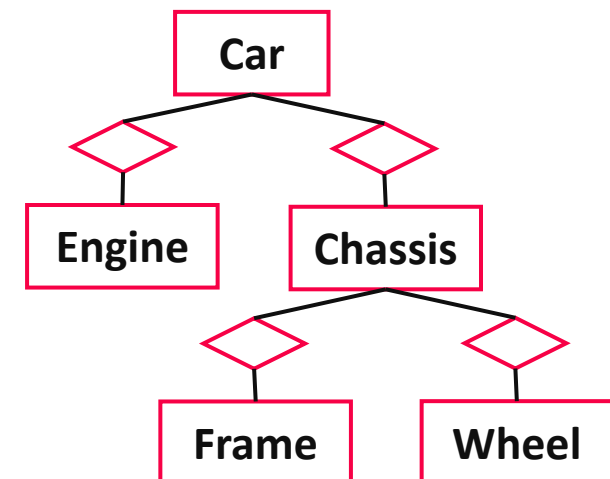
## Specialization via Subclasses

- **Tree of specialized entity types** (no multi-inheritance)
- Graphical symbol: triangle (or hexagon, or subset)
- Each entity of subclass is entity of superclass, but not vice versa



## Aggregation (composition, not specialization)

- **#1: Recursive relationship types**, or
- **#2: Explicit tree of entity** and relationship types
- Design choice: number of types known and finite, and heterogeneous attributes

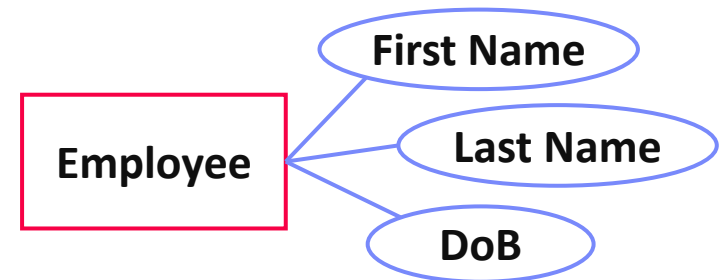


## Beware: **Simplicity is key**

# Types of Attributes

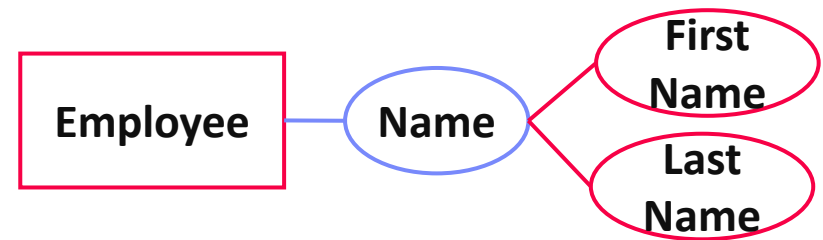
## Atomic Attributes

- Basic, single-valued attributes



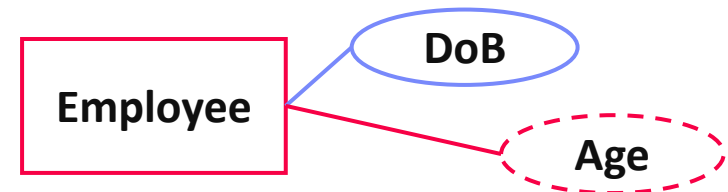
## Composite Attributes

- Attributes as structured data types
- Can be represented as a hierarchy



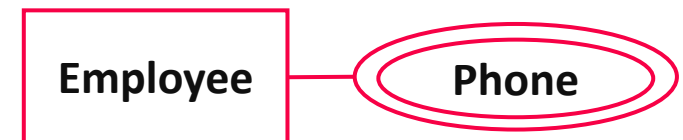
## Derived Attributes

- Attributes derived from other data
- Examples: Number of employees in dep, employee age, employee yearly salary



## Multi-valued Attributes

- Attributes with list of homogeneous entries






# Excursus: Influence of Chinese Characters?



*“What does the Chinese character construction principles have to do with ER modeling? The answer is: both Chinese characters and the ER model are trying to model the world – trying to use graphics to represent the entities in the real world. [...]”*

[Peter Pin-Shan Chen: Entity-Relationship Modeling: Historical Events, Future Trends, and Lessons Learned. **Software Pioneers 2002**]

- Chinese characters representing real-world entities

<u>Original Form</u>	<u>Current Form</u>	<u>Meaning</u>
	日	Sun
	月	Moon
	人	Person

- Composition of two Chinese characters

日 (sun) + 月 (moon) = 明 (Bright/ Brightness by light)

# Design Decisions

**Avoid redundancy**  
**Avoid unnecessary complexity**

- **Meta-Level:**

- Which notations to use (Chen, Modified Chen, (min,max)-notation)?

- **Entities**

- What are the entity types (entity vs relationship vs attribute)?
- What are the attributes of each entity type?
- What are key attributes (one or many)?
- What are weak entities (with partial keys)?

- **Relationships**

- What are the relationship types between entities (binary, n-ary)?
- What are the attributes of each relationship type?
- What are the cardinalities?

- **Attributes**

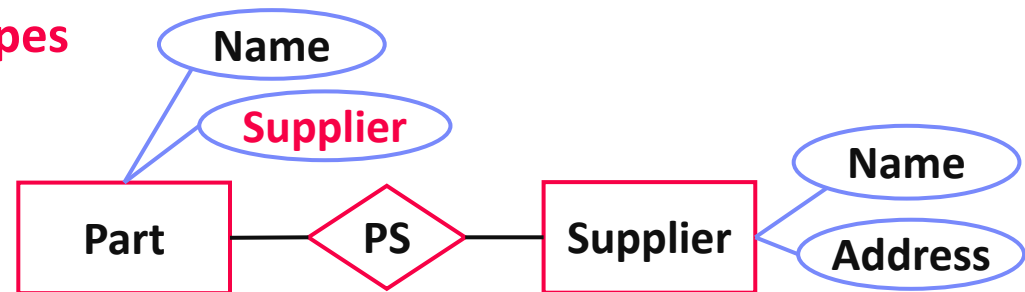
- What are composite, multi-valued, or derived attributes?

# Design Decisions – Examples of **Poor** Choices

- #1 Overuse of **weak entity types**

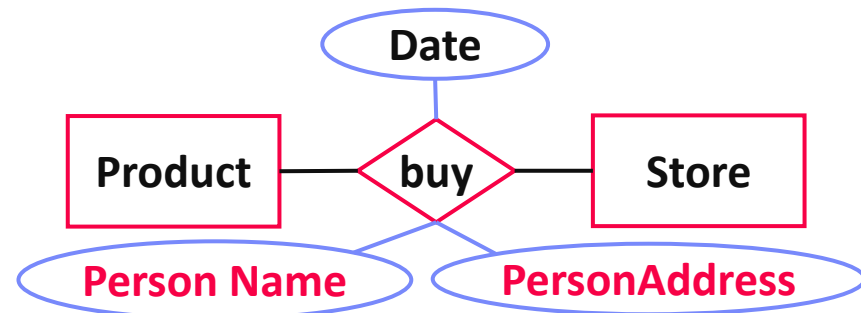
- #2 Redundant attributes

- Redundant supplier name** in Part and Supplier



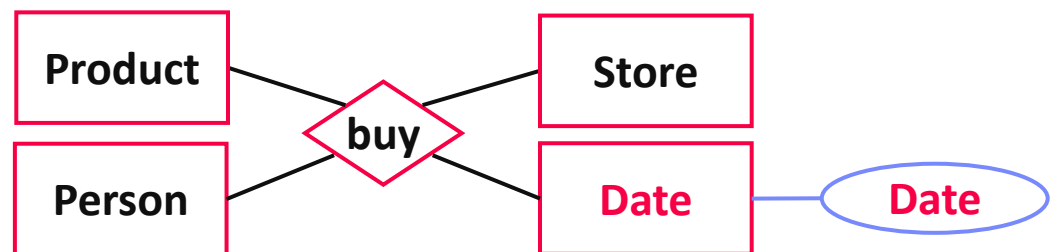
- #3 Repeated information

- Missing person entity type**  
→ redundancy per purchase



- #4 Unnecessary Complexity

- Unnecessary entity type Date**
  - Avoid single-attribute entity types unless in many relationships



# A UniversityDB Example

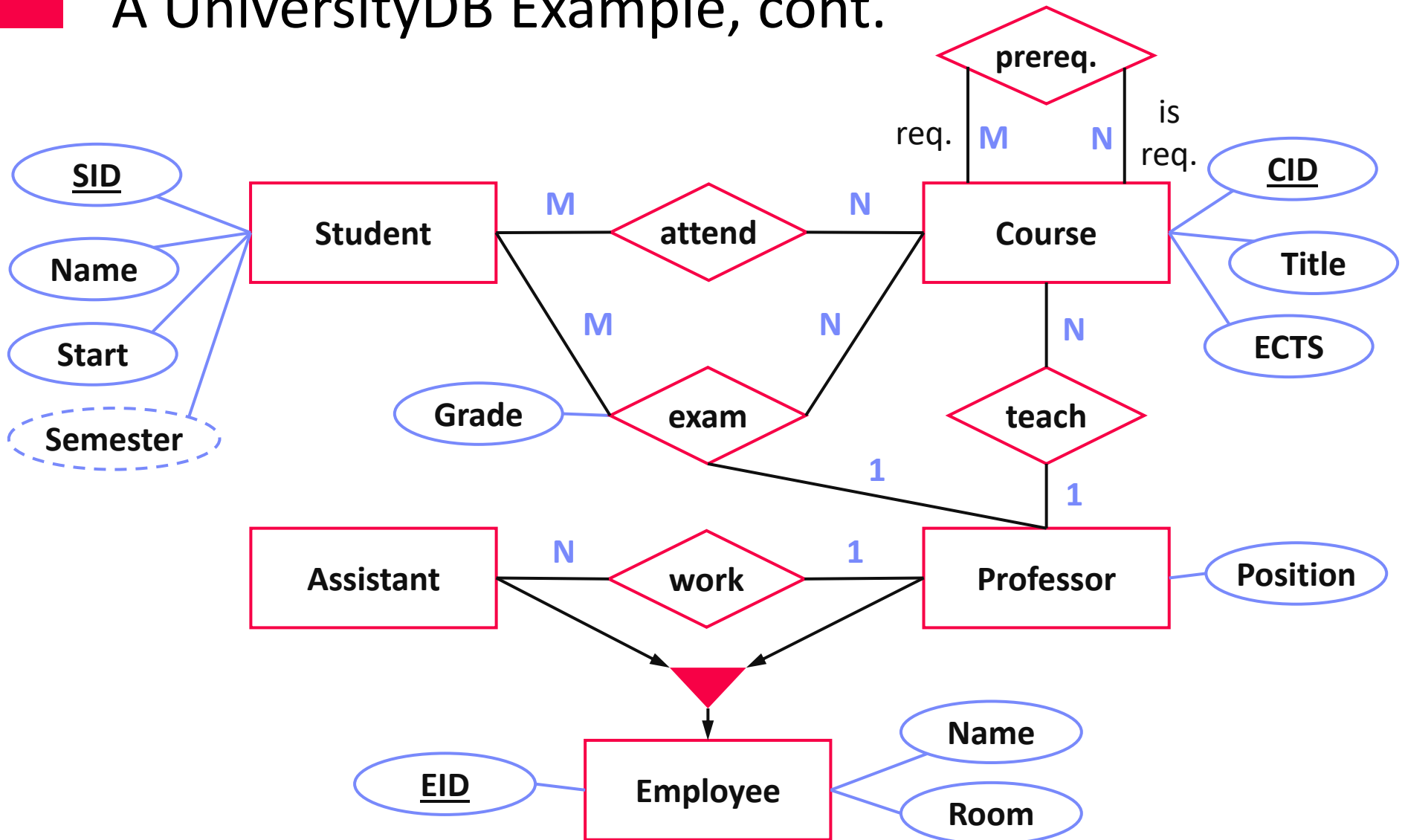
## ▪ Discourse of Real Mini World

- **Students** (with SID, name, and semester) attend **courses** (CID, title, ECTS), and take graded exams per course
- **Professors** teach courses and have positions, **assistants** work for professors
- A course may have another course as prerequisites
- Both professors and assistants are university **employees** (EID, name, and room number); professors also have a position

## ▪ Task: **Create an ER diagram in Chen notation**

- Include entity types, relationship types, attributes, and generalizations
- Mark primary keys, roles for recursive relationships, and derived attributes

# A UniversityDB Example, cont.



# Exercise 01 – Data Modeling

Published: **Oct 12, 2020**

Deadline: **Nov 03, 2020**



# Exercises: The Movies Dataset

## Dataset

- Derived (extracted, cleaned) from **The Movies Dataset** for movies year  $\geq 2011$
- Note: Still in process of data cleaning**
- Clone or download your copy from <https://github.com/tugraz-isds/datasets.git>
- Find CSV files in <datasets>/movies

## Exercises

- 01** Data modeling (relational schema)
- 02** Data ingestion and SQL query processing
- 03** Physical design tuning, query processing, and transaction processing
- 04** Large-scale data analysis (distributed query processing and ML model training)

### Movies

#### Overview

The following dataset was derived from a larger movie dataset - [The Movies Dataset](#). The original dataset contains metadata of over 45,000 movies as well as 26 million ratings from 270,000 users for all of the movies, which have been collected from [TMDB](#) and [GroupLens](#), respectively. From this large collection, data of around 12,000 movies from the past 10 years were extracted and stored as three denormalized CSV files (with ; delimiter and simplified structure without the need for quoting): `Movies.csv`, `Persons.csv` (cast members) and `Ratings.csv` (given to the movies by various users).

#### Structure:

`Movies.csv`: The movies file contains metadata on over 12,000 different movie titles. The datapoints included in the file are:

- `MovieID`: An ID that uniquely identifies each movie.
- `OriginalLanguage`: The language of the original movie title, denoted in its [ISO 639-1](#) language code equivalent.
- `OriginalTitle`: The title of the movie in its original language.
- `EnglishTitle`: The english equivalent of the original title.
- `Budget`: The amount of money invested into making the movie.
- `Revenue`: The amount of money generated by the movie.
- `Homepage`: A link to the movies website.
- `Runtime`: The duration of the movie in minutes.
- `ReleaseDate`: The date on which the movie was/will be released.
  - Format: yyyy-mm-dd
- `Genres`: A list of genres that the movie is categorized under.
  - Format: genre|genre2|...|genreN
- `CastID`: A list of 24 character long IDs, belonging to the movies cast members.
  - Format: castid1|castid2|...|castidN
- `ProductionCompanies`: A list of production companies involved with the movie.
  - Format: company1|company2|...|companyN
- `ProductionCountries`: A list of countries in which the movie was filmed, paired with their [ISO 3166-1](#) country code.
  - Format: code1+country1|code2+country2|...|codeN+countryN
- `SpokenLanguages`: A list of languages spoken in the movie, paired with their [ISO 639-1](#) language code.
  - Format: code1+language1|code2+language2|...|codeN+languageN

### Movies.csv

The following is an excerpt from the `Movies.csv` file which is representative of the dataset's structure:

```
MovieID,OriginalLanguage,OriginalTitle,EnglishTitle,Budget,Revenue,Homepage,Runtime,ReleaseDate,Genres,CastID,ProductionCompany
136558,en,Kingdom Come,Kingdom Come,,http://www.kingdomcomefilm.com/#kingdomcome,88,2011-09-01,Comedy,...
118423,fr,Camille Claudel 1915,Camille Claudel 1915,3512454,115086,,95,2013-03-13,Drama,Canada|Arte France Cinéma|3B Product
62775,fi,Havukka-Ahon Ajattelijat,Havukka-Ahon Ajattelijat,2223000,,2018-01-15,Comedy|Drama,,,fi-suomi
12477,en,When in Rome,When in Rome,,36699403,,91,2018-01-20,Fantasy|Comedy|Romance,Krasnoff Foster Productions|Touchstone|P
12081,en,Edge of Darkness,Edge of Darkness,80000000,74001339,,117,2010-01-29,Crime|Drama|Mystery|Thriller,,Icon Productions|I
37034,sh,Su Qi-Er,True Legend,20000000,,115,2010-02-09>Action|Fantasy,,Shanghai Film Group|Focus Features|EKO Film|Eko Mar
```

`Persons.csv`: The persons file contains information about the cast members and their characters in the movies. The datapoints included in the file are:

- `MovieID`: An ID referencing the movie in which the character appeared.
- `CastID`: An ID that uniquely identifies each character played by the cast member.
- `Name`: The name of the cast member.
- `Gender`: The gender of the cast member (1 = female, 2 = male)
- `Character`: Full name of the movie's character.

### Persons.csv (Cast)

The following is an excerpt from the `Persons.csv` file which is representative of the dataset's structure:

```
MovieID,CastID,Name,Gender,Character
136558,52fedc18c3a368484e1a6d23,Daniel Gillies,2,Himself
136558,52fedc18c3a368484e1a6d27,Rachel Leigh Cook,1,Meriel
118423,52fedc06c3a36847f81e4625,Juliette Binoche,1,Camille Claudel
118423,52fedc06c3a36847f81e4629,Jean-Luc Vincent,Paul Claudel
62775,52fedc06c3a368484e09786f,Kari Lehtinen,2,Konsta Pykkönen
62775,52fedc06c3a368484e097873,Tomi Korpela,1,Misteri Kronberg
```

`Ratings.csv`: The ratings file contains information about the ratings given to the movies by different users. The datapoints included in the file are:

- `UserID`: An ID identifying the user that published the rating.
- `MovieID`: An ID referencing the movie that had been rated.
- `Rating`: The user's rating of the movie on a scale from 1.0 to 5.0.
- `Timestamp`: Timestamp at which the user's rating was published.

### Ratings.csv

The following is an excerpt from the `Ratings.csv` file which is representative of the dataset's structure:

```
UserID,MovieID,Rating,Timestamp
1,58559,4,0,1425942007
1,58621,5,0,1425941392
7,58559,5,0,1486233675
7,68744,1,5,1486233974
11,58559,4,5,1216707182
15,48539,3,5,1346802547
```

# Overview Exercise 1 Tasks

[[https://mboehm7.github.io/teaching/ws2021\\_dbs/01\\_ExerciseModeling.pdf](https://mboehm7.github.io/teaching/ws2021_dbs/01_ExerciseModeling.pdf)]

- **Task 1.1: ER Modeling (movies, cast, ratings)**
  - Create an ER diagram in Modified Chen (MC) notation
  - <https://github.com/tugraz-isds/datasets/tree/master/movies>
  
- **Task 1.2: Mapping ER Diagram into Relational Model**
  - Create a relational schema in 3NF for the ER diagram from Task 1.1
  - **FDs:** LangCode → Language, CountryCode → Country, multi-valued attributes
  - a) text-based schema, **OR** b) SQL DDL script
  
- **Task 1.3: Relational Normalization**
  - Explain why the schema from Task 1.2 is in third normal form (3NF)
  
- **Expected result (for all three subtasks)**
  - **DBExercise01\_<studentID>.zip**



**Don't get your own  
studentID wrong**

# Summary and Q&A

## ■ Summary

- DB Design lifecycle from requirements to physical design
- Entity-Relationship (ER) Model and Diagrams

## ■ Importance of Good Database Design

- Poor database design → **development and maintenance costs**, as well as performance problems
- Once data is loaded, **schema changes very difficult** (data model, or conceptual and logical schema)

## ■ Exercise 1: Data Modeling

- Published Oct 12, 2020; deadline: Nov 03, 2020
- **Recommendation:** start with task 1.1 this week; ask questions in upcoming lectures or on news group

## ■ Next lecture: **03 Data Models and Normalization** [Oct 19]