

Data Management

02 Conceptual Design

Matthias Boehm

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

Last update: Mar 09, 2020

Announcements/Org

■ #1 Video Recording

- Link in [TeachCenter](#) & [TUBE](#) (lectures will be public)



■ #2 Course Registrations SS20

- Data Management (lectures/exercises): **490/485**
- Databases (combined lectures/exercises): **97**

Total:

587

■ #3 CS Talks x7 (**Mar 10, 5pm**, Aula Alte Technik)

- [Claudia Müller-Birn](#) (Freie Universität of Berlin)
- Title: [Collaboration is Key – Human-Centered Design of Computational Systems](#)



■ #4 Study Abroad Fair (**Mar 18, 10am-3pm**, INF 25d HS i4)

- Info booths and short presentations on study abroad programs (e.g., exchange, research, summer)



Announcements/Org, cont.

- #5 Catalyst Coding Contest (**Apr 03, 3-8pm**)
 - Hosted by: **IT Community Styria**
 - Online or in-person (teams/individuals)
 - INF 18, HS i1 (117 seats)
 - <https://register.codingcontest.org/>



Agenda

- DB Design Lifecycle
- ER Model and Diagrams
- Exercise 01 – Data Modeling (preview)



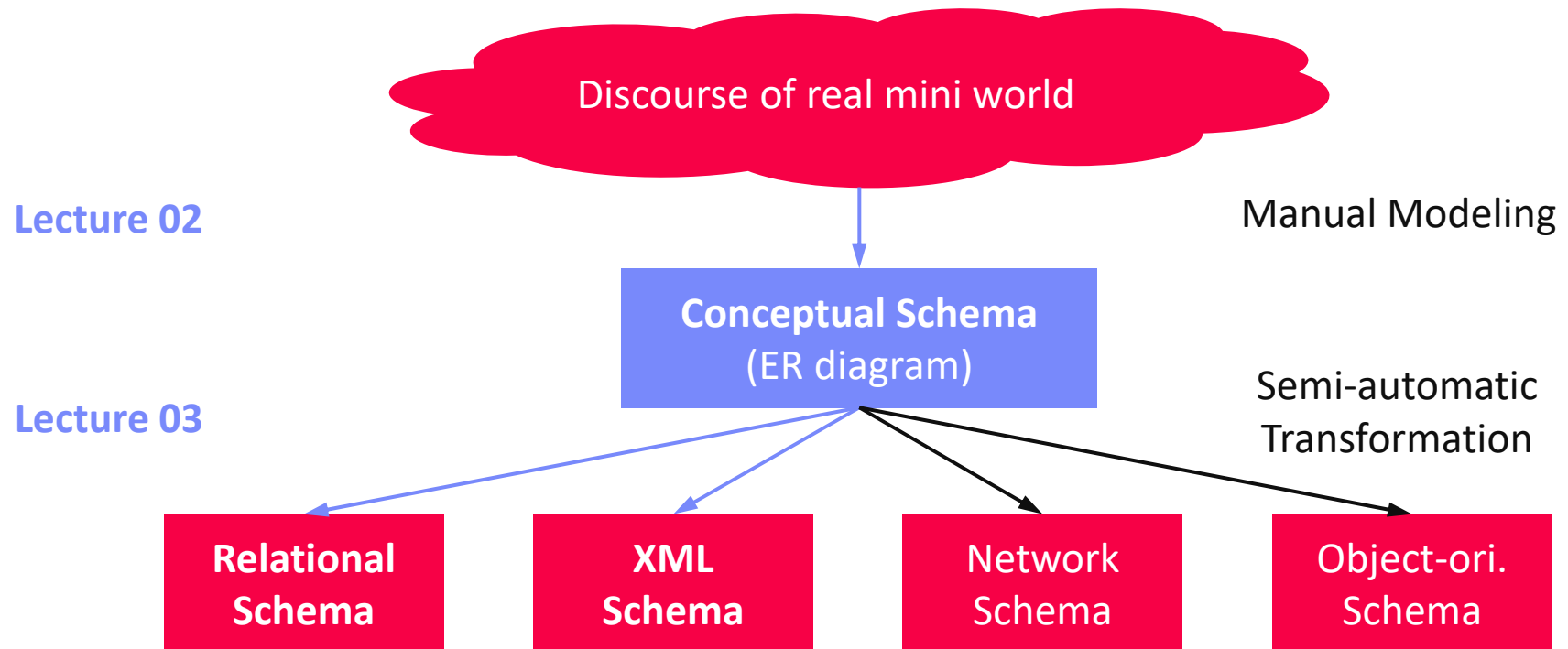
[**Credit:** Alfons Kemper, André Eickler: Datenbanksysteme - Eine Einführung, 10. Auflage. De Gruyter Studium, de Gruyter Oldenbourg 2015, ISBN 978-3-11-044375-2, pp. 1-879]

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

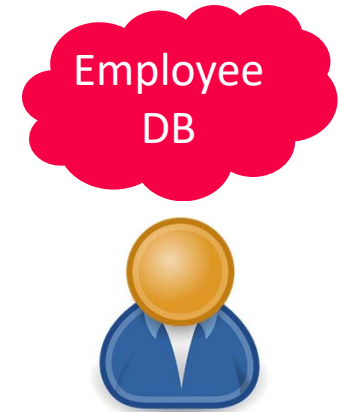
Partly covered
in part B

- 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]

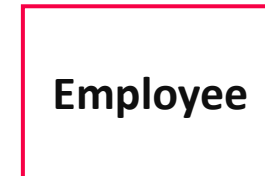
[Peter P. Chen: The Entity-Relationship Model: Toward a Unified View of Data. **VLDB** 1975]



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.

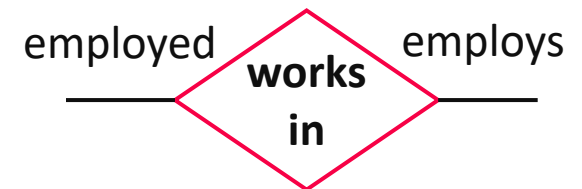
■ Keys

- Attributes that uniquely identify an entity
- Every entity type must have such a key
- Natural or surrogate (artificial) keys



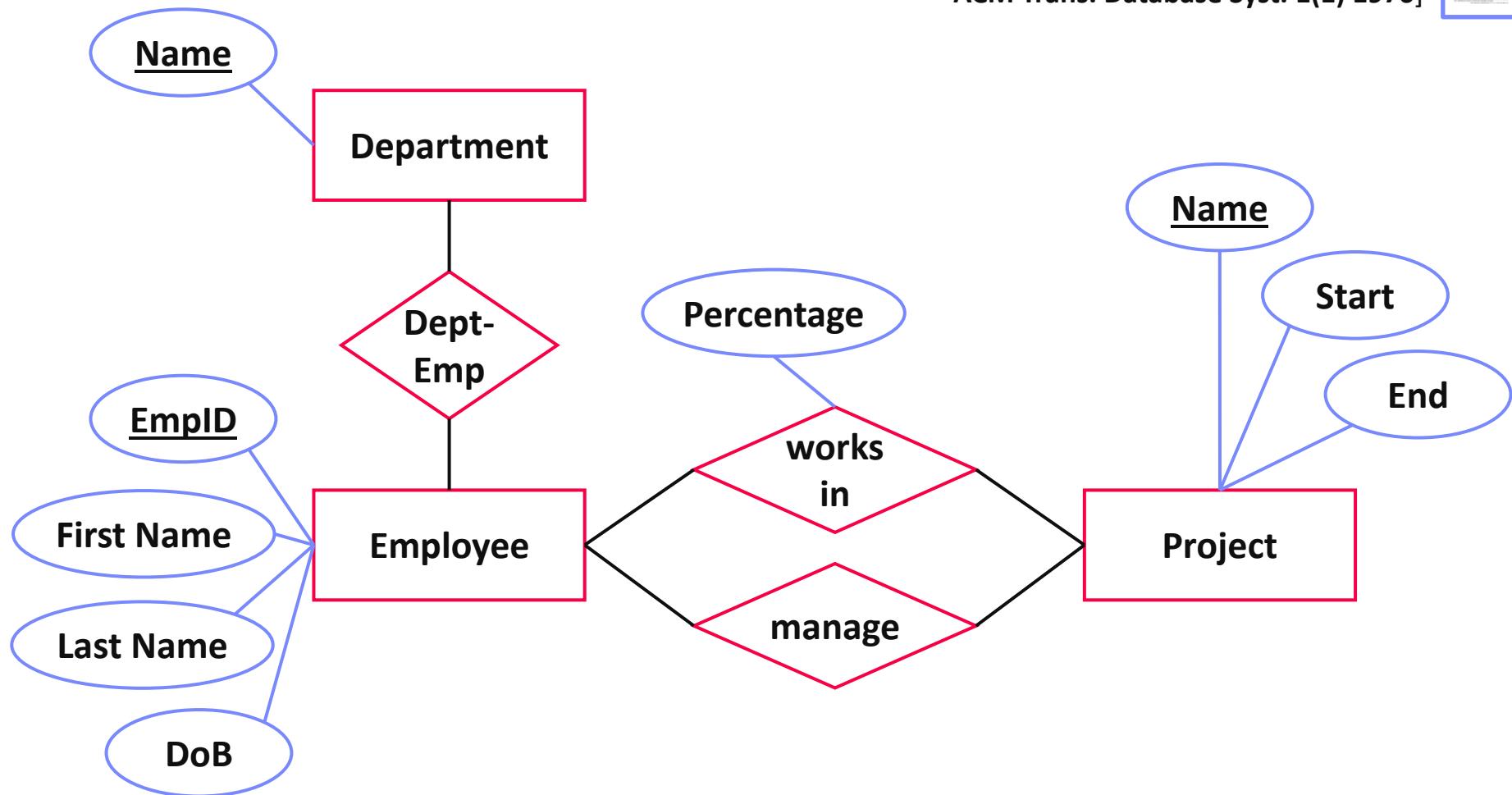
■ Role

- Optional description of relationship types
- Useful for recursive relationships



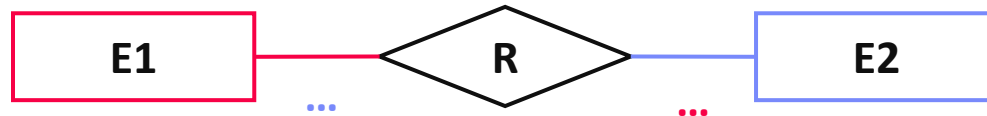
An EmployeeDB Example

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



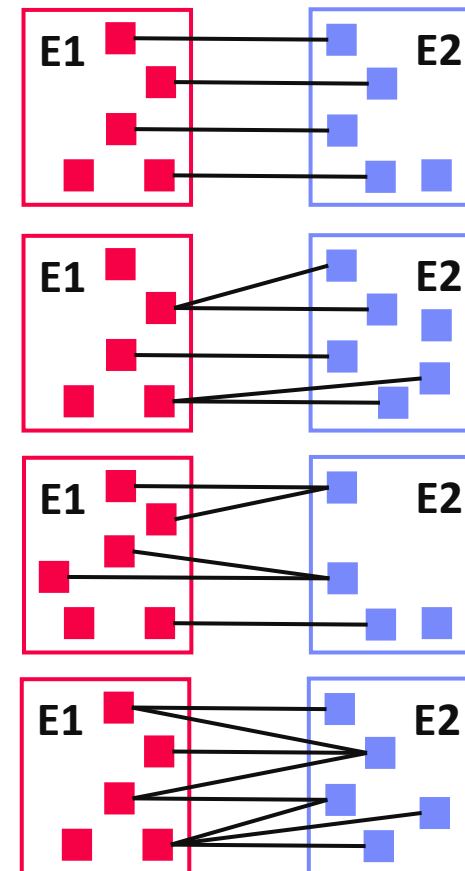
Multiplicity/Cardinality in Chen Notation

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



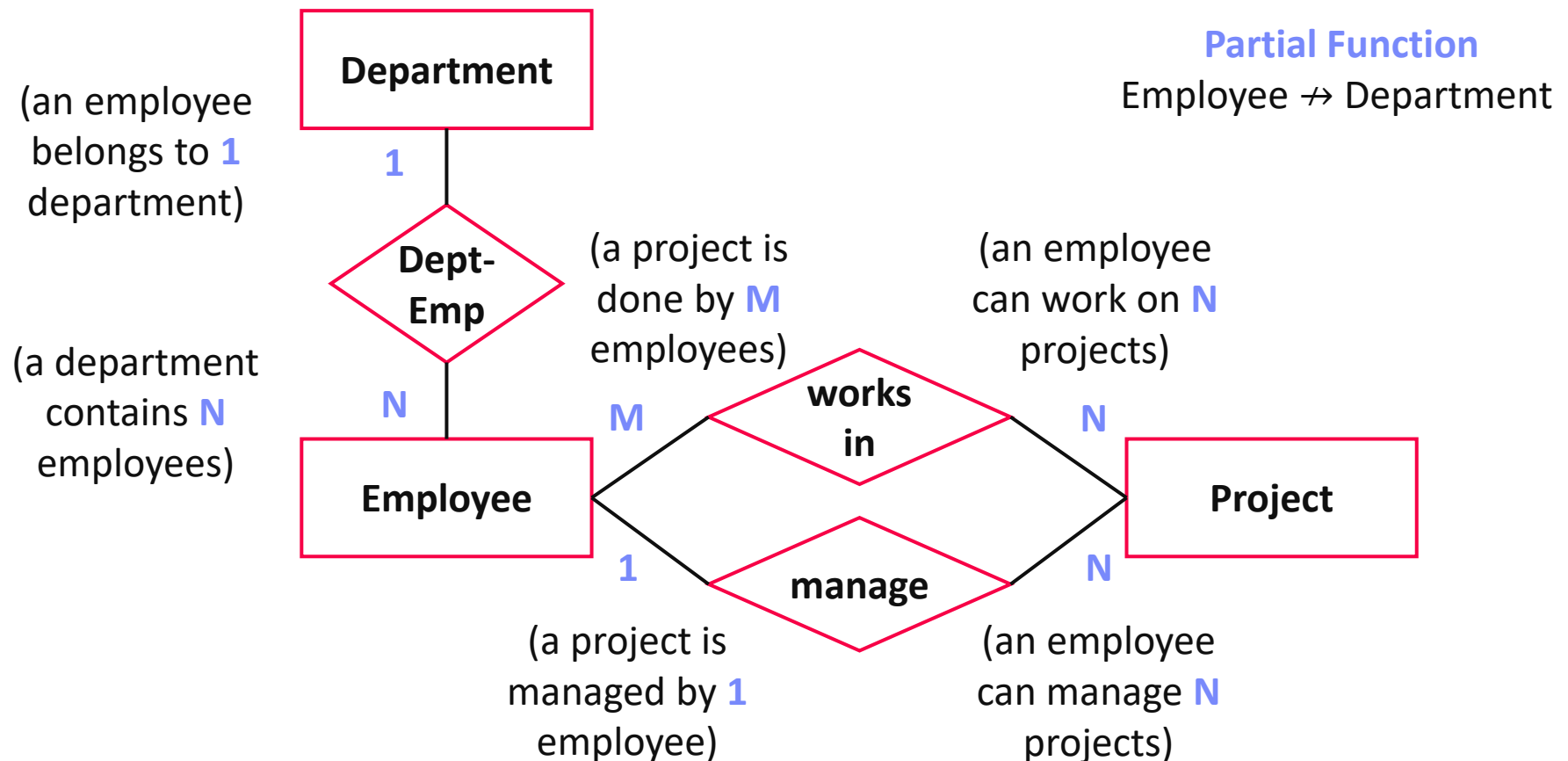
$$R \subseteq E1 \times E2$$

- **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]



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

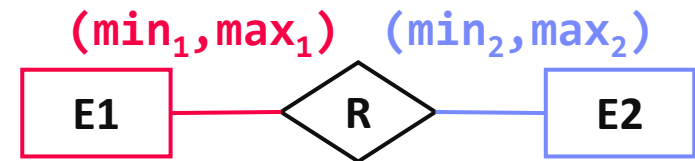
$$\frac{n \cdot (n + 1)}{2}$$

- **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**

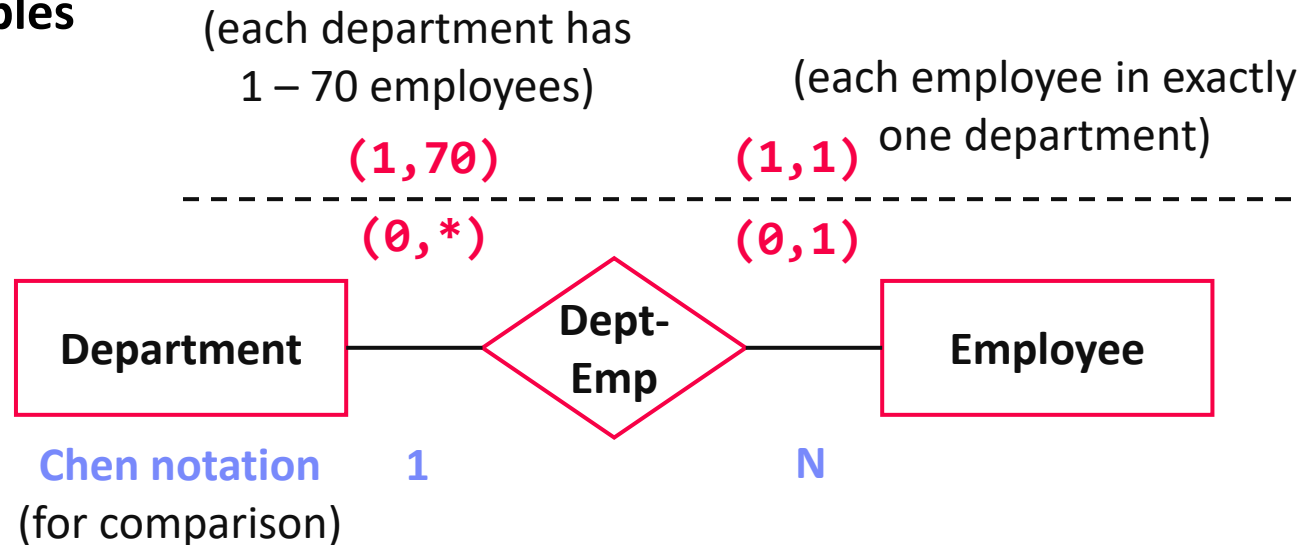
(min,max)-Notation

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

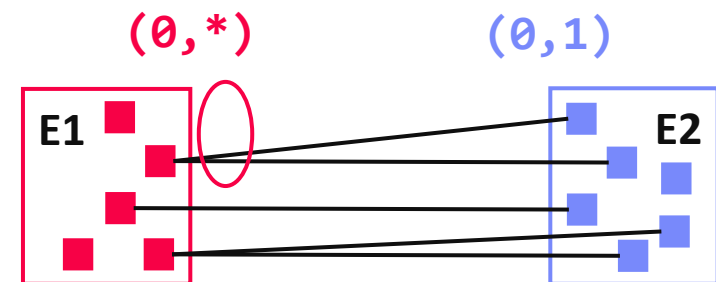


(min,max)-Notation, cont.

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

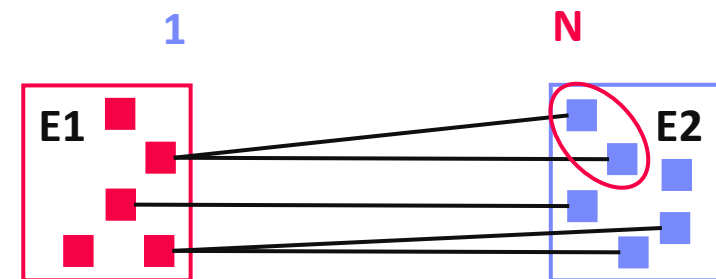
- **Understanding (min, max)-Notation**

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



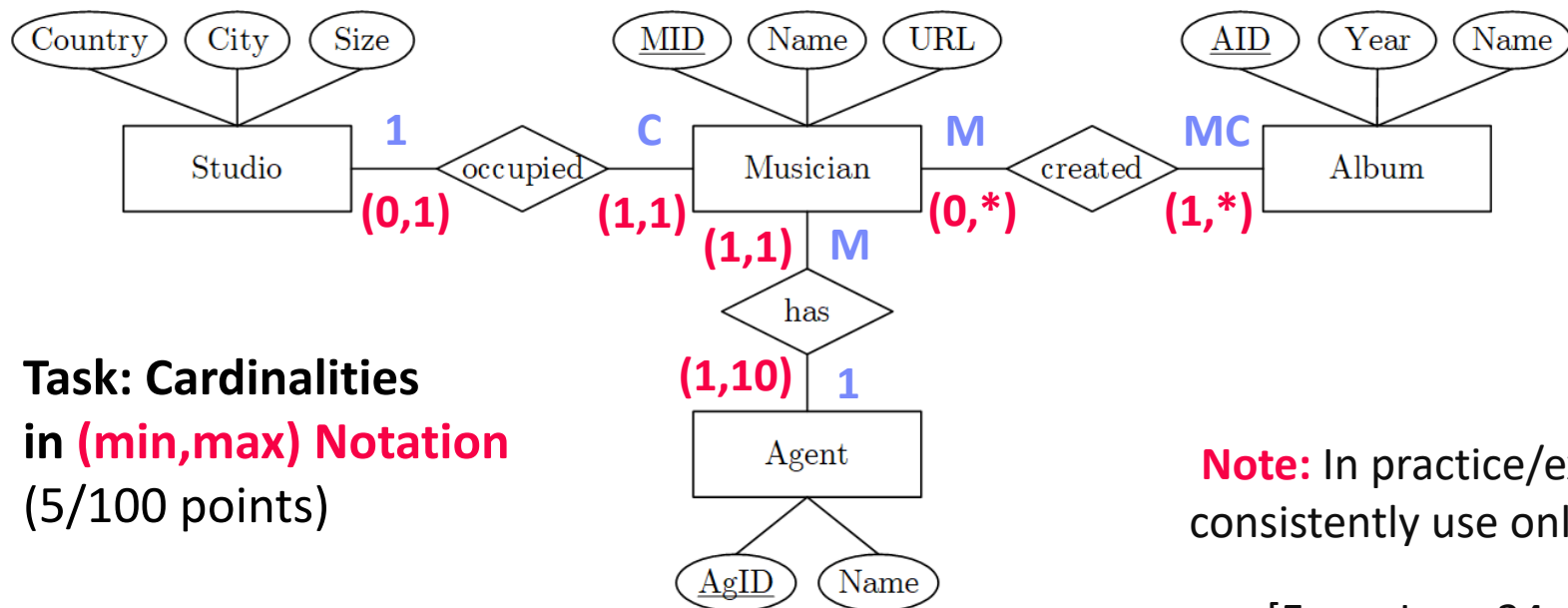
- **Understanding Chen- / Modified-Chen-Notation**

- 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]

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

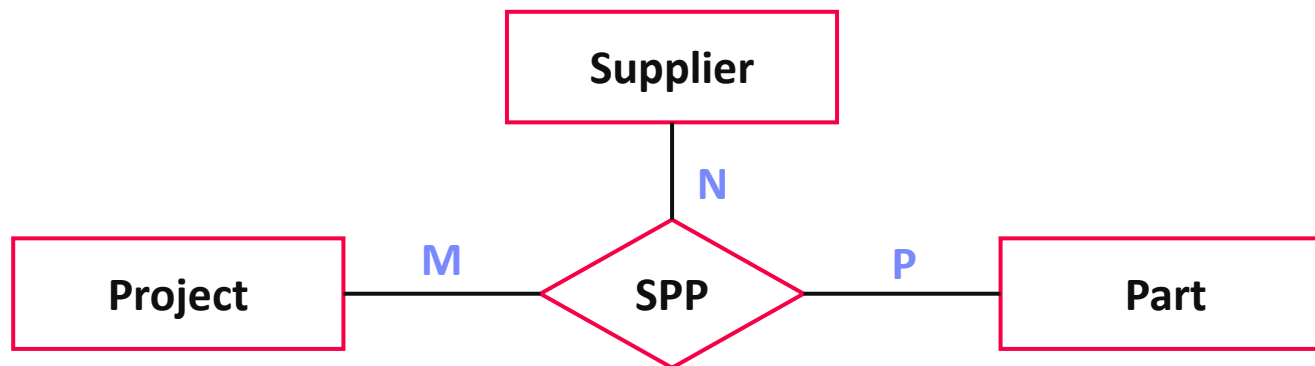
- Dependents of an employee (spouse, children)
- Rooms of a building



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

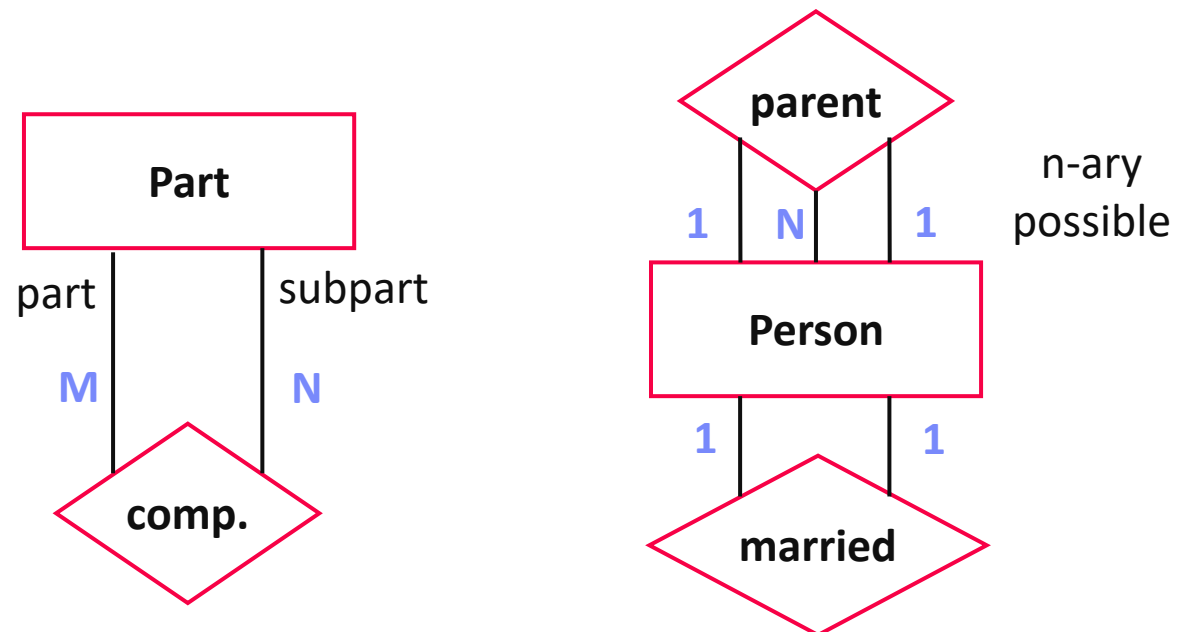
- 1 Project and 1 Supplier → supply **P** parts
- 1 Project and 1 Part → supplied by **N** suppliers (**1 instead of N?**)
- 1 Supplier and 1 Part → supply for **M** projects

Recursive Relationships

■ Definition

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

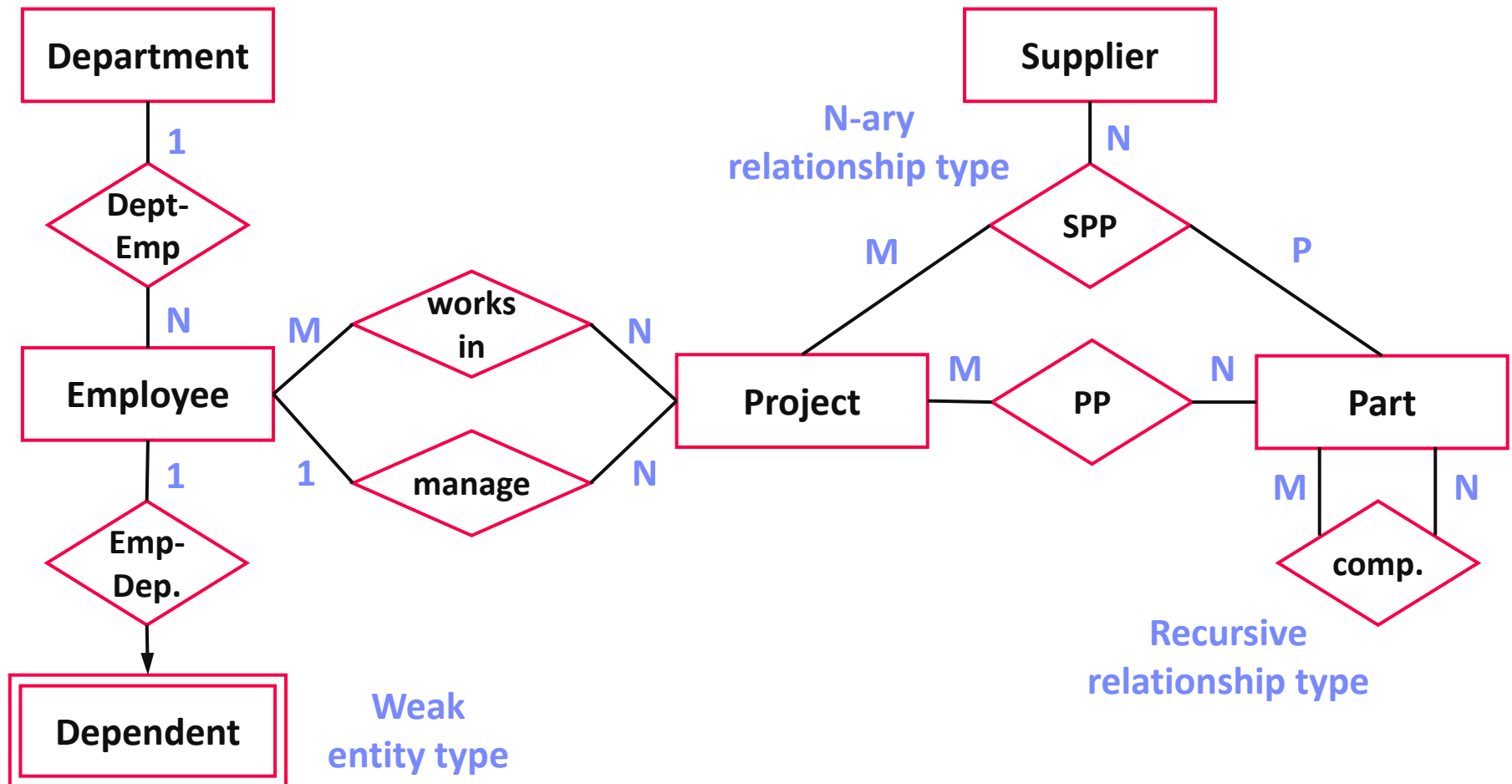
■ Examples



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

An EmployeeDB Example, cont.

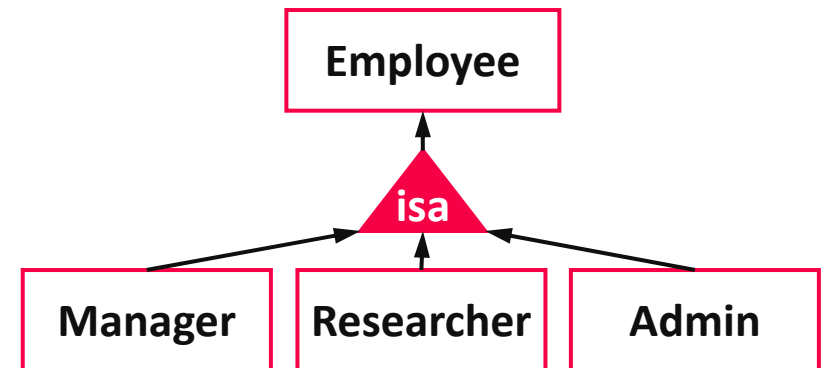
[Peter P. Chen: The Entity-Relationship Model - Toward a Unified View of Data.
ACM Trans. Database Syst. 1(1) 1976]



Specialization and Aggregation

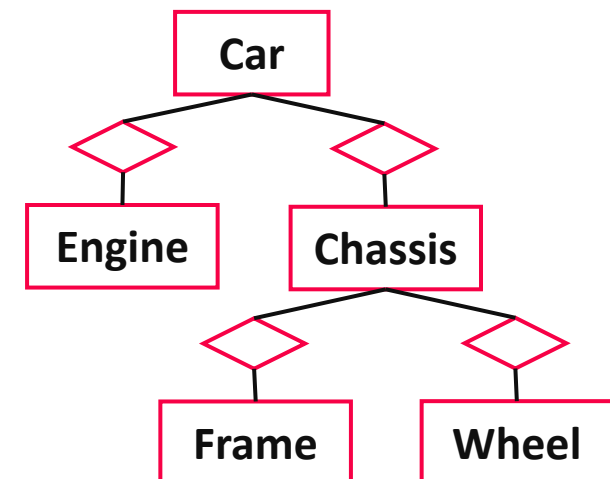
■ Specialization via Subclasses

- **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

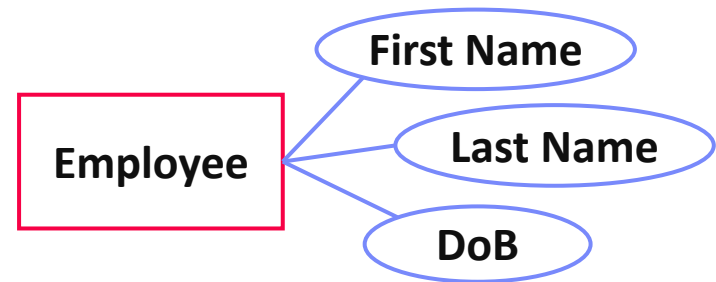


■ Beware: **Simplicity is key**

Types of Attributes

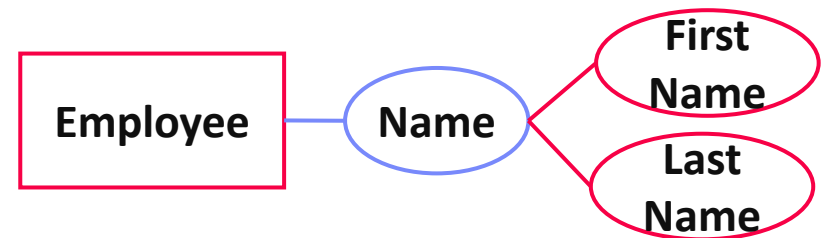
Atomic Attributes

- Basic, single-valued attributes



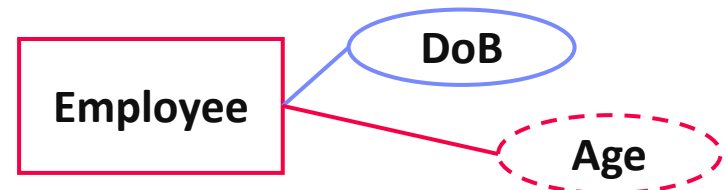
Composite Attributes

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



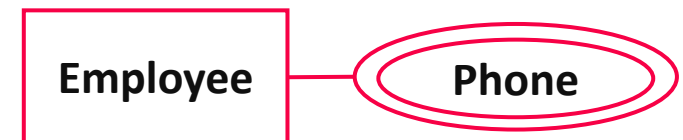
Derived Attributes

- 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>
		Sun
		Moon
		Person

- Composition of two Chinese characters

 (sun) +  (moon) =  (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

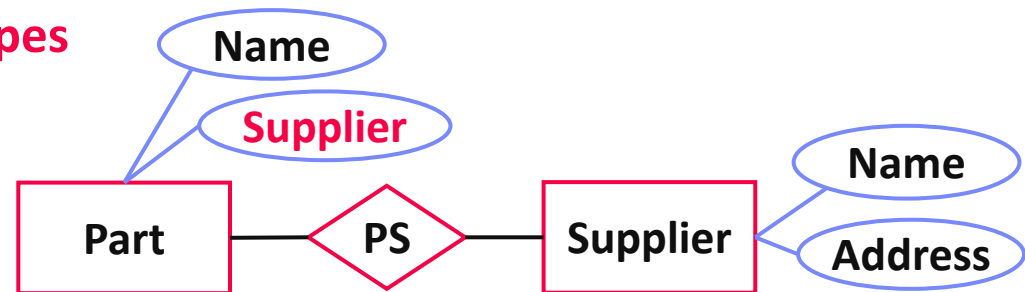
- What are composite, multi-valued, or derived attributes?

Design Decisions – Examples of **Poor** Choices

■ #1 Overuse of **weak entity types**

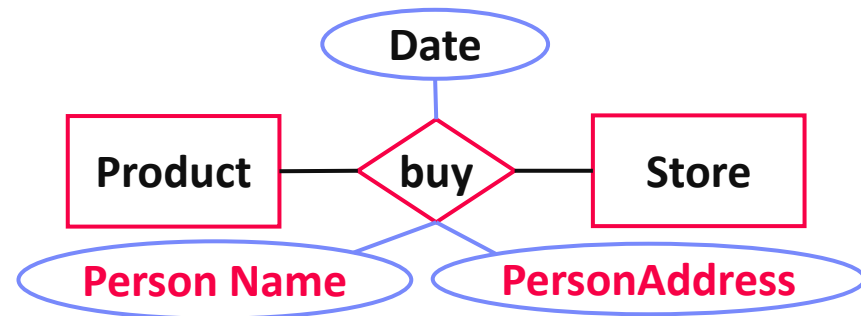
■ #2 Redundant attributes

- **Redundant supplier name** in Part and Supplier



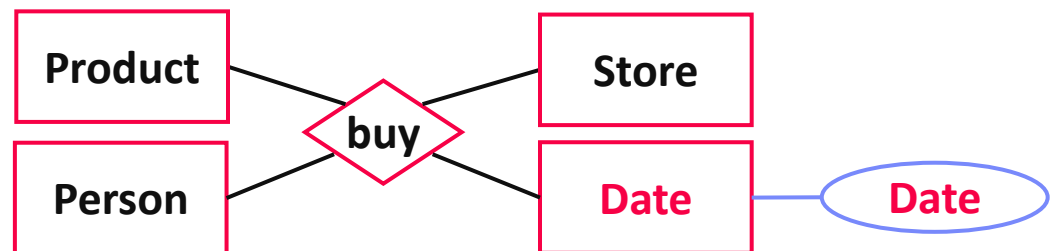
■ #3 Repeated information

- **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

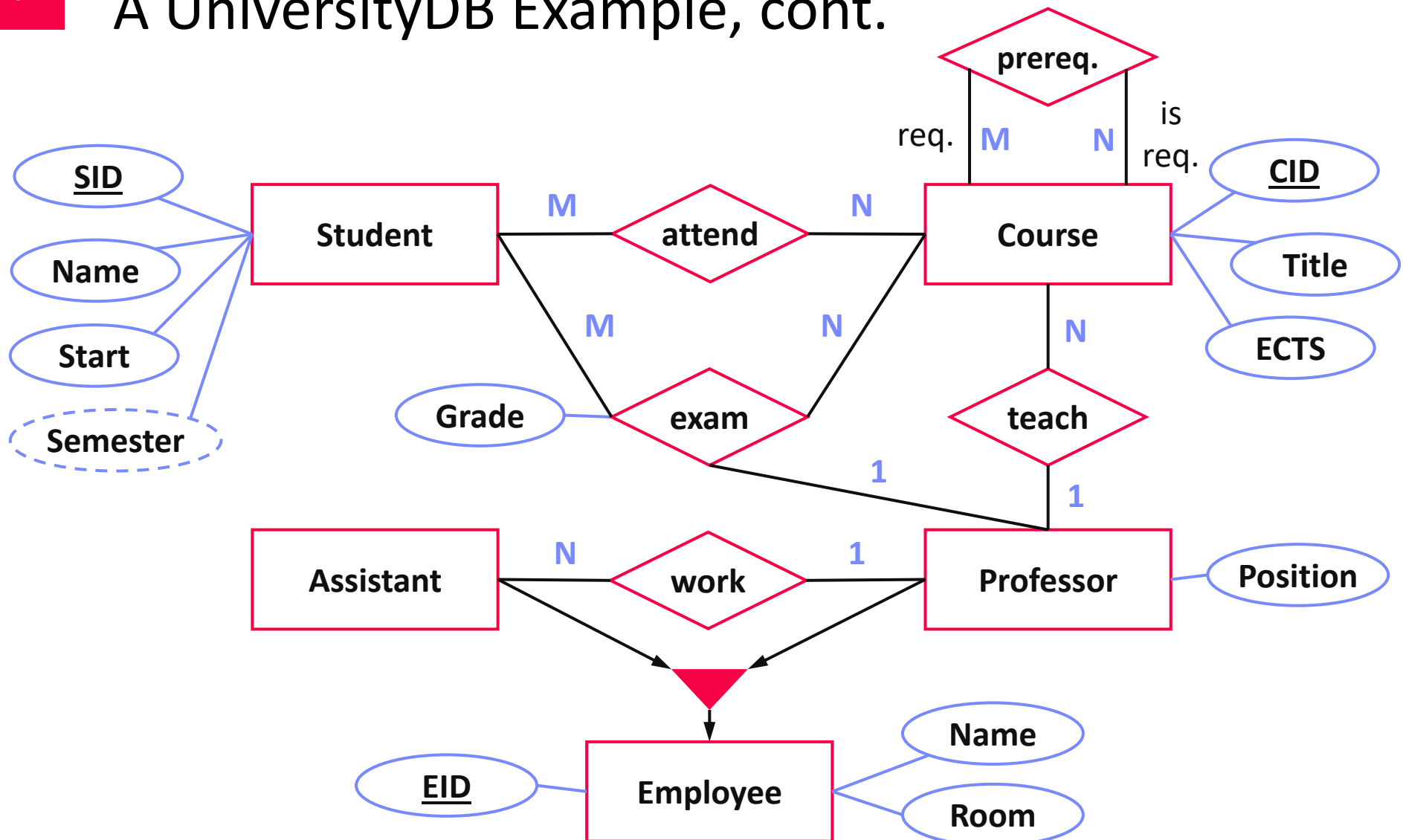
- **Discourse of Real Mini World**

- **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: **Mar 13, 2020**

Deadline: **Mar 31, 2020**

Exercises: DBLP Publications

■ Dataset

- CC0-licensed, derived (extracted, cleaned) from **DBLP** (<https://dblp.org> Feb 1, 2020) for publication year ≥ 2011
- **Note: Still in process of data cleaning**
- Clone or download your copy from <https://github.com/tugraz-isds/datasets.git>

■ 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 data ingestions and query processing)

persons.csv: The persons file contains author information including websites. Its detailed structure and examples look as follows.

```
#PID | name | aliases | affiliation | url
A261789|Matthias Boehm 0001|Matthias Böhm 0001|Graz University of Technology, Austria|http://www.tugraz.at/~tugraz01/staff/boehm
A1537639|Stefanie N. Lindstaedt|Stefanie N. Lindstädt|http://www.tugraz.at/~tugraz01/staff/lindstaedt
A977823|Denis Helic||Graz University of Technology, Austria|http://www.tugraz.at/~tugraz01/staff/helic
```

theses.csv: The theses file contains the information of public PhD theses. Its detailed structure and examples look as follows.

```
#TKey | author | title | year | type | school | pages | isbn
T25621|A261789|Cost-based optimization of integration flow control|2011|Dissertation|Graz University of Technology|150|9783708911111
T30052|A1399369|An Architecture for Fast and General Data Access|2012|Dissertation|Graz University of Technology|150|9783708911111
```

pubs.csv: The pubs file (or better, its individual parts) contains the detailed structure and examples look as follows.

```
#PKey | authors | title | year | type | journal | volume | issue | pages
P519327|A382693:A261789:A261428:A2051042:A69590|MNC: Struempel|2011|Article|ACM Transactions on Database Systems|36(4)|10:1-10:12
P1640801|A261789:A2051042:A2047447:A472485:A261428:A38856|Proceedings of the 2011 ACM SIGMOD Conference on Management of Data|2011|Conference|ACM Press|10:1-10:12
P12485|A1399369:A1703306:A1416241:A557115:A650354:A863102|Proceedings of the 2011 ACM SIGMOD Conference on Management of Data|2011|Conference|ACM Press|10:1-10:12
```

confs.csv: The confs file contains the information on conference proceedings. Its detailed structure and examples look as follows.

```
#CKey | title | editors | year | isbn
C8036|Proceedings of the 2019 International Conference on Database Systems|2019|9781450380361
C76|Proceedings of the 9th USENIX Symposium on Networked Systems|1993|0-913150-00-0
```

Overview Exercise 1 Tasks

- **Task 1.1: ER Modeling (authors, publications)**
 - Create an ER diagram in Modified Chen (MC) notation
 - https://github.com/tugraz-isds/datasets/tree/master/dblp_publications
- **Task 1.2: Mapping ER Diagram into Relational Model**
 - Create a relational schema for the ER diagram from Task 1.1
- **Task 1.3: Relational Normalization**
 - Bring the relational schema from Task 1.2 into third normal form (3NF)
- **Task 1.4: Extra Credit**
- **Expected result (for all three subtasks)**
 - **DBExercise01_<studentID>.pdf**



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

Conclusions 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 Mar 13, 2020; deadline: Mar 31, 2020
- **Recommendation:** start with task 1.1 this weekend; ask questions in upcoming lectures or on news group

■ Next lecture (Mar 16): **03 Data Models and Normalization**