# Systems Design and Security

## Part 7: Database Design

http://staffwww.dcs.shef.ac.uk/people/A.Simons/

Home ⇒ Teaching ⇒ Lectures
⇒ COM2008/COM3008

# Bibliography

- Database Systems
  - T Connolly and C Begg, Database Systems – a Practical Approach to Design, Implementation and Management, 6th ed., Pearson, 2014.
  - C J Date, An Introduction to Database Systems, 8th ed., Pearson, 2003.

- UML Profile for Databases (not official)
  - D Gornik, UML Data Modelling Profile, IBM/Rational Software TP162, May, 2002.
  - S Ambler, A UML Profile for Data Modelling, Agiledata.org, 2009.

# Outline

- Evolution of databases
- Relational data model
- UML profile for databases
- Entity-relationship modelling
- Traditional table normalization
- Normal and pre-normal forms

Reading:  Date chapters 10, 11, 12, 13;
          Connolly and Begg, chapters 11, 12, 13, 14

# Business Data

- Valuable
  - data forms the core of business operations
  - customers, suppliers, purchasing and sales, …
  - often, provides the commercial advantage
- Persistent
  - data survives many executions of the program
  - possibly accessed by many programs, many users
- Protected
  - must be kept safe, even if the system fails
  - must be protected from unauthorised access

# Database Systems

- Early data storage
    - 1880s – punched cards, influenced CODASYL strategy
    - 1950s – magnetic tape, influenced IBM's IMS
- Navigational model
    - 1960s – CODASYL, networks of records, pointer-following
    - IBM IMS, hierarchical records, influences XML today
    - 1990s – XML databases, hierarchical, text-based
- Relational model
    - 1970s – Edgar Codd, relational algebra, platform-independent
    - Relational databases, table-based, most efficient searching
    - Structured Query Language (SQL) – common query language
    - 1980s – object-oriented model, richer datatypes, procedures
    - 1990s – object-relational mapping, OQL, …

# Storage Issues

- Programs in memory
    - objects have rich datatypes, eg: Account, Holder
    - data structure is an arbitrary connected graph
    - structures may be extended dynamically (lists of pointers)
    - navigation is by object reference (memory pointer)
- Database files on disk
    - simply-typed data, eg: Integer, String, Date, Money
    - data structure is a fixed, predefined set of tables
    - data tables have a fixed width, cannot grow/shrink to accommodate varying lists of references
    - navigation is by searching according to key values

# Relational Databases

- **Original Purpose**
  - to eliminate redundancy of stored information
    - don't store same information in several places
    - eliminate blank fields, repeated groups of fields
  - to minimize dependency between data items
    - easy to search for data – optimal links via keys
    - easy to insert, update, delete single items
    - fewer cascading effects (viz. knock-on updates)
- **Structure of Data**
  - logically a set of tables, indexed by row $\times$ column
  - physically a set of files containing many (fixed-length) records

The
University
Of
Sheffield.

# Database Tables

Current Account Table

| number | balance | overDr |
|--------|---------|--------|
| 0214537 | 323.50 | -100 |
| 0773465 | 443.97 | -100 |
| 1334890 | -27.68 | -500 |

- **Table columns**
  - represent attributes
  - have simple types
  - a column has one type
  - a type has a fixed size

- **Table rows**
  - represent objects
  - all attribute values
  - row has mixed types

Holder Table

| custID | forename | surname | addrID |
|--------|----------|---------|--------|
| 235 | Inderpal | Singh | S104DP12 |
| 673 | Sarah | Wilson | S57AA297 |
| 589 | Tariq | Al Harq | S116SQ40 |

# Database Profile

«table»
**Holder**

PK holderID : Integer
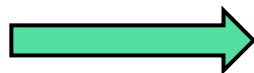　forename : String[30]
　surname : String[30]
FK houseNumber : String[10]
FK postcode : String[10]

primary and foreign keys are attributes that serve a special row-identifying purpose

- UML Profiles
  - a profile is an extension to standard UML notation
  - database profile is not yet universal; but widespread
- Data tables
  - use the «table» label to indicate a data table
  - first box defines columns
  - second box is for triggers, integrity check functions
- Key attributes
  - use PK for primary key
  - use FK for foreign key

# Primary Key

| «table» |
| --- |
| CurrentAccount |
| PK number : Integer<br>  balance : Money<br>  overdraftLimit : Money |
|  |

| «table» |
| --- |
| Address |
| PK houseNumber : String[10]<br>  streetName : String[30]<br>  cityName : String[30]<br>PK postcode : String[10] |
|  |

- Simple key
  - a single column, having a
  - unique value for each row
- Compound key
  - several columns, having a
  - unique value-combination for each row
- Entity integrity
  - PK may not be null – must exist for each row
  - PK must identify a single row – no duplicate keys

The University Of Sheffield.

# Surrogate Key

«table»
**Holder**

PK forename : String[30]
PK surname : String[30]
...

«table»
**Holder**

PK holderID : Integer
  forename : String[30]
  surname : String[30]
...

- Problem
  - forename and surname – can't be a compound PK
  - because of the possibility of duplicate names
- Surrogate key
  - create an artificial attribute
  - only if no other suitable key can be found
  - or if compound key is too large, > 3 rows
  - but don't use for all tables! (wasteful)

The
University
Of
Sheffield.

# Foreign Key

| «table» |
| --- |
| Holder |
| PK holderID : Integer<br>    forename : String[30]<br>    surname : String[30]<br>FK houseID : String[10]<br>FK postcode : String[10] |
|  |

1..*   householder

1   domicile

| «table» |
| --- |
| Address |
| PK houseID : String[10]<br>    streetName : String[30]<br>    cityName : String[30]<br>PK postcode : String[10] |
|  |

- Foreign key
  - a table has one FK for each related table
  - FK from one table refers to the PK in another table
  - value of FK from one table
  - equals the value of the PK in another table
  - may be null, or repeated
- Referential integrity
  - if FK is not null, then
  - row with same PK must exist in the related table

The University Of Sheffield.

# Navigation



```
        «table»
        BookCopy
─────────────────────────
PK copyID : String[30]
FK isbn : Integer
─────────────────────────

```

1..*  copies

```
        «table»
        BookTitle
─────────────────────────
PK isbn Integer
   name : String[50]
   author : String[30]
─────────────────────────
```

directed association

1  title

- Linking tables
  - tables are linked in one direction only
  - link from the "many" to the "one" side – why is this?
- Directed association
  - arrow shows direction of navigation in UML
  - arrow points from the "many" to the "one" side
  - the "many" side has the FK
  - refers to PK of the "one"
  - matched PK, FK will have the same values

# Data Normalization

- Relational Data Analysis
  - Edgar Codd, Ray Boyce – complicated set of rules for transforming arbitrary data tables into Normal Form (NF)
  - based on the detailed analysis of attribute dependency
  - 1NF, 2NF, … 5NF, 6NF (typically 3NF, 4NF required)
- Entity-Relationship Modelling
  - Peter Chen – simple diagram-based technique for normalization
  - based on simplifying object relationships and multiplicities
  - achieves 3NF and sometimes better
- Event-Driven Design
  - Monique Snoeck, Anthony Simons – event table and graph-building technique linking according to existence dependency
  - constructs a data model already in 3NF, sometimes better

# Relationship Types

- One-to-One (1:1)
    - eg: every student has a unique UCard (1:1)
    - eg: each bank branch has a single address (1:1)
- One-to-Many (1:M)
    - eg: an address contains many householders (1:1..*)
    - eg: a person optionally has a driving licence (1:0..1)
    - eg: a woman optionally has many children (1:0..*)
    - eg: a person may optionally also be a student (1:0..1)
- Many-to-Many (M:N)
    - eg: many modules are offered on many degrees (1..*:1..*)
    - eg: books are loaned optionally to a borrower (0..*:0..1)
    - eg: many borrowers optionally reserve many books (0..*:0..*)

# Minimize Data Dependency

- Data design goals
  - remove duplicated data, redundant paths (for space efficiency)
  - optimise table structure for the independent insertion, update and deletion of single data items (reduces cascading effects)
- Minimization technique (ERM)
  - merge tables that are in 1:1 relationship (removes need for join)
  - link tables that are in 1:M relationship (from M→1)
  - introduce new linker tables to encapsulate each M:N relationship, yielding two new 1:M relationships (from M→1)
  - find redundant search paths navigating to the same sets of rows and delete the shorter path

The University Of Sheffield.

# Merge Required

«table»
## Student

PK regNo : Integer
FK degree : String[10]

1 student

«table»
## UCard

PK cardNo : Integer
  startDate : Date
  endDate : Date

1 card

1..* cohort

M:1 relationship

1 enrol

1:1 relationship

in 3NF, a 1:1 relationship
must be eliminated

«table»
## Degree

# Merged Tables

«table»
**Student**

PK regNo : Integer
FK degree : String[10]
cardNo : Integer
startDate : Date
endDate : Date

satisfies
3NF

0..*
cohort

1 enrol

«table»
**Degree**

May need to rename
merged attributes, if
there's a name clash

- **Which table?**
  - keep the stronger concept
  - merge the attributes of the deleted table
  - transfer any associations from the deleted concept
- **Which key?**
  - possibly several candidate keys for the PK
  - preserve one PK only
  - here: cardNo demoted to a dependent attribute

# Linker Required

M:N relationship

possibly has an
association class

In 3NF, a M:N relationship
must be eliminated

| «table»<br>BookCopy | 0..6  borrows ————————— 0..1  borrowedBy | «table»<br>Borrower |

1..*  copies

| Loan |
| --- |
| issueDate : Date<br>dueDate : Date |
| |

1  title

| «table»<br>BookTitle |

1:M relationship

# Linker Added

first 1:M relationship

care with multiplicities

second 1:M relationship

0..1  borrowedBy

| «table» Loan |
|---|
| issueDate : Date<br>dueDate : Date |
|  |

0..6  borrows

1  copy

1  borrower

| «table» BookCopy |
|---|

1..*  copies

| «table» Borrower |
|---|

1  title

| «table» BookTitle |
|---|

the association is promoted
to a separate linker table

the linker relates to exactly
one of each linked object

# Linker Keys

The linker's FKs are the PKs of each related object

The linker's FKs are also its own compound PK

**«table»**
**Loan**

PFK memberID : Integer
PFK copyID : String
issueDate : Date
dueDate : Date

satisfies 3NF

0..1  borrowedBy

0..6  borrows

1 | copy

1 | borrower

**«table»**
**BookCopy**

PK copyID : String
...

A linker is uniquely identified by the two related objects

**«table»**
**Borrower**

PK memberID : Integer
...

The University Of Sheffield.

# Lab 1: Normalize the Library

- Complete the normalization of the tables
  - use the information model developed earlier
  - consider the relationship where borrowers reserve books
  - care with the roles and multiplicities of linkers

Run a Poll

- Identify all navigation paths through the data
  - identify all primary keys (compound? surrogate?)
  - supply all required foreign keys (on the many-side)
  - how does this affect generalisation relationships?

# Generalisation

A subclass always
relates optionally
to its superclass

No candidate key to
identify a person

| Person |
| --- |
| forename : String<br>surname : String<br>gender : Character |
|  |

1            1

Subclasses depend
on the superclass

0..1            0..1

| Student |
| --- |
| regNo : Integer {id}<br>degID : String |
|  |

| Lecturer |
| --- |
| empNo : Integer {id}<br>office : String |
|  |

Some candidates,
but need one key

# Strong/Weak Entities

Standard solution maps all classes to tables in DB (but see later)

Supertype is a strong entity with its own surrogate key

Subtype is a weak entity

Subtype is a weak entity

**«table»**
**Person**

PK personID
   forename : String
   surname : String
   gender : Character

1

1

0..1

0..1

**«table»**
**Student**

PFK personID : Integer
   regNo : Integer
FK  degID : String

**«table»**
**Lecturer**

PFK personID : Integer
   empNo : Integer
   office : String

The FK is also a sufficient PK

Other keys are demoted to dependents

# Composition

The whole is **materially indivisible** from its parts

The parts depend on the whole – **cannot exist** independently

**Book**

title : String
author : String
isbn : String {id}
publisher : String

**Chapter**

number : Integer {id}
title : String

1..* chapters

Cascading deletion: deleting the whole also deletes the parts

# Strong/Weak Entities

Composite is a **strong** entity with a natural key

```
          «table»
           Book

   title : String
   author : String
PK isbn : String
   publisher : String
```
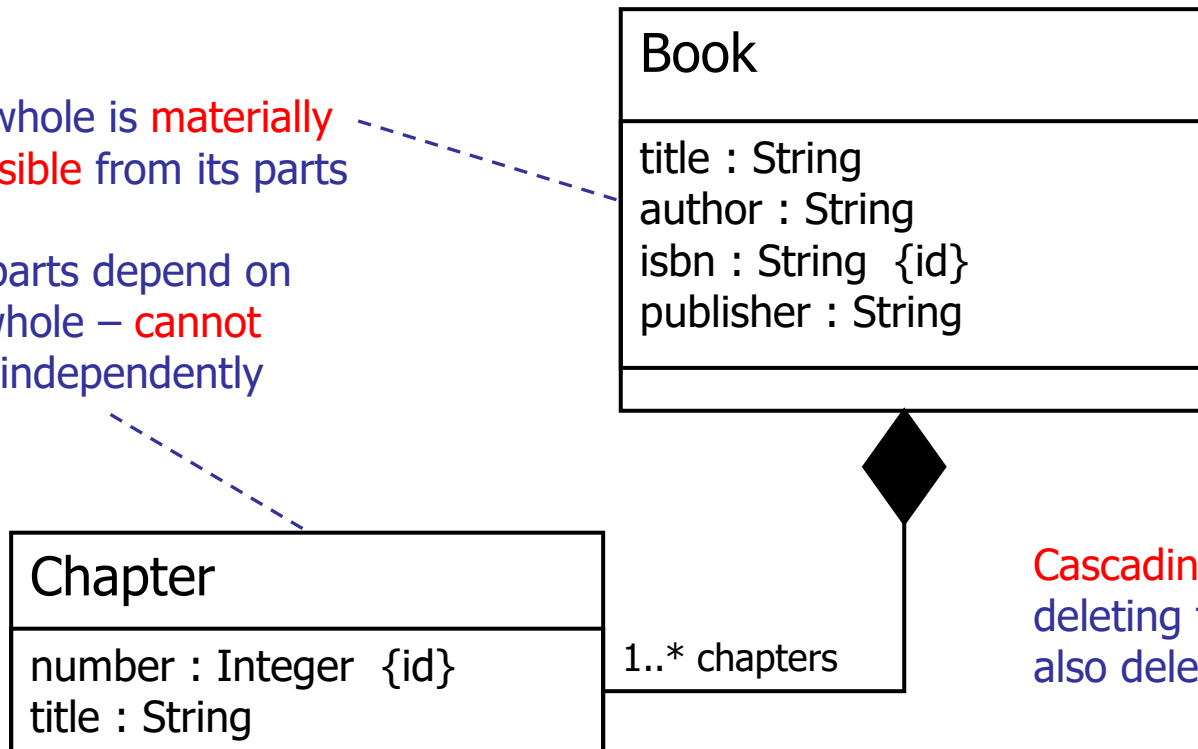
Component is a **weak** entity, cannot be identified just by its **weak key** number

```
          «table»
          Chapter

PFK isbn : Integer
PK  number : Integer
     title : String[30]
```

1..*

chapters

1

Otherwse, standard M:1 linkage

Requires a compound key also including FK

# Aggregation

Aggregation is always an optional assembly depending on its parts

Can disassemble and reassemble …

optional-to-many needs a linker

**Vehicle**

serialNumber : Integer {id}
brand : String

0..1

0..1

optional-to-one

4 wheels

**Wheel**

diameter : Integer
tyrePressure : Real

1 engine

1 chassis

**Engine**

engineNo : Integer {id}
cubicCapacity : Integer

**Chassis**

chassisNo : Integer {id}

…and the parts may exist independently

# Table Design

...except where a linker is required

**«table» Vehicle**

PK serialNumber : Integer
   brand : String
FK engineNo : Integer
FK chassisNo : Integer

**«table» VehicleHasWheels**

PFK serialNumber : Integer
PFK wheelID : Integer

1

0..4

Whole has many FKs referring to parts ...

0..1   0..1

0..1

1

1

1

**«table» Engine**

PK engineNo: Integer
   cubicCapacity: Integer

**«table» Chassis**

PK chassisNo: Integer

**«table» Wheel**

PK wheelID: Integer
   diameter : Real
   tyrePressure : Real

The University Of Sheffield.

# Traditional Normalization

- Start with arbitrary table, subdivide
    - start: data in one table, maybe with list-valued columns
    - next: break up lists by replicating atomic data in many rows
    - next: split up into many tables to remove replicated data
    - next: make sure nonkey attributes depend fully on the key
- Driven by functional dependency
    - eg: customer's name, age depend on which customer
    - viz:  custID $\rightarrow$ name, custID $\rightarrow$ age, but not age $\rightarrow$ name
    - transitive, eg deptID in:  regNo $\rightarrow$ degID $\rightarrow$ deptID
    - partial, eg modID $\rightarrow$ lectID in:  R(lectID, modID, regNo)
    - multivalued, eg product set in:  company $\rightarrow\rightarrow$ product

# Normal Forms

- 1NF (Codd) – no lists or repeating data groups (atomicity rule) and find a primary key for each row

- 2NF (Codd) – and no nonkey attributes depend on part of a compound primary key (full functional dependence on the key)

- 3NF (Codd) – and no nonkey attributes depend on other nonkey attributes (don't depend transitively on the primary key)

  - BCNF (Boyce/Codd; Heath) – and every attribute (or compound) on which other attributes fully depend is also a candidate key

- 4NF (Fagin) – and no multivalued dependency anomalies (no dependent set of values depends on part of a compound key)

- 5NF (Fagin) – and every join dependency is implied by the candidate keys (project/join NF, rare case, hard to explain)

# First Normal Form

Enrolment (non-normal)

non-atomic!
list-valued data!

| regNo | degID | modIDs |
|-------|-------|--------|
| 0214537 | G402 | COM101, COM102, ACS101... |
| 0773465 | G650 | COM101, COM103, EEE105, ... |
| 1334890 | G600 | |

...at the cost of replicating each regNo and degID

now atomic, lists split up...

PK = regNo + degID + modID

Enrolment (1NF)

| regNo | degID | modID |
|-------|-------|-------|
| 0214537 | G402 | COM101 |
| 0214537 | G402 | COM102 |
| 0214537 | G402 | ACS101 |
| 0773465 | G650 | COM101 |
| 0773465 | G650 | COM103 |

# Second Normal Form

Examinations (1NF)

...so, split into two tables

| regNo | modID | lectID | marks |
|-------|-------|--------|-------|
| 0214537 | COM101 | GJB | 72 |
| 0773465 | COM101 | GJB | 65 |
| 1334890 | COM103 | SDN | 67 |

Teaches (2NF)

| modID | lectID |
|-------|--------|
| COM101 | GJB |
| COM103 | SDN |

PK = regNo, modID

while each mark depends
on both regNo, modID,
the lecturer depends only
on modID! ...

| regNo | modID | marks |
|-------|-------|-------|
| 0214537 | COM101 | 72 |
| 0773465 | COM101 | 65 |
| 1334890 | COM103 | 67 |

Examinations (2NF)

The University Of Sheffield.

# Third Normal Form

Registration (2NF)

| regNo | degID | deptID |
|-------|-------|--------|
| 0214537 | G404 | COM |
| 0343662 | G404 | COM |
| 1754532 | H400 | MEC |

PK = regNo

while each degree taken depends directly on the student regNo, the home deptID depends on the degID, only indirectly on the PK!

Owner (3NF)

| degID | deptID |
|-------|--------|
| G404 | COM |
| H400 | MEC |

| regNo | degID |
|-------|-------|
| 0214537 | G404 |
| 0343662 | G404 |
| 1754532 | H400 |

Registration (3NF)

...eliminate transitive dependencies in 3NF

The University Of Sheffield.

# Fourth Normal Form

Sells (3NF)

| company | product | country |
|---------|---------|---------|
| IBM | PC | France |
| IBM | Server | Italy |
| DEC | PC | France |

PK = company, product, country

multivalued dependency: set of products and set of countries depend on company, but are independent of each other

...so decompose into two tables

Exports (4NF)

| company | country |
|---------|---------|
| IBM | France |
| IBM | Italy |
| DEC | France |

| company | product |
|---------|---------|
| IBM | PC |
| IBM | Server |
| DEC | PC |

Makes (4NF)

# How Normal?

- Typical levels
  - aim for at least 3NF – good separation of data
  - 4NF needed where multivalued dependency exists
- What to check
  - insertion anomaly – can't insert X until Y is also supplied
  - update anomaly – affects some, not all rows (duplication)
  - deletion anomaly – delete X causes unintended loss of Y
- Denormalize?
  - to increase efficiency of searches (fewer joins)
  - to decrease fragmentation of data (fewer tables)

# Lab 2: Normalize Orders

| orderID | city | productID | quantity |
|---------|------|-----------|----------|
| S1 | London | P1 | 100 |
| S1 | London | P2 | 200 |
| S2 | Paris | P1 | 200 |
| S2 | Paris | P2 | 300 |
| S3 | Paris | P2 | 300 |
| S4 | London | P2 | 400 |
| S4 | London | P4 | 500 |
| S4 | London | P5 | 700 |

- First: work out the functional dependencies
- Second: split the tables, aiming for at least 3NF

Run a Poll

# Ternary Association

«table»
**Student**

PK regNo : Integer
   cardNo : Integer
   startDate : Date
   endDate : Date

1..*
cohort

higher-order association

«table»
**Degree**

PK courseID : String
   name : String

1
enrol

1 | home

«table»
**Department**

PK deptID: String
   name : String

registration is ternary: relates student, degree, department

multiplicities are assigned to each end when instances are fixed at all other ends

# Ternary Linker



promote to a linker table

PKs are tricky – only need regNo, courseID to uniquely identify a registration

linker relates exactly one of each object

**«table»**
**Registration**

PFK regNo : Integer
PFK courseID: String
FK deptID: String

1

enrol

1..*

cohort

dependency problems?

1 student

1..* register

1 department

1 degree

**«table»**
**Student**

regNo : Integer «PK»
cardNo : Integer
startDate : Date
endDate : Date

**«table»**
**Department**

deptID: String «PK»
name : String

**«table»**
**Degree**

courseID : String «PK»
name : String

# Make Binary

«table»
**Student**

PK regNo : Integer
   cardNo : Integer
   startDate : Date
   endDate : Date
FK courseID : String
FK deptID: String

safer to construct only binary associations – makes dependencies easier to control

«table»
**Degree**

PK courseID : String
   name : String
FK deptID: String

1..*  cohort     enrol  1

1..*  course

«table»
**Department**

PK deptID: String
   name : String

1..*  students    1  department

1  home

# Redundant Paths

«table»
**Student**

PK regNo : Integer
    cardNo : Integer
    startDate : Date
    endDate : Date
FK courseID : String
FK deptID: String

two paths exist from Student to Department - do they lead to the same set of instances?

«table»
**Degree**

PK courseID : String
    name : String
FK deptID: String

1..*    cohort                    enrol    1

1..*  course

1..*  students

«table»
**Department**

PK deptID: String
    name : String

if so, must delete one path...

1  department                    home  1

...but which path is redundant?

# Redundancy Eliminated

«table»
## Student

PK regNo : Integer
cardNo : Integer
startDate : Date
endDate : Date
FK courseID : String

always preserve the
longer path and delete
the shorter path

1..*                                    1
cohort                            enrol

«table»
## Degree

PK courseID : String
name : String
FK deptID: String

saves space on FKs

can still navigate to the
related Department

«table»
## Department

PK deptID: String
name : String

1..*  course

1
home

satisfies
4NF

# Involuted Association

«table»
**Employee**

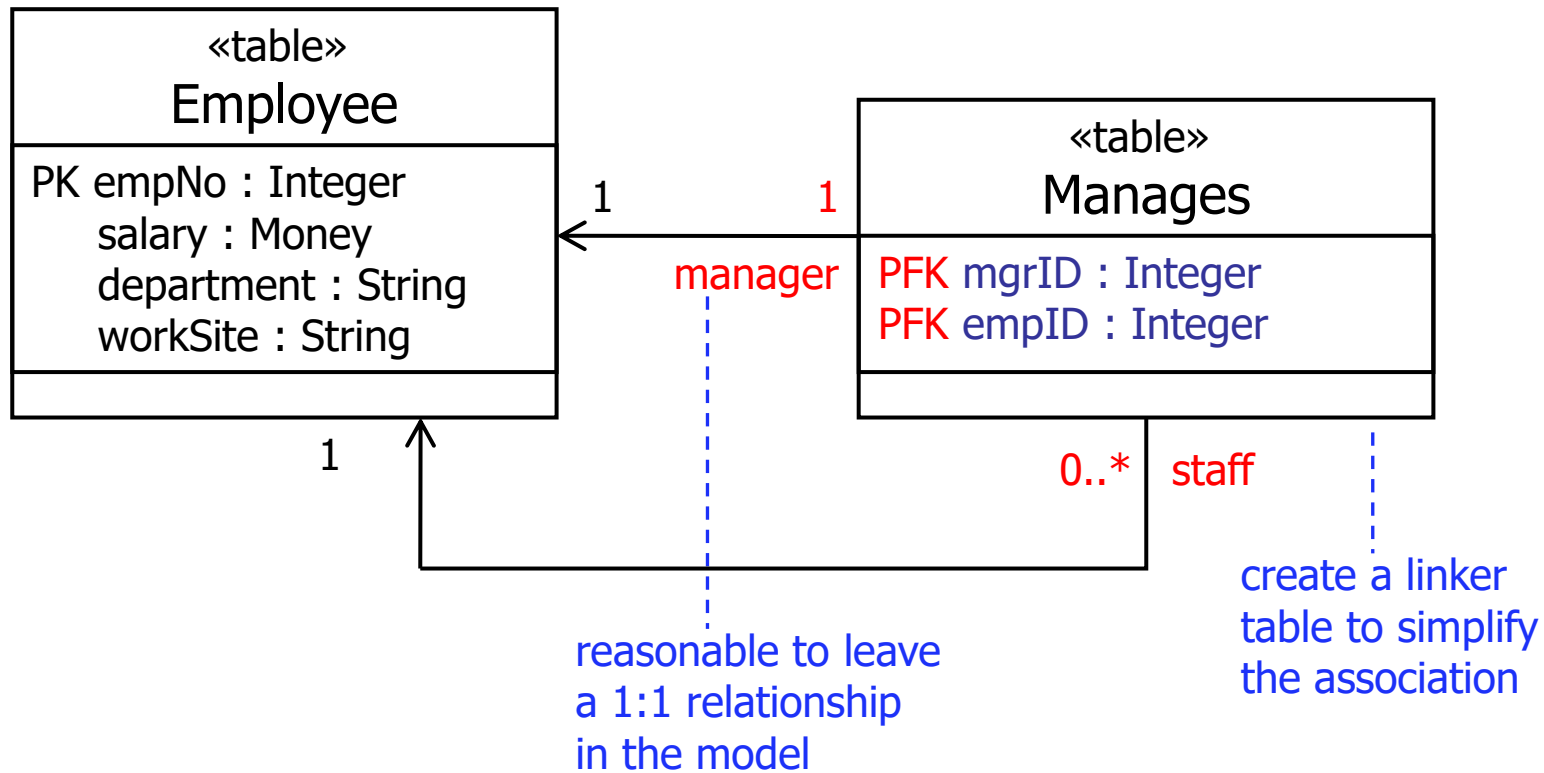| |
|---|
| PK empNo : Integer |
| salary : Money |
| department : String |
| workSite : String |

0..*

staff

1 ↑ manager

involuted association

an involuted association relates a table back to itself

typical in human relationships, eg: marriage, employment

cannot easily link through FKs, since not all Employees are managed (null values)
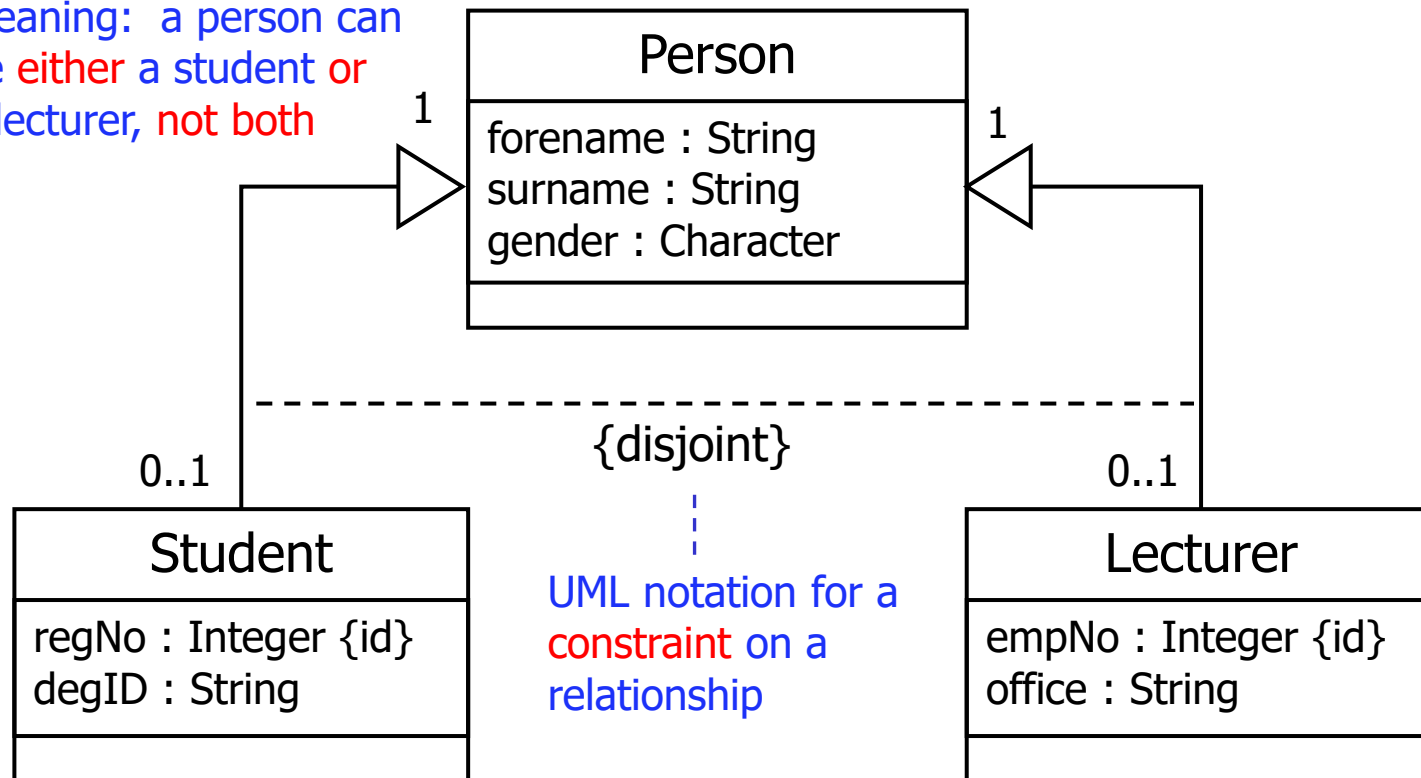
impossible for M:N anyway

# Involuted Linker

«table»
**Employee**

PK empNo : Integer
   salary : Money
   department : String
   workSite : String

«table»
**Manages**

PFK mgrID : Integer
PFK empID : Integer

1     1

manager

1

0..* staff

reasonable to leave
a 1:1 relationship
in the model

create a linker
table to simplify
the association

# Disjoint Semantics

meaning: a person can be either a student or a lecturer, not both

**Person**

forename : String
surname : String
gender : Character

1           1

0..1        {disjoint}        0..1

**Student**

regNo : Integer {id}
degID : String

UML notation for a constraint on a relationship
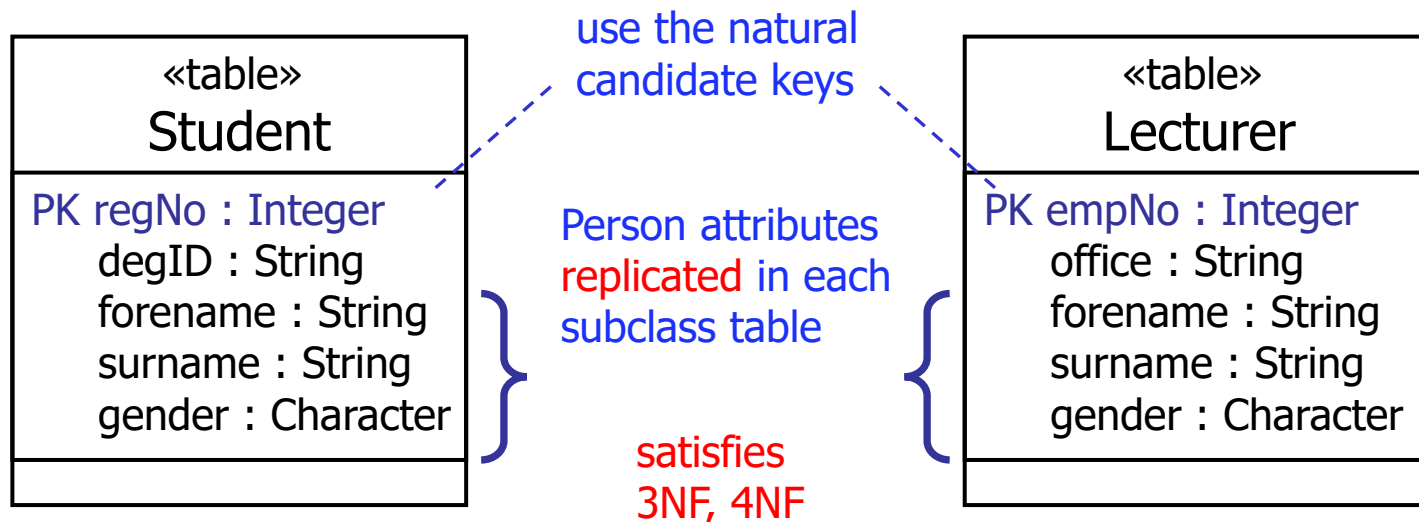
**Lecturer**

empNo : Integer {id}
office : String

# Partitioned Tables

if a person can be
either a student or
a lecturer, not both...
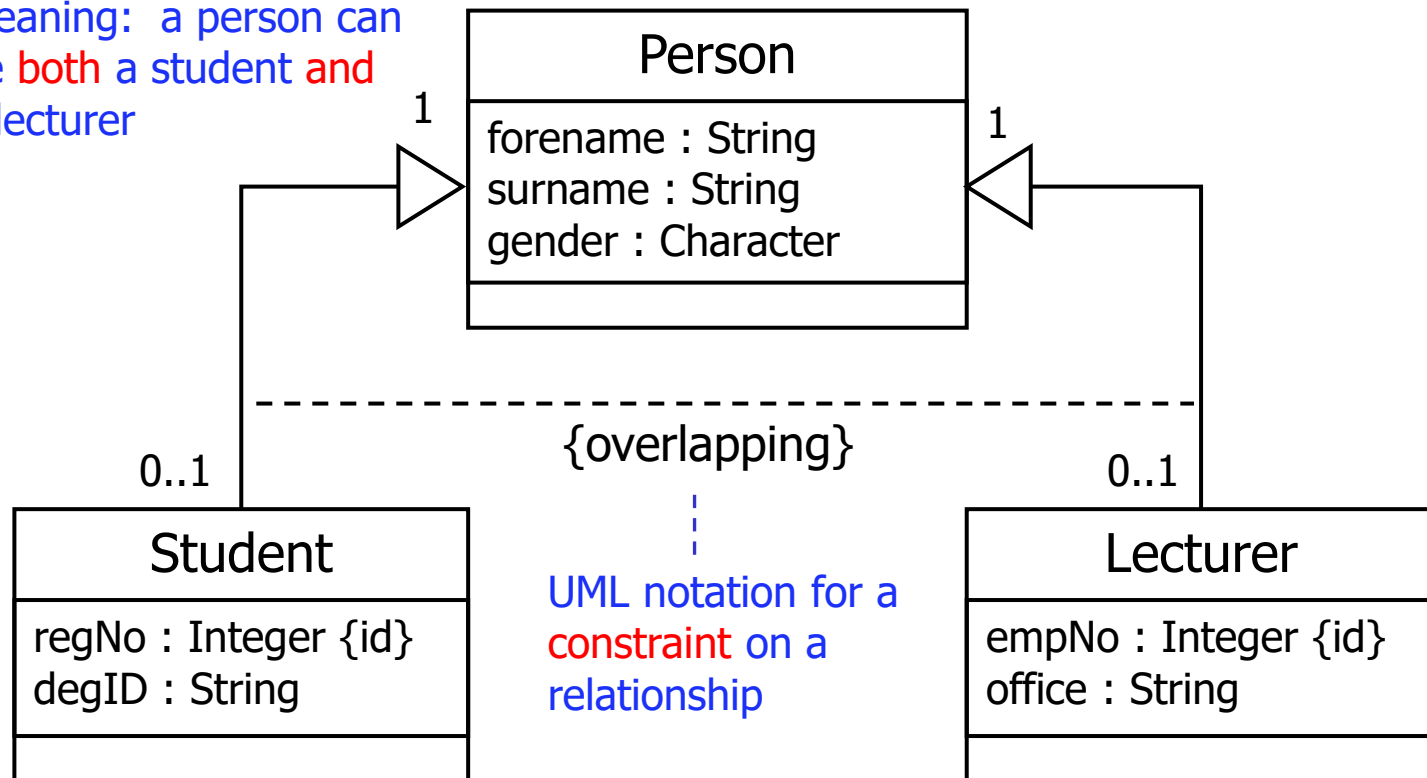
...requires only 2
tables in the DB
(reduces joins)...

...but cannot represent a
person who is both a
student and a lecturer

use the natural
candidate keys

| «table» Student |
| --- |
| PK regNo : Integer |
|   degID : String |
|   forename : String |
|   surname : String |
|   gender : Character |
| |

Person attributes
replicated in each
subclass table

satisfies
3NF, 4NF

| «table» Lecturer |
| --- |
| PK empNo : Integer |
|   office : String |
|   forename : String |
|   surname : String |
|   gender : Character |
| |

# Overlapping Semantics

meaning: a person can be **both a student** and a lecturer

## Person

forename : String
surname : String
gender : Character

1

1

0..1

## Student

regNo : Integer {id}
degID : String

0..1

## Lecturer

empNo : Integer {id}
office : String

{overlapping}

UML notation for a constraint on a relationship

# Denormalized Table

surrogate key for
Person instances
(generated key)

**«table»**
**Person**

| |
|---|
| PK personID |
|     forename : String |
|     surname : String |
|     gender : Character |
|     regNo : Integer |
| FK degID : String |
|     empNo : Integer |
|     office : String |

Student attributes
may be null!

Lecturer attributes
may be null!

folds all subclass attributes
into the superclass

reduces joins between
tables to reconstruct whole
persons...

...but wasteful in space,
because of many null fields

satisfies only 2NF
due to transitive
dependency on
regNo, empNo

The
University
Of
Sheffield.

# Denormalization?

- Reduces number of tables
    - eg: flatten generalisation structures into one table – but at the cost of wasted space (null values)
    - eg: flatten composition structures into one table – at the cost of replicating the part-data many times in the whole
    - saves reconstructing whole objects in memory from fragmented table data – fewer joins
- Optimizes access paths
    - eg: retain redundant access paths through the data – but at the cost of additional FKs, risk of inconsistency
    - saves time when searching – again, fewer joins

The University Of Sheffield.

# Summary

- Database systems have evolved from the network model, via the hierarchical model, to the relational model
- The relational model stores data in tables, whose rows denote objects and whose columns denote their simple attributes.
- Each row must be uniquely identifiable by a primary key
- Navigation is unidirectional, from the many to the one, using a foreign key to identify the related row in another table
- Data must be normalised to 3NF, 4NF to support insert, update and delete without duplication, or cascading effects
- Entity-relationship modelling is one approach to normalising data, based on optimising multiplicities
- Traditional table normalisation is another approach based on analysing functional dependency