# Entrepôts de Données et Big-Data

# Équipe Pédagogique

- Anne-Muriel Arigon - UM, LIRMM - <u>Responsable UE</u>
  (<u>chifolleau@lirmm.fr</u>)

- Federico Ulliana - INRIA, UM
  (<u>federico.ulliana@inria.fr</u>)

  `Special Guest :`

- Christophe Menichetti - HPE (Hewlett Packard), IBM

# Objectif du Module

*vers le*

## BIG-DATA ANALYTICS

*Plus concrètement, présenter*

     *1. les techniques de modélisation*

     *2. les systèmes et technologies*

*conçues pour l'analyse des données massives*

# Programme

1. Rappels

3. Préliminaires : stockage et optimisation

2. Entrepôts de Données

3. Hadoop/Map-Reduce

4. **Séminaires** : Défis, Cas d'Utilisation, Solutions Big-Data Analytics & IA/ML

# Database Evolution History

**1960s**
First Computerized Database Models

Hierarchical Model (IMS)

Network Model (CODASYL)

**1970s**
The Dawn of the Database
- The relational model and its language SQL emerge
- The disruptive model causes the demise of other models

**1970** E.F CODD Writes a Paper on the Relational Database Model

System R by IBM (SQL)
Ingres (QUEL)

Single Instance Relational Database

Entity Relational Database

ORACLE
1st commercially available RDBMS

IBM DB2

**1980s**
An Industry Develops
- SQL becomes the de-facto standard
- Commercial offerings from IBM, Oracle grow market
- Other data models enter the scene, without much traction

SAP Sybase
Informix

Object Oriented Database

Distributed SQL Data Warehouse
Teradata

Dawn of the Internet

**1990s**
Technology Shifts
- Data explodes with the Internet age
- Single server SQL databases run into resource problems
- Business Intelligence and Analytics move out of transactional database

BUSINESS INTELLIGENCE

ANALYTICS

Application Layer

New Distributed SQL Data Warehouse

**2000s**
New Players Emerge
- Data variety, velocity and volume increase
- New analytics SQL databases are introduced
- NoSQL databases fill the gap for processing unstructured data
- Hadoop gains traction for analyzing petabytes of data

NOSQL

Distributed SQL

Netezza
Microsoft
Aster
Oracle
SAP
IBM
Vertica

MPP

HADOOP

**Today**
Databases Adapt and Evolve
- Businesses require real-time analytics on operational data
- Scale-up SQL proves too costly, but scale-out removes resource constraint
- Scale-out provides real time analytics with high volume transactions
- Google and Clustrix are pioneers in this space

Google BigTable
CouchDB
MongoDB
Cassandra
Redis

Clustrix
NuoDB
VoltDB

Google Spanner

Oracle In-memory option

MASSIVELY PARALLEL PROCESSING
Clustrix adds MPP
SAP HANA In-memory

MPP

REAL-TIME ANALYTICS

SQL HIVE

**The Future**
Businesses Advance with Database Innovations
- Single node SQL gets replaced by scale-out SQL
- Data warehouse type analytics will become available in real-time database
- Businesses gain a significant edge and increased agility

COLUMNAR

REPLACED?

ALL ANALYTICS?

**Winning Database Platforms**

NOSQL DATABASE

DISTRIBUTED SQL

HADOOP

Source: Robin Purohit

# MCC 40%TD + 60%Examen

- TD obligatoires          TD facultatifs
    - Rappels                     - Optimisation
    - Entrepôts D. 1          - Entrepôts D. 2
                                        - Map/Reduce

- Mini-projet
    – Rendu intermédiaire          (document)
    – Rendu final       (document, code, transparents)

- **Présence obligatoire** séminaires C. Menichetti

# Divers

- TD
    - salles 36.207 (F. Ulliana) + 5.130 (Anne-Muriel Arigon/Julie Cailler)
    - éléments de correction dispensés en cours
    - binômes ou monômes (<u>sans exceptions</u>)

- Consultez la page Moodle dédiée à l'UE

- Access comptes Oracle *possible à distance* (SSH)
    - instructions sur Moodle

- Signaler les dysfonctionnements des serveurs Oracle
    - ENT --> Assistance --> Centre de Services --> Déclarer un Incident

- Nous écrire pour tout type de problème ou question
    - <u>chifolleau@lirmm.fr</u>,<u>federico.ulliana@inria.fr</u>

# Objectifs Pour Aujourd'hui

- Rappels
  - UML
  - Modèle Relationnel
  - SQL

- Stockage
  - "the database game" (clickbait)

- TD
  - Modélisation UML, SQL (facile)
  - Rendu jeudi 21/09 avant 23h59

# *Rappels*

# Readings

*These slides should be considered simply as "pointers" to the references below.*

**[BD-G]** Bases de données,
Georges Gardarin, 5ème edition 2005
http://georges.gardarin.free.fr/Livre_BD_Contenu/XX-TotalBD.pdf

**[ORA]** Oracle® Database SQL Language
Reference 11g Release 1 (11.1) - 2013
http://docs.oracle.com/cd/B28359_01/server.111/b28286.pdf

**[UML]** Prolegomenes_uml.pdf

**[UML2]** UML2 : de l'apprentissage à la pratique

# Relational Databases & UML

**NB :** *Assumed to be well known from L2/L3, we just recall basic topics.*

1. UML

2. Relational Model

3. SQL

# Levels of Modelling

- Conceptual Model                                    (UML, EA, Merise)
    - defines what the system contains

- Logical Model      (Relational Model, Object Model, Graph)
    - Defines data structures and rules of the system

- Physical Model                                    (SQL, OQL, XML)
    - defines how the system has to be implemented

Peter Chang            Ted Codd            Don Chamberlin

EA – 1976              RM – 1970           SQL – 1974

# BASIC RELATIONAL MODEL THEORY

# The Relational Model

*Everything is a relation*

- **Person**(Bob,42,Paris)

  (can model entities)

- **LiveTogether**(Alice,Bob,Lyon,2010)

  (can model associations)

# Relational Schema

- A set of relations built on a set of attributes, with well defined domains.

**Person**(Name,Age,City)

Name : String      Age: Integer      City : {Lyon,Paris}

# The Model  VS.  The Content

- The idea of representing data using relations is abstract, universal.

- But, as this data is originated from real world interactions (eg., trading, social), all forms of weak and strong correlations are found in it.

# Functional Dependency

*A set of attributes **A** determining a set of attributes **B***

```
(name, surname)    ------------>    birthday
                    determine

        city       ------------>    (population, state)
                    determine
```

# Key(s)

*minimal set of attributes determining a whole tuple*

**LiveTogether**(<u>Person1</u>, <u>Person2</u>, City, <u>Date</u>)

<u>Person1</u> , <u>Person2</u> , <u>Date</u>       --------> City

key                    determines

(ex)    **LiveTogether**(Alice,Bob,Lyon,2010)

Strongly recommended in systems (efficiency, integrity)

# Data dependencies were undesirable

- Except for **keys** and referential **integrity constraints**
  - beside these cases, they just bring redundancy

- Database **normalization** eliminated dependencies

# Normal Forms

Normal-Forms are <u>guidelines</u> for modeling.

Normal forms came <u>after</u> the relational model.

Their definition is motivated by design <u>mistakes</u>.

*So, let's find the right place for the attributes !*

# Normal Forms : 2NF

- **2NF** : non-key attributes fully-dependent from the key

**FournisseurPiece**(<u>Name</u>, <u>Article</u>, Address, Price)        **NO**

(Article --/--> Address)

⬇        (decomposition)

**FournisseurPiece**(<u>Name</u>, <u>Article</u>, Price)

**Fournisseur**(<u>Name</u>, Address)

# Normal Forms : 3NF

[BD-G] chapter VI section 6.3

- 3NF : no dependencies between non-key attributes

**Person**( <u>ID</u> , Name, City, *CityPopulation* )                **NO**

ID -----> City -----> CityPopulation

⬇                        (decomposition)

**Person**( <u>ID</u> , Name, City)        **Place**(City, *CityPopulation)*

# Normal Forms

[BD-G] chapter VI section 6.3

- 3NF respected by **most** *"transactional-database"* you will find in **any** real world company
  - it allows to fix common data redundancy problems
  - also, every schema can be normalized in 3NF

- but sometimes only 2NF …
  - not (always) by mistake of database administrator
  - conscious choice : pay the price of redundancy (and update) to gain in performances (more later)

- 3NF = 2NF + no indirect dependencies          (simplification)
- 2NF = 1NF + no partially dependent attributes     (simplification)
- 1NF =  ?

- Are there transactional databases which do not respect the 1NF ?          (Non-first-normal-form N1NF)

# Normal Forms

[BD-G] chapter VI section 6.3

- 3NF respected by **most** *"transactional-database"* you will find in **any** real world company
  - it allows to fix common data redundancy problems
  - also, every schema can be normalized in 3NF

- but sometimes only 2NF …
  - not (always) by mistake of database administrator
  - conscious choice : pay the price of redundancy (and update) to gain in performances (more later)

- 3NF = 2NF + no indirect dependencies          (simplification)
- 2NF = 1NF + no partially dependent attributes     (simplification)
- 1NF =    only atomic values in a tuple

- Are there transactional databases which do not respect the 1NF ?          (Non-first-normal-form N1NF)
  - Yes ! Object-oriented, XML, many in the NoSQL family (using JSON as data model)

# Normal Forms : Remarks

- 2NF & 3NF respected in practically any information system using a relational database

- Stronger normal forms (BCNF, 4NF, 5NF) are less employed (avoid rarer mistakes; not always achievable)

- Exceptions to normalizations are tolerated to save joins (at the price of redundancy)
  - **we will see this for datawarehouses**

# SQL : SURVIVAL KIT

# SQL

- Structured Query Language

- Declarative (logical) Language : tell what you want from relations, not what to do with them.
  - This is the main difference with C and Java, *not only* the fact that we deal with data.

- In SQL terminology, a relation is called "table".

# SQL

- DDL (Data Definition Language)
  - CREATE/ALTER  structures (table, view, index)


- DML (Data Manipulation Language)
  - UPDATE/INSERT/DELETE content


- DQL (Data Query Language)
  - SELECT data

# SQL

- The purpose of the following slides is just to recall you the basic syntax.

- Check the documentation for more details or advanced features !

# Create Table

```
CREATE TABLE Employee (

 id    NUMBER,
 name VARCHAR2(50),
 birthday DATE

)
```

# Datatypes

Why do we need datatypes ?

- To associate fixed set of properties to attributes.

- This improves the database:
  - coherence : type-checking operations
    - cannot sum two strings

  - efficiency : a datatype has its own best storage
    - BLOB vs integers

# ORACLE Built-in Datatypes

- Character
- Numeric
- Date/Time
- Large Object

Complete list : see [ORA] table 2-1

# Value Constraints

- Why do we need constraints ?

- To restrict the values in a database and ensure the data integrity
  - ex : No employee without an ID

# Not Null

- Prohibits a database value from being null

```
CREATE TABLE Employee (

id         NUMBER NOT NULL,
name       VARCHAR2(50),
birthday   DATE

)
```

# Unique

- Prohibits multiple rows from having the same value (but allows them to be null)

```
CREATE TABLE Employee (

id          NUMBER UNIQUE,
name        VARCHAR2(50),
birthday    DATE

)
```

# Primary Key

- Combines a NOT NULL constraint and a UNIQUE constraint in a single declaration

```
CREATE TABLE Employee (

id          NUMBER PRIMARY KEY,
name        VARCHAR2(50),
birthday    DATE

)
```

# Primary Key : Multiple Attributes

[ORA] section 8-20

- Combines a NOT NULL constraint and a unique constraint in a single declaration

```
CREATE TABLE Employee (
id          NUMBER,
name        VARCHAR2(50),
birthday    DATE,

PRIMARY KEY (name,birthday)
)
```

# Foreign key

*one (or more)* **attributes** *which correspond to the* **key** *of another relation*

**Employee**( <u>ID</u>, Name, **Department_id**)

**Departement**( <u>Dept_ID</u>, Name )

# Foreign Key

Requires values in one table to match values in another table.

```
CREATE TABLE Employee (

    id          NUMBER,
    department  NUMBER

    FOREIGN KEY department
                 REFERENCES Dept(dept_id)
    )
```

# Check

Requires a value to satisfy with a specified condition

```
CREATE TABLE Employee (
        id              NUMBER,
        department      NUMBER,
        office          VARCHAR2(10)

        CHECK (
           office IN
               ('DALLAS','BOSTON', 'PARIS','TOKYO')
        )
)
```

Oracle does not verify mutually exclusive conditions (eg. AGE>1 AND AGE<0)

# ALTER TABLE

- Add a new column
  - `ALTER TABLE Employee ADD (office VARCHAR2(20));`

- Modify an existing column
  - `ALTER TABLE Employee MODIFY (office NUMBER);`

- Define a default value for the new column
  - `ALTER TABLE Employee MODIFY office DEFAULT 'Corridor';`

- Drop a column
  - `ALTER TABLE Employee DROP (office);`

# DELETE    TRUNCATE        DROP

| | removes | | |
|---|---|---|---|
| | rows | table | rollback |
| DELETE | ✓ | ✗ | ✓ |
| TRUNCATE | ✓ | ✗ | ✗ |
| DROP | ✓ | ✓ | ✗ |

- TRUNCATE = DROP + CREATE TABLE

# INSERT

INSERT

INTO Employee

VALUES

     ( 'Bob',

     **TO_DATE(**

     '03-OCT-1972','DD-MON-YYYY'**)**

     )

*TO_DATE* converts a character/numeric to a date

[ORA] 2-49/50

# SELECT     FROM

SELECT
    *TO_CHAR*(birthday,'MM-DD-YYYY')

FROM
    Employee

*TO_CHAR*  converts a numeric to a character

# SELECT FROM WHERE

[BD-G] chapter VII section 3.2 [ORA] section 2-50

```
SELECT
   name

FROM
   Employee

WHERE
   birthday >
      TO_DATE('01-10-1970','DD-MM-YYYY')
```

# JOINS

```
SELECT
  Employee.name, Department.name
FROM
  Employee, Department
WHERE
  Employee.dept = Department.id
```

# JOINS

*Employee*

| name | dept |
|------|------|
| Alice | dep1 |
| Bob | dep2 |
| Eddy | dep1 |

*Department*

| id | name |
|------|------|
| dep1 | Sales |
| dep2 | Production |

*Emp_Join_Dep*

| Employee.name | Department.name |
|---------------|-----------------|
| Alice | Sales |
| Bob | Production |
| Eddy | Sales |

# Group By

```
SELECT
   dept, count(*) as N
FROM
   Employee
GROUP By
   dept
```

*Employee*

| name | dept |
|------|------|
| Alice | dep1 |
| Bob | dep2 |
| Eddy | dep1 |

*Agg_Emp*

| dept | N |
|------|---|
| dep1 | 2 |
| dep2 | 1 |

*Group by queries are at the essence of ANALYTICS.*

# Summing Up

- Relations
  - Functional dependencies, Normal-forms

- SQL
  - CREATE, INSERT, DELETE, SELECT, GROUP-BY

# CONCEPTUAL MODELING WITH ENTITY-RELATIONSHIP DIAGRAMS

*To minimize the number of graphical notations introduced across the various programs of this training, we take an alternative model.*

# CONCEPTUAL MODELING WITH ~~ENTITY-RELATIONSHIP DIAGRAMS~~

# CONCEPTUAL MODELING WITH UML

# UML (Unified Modeling Language)

- Universal **graphical** modeling language designed to model objects, associations, time events, system states

  - Main goal is to ease prototyping

- Good news : UML is a rich model and we can also use it to model data !

# UML (Unified Modeling Language)

**graphical**

means :

- flexible

- error-prone

    - improper use of notation, wrong modelling

- no debugger !

# UML : Plan

- UML Class Diagram

  - Basic constructs that can be used to model Relational Databases in UML

  - Real Object-oriented features

# Class

- Set of elements sharing common properties
  - ex. peoples, animals, cars

- UML : draw a labelled box

| Employee |
|----------|
|          |

# Class : Attributes

- Attributes denote the properties of class objects
  - They are usually typed

- Write the attributes below the class name

| Employee |
| --- |
| name : string<br>age : int |

# Operations can specify

- Visibility (**+** public) (**-** private) (**#** protected)
- Return-type (*optional*; can be undefined)
- Multiplicity of parameter/return-type (*optional*)

| Employee |
| --- |
| name : string<br>dept : int |
| **+**setDepartment(<br> int dept **[1]**) : **bool [1]** |

58

# Instances

- The elements of a class
  - ex. the employee Bob

- Related to their class by a directed edge

| p1:Employee |
| --- |
| name : Bob<br>age : 42 |

| Employee |
| --- |
| name : string<br>age : int |

<<is_instance_of>>

# Binary Associations

- General relationships between elements of two classes
  - ex. an employee works for a department
  - a concrete instance of association is called <u>link</u>

- Binary association : *undirected* edge between classes

| Employee | Department |
|----------|------------|
|          |            |

works_for

# Associations : Links

- Relationships between instances of classes

| p1:Employee |
|---|
| name : Bob<br>age : 42 |

| d1:Department |
|---|
| area : sales<br>head : Alice |

works_for

# Associations : a tricky notation

- Any association is specified by three things

1. Its name

1. The cardinality constraints of the class elements

2. The role of the classes in the association (optional)

# Associations : Name (1)

| Employee | | Department |
|---|---|---|
| | works_for | |

- The name is a label placed in the middle of the edge

# Associations : Cardinality (2)



- This means that one instance of A participates in the association with **N** elements of the other class

# Associations : Cardinality (2)

| Employee | | 1 | Departement | |
|---|---|---|---|---|
| | | | | |

This means that an employee works for
exactly **1** department

# Associations : Cardinality (2)

Employee | 1..* | Departement

This means that an employee can work for
more than **1** department  (but at least one)

# Associations : Cardinality (2)

| Employee | 1..* | Departement |

This means that a department has at least one employee, with no upper-limit.

# Cardinality Specification

Employee ——— works_for

| Cardinality | Entity | Description |
|---|---|---|
| **1** | Department | exactly one |
| **0..1** | Department | zero or one |
| **0..*** | Department | zero or N (default) |
| **1..*** | Department | at least one |
| **2..4** | Department | from two to four |
| **2,4** | Department | two or four |

# Associations : Roles (2)



- This means that in the association an instance of $A$ plays **role**$_A$

# Putting everything together

| Employee |
|---|
|  |

1..*                                1

| Department |
|---|
|  |

works_for

employed          affiliation

# Reflexive Associations



Employee
1
supervisor
supervised
0..*
supervises

# Equivalent Formulation

Employee    0..*

supervised

supervisor    1

supervises

# Notation for Attributes in Associations

[UML] section 4 [BD-G] chapter XVII section 2



Employee     works_for     Department

starting_date : date

# Associations, Associations, Associations…

Often, a relationship between two entities combines:
- *association* features
    - the fact that two or more things are linked
- *representation* features
    - the details about this association


- Ex: a working contract can be seen
    - as a relationship
    (an association between an employee and a company)
    - as an entity (representation of a legal concept)


This duality is the source of most design problems !
- Recognize your own modelling choices !

# Choice 1



| |
|---|
| |
| starting_date : date |

| Company |
|---|
| |

| Employee |
|---|
| |

works_for

# Choice 2



Contract
starting_date : date

employer

Company

employee

Employee

# Associations Participating in Other Associations

| Employee |
|---|
|  |

works_for

| Company |
|---|
|  |

| |
|---|
| starting_date : date |

with_contract

| ContractType |
|---|
|  |

# Notation for 3-ary Associations

[UML] section 4.1 [UML2] section [3.3.4]

# Summing Up

UML for Relational Databases

[UML] section 4 [BD-G] chapter XVII section 2

- Basic UML
  - Classes, binary associations, n-ary associations
  - Can model relational databases

- The construct of the language we have seen so far are enough to model relational databases

# SubClass

- A subset of the instances of a Class
  - ex. every employee is a person

- A subclass <u>inherits</u> all superclass properties or <u>redefines</u> them.

| Person |
| --- |
|  |

| Employee |
| --- |
|  |

80

# SubClass =
# Generalization/Specialization
[UML] section 5

# Association modeled by a Class

[UML] section 4.3



- Consequence : such class can have proper attributes, operations and also other associations
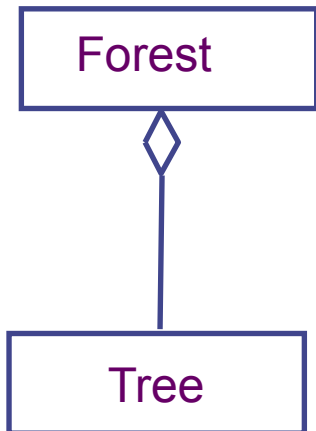
# Part-Of (binary) Association

- "A forest is made of trees"

# Aggregation

- Consequence 1 : a tree can exist even without a forest

```
┌─────────────┐
│   Forest    │
└─────────────┘
      ◇
      │
      │
┌─────────────┐
│    Tree     │
└─────────────┘
```
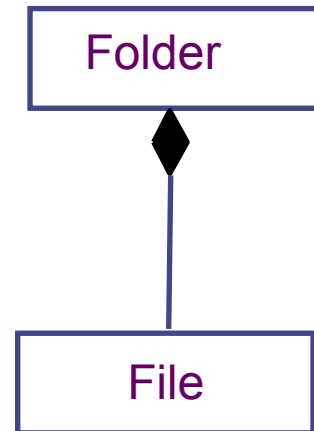
# Part-Of (binary) Association

- "A folder is made of files"

# Composition

- Consequence 1 : a file <u>cannot</u> exists without a folder

```
┌──────────────┐
│    Folder    │
└──────────────┘
        ◆
        │
        │
┌──────────────┐
│     File     │
└──────────────┘
```

# Aggregation          Composition

- Consequence 1 : a tree can exist even without a forest
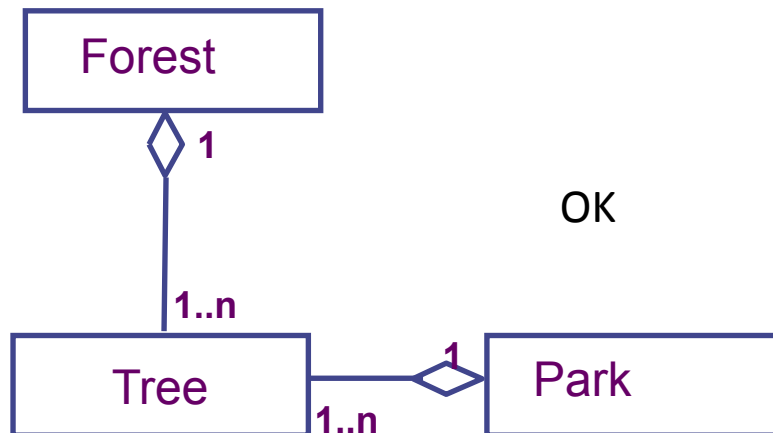
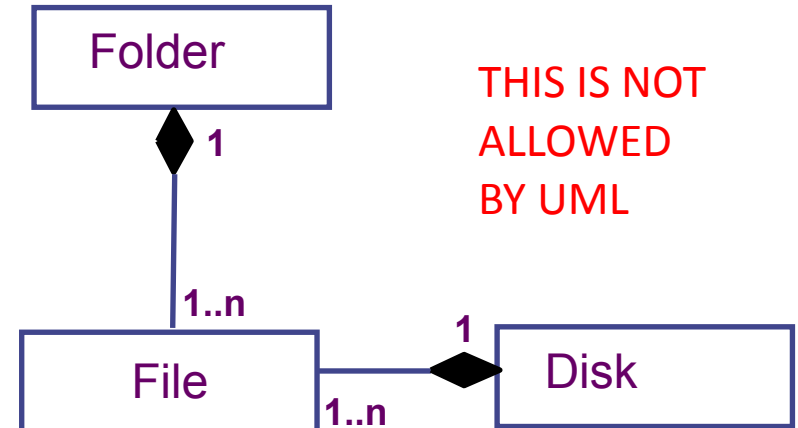- Consequence 1 : a file <u>cannot</u> exists without a folder

# Aggregation

# Composition

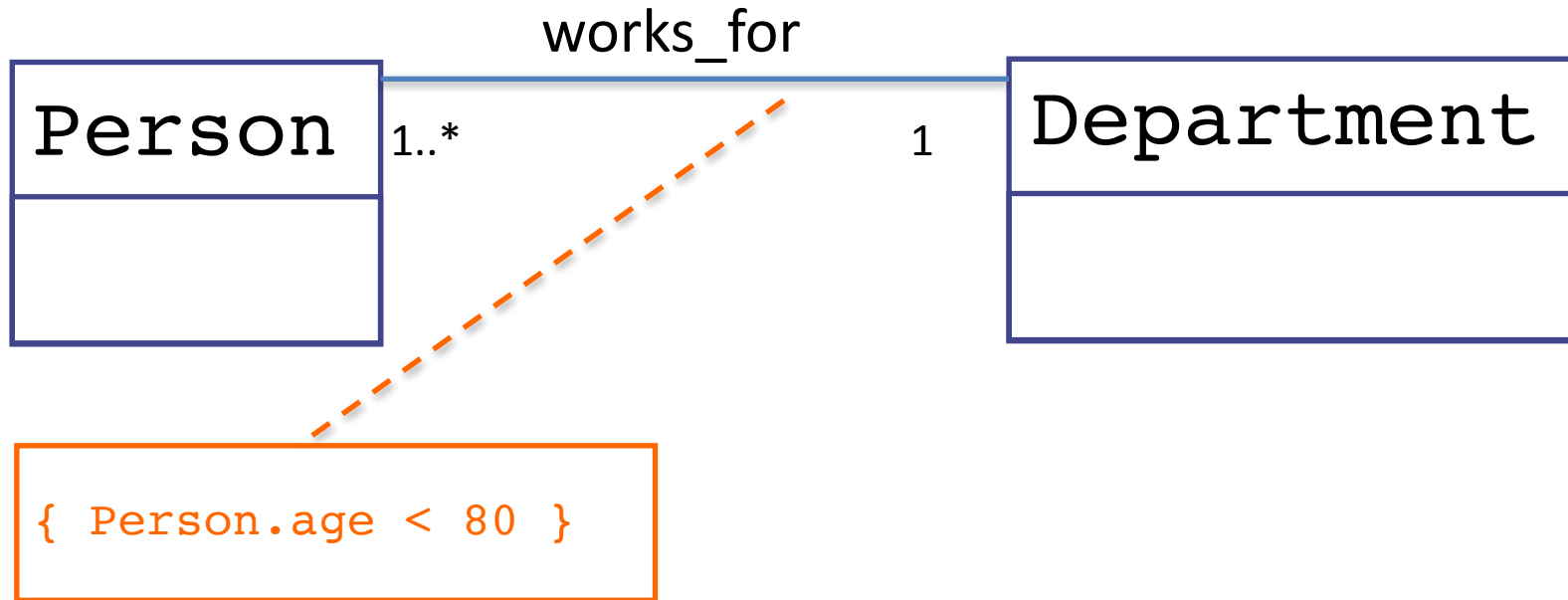- Consequence 2 : a tree can be part of both a forest and a park

- Consequence 1 : a file <u>cannot</u> exists without a folder



Forest

◇ 1

OK

1..n

Tree — ◇ 1 — Park

1..n

Folder

◆ 1

THIS IS NOT ALLOWED BY UML

1..n

File — ◆ 1 — Disk

1..n

# Constraints on Associations

works_for

Person    1..*          1    Department

{ Person.age < 80 }

- Consequence : only people whose age is less than 80 can participate in the association

# Summing Up

- Instances/Links/Operations
- Subclasses
- Aggregation and Composition