

# Semistructured Data and Mediation

**Prof. Letizia Tanca**

Dipartimento di Elettronica e Informazione  
Politecnico di Milano



# SEMISTRUCTURED DATA

FOR THESE DATA THERE IS SOME FORM OF STRUCTURE, BUT IT IS **NOT** AS

- **PRESCRIPTIVE**
- **REGULAR**
- **COMPLETE**

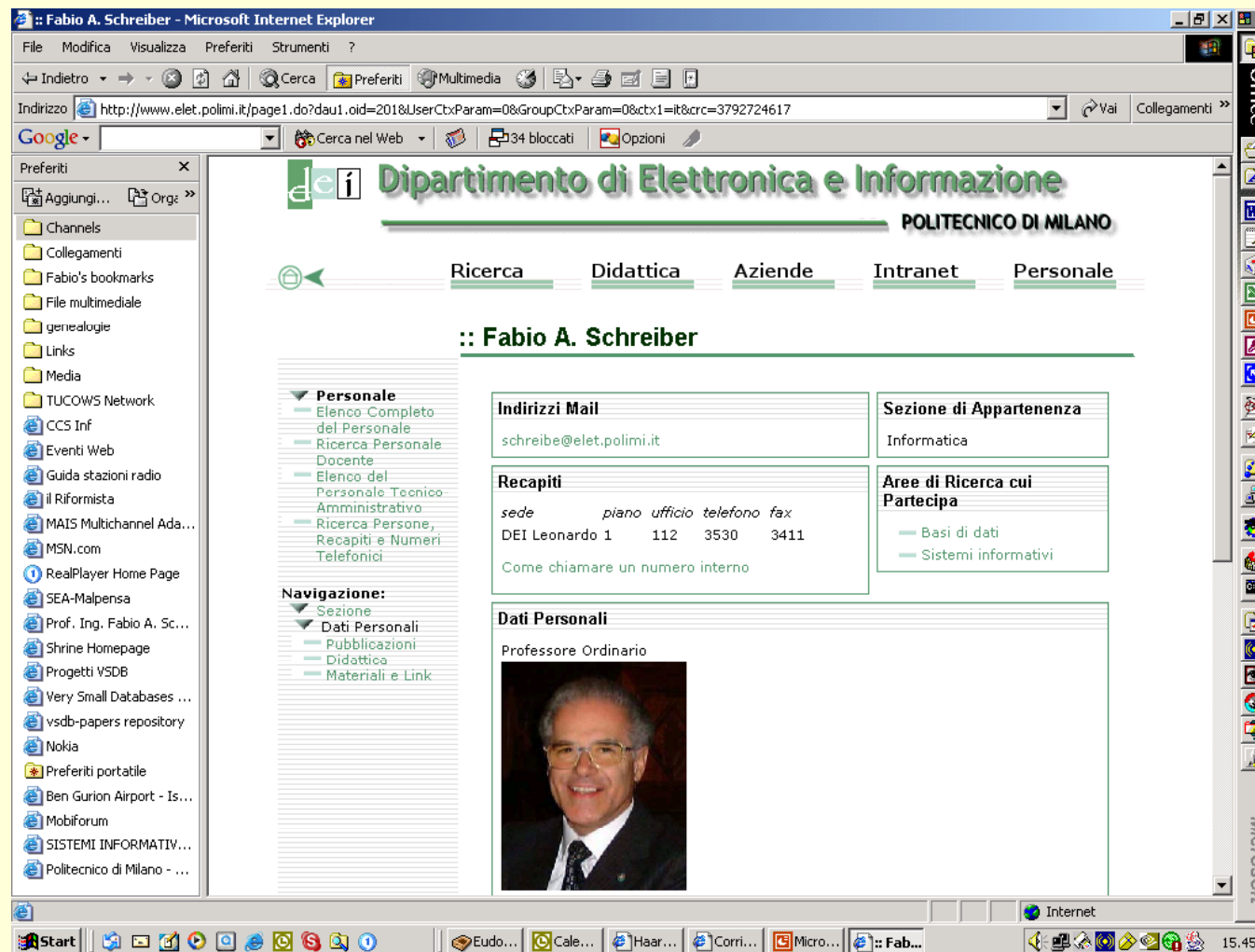
AS IN TRADITIONAL DBMSs

## **EXAMPLES**

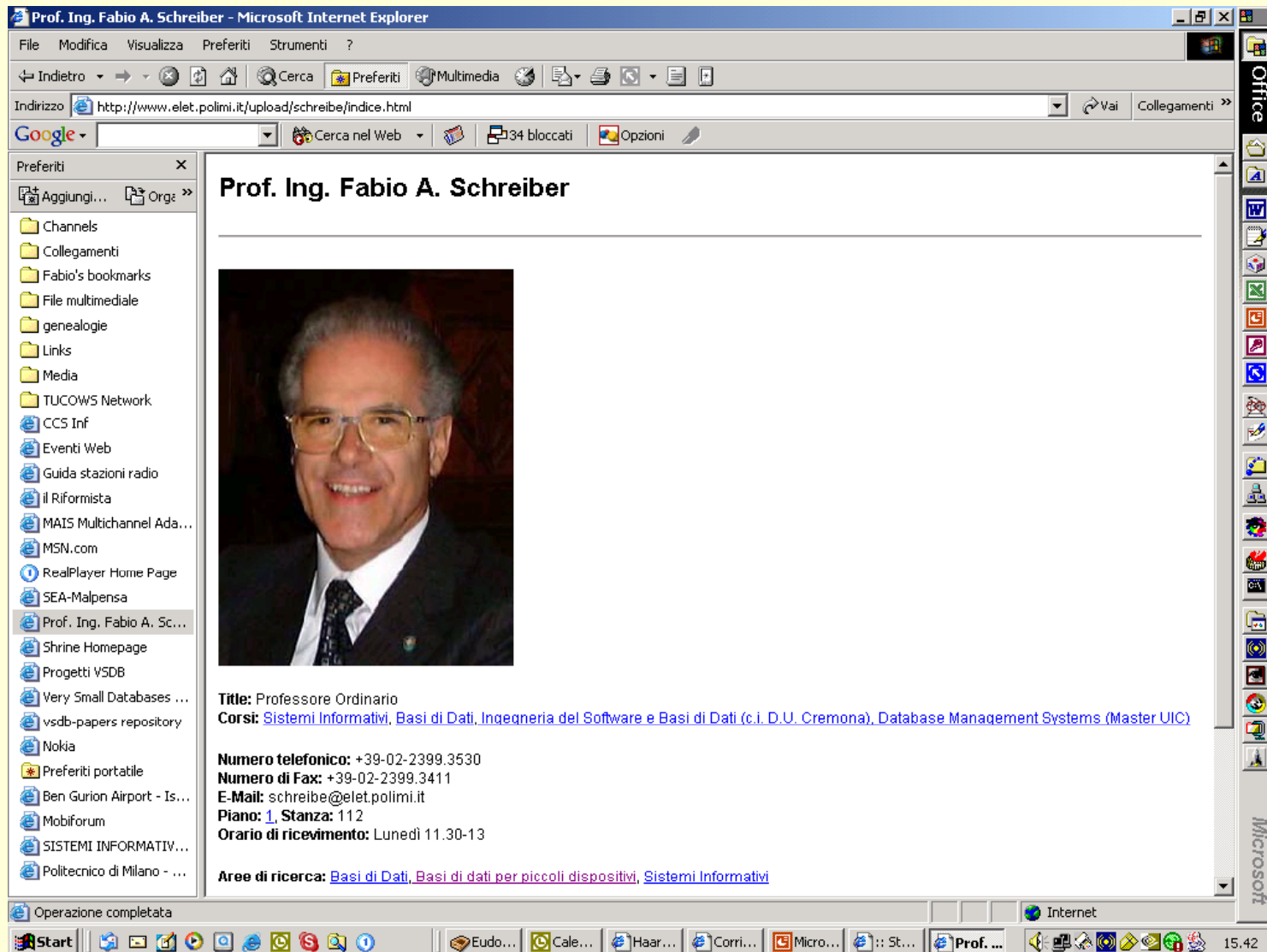
- WEB DATA
- XML DATA
- BUT ALSO **DATA DERIVED FROM THE INTEGRATION OF HETEROGENEOUS DATASOURCES**

# AN EXAMPLE OF SEMISTRUCTURED DATA

## a page produced from a database

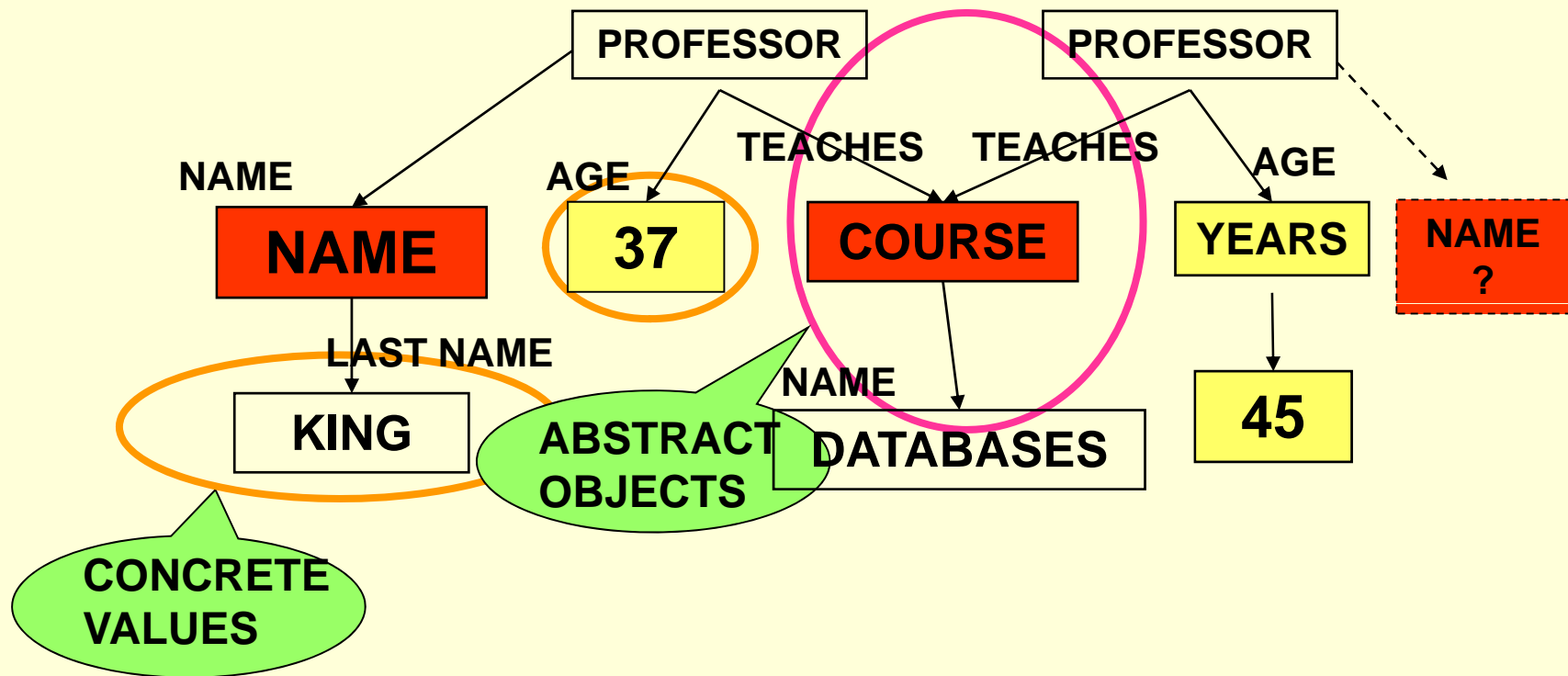


# AN EXAMPLE OF SEMISTRUCTURED DATA



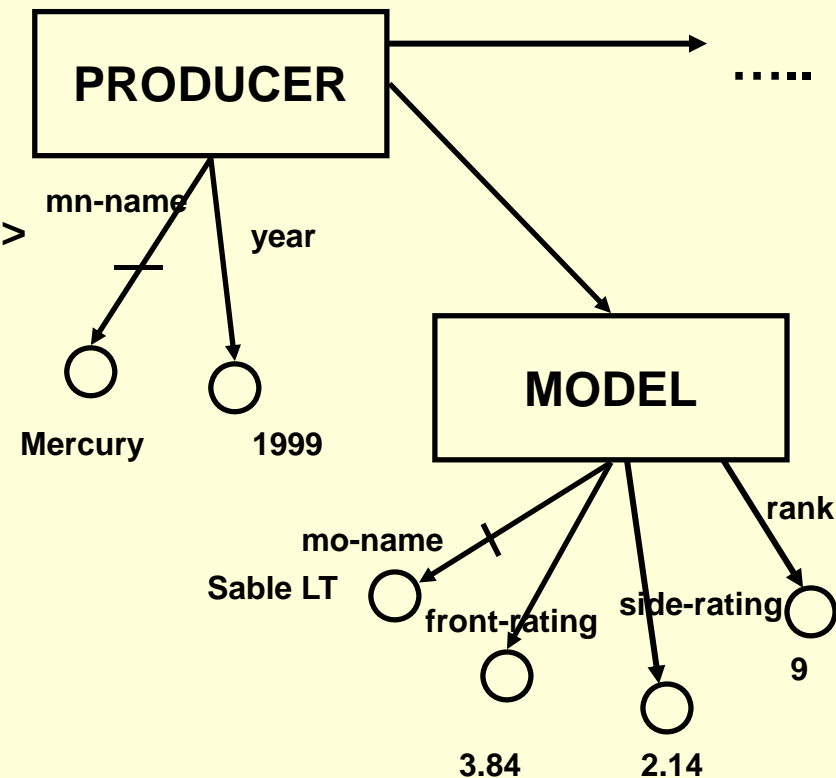
# AN EXAMPLE OF SEMISTRUCTURED DATA

GRAPH-BASED REPRESENTATION: THE IRREGULAR DATA STRUCTURE APPEARS MORE CLEARLY



# A SIMPLE XML DOCUMENT WITH ITS GRAPH BASED REPRESENTATION

```
<producer>
  <mn-name>Mercury</mn-name>
  <year>1999</year>
  <model>
    <mo-name>Sable LT</mo-name>
    <front-rating>3.84</front-rating>
    <side-rating>2.14</side-rating>
    <rank>9</rank>
  </model>
  .....
</producer>
```



# INFORMATION SEARCH IN SEMISTRUCTURED DATABASES

- WE WOULD LIKE TO:

- INTEGRATE
- QUERY
- COMPARE

DATA WITH DIFFERENT STRUCTURES  
ALSO WITH SEMISTRUCTURED DATA,  
JUST AS IF THEY WERE ALL  
STRUCTURED

# DYNAMIC INTEGRATION OF SEMISTRUCTURED DATABASES

AN OVERALL DATA REPRESENTATION SHOULD BE **PROGRESSIVELY BUILT**, AS WE DISCOVER AND EXPLORE NEW INFORMATION SOURCES



# SEMISTRUCTURED DATA MODELS

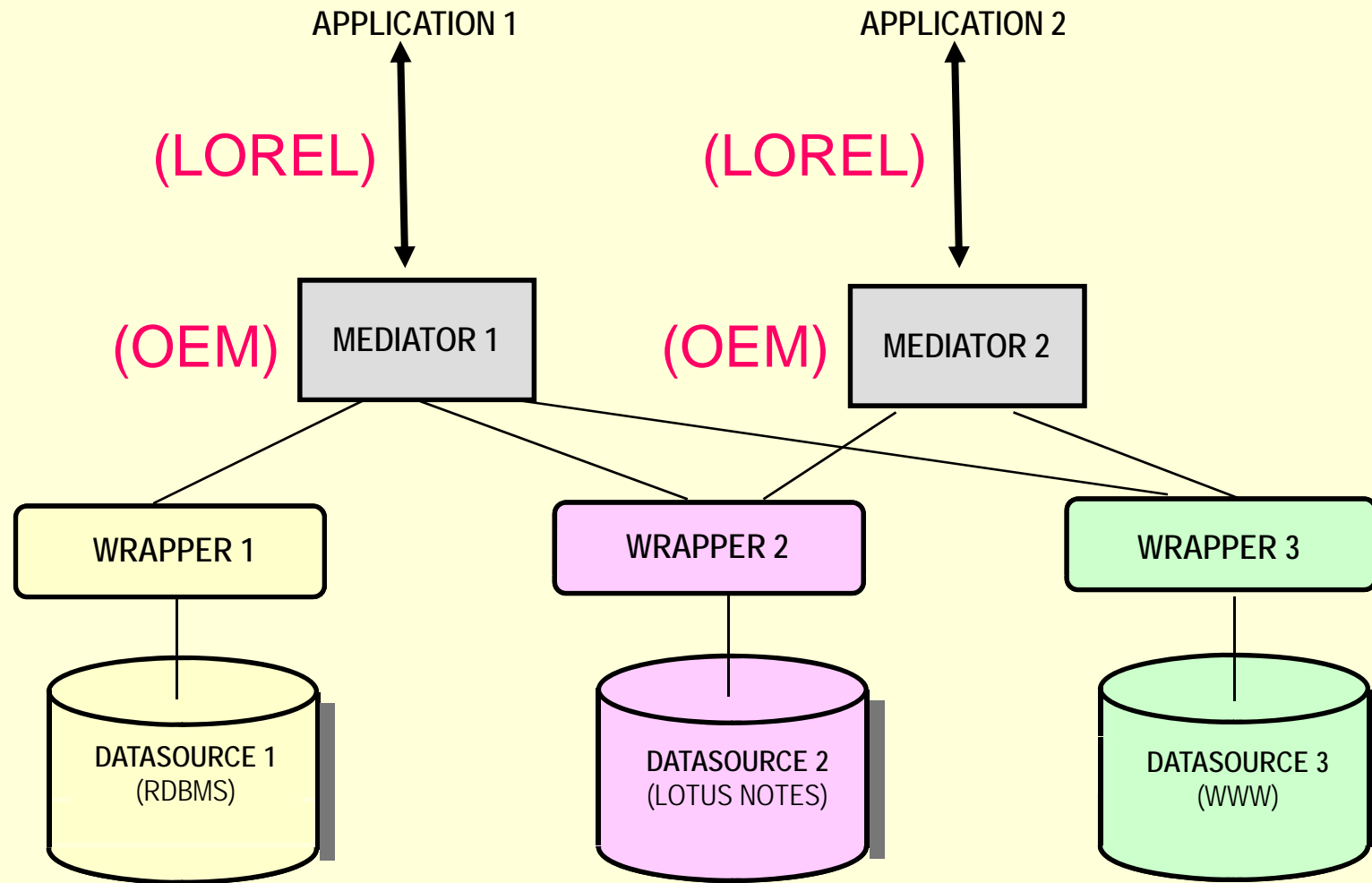
- BASED ON
  - TEXT
  - TREES
  - GRAPHS
    - LABELED NODES
    - LABELED ARCS
    - BOTH
- THEY ARE ALL DIFFERENT AND DO NOT  
LEND THEMSELVES TO EASY  
INTEGRATION

# Recall: MEDIATORS

The term mediation includes:

- the **processing** needed to make the interfaces work
- the **knowledge structures** that drive the transformations needed to transform data to information
- any **intermediate storage** that is needed (Wiederhold)

# TSIMMIS



# Mediator-based approach

## IN TSIMMIS:

- UNIQUE, GRAPH-BASED DATA MODEL
- DATA MODEL MANAGED BY THE MEDIATOR
- WRAPPERS FOR THE MODEL-TO-MODEL TRANSLATIONS

# OEM (Object Exchange Model) (TSIMMIS)

- Graph-based
- Does not represent the schema
- Directly represents data : self-descriptive

`<temp-in-farenheit,int,80>`

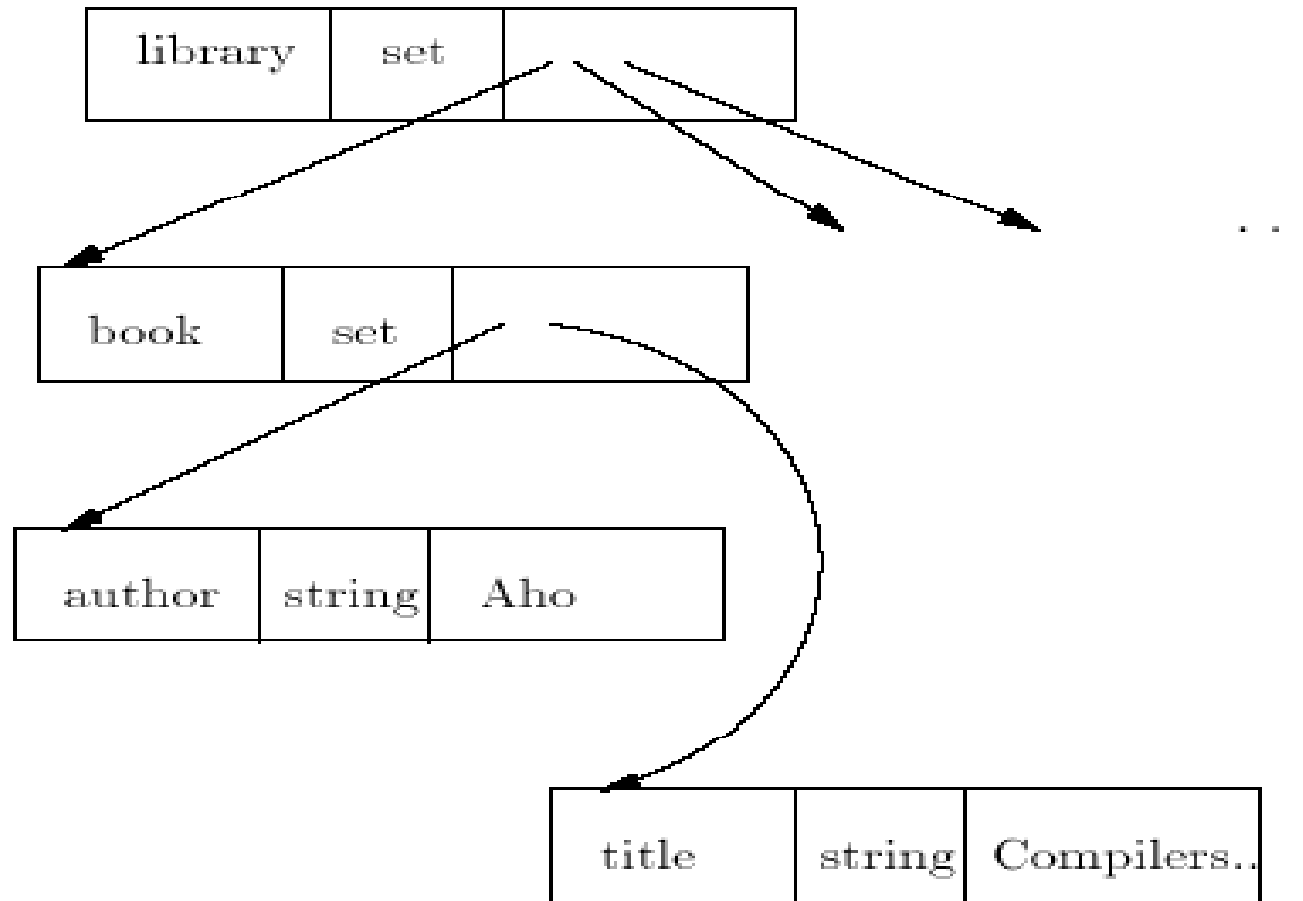
# Object structure

`<Object-id,label,type,value>`

## Nested structure

```
{set-of-temps, set, {cmp1, cmp2}}  
  cmp1: {temp-in-Fahrenheit, int, 80}  
  cmp2: {temp-in-Celsius, int, 20}
```

# OEM (Object Exchange Model) (TSIMMIS)



# Typical complications when integrating

- Each mediator is specialized into a certain domain (e.g. weather forecast), thus
- Each mediator must know **domain metadata** , which convey the data semantics
- On-line duplicate recognition and removal (no designer to solve conflicts at design time here)



# Query formulation

*“Find books authored by Aho”*

```
select library.book.title  
where library.book.author = "Aho"  
from library (if more than one root is available)
```

OK, but if this query must be produced at run-time, how does the user (or the system, if a transformation has to be applied) know that:

- A node *library* exists, which contains nodes *book*, which in turn contain fields *author* and *title*
- TSIMMIS uses the *Dataguide*: a-posteriori schema, progressively built while exploring the data sources

# TSIMMIS's language is *LOREL*

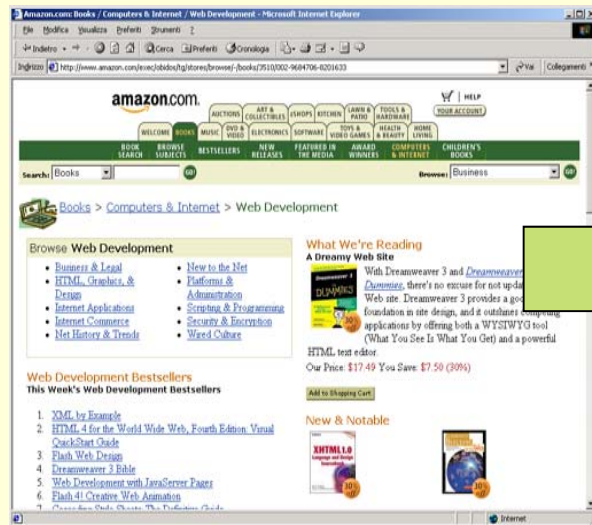
- *L*ightweight *O*bject *RE*pository *L*anguage
- Object-based
- Similar to OQL, with some modifications appropriate for semistructured data

```
select library.book.title
where library.book.author = "Aho"
from library
```

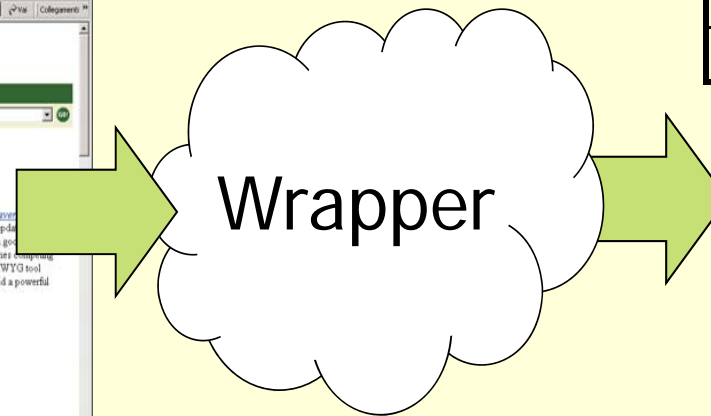
# WRAPPERS (translators)

- Convert queries into queries/commands which are understandable for the specific data source
  - they can **extend** the query possibilities of a data source
- Convert query results from the source's format to a format which is understandable for the application

# WRAPPERS



# HTML page



BookTitle	Author	Editor
The HTML Sourcebook	J. Graham	...
Computer Networks	A. Tannenbaum	...
Database Systems	R. Elmasri, S. Navathe	...
Data on the Web	S. Abiteboul, P. Buneman, D. Suciu	...

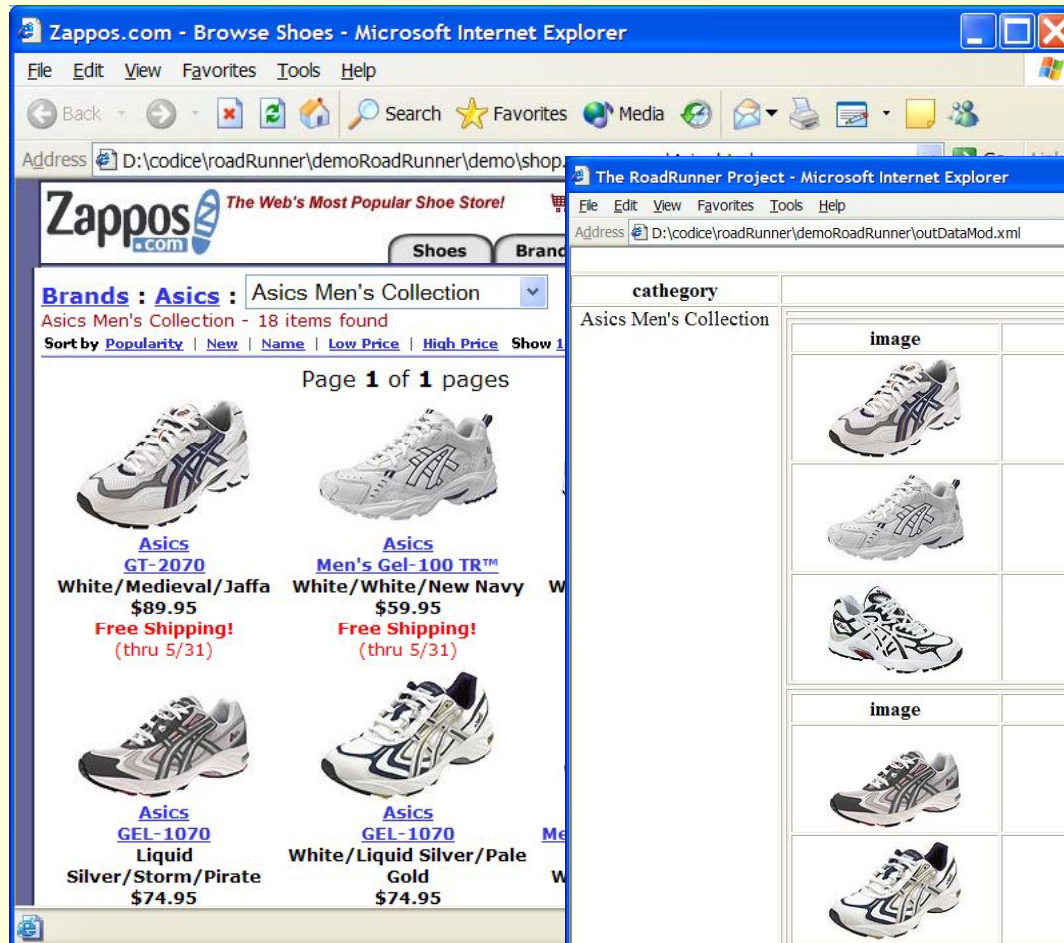
database table(s)  
(or XML docs)

# Example:

## extraction of information from HTML docs







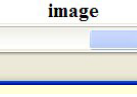

- Information extraction
  - Source Format: plain text with HTML tags (no semantics)
  - Target Format: relational table (possibly nested,  $NF^2$ ) or XML (we add *structure*, i.e. *semantics* )
- Wrapper
  - Software module which performs an *extraction step*
  - Intuition: use extraction rules which exploit the *marking tags*

# A complex extraction process



20-30KB IN HTML

The screenshot shows a data extraction tool window titled 'The RoadRunner Project - Microsoft Internet Explorer'. The address bar shows the file path 'D:\codice\roadRunner\demoRoadRunner\outDataMod.xml'. The tool displays the HTML content of the Zappos.com page, which is a table with 5 columns: image, brand, model, descr, and price. The table contains 10 rows of data, representing the first 10 items from the Zappos.com page.

image	brand	model	descr	price
	Asics	GT-2070	White/Medieval/Jaffa	\$89.95
	Asics	Men's Gel-100 TR™	White/White/New Navy	\$59.95
	Asics	GEL-MC PLUS® V	White/White/Russet	\$99.95
	Asics	GEL-1070	Liquid Silver/Storm/Pirate	\$74.95
	Asics	GEL-1070	White/Liquid Silver/Pale Gold	\$74.95
	Asics	Men's GEL-Foundation III	White/Cinder/Blaze	\$79.95
	Asics	GT-2070	White/Medieval/Jaffa	\$89.95
	Asics	Men's Gel-100 TR™	White/White/New Navy	\$59.95
	Asics	GEL-MC PLUS® V	White/White/Russet	\$99.95
	Asics	GEL-1070	Liquid Silver/Storm/Pirate	\$74.95

>10 attributes  
with nesting

# Problems

- Web sites change very frequently
- A layout change may affect the extraction rules
- Human-based maintenance of an ad-hoc wrapper is very expensive
- Better: *automatic wrapper generation*

## *Automatic wrapper generation...*

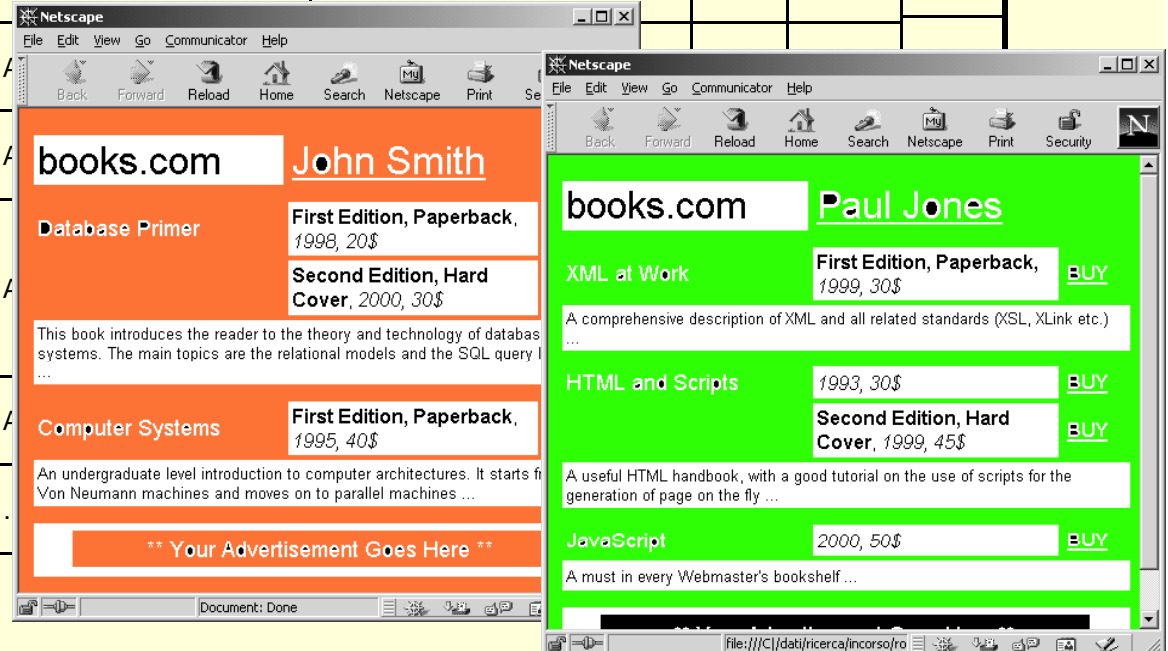
- We can only use them when pages are *regular* to some extent
- OK when:
  - Many pages sharing the same structure
  - e.g. pages are dynamically generated from a DB

→ *data intensive web sites*



# Online library

Name	Books				
	Title	Description	Editions		
			Details	Year	Price
John Smith	Database Primer	This book ...	First Edition, Paperback	1998	20\$
			Second Edition, Hard Cover	2000	30\$
	Computer Systems	A ...			
Paul Jones	XML at Work	A ...			
	HTML and Scripts	A ...			
	JavaScripts	A ...			
...	...	...			



# The ROAD RUNNER project

- Page Class
    - *It is the collection of all pages generated by the same script from a common dataset*
  - Schema Derivation
    - *Given a set of HTML sample pages, belonging to the same class, find the underlying dataset structure (database schema)*
- Solution: **Wrapper Generator**
- Underlying dataset structure
  - Extraction rules

A	B				
	C	D	E		
			F	G	H
John Smith	Database Primer	This book ...	First Edition, ...	1998	20\$
			Second Edition, ...	2000	30\$
	Computer Systems	All under			40\$

```
<HTML><BODY><TABLE>
  <TR>
    <TD><FONT>books.com<FONT></TD>
    <TD><A>John Smith</A></TD>
  </TR>
  <TR>
    <TD><FONT>Database Primer<FONT></TD>
    <TD><B>First Edition, Paperback</B></TD>,
  ...

```

<HTML><BODY><TABLE>	30\$
<TR>	
<TD><FONT>books.com<FONT></TD>	30\$
<TD><A>Paul Jones</A></TD>	
</TR>	45\$
<TR>	
<TD><FONT>XML at Work<FONT></TD>	
<TD><B>First Edition, Paperback</B></TD> ,	50\$
...	

```
<HTML><BODY><TABLE>
<TR>
  <TD><FONT>books.com<FONT></TD>
  <TD><A> #PCDATA</A></TD></TR>
(<TR>
  ( <TD><FONT> #PCDATA <FONT></TD>
    ( <TD><B> #PCDATA </B> )? </TD>
  ( <TR><TD><FONT> #PCDATA <FONT>
    <B> #PCDATA </B> </TD></TR> )+
)+...
```

## Target Schema

```
SET (
  TUPLE (A : #PCDATA;
    B : SET (
      TUPLE ( C : #PCDATA;
        D : #PCDATA;
        E : SET (
          TUPLE ( F : #PCDATA;
            G : #PCDATA;
            H : #PCDATA))))
```

**21° oct**

# Model Management approach

(Atzeni, Bernstein, others...)

- Given two different data models (e.g. OO and relational, or XML and OO, etc.) (when datasources are at least semistructured)
- Define general mappings from one model into another, which allow to
  - Map SQL schema to XML schema
  - Map data source to data warehouse
  - Map OO classes to data source tables, ...
- To this end, one possibility is to use a **metamodel**

# METAMODEL

- A METAMODEL IS AN ABSTRACT MODEL FOR THE SPECIFICATION OF CONCRETE MODELS
- TWO TYPES OF METAMODELS:
  1. GENERAL ENTITIES WHOSE SPECIALIZATIONS BECOME OBJECTS IN THE TARGET MODEL, E.G. GSMM
  2. ENTITIES DESCRIBING THE OBJECTS OF THE TARGET MODEL, E.G. GEOGRAPHIC DATA FILES (GDF)

# Using a metamodel for integrated data representation

- TRANSLATION OF DIFFERENT MODELS INTO A UNIQUE FORMALISM
- EASY A-PRIORI COMPARISON BETWEEN THE DIFFERENT MODEL'S FEATURES
- AUTOMATIC TRANSLATION DICTATED BY THE REPRESENTATION RULES OF THE CONCRETE MODEL INTO THE METAMODEL

# GSMM

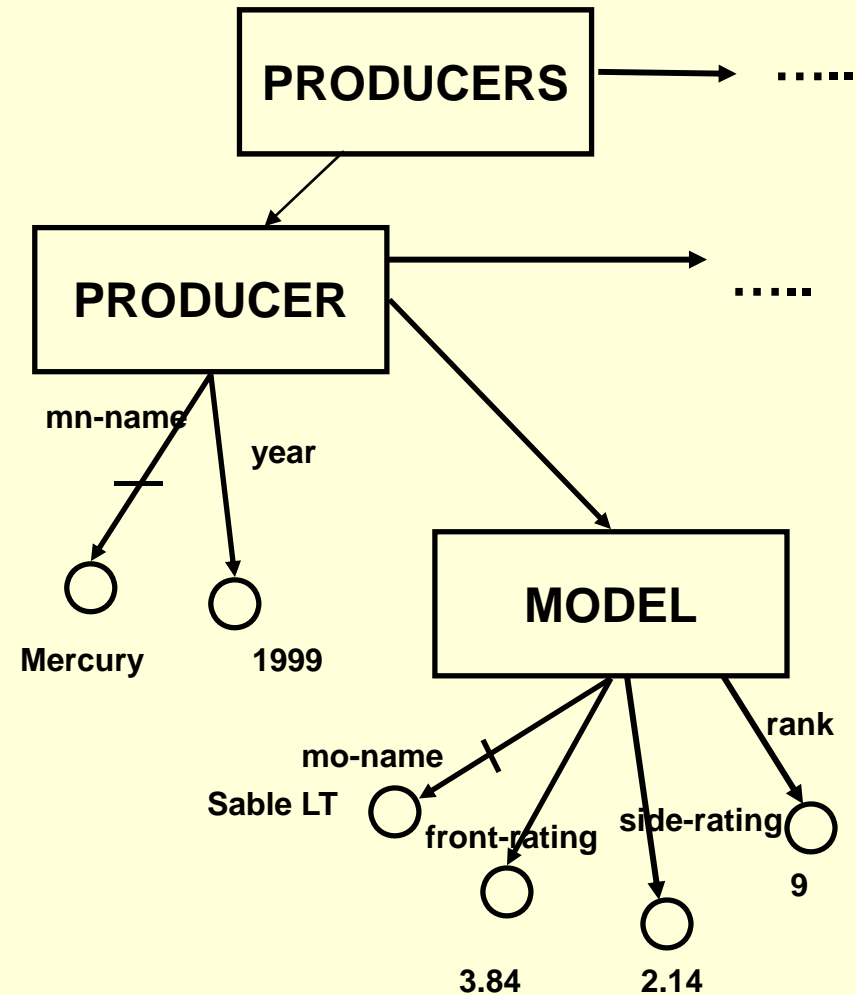
(General Semistructured Meta-Model)

- GRAPH-BASED
- DOES NOT REPRESENT THE SCHEMA: **SELF-DESCRIPTIVE** AS TSIMMIS
- GSSM REPLACES THE CONCEPT OF **SCHEMA** WITH THAT OF **CONSTRAINT**
- INTRODUCTION OF **FLEXIBILITY** IN DATA REPRESENTATION

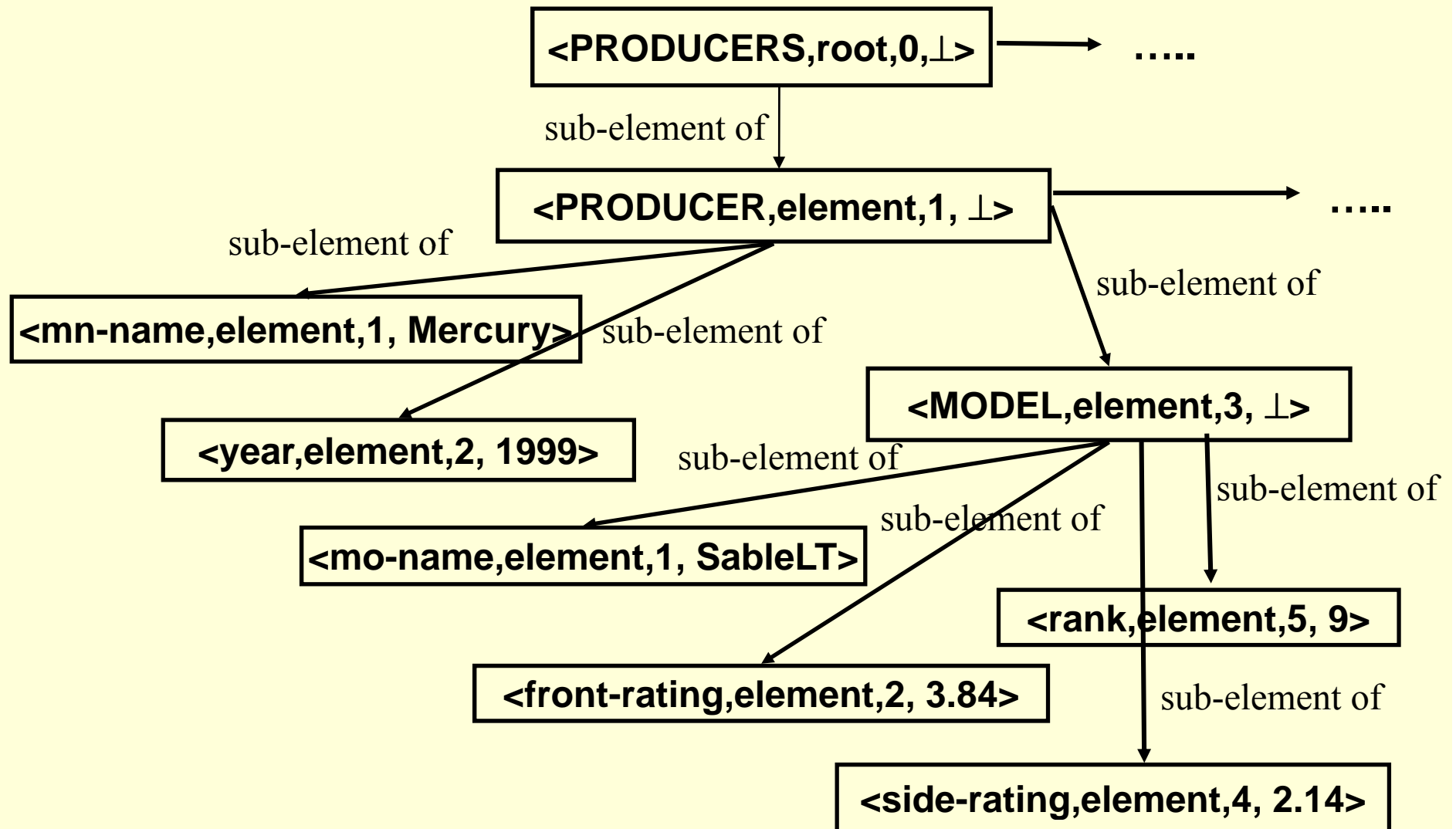


# An XML document...

```
<producers>
  <producer>
    <mn-name>Mercury</mn-name>
    <year>1999</year>
    <model>
      <mo-name>Sable LT</mo-name>
      <front-rating>3.84</front-rating>
      <side-rating>2.14</side-rating>
      <rank>9</rank>
    </model>
    .....
  </producer>
  ...
</producers>
```



# Its representation in GSMM



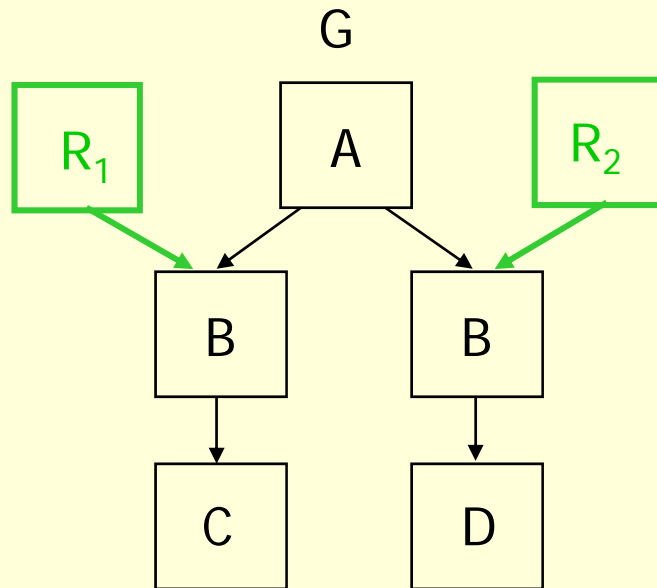
# GSL

## General Semistructured Language

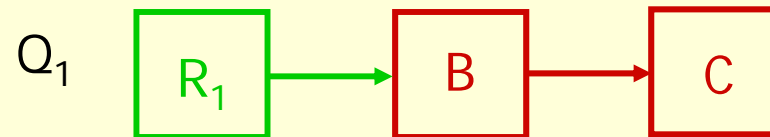
RULE BASED, WHERE

- A RULE IS A **COLORED GRAPH**
  - RED SOLID FOR POSITIVE PREMISES
  - RED DASHED FOR NEGATIVE PREMISES
  - GREEN SOLID FOR POSITIVE CONSEQUENCES
- RULES REPRESENT:
  - **QUERIES**
  - **CONSTRAINTS** (EMPTY CONSEQUENCE)

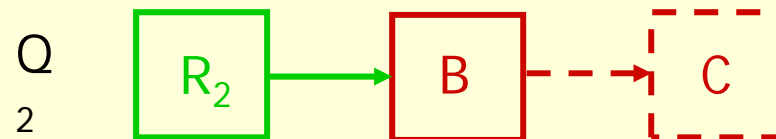
# GSL QUERIES



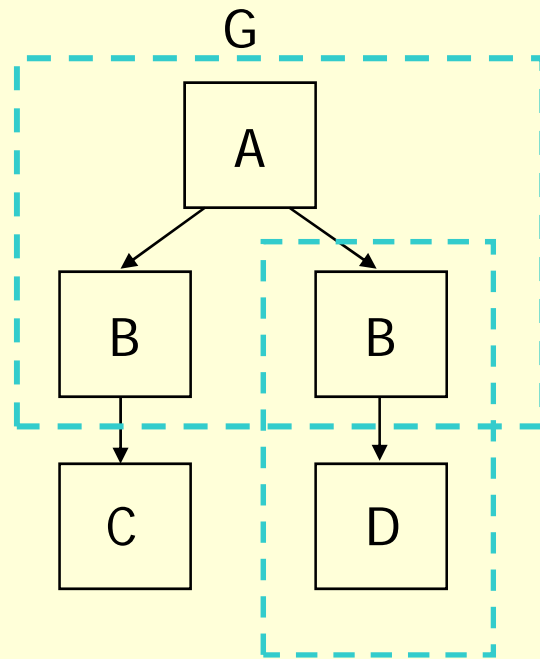
$Q_1$ : Find all the  $B$  nodes which have at least a  $C$  child and link them to a node  $R_1$



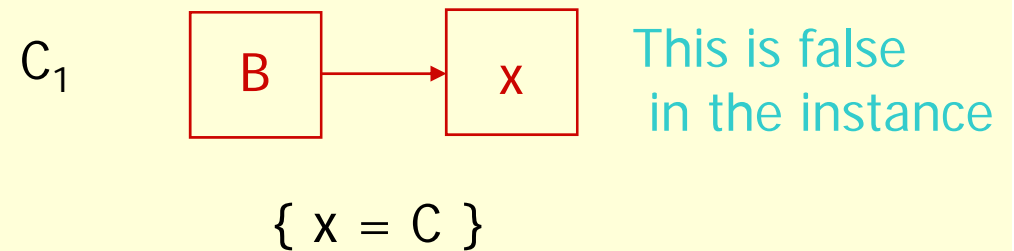
$Q_2$ : Find all the  $B$  nodes that do not have children  $C$  and link them to a node  $R_2$



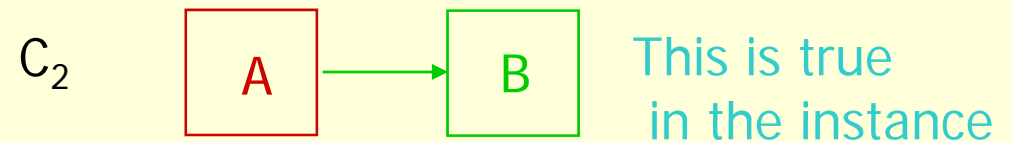
# CONSTRAINTS IN GSL



$C_1$ : Whenever a *B* node has a child, this is a *C* node



$C_2$ : Each *A* node has at least a *B* child



# APPLICATION OF THE GSMM METAMODEL

- THE METAMODEL ALLOWS THE CONSTRUCTION OF A GENERIC GRAPH WHICH REPRESENTS AN INSTANCE OF A CONCRETE MODEL
- PLUS **CONSTRAINTS** (also represented **by means of graphs**):
  - HIGH LEVEL, OR META- CONSTRAINTS – these dictate the syntactic rules of the object data model
  - LOW LEVEL CONSTRAINTS – as usual, these dictate the application domain semantics
- IT IS A METAMODEL **OF THE FIRST KIND**

# XML expressed in GSMM

- EACH **NODE LABEL** HAS CARDINALITY 4:  
 $n_i = \langle Ntag_i, Ntype_i, Norder_i, Ncontent_i \rangle$   
*Ntype<sub>i</sub>* says whether the node is the root, an element, a text, ...  
*Norder<sub>i</sub>* represents the node's position as a child of its parent node  
*Ncontent<sub>i</sub>* may assume a value of type PCDATA or value  $\perp$
- EACH **EDGE LABEL** HAS CARDINALITY 1:  
 $e_j = \langle (nh, nk), EL_j \rangle$   
 $EL_j = \langle Etype_j \rangle$   
 $Etype_j \in \{ \text{attribute of, sub-element of} \}$

# CONSTRAINT REPRESENTATION IN XML

<TAG1,TYPE1,ORDER1,CONTENT1>



<E-type>

<TAG2,TYPE2,ORDER2,CONTENT2>

{ E-type=SubElement-of  $\rightarrow$  TYPE1=element  
 $\wedge$  (TYPE2=element  $\vee$  TYPE2=text),  
 E-type=Attribute-of  $\rightarrow$  TYPE1=element  
 $\wedge$  TYPE2=attribute }

**HIGH LEVEL, OR META-CONSTRAINT:**  
 THE TYPE OF AN ARC DEPENDS ON THE  
 TYPE OF THE DESTINATION NODES

<TAG1,TYPE1,ORDER1,CONTENT1>



<E-type>

<TAG2,TYPE2,ORDER2,CONTENT2>

{ TYPE1=root  $\wedge$  CONTENT1= $\perp$  }

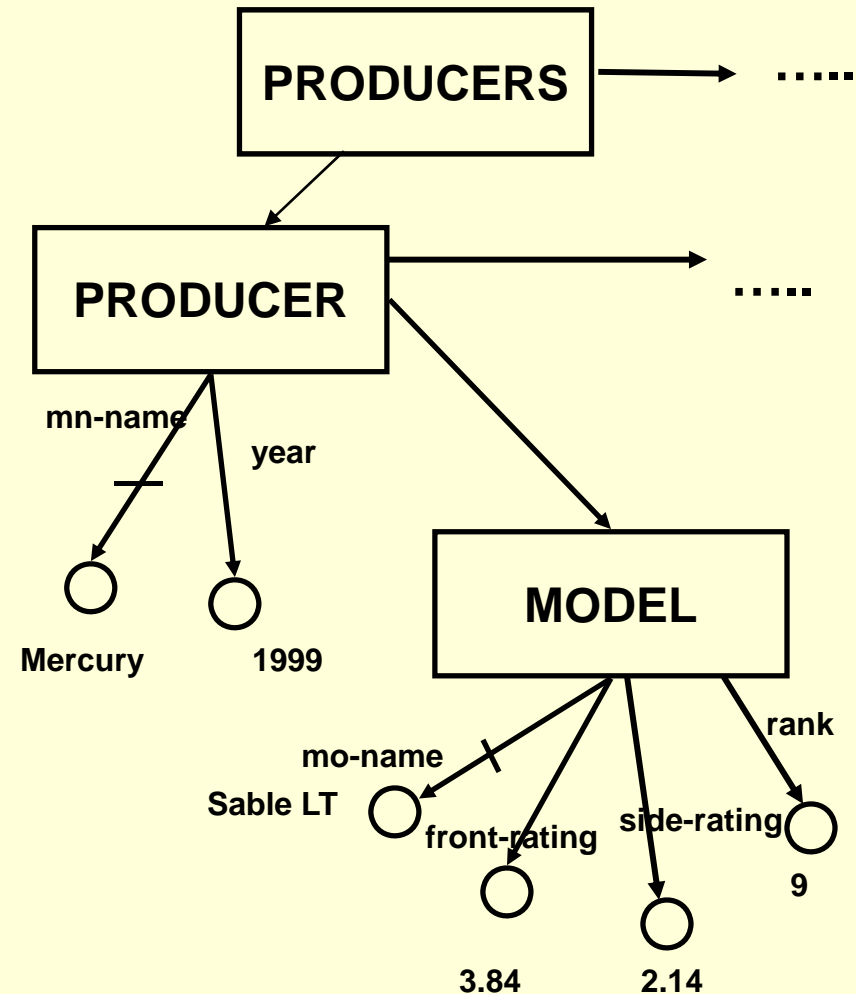
**HIGH LEVEL, OR META-CONSTRAINT:**  
 THE GRAPH ROOT HAS TYPE *root* AND  
 UNDEFINED CONTENT

HIGH LEVEL CONSTRAINTS CONTRIBUTE TO THE  
 DEFINITION OF THE STRUCTURE OF ANY XML DOCUMENT

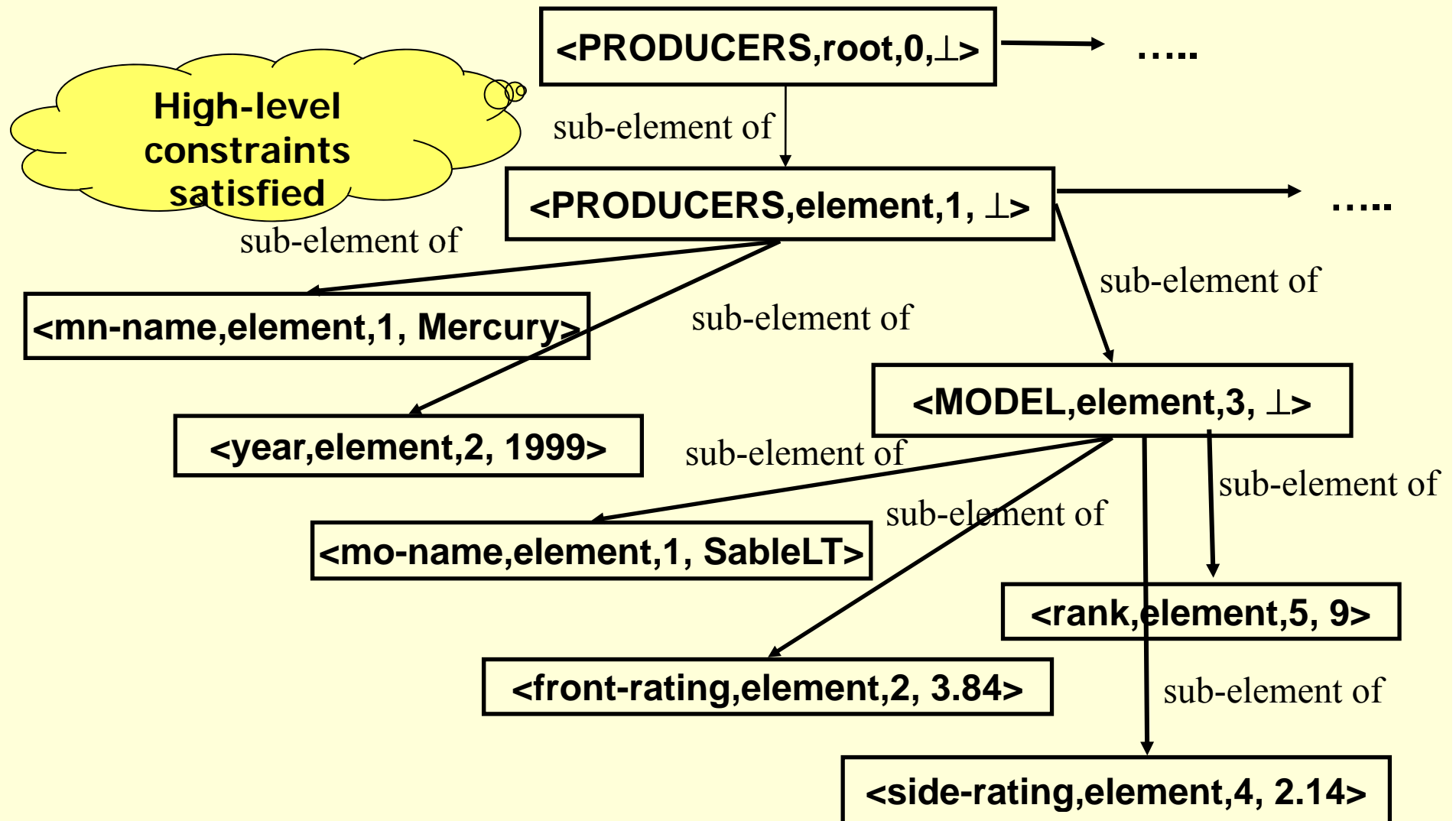


# An XML document...

```
<producers>
  <producer>
    <mn-name>Mercury</mn-name>
    <year>1999</year>
    <model>
      <mo-name>Sable LT</mo-name>
      <front-rating>3.84</front-rating>
      <side-rating>2.14</side-rating>
      <rank>9</rank>
    </model>
    .....
  </producer>
  ...
</producers>
```



# Its representation in GSMM



# LOW-LEVEL, OR APPLICATION-LEVEL, CONSTRAINTS

<producer,element,ORDER1,CONTENT1>



<Subelement-of>

<model,element, ORDER2,CONTENT2>

<producer,element,ORDER1,CONTENT1>



<Subelement-of>

<year,element,ORDER2,CONTENT2>

{ CONTENT2>1990 }

## Low-level constraint:

For each PRODUCER element there is at least one MODEL element

## Low-level constraint:

YEAR must be greater than 1990

**LOW-LEVEL CONSTRAINTS  
EXPRESS THE DOCUMENT'S SEMANTICS**

# Geographic Data Files (GDF)

- GDF is a standard, **used for describing roadmaps**
  - CITY STREET REGISTER
  - STREET SIGN ARCHIVE
  - TRAFFIC CONTROL SYSTEMS
  - CAR NAVIGATOR SYSTEMS (GPS)
  - .....
- IT IS A METAMODEL **OF THE SECOND KIND**

# GDF STRUCTURE

- DATASET OF 82 ASCII CHARACTERS RECORDS
  - ENTITIES
  - ATTRIBUTES
  - RELATIONS
- 11 INTEREST THEMES
  - ROADS AND FERRIES
  - BRIDGES AND TUNNELS
  - RAILROADS
  - RIVERS
  - PUBLIC TRANSPORTATION
  - ADMINISTRATION AREAS
  - .....
- 3 DESCRIPTION LEVELS FOR EACH THEME

# GDF

- **LEVEL 0 – TOPOLOGY**
  - A PLANAR GRAPH:
    - POINT
    - ARC
    - NODE
    - POLYGON
  - EACH ELEMENT IS UNIQUELY IDENTIFIED BY AN *ID*
  - TOPOLOGICAL RULES GOVERN INTERNAL STRUCTURE AND RELATIONSHIPS AMONG ELEMENTS

# GDF LEVELS

- **LEVEL 1 – ELEMENTARY ENTITIES**

THIS IS THE BASIS FOR THE **CITY STREET REGISTER**

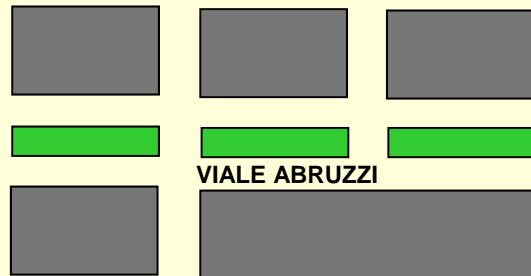
- ROAD ELEMENT
- JUNCTION
- TRAFFIC AREA

- **LEVEL 2 – COMPLEX ENTITIES**

THIS IS THE BASIS FOR THE **GEOGRAPHIC INFORMATION SYSTEMS**

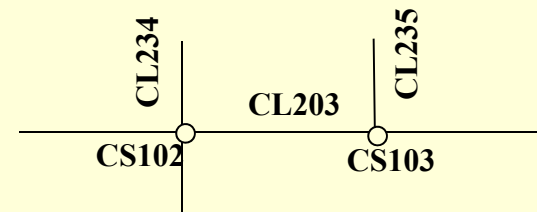
- STREET
- INTERSECTION

# GDF LEVELS

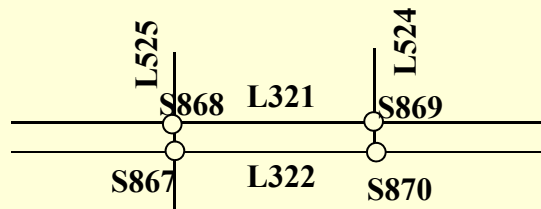


**CITY MAP**

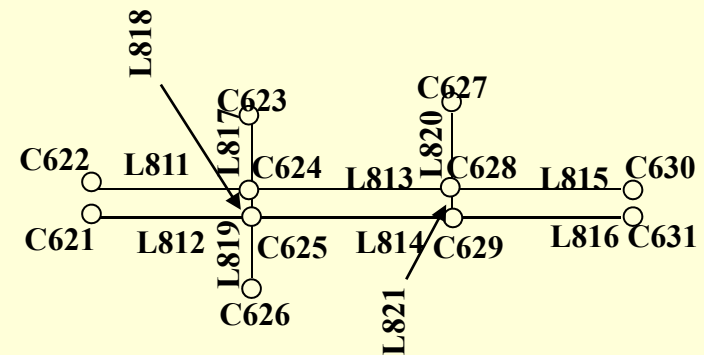
**LEVEL 2**



**LEVEL 1**

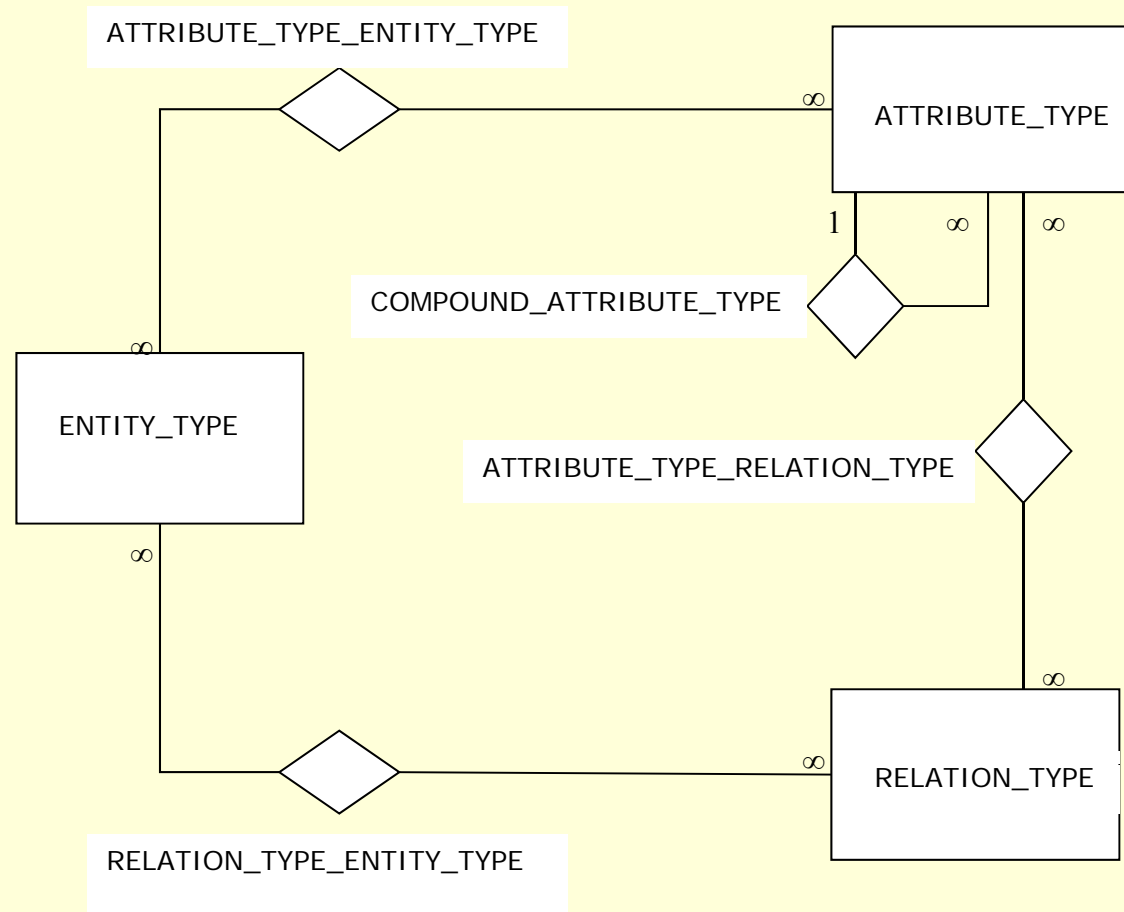


**LEVEL 0**

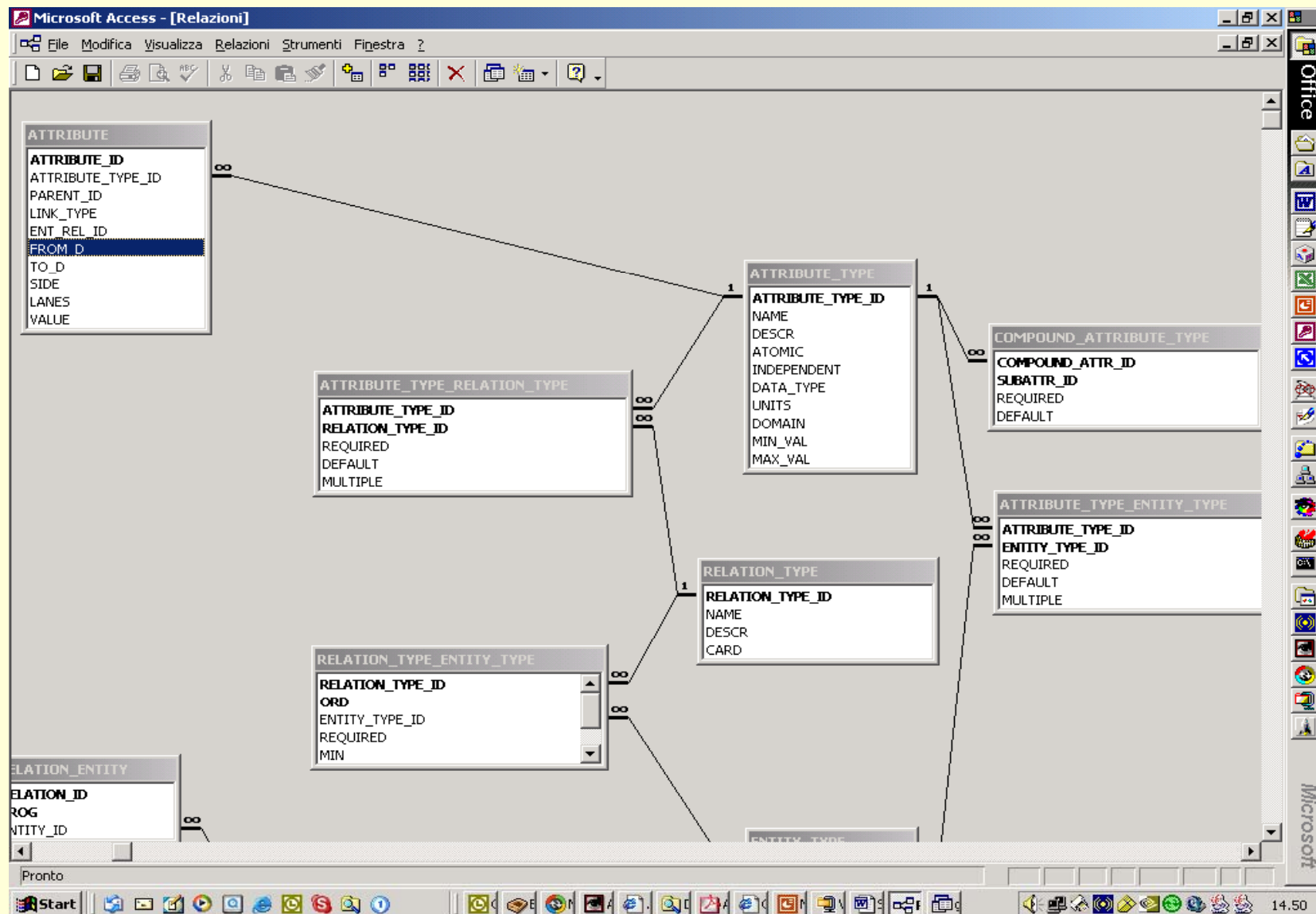




# META-ENTITY DEFINITION



# GDF META-SCHEMA



# AN ORTHOGONAL APPROACH: A STANDARD FOR METADATA

## MULTIMEDIA CONTENT DESCRIPTION INTERFACE

(MPEG-7) <http://www.tilab.com/mpeg>

- PROVIDES A SET OF TOOLS FOR DESCRIBING MULTIMEDIA CONTENT
- XML-SCHEMA BASED
- NOT BASED ON SPECIFIC APPLICATION DOMAINS, IT ENABLES EASY EXCHANGE AND REUSE OF MULTIMEDIA CONTENTS
  - SPEECH → TEXT
  - IMAGES ↔ TEXT
- IT IS APPLIED IN REAL-TIME AS WELL AS IN NON-REAL-TIME SITUATIONS:
  - OFF-LINE CONTENT STORAGE
  - ON-LINE CONTENT STORAGE
  - STREAM (NO PERMANENT STORAGE)

# MPEG-7 TOP LEVEL ELEMENTS

