

Semistructured Data & XML
(Summer Term 2011)

(c) Prof Dr. Wolfgang May
Universität Göttingen, Germany

`may@informatik.uni-goettingen.de`

Advanced Course in Informatics; 3+1 hrs/week, 6 ECTS Credit Points

A comprehensive German-English dictionary can e.g. be found at

<http://dict.leo.org/>

TASKS IN INFORMATICS

1. Implementing a proposed solution: a job
requires: good knowledge of common tools
2. Designing solutions: an interesting task
requires: solid knowledge of up-to-date concepts
3. Development of concepts: a fascination
requires: deep understanding and analysis of existing concepts

XML is a good example for all of them.

AIMS OF THE COURSE

- knowledge of the concepts of the XML-World, practical experiences
⇒ application-oriented
(requires also to work on your own)
- backgrounds why XML developed, and why it is as it is
⇒ understanding of concepts und developments
- underlying meta-concepts
⇒ as an example of “Informatics” as a whole

OVERVIEW

- other talk “Introduction to XML” ...

Chapter 1

Introduction

CONTEXT AND OVERVIEW

- Databases are used in many areas ... economics, administration, research ...
- originally: storage of information
late 60s: Network Data Model, Hierarchical Model
70s: Relational model, SQL – Lecture “Introduction to Databases”
- evolution: information systems, combination of databases and applications, distributed databases, federated databases, interoperability, data integration
- today: Web-based information systems, electronic data interchange
→ new challenges, semistructured data, XML
- tomorrow: Semantic Web etc.

1.1 Data Models

A *data model* defines the modeling constructs that can be used for modeling the application domain.

- *conceptual modeling*: application-oriented model
 - Entity-Relationship-Model (1976, only static concepts: entities and relationships, graphical)
Lecture: Introduction to Databases
 - Unified Modeling Language (UML 1.0: 1996)
comprehensive graphical formalism for modeling of processes, based on the object-oriented idea: classes and objects, properties, behavior (states and actions).
Lecture: Software Engineering
- *logical data models*: (e.g. relational model)
serve as *abstract data types* for implementations:
 - definitions of operations and their semantics, e.g. relational algebra
 - corresponding languages (as *application programming interfaces*): e.g. SQL
- *physical data models*: the implemented structures.

Data Model: Database Schema and Database State

Usually, for a database (for both, conceptual and logical models), its *schema* and its *state* are considered:

Database schema: the schema contains the *metadata* about the database, i.e., it describes the structure (in terms of the concepts of the data model).

The set of legal states is also described in metadata (e.g., by integrity constraints).

Database state: the state of a database is given as the currently stored information. It describes all objects and relationships that exist in the application at a given *timepoint*.

The database state changes over the time (representing changes in the real world), whereas the database schema is in general unchanged.

Logically spoken, the database state is an *interpretation* of the structure that is determined by the metadata.

Languages for Logical Data Models: In general, a language for operating on a data model consists of

- **Data Definition Language (DDL)** for schema definitions,
- **Data Manipulation Language (DML)** for manipulating and querying database states.

LOGICAL/IMPLEMENTATION DATA MODELS

... there are many different data models.

Basically, all database approaches are grounded on the concept of a “data item” (german: “Datensatz”).

- logical data models and implementation models
 - network data model (IDS (General Electric) 1964; CODASYL Standard 1971), hierarchical data model (IMS (IBM) 1965); data records,
 - relational model (Codd 1970), SQL (IBM System R 1973; products since 1979 (Oracle), ISO SQL Standard 1986); tuples
 - object-oriented model (ODMG 1993; OQL); objects
- document-data model (SGML)
- semistructured data models, XML; nodes: elements, attributes, text
 - why?
 - evolution and current situation

1.2 Relational Model

- *relational model* by E.F. Codd (1970, IBM San Jose): mathematical foundation: set theory
- only a single structural concept: *relation* for entities/objects and relationship types (note that the notions “entity” and “relationship” from the ER model [1976] were not yet defined!)
- properties of entities/objects and relationship types are represented by *attributes*
- a relation schema consists of a name and a set of attributes
Continent: Name, Area
- each attribute is associated with a *domain* that contains all legal values of the attribute. Attributes can also have *null values*:
Continent: Name: VARCHAR(25), Area: NUMBER
- a **(relational) database schema** is given by a (finite) set of (relation)schemata:
Continent: ...; Country: ...; City: ...; encompasses: ...

RELATIONS

- a (*database*) state associates a *relation* with each *relation schema*.
- the elements of a relation are called *tuples*.
Each tuple represents an object or a relationship:
(Name: Asia, area: 4.5E7)

Example:

Continent	
<u>Name</u>	Area
VARCHAR(20)	NUMBER
Europe	9562489.6
Africa	3.02547e+07
Asia	4.50953e+07
America	3.9872e+07
Australia	8503474.56

Relations: Example

Continent	
<u>Name</u>	Area
Europe	9562489.6
Africa	3.02547e+07
Asia	4.50953e+07
America	3.9872e+07
Australia	8503474.56

Country				
<u>Name</u>	<u>code</u>	Population	Capital	...
Germany	D	83536115	Berlin	
Sweden	S	8900954	Stockholm	
Canada	CDN	28820671	Ottawa	
Poland	PL	38642565	Warsaw	
Bolivia	BOL	7165257	La Paz	
..	

encompasses		
<u>Country</u>	<u>Continent</u>	Percent
VARCHAR(4)	VARCHAR(20)	NUMBER
R	Europe	20
R	Asia	80
D	Europe	100
...

- ... with referential integrity constraints
- abstract datatype for this model: relational algebra
- application interface: SQL

QUERY LANGUAGE: SQL

- Since 1973 “SEQUEL – Structured English Query Language” in IBM System R (E.F. Codd (Turing Award 1981), D. Chamberlin (2001: co-designer of XQuery)) etc.; Research-only (IBM continued to sell only IMS until SQL/DS (1980), DB2 (1983))
Stories: http://www.mcjones.org/System_R/SQL_Reunion_95/
<http://www.nap.edu/readingroom/books/far/ch6.html>
- 1974 INGRES (UC Berkeley, M. Stonebraker; NSF funding), QUEL language, open-source.
Led to the products INGRES (“Relational Technology Inc.” 1980, QUEL; since 1986 with SQL), INFORMIX (1981; since 1984 with SQL), SYBASE (1984, since 1987 with SQL)
- Oracle: founded in 1977 as “Relational Software” (L. Ellison worked before on a consultant project for CIA who wanted to use SEQUEL), 1983 renamed to “Oracle”.
Product: 1979 Oracle V2 (SQL), first commercial relational DB system.
- Standard SQL: 1986 ANSI/ISO (least common denominator of existing products); SQL-1 1989 (Foreign Keys, ...); SQL-2 1992 (multiple result tuples in subqueries, SFW in FROM, JOIN syntaxes, ...); SQL-3 1999 (PL/SQL etc) ...
- 1995: 80% of all databases use the relational model and SQL

QUERY LANGUAGE: SQL

```
SELECT name, percent  
FROM country, encompasses  
WHERE country.code = encompasses.country  
AND encompasses.continent = 'Europe';
```

- intuitive to understand,
- *clause-based, declarative* language,
- *set-oriented, closed*: result of (nearly) each expression is again a relation,
- *orthogonal constructs*, can be nested (nearly) arbitrarily,
- *functional programming paradigm*: each SQL query is a function that maps relations to another relation. Such functions can be nested.

... so far the things you have learnt in “Databases” about the relational model and SQL.

1.3 Concepts and Notions

- the relational model is a *data model*.
- (relational) databases follow a 3-level architecture:
 - *physical level/schema*: actual storage of tables in files, as sequenced records, with length indicators etc; additional index files, and allocation tables.
 - *logical level/schema*: *user level*.
Relational model (*logical data model*) with given database schema (table names, attributes, keys, foreign keys etc), relational algebra, SQL (*database language*).
Abstract, *declarative*, *set-oriented* language, distinguished notions of schema and state.
Internal: mapping to physical schema. Admin can change the physical schema and adapt the mapping without effecting the logical schema.
 - *external level (optional)*: possible views, given by SQL queries.
A *view* is (any kind of) a mapping from underlying “base” data to derived information.
- note: SQL is the only language with which users work on relational data. Relational data exists only inside databases.

CONCEPTS: PREVIEW

- network data model: mainly a physical data model; "logical" model on a very low level of abstraction.
No database language, only some data-management-oriented operations extending a common programming language.
- relational model: abstract/logical data model, relational algebra, declarative, set-oriented query+update language.
- early semistructured data models (OEM, F-Logic etc.): not comparable, separate experiments how to extend functionality without losing the advantages from relational databases and SQL.
- for XML there are several languages ("views" can also be defined in several ways), and XML exists also as a data structure used in non-database tools.

1.4 Aside: Really Declarative Languages ...

SQL is already called “declarative”: express what, not how.

But there is an even more declarative language family: *logic-based* languages.

Relational Calculus, Datalog

- Facts (tuples) are called “atoms”:
`country(“Germany”, “D”, 83536115, 356910, “Berlin”, “Berlin”),`
`city(“Berlin”, “Berlin”, “D”, 3472009), etc.`
- queries are given as “patterns” with free variables:
`?- country(N,C,Pop,Area,CapProv,Capital).`
yields a set of *answer bindings* for the variables N,C,Pop,Area,CapProv,Capital.
- Projection via *don't care* variables:
`?- country(N,_C,Pop,_Area,_,_).`
yields a set of *answer bindings* for the variables N and Pop.
- Selection: `?- country(“Germany”, “D”, Pop, Area,_,_).` binds only Pop and Area.

Relational Calculus (cont'd)

- Selection as Conjunction:

?- country(N, C, Pop, _,_,_), Pop > 1.000.000. binds N, C, Pop

?- country(N, _, _Pop, _,_,_), _Pop > 1.000.000. returns only the set of names of countries with more than 1000000 inhabitants.

- Joins as conjunctions:

?- country(N,_C,_,_,Area,_,_), encompassed(_C,Cont,Perc), continent(Cont, ContArea).

?- country(_, "D",_,_,_,CapProv,Capital), city(Capital, CapProv, "D", Pop, _, _)

Datalog

- Views as “derived/virtual relations”:

ancestor(X,Y) :- parent(X,Y).

ancestor(X,Y) :- ancestor(X,Z), parent(Z,Y).

?- ancestor(X,Y).

can express e.g. transitive closure (recursive rules and fixpoint semantics).

- flows_transitive(Name, Sea) :- river(Name, _, Sea, _).

flows_transitive(Name, Sea) :- river(Name, River, _ , _), flows_transitive(River, Sea).

?- flows_transitive(River, Sea).

[see F-Logic/transitive-rivers.flp]

REMARK

The term “Database” does not only mean the relational model & SQL, but is a general notion:

- persistent storage
- mass data
- multiuser concepts
 - access control/safety
- transaction concepts
 - correctness/consistency
 - safety: rollback, recovery

... all these are *general concepts* that apply for the network data model, the relational model, the object-oriented model, and also for XML databases.

Chapter 2

Database Concepts and Extensions

- The notion of “semistructured data (SSD)” has mainly been coined by the “TSIMMIS” project (The Stanford-IBM Manager of Multiple Information Sources, 1995-2000; persons: J. Ullman, H. Garcia-Molina, J. Widom, Y. Papakonstantinou).
- The problem has already been investigated before in several areas and projects.

WHY SEMISTRUCTURED DATA?

... mainly two requirements:

1. *data integration* from different sources

(late 1980s/early 1990s):

- increasing networks
- combination of contents of several databases
 - * multi-database-systems
 - * federated database systems
 - * different schemata
 - * mostly only different relational schemata,
 - partially also under the aspect of integration of metadata into the DML – this aspect is originally independent from semistructured data.
 - * sometimes different data models (“legacy”-databases according to earlier data models)
 - * since mid-90s increasingly data from the Web

WHY SEMISTRUCTURED DATA?

2. storage of “unregular” data:

no fixed/homogeneous/known schema, many null values (e.g. biochemistry)

- data exchange (B2B); standard formats e.g. for suppliers in automobile industry
- partly also full-text portions,
- management of document content
 - * coarsely structured
SGML (special form: HTML)
- annotated binaries (pictures, films, etc.)
- mixed forms between databases and documents
 - * collections: (tax) laws, partially in SGML
 - * health care and clinical information systems

THE EVOLUTION TOWARDS XML

The evolution in the area of semistructured data and XML combined concepts, experiences and developments from many previous approaches:

- network data model, hierarchical model (“legacy”-databases),
- relational databases,
- object-oriented databases,
- distributed and federated databases,
- data integration (purely relational environments, or mixed ones),
- document management.

Different lines of evolution have been brought together with XML & friends:

⇒ (nearly) nothing new, but a perfect combination!

Textbook on “Databases” in general (but without document management and XML):

- R.Elmasri, S.Navathe: “Foundations of databases”/“Grundlagen von Datenbanksystemen”. Pearson Studium, 3rd edition, 2002.

2.1 Early Databases: the Network Data Model

Situation 1960: first primitive “high-level” programming languages for “calculations”

- FORTRAN 1957: “formula translator”
- COBOL 1959: “common business-oriented language”

Goal: somehow store and organize lots of data:

- first development in the database system *IDS (Integrated Data Store)* at *General Electric* (Bachman & Williams, 1964; Turing-Award 1973)
- specification of the “[Conference on Data Systems Languages Data Base Task Group \(CODASYL\)](#)”, 1971.
- products: e.g. VAX-DBMS (Digital Equipment)

NETWORK DATA MODEL

- data is stored in *data records*,
- classified by *data record types*, with attributes (name and datatype to be specified).

Country				
Name	Code	Population	Area	...

City		
Name	Population	...

- Sample data records:

"Germany"	"D"	83536115	356910	...
-----------	-----	----------	--------	-----

"Berlin"	"3472009"	...
----------	-----------	-----

- So far, the same as the mapping of entity types in the relational model.
- difference: the *organization* of the records (and their relationships) in the database ...

RELATIONSHIPS: SET TYPES

Relationships are represented as sets: “all B that are in a given relationship with a certain A”
(E.g. all cities in a given country)

Definition of set types:

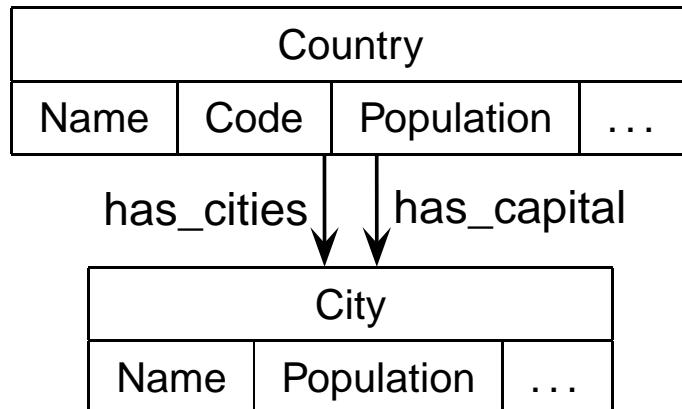
- name of the set type
- owner record type (“owner”; “where the relationship starts from”)
- member record type (“member”)

A set instance then represents the relationship for exact one owner of type “A”. Each instance of a set consists of

- a data record of the owner record type
- an *ordered* set of data records of the member record type
- **comparison with XML:** parent-children relationship, ordered children, but all of the same type
- intuition: not as a set, but as a wire that is fixed at the owner and then pulled through all members.

RELATIONSHIPS: SET TYPES

Graphical representation:
"Bachman Diagram"



has_cities:	Germany	D	83536115	...
Berlin			3472009	...
Hamburg			1705872	...
Frankfurt			652412	...
			⋮	

has_capital:	Germany	D	83536115	...
Berlin			3472009	...

- similar sets for France//Paris/Lyon/Marseille/... and France//Paris
- a member record can belong to only one instance of a set of *each* set type (thus, only 1:N-relationships can be modeled directly)
- n:m relationships: later

ENTRY POINTS

- system-owned instances of a set serve as entry points
(e.g. an instance of a set “countries” whose members are the country-data records)

ACCESS OPERATIONS

Access to (and navigation through) the database only via sets.

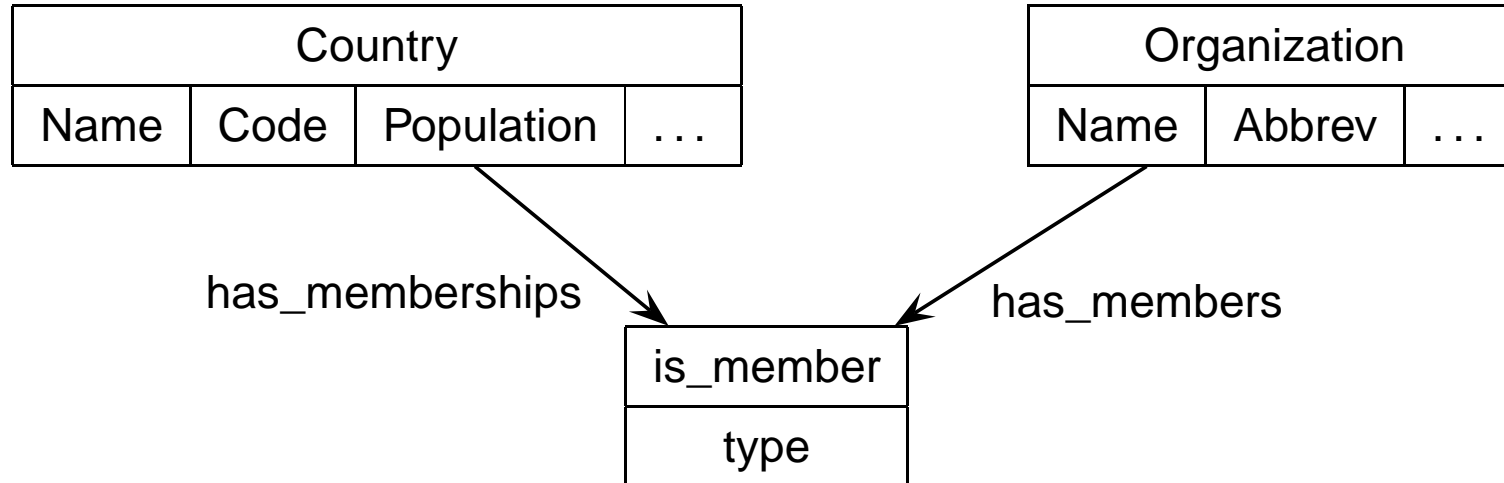
Actually, this is again an *abstract datatype*:

- access to the attribute values of a record,
 - an *iterator* (first, next) for traversing the relationships ,
 - a selector “find_owner” for inverse relationships.
- ⇒ the user does not explicitly work with pointers or identifiers, but already uses the semantic notions of the data model.

N:M-RELATIONSHIPS

Cannot be represented by a single set type (analogously for attributed relationships).

- split into a 1:M and an inverse N:1-relationship
Problem: consistency maintenance (symmetry!)
- introduce an auxiliary data record type that represents the relationship, and two set types:



- later, there is a mapping from the ER model (1976) to the network model.

ORGANISATION OF THE SET TYPES

Each data record contains reference entries for each set type where it belongs to (either as owner or as member):

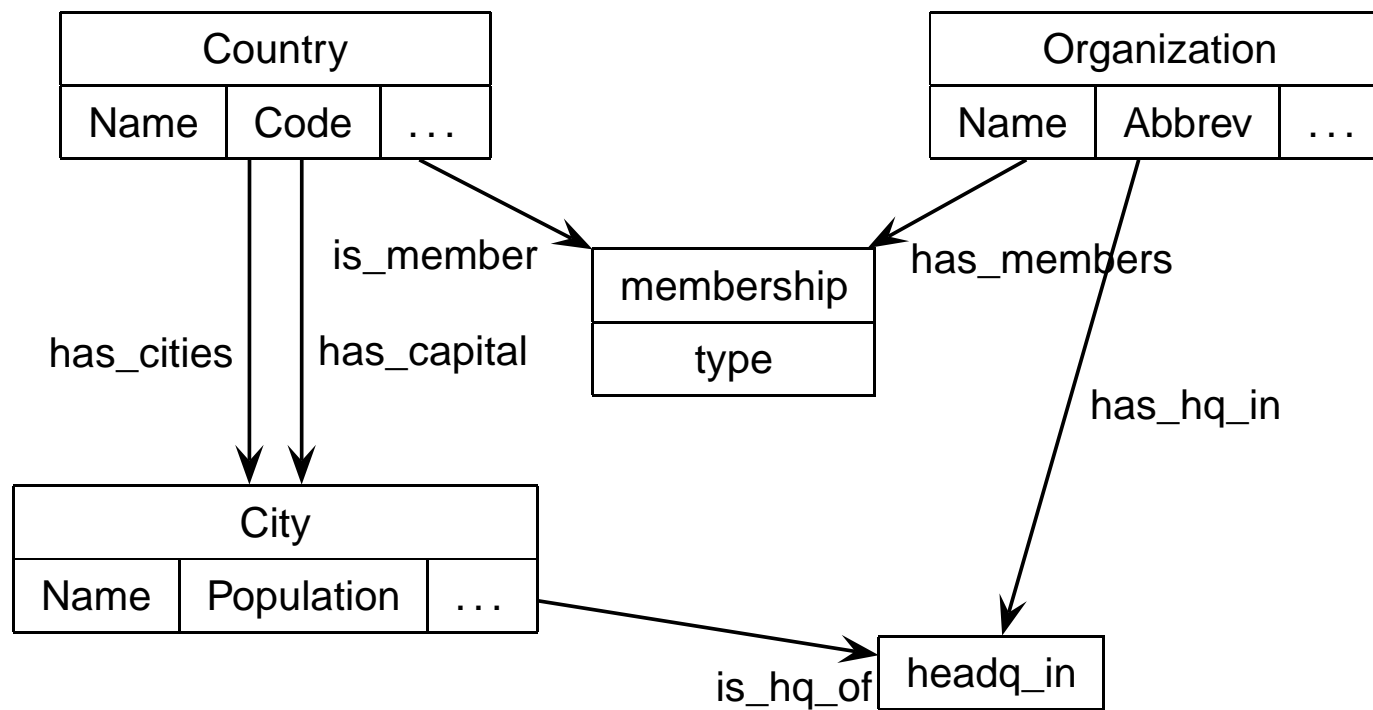
- as owner: a “first”-reference, labeled with the name of the set type, pointing to the first member record
- as member:
 - a “next”-reference, labeled with the name of the set type, pointing to the next member record
 - additionally a labeled backwards pointer to the owner of the set instance
 - a labeled null pointer if there exists no first/next element.

Exercise 2.1

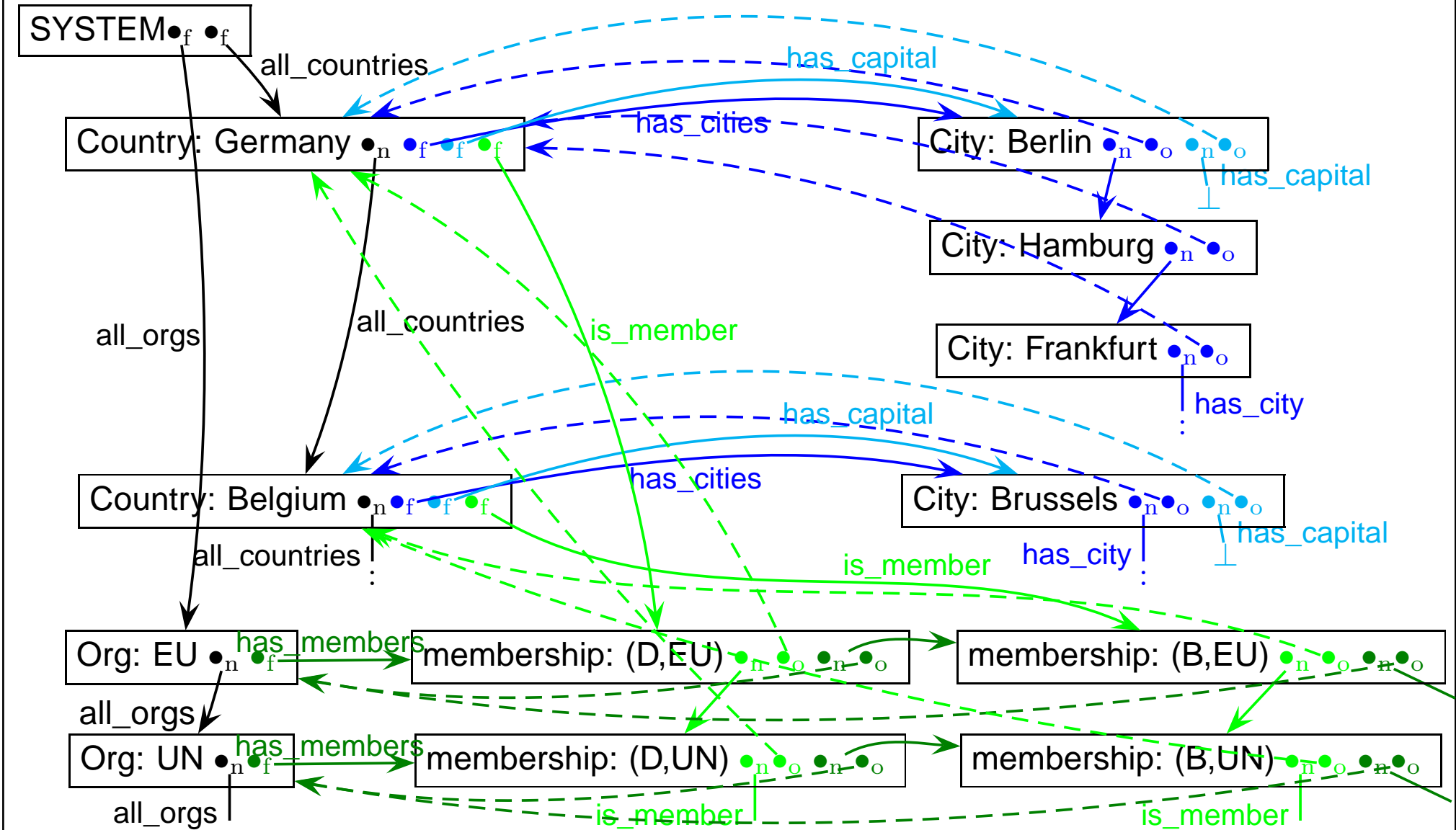
- a) Visualize the model by drawing some country, city and organization data records.
- b) Consider the “has_headq”-relationship that describes that organizations have their headquarters in a city.

□

SOLUTION: NETWORK SCHEMA DIAGRAM



SOLUTION: INSTANCE LEVEL



DATA DEFINITION LANGUAGE: EXAMPLE

RECORD NAME IS country

DUPLICATES ARE NOT ALLOWED FOR Code

Name TYPE IS CHARACTER 20

Code TYPE IS CHARACTER 4

Population TYPE IS NUMERIC INTEGER

Area TYPE IS NUMERIC INTEGER

RECORD NAME IS city

Name TYPE IS CHARACTER 25

Population TYPE IS NUMERIC INTEGER

SET NAME IS all_countries

OWNER IS SYSTEM

MEMBER IS country

SET NAME IS has_cities

OWNER IS country

MEMBER IS city

QUERY AND DATA MANIPULATION LANGUAGE

- record-at-a-time DML
- based on *iterators* (common design pattern/interface, e.g. in Java!) over sets
 - commands for navigation, access and data manipulation
 - embedded into a host language (COBOL, PL/I, later ... Pascal, C)
- “Current of” (cf. PL/SQL: “cursor”) that points to an instance of a record/set type in the DB
 - current of each record type
 - current of each set type (pointing on either the owner or one of the member records)
 - current of run unit (CRU): the record most recently accessed – any record type
- UWA (User Work Area) in the programming language runtime environment
 - one variable for each record type (auto-defined from the schema)
 - current of ... can be “fetched” into the corresponding UWA record

Retrieval and Navigation Commands

Query answering consists of stepwise navigation, carefully tracing currency indicators, and fetching tuples to the UWA:

- Retrieval: move the CRU into the corresponding UWA record,
- Navigation: navigate by using iterators and currency indicators to specific records and set owners/members.

Search for a Record of a Record Type

- FIND ANY <data record type> [USING <UWA.field.list>]
- FIND DUPLICATE <data record type> [USING <UWA.field.list>]
- tests/loops can be programmed by IF/WHILE DBSTATUS=0 // 0: successfully found
- FIND sets all current of record/set type in which the record participates to that record.
Can be avoided with RETAINING clause.

```
UWA.city.name = "Santiago";  
FIND ANY city USING name;  
// sets also current of city indicator  
while DBSTATUS=0 do begin  
    GET city    // fetches data record into UWA.city  
    if UWA.city.population > 1.000.000 then writeln (UWA.city.name|UWA.city.population);  
    FIND DUPLICATE city USING name;  
end;
```

- How to print out the city name and the country where it is located?
Needs the "owner" of the city wrt. "has_cities".

Search for a Record in a Set Type

- FIND (FIRST | NEXT | PRIOR | LAST) WITHIN <set type> [USING <UWA.field.list>]
- FIND OWNER WITHIN <set type>
- starts always from the current of this set (which is implicitly set when the CRU points to a suitable record type)

```
UWA.country.name = "Belgium";  
FIND ANY country USING name;  
FIND FIRST city WITHIN has_capital  
GET city    // fetches data record (Brussels) into UWA.city  
writeln (UWA.city.name);  
FIND OWNER WITHIN in_province  
GET province    // fetches data record (Brabant) into UWA.province  
writeln (UWA.province.name);
```

- Joins are only possible via navigation and loops in the host language.

Exercise 2.2

Write a program that outputs all organizations that have their headquarter in the capital of one of their member countries. Compare with the equivalent SQL query against Mondial. □

UPDATES

Updates on Data Records

STORE, ERASE, MODIFY (of the current data record)

Updates on Sets

CONNECT, DISCONNECT, RECONNECT (for the current data record wrt. a set)

HIERARCHICAL DATA MODEL

- In general very similar: parent-child-relationships define a tree structure; additionally, “virtual” parent-child-relationships.
- Systems: IMS (IBM & Rockwell International, 1969 for NASA Apollo), Adabas (Software AG, 1969), etc ...

SOLUTION

```
// not tested
find any organization // sets current of has_headq, current of has_members
while ok do
{ get organization // current organization into UWA
  find first headq_in within has_headq_in // auxiliary record hq(org,cty)
  find owner within is_headq_of // is a city
  find owner within has_capital // is a country
  if ok then // city is a capital
  { get country // UWA.country now holds this country
    found = 0;
    find first membership within has_members
      // starts from the organization
      // points to an auxiliary membership record m(org,c)
    while ok & not found do
    { find owner within is_member using code // UWA.country.code
      // check if the owner country is the same as in UWA
      if ok then { println(UWA.organization.name); found = 1;}
      find next membership within has_members
    }
  }
  find duplicate organization // next organization
}
```

THE SAME IN SQL

```
SELECT name
FROM organization org
WHERE (city,country) IN (SELECT capital, code
                        FROM country
                        WHERE code IN (SELECT country
                                     FROM is_member
                                     WHERE organization = org.abbreviation))
```

```
SELECT organization.name
FROM organization, is_member, country
WHERE organization.abbreviation = is_member.organization
   AND is_member.country = country.code
   AND organization.city = country.capital
   AND organization.country = country.code
```

```
SELECT organization.name
FROM organization, country
WHERE organization.city = country.capital
   AND organization.country = country.code
   AND (abbreviation, code) IN (SELECT organization, country
                                FROM is_member)
```

CONCLUSION

- importance decreased rapidly since SQL came up (1979), in the meantime it is only present in “legacy systems”.
- no underlying theory (required as a base for normalization and optimization)
- only **procedural**, **(data-model-level) navigation-** and **record-oriented query language**, non-declarative, needs to be embedded into a host language (COBOL, PL/I, Pascal, C).
- not possible to state ad-hoc queries.
Error-prone due to behavior of currency indicators.
- nevertheless, the idea of **navigation** and **parent-child-relationships** between data records is elegant (no problems with referential integrity).
These concepts came up again in later approaches ... with high-level navigation!
- **graph data model**, “**node + edge-labeled**”
- especially, ordered “child data records” are used again in XML. Then, there is
 - the DOM as an abstract datatype (stepwise, record-oriented),
 - XPath/XQuery as a *declarative*, set-oriented high-level language.

2.2 Object-Oriented Databases

Mid-80s: Object-orientation

- object-oriented design and modeling (UML)
- object-oriented programming (C++)

Application programs are developed and programmed in an object-oriented way.

- “impedance mismatch” between **tuple-based** SQL databases and the **object-oriented** data structures of the programming languages.

Goals:

- make objects of the application programs persistent
- bring **object-orientation into the DBMS**
 - class hierarchy and inheritance, polymorphism
 - implementation and encapsulation of behavior

FURTHER INFLUENCES

- Networks: Internet and Intranets
 - **Data exchange** and **interoperability**
 - CORBA (1989) “Common Object Request Broker Architecture” (standardized by OMG – Object Management Group; predecessor of Web Services):
 - central ORB bus where services can connect
 - service registry (predecessor of WSDL and UDDI ideas)
 - description of service interfaces in object-oriented style (IDL - interface description language, similar to C++ declarations)
 - exchanging objects between services
- ⇒ requires a format for exchanging data (DB: between databases)

In this lecture, OODBS are only discussed shortly to sketch the central ideas. An extended lecture can be found in “Information Systems”, available at <http://user.informatik.uni-goettingen.de/~may/Lectures>.

LIFETIME OF OBJECTS

- Object-oriented programming language: Objects are created during runtime of an application program, and they are destroyed when the program terminates.

Objects in OO Database Systems

- persistent: objects that are created by an activity, and then they are stored in the database system and survive also the termination of the activity that created it (until they are explicitly destroyed by another activity)
- transient: objects that are only needed temporarily for executing an activity. They exist only as long as the application is actually active, and they are only managed by the runtime environment of the programming language.

Lifetime of Objects

- Relational DBMS: all SQL types have only persistent instances that are stored in the DBMS. All non-SQL types (i.e., types of the host language) have only transient instances, these are destroyed with the termination of the application-program (= when the host language is left).

Persistent objects can only be manipulated/used by SQL, while transient objects can only be manipulated/used by the host language.

⇒ “impedance mismatch”.

- ODBMS: object types of the DBMS and of the application coincide. They can have parallel and transient instances at the same time.

For persistent and transient objects the same programming language and the same operations are used.

- **comparison with XML**: XML nodes can also be processed uniformly in the runtime environment and stored in a database. The DOM-API can be used in both cases.

OBJECT-ORIENTED DATA MODEL

- describes only the (*database*) state (*attributes, relationships, class membership and class hierarchy*), not the behavior,
- representation of the current state of the application-domain,
- corresponding conceptual modeling language: UML (see Software Engineering)
- **more expressive than the relational model/ER-model**
- (behavior of objects is integrated into the data manipulation language)

OO-DBMS

Standardization activities similar to the standardization of relational databases:

Success of the relational database systems:

- not only by the simple, high-level data model,
- but also due to the standardization: SQL (at least after some time)
 - portability
 - interoperability

ODMG: Object Database Management Group

- founded 1991
- Architecture of OODBMS, DDL, query language (OQL), data formats
- ODMG-1.0 standard (1993)
- ODMG-2.0 standard (1997)
- ODMG-3.0 standard (2000); incremental changes

Literature: Cattell et al; Object Database Management (ODMG, 1993/1997/2001)

ODMG: OBJECT DATABASE MANAGEMENT GROUP

- Voting members: organizations/companies, who commercially work at an ODBMS, among others JavaSoft, Windward Solutions, Lucent Technologies, Unidata, GemStone, ObjectDesign, Versant, ...
Reviewer members: Organizations who have a material interest in the work of ODMG.
- not the goal to define identical products, but to obtain source code portability (cf. Java, SQL, later also XML).
- enough freedom to define own properties and targets of products:
 - performance, optimization, (price)
 - support of certain programming languages,
 - functionality dedicated to special application areas (multimedia, CAD, ...), predefined types
 - integrated programming environments, design tools ...

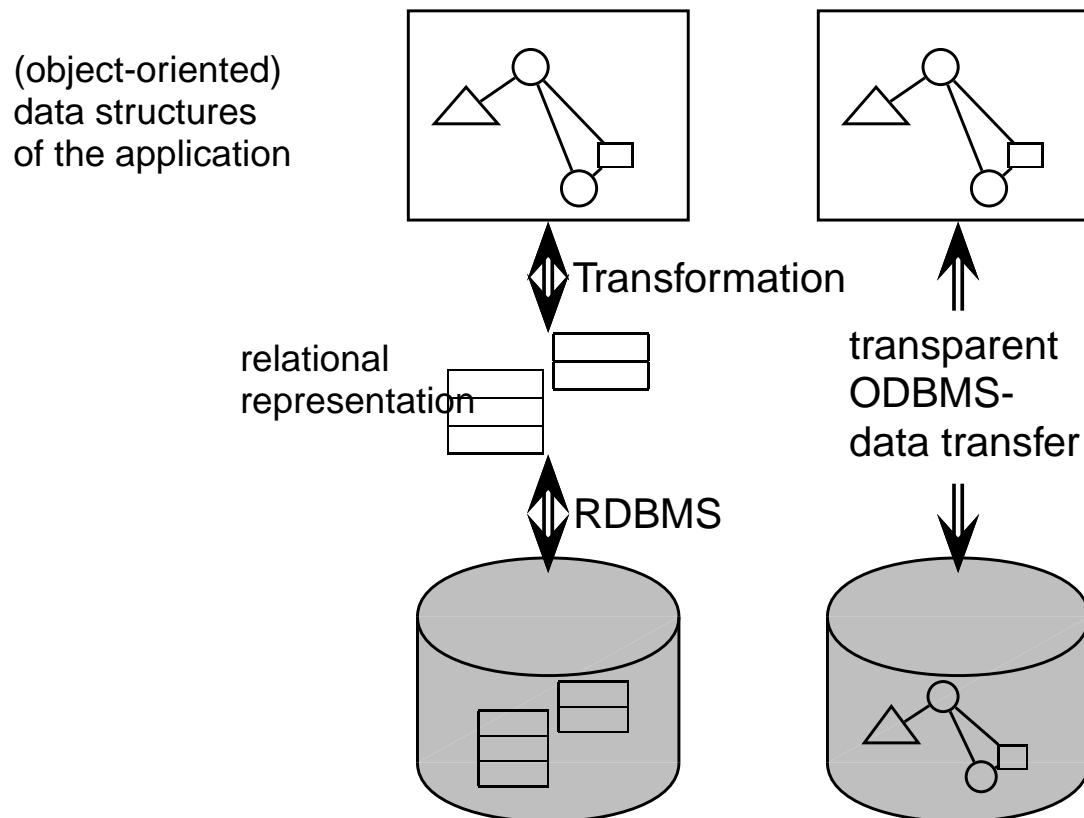
ARCHITECTURE OF ODBMS

- Different from “classical” relational DBMS:
SQL: high-level language for data manipulation,
applications are then written in other programming languages (cf. embedded approaches).
- ODBMS/ODM: transparent integration of DBMS functionality (persistence, multiuser, recovery) into application programming language (cf. Persistent Java).
The objects of the application are simply stored in the database.
- no separate DML necessary. The application-level programming language *is* the DML.
- There is also a *set-oriented, declarative query language*
(the impedance mismatch between variable-orientation and set-orientation remains):
OQL
- no transformation between the (logical) database representation and the representation in the programming language (cf. datatype conversion in JDBC).

ARCHITECTURES

ODMG is concerned with two types of products:

- Object Database Management Systems (ODBMSs) store the objects directly,
- Object-to-Database Mappings (ODMs) convert objects and store them in a relational (or any other) representation.



Remark:

There are similar approaches for XML databases.

ODMG-STANDARD

A standard that consists of several languages for implementation-level specification of object-oriented systems.

COMPONENTS OF THE ODMG STANDARD

- Object specification languages/data model
 - Object Definition Language (ODL)
 - Object Interchange Format (OIF)
- Object Query Language (OQL) – based on SQL
- C++/Smalltalk/Java Language Binding
 - specifies how to work with persistent objects in the target languages.

2.2.1 ODL: Object Definition Language

- Data definition language for object types:
- not a programming language, but only a language for definition of object specifications,
- characterizes object types (class hierarchy, properties and relationships)
- extends IDL (Interface Definition Language) from the OMG/CORBA (1989/1990) standard (which is in course closely related to the declaration commands in Java)

DATA TYPES: LITERALS

Literals are only values, they have no *object identity*.

Atomic literals

- long, short, unsigned long/short, float, boolean, char, string,
- enumeration {...} (“type generator”)
Z.B. `enum Weekday {Sunday, Monday, . . . , Saturday}`

Structured Literals

- predefined types: `date`, `interval`, `time`, `timestamp`
(additionally to *actual object types* Date, Interval, Time, Timestamp)
- user-defined structural types, e.g. `address` or

```
struct geoCoord { real    longitude;  
                  real    latitude; }
```

Collection literals

- `set<t>`, `bag<t>`, `list<t>`, `array<t>`, `dictionary<t>`
(additionally to the *actual collection types* Set, Bag, List, Array, Dictionary)

CLASSES

... are used to define and categorize complex object types.

Classes define the *signature* of their instances (the implementation does not belong to the object model):

```
class <name> { <attribute-defs>;  
    <relationship-defs>;  
    <operation-defs>;  
}
```

<attribute-def> ::= attribute <domain-type> <attribute-name>

```
class City { attribute string name;    % attributes ...  
    attribute number population;  
    attribute geoCoord coordinates;  
    relationship Country in_country;  % ... and relationships  
}
```

RELATIONSHIPS

- relationships are defined in course of the definition of classes.
- in UML and ODMG, only *binary* relationships are allowed.
- *bidirectional* and *inverse* relationships can be specified. Inverse relationships exist in UML, and later again in the Semantic Web languages (OWL).
- one-to-one / one-to-many / many-to-many-relationships.

```
class <name> {  
    <attribute-defs>;  
    <relationship-defs>;  
    <operation-defs>; }
```

```
<relationship-def> ::= relationship <target_of_path> <relationship-name>  
                        inverse <domain-type> :: <relationship-name'>
```

```
<target_of_path> ::= <domain-type> |  
                    <collection type> <<domain-type>>
```

- <collection type> for -to-many-relationships

RELATIONSHIPS

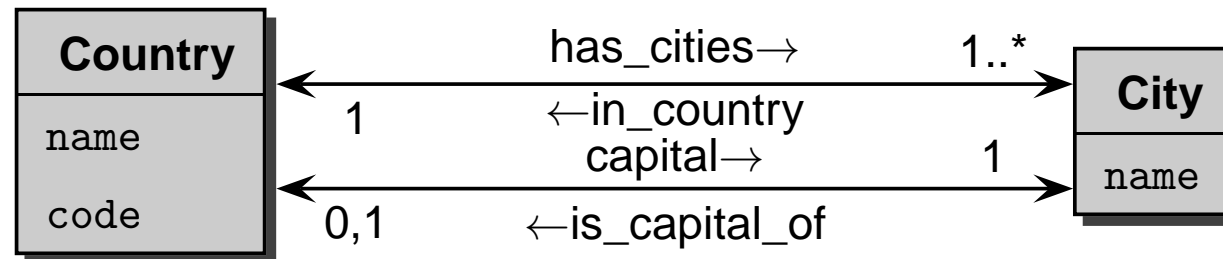
```
class <name> { <attribute-defs>;
    <relationship-defs>;
    <operation-defs>; }
```

```
<relationship-def> ::= relationship <target_of_path> <relationship-name>
    inverse <domain-type> :: <relationship-name>
```

```
<target_of_path> ::= <domain-type> |
    <collection type> <<domain-type>>
```

```
class Country { attribute string name;
    relationship City capital inverse City::is_capital_of;
    relationship set<City> has_cities inverse City::in_country; }
```

```
class City { attribute string name;}
```



RELATIONSHIPS

- the instance level can be represented as a graph:
 - nodes: objects; nodes have labels (names of the object types) and an ID
 - edges: relationships; edges have labels (names of relationships)
- ODBMS is responsible for maintaining referential integrity:
If an object is deleted, all relationships with/to it must also be deleted.
- relationships define *access paths*, e.g. `Germany.capital` for *navigation* through the graph.
- *graph-data model*, “node + edge-labeled”
- The `set<...>` is very similar to the *set-oriented* representation of set-valued relationships from the network data model (→ handled by iterators)
- the query language OQL solves this problem SQL-like in a *declarative* way (see later).

Exercise 2.3

Visualize an excerpt of the Mondial database as an object graph. □

2.2.2 Object Interchange Format (OIF)

- dump the database state to one or more files (cf. *export* in ORACLE)
- specification language for persistent objects and their states
- OIF file contains for each object its type, its attributes and values, and its relationships to other objects;
- the database schema (class definitions *and* class hierarchy) is *not* represented in OIF!

OBJECT INTERCHANGE FORMAT (OIF)

- Simplest form: only the class membership

```
<object> <class> {}  
Germany Country {}  
Berlin City {}
```

- attribute values are enumerated in braces:

```
Germany Country {name "Germany", area 356910, ...}
```

- structured attributes: nested brace structures

```
struct geoCoord { real longitude;  
                  real latitude; }
```

```
class City { attribute string name;  
             attribute geoCoord coordinates;  
             relationship Country in_country; }
```

```
Berlin City {name "Berlin", in_country Germany,  
             coordinates {longitude 13.3, latitude 52.45} }
```

OBJECT INTERCHANGE FORMAT (OIF)

- Collections, set-valued relationships:

```
class Country {  
    attribute string name;  
    relationship set<City> has_cities;}
```

```
Germany Country  
    {name "Germany", capital Berlin,  
    has_cities {Berlin, Frankfurt, Freiburg, ...} }
```

- cyclic references: no problem.
- attributed relationships (e.g. border) cannot be represented directly

⇒ OIF is already a self-describing data format!

2.2.3 OQL (overview)

- Query language of the ODMG standards (Object Query Language)
- similar to SQL:
SELECT - FROM - WHERE - clause, extended by complex objects, object-identity, path expressions, polymorphism, operation calls and late binding.
- but: functional language (like SQL), fully orthogonal (in SQL not completely)
- no explicit UPDATE statement: instead, object methods are used
- not Turing complete (cf. SQL/transitive closure)
- OQL can be embedded into suitable object-oriented programming languages (C++, Java, Smalltalk). Results of queries (collections!) are then processed by iterators.

EXTENTS

SQL: SELECT ... FROM <relation> ...

What corresponds to a *relation* in an ODBMS ?

⇒ *Extension*: set of all instances of a class (similar to system-owned sets in NWDBMS).

Extensions are defined in ODL together with the class declaration:

```
class <name> (extent <extent_name>)
```

```
{ <attribute-def>;  
  <relationship-def>;  
  <operations-def>; }
```

```
class Country (extent Countries)
```

```
{ attribute string name;  
  relationship City capital;  
  set<string> languages;  
  ... }
```

QUERIES

Queries against the database are expressed with the SELECT statement, with the same simple basic structure as in SQL:

```
SELECT <expression>  
FROM <extents>  
WHERE ...
```

```
SELECT c.population  
FROM Countries c  
WHERE c.name = 'Germany'
```

- with an iterator variable (here: c) – cf. SQL Aliasing

Similar to SQL:

- DISTINCT, aggregate functions: COUNT, SUM, . . . , set functions: UNION, INTERSECT, EXCEPT (MINUS)

QUERIES

- SQL: all results of queries of the form
 SELECT a,b,c FROM ...
are virtual “relations” (i.e. sets of tuples),
- OQL: the result is a virtual set of objects,
- in most cases an (implicit) collection.

```
SELECT c.capital  
FROM Countries c
```

Result is of the type

```
collection <City>
```

⇒ queries can be nested arbitrarily (like in SQL)

QUERIES

in case that the result has more than one attribute (e.g. with SELECT *), a

bag <struct{...}>

is automatically generated:

```
SELECT c.name, c.population  
FROM Countries c  
WHERE c.name = 'Germany'
```

Result is of the type

bag <struct {string name; number population}>

COMPLEX RESULTS

- bags (here: set-valued relationship) can be handled as a whole,
- by **explicit** generation of a struct, the properties of the result can be renamed:

```
SELECT struct(name: c.name,  
             cities: c.has_cities)  
FROM Countries c
```

result is of the type

```
collection <struct {string name;  
                  collection<city> cities}>
```

How can something *in the collection* be selected?

... straightforwardly: apply a SELECT statement to the collection:

```
SELECT struct(name: c.name,  
             cities: (SELECT ct  
                     FROM c.has_cities as ct  
                     WHERE ct.population > 1000000))  
FROM Countries c
```

- Traversing the relationship *has_cities* by a *path expression* in the query
- nested SELECT in the SELECT statement: the inner SELECT ranges over the (virtual) set *c.has_cities* of instances of type *set<city>*.
- the inner SELECT is evaluated separately for every result (i.e. for each instance *c*) of the outer SELECT.

PATH EXPRESSIONS

for navigation along *scalar* relationships:

```
SELECT name: c.name,  
       cpp: c.capital.province.population  
FROM Countries c
```

SELECT IN THE FROM-CLAUSE

Navigation along *set-valued* relationships:

not allowed:

```
SELECT name: c.has_cities.name,  
       pop: c.has_cities.population  
FROM Countries c  
WHERE c.name = "Germany"
```

`has_cities` is a *set of cities*, thus, the method *population* cannot be applied (to the set).

This can be done e.g. by a SELECT statement in the FROM-clause:

```
SELECT name: cty.name,  
       pop: cty.population  
FROM (SELECT c.has_cities  
      FROM Countries c  
      WHERE c.name = "Germany") as cty
```

CORRELATED JOINS

... do the above example even better:

```
SELECT name: cty.name,  
       pop: cty.population  
FROM Countries c, c.has_cities cty  
WHERE c.name = "Germany"
```

This would be a nice feature also in SQL ... the right side of the join is computed dependent on the left one.

⇒ asymmetric joins that express nested iteration in a declarative way

⇒ not aligned with the relational algebra

OQL: FUNCTIONAL LANGUAGE CONCEPT

SQL:

- declarative, relational algebra as theoretical base,
- somewhat ad-hoc language (around SELECT – FROM – WHERE),
- not completely orthogonal composition (aggregate functions, method applications)

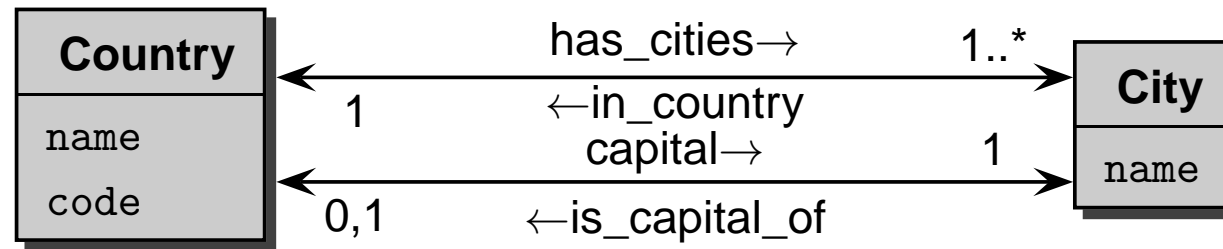
OQL:

- orthogonal composition rules: operators can be nested as long as the type system is not violated
- functional concept, includes the *simple* queries in SQL syntax.
- result of a query is always a
collection()
- can be processed in the same way as an extension (*intensional* part of the database).

CONCLUSION

- Object-oriented databases have not been accepted by the market.
- Products: ObjectStore, Adabas, O2, GemStone, Poet, ...
Some of them served as the base for the first commercial XML database systems (Excelon, Tamino [Software AG]).
- Object-relational extensions to SQL and relational systems (SQL-3-Standard): evolutionary instead of revolutionary development.
- **graph data model**, “node + edge-labeled”
- **set-oriented** (extends similar to relations) and **navigation-based** access, integrated in a **declarative** language.
Problems with navigating along set-valued properties.
- OQL as a **functional language** with **fully orthogonal constructs** and the possibility to **generate structures in the SELECT-clause**.
The XML-Query language XQuery will be very similar ...
- OIF as **self-describing ASCII-based data exchange format**, but still with a **fixed schema**.

2.2.4 Analysis: 1:n-Relationships

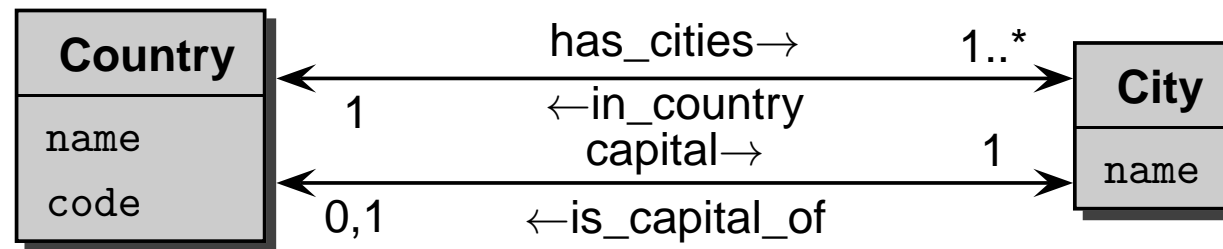


```
class Country { attribute string name;
                relationship City capital inverse City::is_capital_of;
                relationship set<City> has_cities inverse City::in_country; }
```

```
class City { attribute string name;}
```

- correct: `germany.capital.name`
- not correct: `germany.has_cities.name`
- translation to `set<City>` “country is in relation with a set of cities” is a tribute to **programming language influence**: must be something that exists in programming languages and that can be bound to a single variable.
“set-valued” – one answer which is a set.
- applying “.name” to a set is obviously not correct.

ALTERNATIVE TRANSLATION

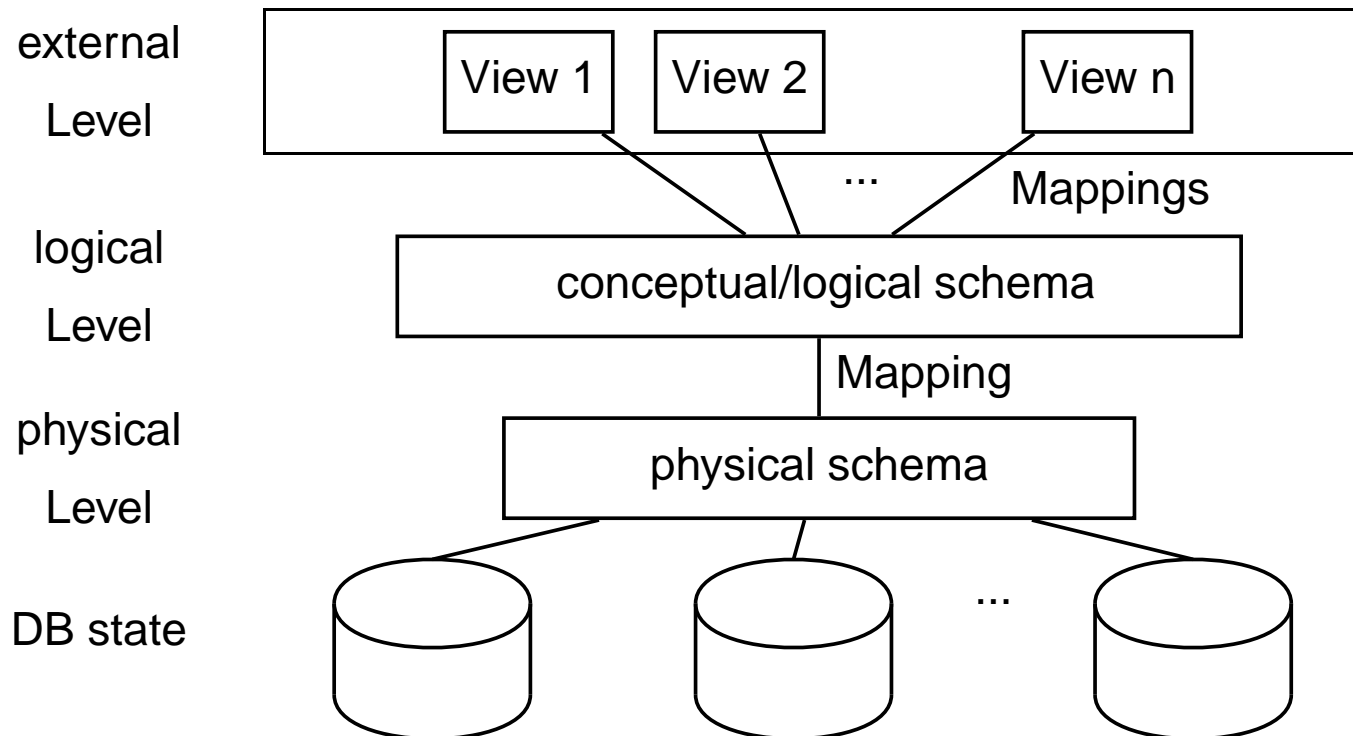


- database style: “country is in relation with multiple cities”
“**multi-valued**” – a set of *answers*, each of them is a city,
- “**set of answers**” is a meta-concept of the query language, not of the underlying programming language,
- applying “.name” to a set of answers can be defined by the semantics of the query language!
- “Modern” query languages change to multivalued semantics:
 - F-Logic (1989, see later): [germany.has_cities.name](#),
 - XPath (for XML, 1998): [//country\[name=“Germany”\]/city/name](#),
 - semantics of path expressions stays within the semantics of the query language.

2.3 Data Integration and Metadata Queries: SchemaSQL

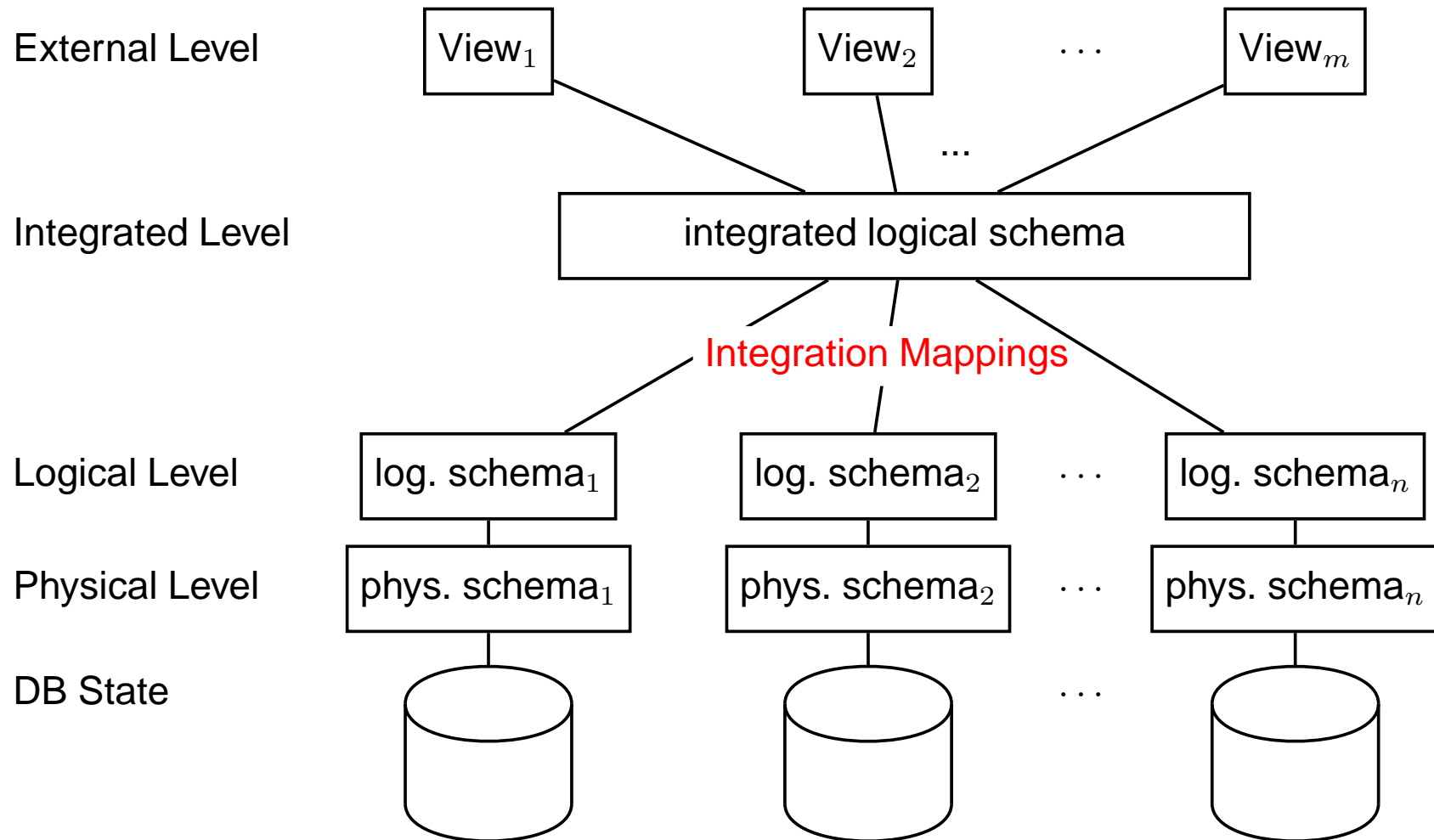
2.3.1 Introduction

- So far: single databases
- according to the classical 3-level architecture



MULTIDATABASE SYSTEMS AND FEDERATED DATABASE SYSTEMS

- providing a common, integrated view over several databases



DATA INTEGRATION AND METADATA QUERIES IN SQL: SCHEMA SQL

SchemaSQL (Lakshmanan et al. 1996; non-commercial academic system) extends SQL:

- combination of relations and attributes of different (federated) databases.
- uniform handling of data and metadata (by SchemaSQL variables).
- possible domains of variables are the names of the components of a federation, names of the relations of a database, names of the attributes of a relation, tuples of a relation, and values of a column of a relation.
- additionally to the “vertical” aggregations over columns, also “horizontal” aggregations over relations or even tables are possible.

Example

univ-A

salInfo

category	dept	salary
Prof	CS	65,000
Assoc Prof	CS	50,000
Prof	Math	60,000
Assoc Prof	Math	55,000

univ-B

salInfo

category	CS	Math
Prof	55,000	65,000
Assoc Prof	50,000	55,000

univ-C

CS

category	salary
Prof	60,000
Assoc Prof	55,000

Math

category	salary
Prof	70,000
Assoc Prof	60,000

univ-D

salInfo

dept	Prof	Assoc Prof
CS	75,000	60,000
Math	60,000	45,000

2.3.2 Declaration of Variables

... as known from SQL in the FROM-clause: FROM <range> <var>

SQL: FROM <table> <var>

SELECT city.name, city.population

FROM City city

$\langle \text{range} \rangle \in \{ \rightarrow, db \rightarrow, db :: rel, db :: rel \rightarrow, db :: rel.attr \}$.

- \rightarrow : names of the databases of the federation
- $db \rightarrow$: names of the relations of the database db .
- $db :: rel$: tuples of the relation rel of the database db [as in SQL].
- $db :: rel \rightarrow$: names of the attributes of the schema of the relation rel of the database db .
- $db :: rel.attr$: values of the column of the attribute $attr$ of the relation rel of the database db .
- **SchemaSQL: iterated declarations of variables are allowed!**
 \Rightarrow joins not longer symmetrical (cf. OQL).

Declaration of Variables: Tuple- and Domain Variables

- **tuple variables** as known from SQL:

$db :: rel$ ranges over the set of tuples of the relation rel of the database db .

```
SELECT tuple.category, tuple.salary  
FROM univ-C::CS tuple
```

category	salary
"Prof"	60000
"AssocProf"	55000

- **Domain-Variables:**

$db :: rel.attr$ ranges over the set of values of the attribute $attr$ of the relation rel of the database db

```
SELECT cat  
FROM univ-A::salInfo.category cat
```

cat
"Prof"
"AssocProf"

Note: SQL-style `SELECT category FROM univ-A::salInfo` yields the same result – but does not allow to bind the values to a variable (that can be used somewhere else)

Declaration of Variables: Metadata Variables

- $db \rightarrow$ ranges over the relation names of the database db .

SELECT relname FROM univ-C \rightarrow relname

relname
"CS"
"math"

- Nested declarations: Second variable depends on the first one:

SELECT dept, tuple.category, tuple.salary
FROM univ-C \rightarrow dept, univ-C::dept tuple

dept	category	salary
"CS"	"Prof"	60000
"CS"	"AssocProf"	55000
"Math"	"Prof"	70000
"Math"	"AssocProf"	60000

... *integrates* both tables from univ-C in one.

Declaration of Variables

- Variables over names of attributes:

$db :: rel \rightarrow$ ranges over the set of attribute names of the schema of the relation rel of the database db .

SELECT attrname
FROM univ-C::CS \rightarrow attrname

attrname
“category”
“Salary”

- SELECT C: *name* of the attribute,
SELECT T.C: *value* of the respective attribute of the current tuple.

SELECT attrname, univ-C::CS.attrname
FROM univ-C::CS \rightarrow attrname

“category”	“Prof”
“category”	“AssocProf”
“Salary”	60000
“Salary”	55000

Declaration of Variables

- \rightarrow ranges over the names of the databases of the federation.

SELECT dbname FROM \rightarrow dbname

dbname
"univ-a"
"univ-b"
"univ-c"
"univ-d"

- SELECT dbname, relname

FROM \rightarrow dbname, dbname \rightarrow relname

dbname	relname
"univ-A"	"SallInfo"
"univ-B"	"SallInfo"
"univ-C"	"CS"
"univ-C"	"math"
"univ-D"	"SallInfo"

2.3.3 Queries

All departments of Univ-A that pay a higher salary to their professors than the corresponding departments of Univ-B:

```
select A.dept
  - all variables are independent
from   univ-A::salInfo A, univ-B::salInfo B,
       univ-B::SalInfo-> AttB
where  AttB <> "category" and
       A.dept = AttB and
       A.category = "Prof" and
       B.category = "Prof" and
       A.salary > B.AttB.
```

Queries (Cont'd)

Same for C/D:

```
select RelC
  - C depends on RelC
from   univ-C-> RelC, univ-C::RelC C,
       univ-D::salInfo D
where  RelC = D.dept and
       C.category = "Prof" and
       C.salary > D.Prof
```

AGGREGATION

Similar to SQL, there can be aggregation over a variable.

⇒ here also *horizontal* and *blockwise* aggregation possible.

Average salary for each kind of professors over all departments of Univ-B:

```
select T.category, avg(T.D)
from univ-B::salInfo→D, univ-B::salInfo T
where D <> "category"
group by T.category
```

- select the values for D,
- compute the cartesian product with univ-B::salInfo T
- include column T.D
- evaluate, do the grouping, compute the aggregate

D	category	CS	Math	T.D
category	Prof	55,000	65,000	55,000
CS	Prof	55,000	65,000	55,000
math	Prof	55,000	65,000	65,000
category	Assoc Prof	50,000	55,000	50,000
CS	Assoc Prof	50,000	55,000	50,000
math	Assoc Prof	50,000	55,000	55,000

Aggregation

Average salary for each kind of professors over all departments of Univ-C:

```
select T.category, avg(T.salary)
from univ-C→D, univ-C::D T
group by T.category
```

- compute values for D,
- join with tuple variable D T

D	category	salary
CS	Prof	60,000
CS	Assoc Prof	55,000
math	Prof	70,000
math	Assoc Prof	60,000

- grouping
- compute the aggregate

RESTRUCTURING

... as usual via views:

```
create view
BtoA::salInfo(category, dept, salary) as
  select T.category, D, T.D
  from univ-B::salInfo→D, univ-B::salInfo T
  where D <> 'category'
```

creates a virtual database BtoA with a virtual relation salInfo in the same format as A::salInfo.

Restructuring

A to B: number of attributes of the result table depends on the number of departments.

⇒ **Dynamic result schema**

```
create view AtoB::salInfo(category,D) as
select A.category, A.salary
from univ-A::salInfo A, A.dept D
```

Result of the FROM-clause:

A.category	A.salary	A.dept D
Prof	65,000	CS
Assoc Prof	50,000	CS
Prof	60,000	Math
Assoc Prof	55,000	Math

Many-to-one-mapping into a schema of the form
salInfo(category, dept₁, ..., dept_n).

AtoB::salInfo		
category	CS	Math
Prof	65,000	60,000
Assoc Prof	50,000	55,000

2.3.4 Exercise

Create the following view that represents the information of all four databases in a uniform way:

```
create view
globalSchema::salInfo(univ, dept, category, salary) as
[TO BE COMPLETED]
```

SOLUTION

```
create view
globalSchema::salInfo(univ, dept, category, salary) as
  select "univ-A", T.dept, T.category, T.salary
  from univ-A::salInfo T
union
  select "univ-B", D, T.category, T.D
  from univ-B::salInfo T, univ-B::salInfo→D
  where D<>"category"
union
  select "univ-C", T, T.category, T.salary
  from univ-C→D, univ-C::D T
union
  select "univ-D", T.dept, C, T.D
  from univ-D::salInfo T, univ-D::salInfo→C
  where C<>"dept"
```

2.3.5 Query Evaluation

Federation System Table (FST): meta-information about the component databases, i.e. names of the databases, relations, attributes, or other statistical information that is useful for query evaluation (similar to the Data Dictionary in SQL).

Variable Instantiation Tables (VIT): contain the possible variable bindings during the evaluation (meta level).

Input: a *SchemaSQL* query

Output: bindings of the variables of the SELECT-clause of the query

Evaluation: two phases:

1. generation of the VITs according to the variables in the FROM-clause. For this, SQL queries are stated against the local databases and against the FST.
2. rewriting of the *SchemaSQL* query into an equivalent query using the VITs (Dynamic SQL). This query is then evaluated by the *resident SQL server*.

EVALUATION: EXAMPLE

```

select RelC
from univ-C → RelC, univ-C::RelC C, univ-D::salInfo D
where RelC = D.dept and C.category = "Prof" and C.salary > D.Prof
    
```

Bindings for meta-variables (query against an *FST*):

VIT_{RelC}
RelC
CS Math

Bindings for tuple variables (queries against component-DBS):

VIT_C (depends on Rel_C)		
RelC	category	salary
CS	Prof	60,000
CS	Assoc Prof	55,000
Math	Prof	70,000
Math	Assoc Prof	60,000

VIT_D		
Dept	Prof	AssocProf
CS	75,000	60,000
Math	60,000	45,000

Evaluation: Example

... again the query:

```
select RelC
from univ-C→ RelC, univ-C::RelC C,
     univ-D::salInfo D
where RelC = D.dept and
       C.category = "Prof" and
       C.salary > D.Prof
```

Query evaluation via standard SQL over the *VIT's*.

```
select VIT_RelC.RelC
from VIT_RelC, VIT_C, VIT_D
where VIT_C.RelC = VIT_RelC.RelC      % Correlation RelC, C
     and VIT_RelC.RelC = VIT_D.dept
     and VIT_C.category = "Prof"
     and VIT_C.salary > VIT_D.Prof
```

EXERCISE: SCHEMA-SQL

Describe the evaluation of the query given on Slide 76 with its FST and VITs.

Solution

VIT_{dbname}	$VIT_{relname}$	
dbname	dname	relname
univ-A	univ-A	salInfo
univ-B	univ-B	salInfo
univ-C	univ-C	CS
univ-D	univ-C	math
	univ-D	salInfo

```
SELECT  $VIT_{dbname}$ .dbname,  $VIT_{relname}$ .relname  
FROM  $VIT_{dbname}$ ,  $VIT_{relname}$   
WHERE  $VIT_{dbname}$ .dbname =  $VIT_{relname}$ .relname
```

2.3.6 Example: Integration of Stock Exchange Data

Frankfurt::Quota		
Date	Name	Price
3.3.93	sun	150
3.3.93	dc	151
3.3.93	b.u.	160
4.3.93	sun	153
4.3.93	dc	154
4.3.93	b.u.	163

Tokyo::Quota			
Date	sun	dc	fuji
3.3.93	150	151	140
4.3.93	153	154	140

Sydney::3.3.	
Name	Price
sun	150
dc	151
kiwi	130

Sydney::4.3.	
Name	Price
sun	153
dc	154
kiwi	135

New York::sun	
Date	Price
3.3.93	150
4.3.93	153

New York::dc	
Date	Price
3.3.93	151
4.3.93	154

New York::msoft	
Date	Price
3.3.93	148
4.3.93	74

Possible extension:
Euro vs. Dollar vs. Yen

EXERCISE: SCHEMA-SQL

- Formulate the “On which days had which stocks the price of 150 \$?” for the schemata given on Slide 91.
- In commercial database systems, the schema information is stored in the *Data Dictionary* (cf. the following excerpts of table definitions of the data dictionary):

```
SQL> desc sys.user_tables;
```

Name	Null?	Type
-----	-----	----
TABLE_NAME	NOT NULL	VARCHAR2(30)

```
SQL> desc sys.user_tab_columns;
```

Name	Null?	Type
-----	-----	----
TABLE_NAME	NOT NULL	VARCHAR2(30)
COLUMN_NAME	NOT NULL	VARCHAR2(30)
DATA_TYPE		VARCHAR2(30)

Describe how the above queries can be formulated in an environment where SQL is embedded into a procedural programming language (e.g. embedded-SQL or PL/SQL) (Pseudocode).

SOLUTION: SCHEMA-SQL

- SELECT Date, Name

```
FROM Frankfurt::Quota
```

```
WHERE Price=150;
```

```
SELECT Date, AttrName
```

```
FROM Tokyo::Quota.Date, Tokyo::Quota → AttrName
```

```
WHERE AttrName ≠ 'Date' AND Price=150;
```

```
SELECT NewYork::TabName.Date, TabName
```

```
FROM NewYork → TabName
```

```
WHERE Price=150;
```

```
SELECT TabName, Sydney::TabName.Name
```

```
FROM Sydney → TabName
```

```
WHERE Price=150;
```

- Information from the Data Dictionary is only needed for Tokyo, New York and Sydney.

SOLUTION: SQL

Algorithm for SQL in a procedural environment (database Tokyo):

- Store the result of

```
SELECT ColumnName
FROM Tokyo.user_tab_columns
WHERE ColumnName ≠ 'Date';
```

(result: the names of the companies) and for each result <cn> execute the query

```
SELECT Date, <cn>
FROM Tokyo.Quota
WHERE <cn>= 150;
```

and collect all results.

Solution: SQL

- database “New York”: store the result of

```
SELECT TableName
FROM user_tables
WHERE
  ( SELECT ColumnName
    FROM user_tab_columns UTC
    WHERE UTC.TableName=TableName = {Date,Price});
```

(the comparison of sets must be formulated in SQL) and for each result <tn> evaluate the query

```
SELECT Date, <tn>
FROM <tn>
WHERE Price = 150;
```

and collect all results.

Problem: SQL statements must be generated *dynamically*: the results of the first query are used in the second statement.

SOLUTION: DYNAMIC SQL

This is e.g. possible in Oracle by using the DBMS_SQL-Package (to be used with PL/SQL), which allows to generate SQL statements at runtime:

```
create procedure findnumber as
declare
  cursor col_cursor is
    select column_name, data_type
    from sys.user_tab_columns
    where table_name = upper('&&table_name')
    order by column_id;
  lv_column_name  sys.user_tab_columns.column_name%TYPE;
  lv_column_typ   sys.user_tab_columns.data_type%TYPE;
  lv_rowid        varchar2(20);
  rows_processed  number;
  loop_count      number;
  stmtnt          varchar2(2000);
  doublecur       BINARY_INTEGER;
  execute_feedback INTEGER;
  type colname_typ is table of lv_column_name%TYPE
    index by binary_integer;
  type rowid_typ is table of lv_rowid%TYPE
```

```

        index by binary_integer;
colname_table      colname_typ;
empty_colname      colname_typ;
rowid_table        rowid_typ;
empty_rowid_table  rowid_typ;

begin
  DBMS_OUTPUT.ENABLE(10000);
  rows_processed := 0;
  -- Search for attributes with datatype "Number"
  open col_cursor;
  loop
    fetch col_cursor into
      lv_column_name, lv_column_typ;
    exit when col_cursor%notfound;
    IF lv_column_typ='NUMBER' THEN
      rows_processed := rows_processed+1;
      colname_table (rows_processed)
        := lv_column_name;
    END IF;
  end loop;
  close col_cursor;

  -- Initialize query statement
  stmtnt := 'select rowid from '
    || '&&table_name '

```

```

    || 'where ' ;

-- generate the query iteratively
loop_count := 1;
WHILE loop_count <= rows_processed
loop
    stmtnt := stmtnt
        || colname_table(loop_count)
        || ' = &&Price';
    if loop_count < rows_processed
    then
        stmtnt := stmtnt || ' or ' ;
    end if;
    loop_count := loop_count + 1;
end loop;
DBMS_OUTPUT.PUT_LINE
    ('Computed Query: ' || stmtnt);

-- execute the generated statement
doublecur := DBMS_SQL.OPEN_CURSOR;
DBMS_SQL.PARSE (doublecur
    ,stmtnt
    ,DBMS_SQL.V7);
DBMS_SQL.DEFINE_COLUMN
    (doublecur, 1, lv_rowid, 20);
execute_feedback := DBMS_SQL.EXECUTE (doublecur);

```

```
-- generate list of all resulting data records and
-- RowIDs
loop
  if DBMS_SQL.FETCH_ROWS (doublecur) = 0
  then
    exit;
  else
    DBMS_SQL.COLUMN_VALUE (doublecur,1, lv_rowid);
    DBMS_OUTPUT.PUT_LINE('RowID: ' ||lv_rowid);
  end if;
end loop;

-- cleaning ...
DBMS_SQL.CLOSE_CURSOR (doublecur);
colname_table := empty_colname;
end;
/
```

Solution: Dynamic SQL

```
SQL> execute find-number;
```

```
Give value for table_name: Tokyo
```

```
Give a value for price: 150
```

```
Generated Query:
```

```
select rowid from Tokyo
where SUN = 150 or DC = 150 or FUJI = 150
```

```
RowID: AAAA2MAADAAAD7nAAA
```

```
SQL> select * from Tokyo
```

```
where rowid='AAAA2MAADAAAD7nAAA';
```

```
03.03.93          150          151          140
```

which must still be postprocessed for obtaining the answer 'sun', 3.3.93.

- Conclusion: SchemaSQL helps to express such queries much shorter and more concise, and it is easier to learn than PL/SQL and DBMS_SQL.

2.3.7 Exercise: Horizontal and blockwise Grouping

- Consider the schemata `univ-B`, `univ-C` and `univ-D`. Give SchemaSQL queries that return for each kind of professors the average salary over all departments.

SOLUTION: HORIZONTAL AND BLOCKWISE GROUPING

- univ-A: same as in standard SQL: vertical aggregation:

```
select T.category, avg(T.salary)
from univ-A::salInfo T      – tuple variable
group by T.category
```

- univ-B: horizontal aggregation
see Slide 81.
- univ-C: aggregation over different tables
see Slide 82.

- univ-D: aggregation over different columns:

```
select T.category, avg(T.C)
from univ-B::salInfo T,  univ-B::salInfo → C
where C <> "dept"
group by C
```

CONCLUSION

- integration of *relational* databases with different schemas
- queries against **metadata**
- **combination of metadata and data**
- data-dependent **generation of schema**

New Features

Generalization of the use of variables:

- SQL: variables only ranging over tuples of a fixed relation,
- SchemaSQL: variables ranging over “everything”: data: tuples, column values
metadata: names of columns, names of relations, even names of databases,
- intuitively simple extension of SQL,
- powerful feature for data integration,
– But: classical query optimization/evaluation not applicable.

Such variables are more (F-Logic) or less extremely (XML: XPath/XQuery) used in Semistructured Data and XML.

Chapter 3

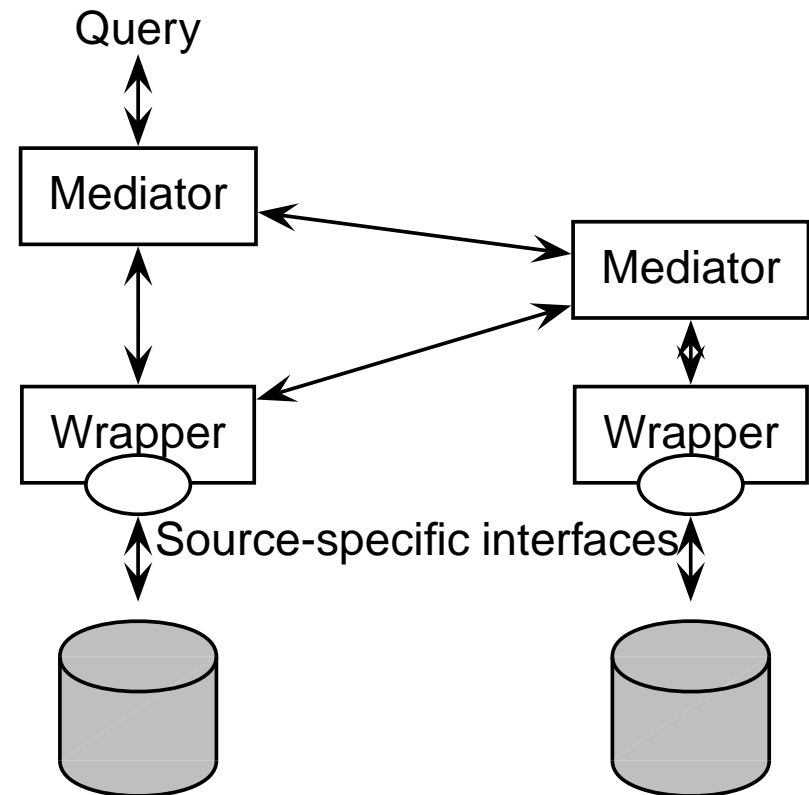
Semistructured Data: Early Approaches

- Data integration
 - different, autonomous data sources
 - *different data models and schemata*
 - more advanced than the approach of SchemaSQL
- Knowledge representation, data exchange
 - schema- and meta-information inside the data
 - examples: KIF (Knowledge Interchange Format), F-Logic
 - up to ontology management (“Semantic Web”)
- Management of data for presentation on the Web
- Extraction of data from the Web

SSD FOR DATA INTEGRATION/DATA EXCHANGE:

Wrapper/Mediator-Based Architectures

- *Mediator* (Vermittler): between users and data sources (Middleware),
- *Wrapper (Translator)*: provides homogeneous access to heterogeneous sources (especially for information extraction from the Web: programming of wrappers for Web pages and then collect the data)



WRAPPER/MEDIATOR-BASED ARCHITECTURES

- sources: databases, interfaces to databases via forms (e.g. library search), search engines, simple Web pages
- each relevant Web source is associated to a wrapper
- mediator contains knowledge about the accessible sources
- mediators can be composed hierarchically

Virtual Approach

The users state queries against the upper level mediator (“external view”) which translates the queries against lower mediators and wrappers. Wrappers answer the queries from the sources. Mediators combine the answers and return them.

Materialized Approach

An integrated view of all data is completely materialized (and maintained). Users state their queries against the materialized database that directly answers them.

REQUIREMENTS FOR DATA INTEGRATION

- upper mediator level: a target data format
- interfaces between wrappers/mediators
 - a common data exchange format
 - a common query language/mechanism
- wrapper level: mapping from sources into the common format

Target Data Model and Languages

- **flexible and extensible**
 - “copy all properties of object X from data source A”
 - extensible to additional sources
 - different source data models and schemata
- **handling metadata and content in combination**
- **self-describing data !?**

3.1 TSIMMIS

(The Stanford-IBM Manager of Multiple Information Sources, 1995-2000)

Persons: J. Ullman, H. Garcia-Molina, J. Widom, Y. Papakonstantinou, etc.

Goal (several subprojects): construction of means for a consistent and efficient integrated access to information sources:

- Heterogeneous information sources
 - databases
 - Web pages
- ⇒ often no explicit schema known/present
- ⇒ mapping to a *common data model*:
Object Exchange Model (OEM)

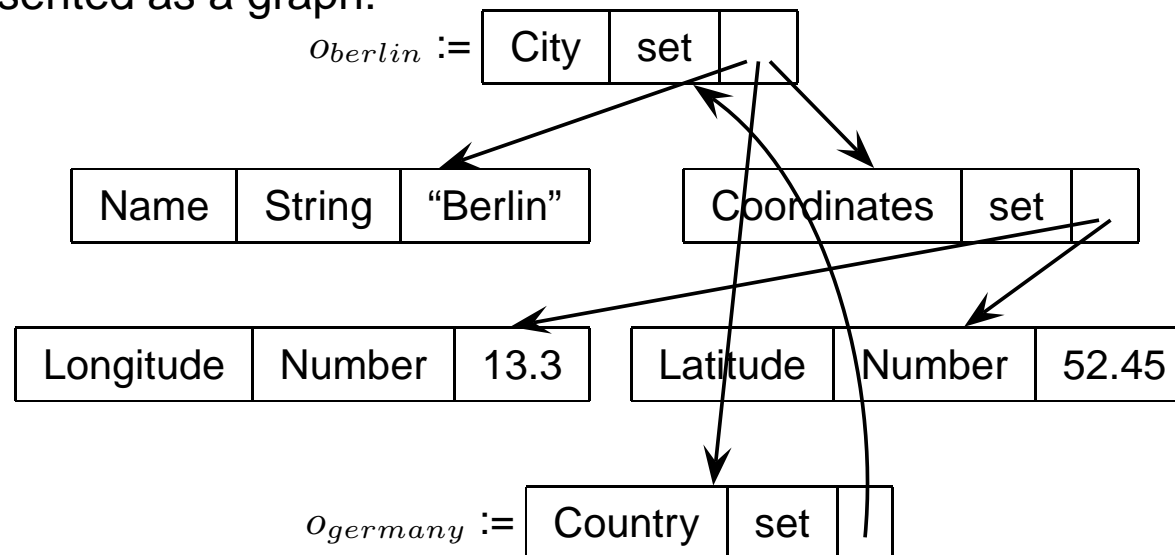
TSIMMIS: Concept

“Virtual” approach:

- users state queries against a mediator
- mediator forwards the subqueries to lower mediators or wrappers
- wrappers are programs that (logically) transform the objects of the data source into OEM and then answer the basic queries
- results of the wrappers are returned in OEM format to the mediator
- mediator integrates the results of the sources
- mediators can be composed hierarchically

3.1.1 OEM: Object Exchange Model

- very simple, “*self-describing*” object model
- knows only object identity and nesting as concepts:
- each object has an object-ID, a label (\sim class), a (data)type and a value,
- values of complex types are sets of references to sub-objects
- labels: “self-describing data”
- top-level objects with semantic object identifiers as entry points (cf. OQL)
- can be represented as a graph:



OEM: EXAMPLE

Source 1: CIA World Factbook

Wrapper cia exports OEM objects as follows:

<&cont1, continent, set, {&a1, &n1, &c1, &c2, &c3, ...}>

<&n1, name, string, 'Europe'>

<&a1, area, number, 9562488>

<&c1, country, set, {&cn1, &cc1, &ca1, &cp1, &cap1}>

<&c2, country, set, {&cn2, &cc2, &ca2, &cp2, &cap2}>

<&cn1, name, string, 'Germany'>

<&cc1, code, string, 'D'>

<&ca1, area, number, 356910>

<&cp1, population, number, 83536115>

<&cap1, name, string, 'Berlin'>

<&cn2, name, string, 'Sweden'>

<&cc2, code, string, 'S'>

<&ca2, area, number, 449964>

<&cp2, population, number, 8900954>

<&cap2, name, string, 'Stockholm'>

OEM: Example

Source 2: Global Statistics

Wrapper gs exports OEM objects as follows:

```
<&cont1, continent, set, {&a1, &n1, &c1, &c2, &c3, ...}>  
<&n1, name, string, 'Europe'>  
<&c1, country, set, {&cn1, &ct11, &ct12, &ct13, ..., &prov11, &prov21,...}>  
<&c2, country, set, {&cn2, &ct21, &ct22, &ct23, ..., &prov12, &prov22,...}>  
  
<&cn1, name, string, 'Germany'>  
<&ct11, city, set, {&ctn11, &ctp11, &prov11}>  
<&prov11, province, set, {&pn11, &pa11, &pp11, &ct11}>  
  
<&ctn11, name, string, 'Stuttgart'>  
<&ctp11, population, number, 588482>  
  
<&pn11, name, string, 'Baden-Württemberg'>  
<&pa11, area, number, 35742>  
<&pp11, population, number, 10272069>
```

OEM

- another version of OEM has been presented that additionally allows for labeled edges; e.g. for *capital*-edges from a country to a city.

Exercise 3.1

Visualize an excerpt of the Mondial database with some countries, cities, continents and organizations as an OEM graph. □

... a very simple data model.

- how to query it?
- generally, the network model language could be used for navigating ...
- ... but in the meantime, *declarative* languages had been invented:
 - clause-based: SQL-style
 - logic-based: Datalog-style

3.1.2 TSIMMIS: Languages

Mediators are programmed in **MSL** (Mediator Specification Language; a rule-based query language for OEM):

MSL rules (cf. Prolog, Datalog):

head(Vars) :- body(Vars,databases)

- head and body consist of expressions over *patterns* of the form

<oid label type value>

or

<oid label value>

or

<label value>

- *value* can be set-valued; in this case, the set consists itself of expressions of the form *<...>* .
- objects of different sources are identified by *<object>@source*.
- *body*: pattern that must be satisfied by suitable variable bindings,
- *head* describes the structure of the OEM object that is generated.

MSL

The country whose code is “D”:

```
?- <C country {<code "D">}>@cia
```

- the query generates a set-valued result object, whose sub-objects are the individual answers:

```
result: <answer {&c1}>
```

The names of all south-american countries in which there is a city with name “Santiago”:

```
<countryname N>:-  
    <continent {<name "South America"> <country  
        {<name N> <city {<name "Santiago">}}>}}>@gs
```

```
<&obj42, answer, set, {<countryname "Chile">,  
    <countryname "Paraguay">,  
    <countryname "Argentina">}>
```

EXAMPLE

Mediator med accesses wrappers cia and gs.

Query: all cities that are stored in gs whose names are mentioned in cia as names of capitals.

Mediator rule:

```
<capital {<name Cap> <country CN> R }>@med :-  
  <country {<name CN> <capital Cap>}>@cia  
  AND <country {<name CN> <city {<name Cap> | R}>}>@gs
```

- R is bound to the remaining sub-objects of the resulting city-objects.

⇒ object creation in the rule head *obj@med* (cf. Views)

exported object e.g. `<&cap, capital, set, {&ctn12, &c1, &ctp12, &ctprov12}>@med`

- additionally: external predicates (string concatenation, substring, comparisons etc).
- variables can be bound to labels and to values
- *syntactically* 2nd order

⇒ queries against the logical schema are possible.

3.1.3 TSIMMIS: User Queries Against OEM

Users can state queries in MSL or LOREL:

- MSL: see above – rule- and pattern-based language
- **LOREL (Lightweight Object Repository Language):**
(LORE: DBMS for OEM data model, Stanford)
clause- and path-based language (based on OQL)

Lorel

- Entry points are named constants (e.g. europe) or extents (e.g. countries)
- result of each query is a collection of OEM objects.

Semantics of 1:N-relationships

multi-valued instead of *set-valued* semantics:

- **germany.city** yields *multiple unary* answers instead of (as in ODMG) *one set-valued* answer.
- **germany.city.name** yields the names of all these cities.

LOREL

clause-based SQL/OQL-style language: SELECT - FROM - WHERE

```
% all europ. capitals:
```

```
select europe.country.capital % note: multivalued semantics
```

```
% the country with the code "S":
```

```
select c
```

```
from europe.country c
```

```
where c.code = "S"
```

```
% South-american countries such that ...
```

```
select c
```

```
from southamerica.country c
```

```
where c.city.name = "Santiago"
```

implicit existential semantics: ... if there is any city whose name is Santiago.

3.2 Frame-Based Models

- objects are represented by *object frames*,
- Frame contains *slots* for storing properties (“signature” of the slots *can* be given by a schema, but not necessarily (\Rightarrow semistructured data))
- Slots can be literal-valued or object-valued (*internal* storage by references) for describing attributes and relationships.

A SELF-DESCRIBING OO-DATA MODEL: F-LOGIC

(M. Kifer, G. Lausen, 1989/1995; SIGMOD Test of Time Award 1999)

- full object-orientation (class hierarchy, inheritance)
- objects have properties
- metadata (class names, method names) as first-class-members of the data model
- “frame-based” model
- deductive language (Prolog-style)

F-LOGIC: DATA MODEL AND SYNTAX

- is-a relationship:

$\langle \text{object} \rangle : \langle \text{class} \rangle$

- subclass-relationship:

$\langle \text{class} \rangle :: \langle \text{class} \rangle$

- properties:

$\langle \text{object} \rangle [\langle \text{property} \rangle \rightarrow \langle \text{object/value} \rangle]$ (scalar)

$\langle \text{object} \rangle [\langle \text{property} \rangle \rightarrow \{ \langle \text{set-of-objects/values} \rangle \}]$ (set-valued/multi-valued)

Analogously with parameters:

$\langle \text{object} \rangle [\langle \text{property} \rangle @ (\langle \text{list-of-objects/values} \rangle) \rightarrow \langle \text{object/value} \rangle]$

$\langle \text{object} \rangle [\langle \text{property} \rangle @ (\langle \text{list-of-objects/values} \rangle) \rightarrow \{ \langle \text{set-of-objects/values} \rangle \}]$

- inheritable properties:

$\langle \text{class} \rangle [\langle \text{property} \rangle \bullet \rightarrow \langle \text{object/value} \rangle]$ (scalar)

$\langle \text{class} \rangle [\langle \text{property} \rangle \bullet \rightarrow \{ \langle \text{set-of-objects/values} \rangle \}]$ (set-valued)

nonmonotonic inheritance semantics with overriding.

F-LOGIC: EXAMPLE

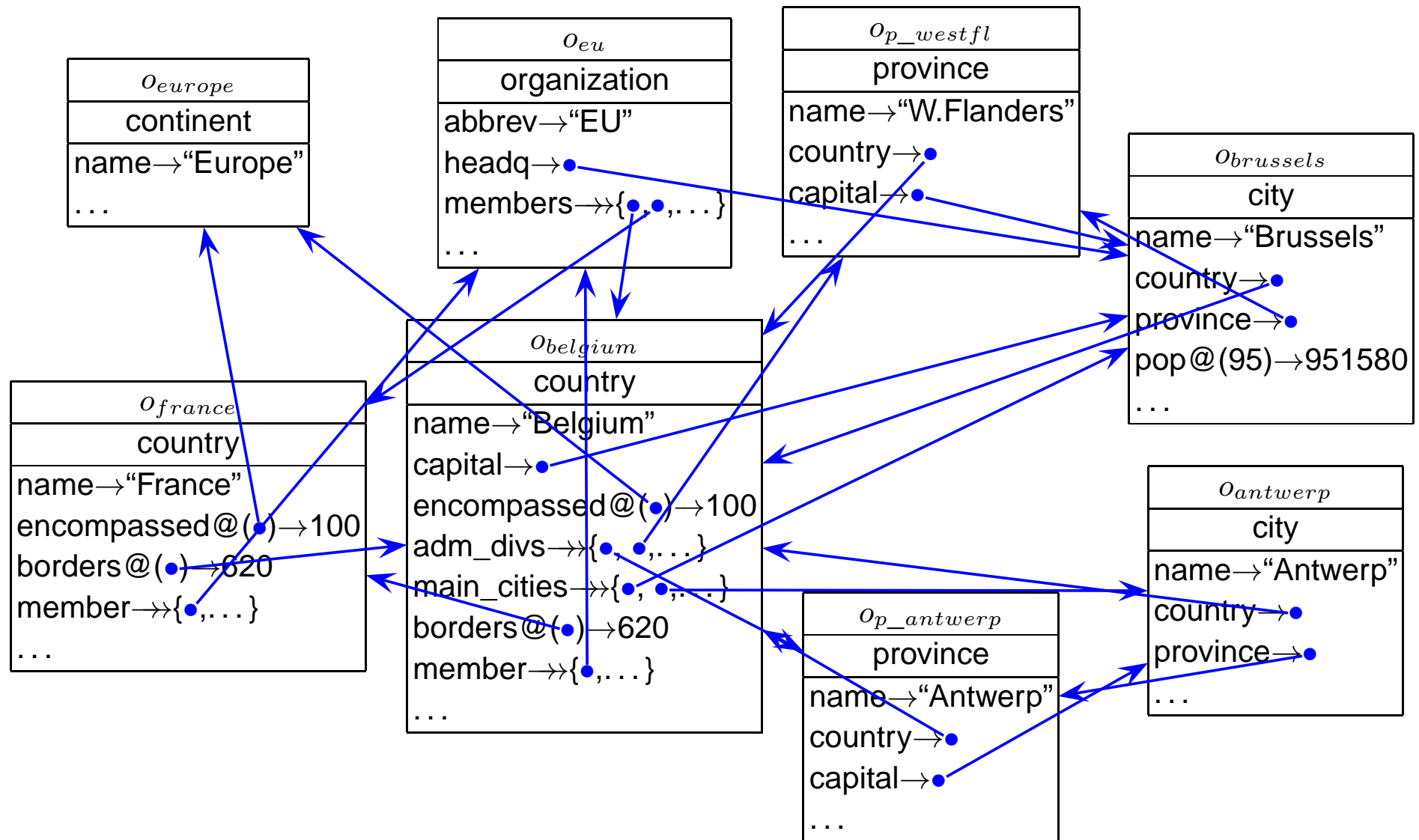
Obelgium : *country*[name→“Belgium”; car_code→“B”;
capital→*Obrussels*; independence→“04 10 1830”;
total_area→30510; population→10170241;
encompassed@(*Oeurope*)→100; pop_growth→0.33;
adm_divs→→{*Op_antwerp*, *Op_westfl*, . . . };
main_cities→→{*Obrussels*, *Oantwerp*, . . . };
borders@(*Ofrance*)→620; borders@(*Ogermany*)→167;
borders@(*Oluxembourg*)→148; borders@(*Onetherlands*)→450].

Obrussels : *city*[name→“Brussels”; country→*Obelgium*;
province→*Op_westfl*; population@(95)→951580].

Op_westfl : *prov*[name→“West Flanders”; country→*Obelgium*;
capital→*Obrussels*; area→3358; population→2253794].

Oeu : *org*[abbrev→“EU”; name→“European Union”;
established→“07 02 1992”; headq→*Obrussels*;
members@("member")→→{*Obelgium*, *Ofrance*, . . . };
members@("applicant")→→{*Ohungary*, *Oslovakia*, . . . }].

REPRESENTATION AS A GRAPH



PROPERTIES

- arbitrary properties of an objects can be stored – simply generate and fill a slot.
self-describing data model
- null values need not to be stored explicitly, slots are simply not filled.
- *navigation* with *path expressions* in combination with *patterns*:
 - population of the province where the capital of Belgium is located:
?- C:country[name→"Belgium"].capital.province[population →P].
 - all names of cities in Belgium:
?- C:country[name→"Belgium"]..main_cities[name→N].
 - sum of population of all provinces of Belgium:
?- Z = sum{X | C:country[name→"Belgium"]..province[population→X]}.
 - can be nested in complex conditions:
?- C:country[encompassed@(_Cont[name→"Europe"])→_X; member→_O[name→"EU"];
population→CP; area→CA]..city[population→CitP; name→N],
CP > 1000000, CA > 100000, CitP > 0.25 * CP.

F-LOGIC: SIGNATURE

The signature can also be formalized in F-Logic:

- properties:

$\langle \text{type} \rangle [\langle \text{property} \rangle \Rightarrow \langle \text{type} \rangle]$ (scalar)

$\langle \text{type} \rangle [\langle \text{property} \rangle \Rightarrow \Rightarrow \langle \text{type} \rangle]$ (set-valued)

analogously with parameters:

$\langle \text{type} \rangle [\langle \text{property} \rangle @ (\langle \text{list-of-types} \rangle) \Rightarrow \langle \text{type} \rangle]$

- $\text{country}[\text{name} \Rightarrow \text{string}; \text{capital} \Rightarrow \text{city};$
 $\text{total_area} \Rightarrow \text{number}; \text{population} \Rightarrow \text{number};$
 $\text{encompassed} @ (\text{continent}) \Rightarrow \text{number};$
 $\text{adm_divs} \Rightarrow \Rightarrow \text{province}; \text{main_cities} \Rightarrow \Rightarrow \text{city};$
 $\text{borders} @ (\text{country}) \Rightarrow \text{number}]$.

$\text{city}[\text{name} \Rightarrow \text{string}; \text{country} \Rightarrow \text{country};$

$\text{province} \Rightarrow \text{province}; \text{population} @ (\text{number}) \rightarrow \text{number}]$.

- queries against metadata: $?- X:\text{country}[M \rightarrow (_V:C)]$ or $?- \text{country}[M \Rightarrow C]$.

F-LOGIC AS A PROGRAMMING LANGUAGE

- object-oriented logic
- *deductive* database language (i.e. Prolog-style) with fixpoint semantics

`<head> :- <body> .`

`?- <query> .`

- e.g., transitive closure can be expressed:

`R[transitive_flows→S] :- R:river[to@(sea)→S].`

`R1[transitive_flows→S] :- R1:river[to@(river)→R2], R2:river[transitive_flows→S].`

Implementations

- The FLORID system (F-Logic Reasoning in Databases; C++):
<http://www.informatik.uni-freiburg.de/~dbis/florid>
- FLORA/FLORA II; XSB-Prolog: <http://flora.sourceforge.net/>

... current use: reasoning in the Semantic Web.

DATA INTEGRATION FROM THE WEB WITH FLORID

- warehouse approach
- direct mapping from HTML trees to F-Logic
- queries against Web pages possible
- wrapper + mediator functionality programmed by F-Logic rules

1998: Generation of the MONDIAL database from the Web

- F-Logic wrappers for several Web pages
 - CIA World Factbook Country Pages
 - CIA World Factbook Organizations Pages
 - “Global Statistics”
 - some smaller sources + the original Karlsruhe TERRA database
- ⇒ materialized F-Logic representation of each source
- F-Logic integration program that *stepwise* materializes an integrated database
- advantages of the warehouse approach for complex integration tasks (basic rules + exceptional and refining rules)

3.3 Summary: Database Aspects (1995)

- **Integration** of data from different, heterogeneous sources:
 - relational distributed/federated databases: metadata queries and -integration (SchemaSQL)
 - arbitrary sources (Tsimmis): metadata queries, **general, flexible data model**
 - sources can also be Web pages (HTML) (Tsimmis, Florid)
- semistructured nature of data
 - no fixed schema
 - implicit null values
 - easily extensible (adding new properties)
 - object as a collection of properties
 - self-describing data + metadata
- languages for semistructured data
 - metadata queries
 - generation of structure

⇒ **flexible**

QUERY LANGUAGES

- declarative
- clause-based; iterator variables and SELECT-FROM-WHERE-clause
(SQL, OQL, Lorel)
- logical/deductive; binding of variables by patterns; constructive semantics of rule heads
 - patterns as terms
(WSL/MSL)
 - patterns as extended path expressions
(F-Logic)
 - as programming languages: fixpoint semantics
- multivalued vs. set-valued semantics
- rule-based (more or less explicit):
 - binding of variables in the “body” (SQL/OQL: FROM-WHERE clause)
 - result generation in the head

F-LOGIC [1989] AS A PREDECESSOR OF XML, RDF, OWL

- semistructured (\Rightarrow XML, RDF)
- self-describing (\Rightarrow XML, RDF)
- data model
 - database model: complex objects, properties, relationships (\Rightarrow XML, RDF)
 - but also knowledge representation model with *built-in reasoning* (\Rightarrow OWL)
 - optional schema information (\Rightarrow XSD, RDFS, OWL)
- query language
 - navigation, path expressions with predicates and multivalued semantics (\Rightarrow XPath)
- derivation rules (\Rightarrow OWL + Rules [SWRL?])

RDF: Resource Description Format, 1997, see Lecture “Semantic Web”

OWL: Web Ontology Language, 2002 [OIL: 2000], see Lecture “Semantic Web”

3.4 Situation 1996

- Experiences with SQL (and ODMG/OQL) as database languages
 - standardization vs. products
- document management with SGML (Structured Generic Markup Language), CSS (Cascading Stylesheets) and DSSSL
- data exchange/access via internet/Web:
 - homogeneous solution necessary
 - availability of documents and data in HTML:
 - * very simple variant of SGML
 - * “native” HTML data (handwritten)
 - * mapping of SGML (document management) to HTML (publication) by CSS
 - * HTML-Web-Servers over relational databases

⇒ “Global” approach coordinated by the W3C (World Wide Web Consortium):
development of a data model (+ language), that can handle (legacy-)databases,
documents and Web (=HTML)

THE W3C (WORLD WIDE WEB CONSORTIUM)

- <http://www.w3.org>.
- founded in 1994 for developing common protocols and languages for the World Wide Web and to ensure interoperability of applications in the Web.
(Tim Berners-Lee, MIT, CERN)
- following the principles of OMG/ODMG who developed the CORBA and ODL/OIF/OQL standards
- members: companies and research institutes
- definition of working groups
- notes → working drafts → recommendations
- not only XML, but also many other Web-related issues

3.5 Documents: SGML and HTML

- Structuring (und presentation) are called (*logical and optical*) “*markup*”.
(document = content + markup)
- SGML (Standard Generalized Markup Language),
development (IBM) since 1979, standard 1986.
structuring and markup of documents, widely used in publishing.
- for publishing in the Web:
HTML (Hypertext Markup Language), development since 1989 (CERN), standard 1991.

⇒ **HTML is an SGML application** with a fixed syntax

(tags, attributes, later: DTD).

goal: optical markup, as a side effect also some structuring of the documents (cf. <P>-Tag).

- SGML much more flexible than HTML → more complex → not suitable for browsers (HTML allows for efficient and fault-tolerant parsing)
- SGML sources can be transformed to HTML by stylesheets (CSS: Cascading Style Sheets).

Chapter 4

XML (Extensible Markup Language)

Introduction

- SGML *very* expressive and flexible
HTML very specialized.
- Summer 1996: John Bosak (Sun Microsystems) initiates the XML Working Group (SGML experts), cooperation with the W3C.

Development of a subset of SGML that is simpler to implement and to understand

<http://www.w3.org/XML/>: the homepage for XML at the W3C

⇒ XML is a “stripped-down version of SGML”.

- for understanding XML, it is not necessary to understand everything about SGML ...

HTML

let's start the other way round: HTML ... well known, isn't it?

- tags: pairwise opening and closing: `<TABLE> ... </TABLE>`
- “empty” tags: without closing tag `
`, `<HR>`
- `<P>` is *in fact* not an empty tag (it should be closed at the end of the paragraph)!
- attributes: `<TD colspan = “2”> ... </TD>`
- empty tags with attributes:
``
- content of tag structures: `<TD>123456</TD>`
- nested tag structures: `<TH>Name</TH>`
``
`Homepage of the IFI`

⇒ hierarchical structure

- Entities: `ä = ä` `ß = ß`

HTML

- browser must be able to interpret tags
→ semantics of each tag is fixed for all (?) browsers.
- fixed specifications how tags can be nested
(described by a DTD (Document Type Definition))

```
<body><H1>... </H1><H2>... </H2>  
    <P> ... </P>  
    <H2>... </H2>  
    <P> ... </P>  
    <H1>... </H1><H2>... </H2>  
    <P> ... </P>  
  
</body>
```

- analogously for tables and lists ...
- reality: people do in general not adhere to this structure
 - closing tags are omitted
 - structuring levels are omitted
- parser has to be fault-tolerant and auto-completing

KNOWLEDGE OF HTML FOR XML?

- intuitive idea – but only of the *ASCII representation*
- this is *not a data model*
- no query language
- only a very restricted viewpoint:
HTML is a markup language for browsers
(note: we don't “see” HTML in the browser, but only what the browser makes out of the HTML).

Not any more.

GOALS OF THE DEVELOPMENT OF XML

- XML must be directly usable and transmitted in the internet (Unicode-Files),
- XML must support a wide range of applications,
- XML must be compatible with SGML,
- XML documents must be human-readable and understandable,
- XML documents must be easy to create,
- it must be easy to write programs that evaluate/process/parse XML documents.

DIFFERENCES BETWEEN XML AND HTML?

- Goal: *not browsing*, but representation/storage of (semistructured) data (cf. SGML)
- SGML allows the definition of new tags according to the application semantics; each SGML application uses its own *semantic tags*.
These are defined in a DTD (Document Type Definition).
- HTML is *an* SGML application (cf. `<HTML>` at the beginning of each document `</HTML>`), that uses the DTD “HTML.dtd”.
- In XML, (nearly) arbitrary tags can be defined and used:

```
<country> ... </country>  
<city> ... </city>  
<province> ... </province>  
<name> ... </name>
```
- These *elements* represent objects of the application.

XML AS A META-LANGUAGE FOR SPECIALIZED LANGUAGES

- For each application, it can be chosen which “notions” are used as element names etc.:
⇒ document type definition (DTD)
- the set of allowed element names and their allowed nesting and attributes are defined in the DTD of the document (type).
- the DTD describes the *schema*
- XML is a *meta-language*, each DTD defines an own language
- for an application, either a new DTD can be defined, or an existing DTD can be used
→ standard-DTDs
- HTML has (as an SGML application) a DTD

EXAMPLE: MONDIAL

```
<mondial>
```

```
:
```

```
<country code="D" capital="city-D-Berlin" memberships="EU NATO UN ...">
```

```
<name>Germany</name>
```

```
<encompassed continent="europe">100</encompassed>
```

```
<population year="1995">83536115</population>
```

```
<ethnicgroup name="German">95.1</ethnicgroup>
```

```
<ethnicgroup name="Italians">0.7</ethnicgroup>
```

```
<religion name="Roman Catholic">37</religion>
```

```
<religion name="Protestant">45</religion>
```

```
<language name="German">100</language>
```

```
<border country="F" length="451"/>
```

```
<border country="A" length="784"/>
```

```
<border country="CZ" length="646"/>
```

```
:
```


Example: Mondial (Cont'd)

```

:
<province id="prov-D-berlin" capital="city-D-berlin">
  <name>Berlin</name>
  <population year="1995">3472009</population>
  <city id="city-D-berlin">
    <name>Berlin</name> <population year="1995">3472009</population>
  </city>
</province>
<province id="prov-D-baden-wuerttemberg" capital="city-D-stuttgart">
  <population year="1995">10272069</population>
  <name>Baden Wuerttemberg</name>
  <city id="city-D-stuttgart">
    <name>Stuttgart</name> <population year="95">588482</population>
  </city>
  <city id="cty-D-mannheim"> ... </city>
:
</province>
:
</country>
:
</mondial>
```

CHARACTERISTICS:

- hierarchical “data model”
- subelements, attributes
- references
- ordering? documents – yes, databases – no

Examples can be found at

<http://dbis.informatik.uni-goettingen.de/Mondial/#XML>

XML AS A DATA MODEL

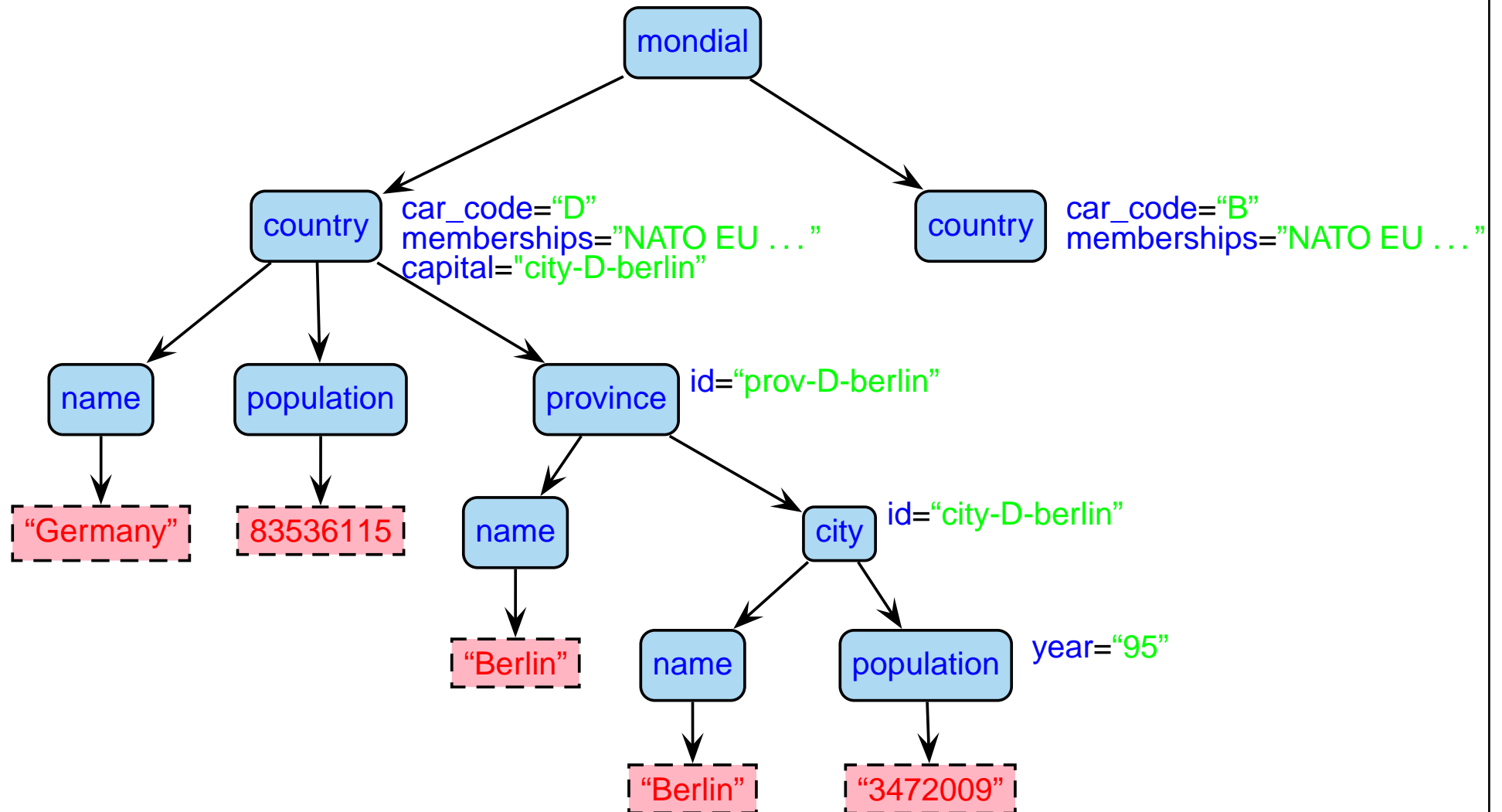
XML is much more than only the ASCII representation shown above as known from HTML
(see also introductory talk)

- abstract data model (comparable to the relational DM)
- abstract datatype: DOM (Document Object Model) – see later
- many concepts around XML
(XML is *not* a programming language!)
 - higher-level declarative query/manipulation language(s)
 - notions of “schema”

4.1 Structure of the Abstract XML Data Model (Overview)

- for each document there is a document node which “is” the document, and which contains information about the document (reference to DTD, doctype, encoding etc).
- the document itself consists of nested *elements* (tree structure),
- among these, exactly one *root element* that contains all other elements and which is the only child of the document node.
- elements have an element type (e.g. Mondial, Country, City)
- element content (if not empty) consists of text and/or *subelements*. These *child nodes* are ordered.
- elements may have *attributes*. Each attribute node has a name and a value (e.g. (car_code, “D”). The attribute nodes are unordered.
- *empty elements* have no content, but can have attributes.
- a *node* in an XML document is a logical unit, i.e., an element, an attribute, or a text node.
- the allowed structure can be restricted by a schema definition.

EXAMPLE: MONDIAL AS A TREE



EXAMPLE: MONDIAL AS A NESTED STRUCTURE

mondial

country car_code="D" memberships="EU NATO ..." capital="city-D-berlin"

name "Germany"

population "83536115"

province id="prov-D-berlin"

name "Berlin"

city id="city-D-berlin"

name "Berlin"

population year="1995" "3472009"

country car_code="B" memberships="EU NATO ..."

OBSERVATIONS

- there is a global order (preorder-depth-first-traversing) of all element- and text nodes, called *document order*.
- actual text is only present in the **text-nodes**
Documents: if all text is concatenated in document order, a pure text version is obtained.
Exercise: consider an HTML document.
- element nodes serve for structuring (but do not have a “value” for themselves)
- attribute nodes contain values whose semantics will be described in more detail later
 - attributes that describe the elements in more detail (e.g. `td/@colspan` or `population/@year`)
 - IDs and references to IDs
 - can be used for application-specific needs

4.2 XML ASCII Representation

- Tree model and nested model serve as abstract datatypes (see later: DOM)

data exchange? how can an XML document be represented?

- a relational DB can be output as a finite set of tuples (cf. relational calculus)
country("Germany", "D", 83536115, 356910, "Berlin", "Berlin")
or
country(Name: "Germany", Code: "D", Population: 83536115, Area: 356910,
Capital: "Berlin", CapitalProvince: "Berlin")
- object-oriented databases: OIF (Object Interchange Format)
- OEM-tripels, F-Logic-frames
- XML?
Exporting the tree in a *preorder-depth-first-traversing*.
The node types are represented in a specified syntax:
⇒ XML as a representation language

ASCII: XML AS A REPRESENTATION LANGUAGE

- elements are limited by
 - opening `<Country>` and
 - closing *tags* `</Country>`,
 - in-between, the *element content* is output recursively.

- Element content consists of text

`<Name> United Nations </Name>`

- and *subelements*:
`<Country> <City> ... </City>`
`<City> ... </City>`
`</Country>`

- *attributes* are given in the opening tag:

`<Country car_code="D"> ... </Country>`

where attribute values are always given as strings, they do not have further structure. The difference between value- and reference attributes is not visible, but is only given by the DTD.

- *empty elements* have only attributes: `<border country="F" length="451"/>`

XML AS A REPRESENTATION LANGUAGE: GRAMMAR

The language “XML” defined as above can be given as an EBNF grammar:

```
Document ::= Element
Element  ::= "<" ElementName Attribute* ">" Content "</" ElementName ">"
           | "<" ElementName Attribute* "/>"
Content  ::= (Element | Text)+
Text     ::= characters including whitespace
Attribute ::= AttributeName "=" AttributeValue "\""
ElementName, AttributeName ::= character string with some restrictions
AttributeValue ::= characters including whitespace
```

- note that this grammar does not guarantee that the opening and closing tags match!
- instead of `'`, also the usual `"` are allowed
- strict adherence to these rules (closing and empty elements) is required.
- an XML instance given as ASCII is *well-formed*, if it satisfies these rules.
- “XML parsers” process this input.

XML PARSER

- an *XML parser* is a program that processes an XML document given in ASCII representation according to the XML grammar, and generates a result:
 - **correctness**: check for well-formedness (and adherence to a given DTD)
 - **DOM-parser**: transformation of the XML instance into a DOM model (implementation of the **abstract datatype**; see later).
 - **SAX-parser**: traversing the XML tree and generation of a sequence of “events” that *serialize* the document (see later).
- XML parsers are required to accept only well-formed instances.
 - simple grammar, simple (non-fault-tolerant) parser
 - HTML: fault-tolerant parsers are much more complex (fault tolerance wrt. omitted tags is only possible when the DTD is known)
- each XML application must contain a parser for processing XML instances in ASCII representation as input.

XML PARSING IN THE GENERAL CASE

- ElementName is a separate production and

Element ::= “<” ElementName Attribute* “>” Content “</” ElementName “>”
| “<” ElementName Attribute* “/>”

does not guarantee matching tags

⇒ not context-free!

- Nevertheless, context-free-style parsing with push-down-automaton *without fixed stack alphabet* possible:
 - for every opening tag, put ElementName on the stack
 - for every closing tag, compare with top of stack, pop stack.

⇒ linear-time parsing

- Exercise: give an automaton for parsing XML and describe the handling of the stack (solution see Slide 179).

VIEWING XML DOCUMENTS?

- as a file in the editor
 - emacs with xml-mode
 - Linux/KDE: kxmleditor
- browser cannot “interpret” XML
(in contrast to HTML)
- with “show source” in a browser:
current versions of most browsers show XML in its ASCII representation with indentation and allow to open/close elements/subtrees.
- but, in general, XML is not intended for viewing:
→ transformation to HTML by XSLT stylesheets
(see later)

4.3 Datatypes and Description of Structure for XML

- relational model: atomic data types and tuple types
- object-oriented model: literal types and object types, reference types

Data Types in XML

- data types for text content
- data types for attribute values
- element types (as “complex objects”)
- somewhat different approaches in DTD (document-oriented, coarse) and XML Schema (database-oriented, fine)

DOCUMENT TYPE DEFINITION – DTD

- the set of allowed tags and their nestings and attributes are specified in the **DTD** of the document (type).
- the idea of the DTD comes from the SGML area
 - meets the requirements for describing document structure
 - does not completely meet the requirements of the database area
→ **XML Schema** (later)
 - simple, and easy to understand.
- the DTD for a document type *doctype* is given by a grammar (context-free; regular expression style) that characterizes a class of documents:
 - **what elements** are allowed in a document of the type *doctype*,
 - **what subelements** they have (element types, order, cardinality)
 - **what attributes** they have (attribute name, type and cardinality)
 - additionally, “entities” can be defined (they serve as constants or macros)

DATA TYPES OF DTDs

- text content: PCDATA – parsed character data
it is up to the application to distinguish between string data and numerical data
- data types for attribute values:
 - CDATA: simple strings
 - NMTOKEN: string without blanks
 - NMTOKENS: a list of tokens, separated by blanks
 - ID: like NMTOKEN, each value must be unique in the document
 - IDREF: like NMTOKEN, each value must occur in the same document as an ID value
 - IDREFS: the same, multivalued
- element types: definition of structure in the style of regular expressions.

DTD: ELEMENT TYPE DEFINITION – STRUCTURE OF THE ELEMENT CONTENTS

<!ELEMENT *elem_name struct_spec*>

- EMPTY: empty element type,
- (#PCDATA): text-only content
- (*expression*): expression over element names and combinators (same as for regular expressions). Note that the expression must be deterministic.
 - “,”: sequence,
 - “|”: (exclusive-)or (choice),
 - “*”: arbitrarily often,
 - “+”: at least once,
 - “?”: optional
- (#PCDATA|*elem_name*₁|...|*elem_name*_{*n*})*
mixed content, here, only the types of the subelements that are allowed to occur together with #PCDATA can be specified; no statement about order or cardinality.
- ANY: arbitrary content

Element Type Definition: Examples

- from HTML: images have only attributes and no content
`<!ELEMENT img EMPTY >`
- from Mondial:
`<!ELEMENT country (name, encompassed+, population*,
ethnicgroup*, religion*, border*,
province+ | city+)>`
`<!ELEMENT name (#PCDATA)>`
- for text documents:
`<!ELEMENT Section (Header,
(Paragraph|Image|Figure|Subsection)+,
Bibliography?)>`
- Element type definitions by **regular expressions**
⇒ can be checked by **finite state automata**

DTD: ATTRIBUTE DEFINITIONS

- General: an element contains at most one attribute of every attribute name.
- details about allowed attribute names and their types are specified in the DTD.

```
<!ATTLIST  elem_name
           attr_name1  attr_type1  attr_constr1
           :             :             :
           attr_namen  attr_typen  attr_constrn>
```

- *attr_type_i*: value/reference attribute and scalar/multi-valued
 - CDATA: arbitrary text.
 - NMTOKEN: scalar, token-content (text without blanks).
 - NMTOKENS: multi-valued, token-content.
 - (*const*₁ | ... | *const*_{*k*}): scalar, from a given domain.
 - ID: distinguished scalar attribute, token-content, unique in the whole document.
 - IDREF: scalar, its value is a token that occurs as a value of an ID attribute in the same document (reference).
 - IDREFS: multi-valued reference attribute.

DTD: Attribute Definitions (cont'd)

```
<!ATTLIST  elem_name
            attr_name1  attr_type1  attr_constr1
            :              :              :
            attr_namen  attr_typen  attr_constrn>
```

- *attr_constr*_{*i*}: minimal cardinality
 - #REQUIRED: attribute must be present for each element of this type.
 - #IMPLIED: attribute is optional.
 - *default*: Default-value (non-monotonic value inheritance).
 - #FIXED *value*: attribute has the same (given) value for each element of this type (monotonic value inheritance).

DTD: ATTRIBUTE-DEFINITIONS (EXAMPLES)

<!ATTLIST Country

| | | |
|-------------|----------|------------|
| Code | ID | #REQUIRED |
| Capital | IDREF | #REQUIRED |
| Memberships | IDREFS | #IMPLIED |
| Products | NMTOKENS | #IMPLIED > |

<!ATTLIST desert

| | | |
|---------|------------------|----------------|
| id | ID | #REQUIRED |
| Type | (sand,rocks,ice) | 'sand' |
| Climate | NMTOKENS | #FIXED 'dry' > |

- when an XML parser reads an XML instance and its DTD, it fills in default and fixed values.

DTD AND XML INSTANCES

- Each DTD defines an own markup language (i.e., an XML application – HTML is one, Mondial is another).
- an XML instance has a *document node* (which is not the root node, but even “superior”) that contains among other things information about the DTD.
(see next slides ...)
- the root element of the document must be of an element type that is defined in the DTD.
- an XML instance is *valid* wrt. a DTD if it satisfies the structural constraints specified in the DTD.
Validity can be checked by an extended finite state automaton in linear time.
- XML-instances can exist without a DTD (but then, it is not explicitly specified what their tags “mean”).

XML DOCUMENT STRUCTURE: THE PROLOG

The *prolog* of an XML document in ASCII-representation contains additional information about the document (associated with the document node):

- XML declaration (with optional attributes)

`<? xml version="1.0" encoding="utf-8"?>`

`encoding="ISO-8859-1"` allows additionally German "Umlauts".

- document type *declaration*: indication of the document type, and where the document type *definition (DTD)* can be found.

– `<!DOCTYPE name {SYSTEM|PUBLIC public-id} url>`

SYSTEM *url*: own document type,

name: one of the element names given in the DTD

`<!DOCTYPE Mondial SYSTEM "mondial-2.0.dtd">`

PUBLIC *public-id url*: standard document type (e.g. XHTML), or

– `<!DOCTYPE name [dtd]>`

with DTD directly included *in* the document.

- then follows the document content (i.e., the root node with the document body as its content).

EXAMPLE: MONDIAL

```
<?xml version="1.0" encoding="utf-8"?>
<!DOCTYPE mondial SYSTEM "mondial-2.0.dtd">
<mondial>
  <country car_code="AL" area="28750" capital="cty-cid-cia-Albania-Tirane"
    memberships="org-BSEC org-CE org-CCC ...">
    <name>Albania</name> <population>3249136</population>
    <encompassed continent="europe" percentage="100"/>
    <ethnicgroups percentage="3">Greeks</ethnicgroups>
    <ethnicgroups percentage="95">Albanian</ethnicgroups>
    <border country="GR" length="282"/> <border country="MK" length="151"/>
    <border country="YU" length="287"/>
    <city id="cty-cid-cia-Albania-Tirane" is_country_cap="yes" country="AL">
      <name>Tirane</name>
      <longitude>10.7</longitude> <latitude>46.2</latitude>
      <population year="87">192000</population>
    </city>
  :
</country>
:
</mondial>
```


TOOL: XMLLINT

`xmllint` is a simple tool that allows (among other things – see later) to validate a document (belongs to libxml2):

- `man xmllint`: lists all available commands
- currently, we are mainly interested in the following:
`xmllint -loadtd -valid -noout mondial-europe.xml`
validates an XML document wrt. the DTD given in the prolog.

XMLLINT: Further Functionality (see later)

XMLLINT can be used to “visit” the document, and to walk through it:

- `call xmlint -loadtd -shell mondial-europe.xml.`

Then, one gets a “navigating shell” “inside” the XML document tree (very similar to navigating in a UNIX directory tree):

- `validate`: validates the document
- `cd xpath-expression`: navigates into a node (the XPath expression must uniquely select a single node)
relativ: `cd country[1]`
absolut: `cd //country[@car_code="D"]`
- `pwd`: gives the path from the root to the current position
- `cat`: prints the current node
- `cat xpath-expression`
`cat ../city/name`
- `du xpath-expression` lists the content of the node that is selected by `xpath-expression` (starting from the current node)
- `dir xpath-expression` prints the node type and attributes of the selected node

Example: "Books" from W3C XML Use Cases

```
<?xml version="1.0" encoding="utf-8"?>
<!DOCTYPE bib SYSTEM "books.dtd">
<!-- from W3C XML Query Use Cases -->
<bib>
  <book year="1994">
    <title>TCP/IP Illustrated</title>
    <author><last>Stevens</last><first>W.</first></author>
    <publisher>Addison-Wesley</publisher>
    <price>65.95</price>
  </book>
  <book year="1992">
    <title>Advanced Programming in the Unix environment</title>
    <author><last>Stevens</last><first>W.</first></author>
    <publisher>Addison-Wesley</publisher>
    <price>65.95</price>
  </book>
  <book year="2000">
    <title>Data on the Web</title>
    <author><last>Abiteboul</last><first>Serge</first></author>
    <author><last>Buneman</last><first>Peter</first></author>
    <author><last>Suciu</last><first>Dan</first></author>
    <publisher>Morgan Kaufmann Publishers</publisher>
    <price>39.95</price>
  </book>
  <book year="1999">
    <title>Economics of ... for Digital TV</title>
    <editor>
      <last>Gerbarg</last><first>Darcy</first>
      <affiliation>CITI</affiliation>
    </editor>
    <publisher>Kluwer Academic Publishers</publisher>
    <price>129.95</price>
  </book>
</bib>
```

[see XML-DTD/books.xml]

Exercise: Generate a DTD for the above XML

... do it step-by-step, using a validator:

- for all element types:
 <!ELEMENT *name* ANY>
- declare <!ATTLIST *name* ...> where needed
- validate
- stepwise refinement of content models ...
- ... blackboard demonstration ...
- solution see Slide 175

DATA-CENTERED VS. DOCUMENT-CENTERED XML DOCUMENTS

Data-Centered XML Documents

- very regular structure with “data fields”
- only some text
- no naturally induced tree structure

Document-Centered XML Documents

- tree structure with much text (text content is the text of the document)
- non-regular structure of elements
- logical markup of the documents
- annotations of the text by additional elements/attributes

Semistructured XML Documents

- combine both (e.g. medical information systems)

SUBELEMENTS VS. ATTRIBUTES

When designing an XML structure, often the choice of representing something as subelement or as attribute is up to the designer.

Document-Centered XML

- the concatenation of the whole text content should be the “text” of the document
- element structures for logical markup and annotations
- attributes contain additional information *about* the structuring elements.

Data-Centered XML

- more freedom
- attributes are unstructured and cannot have further attributes
- elements allow for structure and refinement with subelements and attributes
- using DTDs as schema language allows the following functionality only for attributes:
 - usage as identifiers (ID)
 - restrictions of the domain
 - default values(XML Schema and XLink allow many more things)

EXAMPLES AND EXERCISES

- The MONDIAL database is used as an example for practical experiments.
See <http://dbis.informatik.uni-goettingen.de/Mondial#XML>.
- many W3C documents base on examples about a literature database (book, title, authors, etc.).
- each participant (possibly in groups) should choose an *own* application area to set up an own example and to experiment with it.
 - from the chosen branch of study?
 - database of music CDs
 - lectures and persons at the university
 - exams (better than FlexNever?)
 - calendar and diary
 - other ideas ...

Exercise: Define a DTD and generate a small XML document for your chosen application.

EXERCISES

- Validate your example document with a suitable prolog and internal DTD.
- put your DTD publicly in your public-directory and validate a document that references this DTD as an external DTD.
- take a DTD+url from a colleague and write a small instance for the DTD and validate it.

DATA EXCHANGE WITH XML

For *Electronic Data Interchange (EDI)*, a commonly known+used DTD is required

- producers and suppliers in the automobile industry
- health system, medical area
- finance/banking

PROCEEDING

Usually, XML data is exchanged in its ASCII representation.

- XML-Server make documents in the ASCII representation accessible (i.e., as a stream or as a textfile)
- applications *parse* this input (linear) and store it internally (DOM or anything else).

4.3.1 Aside: XML Parsing

... one of the objectives of this lecture is also to show the applications and connections of basic concepts of CS ...

- XML/DTD: content models are regular expressions
 - ⇒ can be checked by finite state automata
 - design one automaton for each `<!ELEMENT ...>` declaration
 - design a combined automaton for validating documents against a given DTD
 - extension to attributes: straightforward (when processing opening tags, dictionary-based)
 - checking for well-formedness and validity in linear time
 - * with a DOM parser: during generation of the DOM
 - * with a SAX parser: streaming, on the fly
 - * using a DOM instance: depth-first traversal
 - without a DTD: requires a push-down automaton (remembering opening tags); still linear time
 - checking well-formedness
 - generating a DOM instance, or on-the-fly (SAX)

FINITE STATE AUTOMATA FOR VALIDATION

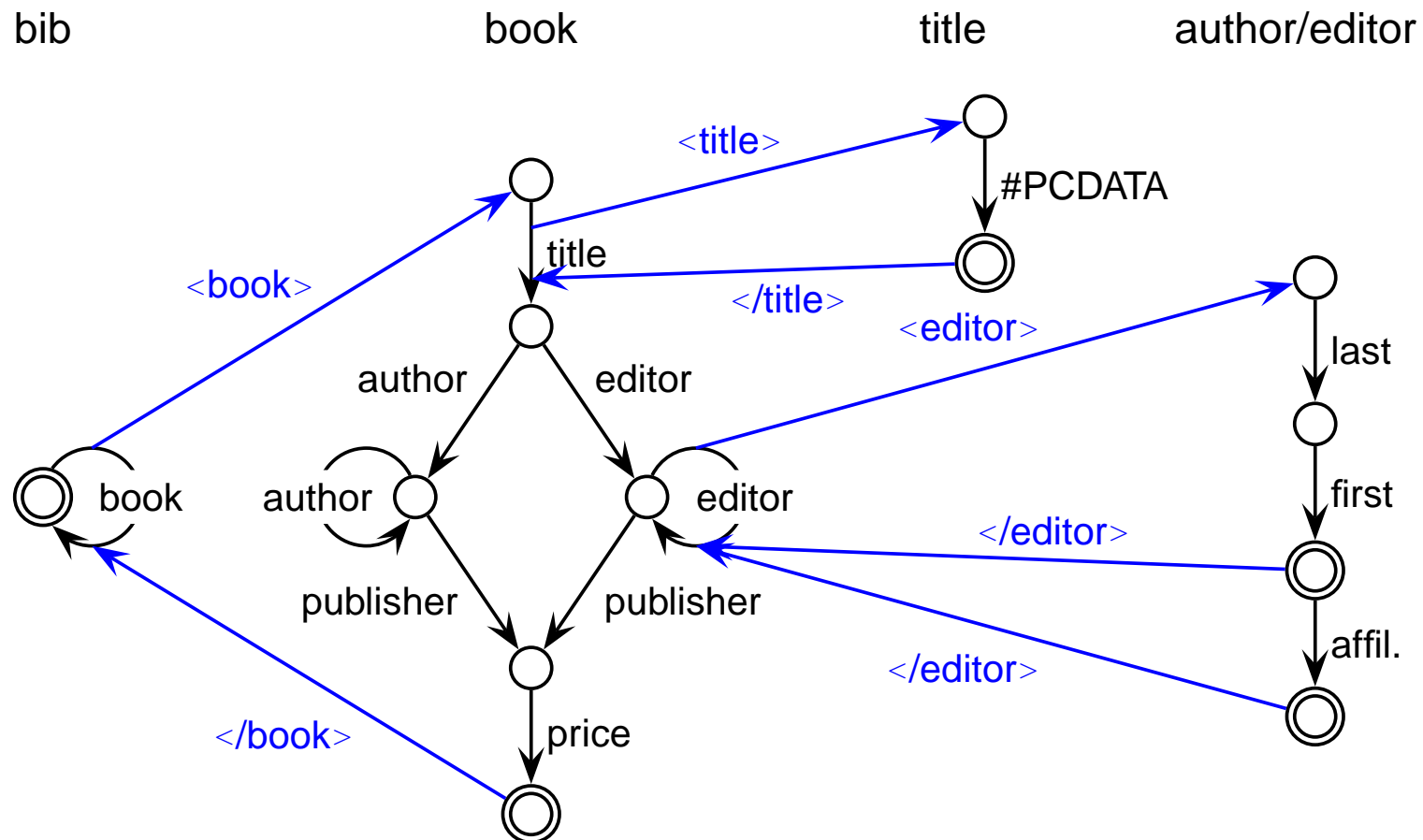
EXAMPLE: BOOKS.DTD

Consider the “books” example:

```
<!ELEMENT bib (book*)>
<!ELEMENT book (title, (author+ | editor+), publisher, price)>
<!ATTLIST book year CDATA #REQUIRED>
<!ELEMENT title (#PCDATA)>
<!ELEMENT author (last, first, affiliation?)>
<!ELEMENT last (#PCDATA)>
<!ELEMENT first (#PCDATA)>
<!ELEMENT publisher (#PCDATA)>
<!ELEMENT editor (last, first, affiliation?)>
<!ELEMENT price (#PCDATA)>
<!ELEMENT affiliation (#PCDATA)>
```

Finite State Automata

- individual automata for element content models
(recall that the content model must be deterministic)
- + detailed by nesting (jumping on opening/closing tags)



XML GRAMMAR IN PRESENCE OF A DTD

Consider the grammar from Slide 150:

- Element names known from a DTD: context-free grammar

Document	::=	Element
Element	::=	"<bib" Attribute* ">" Content "</bib>"
Element	::=	"<book" Attribute* ">" Content "</book>"
:	:	:
Content	::=	(Element Text)+
Text	::=	<i>characters</i>
Attribute	::=	AttributeName "=" AttributeValue "'"
AttributeValue	::=	<i>characters</i>

- there is even a regular grammar, see above automata, but this is not derived from the XML EBNF.

XML GRAMMAR IN GENERAL

- no DTD present/element names not known:

Consider the grammar from Slide 150:

- ElementName is a separate production and

$$\text{Element} ::= \text{"<" ElementName Attribute* ">" Content "</" ElementName ">"}$$
$$| \text{"<" ElementName Attribute* "/>"}$$

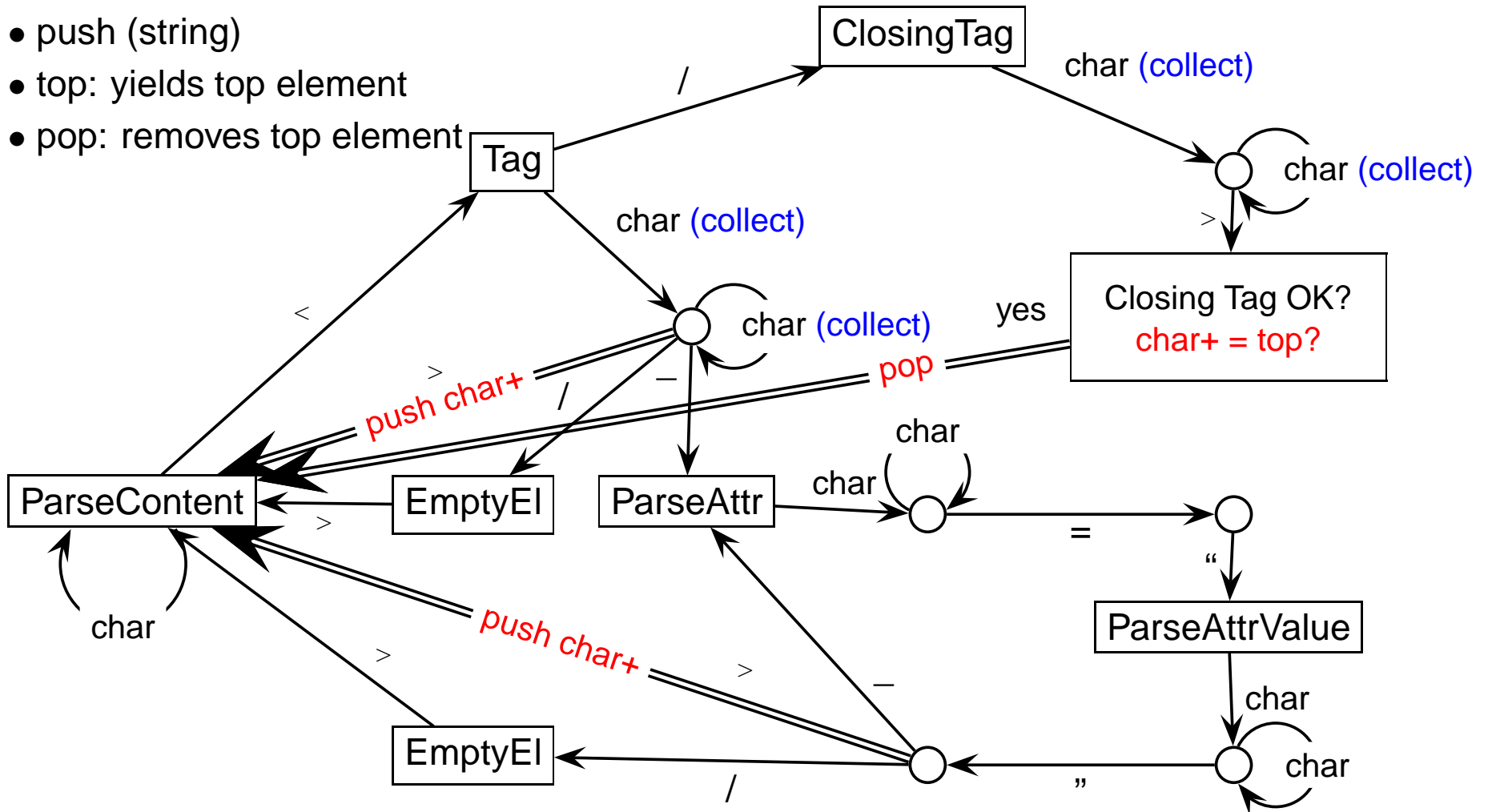
does not guarantee matching tags.

- Nevertheless, context-free-style parsing with push-down-automaton *without fixed stack alphabet* possible:
 - for every opening tag, put ElementName on the stack
 - for every closing tag, compare with top of stack, pop stack.
- Automaton: see next slide.

XML GRAMMAR IN GENERAL

Stack Commands:

- push (string)
- top: yields top element
- pop: removes top element



4.4 Example: XHTML

- XML documents that adhere to a strict version of the HTML DTD
- Goal: browsing, publishing
- DTD at <http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd>
(note that the DTD requires also some entity files)
- Validator at <http://validator.w3.org/>
- Example at ... DBIS Web Pages
- only the text content is shown in the browser, all other content describes *how* the text is presented.
- no logical markup of the documents (sectioning etc), but
- only optical markup (“how is it presented”).

Exercise

Design (and validate) a simple homepage in XHTML, and put it as `index.html` in your public-directory.

4.5 Miscellaneous about XML

4.5.1 Remarks

- all letters are allowed in element names and attribute names
- text (attribute values and element content) can contain nearly all characters. Western european umlauts are allowed if the XML identification contains `encoding="UTF-8"` or `encoding="ISO-8859-1"` etc.
- comments are enclosed in `<!-- ... -->`
- inside XML content,
`<![CDATA[...]]>`
(*character data sequences*) can be included that are not parsed by XML parsers, but which are copied character-by-character.

4.5.2 Entities

Entities serve as macros or as constants and are defined in the DTD. They are then accessible as “&entityname;” in the XML instance and in the DTD:

<!ENTITY *entity_name* *replacement_text*>

- additional special characters, e.g. ç:

DTD: <!ENTITY ccedilla “ç”>

XML: president=“Françla;ois Mitterand”

- reserved characters can be included as *references to predefined entities*:

< = < (less than), > = > (greater than)

& = & (ampersand), *space* = , *apostroph* = ', *quote* = "

ä = ä, ..., Ü = Ü

<name>Düsseldorf </name>

- characters can also be given directly as character references, e.g. (Space),  (CR).

Entities (cont'd)

- global definitions that may change can be defined as constants:

DTD: `<!ENTITY server "http://www.informatik.uni-goettingen.de">`

XML: `<url> &server;/dbis </url>`

- macros that are needed frequently:

DTD: `<!ENTITY european`

`"<encompassed continent='europe'>100</encompassed">`

XML: `<country car_code="D">`

`<name >Germany</name>`

`&european; ...`

`</country>`

- note: single and double quotes can be nested.

PARAMETER ENTITIES

Entities that should be usable only in the DTD are defined as *parameter entities*:

- macros that are needed frequently:

```
<!ENTITY % namedecl "name CDATA #REQUIRED">
```

```
<!ATTLIST City
```

```
    %namedecl;
```

```
    zipcode ID #REQUIRED>
```

- define enumeration types:

```
<!ENTITY % waters "(river|lake|sea)">
```

```
<!ATTLIST City_located_at
```

```
    type %waters; #REQUIRED
```

```
    at IDREF #REQUIRED>
```

ENTITIES FROM EXTERNAL SOURCES

Entity “collections” can also be used from external sources as *external entities*:

```
<!ENTITY entity_name SYSTEM “url”>
```

is an entity that stands for a remote resource which itself defines a set of entities by

```
<!ENTITY entity_name’ replacement_text>
```

e.g. a set of technical symbols:

```
<!ENTITY % isotech SYSTEM  
    “http://www.schema.net/public-text/ISOtech.pen”>  
%isotech;
```

the reference %isotech; makes then all symbols accessible that are defined in the external resource.

This can be iterated for defining “style files” that collect a set of external resources that are used by an author.

4.5.3 Integration of Multimedia

- for (external) non-text resources, it must be declared which program should be called for showing/processing them. This is done by **NOTATION** declarations:

```
<!NOTATION notation_name SYSTEM "program_url">
```

```
<!NOTATION postscript SYSTEM "file:/usr/bin/ghostview">
```

- the entity definition is then extended by a declaration which notation should be applied on the entity:

```
<!ENTITY entity_name SYSTEM "url"
```

```
    NDATA notation_name>
```

```
<!ENTITY manual SYSTEM "file:/.../name.ps"
```

```
    NDATA postscript>
```

- the *application program* is then responsible for evaluating the entity and the NDATA definition.
- XLink will later present another mechanism for referencing resources.

4.6 Summary and Outlook

XML: “basic version” consists of DTD and XML documents

- tree with additional cross references
- hierarchy of nested elements
- order of the subelements
 - documents: 1st, 2nd, . . . section etc.
 - databases: order in general not relevant
- attributes
- references via IDREF/IDREFS
 - documents: mainly cross references
 - databases: part of the data (relationships)
- XML model similar to the network data model:
relationships are mapped into the structure of the data model
 - the basic explicit, stepwise navigation commands of the network data model have an equivalent for XML in the **DOM**-API (see later), but
 - XML also provides a declarative, high-level, set-oriented language.

REQUIREMENTS

- Documents: logical markup (Sectioning etc.)
presentation on Web pages in (X)HTML? – transformation languages
- databases: structuring of data;
several equivalent alternatives
query languages?
presentation on Web pages in (X)HTML? – transformation languages
- application-specific formats
e.g. XHTML: browsing
DTDs are induced by the application-programs
Web-Services: WSDL, UDDI; CAD; ontology languages; . . .
transformation between different XML languages
application-programs must “understand” XML internally

FURTHER CONCEPTS OF THE XML WORLD

Extensions:

- namespaces: use of different DTDs in a database (see Slide 223)
- APIs: DOM, SAX
- theoretical foundations
- query languages: XPath, XML-QL, Quilt, XQuery
- stylesheets/transformation languages: CSS, DSSSL, XSL
- better schema language: XML Schema
- XML with inter-document handling: XPointer, XLink

4.7 Recall

- XML as an abstract *data model*
 - cf. relational DM
 - XML now has become less abstract: creation of instances in the editor, validating, viewing ...
- a data model needs ... implementation? theory?
- ... first, something else: *abstract datatype, interface(s)*
 - constructors, modifiers, selectors, predicates (cf. Info I)
- here: “two-level model”
 - as an ADT (programming interface): Document Object Model (DOM): detailed operations as usual in programming languages (Java, C++).
 - as a database model (end user interface; declarative): import (parser), *queries*, updates
- theory: formal specification of the semantics of the languages, other issues are the same as in classical DB theory (transactions etc.).

Chapter 5

Query Languages: XPath

- Network Data Model: no query language
- SQL – only for a flat data model, but a “nice” language
(easy to learn, descriptive, relational algebra as foundation, clean theory, optimizations)
- OQL: SQL with object-orientation and path expressions
- Lorel (OEM): extension of OQL
- F-Logic: navigation in a graph by path expressions with additional conditions
descriptive, complex.

REQUIREMENTS ON AN XML QUERY LANGUAGE

- suitable both for databases and for documents
- declarative: binding variables and using them
 - rule-based, or
 - SQL-style clause-based (which is in fact only syntactic sugar)
- binding variables in the rule body/selection clause:
suitable for complex objects
 - navigation by path expressions, or
 - patterns
- generation of structure in the rule head/generating clause

EVOLUTION OF XPATH

- when defining a query language, constructs are needed for addressing and accessing individual elements/attributes or sets of elements/attributes.
- based on this *addressing mechanism*, a clause-based language is defined.

Early times of XML (1998)

different navigation formalisms of that kind:

- XSL Patterns (inside the stylesheet language)
- XQL (XML Query Language)
- XPointer (referencing of nodes/areas in an XML document)

used all the same basic idea with slight differences in the details:

- paths in UNIX notation
- conditions on the path

`/mondial/country[@car_code="D"]/city[population > 100000]/name`

5.1 XPath – the Basics

1999: specification of the navigation formalism as *W3C XPath*.

- Base: UNIX directory notation
 - in a UNIX directory tree: `/home/dbis/Mondial/mondial.xml`
 - in an XML tree: `/mondial/country/city/name`

Straightforward extension of the URL specification:

<http://.../dbis/Mondial/mondial.xml#mondial/country/city/name> [XPointer until 2002]

[http://.../dbis/Mondial/mondial.xml#xpointer\(mondial/country/city/name\)](http://.../dbis/Mondial/mondial.xml#xpointer(mondial/country/city/name)) [XPointer now]

- W3C: XML Path Language (XPath), Version 1.0 (W3C Recommendation 16. 11. 1999)
<http://www.w3.org/TR/xpath>
- W3C: XPath 2.0 and XQuery 1.0 (W3C Recommendation 23. 1. 2007)
<http://www.w3.org/TR/xquery>
- Tools: see Web page
 - XML (XQuery) database system “eXist”
 - lightweight tool “saxonXQ” (XQuery)

XPATH: NAVIGATION, SIMPLE EXAMPLES

XPath is based on the UNIX directory notation:

- `/mondial/country`
addresses all country elements in MONDIAL,
the result is a set of elements of the form
`<country code="..."> ... </country>`
- `/mondial/country/city`
addresses all city elements, that are direct subelements of country elements.
- `/mondial/country//city`
addresses all city elements that are subelements (in any depth) of country elements.
- `//city`
addresses all city elements in the current document.
- wildcards for element names:
`/mondial/country/*/city`
addresses all city elements that are grandchildren of country elements
(different from `/mondial/country//city` !)

... and now systematically:

XPATH: ACCESS PATHS IN XML DOCUMENTS

- Navigation paths

/step/step/.../step

are composed by individual navigation steps,

- the result of each step is a set of nodes, that serve as input for the next step.
- each step consists of

*axis::nodetest[condition]**

- an axis (optional),
 - a test on the type and the name of the nodes,
 - (optional) predicates that are evaluated for the current node.
- paths are combined by the “/”-operator
 - additionally, there are function applications
 - the result of each XPath expression is a *sequence* of nodes or literals.

XPATH: AXES

Starting with a *current node* it is possible to navigate in an XML tree to several “directions” (cf. xmllint’s “cd”-command).

In each navigation step

path/axis::nodetest[condition]/path

the *axis* specifies in which direction the navigation takes place. Given the set of nodes that is addressed by *path*, for *each* node, the step is evaluated.

- Default: child axis: *child::country* \equiv *country*.
- Descendant axis: all sub-, subsub-, ... elements:

country/ancestor::city

selects all city elements, that are contained (in arbitrary depth) in a country element.

Note: *path//city* actually also addresses all these city elements, but “//” is *not* the exact abbreviation for “/descendant::” (see later).

XPATH: AXES

... another important axis:

- attribute axis:

`attribute::car_code` \equiv `@car_code`

wildcard for attributes: `attribute::*` selects all attributes of the current context node.

- and a less important:

self axis: `self::city` \equiv `./city`

selects the current element, *if* it is of the element type city.

for the above-mentioned axes there are the presented abbreviations. This is important for *XSL patterns* (see Slide 314):

XSL (match) patterns are those XPath expressions, that are built *without* the use of “axis:” (the abbreviations are allowed).

XPATH: AXES

Additionally, there are axes that do not have an abbreviation:

- parent axis: `//city[name="Berlin"]/parent::country`
selects the parent element of the city element that represents Berlin, *if* this is of the element type country.
(*only* the parent element, not all ancestors!)
- ancestor: all ancestors:
`//city[name="Berlin"]/ancestor::country` selects all country elements that are ancestors of the city element that represents Berlin (which results in the Germany element).
- siblings: `following-sibling::...`, `preceding-sibling::...`
for selecting nodes on the same level (especially in ordered documents).
- straightforward: “descendant-or-self” and “ancestor-or-self”.
Note: The popular short form `country//city` is defined as `country/descendant-or-self::node()/city`.
This makes a difference only in case of *context functions* (see Slide 219).

XPATH: AXES FOR USE IN DOCUMENT-ORIENTED XML

- following: all nodes after the context node in document order, excluding any descendants and excluding attribute nodes
- preceding: all nodes that are before the context node in document order, excluding any ancestors and excluding attribute nodes and namespace nodes

Note: For each element node x , the ancestor, descendant, following, preceding and self axes *partition* a document (ignoring attribute nodes): they do not overlap and together they contain all the nodes in the document.

Example:

Hamlet: what is the next speech of Lord Polonius after Hamlet said “To be, or not to be”?
(note: this can be in a subsequent scene or even act)

Exercise:

Provide equivalent characterizations of “following” and “preceding”

- i) in terms of “preorder” and “postorder”,
- ii) in terms of other axes.

XPATH: NODETEST

- The *nodetest* constrains the node type and/or the names of the selected nodes
- “*” as wildcard: `//city[name="Berlin"]/child::*`
returns all children.
- test if something is a node: `//city[name="Berlin"]/descendant::node()`
returns all descendant nodes.
- test if something is a node: `//city[name="Berlin"]/descendant::element()`
returns all descendant *elements* (note: not the text nodes).
- test if something is a text node: `//city[name="Berlin"]/descendant::text()`
returns all descendant *text nodes*.
`//city[name="Berlin"]/population/text()`
returns the text contents of the population element.
- test for a given element name:
`//country[name="Germany"]/descendant::element(population)`
or short form:
`//country[name="Germany"]/descendant::population`
returns all descendant *population elements*.

XPATH: TESTS

In each step

path/axis::nodetest[condition]/path

condition is a predicate over XPath expressions.

- The expression selects only those nodes from the result of *path/axis::nodetest* that satisfy *condition*. *condition* contains XPath expressions that are evaluated relative to the current *context node* of the respective step.

`//country[@car_code="D"]`

returns the country element whose `car_code` attribute has the value "D"

- When comparing an element with something, the `text()` method is applied implicitly:

`//country[name = "Germany"]` is equivalent to

`//country[name/text() = "Germany"]`

- If the right hand side of the comparison is a number, the comparison is automatically evaluated on numbers:

`//country[population > 1000000]`

XPATH: TESTS (CONT'D)

- boolean connectives “and” and “or” in *condition*:

`//country[population > 100000000 and @area > 5000000]`

`//country[population > 100000000 or @area > 5000000]`

- boolean “not” is a *function*:

`//country[not (population > 100000000)]`

- XPath expressions in *condition* have existential semantics:

The *truth value* associated with an XPath expression is *true*, if its result set is non-empty:

`//country[inflation]`

selects those countries that have a subelement of type inflation.

⇒ formal semantics: a path expression has

- a semantics as a result set, and
- a truth value!

XPATH: TESTS (CONT'D)

- XPath expressions in *condition* are not only “simple properties of an object”, but are path expressions that are evaluated wrt. the current context node:

`//city[population/@year='95']/name`

- Such comparisons also have existential semantics:

`//country[//city/name='Cordoba']/name`

returns the names of all countries, in which a city with name Cordoba is located.

`//country[not (//city/name='Cordoba')]/name`

returns the names of those countries where no city with name Cordoba is located.

Remark:

Note that `descendant::city` (relative) and `//city` (absolute) have different effect:

`//country[//city/name='Cordoba']/name`

returns the names of *all* countries (the filter just checks if there is *some* city with name Cordoba in the document).

XPATH: EVALUATION STRATEGY

- Input for each navigation step: A set of nodes (*context*)
- each of these nodes is considered separately for evaluation of the current step
- and returns zero or more nodes as (intermediate) result.
This intermediate result serves as context for the next step.
- finally, all partial results are collected and returned.

Example

- conditions can be applied to multiple steps

```
//country[population > 10000000]  
  //city[@is_capital and population > 1000000]  
    /name/text()
```

returns the names of all cities that have more than 1,000,000 inhabitants and that are the capital of a country that has more than 10,000,000 inhabitants.

ABSOLUTE AND RELATIVE PATHS

So far, conditions were always evaluated only “local” to the current element on the main navigation path.

- Paths that start with a name are *relative* paths that are evaluated against the current context node (used in conditions):

```
//city[name = “Berlin”]
```

- Semijoins: comparison with results of independent “subqueries”:
Paths that start with “/” or “//” are absolute paths:

```
//country[population > //country[@car_code='B']/population]/name
```

returns the names of all countries that have more inhabitants than Belgium

- conflict between “//” for absolute paths and for descendant axis:

```
//country[//city/name=“Berlin”]
```

(equivalent: `//country[descendant::city/name=“Berlin”]`)

can be used for starting a relative path.

XPATH: FUNCTIONS

Input: a node/value or a set of nodes/values.

Result: in most cases a value; sometimes one or more nodes.

- dereferencing (see Slide 209)
- access to text value and node name (see Slide 212)
- aggregate functions `count(node_set)`, `sum(node_set)`

`count(/mondial/country)`

returns the number of countries.

- context functions (see Slide 219)
- access to documents on the Web:

`doc("file or url")/path`

`doc('http://www.dbis.informatik.uni-goettingen.de/index.html')//text()`

(for querying external HTML documents, consider use of namespaces as described on Slide 230 - nodetests work only with namespace!)

- see W3C document *XPath/XQuery Functions and Operators*

IDREF ATTRIBUTES

- ID/IDREF attributes serve for expressing cross-references
- SQL-style: references can be resolved by semi-joins:
(similar to foreign keys in SQL)

`//city[@id = //organization[abbrev="EU"]/@headq]`

SQL equivalent (uncorrelated subquery):

```
SELECT *
FROM city
WHERE (name, country, province) IN
      (SELECT city, country, province
       FROM organization
       WHERE abbrev = 'EU')
```

... not a really elegant way in a graph-based data model ...

XPATH: DEREFERENCING

Access via “keys”/identifiers

The function `id(string*`) returns all elements (of the current document) whose id’s are enumerated in *string**

- `id(“D”)` selects the element that represents Germany
(`country/@car_code` is declared as ID)
- `id(//country[car_code=“D”]/@capital)`
yields the element node of type `city` that represents Berlin.

This notation is hard to read if multiple dereferencing is applied, e.g.

```
id(id( id(//organization[abbrev='IOC']/@headq)/@country)/@capital)/name
```

Alternative syntaxes:

```
//organization[abbrev='IOC']/id(@headq)/id(@country)/id(@capital)/name
```

```
//organization[abbrev='IOC']/@headq/id(./@country/id(./@capital/id(./name
```

XPath: Dereferencing (Cont'd)

Analogously for multi-valued reference attributes (IDREFS):

- `//country[@car_code="D"]/@memberships`
returns "org-EU org-NATO ..."
- `id(//country[@car_code="D"]/@memberships)`
`//country[@car_code="D"]/id(@memberships)`
returns the set of all elements that represent an organisation where Germany is a member.
- `id(//organization[abbrev="EU"]/members/@country)`
`//organization[abbrev="EU"]/members/id(@country)`
returns all countries that are members (of some kind) in the EU.

Aside: Dereferencing by Navigation [Currently not supported]

Syntax:

attribute::nodetest⇒*elementtype*

Examples:

- `//country[car_code="D"]/@capital⇒city/name`
yields the element node of type city that represents Berlin.
- `//country[car_code="D"]/@memberships⇒organization`
yields elements of type organization.
- Remark: this syntax is not supported by all XPath Working Drafts:
 - XPath 1.0: no
 - has originally be introduced by Quilt (2000; predecessor of XQuery)
 - XPath 2.0: early drafts yes, later no
 - announced to be re-introduced later ...

XPATH: STRING() FUNCTION

The *function* `string()` returns the string value of a node:

- straightforward for elements with text-only contents:

`string(//country[name='Germany']/population)`

Note: for these (and only for these!) nodes, `text()` and `string()` have the same semantics.

- for attributes: `//country[name='Germany']/string(@area)`

Note: an attribute node is a name-value pair, not only a string (will be illustrated when constructing elements later in XQuery)!

free-standing attribute nodes as result cannot be printed!

- the `string()` function can also be appended to a path; then the argument is each of the context nodes: `//country[name='Germany']//name/string()`

- the string value of a subtree is the concatenation of all its text nodes:

`//country[@name='Germany']/string()`

Note: compare with `//country[@name='Germany']//text()` which lists all text nodes.

- `string()` **cannot** be applied to node sequences: `string(//country[name='Germany']//name)` results in an error message.

(see W3C XPath and XQuery Functions and Operators).

XPATH: SOME MORE DETAILS ON COMPARISONS

- in the above examples, all predicate expressions like `[name="Berlin"]` or `[@car_code="D"]` always *implicitly* compare the string value of nodes, e.g., here the string values of `<name>Berlin</name>` or `attribute: (car_code, "D")`.

Usage of Numbers

- comparisons using `>` and `<` and a number literal given in the query implicitly cast the string values as *numeric* values.

```
//city[population > 200000]
```

returns the all cities with a population higher than 200,000.

```
//city[population > '200000']
```

returns the all cities with a population *alphabetically* "bigger" than 200,000, e.g., 3500, but not 1,000,000!

```
//city[population > //city[name="Munich"]/population]
```

does *not* recognize that numerical values are meant:

All cities with population lexically bigger than "1244676" are returned.

```
//city[population > //city[name="Munich"]/population/number()]
```

It is sufficient to apply the `number()` casting function (see later) to one of the operands.

XPATH: COMPARISON BETWEEN NODES

Usage of Node Identity

- as seen above, the “=” predicate uses the string values of nodes.

In most cases, this is implicitly correct:

Consider the following query: “Give all countries whose capital is the headquarter of an organization”:

```
//country[id(@capital)=//organization/id(@headq)]/name
```

Compares the overall string values of city elements, e.g., “Brussels 4.35 50.8 951580”.

- but for empty nodes, the result is not as intended ...

Comparison by Node Identity: “a is b”

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE mondial-simple SYSTEM "mondial-simple.dtd">
<mondial-simple>
  <country car_code="D" capital="Berlin"/>    <city name="Berlin"/>
  <country car_code="B" capital="Brussels"/>  <city name="Brussels"/>
  <organization name="EU" headq="Brussels"/>
</mondial-simple>                                [Filename: XPath/node-comparison.xml]
```

- the query `//country[id(@capital)=//organization/id(@headq)]/string(@car_code)` yields “D” and “B”.
- Comparison by *node identity* is done by “is”:
`//country[id(@capital) is //organization/id(@headq)]/string(@car_code)`
 - “is” is only provided since XPath 2.0
 - “is” allows only one node as argument, not a node sequence
(\Rightarrow XQuery: not something bound by “let \$x := *node sequence*”)
- Aside: “deep equality” of nodes can be tested with the predicate `deep-equal(x, y)`.
(by this, two subtrees can be checked to have the same structure+contents)

XPATH: PREDICATES AND OPERATIONS ON STRINGS

- `concat(string, string, string*)`
- `startswith(string, string)`
`//city[starts-with(name,'St.')] /name`
- `contains(string, string)`
`//city[contains(name,'bla')] /name`
- `substring-before(string, string, int?)`
- `substring-after(string, string, int?)`
- `substring(string, int, int)`: the substring consisting of i_2 characters starting with the i_1 th position.

XPATH: NAME FUNCTION

- the function `name()` returns the element name of the current node:
 - `name(//country[@car_code='D'])` or `//country[@car_code='D']/name()`
 - `//*[name='Monaco' and not (name()='country')]` yields only the city element for Monaco.

XPATH: IDREF FUNCTION

- the function `idref(string*)` returns all nodes that have an IDREF value that refers to one of the given strings (note that the results are attribute nodes):
`idref('D')/parent::* /name` yields the name elements of all “things” that reference Germany.

FUNCTIONS ON NODESETS

- Aggregation: `count(nodeset)`, `sum(nodeset)`, analogously `min`, `max`, `avg`
`sum(//country[encompassed/id(@continent)/name="Europe"]/population)`
`count(//country)`
all numeric functions implicitly cast to numeric values (double).
- removal of duplicates:
 - recall that the XPath strategy works on *sets of nodes* in each step - duplicate *nodes* are automatically removed:
`//country/encompassed/id(@continent)/name`
 - function `distinct-values(nodeset)`:
takes the *string values* of the nodes and removes duplicates:
`doc('hamlet.xml')//SPEAKER`
returns lots of `<SPEAKER>...</SPEAKER>` *nodes*.
`distinct-values(doc('hamlet.xml')//SPEAKER)`
returns only the different (text) *values*.
- and many more (see W3C XPath/XQuery Functions and Operators).

XPATH: CONTEXT FUNCTIONS

- All functions retain the order of elements from the XML document (document order).
- the `position()` function yields the position of the current node in the current result set.

`/mondial/country[position()=6]`

Abbreviation: `[x]` instead of `[position()=x]`; `[-1]` yields the last node:

`/mondial/country[population > 1000000][6]`

selects the 6th country that has more than 1,000,000 inhabitants (in document order, not the one with the 6th highest population!)

`/mondial/country[6][population > 1000000]`

selects the 6th country, if it has more than 1,000,000 inhabitants.

- the `last()` function returns the position of the last elements of the current sub-results, i.e., the size of the result.

`//country[position()=last()]`

XPATH: CONTEXT FUNCTIONS (CONT'D)

- consider again the “//” abbreviation (cf. Slide 199):
 - `/mondial/descendant::city[18]` selects the 18th city in the document,
 - `/mondial/descendant-or-self::node()/city[18]` selects each city which is the 18th child of its parent (country or province).
(note that some implementations are buggy in this point ...)
- Example queries against `mondial.xml` and `hamlet.xml`.

XPATH: FORWARD- AND BACKWARD AXES

- the result of each query is a *sequence of nodes*
- document order (and final results): forward
- context functions: forward or backward
- all axes enumerate results starting from the current node.
 - forward axes: child, descendant, following, following-sibling
 - backward axes: ancestor, preceding, preceding-sibling
 - `//table/preceding-sibling::h4//text()`
selects all preceding h4 elements (section headers).
The result is -as always- output in document order
 - `//table/preceding-sibling::h4[1]//text()`
selects the last preceding section header (context function on backward axis)
 - undirected: self, parent, attribute (and namespace)
- only relevant for queries against document-oriented XML.

EXTENSIONS WITH XPATH 2.0

- further string- and aggregate functions
- more complex path constructs (alternatives, parentheses)
`(//city|//country)[name='Monaco']`
`/mondial/country/(city|(province/city))/name`
- extended subscript operator:
`//country[population > 1000000][-3]`
`//country[population > 1000000][5-10]`
`//country[population > 1000000][1,5-10,-3]`
- ANY and ALL semantics for *condition*:
`//country[ALL city/population > 1000000]`
`//country[ANY city/population > 1000000]`
(countries where all/at least one city has more than 1000000 inhabitants)
- extending the language to more than usual navigation ...
- alignment of the whole XML world (XPath, XQuery) with datatypes (data model and XML Schema)

5.2 Aside: Namespaces

The names in an XML instance (i.e., tag names and the attribute names) actually consist of two parts:

- localpart + namespace (which can be empty, as in the previous examples)

Use of Namespaces

- a namespace is similar to a language: defining a set of names and sometimes having a DTD (if intended as an XML vocabulary).
- e.g. “mondial:city”, “bib:book”, “xhtml:tr” “dc:author”, “xsl:template” etc.
- used for distinguishing coinciding element names in different application areas.
- each namespace is associated with a URI (which can be a “real” URL), and abbreviated by a *namespace prefix* in the document.
- e.g., associate the namespace prefix `xhtml` with url `http://www.w3.org/1999/xhtml`. these things will become clearer when investigating the RDF, RDFS, and Semantic Web Data Models.

USAGE OF NAMESPACES IN XML DOCUMENTS

- each element can have (or can be in the scope of) multiple *namespace declarations* (represented by a node in the data model, similar to an attribute node).
- namespace declarations are inherited to subelements
- the element/tag name and the attribute names can then use one of the declared namespaces.

By that, every element can have one *primary namespace* and “knows” several others.

Alternatives:

1. node has no namespace (e.g. mondial),
 2. document declares a default namespace (for all elements (not the attributes!) that do not get an explicit one (often in XHTML pages)),
 3. elements have an explicit namespace (multiple namespaces allowed in a document; e.g. an XSL document that operates with XHTML markup and “mondial:” nodes).
- (2) and (3) are semantically equivalent.

... see next slides.

EXPLICIT NAMESPACE IN AN XML DOCUMENT

```
<xh:html xmlns:xh="http://www.w3.org/1999/xhtml">
  <xh:body>
    <xh:h3>Header</xh:h3>
    <xh:a href="http://www.informatik.uni-goettingen.de">IFI</xh:a>
  </xh:body>
</xh:html>
```

[Filename: XML-DTD/xhtml-expl-namespace.xml]

- Note: attribute is not in the HTML namespace!

This is actually already not XPath, but a simple XQuery query:

```
declare namespace ht = "http://www.w3.org/1999/xhtml";
/ht:html//ht:a/string(@href)
```

[Filename: XPath/xhtml-query.xq]

- Note: the namespace *must* be used in the query, i.e., “ht:html” is different from just “html”
- more accurate, it means something like `<{http://www.w3.org/1999/xhtml}html>...</...>` since not the chosen namespace prefix matters, but only the URI assigned to it.

TWO EXPLICIT NAMESPACES IN AN XML DOCUMENT

- “Dublin Core” defines a vocabulary for metadata description of resources (here: of XML documents); cf. <http://dublincore.org/documents/dces/>

```
<xh:html xmlns:xh="http://www.w3.org/1999/xhtml"
         xmlns:dc="http://purl.org/dc/elements/1.1/">
  <dc:creator>John Doe</dc:creator>
  <dc:date>1.1.2000</dc:date>
  <xh:body> ... </xh:body> </xh:html>
```

[Filename: XML-DTD/xhtml-expl-namespaces.xml]

```
declare namespace ht = "http://www.w3.org/1999/xhtml";
declare namespace dc = "http://purl.org/dc/elements/1.1/";
/ht:html//dc:creator/text()
```

[Filename: XPath/xhtml-dc-query.xq]

- the document is *not* valid wrt. the XHTML DTD since it contains additional “alien” elements.
(combination of languages is a problem in XML – this is better solved in RDF/RDFS)
- in RDF, dc:creator from above expands to the URI <http://purl.org/dc/elements/1.1/creator>.

DEFAULT NAMESPACES IN AN XML DOCUMENT

- a Default Namespace can be assigned to an element (and inherited to all its subelements where it is not overwritten):

```
<html xmlns="http://www.w3.org/1999/xhtml"
      xmlns:dc="http://purl.org/dc/elements/1.1/">
  <dc:creator>John Doe</dc:creator>
  <date xmlns="http://purl.org/dc/elements/1.1/">1.1.2000</date>
  <body> ... </body> </html>
```

[Filename: XML-DTD/xhtml-def-namespaces.xml]

```
declare namespace ht = "http://www.w3.org/1999/xhtml";
declare namespace dc = "http://purl.org/dc/elements/1.1/";
/ht:html/dc:date/text()
```

[Filename: XPath/xhtml-dc-def-query.xq]

NAMESPACES AND ATTRIBUTES

- Namespaces are *not* inherited to attributes in any case. If an attribute should be associated with a namespace, this *must* be done explicitly:

```
<ht:html xmlns:ht="http://www.w3.org/1999/xhtml">
  <ht:body>
    <ht:a href="1+" ht:href="2-">IFI</ht:a>
    <x:a xmlns:x="http://www.w3.org/1999/xhtml" href="3+" x:href="4-">IFI</x:a>
    <a xmlns="http://www.w3.org/1999/xhtml" href="5+" ht:href="6-">IFI</a>
  </ht:body> </ht:html>
```

[Filename: XML-DTD/namespaces-attr.xml]

```
declare namespace ht = "http://www.w3.org/1999/xhtml";
/ht:html//ht:a/@href/string()
```

[Filename: XPath/namespaces-attr-query.xq]

- the “HTML-correct” attributes “1+”, “3+”, and “5+” are returned,
- the query `/ht:html//ht:a/@href/string()` returns the “wrong” attributes “2-”, “4-”, and “6-”.

DECLARING NAMESPACES IN THE DTD DOCUMENT

- introduce default namespace in the DTD as attribute of the root element (e.g. in XHTML):

```
<!ELEMENT html (head, body)>
<!ATTLIST html
  xmlns %URI; #FIXED 'http://www.w3.org/1999/xhtml' >
```

- XHTML instance:

```
<html xmlns="http://www.w3.org/1999/xhtml"> <body> ... </body></html>
```

- introduce explicit namespaces as attribute of the root element (e.g. in XHTML):

```
<!ELEMENT html (head, body)>
<!ATTLIST html xmlns:xh %URI; #FIXED 'http://www.w3.org/1999/xhtml' >
```

This is used with RDF/XML in the Semantic Web

EXAMPLE: QUERYING XHTML IN PRESENCE OF NAMESPACES

XHTML DTD at <http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd> contains:

```
<!ELEMENT html (head, body)>
<!ATTLIST html id ID #IMPLIED
              xmlns %URI; #FIXED 'http://www.w3.org/1999/xhtml'>
```

Sample XHTML files:

- DBIS Web pages:

```
declare namespace h = "http://www.w3.org/1999/xhtml";
doc('http://www.dbis.informatik.uni-goettingen.de/')//h:li/h:a/@href/string()
```

[Filename: XPath/web-queries.xq]

- DBIS WWW2002 paper: in the local exist at /db/xmlcourse/xlink.htm

```
declare namespace ht = "http://www.w3.org/1999/xhtml";
doc('/db/xmlcourse/xlink.htm')//ht:h1
```

[Filename: XPath/exist-xhtml-query.xq]

DECLARING A DEFAULT NAMESPACE IN XQUERY

XQuery allows to declare default namespaces for elements and for functions:

- are then added to each element and function step, respectively;
- not for attributes (recall that namespaces from elements are not inherited to attributes).
(cf. Slide 228)

```
declare default element namespace "http://www.w3.org/1999/xhtml";  
/html//a/@href/string()
```

[Filename: XPath/namespaces-default-query.xq]

- the “HTML-correct” attributes “1+”, “3+”, and “5+” are returned,
- the equivalent query is `/h:html//h:a/@href/string()`.

5.3 XPath: The Limits

- addressing only sets of nodes
- not “give all pairs of ...”
- the highest mountain in Africa:

```
doc('mondial.xml')//mountain[
  id(id(located/@country)/encompassed/@continent)/name='Africa'
  and
  not (height <
    //mountain[
      id(id(located/@country)/encompassed/@continent)/name='Africa']/height)]
/name
```

[Filename: XPath/highestmountain.xq]

... comparison only by semijoins in the condition.

- for *each* continent, give the highest mountain?
not possible: two properties of the same object (height, continent) must be compared independently → **requires variable binding**

5.4 XPath: Conclusion

What can XPath do?

Comparison with relational operators

- selection: yes (selection of values and of (sub)structures)
- projection/reduction: no. Only complete nodes can be selected
- join/combination: no. Only semi-joins can be expressed in the conditions

Other functionality:

- correlated subqueries: inside the conditions as semijoins
 - restructuring of the results: no
 - only following a “main path” for navigating to nodes (including semijoins)
- ⇒ only a fragment of a query language for addressing nodes.
- compared with SQL, XPath is only a unary “FROM” clause!
 - XQL (Software AG, 1998/1999) for some time followed (as one of the predecessors of XPath) an approach to add join variables and constructs for projection and restructuring/grouping to the path language.

IMPORTANCE OF XPATH IN THE XML-WORLD

- addressing mechanism for nodes in XML documents
- navigation in the tree structure
- serves as base for different concepts:
 - XQuery
 - XSL/XSLT: stylesheets, transformation language
 - other query languages
 - XML Schema
 - XPointer/XLink

Chapter 6

XML Query Languages

- XPath is not a query language:
selects only sets of nodes
- additional functionality of query languages:
 - composition of tuples/structures from several nodes of a path
 - joins
 - dereferencing
 - * via joins
 - * via direct resolving of IDs (seen as values)
 - * via dereferencing of ID attributes
 - aggregations
 - formatting and restructuring of results
 - operations on the order of nodes!

XML QUERY LANGUAGES

Collected experiences from SQL, OQL, OEM/WSL/MSL, F-Logic and some more ...

- predecessors of XPath: XSL Patterns/XPointer/XQL (1998)
- XQL extended the early “basic form” to a query language
 - adding several constructs to the path expressions
 - increasingly complicated
 - still not sufficiently expressive
 - showed the limits and requirements
- XML-QL (1998): pattern-matching-based “extraction language”
 - not path-based, but XML-pattern/template-based binding of variables
 - semantics by a clause-construct
 - generation and structuring of the result by an XML pattern with variables

XML QUERY LANGUAGES (CONT'D)

- Quilt (2000): SQL-style extension of XPath
 - binding of variables by XPath expressions
 - nested loops by “for”-clauses
 - additional conditions in a “where”-clause
 - structuring of the result by a “return”-clause
- XQuery (2001): “official” version of Quilt
 - W3C Working Draft XQuery first version from 15 February 2001
 - XQuery 1.0: W3C Recommendation since 23.1.2007
 - <http://www.w3.org/TR/xquery/>

6.1 XQL

XQL (XML Query Language; 1998) is a simple query language based on early constructs of XPath:

- all XPath expressions that can be expressed without the use of “axis::” (cf. Slide 198 - axes have only been added later).
- text() was a function,
- function applications have been expressed by “!” at the end of the path expression:
`//country/name!text()`

Further querying functionality was integrated syntactically into the path expressions.

XQL: BOOLEAN OPERATIONS AND SET OPERATIONS

- q_1 union q_2 , $q_1 \mid q_2$
- q_1 intersect q_2
- $q_1 \sim q_2$ (union, in case that both are non-empty)
- q_1 or q_2
- q_1 and q_2

XQL: RETURN OPERATORS (PROJECTIONS ON THE PATH)

Operators that output the node that is addressed at the given position:

?: the complete node is added to the output structure (including attributes and subelements)

?: only the element "hull" is added to the output

- `country/city[@isCountryCap]/name`

```
<name>Berlin</name>
```

```
<name>Rome</name>
```

- `country?/city[@isCountryCap]/name`

```
<country> <name>Berlin</name> </country>
```

```
<country> <name>Rome</name> </country>
```

- `country?[@car_code?]/city[@isCountryCap]/name`

```
<country car_code="D"> <name>Berlin</name> </country>
```

```
<country car_code="I"> <name>Rome</name> </country>
```

- `country?[@car_code?]/city?[@isCountryCap]/name!text()`

```
<country car_code="D"> <city>Berlin</city> </country>
```

```
<country car_code="I"> <city>Rome</city> </country>
```

XQL: GROUPING

- copy a part of the original document structure:

path₁ { path₂ }

- without grouping:

`country?[@car_code?]/city?/name!text()`

```
<country car_code="D"> <city>Berlin</city> </country>
```

```
<country car_code="D"> <city>Hamburg</city> </country>
```

```
<country car_code="D"> <city>Munich</city> </country>
```

- with grouping:

`country?[@car_code?] {/city?/name!text()}`

```
<country car_code="D">
```

```
  <city>Berlin</city>
```

```
  <city>Hamburg</city>
```

```
  <city>Munich</city>
```

```
</country>
```

XPATH: SEMIJOINS ARE POSSIBLE

- Semi-joins via subqueries in the condition:

$$\pi[A](r \bowtie s), \quad A \subset \text{attr}(r)$$

Query: name of the continent where Germany is located:

```
/mondial/continent[@id =  
    /mondial/country[@car_code="D"]  
    /encompassed/@continent]  
/name!text()
```

Problems

- full joins with join conditions not possible
- no restructuring/generation of answer structure

XQL: JOINS

Asymmetric full joins expressed by *correlating variables* and “alternative”-construct:
Filters may contain variable assignments of the form

$[\$var := expr]$

that are then used in another condition

$[expr' = \$var]$

$//organization?[\$s := @headq] \{ name?? \mid abbrev?? \mid member?? \mid //city[@id=\$s]?? \}$

```
<organization>
  <name>European Union</name>
  <abbrev>EU</abbrev>
  <member type="member" country="GR F E A D I B L NL DK SF S IRL P GB"/>
  <member type="membership applicant"
    country="AL CZ H SK LV LT PL BG RO EW M CY"/>
  <city> <name>Brussels</name> ... </city>
</organization>
```

Equivalent:

$//organization?[\$s := @headq \text{ and } name?? \mid abbrev??] \{ member?? \mid //city[@id=\$s]?? \}$

XQL: CONCLUSION

- Ad-hoc-constructs (in different versions)
- insufficient restructuring functionality
 - tree structure of the input is in principle retained
- insufficient join functionality
- no clear semantics for the result format
- queries cannot be nested (cf. SQL, OQL: results are again relations); here is even no notion of a subquery
- one of the reasons: no variable concept
- implemented and used up to 2002 in the “Tamino” system of Software AG.

6.2 Query Languages: Requirements

Requirements on XML Query Languages [David Maier and W3C XML Query Requirements]

- closedness: output must be XML
- orthogonality/composability: everywhere where a set of XML elements is required, also a query is allowed.
- clean definition and nesting of operations: selection, extraction/projection, restructuring, combination/join, fusion of elements,
- applicable without presence of schema, but can use a schema,
- retaining the order of nodes,
- [queries should have an XML representation, especially, XML documents should be able to contain embedded queries]
- resolving of XPointer and XLink
- formal semantics: deriving structure of the result, equivalence and query containment

6.3 XML-QL

- <http://www.w3.org/TR/1998/NOTE-xml-ql-19980819/>

- simple, pattern-based XML query language:

WHERE *xml-pattern* IN *url* CONSTRUCT *result*

- usage of variable bindings:

xml-pattern contains variables that can be used in *result*,

- declarative,
- “relationally complete”, i.e., joins can be expressed.

Example:

```
WHERE <country car_code=$id>
      <name>$name</name>
      </country>
IN "http://www.../mondial.xml"
CONSTRUCT <country car_code=$id name=$name/>
```

XML-QL: JOINS

Joins are expressed as a list of

WHERE (*expr*₁ IN *doc*₁, ... , *expr*_{*n*} IN *doc*_{*n*})-clauses:

- equijoin inside a document:

WHERE

```
<country car_code=$c></country>
```

```
IN mondial.xml,
```

```
<organization abbrev=$org>
```

```
  <members type=$type country=$c/>
```

```
</organization>
```

```
IN mondial.xml,
```

```
CONSTRUCT ...
```

XML-QL: JOINS

- Joins that combine multiple documents:

WHERE

```
<city name=$name1>  
IN http://www.../europe.xml,  
<city name=$name2>  
IN http://www.../america.xml,  
<connection from=$name1 to=$name2>  
IN http://www.../lufthansa.xml
```

CONSTRUCT

```
<connection>  
  <from continent="europe" city="name1"/>  
  <to continent="america" city="name2"/>  
</connection>
```

XML-QL: NESTED QUERIES

WHERE *xml-pattern* IN *url* CONSTRUCT *result*

- *result* can contain nested WHERE ... IN ... CONSTRUCT statements.

FURTHER FUNCTIONALITY

- tag-variables: WHERE <\$tag> ... </>
- regular path expressions: instead of XPath's "//", <*> ... </> is used.

DATA MODEL

- Graph-based: XML-tree with IDREF edges:

WHERE

```
<country car_code=$cc>
```

```
  <capital><name>$name</name><population>$pop</population></capital>
```

```
</country>
```

IN ... CONSTRUCT ...

XML-QL: CONCLUSION

- clause-based high-level language
- selection and construction pattern-based (by binding variables in the patterns; similar to Logic Programming)
- join conditions: not in a WHERE clause, but implicitly expressed by *join variables* (like in Logic Programming)
- graph data model; no difference between tree edges and reference edges
- has been implemented
- used in different projects (e.g. MIX – Mediation in XML; UC San Diego 1999/2000)
 - allows for access and combination of different HTML/XML-sources in a query.

6.4 SQL, OQL etc.

- set-oriented (sets of tuples or objects) language
- implicit iteration over sets:
SELECT ... FROM *relation-or-extent* **c**
- variable **c** ranges over *data items*
- join: use several such variables and correlate them
- WHERE and SELECT part: use these variables

Similar constructs for XML?

- variables range over sets of nodes
- ... sets of nodes can be addressed by XPath
- straightforward and intuitive:

for **\$c** in *//country*

where **\$c**/population > 1000000

return **\$c**/name/text()

6.5 Quilt and XQuery

- a Quilt is a “Flickenteppich” ...
- IBM, Software AG, INRIA; literature: WebDB2000-Workshop
- Structure similar to [SQL/OQL: clause-based, functional language](#) (arbitrary nesting of FLWR expressions allowed),
- Use of variables similar to SQL/OQL,
- based upon [XPath \(previously XQL/XSL Patterns\)](#) in the [selection part](#) and upon [XML-QL](#) (XML patterns) in the [construction part](#):

- For **L**et **W**here **R**eturn-clauses

for variable in xpath-expr // from XQL/XPath and XML-QL

let additional_variable := xpath-expr

where condition

return xml-expr // from XML-QL

- has been moved into W3C’s “XML Query” in 2001 with only small changes.
- Remark: XQuery is case-sensitive.

ALL KEYWORDS MUST BE WRITTEN WITH **non-capital** LETTERS!

XQUERY: EXAMPLE

- for-clause: binding of variables (cf. SQL: FROM)
- where-clause: evaluation of conditions
- return-clause: generation of the result (cf. SQL: SELECT)

```
<result>
  { for $c in /mondial/country
    where $c/population > 10000000
    return <bigcountry>
      { $c/name }
      <area> { string($c/@area) } </area>
    </bigcountry>
  } </result>
```

[Filename: XQuery/first-example.xq]

generates

```
<result><bigcountry><name>France</name><area>547030</area></bigcountry>
  <bigcountry><name>Spain</name><area>504750</area></bigcountry>
  :
</result>
```

ASIDE: TOOLS – XQUERY AS A DATABASE AND WEB QUERY LANGUAGE

XML Databases

- local repository of XML documents
- adding documents to the Database
- access only against locally stored documents
- presence of access paths like indexes etc
- manipulation of documents

Tool: a local eXist repository; see Web page

paths: [/mondial//country/name](#) or [doc\('/db/xmlcourse/mondial.xml'\)/mondial/country/name](#)

Queries against the Web

- querying the whole Web
- documents not locally stored; only on-the-fly-indexing possible
- access to remote documents by their url

Tool: saxonXQ; see Web page

paths: [doc\('filename or url'\)//country/name](#)

XQUERY: OVERVIEW OF FUNCTIONALITY

- for-clause: defines nested loops where each of the variables runs over the set of selected values
- variables in XPath expressions: bound in for/let (or by surrounding statements), they are used as starting points for paths and in conditions
- joins:
 - multiple variables in a for-clause:
for $\$var_1$ in $doc_1/path_1$, ..., $\$var_n$ in $doc_n/path_n$
 - correlated definition of the variables in the for-clause:
for $\$var_1$ in $doc_1/path_1$, $\$var_2$ in $\$var_1/path_2$, ...
- let-clause for definition of “constants”:
let $\$var := expr$
binds $\$var$ to the *whole result* of $expr$ (in general, a node set).
- nested/iterated for-let-for-let-clauses allowed
- generation of nested and grouped structures:
the return-clause may contain further FLWR-clauses (which can contain variables from the outer clause).

SIMPLEST XQUERY QUERIES: XPATH

- Each XPath query is also an XQuery query
result: a sequence of nodes or literal values

```
doc('mondial.xml')//country/name
```

Note: different behavior when returning attribute nodes!

```
doc('mondial.xml')//country/@area
```

XQUERY: FOR-CLAUSE

for $\$var = xpath\text{-}expr$

- iterates over the result of *xpath-expr*

```
for $c in /mondial//country/name  
return $c
```

[Filename: XQuery/for-example.xq]

XQUERY: RETURN-CLAUSE

Output of all statements must be XML.

- simple case: content of a variable

```
for $c in /mondial//country/name  
return $c
```

- and generation of structured results (cf. OQL)

Generation of Structures

- literal XML
- computed element- and attribute constructors (later)

Use of Computed Values/Structures

- enclosed between “{” ... “}”
- evaluation of variables and XPath expressions
- nested FLWR-clauses

RETURN-CLAUSE: CONSTRUCTION OF RESULT ELEMENTS

- literal XML, values of variables and results of XPath expressions

```
<html><table>
<tr><th>Name</th><th>Area</th><th>Population</th></tr>
{ for $c in /mondial/country
  return
    <tr><td>{$c/name/text()}</td>
      <td>{string($c/@area)}</td>
      <td>{$c/population/text()}</td>
    </tr>
}
</table></html>
```

[Filename: XQuery/table-example.xq]

returns one table row for each country.

XQUERY: FOR-CLAUSE

Multiple Variables in a For-Clause

- cartesian product
(cf. FROM-clause in SQL)

```
for $c in /mondial//country,  
    $o in /mondial//organization  
where $c/@capital = $o/@headq  
return  
    <answer>  
        <country>{$c/name/text()}</country>  
        <organization>{$o/name/text()}</organization>  
    </answer>
```

[Filename: XQuery/cartesian-example.xq]

- compare where clause with equivalent
where \$c/id(@capital) is \$o/id(@headq)
on node level (“=” would also be correct here, taking the string value of the nodes).

XQUERY: FOR-CLAUSE

Multiple Variables in a For-Clause

- “correlated” Join
(cf. FROM-clause in Schema-SQL and OQL)
- subset of the cartesian product

```
for $c in /mondial/country,  
    $p in $c/province  
return  
  <answer>  
    <country>{$c/name/text()}</country>  
    <prov>{$p/name/text()}</prov>  
  </answer>
```

[Filename: XQuery/correlated-join-example.xq]

RETURN-CLAUSE WITH NESTED FLWR-CLAUSE

- inner query used in the outer return-clause (cf. OQL)

```
for $c in /mondial/country
where $c/province
return
  <answer>
    {$c/name}
    { for $p in $c/province
      return
        <prov>{$p/name/text()}</prov>
    }
  </answer>
```

[Filename: XQuery/nested-flwr-example.xq]

generates for each country that has provinces an <answer> element that contains a <name> element and a sequence of <prov> elements.

LET-CLAUSE

`let $var := xpath-expr`

- does not iterate over the result of *xpath-expr*
- but binds the complete result of *xpath-expr* as sequence of nodes to the variable:

```
for $c in /mondial/country
let $cities := $c//city/name
return
  <country>
    {$c/name}
    {$cities}
  </country>
```

[Filename: XQuery/let-example.xq]

- useful for keeping intermediate results for reuse (often missed in SQL)

WHERE-CLAUSE: CONDITIONS

Similar to XPath's conditions (same predicates etc):

- logical “and” and “or”
- “not(...)” as a boolean function
- Comparisons: “is” for node identity, “<<” and “>>” for document order, “follows” and “precedes”
- Quantifiers: *where some|every \$var in expr satisfies condition*

```
for $c in /mondial/country
where some $city in $c//city satisfies $city/population > 1000000
return $c/name
```

```
for $c in /mondial/country
where every $city in $c//city satisfies $city/population > 1000000
return $c/name
```

[Filenames: XQuery/some-example.xq and every-example.xq]

USE CASE: JOIN BETWEEN DIFFERENT DOCUMENTS

- doc(...) function to access files (local or from the Web)
- here: join by a subquery

```
<result>
{ for $c in doc(concat('http://www.dbis.informatik.uni-goettingen.de',
                        '/Mondial/mondial-europe.xml'))/mondial/country
  where some $l in doc('hamlet.xml')//LINE
    satisfies contains($l, $c/name)
  return
    <country>
      {$c/name}
    </country>
}
</result>
```

[Filename: XQuery/join-web-documents.xq]

CONDITIONAL EVALUATION AND ALTERNATIVES

- if-then: alternative choice of subelements

if (expr) then expr else expr

```
<result>
  { for $c in /mondial/country
    return
      <country>
        {$c/name}
        {if ($c/province) then $c/province/city else $c/city}
      </country>
    }
</result>
```

[Filename: XQuery/if-else-example.xq]

ORDER OF RESULT SET

XPath: the result is *always* returned in *document order*:

- purely navigational access:

```
//country/city/name
```

- even when a backward axis is used during navigation, the nodes are enumerated in document order:

```
//country[name='Germany']/province[last()]/preceding-sibling::* /name
```

(backward axis is only relevant for context functions in immediate conditions)

- or when id-referencing is used:

```
id(//organization/@headq)/name
```

(note: cities are *not* ordered according to the order of the organizations!)

XQuery: result set is ordered according to for-clause:

```
for $c in //organization
return id($c/@headq)/name
```

let-clause: binds the result set according to the respective order.

SORTING

- order by: *expr order by (expr [ascending|descending])*

```
<result>
  { for $c in //country
    order by $c/name
    return $c/name }
</result>
```

[Filename: XQuery/orderby-example.xq]

- note that the interpreter must be told whether the values should be regarded as numbers or as strings (default: alphanumerical)

```
<result>
  { for $c in //country
    where $c/population > 0
    order by number($c/population)
    return $c/name  }
</result>
```

[Filename: XQuery/orderby-num-example.xq]

GROUPING AND AGGREGATION

- aggregate functions over result sets (avg, sum)
- bind variable with “for”-clause
- assign group with “let” (dependent on the current value in the for-clause) to a variable
- apply aggregate function to a nodeset

```
<result>
{ for $c in /mondial/country
  let $cities := $c//city
  where sum($cities/population) > 10000000
  return
    <answer>
      {$c/name}
      {sum($cities/population)}
    </answer>
}
</result>
```

[Filename: XQuery/aggr-1-example.xq]

AGGREGATION

- aggregation over result of a FLWR subquery
- bind (single) intermediate result by “let”

```
<result>
{ for $c in /mondial/country
  let $maxpop := max( for $citypop in $c//city/population/text()
                    return $citypop )
  return
    <answer>
      {$c/name}
      {$maxpop}
    </answer>
}
</result>
```

[Filename: XQuery/aggr-2-example.xq]

ATTRIBUTES IN THE RETURN-CLAUSE

- note that expressions the form “@bla” return *attribute nodes* - these are (AttrName,value)-pairs:

```
<result>
  {//country[name='Germany']/@car_code}
</result>
```

generates `<result car_code="D"/>`.

- attribute nodes are always added to the surrounding element.
- if only their value is needed, apply `string()`.

```
for $c in /mondial/country
return
  <country>
    {$c/@area}
    {string($c/@car_code)}
  </country>
```

Result:

```
<country area="28750">AL</country>
<country area="131940">GR</country>
:
```

[Filename: XQuery/attribute-example.xq]

COMPUTED ELEMENT- AND ATTRIBUTE NAMES

- explicit constructors
 - element *expr attrs-and-content*
the evaluation of *expr* yields the name of the element, the result of *attrs-and-content* is then inserted as attributes and content
Note: content is a node sequence, separated by “,”
 - attribute *expr expr-value*
the evaluation of *expr* yields the name of the attribute, *expr-value* yields its value.

```
<result>
  { for $c in doc('mondial.xml')//country
    where $c/encompassed
    return
      element { $c/@car_code }
              { attribute {$c/encompassed[1]/@continent} {"yes"},
                $c/name
              }
  } </result>
```

A result node:

```
<B europe="yes">
  <name>Belgium</name>
</B>
```

[Filename: XQuery/computed-constructors-example.xq]

COMPUTED ELEMENT- AND ATTRIBUTE NAMES: ANOTHER EXAMPLE

- the element content can be computed by an XQuery expression (cf. usage of *expr* on the previous slide):

```
<result>
  { for $c in doc('mondial.xml')//country
    where $c/encompassed
    return
      element { $c/@car_code }
        { for $e in $c/encompassed
          return attribute {string($e/@continent)} {"yes"},
            $c/name
        }
  } </result>
```

[Filename: XQuery/computed-constructors-example2.xq]

HANDLING DUPLICATES

- recall from XPath: results (and intermediate results) of XPath expressions are *node sets* in document order
⇒ for $\$x$ in *xpath-expr*, let $\$y := \text{xpath-expr}$
always results in a set (i.e., duplicates removed)
- recall Slide 218 for removal of duplicate *values*: `distinct-values(...)`

```
distinct-values(doc('...')//SPEAKER)
```

How many speeches has each of the speakers in “Hamlet”?

```
for $a in distinct-values(doc('/db/xmlcourse/hamlet.xml')//SPEAKER)
let $n := count(//SPEECH[SPEAKER = $a])
order by $n descending
return
  <answer>
    {$a}
    {$n}
  </answer>
```

[Filename: distinct-values.xq]

- takes only the string values (⇒ no further navigation applicable)

Handling Duplicates in XQuery(cont'd)

- FLWR expressions (e.g., for \$c in ... return \$c) do *not* eliminate duplicates automatically
- `for $o in //organization return $o/id(@headq)`
returns duplicates
- `distinct-values(for $o in //organization return $o/id(@headq))`
returns only the string values
- so it must be done programmatically (often, specific for the given problem: iterate over the target set and do the test in a subquery) – cf. SQL:
`select * from <table-of-entity-tuples> where <condition>`
- or by a generic function – see Slide 281

SPECIALIZED DATATYPES FOR TIME ETC.

The datatypes specified by XML Schema are used in XPath/XQuery (and XSLT)

- Syntax: constructors like `xs:dateTime('syntactical representation')`
- syntactical representations:
 - `xs:dateTime: yyyy-mm-ddThh:mm:ss[.xx][{+|-}hh:mm]`
 - `xs:date: yyyy-mm-dd` and `xs:time: hh:mm:ss[{+|-}hh:mm]`
 - `xs:duration: P[nY][nM][nD][T[nH][nM][n.n]S]`, where *n* can be any natural number
 - `xs:dayTimeDuration`, `xs:yearMonthDuration`: restrictions of `xs:duration`.

```
let $x := xs:dateTime('2009-08-01T13:51:20.99'),
    $y := xs:date('2008-12-31'),
    $t1 := xs:time('12:50:00+01:00'),      (: timezone +1 = Frankfurt :)
    $t2 := xs:time('15:35:00.50-05:00')   (: timezone -5 = New York :)
return <e const="{ $x }" diff="{ $t2 - $t1 }" d="{ $y + xs:yearMonthDuration("P1Y2M") }"
      sum1="{ xs:time("11:12:00") + xs:dayTimeDuration("PT1H75M") }"
      sum2="{ xs:dateTime("2009-01-10T11:12:00") + xs:dayTimeDuration("P3DT26H40M") }"/>
```

[Filename: XQuery/datetime-test.xq]

- resulting diff = "PT8H45M0.5S" (an `xs:duration`), sum1 = "13:27:00" (an `xs:time`), sum2= "2009-01-14T13:52:00" (an `xs:dateTimes`), d= "2010-02-28" (an `xs:date`)

Actual Usage

... is often simple

```
<country car_code="B">  
  <indep_date>1830-10-04</indep_date>  
</country>
```

```
for $c in //country[indep_date < '1900-01-01']  
return concat($c/name, $c/indep_date)
```

[Filename: XQuery/simple-date-example.xq]

note: explicit `[indep_date < xs:date('1900-01-01')]` would be safer.

FUNCTIONS AND OPERATORS

XPath and XQuery Standard Operators

- Recall Slide 212 for `string()` and `name()`, and Slide 209 for `id()`.
- See “W3C XML Query Functions and Operators” for predefined functions, especially concerning dates+times.

Further Operators: EXSLT

Some mathematical functions (`sqrt` etc.) are not supported as builtins in XPath/XQuery, but only as extensions in EXSLT:

- <http://www.exslt.org/>, supported e.g. by saxon.
- `declare namespace math="http://exslt.org/math"`
- use `math:sqrt`, `math:sin`, ...
- “XPath and XQuery Functions and Operators 3.0” will support math functions (WD 2010, not yet finalized [2011])

USER-DEFINED FUNCTIONS

- User defined functions are declared in the prolog:

```
declare function func_name ([$var1, ..., $varn]) [as returnType]  
{  
    expr that uses $var1, ..., $varn  
}
```

- Parameters: *\$var*_{*i*} [*as paramType*], default for parameter and return types is `item()*` (i.e. a sequence of nodes, literals etc.),
- Any sequence type may be used for *paramType* and *returnType* (cf. XML Schema),
- Any XQuery expression is allowed in the function body.

USER-DEFINED FUNCTIONS: EXAMPLES

- A function computing the population density for a given country:

```
declare function local:density ($name as xs:string) as element(density)
{
  for $c in doc('mondial.xml')//country[name=$name]
  let $density := if ($c/@area > 0) then $c/population div $c/@area else 0
  return <density>{$density}</density>
};
local:density('Germany')
```

[Filename: XQuery/function-density.xq]

- Example for a recursive function:

```
declare function local:depth($e as node()) as xs:integer
{
  if (fn:empty($e/*)) then 1
  else fn:max(for $c in $e/* return local:depth($c)) + 1
};
local:depth(/)
```

[Filename: XQuery/function-depth.xq]

USER-DEFINED FUNCTIONS: EXAMPLE

Ignoring the FLWR, XQuery can even be used as a common functional language:

- every (arithmetic + if) expression is a valid XQuery expression

```
(:call saxonXQ faculty.xq x=5 :)  
declare variable $x external;  
declare function local:faculty($n as xs:integer) as xs:integer  
{ if ($n=1) then 1  
  else $n * local:faculty($n - 1)  
};  
local:faculty($x)
```

[Filename: XQuery/faculty.xq]

USER-DEFINED FUNCTION: EXAMPLE

Remove duplicates from a node set (taken from the example Section from W3C XPath/XQuery Functions and Operators):

```
declare function distinct-nodes-stable ($arg as node()*) as node()*
{
  for $a at $apos in $arg
  let $before_a := fn:subsequence($arg, 1, $apos - 1)
  where every $ba in $before_a satisfies not($ba is $a)
  return $a
}
```

PRACTICAL HINTS: OUTPUT

When creating output, most XQuery engines generate the XML declaration, and output “<” and “>” as “<” and “>”, respectively.

Add Doctype Declaration to the Output

- XQuery engines output only the XML structure itself
- how to add the `<!DOCTYPE mondial SYSTEM "mondial.dtd">` preamble?

With saxon, use

```
declare namespace saxon="http://saxon.sf.net/";  
declare option saxon:output "indent=yes";  
declare option saxon:output "doctype-system=mondial.dtd";
```

Generating non-XML code, text, etc

- In case it is intended to generate e.g. LaTeX, N3 or whatever output, the XML declaration and the `<->`-conversion must be avoided.

With saxon, use

```
declare option saxon:output "method=text";
```

OPERATING WITH SEQUENCES

Comparisons are instance-based: if one operand is a sequence, each value is compared:

- ... as we have seen for XPath: `country[./city/name = "Cordoba"]/name`
`country[./city/population > 1000000]/name`
- but somewhat surprising when using a “let”-view:

```
let $europenames := //country[encompassed/@continent="europe"]/name
for $country in //country
where not ($country/name = $europenames)
return $country/name
```

[Filename: XQuery/seq-comparison-example.xq]

outputs all names of non-european countries.

- selection from let-sequences is also instance-based:

```
let $europcountries := //country[encompassed/@continent="europe"]
return $europcountries[@area>300000]/name
```

[Filename: XQuery/seq-selection-example.xq]

OPERATIONS ON NODES AND NODE SEQUENCES

- “=” compares the string-values of nodes, not “correct” if node identity has to be checked
- “is” compares node identity:

```
for $c in //country,  
    $o in //organization  
where $c/id(@capital) is $o/id(@headq)  
return <pair country='\{$c/@car_code\}' org='\{$o/abbrev\}' />
```

[Filename: XQuery/node-comparison.xq]

- “is” is not allowed for sequences:

```
let $caps := //country/id(@capital)  
for $hq in //organization/id(@headq)  
where $hq is $caps      (: not allowed :)  
return $hq/name
```

[Filename: XQuery/nodes-comparison-example-1.xq]

⇒ “A sequence of more than one item is not allowed as the second operand of “is” ”

OPERATIONS ON NODES AND NODE SEQUENCES

- explicit iteration via some:

```
let $caps := //country/id(@capital)
for $hq in //organization/id(@headq)
where some $cap in $caps satisfies $hq is $cap
return $hq/name
```

[Filename: XQuery/nodes-comparison-example-2.xq]

- `index-of(sequence,item) → integer`

```
let $caps := //country/id(@capital)
for $hq in //organization/id(@headq)
where index-of($caps,$hq) (: checks if $caps contains $hq :)
return $hq/name
```

[Filename: XQuery/nodes-comparison-example-3.xq]

- ... or “where \$caps intersect \$hq”, or even shorter:

```
let $caps := //country/id(@capital)
let $hqs := //organization/id(@headq)
return ($hqs intersect $caps)/name
```

[Filename: XQuery/nodes-comparison-example-4.xq]

FLEXIBILITY

For each task, there is a multitude of possible solutions ...

Example: Uncorrelated Subqueries

Names of all countries that are larger than Germany:

- XPath:

```
//country[@area > number(//country[@car_code='D']/@area)]/name
```

- XQuery and SQL: uncorrelated subquery/semijoin

<pre>for \$c in //country</pre>	<pre>SELECT c.name</pre>
<pre>where \$c/@area ></pre>	<pre>FROM country c</pre>
<pre> number(//country[@car_code='D']/@area)</pre>	<pre>WHERE c.area > (SELECT c2.area</pre>
<pre>return \$c/name</pre>	<pre> FROM country c2</pre>
	<pre> WHERE c2.code = 'D')</pre>

- binding the uncorrelated subquery to a variable:

```
let $germanyarea := number(//country[@car_code='D']/@area)
for $c in //country
where $c/@area > $germanyarea
return $c/name
```

EXERCISES

... see Web.

Exercise 6.1

Determine the lowest mountain that is the highest mountain of the continent where it is located.

Solve the problem for the relational Mondial-DB in SQL, and for XML in XQuery. □

XQUERY: CONCLUSION

Design and Functionality

- combines the positive experiences of previous approaches
- avoids their drawbacks
- intuitively clear syntax and semantics
- declarative, orthogonal, functional style: every expression is a function on nodesets that also returns a nodeset
 - explicit, variable-based iteration: “for *var* in *expression*”
 - implicit iteration: “*collection*[*condition*]” or “*collection/path*”
- Theoretical background (see W3C XML Query Formal Semantics; datatypes of the XML Schema and XML Query Data Model)
 - for each expression (and thus also for its result), the formal type (according to the XML Schema datatypes) can be determined.
 - the type of each variable is determined in the same way.
 - formal, denotational semantics of queries:
“what is the answer set of a given expression?”

XQUERY: CONCLUSION (CONT'D)

W3C XML Query Formal Semantics:

- XPath/XQuery is a functional language.
- is built from expressions, rather than statements. Every construct in the language (except for the XQuery query prolog) is an expression and expressions can be composed arbitrarily.
- The result of one expression can be used as the input to any other expression, as long as the type of the result of the former expression is compatible with the input type of the latter expression with which it is composed.
- Another characteristic of a functional language is that variables are always passed by value, and a variable's value cannot be modified through side effects.

XQUERY: CONCLUSION (CONT'D)

- Note: XQueryX provides a syntax that is formulated in XML

Restrictions

- up to now no resolving of XLink/XPointer (see later)
- only a *query language*:
decision of the W3C: first complete XQuery 1.0 as a query language and make it consistent with XML Schema and XML Query Data Model as a “Recommendation”, and then start official thoughts about updates in XQuery 2.0.

GENERAL DESIGN PATTERNS FOR DATABASE QUERY LANGUAGES

SQL, OQL, XML-QL, XQuery (and many others) use the same underlying principle:

- binding variables
- evaluating a condition
- generating a result (which is a set of data items of the underlying data model)

Note: XQL did not follow this idea \Rightarrow restricted expressiveness and clarity

... let's now have a look on one more XML query language

- the underlying principle is the same

\Rightarrow everything else is “just syntax”!

6.6 Further (Academic) Query Languages

XPATHLOG

- Prolog-/Datalog-style (May, DBPL and VLDB 2001; TPLP 2004)
- based on F-Logic
 - path syntax changed from *step.step.step* to *step/step/step*
 - same syntax for conditions as for F-Logic: “[...]” could be reused
 - F-Logic semantics (1989) closely related with XPath semantics
 - new: distinction between attributes/subelements
- Binding of variables at *arbitrary* positions of an expression
- joins as conjunction (as in Prolog/Datalog)

XPATHLOG

- implicit resolving of multi-valued attributes
- implicit resolving of reference attributes

```
?- //country->C[name->N and @membership->O/name->A].
```

- access to signature/metadata

```
?- //country [name="Germany"] /M.
```

```
?- //country [name="Germany"] /@A.
```

- class membership and -hierarchy

```
?- C isa country [name->N] /M.
```

```
?- _C isa country /@A->_O, _O isa X.
```

```
?- country [@M=>C].      % from DTD
```

XPATHLOG: OVERVIEW

- declarative language
- implicit iteration (fixpoint semantics)
- (equi-)join variables
- XPath-style semantics in rule heads for *generation* and *manipulation* of XML data
- first implementation of an update language for XML (Demo VLDB 2001)
generation of XML in rule heads:

$C[\text{density} \rightarrow D] :- C \text{ isa country}[\text{population} \rightarrow P; @\text{area} \rightarrow A], D \text{ is } P \text{ div } A.$

- fixpoint semantics for execution
- can compute transitive closure etc.

$R[\text{tr_flows_into} \rightarrow S] :- R \text{ isa river}, R/\text{to}[@\text{watertype} \rightarrow \text{"seas"}; @\text{water} \rightarrow S].$

$R[\text{tr_flows_into} \rightarrow S] :- R \text{ isa river}, R/\text{to}[@\text{watertype} \rightarrow \text{"river"}; \text{water} \rightarrow R2],$
 $R2[\text{tr_flows_into} \rightarrow S].$

GENERAL DESIGN PRINCIPLES FOR DATABASE QUERY LANGUAGES

SQL, OQL, XML-QL, XQuery (and many others) use the same underlying principle:

- binding variables
- evaluating a condition
- generating a result (which is a set of data items of the underlying data model)

	SQL/OQL	XML-QL	XQuery	XPathLog
variables:	1-step-navig. SQL: flat data model OQL: + path navig.	XML patterns	XPath navig.	XPath navig.+ XPath patterns
conditions:	WHERE clause	Patterns (equality join conds) WHERE clause (comparisons+joins)	XPath fragment (only non-join-conds) WHERE clause (all)	XPath filters (join conds) separate conjuncts (comparisons+joins)

- the underlying Logic Programming fixpoint semantics enables XPathLog to compute the transitive closure
- ... but it does not allow for syntactically nested statements

FURTHER (ACADEMIC) QUERY LANGUAGES

- XML-GL (Comai, Politecnico Milano, 1999): graphical “language”
- Lorel-XML (Stanford Univ., 1999): OQL-style language, migration of Lorel
- YATL-XML (Cluet, INRIA, 2000): term-based language, migration of YATL
- Lixto/Elog (Gottlob, TU Wien, 2001): graphical tool for data extraction from the Web, Datalog-based internals
- Xcerpt, XChange (Bry et al, LMU München, 2002): term- and unification-based language

... many different approaches to the same goal (mainly in Europe).

Overview in (May, TPLP 2004).

Chapter 7

Manipulating XML Data

- XML data in files:
 - usually no changes (except manually or by scripts)
 - transformations XML → HTML etc: XSLT
- XML data in application systems
 - inside the application programming language; mostly by the DOM-API
 - no special data manipulation language necessary (cf. OQL)?
- different proposals
 - pre-XQuery commercial area:
 - * XMLDB: XUpdate (1999)
 - * eXcelon (2000; XUL as extension of XSLT)
 - academic area:
 - * “Updating XML” (Halevy et al, SIGMOD 2001) as an extension to XQuery
 - * XPathLog (May, VLDB 2001): Prolog-style query- and manipulation language

EXTENDING XQUERY WITH UPDATES – CONCEPTS

In the meantime consensus about which operations is reached. Syntax is still open.

- always wrt. a context node
- base operations:
 - delete *node*
 - rename *node* as *name*
 - insert *node/nodes* before|after|into *node*
- combined operations:
 - replace *node* with *node*
 - move *node* before|after|into *node*

7.1 XML:DB Initiative's XUpdate

- XML:DB Initiative founded in late 1999
Goal: interface for storing XML in databases
- Low-level API (Java etc., using DOM + XPath ...)
- an update concept: XUpdate
- Implementation:
dbXML Core XML Database released as Open Source software in Sept. 2000
transferred to the Apache Software Foundation ("Xindice")
- <http://xml-db-org.sourceforge.net/>
(inactive?)
- The XML:DB database API is implemented in several systems:
eXist, X-Hive, Tamino, XML:DB Lexus, ...

... but here we are mainly interested in XUpdate ...

(note that XUpdate (1999) is not related with XQuery (2001))

XML:DB XUPDATE

Situation in 1999: XML, XPath, XSLT [see later], low-level APIs

- Requirement: “The XML Update specification MUST be an XML element”
i.e., the language is itself in XML syntax (like XSLT and XML Schema)
- XUpdate: a very basic description of update operations:
 - which node (elements, attributes)
 - which operation (delete, update value, append/insert to contents)
 - new value (in case of update/append/insert)

Basic structure:

```
<xu:modifications xmlns:xu= “http://www.xmldb.org/xupdate”>  
  <xu:operation select= “xpath-expression”>  
    contents (e.g. new value)  
  </xu:operation>  
</xu:modifications>
```

... submit such an element as a kind of a “message” to the DB and get the update.

XUpdate: Example

```
<?xml version="1.0"?>
<xu:modifications version="1.0" xmlns:xu="http://www.xmldb.org/xupdate">
  <xu:append select="/mondial/country[name='Germany']">
    <xu:element name='localname'>Deutschland</xu:element>
  </xu:append>
</xu:modifications>
```

[Filename: XUpdate/append.xu]

Calling eXist with (see `client.sh -h`)

```
/bin/gen_client.sh -u user -P password -c /db/may -f mondial.xml -X append.xu
```

executes the update.

- `select= "xpath"` is the same as in XSLT (see later), XML Schema etc. – a widely used concept in the XML world.
(if multiple nodes are addressed, each one is modified)
- `<xu:element>` constructor is the same as in XSLT (1998) and later in XQuery's RETURN clause
- analogously insert-before and insert-after.

XUpdate: Examples (Cont'd)

```
<?xml version="1.0"?>
<xu:modifications version="1.0" xmlns:xu="http://www.xmldb.org/xupdate">
  <xu:remove select="/mondial/country[name='Germany']/localname"/>
</xu:modifications>
```

[Filename: XUpdate/remove.xu]

```
<?xml version="1.0"?>
<xu:modifications version="1.0" xmlns:xu="http://www.xmldb.org/xupdate">
  <xu:update select="/mondial/country[name='Germany']/population">
    80000000
  </xu:update>
</xu:modifications>
```

[Filename: XUpdate/update.xu]

XUpdate: Examples (Cont'd)

- get the new value from the database:

```
<?xml version="1.0"?>
<xu:modifications version="1.0" xmlns:xu="http://www.xmldb.org/xupdate">
  <xu:update select="/mondial/country[name='Germany']/population/text()">
    <xu:value-of select="/mondial/country[name='Germany']/@area"/>
  </xu:update>
</xu:modifications>
```

[Filename: XUpdate/update-select.xu]

note: the inner `select` cannot depend on the current node.

```
<?xml version="1.0"?>
<xu:modifications version="1.0" xmlns:xu="http://www.xmldb.org/xupdate">
  <xu:variable name="bla"
    select="/mondial/country[name='Germany']/gdp_total/text()"/>
  <xu:update select="/mondial/country[name='Germany']/population/text()">
    <xu:value-of select="$bla"/>
  </xu:update>
</xu:modifications>
```

[Filename: XUpdate/update-variable.xu]

XUPDATE: CONCLUSION AND COMMENTS

- XML-syntax of the language strongly influenced by XSLT (1998)
 - elements as commands
 - `select="..."` selects nodes to which the command is applied
 - use of variables `select="$variable"` as in XSLT, and later also in XQuery
 - element/command contents specifies what is to be done
 - element generation by literal XML (also in XSLT and later XQuery)
- only very simple functionality
 - no way to compute the inner value,
 - no iteration etc.
- same time: combination with XSLT and XUpdate to XUL (XML Update Language/Updategrams [Excelon]):
XSLT program structures + XUpdate operations, applied to “current node” of XSLT.

7.2 XQuery with Updates 2001

- extend a *declarative query language* with updates
- based on *variable bindings*
- SQL: FROM-WHERE for selecting nodes ...
... that are then modified.
- XQuery: FOR-LET-WHERE for selecting nodes ...
... that are then modified.
- update instead of the return-clause (cf. SQL: UPDATE vs. SELECT)?
- or what?

XQUERY WITH UPDATES – AN EARLY PROPOSAL

- QuiP: the 2001/02 XQuery prototype of Software AG [Diplomarbeit P. Lehti 2001], later integrated into the Tamino system (before: XQL).
- calling `quip filename.xq > bla.xml` wrote the modified XML to a file.

```
update
for $c in document("twocountries.xml")//country
let $area := string($c/@area)
delete $c/@area
insert <area>{$area}</area> after $c/name
rename $c//city[@id=$c/@capital] as capital
replace $c/@car_code with
    attribute code {concat($c/name/text(), ":", string($c/@car_code))}
replace $c/population/text() with
    $c/population/text() * (1 + $c/population_growth div 100)
insert "biggest city" into
    $c//city[population = max(for $citypop in $c//city/population/text()
        return int($citypop))]
```

[Filename: XQuery/update.quip]

XQUERY WITH UPDATES – W3C PROPOSAL

- XQuery reached recommendation state in 2007 ... as a query language still without updates.
- “XQuery Update Facility”, first W3C Working Draft has been published 27 January 2006;
<http://www.w3.org/TR/xqupdate>

New Expressions

do insert *SourceExpr* (([as (first | last)] into) | after | before) *TargetExpr*

do delete *TargetExpr*

do rename *TargetExpr* as *Expr* // *Expr* must result in a *qname*

do replace [value of] *TargetExpr* with *Expr*

- Syntax still changing,
- not implemented in saxonA 9.1, only in (commercial) saxonB 9.1

XQuery with Updates – Transformation Command

- not an update!

1. assign variable(s),
2. update things bound to the variable(s),
3. return something generated from the (updated) variables.

transform copy $\$VarName := Expr$ (, $\$VarName := ExprSingle$)*
modify *UpdateExpr* return *Expr*

Chapter 8

The Transformation Language XSL

8.1 XSL: Extensible Stylesheet Language

- developed from
 - CSS (Cascading Stylesheets) scripting language for transformation of data sources to HTML or any other optical markup, and
 - DSSSL (Document Style Semantics and Specification Language), stylesheet language for SGML. Functional programming language.
- idea: rule-based specification how elements are transformed and formatted *recursively*.
 - Input: XML
 - Output: XML (special case: HTML)
- declarative/functional: **XSLT (XSL Transformations)**

APPLICATIONS

- XML → XML
 - Transformation of an XML instance into a new instance according to another DTD,
 - Integration of several XML instances into one,
 - Extraction of data from an XML instance,
 - Splitting an XML instance into several ones.
- XML → HTML
 - Transformation of an XML instance to HTML for presentation in a browser
- XML → anything
 - since no data structures, but only ASCII is generated, \LaTeX , postscript, pdf can also be generated
 - ... or transform to **XSL-FO (Formatting objects)**.

THE LANGUAGE(S) XSL

Partitioned into two sublanguages:

- functional programming language: **XSLT**
“understood” by **XSLT-Processors** (e.g. xt, xalan, saxon, xsltproc ...)
- generic language for document-markup: **XSL-FO**
“understood” by XSL-FO-enabled **browsers** that transform the XSL-FO-markup according to an internal specification into a direct (screen/printable) presentation.
(similar to LaTeX)
- **programming paradigm: self-organizing tree-walking**
- XSL itself is written in **XML-Syntax**.
It uses the **namespace prefixes** “xsl:” and “fo:”,
bound to `http://www.w3.org/1999/XSL/Transform` and
`http://www.w3.org/1999/XSL/Format`.
- XSL programs can be seen as XML data.
- it can be combined with other languages that also have an XML-Syntax (and an own namespace).

APPLICATION: XSLT FOR XML → HTML

- the prolog of the XML document contains an instruction that specifies the stylesheet to be used:

```
<?xml version="1.0"?>
```

```
<?xml-stylesheet type="text/xsl" href="mondial-simple.xsl"?>
```

```
<!DOCTYPE mondial SYSTEM "mondial.dtd">
```

```
<mondial>    ... </mondial>
```

- if an (XSL-enabled) browser finds an XML document with a stylesheet instruction, then the XML document is processed according to the stylesheet (by the browser's own XSLT processor), and the result is shown in the browser.

(e.g.,

```
http://www.informatik.uni-goettingen.de/Teaching/SSD/XSLT/mondial-with-stylesheet.xml)
```

⇒ click "show source" in the browser

- **Remark: not all browsers support the full functionality (id()-function)**
- in general, for every main "object type" of the underlying application, there is a suitable stylesheet how to present such documents.

8.2 XSLT: Syntax and Semantics

- Each XSL-stylesheet is itself a valid XML document,

```
<?xml version="1.0">
```

```
<xsl:stylesheet version="2.0"
```

```
  xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
```

```
  ...
```

```
</xsl:stylesheet>
```

- contains elements of the namespace **xsl:** that specify the transformation/formatting,
- contains literal XML for generating elements and attributes of the resulting document,
- uses XPath expressions for accessing nodes in an XML document. XPath expressions (mostly) occur as attribute values of `<xsl:...>` elements, (e.g., `<xsl:copy-of select='xpath'>`)
- XSL stylesheets/programs recursively generate a result tree from an XML input tree.

8.2.1 XSLT: Flow Control by Templates

The stylesheet consists mainly of *templates* that specify the instructions *how* elements should be processed:

- `xsl:template`:

```
<xsl:template match="xsl-pattern">  
  content  
</xsl:template>
```

- *xsl-pattern* is an XPath expression without use of “*axis::*” (cf. Slide 198). It indicates for which elements (types) the template is applicable:
a node *x* satisfies *xsl-pattern* if there is some ancestor node *k* of *x*, such that *x* is in the result set of *xsl-pattern* for *k* as context node.
(another selection takes place at runtime when the nodes are processed for actually deciding to apply a template to a node).
- *content* contains the XSL statements for generation of a fragment of the result tree.

TEMPLATES

- `<xsl:template match="city">`
 `<xsl:copy-of select="current()"/>`
`</xsl:template>`

is a template that can be applied to cities and copies them unchanged into the result tree.

- `<xsl:template match="lake|river|sea"> ... </xsl:template>`
can be applied to waters.

- `<xsl:template match="country/province/city"> ... </xsl:template>`
can be applied to city elements that are subelements of province elements that in course are subelements of country elements.

- `<xsl:template match="id('D')"> ... </xsl:template>`
can be applied to the element whose ID is "D".

- `<xsl:template match="city[population > 1000000]"> ... </xsl:template>`
can be applied to city elements that have more than 1000000 inhabitants.

EXECUTION OF TEMPLATES: “TREE WALKING”

- `xsl:apply-templates`:

`<xsl:apply-templates select=“xpath-expr”/>`

- *xpath-expr* is an XPath expression that indicates for which elements (starting from the node where the current template is applied as context node) “their” template should be applied.

Note that elements are processed in order of the final axis of the select expression.

- By `<xsl:apply-templates>` elements inside the content of `<xsl:template>` elements, the hierarchical structure of XML documents is processed
 - simplest case (often in XML → HTML): depth-first-search
 - can also be influenced by the “select” attribute: “tree jumping”
- if all subelements should be processed, the “select” attribute can be omitted.

`<xsl:apply-templates/>`

TEMPLATES

- `<xsl:apply-templates select="country"/>`
processes all country subelements of the current context element.
- `<xsl:apply-templates select="country/city"/>`
processes all city subelements of country subelements of the current context element,
- `<xsl:apply-templates select="/mondial//city[population > 1000000]"/>`
processes all city elements that are contained in Mondial and whose population is more than 1000000,
- `<xsl:apply-templates select="id(@capital)"/>`
processes the element whose ID equals the value of the capital-(reference) attribute of the current context element.

TEMPLATES

- One template must be applicable to the root element for initiating the processing:
 - `<xsl:template match="name_of_the_root_element">`
 - `<xsl:template match="/">`
 - `<xsl:template match="*">`

RULE-BASED “PROGRAMMING”

- local semantics: templates as “rules”
- global semantics: built-in implicit tree-walking combines rules

TEMPLATES: EXAMPLE

Presentation of the country information as a table (→ HTML)

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="2.0">

  <xsl:template match="mondial">
    <html> <body> <table>
      <xsl:apply-templates select="country"/>
    </table> </body> </html>
  </xsl:template>

  <xsl:template match="country">
    <tr><td> <xsl:value-of select="name"/> </td>
      <td> <xsl:value-of select="@car_code"/> </td>
      <td align="right"> <xsl:value-of select="population"/> </td>
      <td align="right"> <xsl:value-of select="@area"/> </td>
    </tr>
  </xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/mondial-simple.xsl]

TEMPLATES: EXAMPLE

Presentation of the country and city information as a table:

```
<xsl:stylesheet version="2.0"
  xmlns:xsl="http://www.w3.org/1999/XSL/Transform">

<xsl:template match="mondial">
  <html><body><table>
    <xsl:apply-templates select="country"/>
  </table></body></html>
</xsl:template>

<xsl:template match="country">
<tr valign="top">
  <td><xsl:value-of select="name"/></td>
  <td><xsl:value-of select="@car_code"/></td>
  <td align="right"><xsl:value-of select="population"/></td>
  <td align="right"><xsl:value-of select="@area"/></td>
  <td valign="top">
    <table><xsl:apply-templates select="./city"/></table>
  </td>
</tr>
</xsl:template>

<xsl:template match="city">
<tr> <td width="100"><xsl:value-of select="name"/></td>
  <td align="right" width="100">
    <xsl:value-of select="population[1]"/>
  </td>
</tr>
</xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/mondial-nested.xsl]

TEMPLATES: EXAMPLE

The following (transformation: XML → XML) stylesheet copies all country and city elements from Mondial and outputs first all country elements, and then all city elements as top-level elements:

```
<xsl:stylesheet version="2.0" xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
  <!-- Template that copies elements -->
  <xsl:template match="city|country">
    <xsl:copy-of select="current()"/>
  </xsl:template>
  <xsl:template match="mondial">
    <!-- apply templates: first countries -->
    <xsl:apply-templates select="/mondial/country">
    <!-- apply templates: then cities -->
    <xsl:apply-templates select="//country/city | //country/province/city"/>
  </xsl:template>
</xsl:stylesheet>
```

TEMPLATES

Difference between:

1. `<xsl:template match="xsl-pattern">`
 content
`</xsl:template>`

2. `<xsl:apply-templates select="xpath-expr"/>`

- `select="..."` is evaluated wrt. the current context node (selects which elements are addressed by the given XPath expression),
- `match="..."` is evaluated wrt. the document structure starting from "below" (checks if the document structure matches with the pattern),
- `xsl:apply-templates` selects nodes for application by its *xpath-expr*, and then the suitable templates are applied,
- the order of templates has no effect on the order of application (document order of the selected nodes).

TEMPLATES

Exercise 8.1

Describe the difference between the following stylesheet fragments:

1. `<xsl:template match="city">`
 `<xsl:copy-of select="current()"`
 `</xsl:template>`
 `<xsl:apply-templates select="//country/city"/>`
 `<xsl:apply-templates select="//country/province/city"/>`
2. `<xsl:template match="country/city">`
 `<xsl:copy-of select="current()"`
 `</xsl:template>`
 `<xsl:template match="country/province/city">`
 `<xsl:copy-of select="current()"`
 `</xsl:template>`
 `<xsl:apply-templates select="//country/city|//country/province/city">`

□

CONFLICTS BETWEEN TEMPLATES

When using non-disjoint match-specifications of templates (e.g. *, city, country/city, city[population>1000000]) (including possibly templates from imported stylesheets), several templates are probably applicable.

- in case that during processing of an `<xsl:apply-templates>`-command several templates are applicable, the one with the most specific match-specification is chosen.
- defined by *priority rules* in the XSLT spec (that also define priorities between incomparable patterns)

OVERRIDING (SINCE XSLT 2.0)

The above effect is similar to *overriding* of methods in object-oriented concepts: always take the most specific implementation

- `<xsl:next-match>`: apply the next-lower-specific rule (among those defined in the same stylesheet)
- `<xsl:apply-imports>`: apply the next-lower-specific rule (among those defined in imported stylesheets (see later))

RESOLVING TEMPLATE CONFLICTS MANUALLY

Process a node with different templates depending on situation:

- associating “modes” with templates and using them in apply-templates

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
    version="2.0">
  <xsl:output method="xml" indent="yes"/>
  <xsl:template match="mondial">
    <xsl:apply-templates select="country[@area>1000000]"/>
    ... and now the second part ...
    <xsl:apply-templates select="country[@area>1000000]" mode="bla"/>
  </xsl:template>
  <xsl:template match="country">
    <firsttime> <xsl:value-of select="name"/> </firsttime>
  </xsl:template>
  <xsl:template match="country" mode="bla">
    <secondtime> <xsl:value-of select="name"/> </secondtime>
  </xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/mondial-modes.xsl]

NAMED TEMPLATES

Named templates serve as macros and can be called by their name.

- `xsl:template` with “name” attribute:

```
<xsl:template name=“name”>  
  content  
</xsl:template>
```

- *name* is an arbitrary name
 - *content* contains xsl-statements, e.g. `xsl:value-of`, which are evaluated against the current context node.
- `xsl:call-template`

```
<xsl:call-template name=“name”/>
```
 - Example: Web pages – templates for upper and left menus etc.

8.2.2 XQuery and XSLT

- both are declarative, functional languages ...
- ... with completely different strategies:
 - XQuery: nesting of the return-statement directly corresponds to the structure of the result
 - XSLT: the nested processing of templates yields the structure of the result.

XSLT

- modular structure of the stylesheets
- extensibility and reuse of templates
- flexible, data-driven evaluation

XQuery

- better functionality for joins (for \$a in ..., \$b in ...)
- XSLT: joins must be programmed explicitly as nested loops (xsl:for-each)

TRANSLATION XSLT → XQUERY

- each template is transformed into an FLWR statement,
- inner template-calls result in nested FLWR statements inside the return-clause
- genericity of e.g. `<apply-templates/>` cannot be expressed in XQuery since it is not known which template is activated

⇒ the more flexible the schema (documents), the more advantages show up for XSLT.

Exercise 8.2

- Give XQuery queries that do the same as `mondial-simple.xsl` and `mondial-nested.xsl`.
- Give an XQuery query that does the same as the stylesheet on Slide 321. □

8.2.3 XSLT: Generation of the Result Tree

Nodes can be inserted into the result tree by different ways:

- literal XML values and attributes,
- copying of nodes and values from the input tree,
- generation of elements and attributes by constructors.

Configuring Output Mode

- recommended, top level element (see xsl doc. for details):
`<xsl:output method="xml|html|xhtml|text" indent="yes|no"/>`
(not yet supported by all XSLT tools; saxon has it)

Generation of Structure and Contents by Literal XML

- All tags, elements and attributes in the content of a template that do not belong to the xsl-namespace (or to the local namespace of an xsl-tool), are literally inserted into the result tree.
- with `<xsl:text> some_text </xsl:text>`, text can be inserted explicitly (whitespace, e.g. when generating IDREFS attributes).

GENERATION OF THE RESULT TREE

Copying from the Input Tree

- `<xsl:copy>contents</xsl:copy>`
copies the current context node (i.e., its “hull”): all its namespace nodes, but *not* its attributes and subelements (note that contents can then be generated separately).
- `<xsl:copy-of select=“xpath-expr”/>`
copies the result of *xpath-expr* (applied to the current context) unchanged into the result tree.
- `<xsl:value-of select=“xpath-expr” [separator=“char”]/>`
generates a text node with the value of *xpath-expr*.
Applied to multiple nodes, the partial results are separated by *char* (default: space).
[note: the latter changed from XSLT 1.0 (apply only to 1st node) to 2.0]

Exercise 8.3

Consider the differences between `<xsl:copy/>`, `<xsl:copy-of select=“current()”/>` and `<xsl:value-of select=“current()”/>`.

In which cases do two commands have the same result?

□

GENERATION OF THE RESULT TREE

Example:

```
<xsl:template match="city">
  <mycity>
    <xsl:value-of select="name"/>
    <xsl:copy-of select="longitude|latitude"/>
  </mycity>
</xsl:template>
```

- generates a mycity element for each city element,
- the name is inserted as #PCDATA content,
- the subelements longitude and latitude are copied:

```
<mycity>Berlin
  <longitude>13.3</longitude>
  <latitude>52.45</latitude>
</mycity>
```

GENERATION OF THE RESULT TREE

For inserting attribute values,

```
<xsl:value-of select="xpath-expr"/>
```

cannot be used *directly*. Instead, XPath expressions have to be enclosed in {...}:

```
<xsl:template match="city">  
  <mycity key="{@id}">  
    <xsl:value-of select="name"/>  
    <xsl:copy-of select="longitude|latitude"/>  
  </mycity>  
</xsl:template>
```


GENERATION OF THE RESULT TREE

Example:

```
<xsl:template match="city">
  <mycity source="mondial"
        country="{ancestor::country/name}">
    <xsl:apply-templates/>
  </mycity>
</xsl:template>
```

- generates a “mycity” element for each “city” element,
- constant attribute “source”,
- attribute “country”, that indicates the country where the city is located,
- all other attributes are omitted,
- for all subelements, suitable templates are applied.

XSLT: GENERATION OF THE RESULT TREE

Generation of Elements and Attributes

- `<xsl:element name="xpath-expr">`
 content
`</xsl:element>`

generates an element of element type *xpath-expr* in the result tree, the content of the new element is *content*. This allows for computing element names.

- `<xsl:attribute name="xpath-expr">`
 content
`</xsl:attribute>`

generates an attribute with name *xpath-expr* and value *content* which is added to the surrounding element under construction.

- With `<xsl:attribute-set name="name"> xsl:attribute* </xsl:attribute-set>`
attribute sets can be predefined. They are used in `xsl:element` by
 `use-attribute-sets="attr-set1 ... attr-setn"`

GENERATION OF IDREFS ATTRIBUTES

- XML source: “border” subelements of “country” with an IDREF attribute “country”:
`<border country=“car_code” length=“...”>`
- result tree: IDREFS attribute `country/@neighbors` that contains all neighboring countries
- two ways how to do this (both require XSLT 2.0)

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="2.0">
<xsl:template match="*"><xsl:apply-templates select="country"/></xsl:template>
<xsl:template match="country">
  <country neighbors1="{border/@country}"> <!-- note: adds whitespace as separator -->
    <xsl:attribute name="neighbors2">
      <xsl:value-of select="border/@country"/> <!-- default separator: whitespace -->
    </xsl:attribute>
  </country>
</xsl:template></xsl:stylesheet>
```

[Filename: XSLT/mondial-neighbors.xsl]

8.2.4 XSLT: Control Structures

... so far the “rule-based”, clean XSLT paradigm with implicit recursive semantics:

- templates: recursive control of the processing

... further control structures inside the content of templates:

- iterations/loops
- branching

DESIGN OF XSLT COMMAND ELEMENTS

- semantics of these commands as in classical programming languages (Java, C, Pascal, Basic, Cobol, Algol)
- Typical XML/XSLT design: element as a command, further information as attributes or in the content (i.e., iteration specification, test condition, iteration/conditional body).

ITERATIONS

For processing a list of subelements or a multi-valued attribute, local iterations can be used:

```
<xsl:for-each select="xpath-expr">  
  content  
</xsl:for-each>
```

- inside an iteration the “iteration subject” is not bound to a variable (like in XQuery as `for $x in xpath-expression`), but
- the current node is that from the `xsl:for-each`, not the one from the surrounding `xsl:template`
- an `xsl:for-each` iteration can also be used for implementing behavior that is different from the templates “matching” the elements (instead of using modes).

FOR-EACH: EXAMPLE

Presentation of the country and city information as a table:

```
<xsl:stylesheet version="2.0"
  xmlns:xsl="http://www.w3.org/1999/XSL/Transform">

<xsl:template match="mondial">
  <html><body><table>
    <xsl:apply-templates select="country"/>
  </table></body></html>
</xsl:template>

<xsl:template match="country">
<tr valign="top">
  <td><xsl:value-of select="name"/></td>
  <td><xsl:value-of select="@car_code"/></td>
  <td align="right"><xsl:value-of select="population"/></td>
  <td align="right"><xsl:value-of select="@area"/></td>
  <td valign="top">
    <table>
      <xsl:for-each select="//city">
        <tr> <td width="100"><xsl:value-of select="name"/></td>
          <td align="right" width="100">
            <xsl:value-of select="population[1]"/>
          </td>
        </tr>
      </xsl:for-each>
    </table>
  </td>
</tr>
</xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/mondial-nested-for-each.xsl]

XSLT: CONDITIONAL PROCESSING

- Simple Test:

```
<xsl:if test="predicate"> content </xsl:if>
```

Example:

```
<xsl:template match="country">  
  <table>  
    <tr> <th colspan="2"> <xsl:value-of select="name"> </th>  
    </tr>  
    <xsl:if test="@area">  
      <tr>  
        <td> Area: </td>  
        <td> <xsl:value-of select="@area"> </td>  
      </tr>  
    </xsl:if>  
    ⋮  
  </table>  
</xsl:template>
```

XSLT: CONDITIONAL PROCESSING

- Multiple alternatives:

```
<xsl:choose>
  <xsl:when test="predicate1">
    content1
  </xsl:when>
  <xsl:when test="predicate2">
    content2
  </xsl:when>
  ...
  <xsl:otherwise>
    contentn+1
  </xsl:otherwise>
</xsl:choose>
```


8.2.5 XSLT: Variables and Parameters

Variables and parameters serve for binding values to names.

VARIABLES

- variables can be assigned only once (in their definition). A later re-assignment (like in C or Java) is not possible.
- variables can be defined as top-level elements which makes them visible in the whole document (as a constant).
- a variable definition can take place at an arbitrary position inside a template - such a variable is visible in all its following siblings, e.g.,
 - a variable before a `<xsl:for-each>` is visible inside the `<xsl:for-each>`;
 - a variable inside a `<xsl:for-each>` gets a new value for each iteration to store an intermediate value.

BINDING AND USING VARIABLES

- value assignment either by a “select” attribute (value is a string, a node, or a set of nodes)

```
<xsl:variable name=“var-name” select=“xpath-expr”/>
```

- or as element content (then, the value can be a tree which is generated dynamically by XSLT)

```
<xsl:variable name=“var-name”>
```

```
  content
```

```
</xsl:variable>
```

- Usage: by select=*“\$var-name”*

Example: Variables

A simple, frequent use is to “keep” the outer current element when iterating by an xsl:for-each:

- Consider the previous “border”-example
- now: generate a table of neighbors

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="1.0">
<xsl:template match="*">
  <table><xsl:apply-templates select="country"/></table>
</xsl:template>
<xsl:template match="country">
  <xsl:variable name="country" select="."/>
  <xsl:for-each select="border">
    <tr>
      <td><xsl:value-of select="$country/@car_code"/></td>
      <td><xsl:value-of select="@country"/></td>
    </tr>
  </xsl:for-each>
</xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/mondial-neighbors-table.xsl]

PARAMETERS

- ... similar to variables
- values are communicated to called templates by parameters,
- the definition of parameters is allowed *only at the beginning* of `xsl:template` elements. The defined parameter is then visible everywhere in the template body.
- the assignment of a value takes place in the calling `<xsl:apply-templates>` or `<xsl:call-template>` element.
- pure call-by-value, no call-by-reference possible.

Remark: since a parameter can be an element with substructures, theoretically, a single parameter is always sufficient.

COMMUNICATION OF PARAMETERS TO TEMPLATES

- Parameters are declared at the beginning of a template:

```
<xsl:template match="...">
  <xsl:param name="param-name"
             select="xpath-expr"/>    <!-- with a default value -->
  :
</xsl:template>
```

- the parameter values are then given with the template call:

```
<xsl:apply-templates select="xpath-expr1">
  <xsl:with-param name="param-name"
                  select="xpath-expr2" />
</xsl:apply-templates>
```

- Often, parameters are propagated downwards through several template applications/calls.

This can be automatized (since XSLT 2.0) by

```
<xsl:param name="param-name" select="xpath-expr" tunnel="yes">
```

Example: Parameters

Generate a table that lists all organizations with all their members. The abbreviation of the organisation is communicated by a parameter to the country template which then generates an entry:

→ next slide

[Filename: orgs-and-members.xsl]

Exercise 8.4

- Extend the template such that it also outputs the type of the membership.
- Write an equivalent stylesheet that does not call a template but works explicitly with `<xsl:for-each>`.
- Give an equivalent XQuery query (same for the following examples). □

EXAMPLE (CONT'D)

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="1.0">
  <xsl:template match="mondial">
    <html><body> <h2>Membership Table</h2>
      <table> <xsl:apply-templates select="organization"/>
    </table></body></html>
  </xsl:template>
  <xsl:template match="organization">
    <tr><td colspan="2"><xsl:value-of select="name"/></td></tr>
    <xsl:apply-templates select="id(members/@country)">
      <xsl:with-param name="the_org" select="name/text()"/>
    </xsl:apply-templates>
  </xsl:template>
  <xsl:template match="country">
    <xsl:param name="the_org"/>
    <tr><td><xsl:value-of select="$the_org"/></td>
      <td><xsl:value-of select="name/text()"/></td></tr>
  </xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/orgs-and-members.xsl]

EXAMPLE/COMPARISON OF MECHANISMS

Example: This example illustrates the implicit and explicit iterations, and the use of variables/parameters

[use file:XSLT/members1.xsl and develop the other variants]

- Generate a list of the form

```
<organization> EU <member>Germany</member>  
                <member>France</member> ... </organization>
```

- using template-hopping [Filename: XSLT/members1.xsl]
- using xsl:for-each [Filename: XSLT/members2.xsl]

- Generate a list of the form

```
<membership organization="EU" country="Germany"/>
```

based on each of the above stylesheets.

- template hopping: requires a parameter [Filename: XSLT/members3.xsl]
- iteration: requires a variable [Filename: XSLT/members4.xsl]

A POWERFUL COMBINATION: VARIABLES AND CONTROL

```
<xsl:variable name="var-name">  
  content  
</xsl:variable>
```

Any structure that is generated in *content* is then bound to the variable.

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="1.0">  
  <xsl:output method="xml" indent="yes"/>  
  <xsl:template match="mondial">  
    <xsl:variable name="berlin">  
      <bla>  
        <xsl:copy-of select="//city[name='Berlin']"/>  
      </bla>  
    </xsl:variable>  
    <xsl:copy-of select="$berlin"/>  
    <xsl:copy-of select="$berlin/bla/city/name"/>  
  </xsl:template>  
</xsl:stylesheet>
```

[Filename: XSLT/var-1.xsl]

Even more powerful

Anything inside the *contents* is bound to the variable – this allows even to generate complex structures by template applications:

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="1.0">
  <xsl:output method="xml" indent="yes"/>

  <xsl:template match="mondial">
    <xsl:variable name="bigcities">
      <xsl:apply-templates select="//city"/>
    </xsl:variable>
    <xsl:copy-of select="$bigcities"/>
    <xsl:copy-of select="$bigcities//name"/>
  </xsl:template>

  <xsl:template match="city">
    <xsl:if test='number(population)>1000000'>
      <xsl:copy-of select="current()"/>
    </xsl:if>
  </xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/var-2.xsl]

EXTERNAL PARAMETERS

Stylesheets can be called with external parameters (e.g., from the shell, or from a Java environment):

- define formal parameters for the stylesheet:

```
<xsl:stylesheet ...>  
  <xsl:parameter name="name1"/>  
  <xsl:parameter name="name2"/>  
  stylesheet contents  
  (parameters used as $namei)  
</xsl:stylesheet ...>
```

- call e.g. (with saxon)

```
saxonXSL -s bla.xml bla.xsl name1=value1 name2=value2
```

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="1.0">  
  <xsl:param name="country"/>  
  <xsl:template match="mondial">  
    <xsl:copy-of select="//country[name=$country]"/>  
  </xsl:template>  
</xsl:stylesheet>
```

[Filename: XSLT/external-param.xsl]

8.2.6 XSLT: Miscellaneous

SORTING

For the set-based XSLT elements

- `xsl:apply-templates` and
- `xsl:for-each`

it can be specified whether the elements should be processed in the order of some key:

```
<xsl:sort select="xpath-expr"  
  data-type = {"text"|"number"}  
  order = {"descending"|"ascending"}/>
```

- “select” specifies the values according to which the nodes should be ordered (evaluated wrt. the node as context node)
- “data type” specifies whether the ordering should be alphanumeric or numeric,
- “order” specifies whether the ordering should be ascending or descending,
- if an “`xsl:apply-templates`”- or “`xsl:for-each`” element has multiple “`xsl:sort`” subelements, these are applied in a nested way (as in SQL).

GROUPING (SINCE XSLT 2.0)

Extends the `<xsl:for-each>` concept to groups:

```
<xsl:for-each-group select="xpath-expr" group-by="local-key">  
  content  
</xsl:for-each-group>
```

Inside the content part:

- current element is the *first* element of the current group
⇒ for accessing/returning the whole group, something else must be used:
- `current-group()` returns the sequence of all elements of the current group (e.g., `current-group()/name` for all their names); can e.g. be used for aggregation
- `current-grouping-key()` returns the current value of the grouping key
- `position()` returns the number of the current group

Grouping (Example)

Example

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="1.0">
  <xsl:output method="xml" indent="yes"/>
  <xsl:template match="mondial">
    <xsl:for-each-group select="country" group-by="encompassed/@continent">
      <continent nr="{position()}">
        <xsl:copy-of select="id(current-grouping-key())/name"/>
        <xsl:copy-of select="current-group()/name"/>
      </continent>
    </xsl:for-each-group>
  </xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/for-each-group.xsl]

Exercise 8.5

Do the same in XQuery (note: use “let” for the group).

□

HANDLING NON-XSLT NAMESPACES IN XSLT

- namespaces used in the queried document (e.g., xhtml)
- namespaces to be used in the generated document
- namespaces used in the XSLT stylesheet (xsd, fn, ...)

Declare the namespaces in the surrounding `<xsl:stylesheet>` element:

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
               xmlns:ht="http://www.w3.org/1999/xhtml"
               version="2.0">
```

tells the XSL processor that the namespace bound to 'http://www.w3.org/1999/xhtml' is denoted by "ht:" in this document.

(and `<ht:body>` is different from `<body>`)

Querying XHTML documents with namespace

```
<!--
  call: saxonXSL -s http://www.dbis.informatik.uni-goettingen.de/index.html
        -xsl xsl-html.xsl
  note: takes some time ...
-->
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
                xmlns:ht="http://www.w3.org/1999/xhtml"
                version="2.0">
<xsl:template match="/">
  <result>
    <xsl:copy-of select="//ht:li"/>
  </result>
</xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/xsl-html.xsl]

USING FUNCTIONS FROM XQUERY FUNCTIONS AND OPERATORS

- the functions and operators from “XQuery Functions and Operators” (e.g., aggregations) are also available in XSLT.
- Namespace: <http://www.w3.org/2005/xpath-functions>

Example

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
               xmlns:fn="http://www.w3.org/2005/xpath-functions"
               version="1.0">
  <xsl:output method="xml" indent="yes"/>
  <xsl:template match="mondial">
    <xsl:for-each select="country">
      <country name="{name/text()}">
        <xsl:copy-of select="fn:sum(../city/population[1])"/>
      </country>
    </xsl:for-each>
  </xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/fn-sum.xsl]

USER-DEFINED FUNCTIONS (SINCE XSLT 2.0)

```
<xsl:function name="local-ns:fname">  
  <xsl:param name="param1"/>  
  :  
  <xsl:param name="paramn"/>  
  contents  
</xsl:function>
```

- the *local-ns* must be declared by `xmlns:local-ns='uri'` in the `xsl:stylesheet` element;
- function can then be used with *n* parameters in `xsl:value-of`, or in any XPath expression.

e.g.,

```
<xsl:value-of select="local-ns:fname(value1, . . . , valuen)"/>
```

ACCESS TO DATA FROM MULTIPLE DOCUMENTS

- using the document()-function from XPath:
- note the use of " and '

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
                xmlns:ht="http://www.w3.org/1999/xhtml"
                version="2.0">
<xsl:template match="/"> <!-- call it for any xml document -->
  <result>
  <xsl:copy-of
    select="document('http://www.dbis.informatik.uni-goettingen.de/index.html')//ht:li"/>
  </result>
</xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/web-queries.xsl]

GENERATION OF MULTIPLE INSTANCES

- controlling output to different files (since XSLT 2.0):
`<xsl:result-document href="output-file-uri">`
- note: generates directories if required.

```
<?xml version="1.0" encoding="utf-8"?>
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="2.0">
  <xsl:output method="xml" indent="yes"/>
  <xsl:template match="*">
    <xsl:result-document href="tmp/countries.xml">
      <countries><xsl:apply-templates select="country"/></countries>
    </xsl:result-document>
    <xsl:result-document href="tmp/organizations.xml">
      <organizations><xsl:apply-templates select="organization"/></organizations>
    </xsl:result-document>
  </xsl:template>
  <xsl:template match="country"><xsl:copy-of select="name"/></xsl:template>
  <xsl:template match="organization"><xsl:copy-of select="name"/></xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/redirected-output.xml]

GENERATION OF MULTIPLE INSTANCES

- also possible with dynamically computed filenames:
generates a file for each country:

```
<?xml version="1.0" encoding="utf-8"?>
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="2.0">
  <xsl:output method="xml" indent="yes"/>
  <xsl:template match="*">
    <xsl:result-document href="tmp/countries.xml">
      <countries><xsl:apply-templates select="country"/></countries>
    </xsl:result-document>
    <xsl:result-document href="tmp/organizations.xml">
      <organizations><xsl:apply-templates select="organization"/></organizations>
    </xsl:result-document>
  </xsl:template>
  <xsl:template match="country"><xsl:copy-of select="name"/></xsl:template>
  <xsl:template match="organization"><xsl:copy-of select="name"/></xsl:template>
</xsl:stylesheet>
```

[Filename:

XSLT/redirected-output-countries.xml]

IMPORT MECHANISMS

XSL-stylesheets can import other stylesheets (i.e., they import their rules):

- `<xsl:include href="url"/>`
for conflict-free stylesheets,
- `<xsl:import href="url"/>`
definitions of the importing document have higher priority than definitions from the imported documents,
the `xsl:import` subelements must precede all other subelements.

Example: DBIS Web pages

- general macros, frames etc. as templates in separate files
- individual page content in XML
- stylesheets generate Web pages from XML content file

8.3 XSL-FO

XSL-FO specifies *formatting objects*, that are added to a result tree and describe the later formatting

- page layout, areas, frames, indentation,
- colors, fonts, sizes,
- structuring, e.g. lists, tables ...

XSL-FO provides similar concepts as known from \LaTeX .

FO-objects are e.g. (Namespace **fo:**) fo:block, fo:character, display-graphic, float, footnote, inline-graphic, list-block, list-item, list-item-body, list-item-label, multi-case, page-number, page-number-citation, region-before/after, region-body, simple-link, table, table-and-caption, table-body, table-caption, table-cell, table-column, table-footer, table-header, table-row.

- Each of these objects has appropriate attributes.

XSL-FO

- result tree contains *formatting objects* elements
- the result tree is then input to a formatter that generates HTML/L^AT_EX/RTF/PDF etc.
- currently only understood by
 - FOP (originally by James Tauber, now by Apache), a Java program that translates XML documents that include XSL-FO-markup to PDF:
<http://xml.apache.org/fop/>
 - Adobe Document Server (XML → PDF)

8.4 XSLT: Language Design in the XML-World

- XSLT is itself in XML-Syntax
- there is a DTD for XSLT: <http://www.w3.org/TR/xslt#dtd>

⇒ Analogously, there is an XML syntax for XQuery: **XQueryX**,
<http://www.w3.org/TR/xqueryx>, W3C Recommendation since 23 January 2007.

- XSLT uses an own *namespace*, **xsl:....**
- there are several further languages of this kind (programming languages, markup languages, representation languages ...):
XLink, XML Schema,
SOAP (Simple Object Access Protocol)
WSDL (Web Services Description Language)
OWL (Web Ontology Language)
DocBook
... lots of application-specific languages.

8.5 Concepts

(cf. Slide 10)

- XML as an abstract *data model* (Infoset) with an *abstract datatype* (DOM) and several implementations (*physical level*),
- High-level *declarative*, *set-oriented* query language *logical level*: XPath/XQuery
- new: XSLT: transformational language
- two possibilities to define *views*:
 - XQuery: views as queries,
 - XSLT: views by transformations, especially XHTML views to the user.

Chapter 9

XPointer and XLink

- Considered up to now: XML as a data model, data representation, queries
- only single, isolated documents

World Wide **Web**

- references between documents,
- links in HTML: point to a document, sometimes to an *anchor* in a document:
<http://user.informatik.uni-goettingen.de/~may/Mondial/mondial.html#XML>
(the target must be prepared in the remote document with ``),
- browser: when clicking on a link, something happens.

What does that mean for XML? – concepts?

- a language for expressing references: XPointer
- a language for specifying the semantics of references: XLink

9.1 XPointer

- links in HTML: point to a document, sometimes to an *anchor* in a document:
<http://user.informatik.uni-goettingen.de/~may/Mondial/mondial.html#XML>
(the target must be prepared in the remote document with ``),
- Goal in XML: express a pointer to *something* in another XML document.
- possibilities to address individual elements, attributes, or also characters in an XML document:
 - element-, attribute-, comment-, processing instruction nodes,
 - all “information” that can be selected on the monitor by “mousing” can also be addressed by an XPointer.
(independent from borders of elements – can start in the middle of an element and end in the middle of another element).
 - each point directly before or after an element can be addressed.

XPOINTER

- XPointer is a *semantical*, not a syntactical (wrt. the target document) concept. XPointers must be transparent against mechanical changes in the target document (i.e., not “point to the 3rd character in the 6th line in the browser”).
- extends the URL concept with the use of XPath:

XPointer = *url#xpointer-expr*

[http://.../Mondial/mondial.xml#xpointer\(descendant::country\[@car_code="D"\]\)](http://.../Mondial/mondial.xml#xpointer(descendant::country[@car_code=)

- “shorthand pointer”: *url#id*
 - as in HTML: [](#) and [](http://filename#bla) addresses the element that has *id* as its ID-value (DTD: value of an attribute declared as ID)
- full form – “xpointer scheme” (there are also other schemes):
url#xpointer(xpointer-expr)
- For this, XPath is extended with some constructs.
- alternative: element() scheme, e.g. [element\(D\)](#), [element\(/1/4/3\)](#), [element\(D/8/3\)](#) (last: third child of the eight child of the element identified by “D”)

XPointer

- every XPath expression is also an XPointer expression
- $xpath\text{-}expr_1/range\text{-}to(xpath\text{-}expr_2)$ is a pointer, that selects an area in an XML document:

`mondial.xml#xpointer(//country[name="Germany"]/city[1]/range-to(..city[6]))`

selects the area from the 1st to the 6th city of Germany in mondial.xml.

(*not* as set of nodes, but as an area. This can e.g. include changing from one province element to another).

- $string\text{-}range(xpath\text{-}expr, string, m, n)$ selects sequences of characters in documents: for each result of $xpath\text{-}expr$, the first occurrence of $string$ is searched, and the characters from positions m - n are “referenced”.

Markup is ignored in this sequence (including attribute values!)

Remark: since we speak about pointers, the result is not a fragment of an XML document, but simply two positions in a document!

XPointer: EXAMPLES

- Addressing via the id-function:

`mondial.xml#xpointer(id("D"))`

shorthand: `mondial.xml#D`

- robust against changes in the XML document structure,
- requires knowledge about the schema definition (ID-declaration)

- “object-oriented” addressing via semantic “keys”:

`mondial.xml#xpointer(//country[name="Germany"])`

`mondial.xml#xpointer(//country[name="Germany"]//city[name="Berlin"])`

9.2 XLink: World Wide Linking

- extended possibility for specifying hyperlinks.
- relationships between resources (documents, elements, ...) resources can also be programs, images, films, etc.
- – Language: “XLink”
 - Namespace xlink:
- uses (naturally) XPointer

Requirement Analysis

- What “kinds” of references are needed?
Is the functionality of HTML’s <a>-tag sufficient?
- semantics of references?
click? and then?
- ... up to now, XLink is officially only investigated for browsing applications.

SEMANTICS OF EXISTING REFERENCE TYPES: HTML

HTML: ``

- specified in the source document, unidirectional, only one target,
- either the whole page, or to a predefined anchor.
- behavior?
 - standard: when clicked, the target page is shown in the current window.
user-activated, “replace”
 - alternative: when clicked, the target page is shown in a new window.
user-activated, “new”
 - alternative: instead of building up a page, another page is shown in the current window (forwarding)
automatically activated, “replace”
 - alternative: when building up a page in the browser, other pages are shown in small, separate windows
automatically activated, “new”

... sufficient for clicking/browsing, but not for a data model.

HTML:

... is also a "link"!

- specified in the source document, unidirectional, only one non-HTML/XML target,
- behavior?
 - standard: when the page is loaded, the image is embedded at the given position.
automatically activated, "embed"
 - alternative: when building up a page in the browser, show pictures in small, separate windows
automatically activated, "new"

SEMANTICS OF EXISTING REFERENCE TYPES: ID/IDREF

ID/IDREF/IDREFS is already a reference mechanism in XML: Simplest kind of references *inside an XML document*:

- unidirectional, internal to the document, one or more targets
- “Activation”?
 - ... when a query is executed (dereferencing; “user-activated”)

... insufficient for a data model, useless for clicking ...

EXAMPLE-SCENARIOS

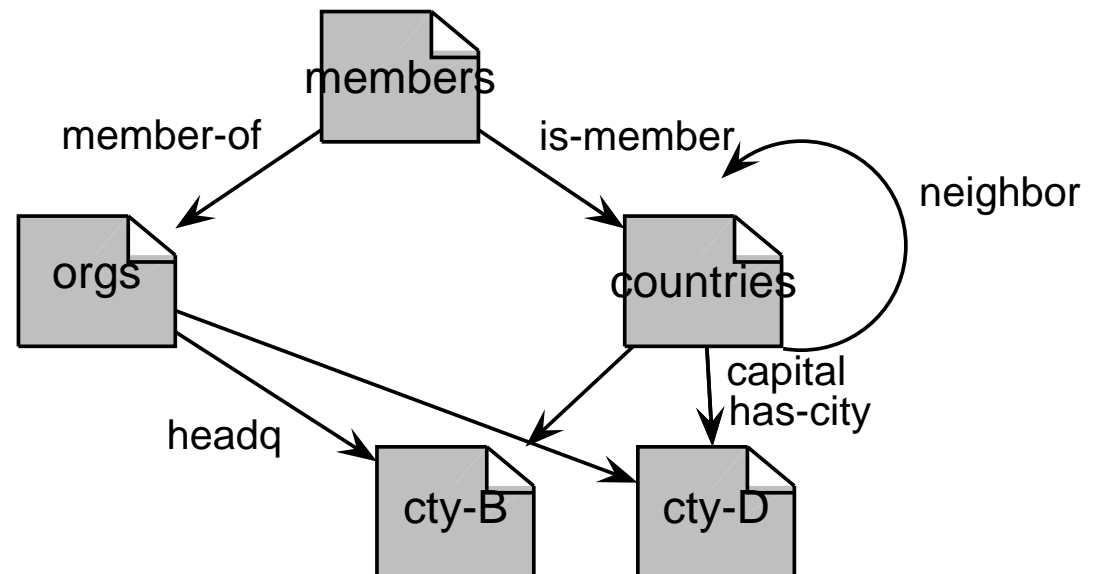
World-Wide-Web

- Web pages
- Hyperlinks
- other kinds of relationships between Web pages

Storage of XML Data in XML (Mondial)

- Distribution over multiple documents

- countries.xml
- cities-*car-code*.xml
(cities and provinces of each country)
- organizations.xml
- memberships.xml



XLINK: BASIC NOTIONS

Resources: XML documents, parts of XML documents, HTML pages, images, films, Web services ...

- local resource: a resource that belongs as a structure to the content of the XLink element itself (or that is the link itself)
- remote resource: a resource that is given by a URI

Examples

- `Göttingen` is a simple link: connects the (local) resource “Göttingen” (string to be clicked) with a (remote) resource located at the URL `www.goettingen.de` (Web page).
- `` is an even simpler link: has no local resource, but points only to a remote one

XLINK: BASIC NOTIONS (CONT'D)

Arcs: directed connections between resources (starting point → endpoint)

- outbound: the starting point is a local resource, the end is a remote resource.
 - ` ...`,
 - country-capital-relationship: a country element is the local resource, and city element is the other, remote, resource.
- inbound: the starting point is a remote resource, the endpoint is a local one.
Inbound-arcs cannot be represented in the same document as their starting point.
- third-party: starting point and endpoint are remote resources.
 - e.g. own linkbase over the Web: each link connects two remote resources (an area of an HTML document with another URL).
 - e.g. memberships of countries in organizations:
 - * each link connects two remote resources, a country and an organisation
 - * n:m-relationship ... see later

XLINK: KINDS OF LINKS AND THEIR SEMANTICS

XLink offers a meta-semantics for describing references that is then used by applications.

- different kinds of references
 - simple: like `...` Or ``
 - links to multiple targets/resources/documents
activate several resources at the same time
DB: a country has several cities
 - the links described above are *inline*-links, i.e., contained in the document itself (outbound arcs).
 - *out-of-line*-links: a user can define connections between (sets of) documents that are owned by somebody else (third-party arcs).
“overlay” own hyperlinks for clicking over the Web
DB: connections between countries and organizations
- timepoint of activation (onLoad, onRequest)
- action (new, replace, embed)

XLINK ELEMENTS

- Element- and attribute names from the *xlink:* namespace
- *Each* element can become a link ...
- ... by adding an *xlink:type* attribute having one of the values defined by XLink, the element is assigned XLink functionality.
- Properties and substructures (chosen from a predefined set of XLink behavior) can then be specified.

<!ELEMENT *linkelement* (*contentmodel*)>

<!ATTLIST *linkelement*

xmlns:xlink CDATA #FIXED "http://www.w3.org/1999/xlink"

xlink:type (simple|extended|locator|resource|arc) #FIXED "..."

... >

- different possibilities for the *content model* (depending on *xlink:type*, some things are required).
- additional attributes (depending on *xlink:type*)

XLINK ELEMENT TYPES

- Definition of arbitrary *link element* types,
- XLink defines 5 basic types of element types:
 - simple: similar to what is known from HTML's ``,
 - extended: multiple targets,
 - locator: is used in extended links for specifying individual remote resources,
 - resource: is used in extended links for specifying individual local resources,
 - arc: is used in extended links for specifying connections between locator- and resource elements.

XLINK -ELEMENTS AND THEIR ATTRIBUTES

Structure

- xlink:type: chooses a link type,
- xlink:href: specification of target(s) (for simple or locator elements) (note that an XPointer can specify multiple targets)
- xlink:label: give names for locators and resources.
- xlink:from, xlink:to: for arcs – references to an xlink:label.

Behavior

- xlink:actuate: specifies, when the link is “activated”:
 - onLoad: when loading/parsing the XML document,
 - onRequest: only when the user does something (e.g. clicking).
- xlink:show: specifies what happens when the link is activated:
 - new: target appears in a new window,
 - replace: the current document is replaced by the target,
 - embed: the target is embedded into the current document.

XLINK ELEMENTS AND THEIR ATTRIBUTES (CONT'D)

Further Attributes (application-specific)

- xlink:title: the user can give a comment about the target.
- xlink:role and xlink:arcrole: allows to group links to roles, e.g., xlink:role="Annotation" for giving additional information to the application.
- xlink:arcrole is mainly used for annotations for mappings to RDF (Resource Description Framework):

object has property *rolename* *value*

(RDF is e.g. used in the Semantic Web)

SIMPLE LINKS

- By setting xlink:type to “simple”, a simple link analogous to the <A>-tag in HTML is defined.
- the xlink:href attribute specifies the target
- xlink:actuate and xlink:show specify the behavior.

Example: <A>-Element in HTML (known behavior)

```
<!ELEMENT A ANY>
<!ATTLIST A xmlns:xlink CDATA #FIXED "http://www.w3.org/1999/xlink"
  xlink:type (simple|extended|locator|resource|arc) #FIXED "simple"
  xlink:href CDATA #REQUIRED
  xlink:actuate (onLoad|onRequest) "onRequest"
  xlink:show (new|embed|replace) "new">
```

A sample element that embeds an HTML fragment when clicking on it:

For getting the whole list, please click

```
<A xlink:href="liste-fragment.html" xlink:show="embed">here</A>.
```

SIMPLE LINK: EXAMPLE

- country/capital as simple link

```
<!ELEMENT country (name,population?, ..., capital,cities, ...)>
<!ELEMENT capital EMPTY>
<!ATTLIST capital xmlns:xlink CDATA #FIXED "http://www.w3.org/1999/xlink"
                xlink:type (simple|extended|locator|resource|arc)
                #FIXED "simple"
                xlink:href CDATA #REQUIRED>
```

In the XML document:

```
<country car_code="D" area="356910">
  <name>Germany</name>
  :
  <capital xlink:type="simple"
    xlink:href="cities-D.xml#xpointer(//city[name='Berlin'])"/>
</country>
```

- the (empty) “capital” element is the local resource, the link is specified between country/capital and Berlin.

Exercise: Define the HTML element with XLink.

XLINK: INLINE EXTENDED LINKS

- Extended links can contain multiple XPointers,
- they have no own href attribute,
- the specification of the targets is done by locator subelements and resource subelements.

```
<!ELEMENT linkelement (... locatorelement* ... resourceelement* ...)>
```

```
<!ATTLIST linkelement
```

```
  xmlns:xlink ... ..
```

```
  xlink:type (simple|extended|locator|resource|arc) #FIXED "extended" >
```

```
<!ELEMENT locatorelement (contentmodel)>
```

```
<!ATTLIST locatorelement
```

```
  xmlns:xlink ... ..
```

```
  xlink:type (simple|extended|locator|resource|arc) #FIXED "locator"
```

```
  xlink:href CDATA #REQUIRED >
```

resourceelement: any element can be made a local resource by xlink:type="resource".

EXAMPLE: EXTENDED LINK IN THE WWW

Pictures in the Web:

```
<!ELEMENT a-loc EMPTY>
<!ATTLIST a-loc xlink:type (...) FIXED "locator"
           xlink:href CDATA>
<ul> <li> <ext-a xlink:type="extended">Goettingen
        <a-loc xlink:type="locator" xlink:href="goe1.jpg">
        <a-loc xlink:type="locator" xlink:href="goe2.jpg">
        <a-loc xlink:type="locator" xlink:href="goe3.jpg">
        </ext-a>
    <li> <ext-a xlink:type="extended">Freiburg
        <bla-loc xlink:type="locator" xlink:href="fr1.jpg">
        <bla-loc xlink:type="locator" xlink:href="fr2.jpg">
        </ext-a>
</ul>
```

when the user (in an XLink-enabled browser) clicks on such an `<ext-a>` element, all corresponding pictures are shown in separate new windows.

EXAMPLE: EXTENDED LINK IN MONDIAL

The country elements can also be regarded as extended links:

- each country element as a whole is a local resource
- with own resources: name etc.
- with references (locators) to other resources.

```
<!ELEMENT country (name,...,capital,encompassed*,border*,cities)>
<!ATTLIST country car_code ID #IMPLIED
                :
                xlink:type CDATA #FIXED 'extended'>
<!ELEMENT capital EMPTY>
<!ATTLIST capital xlink:type CDATA #FIXED 'locator'
                xlink:href CDATA #REQUIRED >
```

Analogously for encompassed, border, provinces, cities.

Exercise 9.1

Give a DTD for these elements (note that encompassed and border have further attributes). □

EXAMPLE (CONT'D)

XML instance:

```
<country xlink:type="extended"
        car_code="D" area="356910" ...>
  <capital xlink:type="locator"
          xlink:href="cities-D.xml#xpointer(//city[name='Berlin'])"/>
  <cities xlink:type="locator"
          xlink:href="cities-D.xml#xpointer(//city)"/>
  :
</country>
```

Exercise 9.2

Complete the XML instance.

□

EXAMPLE (CONT'D)

- the attributes xlink:show and xlink:actuate are not relevant here
- application: not browsing, but queries, transformations etc.
- not considered neither in XPath/XQuery nor in XLink.
- These aspects are investigated in the LinXIS project:
<http://www.dbis.informatik.uni-goettingen/LinXIS>
(see later)
- up to now: only outbound links
 - country/capital is an implicit outbound-arc from the local resource (country) to a city.
 - country/cities references multiple targets (city elements), defines multiple outbound arcs.

XLINK: OUT-OF-LINE-LINKS

Link elements that are not in the document, but in separate documents (i.e., possible to “add” links to other people’s documents):

- expressed by extended link elements with locators and resources; these are equipped with an `xlink:label` attribute.
- in addition to the locator elements (that address the (remote) resources), additional information must be stored:
 - which resources are connected by an arc,
 - and the direction of the connection.

⇒ additional *arc* elements
connect resources/locators by `xlink:from` and `xlink:to` attributes.

XLINK: OUT-OF-LINE-LINKS

- element content allows for subelements of locator element types (as above) and subelements of arc element types that describe relationships between locator elements:

```
<!ELEMENT linkelement ((locatorelement*|resourceelement*|arcelement)*)>
```

```
<!ATTLIST linkelement as above >
```

```
<!ELEMENT locatorelement as above >
```

```
<!ATTLIST locatorelement as above
```

```
  type, href, label: NMTOKEN, title, role >
```

```
<!ELEMENT arcelement (contentmodel)>
```

```
<!ATTLIST arcelement
```

```
  xmlns:xlink CDATA #FIXED "http://www.w3.org/1999/xlink"
```

```
  xlink:type (simple|extended|locator|resource|arc) #FIXED "arc"
```

```
  xlink:from NMTOKEN #IMPLIED
```

```
  xlink:to NMTOKEN #IMPLIED
```

```
  xlink:arcrole CDATA #IMPLIED >
```

- Note: xlink:label is not an ID; there can be several elements with the same label.

EXAMPLE OUT-OF-LINE-LINKS: WEB

Adding an own link base in addition to the ``-entries in HTML Web documents:

- two kinds of locator elements:

1. XPointers on keywords and sentences where the link should start from:

```
<univis xlink:type="locator" xlink:label="univis-ssdxml"
  xlink:href="http://www.univis.goe/somepath#xpointer(//h1[text()='...'])"/>
```

2. URLs of target pages:

```
<teaching xlink:type="locator" xlink:label="ssdxml-homepage"
  xlink:href="http://www.cs.uni-goettingen.de/teaching/ssdxml.html"/>
```

- Arc elements:

```
<arc xlink:type="arc"
  xlink:from="univis-ssdxml" xlink:to="ssdxml-homepage"/>
```

- when using an XLink-enabled browser that has access to the link base, there is a link on the Univis page from the headline to the XML lecture.

EXAMPLE OUT-OF-LINE-LINKS: MONDIAL

Memberships of countries in organizations:

```
<!ELEMENT memberships (country*,organization*,membership*)>
```

```
<!ATTLIST memberships xlink:type (...) #FIXED "extended" >
```

```
<!ELEMENT country EMPTY>
```

```
<!ATTLIST country xlink:type (...) #FIXED "locator"
```

```
xlink:href CDATA #REQUIRED
```

```
xlink:label NMTOKEN #REQUIRED >
```

```
<!ELEMENT organization EMPTY>
```

```
<!ATTLIST organization xlink:type (...) #FIXED "locator"
```

```
xlink:href CDATA #REQUIRED
```

```
xlink:label NMTOKEN #REQUIRED >
```

```
<!ELEMENT membership EMPTY>
```

```
<!ATTLIST membership xlink:type (...) #FIXED "arc"
```

```
xlink:from NMTOKEN #IMPLIED
```

```
xlink:to NMTOKEN #IMPLIED
```

```
membership_type CDATA #REQUIRED >
```

EXAMPLE OUT-OF-LINE-LINKS: MONDIAL

```
<memberships>
  <country xlink:label="B" xlink:type="locator"
    xlink:href="#../countries.xml#xpointer(//country[@car_code='B'])/>
  <country xlink:label="D" xlink:type="locator"
    xlink:href="#../countries.xml#xpointer(//country[@car_code='D'])/>
  <organization xlink:label="org-EU" xlink:type="locator"
    xlink:href="#../organizations.xml#xpointer(//organization[@abbrev='EU'])/>
  <organization xlink:label="org-UN" xlink:type="locator"
    xlink:href="#../organizations.xml#xpointer(//organization[@abbrev='UN'])/>
  <membership xlink:from="B" xlink:to="org-EU" xlink:type="arc"
    membership_type="member"/>
  <membership xlink:from="B" xlink:to="org-UN" xlink:type="arc"
    membership_type="member"/>
</memberships>
```

The diagram illustrates the relationships between XML elements. It shows a sequence of elements: two country elements (labeled 'B' and 'D'), two organization elements (labeled 'org-EU' and 'org-UN'), and two membership elements. The membership elements use 'xlink:from' and 'xlink:to' attributes to reference the country and organization elements respectively. The 'xlink:from' attribute of the first membership element points to the 'xlink:label' attribute of the 'B' country element. The 'xlink:to' attribute of the first membership element points to the 'xlink:label' attribute of the 'org-EU' organization element. The 'xlink:from' attribute of the second membership element points to the 'xlink:label' attribute of the 'B' country element. The 'xlink:to' attribute of the second membership element points to the 'xlink:label' attribute of the 'org-UN' organization element. The 'xlink:href' attributes of the country and organization elements point to their respective locations in an external XML document.

SEMANTICS OF ARCS

- In case that all xlink:label in an extended link element are unique, each arc element stands for the unique relationship given by the xlink:from and xlink:to attributes.
- In case that the labels are not unique, every arc stands for all relationships between pairs of locators that have the corresponding from- and to-labels.
- an arc that has no xlink:to attribute, stands for a connection to each locator (analogously for from).
- an arc that has neither from nor to stands for all possible relationships.

XLINK: USAGE

Browsing: obvious. xlink:show and xlink:actuate

- **W3C Amaya** (<http://www.w3.org/Amaya>): partially understands XLink and is open-source
 - use XLink for annotations to Web pages (→ RDF).
- queries against XML data sources:
 - The *W3C XML Query Requirements* state that the query language must support queries over references. The XLink/XQuery combination does not (yet) satisfy this.
 - behavior of XPath and XLink has not yet been considered in the W3C documents:
 - there is even no *data model* for XLinks
 - currently: requires real programming for resolving XLink elements and evaluating the references dynamically.

9.3 XInclude: Database-Style Use of XPointer

Include-elements are *replaced* by the corresponding included items:

```
<xi:include parse="xml|text" href="url" xpointer="xpointer"/>
```

- no browsing semantics (XHTML: include must be resolved when loading)
- query/database semantics: obvious

```
<country xlink:type="extended" car_code="D" area="356910" ...>  
  <xi:include parse="xml" href="mondial-D-cities.xml" xpointer="//city"/>  
</country>
```

becomes

```
<country xlink:type="extended" car_code="D" area="356910" ...>  
  <city><name>Berlin</name> ... </city>  
  <city><name>Stuttgart</name> ... </city>  
  :  
</country>
```

- resolve inclusion when loading
- resolve inclusion on-demand when querying

9.4 The LinXIS Project – Linked XML Information Sources

Research project in the DBIS group at Göttingen:

... combine XLink with XInclude-style functionality

- extend XLink with data model semantics: insert referenced targets
- **dbxlink**: namespace for specifying the data model

```
<country car_code="D" area="356910" ...>
  <capital xlink:type="simple" dbxlink:transparent="keep-element insert-contents"
    xlink:href="cities-D.xml#xpointer(//city[name='Berlin'])"/>
</country>
```

```
<country car_code="D" area="356910" ...>
  <capital><name>Berlin</name><population>3472009</population></capital>
</country>
```

- **query**: `document(countries.xml)//country[name="Germany"]/capital/population)` resolves the xlink and “jumps” automatically to the target of the reference to `mondial/cities-D.xml#xpointer(//city[name="Berlin"])`
- <http://www.dbis.informatik.uni-goettingen.de/LinXIS>

Chapter 10

XML Schema

10.1 Motivation

- Database area: schema description
 - cf. SQL datatypes, table schemata
 - constraints
- Programming languages: typing – *real* typing – means: theory
 - every expression (query, variable etc) can be assigned with a type
 - structural induction
 - static typechecking for queries/programs/updates
 - validation of resulting structures wrt. target DTD/Schema

XML QUERY FORMAL SEMANTICS: OVERVIEW

Every (query) expression is assigned with a semantics

Static Semantics

Given a static environment, an expression is of a certain type:
(static env.: namespace decl, typedefs, type decls. of variables)

- $statEnv \vdash Expr : Type$

Dynamic Semantics

Given a dynamic environment, an expression yields a certain result:
(dynamic env.: context node, size+position in context, variable bindings)

- $dynEnv \vdash Expr \Rightarrow Value$
(equivalent to “classical” notation: $[[Expr]]_{dynEnv} = Value$)

... both defined by structural induction.

(for short example: show 2.1.5 and “if” in 4.10 of W3C XQFS document)

XML QUERY DATA TYPES

... by examples:

```
define type coordinates { -- in any order
  element longitude of type xs:float &
  element latitude of type xs:float}
```

```
define type city { -- sequence
  attribute country of type xs:string,
  attribute province of type xs:string?, -- optional
  element name of type xs:string,
  element population of type { -- anonymous
    attribute year of type xs:decimal,
    xs:string }
  element coordinates of type { -- anonymous
    element longitude of type xs:float &
    element latitude of type xs:float} }
```

... similar to DTD expressions extended by primitive datatypes

XML QUERY DATA TYPES (CONT'D)

XML type theory:

- operations on types (e.g. “union”): result type of a query that yields either a result of type a or type b
- derivation of new types by
 - additional constraints
 - additional content
- constraints: does the derived result type for some expression guarantee that some conditions hold?
- containment of types: is the derived result type for some expression covered by a certain target type?
(static type checking of programs)
- can e.g. be applied for query and storage optimization, indexing etc.

REQUIREMENTS ON AN XML SCHEMA LANGUAGE

- requirement: a schema description language for the user that is *based* on these types: (usage is optional – XML is self-describing)
- DTD: heritage of SGML; database-typical aspects are not completely supported (datatypes [everything is CDATA/PCDATA], cardinalities); but: order, iteration.
- DTD: syntax is not in XML.

⇒ better formalism for representing schema information

- XML syntax → easy to process
- more detailed information as in the DTD
- database world: datatypes with derived types, constraints etc.

XML SCHEMA: IDEAS

- no actually new concepts (in contrast to the definition of the object-oriented model), but ...
- combination of the power of previous schema languages:
 - datatype concepts from the database area (SQL)
 - idea of complex object types/classes from the OO area
 - structured types from the area of tree grammars (e.g. DTD)
- new in contrast to DTDs: difference between element types and elements: <country> elements are of the type countryType.

⇒ very complex

10.2 XML Schema: Design

Using XML syntax and a verbose formalism, XML Schema uses a very detailed and systematic approach to type definitions and -derivations.

- only a few primitive, atomic datatypes
- other *simple types* are derived from these by *restriction*,
- *complex types* with text-only contents are derived by *extension* from simple types,
- other (*complex types*) are derived by *restriction* from a general *anyType* (cf. class *object* in OO),
- these types are then used for declaring elements.

XML SCHEMA: THE STANDARD

The XML Schema Recommendation (since May 2001) consists of 2 parts:
(note: XML Schema 1.1 work is in progress [2007])

- Part 2: “Datatypes”
 - Definition of *simple types*:
have no attributes and no element content; are used only for text content and as attribute values.
- Part 1 “Structures”:
 - Definition of structured datatypes (*complex types*):
with subelements and attributes; are used as element types.
 - * names/types of the subelements and attributes
 - * order of the subelements
 - Definition of elements using the complex types.
- many syntax definitions
- Part 0: “Primer” (<http://www.w3.org/TR/xmlschema-0/>) explains and motivates the concepts.

USAGE OF XML SCHEMA

- understand concepts and ideas (XML Schema Primer, lecture)
- apply them in practice
- lookup for syntax details in the W3C documents
- make experiences

XML SCHEMA DOCUMENTS

An XML-Schema document consists of

- a preamble and
- a set of definitions and declarations

```
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">  
  content  
</xs:schema>
```

content uses the following kinds of *schema components*:

- datatype definitions (simple types and complex types)
- attribute declarations
- element declarations
- miscellaneous ...

10.3 Datatypes

Datatypes are seen as triples:

- range (possible values, cardinality, equality and ordering),
- lexical representation by *literals* (e.g. 100 also as 1.0E2 and 1.0e+2)
- set of properties

The set of possible values can be specified in different ways:

- extensional (enumeration)
- intensional (by axioms)
- restriction of another domain
- construction from other domains (lists, sets ...)

DATATYPES

Datatypes are characterized by their domain and *properties (called facets)* in multiple independent “dimensions”. The facets describe the differences between datatypes.

The basic facets (present for each datatype) are

- equality,
- order relation,
- upper and lower bound,
- cardinality (finite vs. countable infinite),
- numerical vs. non-numerical.

DATATYPES

- primitive datatypes (predefined)
- generated datatypes (derived from other datatypes). This can happen by aggregation (lists) or restriction of another datatype.
- Primitive predefined types in XML Schema:
 - string (with many subtypes: token, NMTOKEN),
 - boolean (lexical repr.: true, false),
 - float, double,
 - decimal (with several subtypes: integer etc.),
 - duration, time, dateTime, ...
 - base64Binary, hexBinary
 - anyURI (Universal Resource Identifier).
- generated predefined types:
 - integer, [non]PositiveInteger, [non]NegativeInteger, [unsigned](long|short|byte)

XML-SPECIFIC DATATYPES

There are some XML-specific datatypes (subtypes of string) that are defined based on the basic XML recommendation. They are only used for attribute types (atomic and list types):

- NMTOKEN (restriction of string according to the definition of XML tokens),
- NMTOKENS derived from NMTOKEN by list construction,
- IDREF/IDREFS analogously,
- Name: XML Names,
- NCName: non-colonized names,
- language: language codes according to RFC 1766.

CONSTRAINING FACETS

By specifying constraining facets, further datatypes can be derived:

- for sequences of characters: length, minlength, maxlength, pattern (by regular expressions);
- for numerical datatypes: maxInclusive, minInclusive, maxExclusive, minExclusive,
- for lists: length, minLength, maxLength
- for decimal datatypes: totalDigits (number of digits), fractionDigits (number of positions after decimal point);
- enumeration (definition of the possible values by enumeration),

... for a description of all details, see the W3C XMLSchema Documents.

GENERATION OF SIMPLE DATATYPES

Simple datatypes can be derived as `<simpleType>` from others:

Derivation by Restriction

Restriction of a base type (i.e., specification of further restricting facets):

```
<xs:simpleType name="name">  
  <xs:restriction base="simple-type">  
    facets  
  </xs:restriction>  
</xs:simpleType>
```

```
<xs:simpleType name="carcodeType">  
  <xs:restriction base="xs:string">  
    <xs:minLength value="1"/>  
    <xs:maxLength value="3"/>  
    <xs:pattern value="[A-Z]+"/>  
  </xs:restriction>  
</xs:simpleType>
```

Derivation by Restriction

Example:

```
<xs:simpleType name="longitudeType"  
  <xs:restriction base="xs:decimal">  
    <xs:minExclusive value="-180"/>  
    <xs:maxInclusive value="180"/>  
  </xs:restriction>  
</xs:simpleType>
```

```
<xs:simpleType name="latitudeType"  
  <xs:restriction base="xs:decimal">  
    <xs:minInclusive value="-90"/>  
    <xs:maxInclusive value="90"/>  
  </xs:restriction>  
</xs:simpleType>
```

defines two derived simple datatypes.

Remark: Usage of SimpleTypes

- as attributes (details later):

```
<xs:attribute name="car_code" type="carcodeType"/>
```

can e.g. be used in

```
<country car_code="D"> ... </country>
```

- as elements (details later):

```
<xs:element name="longitude" type="longitudeType"/>
```

can e.g. be used in

```
<longitude>48</longitude>
```

Only if also attributes are required, a `<complexType>` must be defined.

Derivation by Restriction: Enumeration

- Enumeration of the allowed values.

Example (from XML Schema Primer)

```
<xs:simpleType name="USState">  
  <xs:restriction base="xs:string">  
    <xs:enumeration value="AK"/>  
    <xs:enumeration value="AL"/>  
    <xs:enumeration value="AR"/>  
    <!-- and so on ... -->  
  </xs:restriction>  
</xs:simpleType>
```

- so far, this functionality is similar to what could be done in SQL by attribute types and integrity constraints.
- additionally:
 - “multi-valued” list types (but still simple types)
 - complex types

Derivation as List Types

```
<xs:simpleType name="name">  
  <xs:list itemType="simple-type">  
    facets      <!-- optional -->  
  </xs:list>  
</xs:simpleType>
```

- *simpleType* must be a non-list type,
- facets of the list (e.g., maxLength, minLength, pattern) can be defined by subelements

Example

Datatype for a list of country codes:

```
<xs:simpleType name="countrylist">  
  <xs:list itemType="carcodeType"/>  
</xs:simpleType>  
  
<xs:attribute name="neighbors" type="countrylist"/>  
  
for <country neighbors="NL L F CH ..."> ... </country>
```

Derivation as Union Types

- Analogously union of sets with `xs:union` and `@xs:memberTypes`.

Component of a data type for postal addresses for US suppliers: send e.g., to D 37075 Göttingen (car code), or CA 94065 Redwood (US State Code)

```
<xs:simpleType name="stateOrCountry">  
  <xs:union memberTypes="carcodeType USState"/>  
</xs:simpleType>
```


10.4 Complex Datatypes

Complex datatypes can be derived from others by `<complexType>`. They describe

- `ContentType`: `simpleContent` or `complexContent` (an already defined `simpleType` or an own `ContentModel`),
- attributes.

Different possibilities:

- by extension from a simple datatype (adding attributes, making an element type out of a simple type)
- by restriction from another datatype (restriction of its components or its structure),
- completely new definition (formally, a restriction of `base="xs:anyType"`)

With these datatypes, element types can be defined later.

COMPLEX DATATYPES

Complex datatypes are defined by the following properties:

- name,
`<xs:complexType name="name"> ... </xs:complexType>`
- kind of content (simple or complex; mixed)
`<xs:simpleContent> ... </>`
`<xs:complexContent [mixed="true"]> ... </>`
- derivation method (extension or restriction),
`<xs:extension base="typename"> ... </>`
`<xs:restriction base="typename"> ... </>`
- attribute declarations
`<xs:attribute name="name" type="typename"/>`
- structure of content model
... a bit more complex ...

COMPLEX DATATYPES: TYPE WITH ATTRIBUTES

Population: text content and an attribute:

```
<population year="1997">130000</population>
```

- take the simpleType for the text content and extend it with an attribute:

```
<xs:complexType name="population">  
  <xs:simpleContent>  
    <xs:extension base="xs:nonNegativeInteger">  
      <xs:attribute name="year" type="xs:nonNegativeInteger"/>  
    </xs:extension>  
  </xs:simpleContent>  
</xs:complexType>
```

- the complexType can have (but does not necessary have) the same name, as the element to be defined later.

COMPLEX DATATYPES: EMPTY ELEMENT TYPES

Border: only two attributes:

```
<border country="F" length="500"/>
```

- take the base type anyType (arbitrary complexContent and attributes) and restrict it:

```
<xs:complexType name="border">  
  <xs:complexContent>  
    <xs:restriction base="anyType">  
      <xs:attribute name="country" type="xs:IDREF"/>  
      <xs:attribute name="length" type="xs:decimal"/>  
    </xs:restriction>  
  </xs:complexContent>  
</xs:complexType>
```

COMPLEX DATATYPES: ARBITRARY ELEMENT TYPES

Newly defined element types with `<complexContent>` are usually defined by

```
<xs:complexType name="...">
  <complexContent>
    <xs:restriction base="anyType">
      <!-- type definition: content model and attributes -->
    </xs:restriction>
  </complexContent>
</xs:complexType>
```

As an abbreviation, it is allowed to omit `<complexContent>` and `<xs:restriction base="anyType">`:

```
<xs:complexType name="border">
  <xs:attribute name="country" type="xs:IDREF"/>
  <xs:attribute name="length" type="xs:decimal"/>
</xs:complexType>
```

COMPLEX DATATYPES: ARBITRARY ELEMENT TYPES

- element types with complex content use (nested) structure-defining elements (called *Model Groups*):
 - `<xs:sequence> ... </xs:sequence>`
 - `<xs:choice> ... </xs:choice>`
 - `<xs:all> ... </xs:all>`
 (“all” with some restrictions - only top-level, no substructures allowed)
- inside, the allowed element types are specified:
`<xs:element name=“name” type=“typename”/>`
- note: even if only one type of subelements is contained, one of the above must be used around it.
- note: the XML Schema definition requires to list the content model specification before the attributes.

COMPLEX DATATYPES: ARBITRARY ELEMENT TYPES

Example

Definition of the type of the <city> elements:

```
<xs:complexType name="city">
  <xs:sequence>
    <xs:element name="name" type="xs:string"/>
    <xs:element name="population" type="population"/>
    <xs:element name="longitude" type="longitude"/>
    <xs:element name="latitude" type="latitude"/>
  </xs:sequence>
  <xs:attribute name="country" type="xs:IDREF"/>
</xs:complexType>
```

... so far, the basic concepts.

10.5 Composing Declarations

FURTHER ATTRIBUTES OF ATTRIBUTE DEFINITIONS

- use (optional, required, prohibited)
default is “optional”
- default (same as in DTD: attribute is added if not given in the document)
- fixed (same as in DTD)

FURTHER ATTRIBUTES OF SUBELEMENT DEFINITIONS

- minOccurs, maxOccurs: default 1.
- <default value=“*value*”/> (bit different from attribute default): if the element is given in a document with empty content, then the default contents *value* is inserted.
In case that an element is not given at all, no default is used.
- <fixed value=“*value*”/>: analogous.

Examples: later.

GLOBAL ATTRIBUTE- AND ELEMENT DEFINITIONS

... up to now, arbitrary element *types* have been defined.

At least, for the root element, a separate element declaration is needed.

- `<xs:attribute>` and `<xs:element>` elements can not only occur inside of `<complexType>` elements, but can also be global.
- as global declarations, they must not contain specifications of `@use`, `@maxOccurs`, or `@minOccurs`.
- global declarations can then also be used in type definitions by `@ref`.
Then, they are allowed to have `@use`, `@maxOccurs` and `@minOccurs`.
- especially useful if the same element type is used several times.

EXAMPLE

```
<xs:element name="city" type="city"/>
<xs:element name="name" type="xsd:string"/>
<xs:attribute name="car_code" type="carcodeType"/>

<xs:complexType name="country">
  <xs:sequence>
    <xs:element ref="name"/> -- minOccurs and minOccurs default 1
    <xs:element ref="city" minOccurs="unbounded"/>
    <xs:element ref="border" minOccurs="0" minOccurs="unbounded"/>
  </xs:sequence>
  <xs:attribute ref="car_code" use="required"/>
</xs:complexType>
<xs:element name="country" type="country"/>

<xs:complexType name="mondial">
  <xs:sequence>
    <xs:element ref="country" minOccurs="unbounded"/>
    :
  </xs:sequence>
</xs:complexType>
```

ANONYMOUS, LOCAL TYPE DEFINITIONS

- Instead of

```
<xs:element name="name" type="typename"/>
```

```
<xs:attribute name="name" type="typename"/>
```

anonymous, local type definitions in the content of such elements are allowed:

```
<xs:complexType name="city">
  <xs:sequence>
    <xs:element name="name" type="xs:string"/>
    <xs:element name="population">
      <xs:simpleContent>
        <xs:extension base="xs:nonNegativeInteger">
          <xs:attribute name="year" type="xs:nonNegativeInteger">
        </xs:extension>
      </xs:simpleContent>
    </xs:element>
  </xs:sequence>
  <xs:attribute name="country" type="xs:IDREF"/>
</xs:complexType>
```

LOCAL DECLARATIONS

- `<complexType>` declarations define local *symbol spaces*, i.e., the same attribute/element names can be used in different complex datatypes with different specifications of result-datatypes (this is not possible in DTDs; cf. country/population and city/population elements)

Using global types:

```
<xs:complexType name="countrypop"> ... without @year ... </xs:complexType>
```

```
<xs:complexType name="citypop"> ... with @year ... </xs:complexType>
```

```
<xs:complexType name="countryType">
```

```
  :
```

```
    <xs:element name="population" type="countrypop"/>
```

```
</xs:complexType>
```

```
<xs:complexType name="cityType">
```

```
  :
```

```
    <xs:element name="population" type="citypop"/>
```

```
</xs:complexType>
```

Local Declarations (Cont'd)

Using local “population” types:

```
<xs:complexType name="countryType">
  <xs:complexType name="pop"> ... without @year ... </xs:complexType>
  :
  <xs:element name="population" type="pop"/>
</xs:complexType>

<xs:complexType name="cityType">
  <xs:complexType name="pop"> ... with @year ... </xs:complexType>
  :
  <xs:element name="population" type="pop"/>
</xs:complexType>
```

ATTRIBUTE GROUPS

Groups of attributes that are used several times can be defined, named and then reused:

```
<xs:attributeGroup name="groupname">  
  attributedefs  
</xs:attributeGroup>
```

```
<xs:complexType name="name" ...>  
  :  
  <xs:attributeGroup ref="groupname"/>  
</xs:complexType>
```

- group definitions can also be nested ...

CONTENT MODEL GROUPS

In the same way, parts of the content model can be predefined:

```
<xs:group name="groupname">
  modelgroupdef
</xs:group>

<xs:complexType name="name" ...>
  :
  <xs:group ref="groupname"/>
</xs:complexType>
```

Exercise 10.1

Use the following group definitions in your MONDIAL schema:

- an attribute group for (country, province) in city, lake, mountain etc.
- a content model group for (longitude, latitude)

□

PRACTICAL ISSUES: XSI:SCHEMALOCATION

In addition to use separate separate .xsd and .xml files (call e.g. saxonXSD bla.xml bla.xsd), the XML Schema can be identified in the XML instance:

- simple things without namespace: the `xsi:noNamespaceSchemaLocation` attribute gives the URI or local file path of the XML Schema file:

```
<mondial xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:noNamespaceSchemaLocation="mondial.xsd"> ... </mondial> <!-- local -->
<mondial xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:noNamespaceSchemaLocation="http://.../mondial.xsd"> ... </mondial>
```

- when a namespace is used: declare the namespace, and the `xsi:schemaLocation` attribute is of the form `xsi:schemaLocation="namespace uri-of-xsd-file"`:

```
<mon:mondial xmlns:mon="http://www.semwebtech.org/Mondial"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.semwebtech.org/Mondial
  http://www.semwebtech.org/Mondial/mondial.xsd">
  ... </mon:mondial>
```

- if a document uses several namespaces, several `xsi:schemaLocations` can be given; also inside of inner elements.

10.6 Integrity Constraints

XML Schema supports three further kinds of integrity constraints (*identity constraints*):

- [unique](#), [key](#), [keyref](#)

that have *very* strong similarities with the corresponding SQL-concepts:

- a name,
- a *selector*: an XPath expression, e.g. `//city`, that describes the set of elements for which the condition is specified (stronger than SQL: [relative to the instance of the element type where the spec is a child of](#)),
- a list of fields (relative to the result of the selector), that are subject to the condition,
- for `keyref`: the name of a key definition that describes the corresponding referenced key.

More expressive than ID/IDREF:

- not only document-wide keys, but can be restricted to a set of nodes (by type, and by subtree),
- multiple fields; can not only contain attributes, but also (textual) element content,
- but not applicable to IDREFS (then, e.g., “D NL B ...” would be seen as a single value).

INTEGRITY CONSTRAINTS

- are subelements of an element type. The scope of them is then each instance of that element type (e.g., allows for having a key amongst all cities of a given country, and keyrefs in that country only referring to such cities)
- document-wide: define them for the root element type.

```
<xs:unique name="..." >  
  <xs:selector xpath="xpath-expr"/>  
  <xs:field xpath="xpath-expr"/> ... <xs:field xpath="xpath-expr"/>  
</xs:unique>
```

```
<xs:key name="name" >  
  <xs:selector xpath="xpath-expr"/>  
  <xs:field xpath="xpath-expr"/> ... <xs:field xpath="xpath-expr"/>  
</xs:key>
```

```
<xs:keyref name="..." refer="name" >  
  <xs:selector xpath="xpath-expr"/>  
  <xs:field xpath="xpath-expr"/> ... <xs:field xpath="xpath-expr"/>  
</xs:keyref>
```

INTEGRITY CONSTRAINTS: EXAMPLE

```
<xs:element name="mondial">
  <xs:complexType>
    <xs:element ref="country" maxOccurs="*" />
    :      <!-- with <border country="..." /> subelements -->
  </xs:complexType>

  <xs:key name="countrykey" <-- key amongst all countries -->
    <xs:selector xpath="country" /> <!-- range: unique amongst all countries -->
    <xs:field xpath="@car_code" /> <!-- is the field @car_code -->
  </xs:key>

  <xs:keyref name="bordertocountry" refer="countrykey">
    <xs:selector xpath="."/ /> <!-- for all border elements, -->
    <xs:field xpath="@country" /> <!-- the @country attr refs to a country -->
  </xs:keyref>
</xs:element>
```

INTEGRITY CONSTRAINTS: EXAMPLE

```
<?xml version="1.0" encoding="UTF-8"?>
<countries-and-cities
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:noNamespaceSchemaLocation="keys.xsd">
  <country code="D">
    <neighbor code="A"/>
    <neighbor code="CH"/>
    <county code="FR" name="Freiburg">
      <neighbor code="VS"/>
      <city>Freiburg</city>
    </county>
    <county code="VS" name="Villingen-Schwenningen">
      <neighbor code="FR"/>
      <city>Villingen</city>
    </county>
    <county code="D" name="Duesseldorf"/>
  </country>
  <country code="CH">
    <neighbor code="D"/>
    <neighbor code="A"/>
    <county code="FR" name="Fribourg">
      <neighbor code="VS"/>
      <city>Fribourg</city>
    </county>
    <county code="VS" name="Valais/Wallis">
      <neighbor code="FR"/>
      <city>Sion</city>
    </county>
    <county code="VD" name="Vaud/Waadts">
      <neighbor code="FR"/>
      <neighbor code="VS"/>
    </county>
  </country>
  <country code="A"/>
</countries-and-cities>
```

[Filename: XMLSchema/keys.xml]

Integrity Constraints: Example (Cont'd)

```
<?xml version="1.0" encoding="UTF-8"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
  <xs:element name="countries-and-cities">
    <xs:complexType>
      <xs:sequence>
        <xs:element ref="country" maxOccurs="unbounded"/>
      </xs:sequence>
    </xs:complexType>
    <xs:key name="countrykey"> <!-- all countries -->
      <xs:selector xpath="//country"/>
      <xs:field xpath="@code"/>
    </xs:key>
    <xs:keyref name="neighbortocountry" refer="countrykey">
      <xs:selector xpath="//country/neighbor"/>
      <xs:field xpath="@code"/>
    </xs:keyref>
  </xs:element>
  <xs:element name="country">
    <xs:complexType>
      <xs:sequence>
        <xs:element ref="neighbor" minOccurs="0" maxOccurs="unbounded"/>
        <xs:element ref="county" minOccurs="0" maxOccurs="unbounded"/>
      </xs:sequence>
      <xs:attribute name="code" type="xs:string"/>
    </xs:complexType>
    <xs:key name="countykey"> <!-- key is local to the country -->
      <xs:selector xpath="//county"/>
      <xs:field xpath="@code"/>
    </xs:key>
    <xs:keyref name="neighbortocounty" refer="countykey"> <!-- local in the country -->
      <xs:selector xpath="//county/neighbor"/>
      <xs:field xpath="@code"/>
    </xs:keyref>
  </xs:element>
  <xs:element name="county">
    <xs:complexType>
      <xs:sequence>
        <xs:element ref="neighbor" minOccurs="0" maxOccurs="unbounded"/>
        <xs:element ref="city" minOccurs="0" maxOccurs="unbounded"/>
      </xs:sequence>
      <xs:attribute name="code" type="xs:string"/>
      <xs:attribute name="name" type="xs:string"/>
    </xs:complexType>
  </xs:element>
  <xs:element name="city" type="xs:string"/>
  <xs:element name="neighbor">
    <xs:complexType>
      <xs:attribute name="code" type="xs:string"/>
    </xs:complexType>
  </xs:element>
</xs:schema>
```

[Filename: XMLSchema/keys.xsd]

USE OF KEY/KEYREF

- allows for local keys and multi-field keys.
- queries can only be stated by joins (as in SQL); then local keys are only of limited use.
- no multivalued keyrefs as in IDREFS. Each reference must be in a separate subelement.

GENERAL ASSERTIONS

Assertions between attributes and/or subelements on instances of an element type:

- children of a complexType declaration,
- `<xs:assert test="xpath-based test"/>`

10.7 Applications and Usage of XML Schema

- (simple) datatype arithmetics and reasoning (numbers, dates)
The simpleTypes and restrictions are also used in the Semantic Web languages RDF/RDFS/OWL.
- specification of allowed structure: validation of documents

The information about a class of documents is also often used:

- derive an efficient mapping to relational storage (cf. Slide 577)
- definition of indexes (over keys and other properties)
- the type definitions can be used to derive corresponding Java Interfaces (JAXB; (cf. Slide 451))
 - get/set methods for properties,
 - automatical mapping between Java and XML serialization,
 - classes that add behavior can then be programmed by extending the interfaces.

Chapter 11

Algorithms and APIs

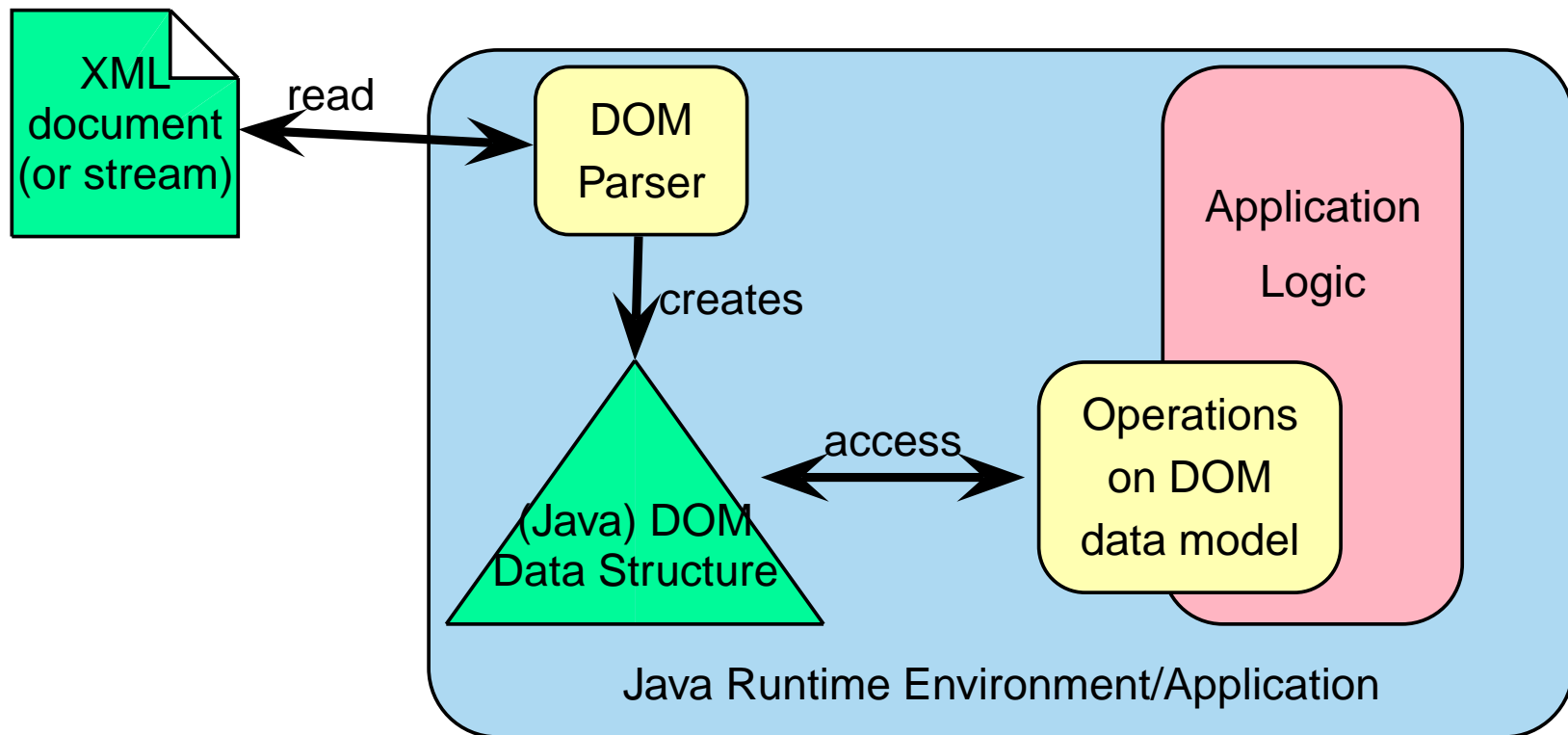
- XML as a data structure:
 - *abstract datatype* with API: DOM
 - (mainly main-memory) implementations; used e.g. in Java applications
 - low-level API with variable-based access
- Databases?
 - high-level API: XPath, XQuery
 - mapping to relational model (Oracle, IBM DB2) or ObjectTypes (Oracle, DB2)
 - “Native” storage: Software AG-Tamino
 - classical database functionality: multiuser, transactions, recovery

Algorithms and APIs (Cont'd)

- XML as Data Exchange Format in Web Services
 - serialize application objects as XML
 - SOAP: generic [not discussed in this course]
 - JAXB: "model-aware" infrastructure
- Stream Processing:
 - XML data transfer as sequence of events
 - SAX (Simple Application Interface for XML), StAX (Streaming API for XML)

11.1 DOM

- DOM (Document Object Model) defines a platform- and language-independent object-oriented *interface* (i.e., an *abstract datatype*) for generating, processing and manipulating XML data.



DOM

- DOM is a specification of an interface/abstract datatype for the XML data model, *not a* data model and *not a* programming language!
- implementations in Java, C++, etc; usually main-memory-based; specialized Java interface definitions:
 - this course: JDOM `jdom.jar`, `org.jdom.*`
 - another alternative: `dom4j`
 - not recommended: `org.w3c.dom.*` (the plain `dom` is an implementation that exists in nearly all programming languages and does not make use of Java's advantages);
- language base of the DOM specification: OMG-IDL
- Main-memory-based: only for relatively small application programs (most of the “lightweight”-tools used in the course are internally based on DOM)

DOM: PRINCIPLES

- only one document in a DOM
- step-by-step-access to the data:
based on variable assignments in the surrounding imperative/object-oriented programming language and on iterators (cf. proceeding in the [network data model](#)):
 - document: represents the complete document,
 - * Query-Methods, e.g. `NodeList getElementByName(string)`
 - class “Node”: `getNodeTypes()`, `getChildren()`, `getFirstChild()`, `getNextSibling()`, `getParentNode()`, ...
 - class “Element”: `getName()`, `getAttributes()`, `getContent()`, ...
 - class “Attribute”: `getName()`, `getValue()`, ...
 - corresponding methods for generating and changing nodes.
- additionally, XPath and XSLT can be applied to instances of Document and Element;
- based on DOM, XPath and XQuery can be implemented (cf. Apache Xerces (XML/DOM)/Xalan (XSLT)/Xindice (DB))
- often inefficient (no indexes, query optimization)

DOM – sample code fragment: Stepwise access

(taken from LanguageElement.java from MARS, using JDOM)

```
// given: Element element;

protected Set<InputVariableDefinition> getInputVariableDefinitions(
    boolean includeJoinVariables)
{ Set<InputVariableDefinition> definitions = new HashSet<InputVariableDefinition>();

    @SuppressWarnings("unchecked")
    List<Element> elements = element.getChildren();
    for (Element e : elements)
    { String elementName = e.getName();
      if (!elementName.equals("Opaque"))
      { String name = e.getAttributeValue("name", "");
        InputVariableDefinition variable = null;
        if (elementName.equals("has-input-variable"))
            variable = new InputVariableDefinition(name, InputVariableDefinition.INPUT);
        else if (elementName.equals("uses-variable") && includeJoinVariables)
            variable = new InputVariableDefinition(name, InputVariableDefinition.USE);
        if (variable != null)
        { String use = e.getAttributeValue("use", "");
          if (use.length() > 0) variable.setUse(use);
          definitions.add(variable);
        }
      }
    }
    return definitions;
}
```

DOM – sample code fragment: XPath

(taken from ServiceRegistry.java from MARS)

- similar to the JDBC statement concept for SQL in Java:

```
public Element getTaskDescr(Element serviceDescr, String task)
{
    Element taskDescr = null;
    try
    {
        XPath xpath = XPath.newInstance(
            "./lsr:has-task-description/lsr:TaskDescription[" +
            "contains(lsr:describes-task/@rdf:resource,$task)]");
        xpath.addNamespace(Namespaces.RDF_NS);
        xpath.addNamespace(Namespaces.MARS_NS);
        xpath.addNamespace(Namespaces.LSR_NS);
        xpath.setVariable("task", task);
        taskDescr = (Element) xpath.selectSingleNode(serviceDescr);
    }
    catch (Exception e) {...}
}
```

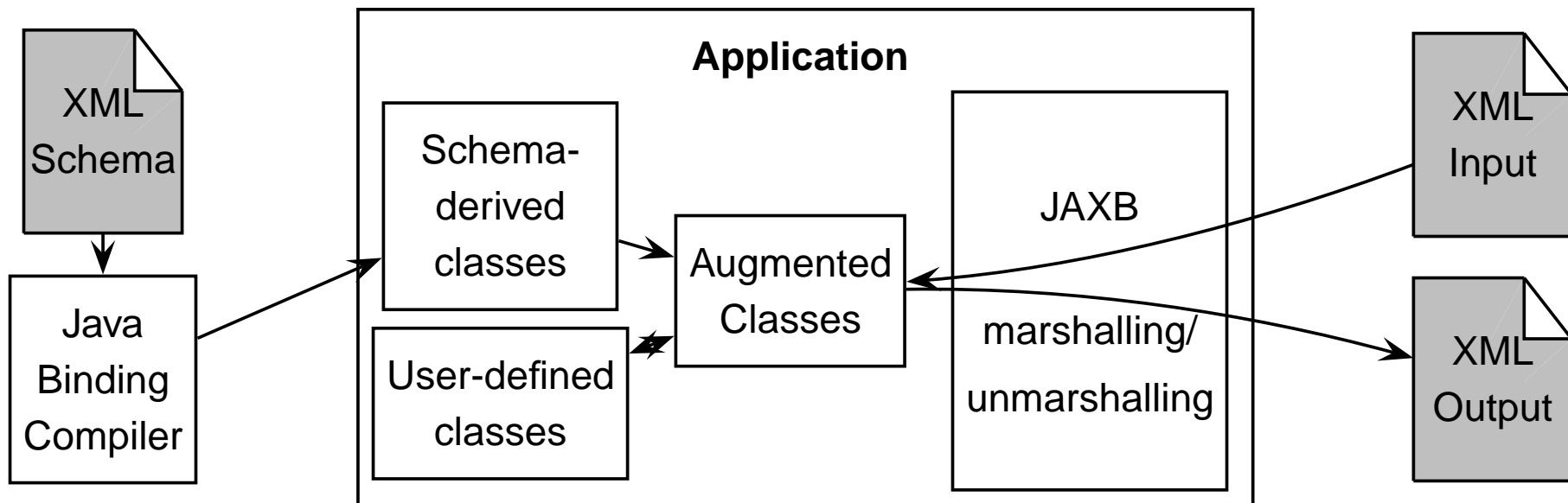
11.2 JAXB - The Java API for XML Binding

- Part of the Java Web Services Developer Pack
- SUN's "Java Web Service Tutorial"
<http://java.sun.com/webservices/tutorial.html>
- XML elements describe objects with properties,
- correspond to classes of an application,
- derive interface with setX/getX methods (= Java Beans) as skeletons for these classes (automatically generated from an XML Schema description),
- user derives classes from these interfaces by adding behavior,
- application logics implemented by using these classes,
- import/export of XML instances of these classes via generic mappings (derived from the XSD).

JAXB ARCHITECTURE

- map XML Schemas to Java classes (get/set methods),
- methods for *unmarshalling* XML instance documents into Java objects,
- methods for *marshalling* Java objects back into XML instance documents.

Architecture



JAXB - EXAMPLE

[Filename: java/JAXB/books.xml]

```
<?xml version="1.0"?>
<Collection>
  <books>
    <book isbn="111-1234">
      <name>Learning JAXB</name>
      <price>34</price>
      <authors>
        <authorName>Jane Doe</authorName>
      </authors>
      <language>English</language>
      <language>French</language>
      <promotion>
        <Discount>10% until March 2003</Discount>
      </promotion>
      <publicationDate>2003-01-01</publicationDate>
    </book>

    <book isbn="112-0815">
      <name>Java Web Services Today and Beyond</name>
      <price>29</price>
      <authors>
        <authorName>John Brown</authorName>
        <authorName>Peter T.</authorName>
      </authors>
      <language>English</language>
      <promotion>
        <Discount>Buy one get Web Services Part 1 free</Discount>
      </promotion>
      <publicationDate>2002-11-01</publicationDate>
    </book>
  </books>
</Collection>
```

- values for `xd:date` and `xs:time` must conform to the syntax required for these XML types (cf. Slide 275)

JAXB - Example: XSD

[Filename: java/JAXB/books.xsd]

```
<?xml version="1.0"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"
  xmlns:jaxb="http://java.sun.com/xml/ns/jaxb" jaxb:version="2.0">

<xs:element name="Collection">
  <xs:complexType>
    <xs:sequence>
      <xs:element name="books">
        <xs:complexType>
          <xs:sequence>
            <xs:element name="book" type="bookType"
              minOccurs="0" maxOccurs="unbounded"/>
          </xs:sequence>
        </xs:complexType>
      </xs:element>
    </xs:sequence>
  </xs:complexType>
</xs:element>

<!-- continue next page -->
```

```

<xs:complexType name="bookType">
  <xs:sequence>
    <xs:element name="name" type="xs:string"/>
    <xs:element name="price" type="xs:long"/>
    <xs:element name="authors" >
      <xs:complexType>
        <xs:sequence>
          <xs:element name="authorName" type="xs:string" minOccurs="1"
            maxOccurs="unbounded"/>
        </xs:sequence>
      </xs:complexType>
    </xs:element>
    <xs:element name="language" type="xs:string" minOccurs="1"
      maxOccurs="unbounded"/>
    <xs:element name="promotion">
      <xs:complexType>
        <xs:choice>
          <xs:element name="Discount" type="xs:string" />
          <xs:element name="None" type="xs:string"/>
        </xs:choice>
      </xs:complexType>
    </xs:element>
    <xs:element name="publicationDate" type="xs:date"/>
  </xs:sequence>
  <xs:attribute name="isbn" type="xs:string" />
</xs:complexType>
</xs:schema>

```

JAXB HowTo

- README file in `java/JAXB/JAXB-README.txt`:

```
## Java Architecture for XML Binding (JAXB)

mkdir myjaxb
cd myjaxb
mkdir classes gen-src
## java6: included in JDK
## earlier java: download jaxb and adjust below HOME:
java -jar JAXB2_20070122.jar
export JAXB_HOME=whereeveritis/jaxb20
export JAXB_LIB=$JAXB_HOME/lib
export JAXB_JAR=$JAXB_LIB/jaxb-api.jar:$JAXB_LIB/jaxb-libs.jar:$JAXB_LIB/jaxb-xjc.jar

$JAXB_HOME/bin/xjc.sh -p JAXBbooks books.xsd -d gen-src
# created classes can the be found in gen-src/JAXBbooks
javac -d classes `find gen-src -name '*.java'`
javac -d classes -classpath classes JAXBbooks.java
java -classpath classes JAXBbooks books.xml

# Syntax for old Java (with jaxb.jar in classpath)
$JAXB_HOME/bin/xjc.sh -p JAXBmondial mondial-jaxb.xsd -d gen-src
javac -d classes -classpath $JAXB_JAR `find gen-src -name '*.java'`
javac -d classes -classpath $JAXB_JAR:classes JAXBmondial.java
java -classpath classes:$JAXB_JAR JAXBmondial mondial-jaxb.xml
```

JAXB: Binding XML Schema to Java Classes

```
<xs:element name="Collection">
```

```
<xs:complexType>
```

```
<xs:sequence>
```

```
<xs:element name="books">
```

minOccurs = maxOccurs = 1 by default

```
<xs:complexType>
```

```
<xs:sequence>
```

```
<xs:element name="book" type="bookType"
```

```
minOccurs="0" maxOccurs="unbounded"/>
```

```
</xs:sequence>
```

```
</xs:complexType>
```

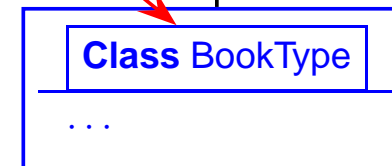
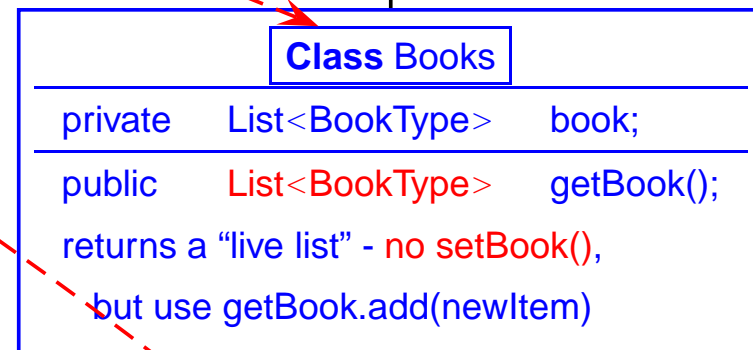
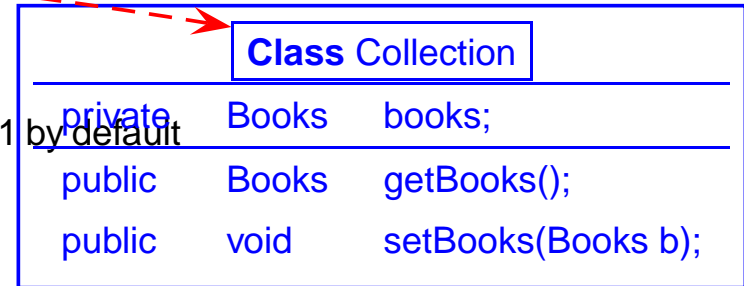
```
</xs:element>
```

```
</xs:sequence>
```

```
</xs:complexType>
```

```
</xs:element>
```

- elements that have complexTypes are mapped to classes (for local type declarations: local classes like Collection.Books),
- elements of simpleTypes and attributes are mapped to instance properties
- multivalued properties are handled by lists; updates not via setXXX(), but via list modifications



JAXB: Binding XML Schema to Java Classes (2)

```

<xs:complexType name="bookType">
  <xs:sequence>
    <xs:element name="name" type="xs:string"/>
    <xs:element name="price" type="xs:string"/>
    <xs:element name="authors" >
      <xs:complexType>
        <xs:sequence>
          <xs:element name="authorName" type="xs:string"
            minOccurs="1" maxOccurs="unbounded"/>
        </xs:sequence>
      </xs:complexType>
    </xs:element>
    <xs:element name="language" type="xs:string"
      minOccurs="1" maxOccurs="unbounded"/>
    <xs:element name="promotion">
      <xs:complexType>
        <xs:choice>
          <xs:element name="Discount" type="xs:string" />
          <xs:element name="None" type="xs:string"/>
        </xs:choice>
      </xs:complexType>
    </xs:element>
    <xs:element name="publicationDate" type="xs:date"/>
  </xs:sequence>
  <xs:attribute name="isbn" type="xs:string"/>
</xs:complexType>
  
```

```

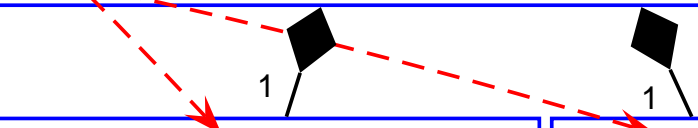
Class BookType
-----
private String    name;
private String    isbn;
private long      price;
private XMLGregorianCalendar publicationDate;
private Authors   authors;
private List      language
private Promotion promotion;
-----
public String    getName(); void setName(String s);
public String    getIsbn(); void setIsbn(String s);
public long      getPrice(); void setPrice(long x);
public XML-G-C. getPublicationDate(); void set(...)
public Authors   getAuthors(); void setAuthors(Authors as);
public List      getLanguage(); no setLanguage()
public Promotion getPromotion(); void setPromotion();
  
```

```

Class Authors
-----
private List<String> authorName;
-----
public List<String> getAuthorName();
live list - no setAuthorName()
  
```

```

Classes Promotion
-----
private String discount;
private String none;
-----
public get/set methods
  
```



JAXB - Example Usage

[Filename: java/JAXB/JAXBbooks.java]

```
import javax.xml.bind.JAXBContext;
import javax.xml.bind.Unmarshaller;
import javax.xml.bind.Marshaller;
import java.util.List;
import javax.xml.datatype.XMLGregorianCalendar;

import java.io.File;
import org.w3c.dom.*;
import javax.xml.transform.dom.*;
import javax.xml.transform.*;
import javax.xml.transform.stream.*;

// import java content classes generated by binding compiler
import JAXBbooks.*;

/**
 * This shows how to use JAXB to unmarshal an xml file
 * Then display the information from the content tree
 */

public class JAXBbooks {

    public static void main (String args[]) {
        try
        {
            JAXBContext jc = JAXBContext.newInstance("JAXBbooks");
            Unmarshaller unmarshaller = jc.createUnmarshaller();

            Collection collection= (Collection)
                unmarshaller.unmarshal(new File( "books.xml"));

            Collection.Books books = collection.getBooks();
            List bookList = books.getBook();
        }
    }
}
```

```

for( int i = 0; i < bookList.size();i++ )
{
    System.out.println("Book details " );
    BookType book =(BookType) bookList.get(i);
    System.out.println("Book Name: " + book.getName().trim());
    System.out.println("Book ISBN: " + book.getIsbn().trim());
    System.out.println("Book Price: " + book.getPrice());
    System.out.println("Book promotion: " +
        book.getPromotion().getDiscount().trim());
    System.out.println("No of Authors " +
        book.getAuthors().getAuthorName().size());

    BookType.Authors authors = book.getAuthors();
    for (int j = 0; j < authors.getAuthorName().size();j++)
    {
        String authorName = (String) authors.getAuthorName().get(j);
        System.out.println("Author Name " + authorName.trim());
    }
    XMLGregorianCalendar date = book.getPublicationDate();
    System.out.println("Date " + date);
    for (int j = 0; j < book.getLanguage().size();j++)
    {
        String language = (String) book.getLanguage().get(j);
        System.out.println("Language " + language.trim());
    }
    // add an element to a live list:
    book.getLanguage().add("Kisuaheli");
    System.out.println();
}

// write the result to an XML file:
Marshaller m = jc.createMarshaller();
DOMResult domResult = new DOMResult();
m.marshal(collection, domResult);
Document doc = (Document) domResult.getNode();
// transformer stuff is only for writing DOM tree to file/stdout
TransformerFactory factory = TransformerFactory.newInstance();
Source docSource = new DOMSource(doc);
StreamResult result = new StreamResult("foo.xml");
Transformer transformer = factory.newTransformer();
transformer.transform(docSource, result);
}catch (Exception e ) {
    e.printStackTrace();
}
}
}

```


JAXB - ANOTHER EXAMPLE

[Filename: java/JAXB/mondial-jaxb.xml]

```
<?xml version="1.0"?>
<mondial>
  <country name="Austria" area="83850" indep_date="1918-11-12" capital="cty-Austria-
    <population>8023244</population>
    <province name="Styria" area="16386">
      <population>1203000</population>
      <city name="Graz">
        <population year="1994-01-01">238000</population>
      </city>
    </province>
    <province name="Salzburg" area="7154">
      <population>501000</population>
      <city name="Salzburg">
        <population year="1994-01-01">144000</population>
      </city>
    </province>
    <province name="Vienna" area="415">
      <population>1583000</population>
      <city name="Vienna" id="cty-Austria-Vienna">
        <population year="1994-01-01">1583000</population>
      </city>
    </province>
  </country>
</mondial>
```

JAXB - Example: XSD

[Filename: java/JAXB/mondial-jaxb.xsd]

```
<?xml version="1.0"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"
  xmlns:jaxb="http://java.sun.com/xml/ns/jaxb" jaxb:version="2.0">
  <xs:element name="mondial">
    <xs:complexType>
      <xs:sequence>
        <xs:element name="country" type="country"
          minOccurs="0" maxOccurs="unbounded"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>

  <xs:complexType name="country">
    <xs:sequence>
      <xs:element name="population" type="populationtype"
        minOccurs="0" maxOccurs="1" />
      <xs:element name="province" type="province"
        minOccurs="0" maxOccurs="unbounded" />
    </xs:sequence>
    <xs:attribute name="name" type="xs:string" use="required"/>
    <xs:attribute name="area" type="xs:integer" use="optional"/>
    <xs:attribute name="car_code" type="xs:ID" use="optional"/>
    <xs:attribute name="indep_date" type="xs:date" use="optional"/>
    <xs:attribute name="capital" type="xs:IDREF" use="optional">
      <xs:annotation> <!-- annotation of the target type <<<<<<< -->
        <xs:appinfo>
          <jaxb:property>
            <jaxb:baseType name="City"/>
          </jaxb:property>
        </xs:appinfo>
      </xs:annotation>
    </xs:attribute>
  </xs:complexType>
```

```

<xs:complexType name="province">
  <xs:sequence>
    <xs:element name="population" type="populationtype"
      minOccurs="0" maxOccurs="1" />
    <xs:element name="city" type="city"
      minOccurs="0" maxOccurs="unbounded" />
  </xs:sequence>
  <xs:attribute name="name" type="xs:string" use="required"/>
  <xs:attribute name="area" type="xs:integer" use="optional"/>
</xs:complexType>

<xs:complexType name="city">
  <xs:sequence>
    <xs:element name="population" type="populationtype"
      minOccurs="0" maxOccurs="unbounded"/>
  </xs:sequence>
  <xs:attribute name="name" type="xs:string" use="required"/>
  <xs:attribute name="id" type="xs:ID"/>
</xs:complexType>

<xs:complexType name="populationtype">
  <xs:simpleContent>
    <xs:extension base="xs:decimal">
      <xs:attribute name="year" type="xs:date" use="optional"/>
    </xs:extension>
  </xs:simpleContent>
</xs:complexType>
</xs:schema>

```

- annotation of the country/@capital IDREFS attribute:
⇒ public City getCapital()
- countries have at most one population subelement, cities may have several ones.

JAXB - Example Usage

[Filename: java/JAXB/JAXBmondial.java]

```
import javax.xml.bind.JAXBContext;
import javax.xml.bind.Unmarshaller;
import java.io.File;
import java.util.List;

// import java content classes generated by binding compiler
import JAXBmondial.*;

/**
 * This shows how to use JAXB to unmarshal an xml file
 * Then display the information from the content tree
 */

public class JAXBmondial {

    public static void main (String args[]) {
        try {
            JAXBContext jc = JAXBContext.newInstance("JAXBmondial");
            Unmarshaller unmarshaller = jc.createUnmarshaller();

            Mondial mondial =
                (Mondial) unmarshaller.unmarshal(new File("mondial-jaxb.xml"));
            List countryList = mondial.getCountry();
            Province prov;
            City city;

            for ( int i = 0; i < countryList.size();i++ )
            {
                Country country = (Country) countryList.get(i);
                System.out.println("Country: " + country.getName() );
                System.out.println("    pop: " +
                    country.getPopulation().getValue());

                // Java knows from the annotation of the IDREF attribute
                // that this is a city

                City c = country.getCapital();
                System.out.println("    cap: " + c.getName());

                List provList = country.getProvince();
                for (int j = 0; j < provList.size() ; j++)
                {
                    464
                    prov = (Province) provList.get(j);
                    System.out.println("        Province name: " + prov.getName() );
                }
            }
        }
    }
}
```

JAXB MAPPING: SUMMARY

- allows for easy and lightweight unmarshalling, bean-based manipulation and marshalling of XML data,
- higher level of abstraction from XML representation, compared with DOM and SAX,
- but **still actually just a way to manipulate XML data without having to know the specific notions of the XML data model.**

Minor Comments

- naming (getBook() for a list etc.) not always intuitive; can be customized by annotations to the XSD;
 - intermediate elements (example: Books, Authors) lead to unnecessary classes; can often be omitted (example: Language elements)
- ⇒ to get a better “modeling”, do not use structures like
Country-hasProvince-Province-hasCity-City
(as in Striped RDF/XML [Semantic Web lecture]; this would generate intermediate classes), but
Country-Province-City.

JAXB INTEGRATION WITH JAVA APPLICATION?

A comfortable usage of the generated classes into an application program is not yet supported:

- means: add application-specific methods
(and properties that would be local to the Java existence of the object)
- define a subclass: `java_xxx` extends `xxx`
 - after unmarshalling, the objects are only instances of `xxx`

⇒ methods of `java_xxx` not applicable
- define class `java_xxx` where `xxx` is a subclass of:
 - useful from the java point of view: extend application class with bean functionality and marshalling
 - cannot be communicated to the JAXB generation of the classes (annotation with `xjc:superClass c` in the xsd does only allow to make all classes subclasses of `c`)

JAXB INTEGRATION WITH JAVA APPLICATION

Delegation

- (manually) write application classes that delegate to the JAXB-generated classes and extend them with application functionality,
- after unmarshalling, traverse the tree and create the “real” objects

⇒ application classes must be manually adapted after schema changes.

Manual editing of generated classes themselves

- edit the generated xxx.java files
- if instance attributes are added, they must also be added either to propOrder, or get an annotation as @XmlAttribute – and then they will be exported when marshalling them

⇒ must be manually redone/adapted after schema changes.

Using Helper Classes

- encode behavior in separate helper classes that provide static methods:
mondialHelper.addProvince(Country,Province)

JAXB - Example Usage with extended class definition

put the following into the generated JAXB/gen-src/JAXBmondial/Country.java and then recompile:

```
// a method for more comfortable manipulation:
public void addProvince(Province p) {
    getProvince().add(p);
}
// a "useful" method:
public void printCityNames() {
    List provList = getProvince();
    for (int j = 0; j < provList.size() ; j++)
    {
        Province prov = (Province) provList.get(j);
        List cityList = prov.getCity();
        for (int k = 0; k < cityList.size() ; k++)
        {
            City city = (City) cityList.get(k);
            System.out.println(city.getName().trim());
        }
    }
}
```

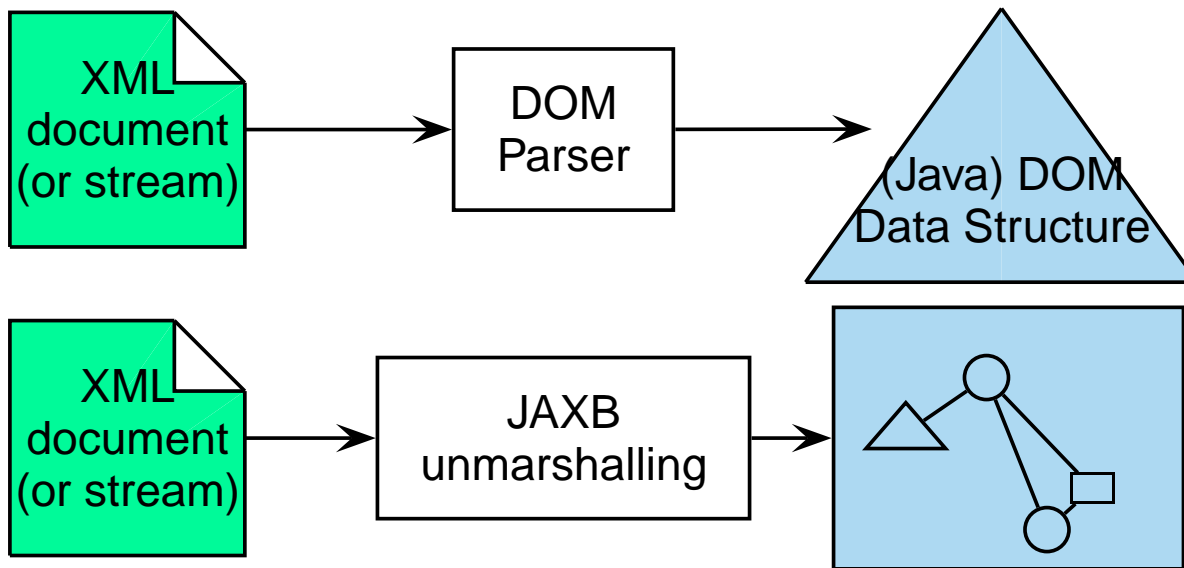
[Filename: java/JAXB/put-into-Country.java]

ASIDE: SOAP (SIMPLE OBJECT ACCESS PROTOCOL)

- Generic “protocol” (nevertheless, HTTP-based)
 - Any object can be serialized in XML, sent, and deserialized (only having the Java class code, without having an XSD).
So far similar to the OIF (Object Interchange Format) of ODMG (cf. Slide 53 ff.).
 - The XML representation is not intended to be processed on the XML level, but only by soap-unpacking it.
 - Bad experience: correct packing/unpacking only between same SOAP implementations.
 - Note: Instances of Java XML (DOM) are not serialized as plain XML, but as SOAP serialization of an instance of the underlying DOM implementation class.
(not intended *for* exchanging XML, but for exchanging objects *by* XML).
- ⇒ when messages are designed to be XML, SOAP is not the right way, but use simple, plain HTTP!
- One does not need to have any knowledge of XML to use soap (actually, knowledge of XML doesn't help).
- ⇒ so it does not fit in this course.

11.3 XML Stream Processing

- reading from a file or from an HTTP connection both is actually reading char by char from a stream
- the stream is parsed, resulting in something that can be used by application:



PARSING: GENERAL CONCEPTS

- **Compiler Construction in general:**

1. input: a sequence of characters
2. lexical analysis ("lexer") extracts the keywords (e.g. "begin", "end", "for") and outputs a sequence of *tokens*
3. syntactical (grammar) analysis: check grammatical structure and generate the *parse tree* (e.g. via automaton)
4. tools: lex & yacc/bison
5. interpreter, optimizer, compiler, visualizer etc. process the parse tree

- **XML:**

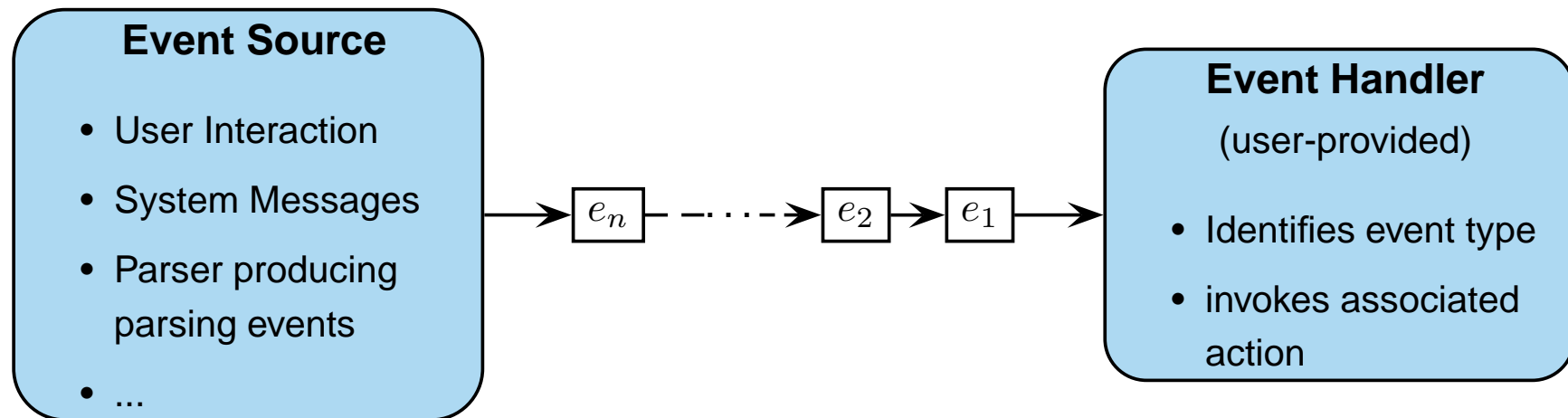
1. lexical: split unicode input sequence into opening tags, closing tags, attributes, PCDATA, processing instructions, etc.
2. syntactical and structural: is it well-formed?
3. processing: build DOM tree, build JAXB structure, visualize, . . .

- the above DOM and JAXB are actually parser+specific processing

⇒ **XML Stream Processing: works on the tokens sequence!**

EVENT-BASED PROCESSING AS A *general Design Pattern*

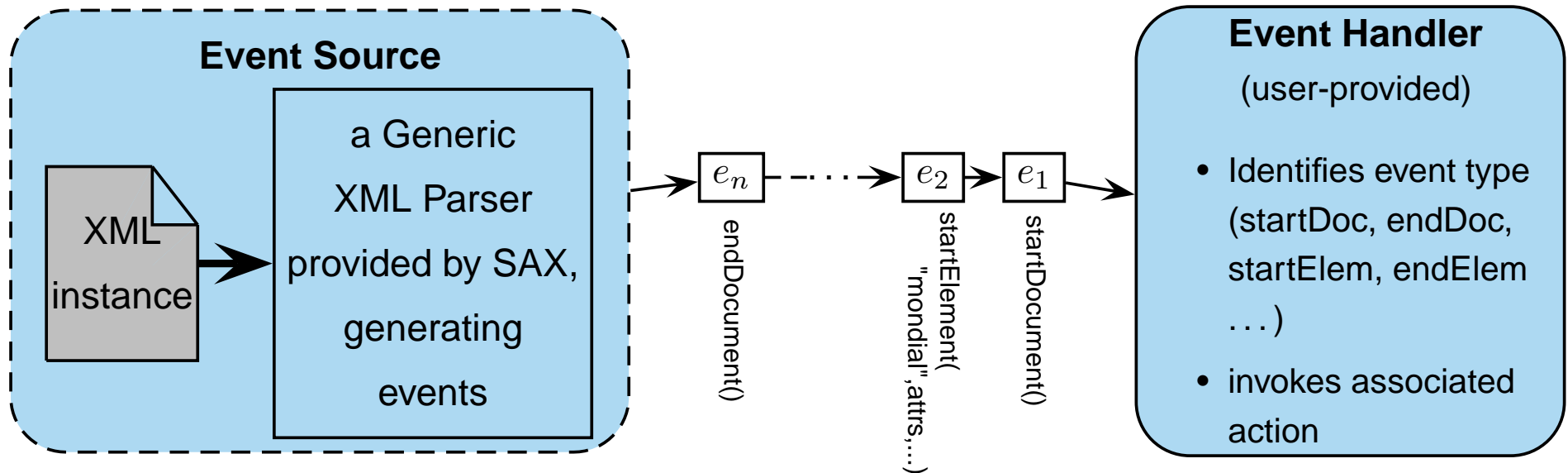
- A stream of (high-level) items that carry some inherent semantics can be seen as a stream of “events”
(in contrast to a simple 0-1-stream, a byte stream or similar low-level streams)



- The application programmer provides the Event Handler implementation, containing actions for each type of event;
- kind of *rule-based*;
- programmer is *not* in charge of the control flow

11.4 Event-based XML Parsing with SAX

- SAX (“The Simple API for XML”) is an *event-based interface/model*



Represents/processes an XML document as a sequence of events (depth-first traversal), e.g.

- `startDocument()`, `endDocument()`
- `startElement(Name, attributesList)` – attributes not split
- `endElement(Name)`
- `characters(string)`

XML PARSING WITH SAX

SAX: parse XML from a file (in general: char stream).

- import classes: `javax.xml.parsers.*`, `org.xml.sax.*`
- a generic XML Parser is parameterized with a *Content Handler* (plus *Error Handler*, *DTD Handler*, and *EntityResolver*) implementation.
- The most trivial Content Handler is the *DefaultHandler* that does nothing: the document is parsed, events are detected, but no action is performed (DTD / XML Schema validation can be switched on).
- Event handler programmed wrt. a “*push API*”.
- Normally, the user-provided Content Handler extends the *DefaultHandler*, overwriting (some of) its Event Methods.
- With the content handler implementation, the user provides “actions” in form of Java code, associated with specific events (and even dependent on context information).
- If during parsing of the XML document, a specific event occurs, the code of the associated action from the content handler is invoked (“*callback*”).

SAX: APPLICATIONS

Only events are signaled: linear processing based on incoming sequence of events.

- ... among many other things, one can generate a DOM tree structure,
- validation according to a DTD (using the automaton as given on Slide 176) in linear time,
- stream-processing of XML input
 - start processing already when input document is not yet complete
 - filtering for elements that are relevant for a given application
 - linear search for something, e.g., names of countries
(Exercise: sketch the behavior of the event handler on relevant events)
 - if the stream is a list of elements of the same structure:
generate a database entry for each element (use JDBC)
- if necessary: application needs to maintain context.

SAX EXAMPLE CODE

Consider a very simple application that

- detects all elements with attributes
 - for each element, output the element's name
 - for each element, output the name-value pairs of its attributes
-

```
>java PrintAttributes mondial.xml > bla.out
>less bla.out
element: country
- attribute: 'car_code' value: 'AL' type: 'ID'
- attribute: 'area' value: '28750' type: 'CDATA'
- attribute: 'capital' value: 'cty-cid-cia-Albania-Tirane' type: 'IDREF'
- attribute: 'memberships' value: 'org-BSEC org-CE org-CCC org-ECE org-EBRD org-EU ...' type: 'IDREFS'
element: encompassed
- attribute: 'continent' value: 'europe' type: 'IDREF'
- attribute: 'percentage' value: '100' type: 'CDATA'
element: ethnicgroups
- attribute: 'percentage' value: '3' type: 'CDATA'
element: ethnicgroups
- attribute: 'percentage' value: '95' type: 'CDATA'
element: religions
- attribute: 'percentage' value: '70' type: 'CDATA'
...
```

Class "PrintAttributes.java":

```
import java.io.IOException;
import javax.xml.parsers.*;
import org.xml.sax.SAXException;
import org.xml.sax.helpers.DefaultHandler;

public class PrintAttributes {
    public static void main(String[] args) {
        if (args.length != 1) {
            System.err.println("usage: PrintAttributes <url>");
            System.exit(1);
        }
        String url = args[0];    // ... prepare a contentHandler:
        DefaultHandler handler = new ContentHandlerPrintAttributes(
            "printing attributes of document at url '" + url + "'");
        SAXParserFactory factory = SAXParserFactory.newInstance();
        try {
            SAXParser parser = factory.newSAXParser();
            parser.parse(url, handler); // <<<<<<<< and now it runs ...
        } catch (IOException e1) {
            e1.printStackTrace();
        } catch (ParserConfigurationException e) {
            e.printStackTrace();
        } catch (SAXException e) {
            e.printStackTrace();
        }
    }
}
```

[see java/SAX/PrintAttributes.java]

Class "ContentHandlerPrintAttributes.java":

```
import org.xml.sax.Attributes;
import org.xml.sax.SAXException;
import org.xml.sax.helpers.DefaultHandler;

public class ContentHandlerPrintAttributes extends DefaultHandler {
    public ContentHandlerPrintAttributes(String message) {
        System.out.println(message);
    }

    // react on opening elements:
    public void startElement(String url, String localName, String qName,
        Attributes attrs) throws SAXException {
        if (attrs.getLength() > 0) {
            String elementName;
            if(qName == null || qName.equals("")) elementName = localName;
            else elementName = qName;
            System.out.println("element: " + elementName);
            for (int i = 0; i < attrs.getLength(); i++) {
                System.out.println(" - attribute: '" + attrs.getQName(i)
                    + "' value: '" + attrs.getValue(i) + "' type: '"
                    + attrs.getType(i)+"'");
            }
            System.out.println();
        }
    }
    // methods for endElement(), startDocument(), endDocument(), characters() omitted
    // all other "parsing events" are ignored in this case
}
```

[see java/SAX/ContentHandlerPrintAttributes.java]

SAX: APPLICATIONS TO XPATH QUERY ANSWERING

Forward queries

XPath-queries like `//country[@car_code='D']/population` can be answered very (time- and memory-)efficient,

- use the sequence of events (linear)
- maintain some context (often LOGSPACE/additional LOGTIME sufficient)

... works only for queries, that contain only forward steps,

General queries

which XPath expressions can be *transformed* in equivalent forward-expressions (and with what efforts)?

- “XPath: Looking forward”; F. Bry et al ; 2002; LMU München

- [theory: complexity, connections to linear temporal logic](#)

For every linear temporal logic formula that uses past and future operators, there is an equivalent formula that uses only future operators

... but in general of exponential size.

11.5 XML Streams/StAX - The Streaming API for XML

Higher abstraction level (than character-based XML) for XML data exchange:

javax.xml.stream (rt.jar)

Reconsider SAX

- on-the fly processing, no in-memory representation for good performance
- idea of “XML Event Stream”: a char stream (File, HTTP) can be converted into an XML Event Stream by an XML parser; see example’s main() method.
- SAX does not make the XML Event Stream accessible, but only via the Event Handler.

Generalization and Abstraction: XML Streams

- XMLEvents: StartDocument, StartElement, Character, EndElement,
- XMLStreamWriter, XMLStreamReader,
- XML Streams also can be connected *directly* as an *abstract* means to exchange XML

SAX AND STAX: APPLICATIONS

Stream-based processing can be applied to XML data on multiple levels:

- low-level applications:
SAX is often used for building a DOM from ASCII XML input: “opening tag with attributes”, “text”, “closing tag” can immediately be translated into the DOM constructors.
- low-level streaming of an XML instance:
answering XPath (forward-axes only) queries; optionally maintaining some context (e.g., stack).
- higher level “application-level events”:
the XML stream is not seen as the traversal of a large instance, but as a sequence of (independent) XML fragments that are seen as application-level events
[RFID applications, time series of stock quotes, RSS feeds]

XML Streams: Application Scenarios

- READ: usage analogous to SAX: process an XML file input as an XML Event Input Stream:
control flow is not passed to the parser (**unlike** SAX), but XML events are accessed using an *iterator*, controlled by the Java program using the StAX API (*Pull-API*).
[Note: iterators are a common design pattern, not only applied to collections, but as we see here also to streams: `init()`, `next()`, ...]
⇒ application code: same as for SAX, only operational embedding done differently.
- WRITE AND READ: streamed data exchange between processed on the XML level

Interfaces XMLStreamWriter, XMLStreamReader

(only some comments; see also following examples)

Reader

- `int event = r.next()` and then switch based on event type
javax.xml.stream.XMLStreamConstants.XX:
START_DOCUMENT, START_ELEMENT, ATTRIBUTE, CHARACTERS, END_ELEMENT, ...
- goal-driven access methods when on START_ELEMENT:
`r.getLocalName()`, `r.getAttributeValue(name)`,
`r.getAttributeCount()`, `getAttributeValue(n)` for iteration,
`r.getElementText()` (reads also the next EndElement from the stream!)
- goal-driven access method when on CHARACTERS: `r.getText()`

Writer

- Writer: `w.writeStartDocument()`, `w.writeStartElement(name), a`
`w.writeAttribute(name, value)`, `w.writeCharacters(text)`: obvious;
- `w.writeEndElement()`: closes the innermost open element;
- `w.writeEndDocument()`: closes all open elements.
- `w.flush()`: force write any data to the underlying output mechanism.

StAX EXAMPLE: EXAM REGISTRATION

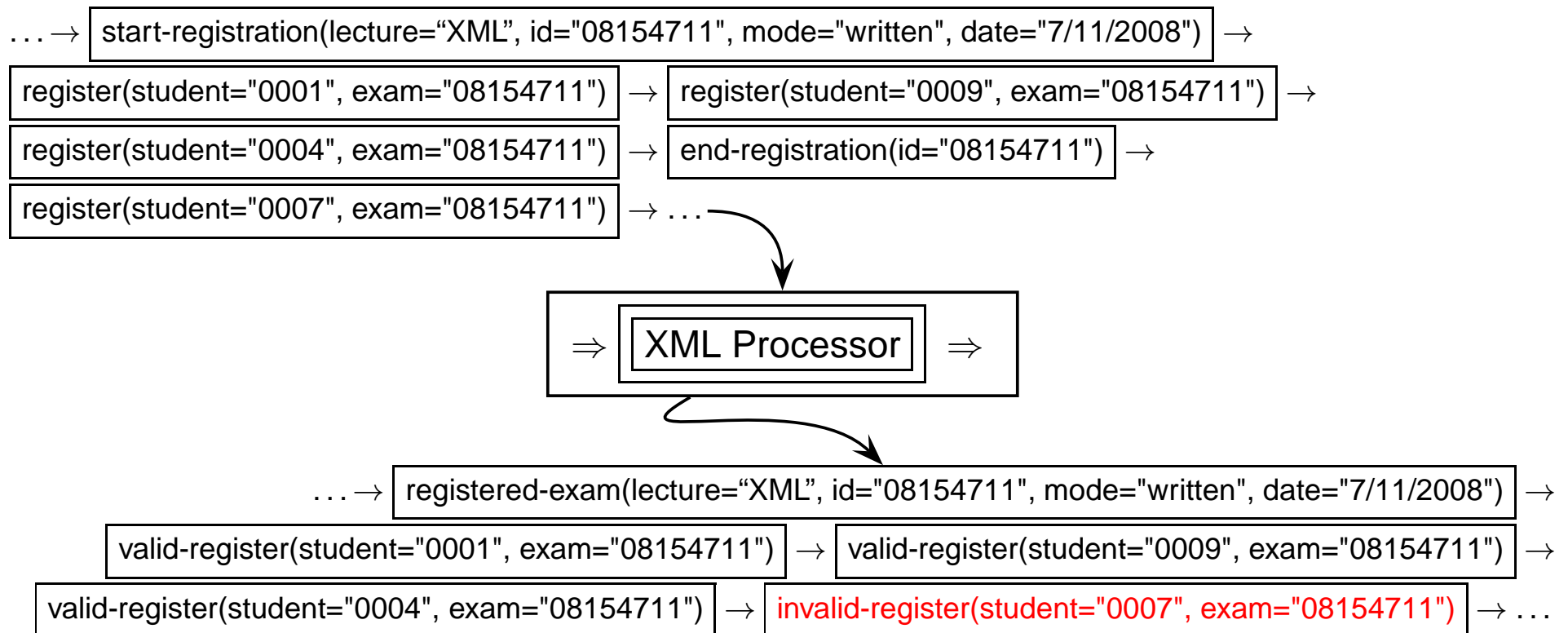
Assume the administration of exams in a student's office ("Prüfungsamt"):

- The *subject* (e.g., "Semi-structured Data and XML") and ID of lectures/exams,
- whether the exam is *written* or *oral*,
- for written exams, the date of the exam,
- for oral exams, a number of dates is given when the single exams are held.
- the registration period *starts* when receiving an incoming XML message
`start-registration`
- the registration period *ends* when receiving an incoming XML message
`end-registration`
- for all students that did (`register`) correctly, the student's relevant details are extracted and written to some output stream (`valid-register`; in the example, we use `STDOUT`.)
- students that register before beginning or after the end of registration, are not accounted for the exam; an XML message `invalid-register` goes to `STDOUT`,

StAX Example: Exam Registration (Cont'd)

- the registration data of the students comes in via a continuous input stream;
- the program should allow the management of registrations for multiple exams at one time (all incoming over the same input stream).

Example stream:



StAX EXAMPLE CONT'D:

Consider the following XML sequence as input stream:

```
<?xml version="1.0" encoding="UTF-8"?>
<stream>
  <register student="0007" exam="08154711"/>
  <start-registration id="08154711" mode="written">
    <subject>Semistructured Data and XML</subject>
    <date>07/11/2008</date>
  </start-registration>
  <register student="0001" exam="08154711"/>
  <register student="0009" exam="08154711"/>
  <start-registration id="12345678" mode="oral">
    <subject>Dental Hygiene</subject>
    <dates>
      <date>17/9/2008</date>
      <date>18/9/2008</date>
    </dates>
  </start-registration>
  <end-registration id="12345678"/>
  <register student="0004" exam="08154711"/>
  <register student="0004" exam="12345678"/>
  <register student="0007" exam="12345678"/>
  <end-registration id="08154711"/>
  <register student="0007" exam="08154711"/>
</stream>
```

[Filename: java/StAX/exam.xml]

StAX EXAMPLE CONT'D (2):

Code for the Exam bean, containing the exam's properties and some constants):

```
import java.util.ArrayList;
import java.util.List;

public class Exam {
    public static final String DATE = "date";
    public static final String SUBJECT = "subject";
    public static final String ID = "id";
    public static final String MODE = "mode";
    public static final String DATES = "dates";
    public static final String STARTOFREG = "start-of-registering";
    public static final String ENDOFREG = "end-of-registering";

    private String id;
    private boolean oral;
    private String subject;
    private String date;
    private List<String> dates;
    private boolean registeringClosed = false;
    private String startOfReg;
    private String endOfReg;

    public Exam(String id, String mode) {
        this.id = id;
        this.oral = "oral".equals(mode);
        this.dates = new ArrayList();
    }
}
```

```

public String getId() { return id; }
public void setDate(String date) { this.date = date; }
public String getDate() { return date; }
public void setDates(List<String> dates) {this.dates = dates; }
public List<String> getDates() { return dates; }
public void setSubject(String subject) { this.subject = subject; }
public String getSubject() { return subject; }
public boolean isOral() { return oral; }
public boolean isWritten() { return (!oral); }

public String getMode() {
    if (oral) return "oral";
    return "written";
}
public boolean isRegisteringClosed() {
    return registeringClosed;
}
public void setRegisteringClosed(boolean registeringClosed) {
    this.registeringClosed = registeringClosed;
}
public String getEndOfReg() {
    return endOfReg;
}
public String getStartOfReg() {
    return startOfReg;
}
public void setStartOfReg(String startOfReg) {
    this.startOfReg = startOfReg;
}
public void setEndOfReg(String endOfReg) {
    this.endOfReg = endOfReg;
}
}

```

[Filename: java/StAX/Exam.java]

StAX EXAMPLE CONT'D (3):

Code for the main parser class, containing the `main` method:

```
import java.io.FileInputStream;
import java.io.FileNotFoundException;
import java.io.OutputStream;
import java.text.DateFormat;
import java.text.SimpleDateFormat;
import java.util.Date;
import java.util.HashMap;
import java.util.Iterator;
import java.util.Map;

import javax.xml.stream.XMLInputFactory;
import javax.xml.stream.XMLOutputFactory;
import javax.xml.stream.XMLStreamConstants;
import javax.xml.stream.XMLStreamException;
import javax.xml.stream.XMLStreamReader;
import javax.xml.stream.XMLStreamWriter;

public class ExamStreamParser {

    FileInputStream inputStream;
    OutputStream outputStream;

    public ExamStreamParser(FileInputStream in, OutputStream out) {
        this.inputStream = in;
        this.outputStream = out;
    }

    public void startParsing() {
        try {
            XMLInputFactory inputFactory = XMLInputFactory.newInstance();
            XMLOutputFactory outputFactory = XMLOutputFactory.newInstance();
            XMLStreamReader parser = inputFactory.createXMLStreamReader(inputStream);
            XMLStreamWriter writer = outputFactory.createXMLStreamWriter(outputStream);
            Exam currentExam = null;
            Map<String,Exam> exams = new HashMap<String,Exam>();
            boolean goOn = true;
```

```

while (goOn) {
int event = parser.next();
switch(event) {
case XMLStreamConstants.END_DOCUMENT:
    parser.close();
    writer.flush();
    writer.close();
    goOn = false;
    break;
case XMLStreamConstants.START_ELEMENT:
    // start-registration and its subelements
    if("start-registration".equals(parser.getLocalName())) {
        currentExam = new Exam(parser.getAttributeValue(null, Exam.ID), parser.getAttributeValue(null, Exam.MOD), Exam.ID);
        currentExam.setStartOfReg(getDate());
        break;
    }
    if(Exam.SUBJECT.equals(parser.getLocalName())) {
        currentExam.setSubject(parser.getElementText()); break;
    }
    if(Exam.DATE.equals(parser.getLocalName())) {
        if(currentExam.isWritten()) currentExam.setDate(parser.getElementText());
        else currentExam.getDates().add(parser.getElementText());
        break;
    }
    if("end-registration".equals(parser.getLocalName())) {
        String examId = parser.getAttributeValue(null, Exam.ID);
        Exam exam = exams.get(examId);
        if(exam == null) {
            System.err.println("no such exam with id '"+examId+"' open for registration!");
            break;
        }
        exam.setEndOfReg(getDate());
        exam.setRegisteringClosed(true);
        break; // no output is generated.
    }
    // register and its subelements
    if("register".equals(parser.getLocalName())) {
        String studentId = parser.getAttributeValue(null, "student");
        String examId = parser.getAttributeValue(null, "exam");
        if(exams.containsKey(examId)) {
            Exam exam = exams.get(examId);
            if(! exam.isRegisteringClosed()) {
                writer.writeStartElement("valid-register");
                writer.writeAttribute("student", studentId);
                writer.writeAttribute("exam", examId);
                writer.writeEndElement();
            } else {
                writer.writeStartElement("invalid-register");
                writer.writeAttribute("student", studentId);
                writer.writeAttribute("exam", examId);
                writer.writeStartElement("message");
                writer.writeCharacters("invalid registration! registration for exam '" + exam.getId()
                    + "' (" + exam.getSubject() + ") has ended on " + exam.getEndOfReg());
                writer.writeEndElement();
            }
        } else {
            writer.writeStartElement("invalid-register");
            writer.writeAttribute("student", studentId);
            writer.writeAttribute("exam", examId);
            writer.writeStartElement("message");
            writer.writeCharacters("invalid registration! exam '"+examId+"' is not (yet?) open for registration.");
            writer.writeEndElement(); writer.writeEndElement();
        }
        writer.writeCharacters("\n");
        break;
    }
}
break;

```

```

case XMLStreamConstants.END_ELEMENT:
    if("start-registration".equals(parser.getLocalName())) {
        exams.put(currentExam.getId(),currentExam);
        writer.writeStartElement("registered-exam");
        writer.writeAttribute(Exam.ID, currentExam.getId());
        writer.writeAttribute(Exam.MODE, currentExam.getMode());
        writer.writeCharacters("\n ");
        writer.writeStartElement(Exam.SUBJECT);
        writer.writeCharacters(currentExam.getSubject());
        writer.writeEndElement();
        writer.writeCharacters("\n ");
        writer.writeStartElement(Exam.STARTOFREG);
        writer.writeCharacters(currentExam.getStartOfReg());
        writer.writeEndElement();
        writer.writeCharacters("\n ");
        if(currentExam.isWritten()) {
            writer.writeStartElement(Exam.DATE);
            writer.writeCharacters(currentExam.getDate());
            writer.writeEndElement();
        } else {
            writer.writeStartElement(Exam.DATES);
            for(Iterator<String> i=currentExam.getDates().iterator();i.hasNext();) {
                writer.writeStartElement(Exam.DATE);
                writer.writeCharacters(i.next());
                writer.writeEndElement();
            }
            writer.writeEndElement();
        }
        writer.writeCharacters("\n");
        writer.writeEndElement();writer.writeCharacters("\n");
        currentExam = null; // it's better to provoke
        // a nullpointer exception than to edit the wrong exam object
        break;
    }
}
} catch (XMLStreamException e) {
    e.printStackTrace();
}
}

private String getDate() {
    DateFormat format = new SimpleDateFormat();
    Date date = new Date();
    return format.format(date);
}

public static void main(String[] args) {
    try {
        ExamStreamParser examStreamParser = new ExamStreamParser(new FileInputStream(args[0]), System.out);
        examStreamParser.startParsing();
    } catch (FileNotFoundException e) {
        e.printStackTrace();
    }
}
}
}

```

[Filename: java/StAX/ExamStreamParser.java]

StAX COMPARISON WITH SAX

SAX: • “Push” API

- Common pattern: methods for each event type, where `startElement()` and `endElement()` contain large `ifs`.

StAX: • “Pull” API

- Common pattern: huge `switch` command whose cases again contain large `if`.
- Performance: no difference
(inside, StAX's `next()` does the same as the SAX builtin stream reader algorithm that call the `EventHandler`)
- The actual code to be written is not much different in both cases.
- SAX maps a unicode input stream directly to the `EventHandler` calls.
- StAX makes the [intermediate abstraction level](#) of XML event streams accessible
 - can easily produce XML output via `XMLStreamWriter` (e.g. to another StAX appl.)
(programmatically, the same feature can be provided by a SAX Event Handler as well using StAX's `xml.stream` classes)

Example: XML Stream Communication

```
import java.io.PipedInputStream;
import java.io.PipedOutputStream;
import java.io.OutputStream;
import javax.xml.stream.XMLOutputFactory;
import javax.xml.stream.XMLStreamWriter;

public class XMLStreamTestWriter implements Runnable
{
    OutputStream outputStream;

    public XMLStreamTestWriter(OutputStream out) {
        this.outputStream = out;
    }

    public void run() {
        try {
            XMLOutputFactory outputFactory = XMLOutputFactory.newInstance();
            XMLStreamWriter writer = outputFactory.createXMLStreamWriter(outputStream);
            writer.writeStartElement("foo");
            int i=1;
            while (i<100) {
                writer.writeStartElement("bla");
                writer.writeCharacters(" " + i);
                writer.writeEndElement();
                System.out.print("Write <bla>" + i + "</bla> ");
                //writer.flush(); // if not uncommented: strictly alternating
                // comment out flush: sleep < 700 causes alternating after blocks of 2...5 elements
                try{ java.lang.Thread.sleep(50); }
                    catch (Exception e) { e.printStackTrace(); }
                i++;
            }
            // writer.writeEndElement(); // close </foo> is done by the next line:
            writer.writeEndDocument(); // docu: closes all tags, but does not send anything else
            writer.flush();
            writer.close();
        } catch (Exception e) { e.printStackTrace(); }
        System.out.println("Writer finished");
    }

    public static void main(String[] args) throws Exception{
        PipedOutputStream pos = new PipedOutputStream();
        PipedInputStream pis = new PipedInputStream();
        pis.connect(pos);
        new Thread (new XMLStreamTestWriter(pos)).start();
        new Thread (new XMLStreamTestReader(pis)).start();
    }
}
```

[Filename: java/StAX/XMLStreamTestWriter.java]

- underlying: connected PipedOutput/InputStream

Example: XML Stream Communication (Cont'd)

```
import java.io.PipedInputStream;
import java.io.InputStream;
import javax.xml.stream.XMLInputFactory;
import javax.xml.stream.XMLStreamConstants;
import javax.xml.stream.XMLStreamReader;

public class XMLStreamTestReader implements Runnable {

    InputStream inputStream;

    public XMLStreamTestReader(PipedInputStream in) {
        this.inputStream = in;
    }

    public void run() {
        try {
            XMLInputFactory inputFactory = XMLInputFactory.newInstance();
            XMLStreamReader parser = inputFactory.createXMLStreamReader(inputStream);
            boolean goOn = true;
            while (goOn) {
                int event = 0;
                try {
                    event = parser.next();
                    switch(event) {
                        case XMLStreamConstants.START_ELEMENT:
                            System.out.println("Read start element " + parser.getLocalName());
                            break;
                        case XMLStreamConstants.CHARACTERS:
                            System.out.println("Read " + parser.getText());
                            break;
                        case XMLStreamConstants.END_ELEMENT:
                            System.out.println("Read end element " + parser.getLocalName());
                            break;
                        case XMLStreamConstants.END_DOCUMENT: // never happens!
                            System.out.println("Read end document");
                            goOn = false;
                        default: System.out.println("Read something else. event: " + event);
                    }
                } catch (Exception e) { parser.close(); goOn = false; }
            }
            parser.close();
            System.out.println("Reader finished");
        } catch (Exception e) { e.printStackTrace(); }
    }
}
```

[Filename: java/StAX/XMLStreamTestReader.java]

11.6 Aside: Use of Date and Time Datatypes

- XML Schema: simple datatypes for dateTime
- represented by String literals in attribute values or text contents
 - xs:dateTime: *yyyy-mm-ddThh:mm:ss[.xx][{+|-}hh:mm]*
 - xs:time: *hh:mm:ss[{+|-}hh:mm]*
 - xs:duration: *P[nY][nM][nD][T[nH][nM][n[.n]S]]*, where *n* can be any natural number
 - xs:dayTimeDuration, xs:yearMonthDuration: restrictions of xs:duration.
- for XQuery handling with specific operations (similar to those known from SQL) cf. Slide 275.
- map to appropriate classes for processing by Java (and by this e.g. with JDBC from and to SQL databases).

ASIDE: DATE AND TIME IN JAVA

Java provides several classes for handling date and time:

- Datatype: `import java.util.GregorianCalendar;`

Create from string representation:

- `import java.text.DateFormat; import java.text.SimpleDateFormat;`

```
public static String datedefaultpattern = "yyyy-MM-dd'T'HH:mm:ss";  
// input: string s (following the XML Schema pattern)  
DateFormat df = new SimpleDateFormat(datedefaultpattern);  
GregorianCalendar typedvalue = df.parse(s);  
// result: typedvalue as an object
```

- see `java.util.GregorianCalendar` for method documentation.
- TODO: check if JAXB declares automatically as date (which class in Java) and whether set-methods of JAXB-created classes automatically convert data.

11.7 Web Services (Overview)

- History: RPC (Remote Procedure Call)
 - call a specific procedure at a specific server
(client stub→marshalling→message→unmarshalling→ server stub→ server).
- History: OMG Standard (Object Management Group) CORBA (1989, “Common Object Request Broker Architecture”; cf. Slides 38 ff.):
 - Middleware, usually applied in an Intranet,
 - central ORB bus where services can connect,
 - service registry (predecessor of WSDL and UDDI ideas)
 - description of service interfaces in object-oriented style
(IDL - interface description language, similar to C++ declarations)
 - exchanging objects between services via OIF (Object Interchange Format)

⇒ RPC abstraction (call abstract functionality) by the ORB as a broker.
- XML-RPC and SOAP+WSDL+UDDI) are XML-based variants of RPC+Corba.
- SOA (“Service-Oriented Architecture”).

HTTP: HYPERTEXT TRANSFER PROTOCOL (OVERVIEW)

- HTTP 0.9 (1991), HTTP 1.0 (1995), HTTP 1.1 (1996).
- Application Layer Protocol, based on a (reliable) transport protocol (usually TCP “Transmission Control Protocol” that belongs to the “Internet Protocol Suite” (IP)) [see Telematics lecture].
- Request-Response Protocol: open connection, send something, receive response (upon completion), close connection
- usually associated with Web Browsing and HTML:
send (HTTP GET) URL, get URL (=resource) contents
⇒ this is already a (very basic) Web Service
also: send HTTP POST URL+Data (Web Forms) get answer
⇒ this is also a (still basic) Web Service; “Hidden Web”
- common protocol used for communication with and between Web Services ...

(JAVA) SERVLET

- a piece of software that should be made available as a Web Service
- implements the methods of the Servlet interface
(Java: `javax.Servlet`, subclasses `GenericServlet`, `HttpServlet`)

WEB (SERVICE|SERVLET) CONTAINER

- a piece of software that extends a Web Server with infrastructure to provide the runtime environment to run servlets as Web Services,
- hosts one or more Web Services that extend the container's base URL,
- the servlets' code must be accessible to the Web Service Container, usually located in a specific directory,
- controls the lifecycle of the servlets: (`init()`, `destroy()`)
- maps the incoming communication from ports via the URLs to the appropriate servlet invocation
Method `service(httpContents)`, mapped to `doGet(httpContents)`,
`doPut(httpContents)`.

ABSTRACTION LEVELS

- a Web Services Container contains several “projects” (eclipse terminology) or “applications”:
 - from the programmer’s view, a “project” is an eclipse project, as a package it is a single .war file, at the end, it is a subdirectory in the container.
Each project has an (internal) name (its directory name in the container), e.g. `project1` or `mondial`.
- Each project consists of one or more servlets:
 - each servlet has an (internal) name (relative to its directory name in the container), e.g. `project1` contains servlets `servlet1a` and `servlet1b`.
 - each servlet’s code is a class that extends `javax.HttpServlet`;

ABSTRACTION LEVELS: URL MAPPING

HTTP connections (GET, POST) received by the server are transparently forwarded to the servlets.

URLs are e.g.

```
http://www.example.org/services/2011/oneservice/handle1
```

```
http://www.example.org/mondial/sqlonline
```

- the Web Service Container has a base url;
`http://www.example.org.`
- the Web Service Container maps *relative paths* to projects (later: by tomcat's `server.xml`), e.g.
`/services/2011/oneservice` to `project1`,
`/mondial` to `mondial`.
- each project's configuration (later: `web.xml`) maps path tails to servlet ids, and servlet ids to servlet classes, e.g.
`/handle1` and `/` to `servlet1a` to `Project1.MyServletA`,
`/handle2` to `servlet1b` to `Project1.MyServletB`.

TOMCAT BASIC INSTALLATION

- Web Server with Web Service Container: Download and install Apache Tomcat
 - can *optionally, but not necessarily* be combined with the Apache Webserver
 - can be installed in the CIP Pool

- set environment variable (catalina is tomcat's Web Service Container)

```
export CATALINA_HOME=~/apache-tomcat-x.x.x
```

- configure server: edit

```
$CATALINA_HOME/conf/server.xml:
```

```
<!-- Define a non-SSL HTTP/1.1 Connector on port 8080 -->
```

```
<Connector port="8080" .../>
```

- start/stop tomcat:

```
$CATALINA_HOME/bin/startup.sh
```

```
$CATALINA_HOME/bin/shutdown.sh
```

- logging goes to

```
$CATALINA_HOME/logs/catalina.out
```

GENERAL SERVLET PREPARATION DIRECTORY STRUCTURE

MyProject project directory (anywhere outside tomcat)

MyProject/build.xml the ant file for compiling and deploying – see later.

MyProject/src the .java (and other) sources

MyProject/lib jars that are needed for building, but should not be copied to the Web Service Container.

Put javax.servlet.jar there (tomcat has own classes for servlets, this would create conflicts).

MyProject/WebRoot roughly, all this content is copied later to the Web Service Container.

MyProject/WebRoot/WEB-INF

the whole content of MyProject/WebRoot except WEB-INF is visible later (e.g., HTML pages can be placed here); the contents of WEB-INF is used by the Web Service Container.

MyProject/WebRoot/WEB-INF/web.xml web application configuration – see later.

MyProject/WebRoot/WEB-INF/classes compiled binary stuff,

MyProject/WebRoot/WEB-INF/lib used jars.

HTTP METHODS GET AND POST

HTTP GET and POST: request-response

HTTP GET should be used only if invocation does *not* change

- Request consists only of URL+parameters:

`http://www.example.org/mondial?type=city&code=D&name=Berlin&province=Berlin`

HTTP POST should be used if it has side effects or changes the state of the Web Service

- Request URL consists only of the plain URL,
- parameters (e.g. queries using forms) or any other information is sent via a stream

⇒ often also queries use POST

Response: always as a stream.

- other HTTP methods PUT (resource), DELETE (resource) are used in REST (Representational State Transfer) “architectures” (e.g. the eXist XML database and document management system uses REST)

SERVLET PROGRAMMING

- servlets handle *incoming* HTTP connections via doGet() (=react-on-get) and doPost() (=react-on-post) methods:

```
public class MyServletA extends HttpServlet
{ public void init(ServletConfig cfg) throws ServletException
  { ... initialization ... }
  protected void doGet(HttpServletRequest request,
                        HttpServletResponse response) throws ServletException
  { ... }
  protected void doPost(HttpServletRequest request,
                        HttpServletResponse response) throws ServletException
  { ... }
}
```

- doGet() and doPost() both read the HttpServletRequest and write the HttpServletResponse object.
- the HttpServletRequest is different between GET (simpler) and POST (including a stream).

SERVLET PROGRAMMING - SOME METHODS

`doGet/doPost(HttpServletRequest req, HttpServletResponse resp)` throws ...

```
{ String path = req.getPathInfo();
```

tail of the URL path, e.g. do

```
if (path.equals("/handle1")) { ... }
```

```
java.util.Map<java.lang.String,java.lang.String[]> req.getParameterMap();
```

(doGet) returns a `java.util.Map` of the parameters of this request.

```
ServletInputStream req.getInputStream(); or
```

```
java.io.BufferedReader req.getReader();
```

(doPost) retrieves the body of the request as binary data using a `ServletInputStream`.

(doPost) retrieves the body of the request as character data using a `BufferedReader`.

```
PrintWriter out = response.getWriter();
```

yields a `Writer` to the response – send character text.

```
ServletOutputStream os = response.getOutputStream();
```

yields an output stream. Don't forget `os.flush()`.

INVOKING A NEW HTTP CONNECTION

Strongly stripped fragment (without try-catch etc.):

```
public static StringBuffer connectAndSend(String url, StringBuffer content) {
    HttpURLConnection.setFollowRedirects(true);
    HttpURLConnection con = (HttpURLConnection) new URL(url).openConnection();
    con.setRequestMethod("POST");
    con.setDoInput(true);
    con.setDoOutput(true);
    con.setRequestProperty("Connection", "keep-alive");
    con.setConnectTimeout(timeout);
    con.setRequestProperty("Content-type", "text/xml");
    write(con.getOutputStream(), content);
    con.getOutputStream().close();
    con.connect();
    StringBuffer response = read(con.getInputStream());
    con.getInputStream().close();
    return response;    }
```


A NOTE ON MULTITHREADING

- servlets can be instantiated by the container permanently or on-demand.
- if multiple requests for the same servlet come in, the servlet container can run multiple threads on the same instance of a servlet.
 - be careful with instance variables,
 - implement mutual exclusion if necessary
- the server can also create (and remove) additional instances of a servlet.

THE PROJECT'S WEB.XML

```
<web-app>
  <!-- Define servlet names and associate them with classfiles -->

  <servlet>
    <servlet-name>servlet1a</servlet-name>
    <servlet-class>org.dbis.project1.MyServletA</servlet-class>
    <init-param> ... </init-param>
  </servlet>
  <servlet> ... </servlet>

  <!-- define mapping of path tails to servlets -->

  <servlet-mapping>
    <servlet-name>servlet1a</servlet-name>
    <url-pattern>/handle1</url-pattern>
  </servlet-mapping>
  <servlet-mapping>
    <servlet-name>servlet1a</servlet-name>
    <url-pattern>/</url-pattern>
  </servlet-mapping>
  <servlet-mapping>
    <servlet-name>servlet1b</servlet-name>
    <url-pattern>/handle2</url-pattern>
  </servlet-mapping>
</web-app>
```

SERVLET DEPLOYMENT

- upon startup, tomcat deploys all servlets that are available in
`$CATALINA_HOME/webapps`
(considering path mappings etc. in `$CATALINA_HOME/server.xml`)

Two alternatives how to make servlets available there:

- create a `project1.war` file (web archive, similar to jar) and copy it into
`$CATALINA_HOME/webapps`.
(e.g. `build.xml` targets "dist" and "deploy")
(tomcat will unpack and deploy it upon startup)
- create a directory `project1`, copy the everything that is in the `WebRoot` directory there.
(e.g. `build.xml` target "deploy")

TOMCAT'S CONF/SERVER.XML

The URL paths to the projects have to be defined.

This is done in the `<Host>` element:

```
<Host>  
  <Context path="/services/2011/oneservice" reloadable="false" docBase="project1"/>  
  :  
</Host>
```

(since the standard path `mondial` for `mondial` is used, no explicit entry is needed)

- `reloadable`: automatically reloads the servlet if the code is changed (e.g. a new `.war`). Should be done only during development.
- the `path` attribute is key. There can be multiple paths that are mapped to the same `docBase`.

SOME COMMENTS

- HTTP connections are Unicode.
- exchanging XML via HTTP must be programmed via String/StringBuffer
 - out: serialize XML,
 - in: put SAX or StAX on the Unicode stream.
- exchanging StAX's XMLStream:
 - out: have an XMLStream, run StAX's next() iterator on it,
 - serialize each event, put it in the HTTP stream
 - in: create a naked StAX XMLStreamReader on the HTTP stream that just provides the XMLStream.

Chapter 12

Between Relational Data and XML

Data integration between “Legacy Systems” and XML databases

- Note: “legacy” now means SQL ...

Mixing everything up ...

Access to data stored in relational databases by

- using an XML environment (e.g., saxon) and mapping relational data from a remote SQL database (e.g. Oracle) to XML, and working with it.
- using XML-world concepts and languages in an SQL environment, e.g. for exchanging data in XML format (again, e.g., Oracle).

(Note that IBM DB2/XML Extender and MS SQL Server provide similar functionality)

12.1 Publishing/Mapping Relational Data in XML

Several *generic* mappings are possible:

Consider country(name: "Germany", code: "D", population: 83536115, area: 356910)

- tables, rows, and subelements

```
<table name="country">
  <row><name>Germany</name><code>D</code>
    <population>83536115</population><area>356910</area></row>
  :
</table>
```

- tuples with subelements

```
<country><name>Germany</name><code>D</code>
  <population>83536115</population><area>356910</area></country>
  :
```

- analogous with XML attributes
- advantage with subelements (vs. attributes): SQL values can also be object types (mapped to XML) and XMLType contents!

Example: HTTP-based access to Oracle

The whole database is mapped (virtually) to XML:

```
<SCHMIDT>    -- user name as root element
  <COUNTRY>  -- all names are capitalized
    <ROW><NAME>Germany</NAME><CODE>D</CODE>
      <POPULATION>83536115</POPULATION><AREA>356910</AREA></ROW>
    :
  </COUNTRY>
  :
</SCHMIDT>
```

Access by extended URL notation:

- URL: *computer:port/oradb/user/tablename/ROW[condition]*
- capitalize user, table and attribute names
- URL must select a single element (whole table, or single row)

```
ap34.ifi...:8080/oradb/DUMMY/COUNTRY           %% show page source
ap34.ifi...:8080/oradb/DUMMY/COUNTRY/ROW [CODE='D']
ap34.ifi...:8080/oradb/DUMMY/COUNTRY/ROW [CODE='D'] /NAME
ap34.ifi...:8080/oradb/DUMMY/COUNTRY/ROW [CODE='D'] /NAME/text ()
```


Generic Mappings (Cont'd)

Up to now: mapping of materialized base tables.

Problem: how to map the result of a query with computed columns?

```
SELECT Name, Population/Area FROM country
```

- tables, rows, and subelements:
the DTD is independent from the relational schema
metadata is contained in the attributes
("JDBC-style" processing of result sets)

```
<table name="country">  
  <row><column name="name">Germany</column>  
    <column name="population/area">83536115</column>  
    <column name="area">234.05473</column>  
  </row>  
  :  
</table>
```

- another "most generic mapping" as (object, property, value) to be discussed later ...

Additionally: often, tools define their own access functionality ...

ACCESS TO SQL DATABASES WITH SAXON-XSLT

[does currently not work]

- uses JDBC technology for remote access (at least for Java XSL tools)
- defines namespace “sql”
- `<sql:connect>` with attributes “database” (JDBC url), “driver” (JDBC driver) returns a JDBC connection object as a value of type “external object” that can be bound to a variable, e.g. `$connection`.

Note: there can be several connections at the same time.

- `<sql:query>` with following attributes allows to state an SQL query whose result is **generically** mapped to XML:
 - connection
 - table: ... the “FROM” clause
 - column: ... the “SELECT” clause
 - where: optional condition
 - row-tag: tag to be used for rows (default: “row”)
 - col-tag: tag to be used for columns (default: “col”)

result is a collection of `<row> ... </row>` elements that can e.g. be bound to a variable.

Administrative Parameters

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="2.0"
  xmlns:sql="java://net.sf.saxon.sql.SQLElementFactory">

  <!-- insert your database details here, or supply them in parameters -->
  <xsl:param name="driver" select="'oracle.jdbc.driver.OracleDriver'"/>
  <xsl:param name="database" select="'jdbc:oracle:thin:USER/PASSWD@IPADDRESS:1521:USER'"/>
  <xsl:param name="user" select="'USER'"/>
  <xsl:param name="password" select="'PASSWD'"/>

  <xsl:variable name="connection" as="java:java.sql.Connection"
    xmlns:java="http://saxon.sf.net/java-type">
    <sql:connect xsl:extension-element-prefixes="sql"
      driver="{ $driver }" database="{ $database }"
      user="{ $user }" password="{ $password }">
      <xsl:fallback>
        <xsl:message terminate="yes">SQL extensions not installed</xsl:message>
      </xsl:fallback>
    </sql:connect>
  </xsl:variable>
</xsl:stylesheet>
```

[Filename: XSLT/sql-administrativa-fake.xml]

Example Access/Query

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="2.0"
  xmlns:sql="java:/net.sf.saxon.sql.SQLElementFactory">
<xsl:include href="sql-administrativa.xsl"/>

<xsl:template match="*">
  <xsl:call-template name="countries"/>
</xsl:template>

<xsl:template name="countries">
  <xsl:variable name="country-table">
    <!-- ** the query's result is bound to the surrounding variable ** -->
    <sql:query xsl:extension-element-prefixes="sql"
      connection="$connection" table="country" column="*" />
  </xsl:variable>
  <xsl:copy-of select="$country-table" />
  <sql:close connection="$connection" />
</xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/sql-query.xsl]

uses a primitive mapping that relies on the order of columns.

Example Access/Query

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="2.0"
  xmlns:sql="java:/net.sf.saxon.sql.SQLElementFactory"
  extension-element-prefixes="sql">
<xsl:include href="sql-administrativa.xsl"/>

<xsl:template match="*">
  <xsl:call-template name="countries"/>
</xsl:template>

<xsl:template name="countries">
  <xsl:variable name="country-table">
    <sql:query connection="$connection" table="country,encompasses"
      where="country.code=encompasses.country"
      column="country.name,encompasses.continent,percentage"
      row-tag="country" column-tag="bla"/>
  </xsl:variable>
  <xsl:copy-of select="$country-table"/>
  <sql:close connection="$connection"/>
</xsl:template>
</xsl:stylesheet>
```

[Filename: XSLT/sql-query2.xsl]

Further Commands

- `<sql:insert>` with attribute
 - attribute: `table="..."`
 - children: `<column name="..." select="..."/>`
value can also be given as contents of the column element; currently always interpreted as a string.
- `<sql:close>` with a “connection” attribute

See also <http://www.saxonica.com/documentation/documentation.html>

12.2 The SQL/XML or SQLX standard

Goal: Coexistence between Relational Data and XML

- mapping relational data to XML
 - by a default mapping (previous section)
 - a (user-defined) XML views over relational data (“XML publishing”)
- storing XML data in relational systems
 - data-centric XML: efficient, several possibilities
 - document-centric XML: problematic

SQL/XML

- draft for an ISO standard since 2003: www.sqlx.org
- an SQL datatype “XMLType” for storing XML as “value” in databases:
Jim Melton (Oracle): *“SQL/XML includes an XML datatype for use inside SQL. This allows the SQL engine to automatically ‘shred’ data-oriented XML documents for storing some elements and contents in one place while keeping other elements in other places. Using indexes and joins, you can then bring this data back together very quickly.”*
 - **theory: an abstract datatype with operators/methods**
(cf. Computer Science I lecture)
 - **API: like (user-defined) object types in object-relational databases**
(cf. SQL database lab)
 - * handled with special methods (constructors, selectors)
 - * can be exported to XML tools, e.g. by JDBC
(either as DOM instance, or *serialized* as ASCII)
 - * libraries provide additional functions for processing XML in PL/SQL.
 - **internal implementation: not seen by the user**
(i) shredding or (ii) storing as LOB (Large Object)

SQL/XML

Making XML data a first-class citizen in relational databases,
seamless flow from relational data to XML and back

SQL to XML

- XML generation from SQL queries (i.e., in the SELECT clause)
(e.g., as packets for data exchange)
- define XMLType views over SQL data

SQL access and manipulation of XML inside the RDB

... use XPath expressions inside SQL (and rise the question what is actually the difference to XQuery):

- storing XML in RDB (e.g. if XML-data exchange packets came in),
- XPath-based extraction of XML content (SELECT clause),
- XPath-based query of XML content (WHERE clause),
- XPath-based update of XML content (in SQL),
- define XPath-based relational views over XML content.

“XML” AS AN SQL DATATYPE

XML/XMLType: an SQL datatype to hold XML data:

CREATE TABLE mondial OF XMLType; use as Row Object Value

CREATE TABLE CountryXML OF XMLType; use as Row Object Values

As **column object type** in relational tables:

```
CREATE TABLE CityXML
  (name XMLType, province VARCHAR2(35), country VARCHAR2(4),
   population XMLType,
   coordinates XMLType);
```

[Filename: SQLX/cityxmltable.sql]

- generation: INSERT INTO *table* VALUES (... , XMLType('XML as ASCII') ...)

```
INSERT INTO CityXML
VALUES(XMLType('<name>Clausthal</name>'), 'Niedersachsen', 'D',
XMLType('<population year="2004">10000</population>'),
XMLType('<coordinates><longitude>10.4</longitude><latitude>51.8</latitude>
</coordinates>'));
```

[Filename: SQLX/cityxmltuple.sql]

HANDLING OF SQL XMLTYPE DATATYPE

- generate it by certain *constructors* (“XML Publishing Functions”)
- storage: chosen by the database
 - “shredding” and distributing over suitable tables (of object-relational object types) (queries are translated into SQL joins/dereferencing)
 - storing it as VARCHAR, CLOB (Character Large Object), or as separate file (the remainder of this section uses CLOB)
 - storing it “natively”
- query it by XPath
- export/exchange in ASCII:
XMLSerialize: a function to serialize an XML value as an ASCII character string:
[XMLSerialize: XMLType → String](#)
- additional methods provided by PL/SQL libraries,
- XML objects can also be used e.g., as documents or as stylesheets, applied to documents (by PL/SQL libraries).

HOW TO GET XMLTYPE INSTANCES

- by the *opaque* constructor
XMLType: STRING → ELEMENT
that generates an XMLType instance from an ASCII string
 - the inverse to Java's `to_string`,
 - nearly all datatypes have such an opaque constructor (e.g., for lists: `list("[1,2,3,4,5]");`);
- generate instances recursively by structural constructors that are closely related to the underlying *Abstract Datatype* (cf. binary trees, lists, stacks in Computer Science I);
- or load them from an XML file (that then actually contains the ASCII serialization and uses the opaque constructor).

LOADING XMLTYPE INSTANCES

- from files:

- tell Oracle where it finds XML files:

```
CREATE OR REPLACE DIRECTORY XMLDIR AS '/tmp/' ;
```

```
CREATE OR REPLACE DIRECTORY XMLDIR AS '/home/bla/...' ;
```

- * Admin only, or after GRANT CREATE ANY DIRECTORY,
 - * on the same computer where the DBMS server (!) is installed.
- copy XML file (e.g. m.xml) into XMLDIR
 - * the XML file must not contain a reference to a DTD!
 - * the file must be publicly readable – `chmod filename 644`
 - load the file into the database (use the xdb_utilities package that must be installed separately)

```
INSERT INTO mondial
```

```
VALUES(xdb_utilities.getXMLfromfile('m.xml','XMLDIR'));
```

```
set long 10000 ;    -- number of characters in the output
```

```
SELECT * FROM mondial;
```

Admin's Note on Creating Directories

```
admin: GRANT CREATE ANY DIRECTORY TO scott;  
scott: CREATE OR REPLACE DIRECTORY XMLDIR AS '/tmp';  
SELECT owner, object_name FROM all_objects WHERE object_type='DIRECTORY';  
SELECT * FROM all_directories;
```

- The directory belongs then to SYS
- there is only one such directory declaration for the whole system
- “delete” it (i.e., make it unknown) `DROP directory XMLDIR`
- allow users to use the directory for reading:

```
GRANT READ on XMLDIR TO scott;
```

LOADING XML FILES: LOCAL SOLUTION

- for importing XML files, our local installation provides a method

```
system.getxml('http-url');
```

```
SELECT system.getxml(  
  'http://www.dbis.informatik.uni-goettingen.de/Teaching/DBP/XML/mondial.xml')  
FROM dual;
```

or, inserting it into a table:

```
INSERT INTO mondial VALUES(  
  system.getxml(  
    'http://www.dbis.informatik.uni-goettingen.de/Teaching/DBP/XML/mondial.xml'));
```

[Filename: SQLX/insertmondial.sql]

- that the XML file must not contain a reference to a DTD!
- the file can e.g. reside in the local homedirectory or anywhere in the Web,
- note that the file must be publicly readable – `chmod filename 644`

```
SET LONG 10000;  
SELECT * FROM mondial;
```

SQL/XML: GENERATING XML BY XML PUBLISHING FUNCTIONS

The SQL/XML Standard defines “XML publishing functions” that act as constructors (the name comes from the fact that they are also used to publish relational data in XML format):

- constructors of the recursively defined abstract datatype “XMLType”,
- create fragments or instances of XMLType,
- usage in the same way as predefined or user-defined functions (e.g., in the SELECT clause),
- output (e.g. in SELECT) in ASCII notation.

Some Theory: the Abstract Datatype

... constructors of the recursively defined abstract datatype “XMLType”:

Sub-datatypes:

- ELEMENT for element nodes
- ATTRIBUTE for attribute nodes
- QNAME for names of elements and attributes
(restriction of STRING without whitespaces etc.)
- STRING for text values (text nodes and attribute values)
- TUPLE(*type*) for a tuple of instances of *type*
- TABLE(*type*) for a table of instances of *type*

Constructors are very similar to those of XQuery (in the return clause), e.g.,:

```
element name attrs-and-content
```

and those of XSLT: `<xsl:element name="..."> content </xsl:element>`

and those of the ... DOM.

(always the same abstract datatype, but expressed with different syntaxes)

SQL/XML PUBLISHING FUNCTIONS: OVERVIEW

Basic constructors:

- XMLType: generates an XMLType instance from an ASCII string (“opaque constructor”)
 $\text{XMLType: STRING} \rightarrow \text{ELEMENT}$
- XMLElement: generates an XML element with a given name and content (either text (simple contents) or recursively created XML (complex contents) or mixed)
 $\text{XMLElement: QNAME} \times (\text{STRING} \cup \text{ELEMENT} \cup \text{ATTRIBUTE})^* \rightarrow \text{ELEMENT}$
 $\text{XMLElement: QNAME} \rightarrow \text{ELEMENT}$ for empty elements
- XMLAttributes: generates a one or more attribute nodes from a sequence of name-value-pairs
 $\text{XMLAttributes: (QNAME} \times \text{STRING)}^+ \rightarrow \text{ATTRIBUTE}^+$

SQL/XML PUBLISHING FUNCTIONS: OVERVIEW (CONT'D)

Further constructors:

- XMLForest: a function to generate a sequence, called a "forest," of XML elements with simple contents from a *sequence* of name-value-pairs
 $\text{XMLForest: (QNAME} \times \text{STRING)}^+ \rightarrow \text{ELEMENT}^+$
(note: the analogue to XMLAttributes for simple elements)
- XMLAgg: a function to group, or aggregate, XML data in a table into a sequence of nodes
 $\text{XMLAgg: TABLE(XMLTYPE)} \rightarrow \text{XMLTYPE}^*$
(note that a table is different from a list as in XMLForest!)
- XMLConcat: a function to concatenate the components of a tuple into a sequence
 $\text{XMLConcat: TUPLE(XMLTYPE}^+) \rightarrow \text{XMLTYPE}^*$
(note that a tuple is also different from a list as in XMLForest!)
- [XMLNamespaces: a function to declare namespaces in an XML element]

CONSTRUCTING XML: ELEMENTS

Basic form: XMLElement

- XMLElement: Name × Element-Body → Element:
 - Element-Body: text or recursively generated (attributes, elements, mixed)

```
SELECT XMLElement("x") FROM DUAL;
```

(note: this result is not correct: <x/> is an empty Element, while <x></x> is an element with the empty string as contents!)

```
SELECT XMLElement("Country", 'bla') FROM DUAL;
```

```
SELECT XMLElement(Country, 'bla') FROM DUAL;
```

- note: using “...” to indicate non-capitalization (otherwise the whole name is capitalized). (note that single and double “...” must be used exactly as in the example).

Elements with Non-Empty Content

- XMLElement: second argument contains the element body (attributes, subelements, text),
- XMLAttributes: list of name-value pairs that generate attributes.

```
SELECT XMLElement("Country",
  XMLAttributes(code AS "car_code", capital AS "capital"),
  name,
  XMLElement("Population",population),
  XMLElement("Area",area))
FROM country
WHERE area > 1000000;
```

[Filename: SQLX/xmlelement.sql]

A result element:

```
<Country car_code="R" capital="Moscow">
  Russia
  <Population>148178487</Population>
  <Area>17075200</Area>
</Country>
```

Optional Substructures

- XML as abstract datatype, *functional* constructors
- semistructured data: flexible and optional substructures

```
SELECT XMLElement("City",
  XMLAttributes(country AS country),
  XMLElement("Name",name),
  CASE WHEN longitude IS NULL THEN NULL
        ELSE XMLElement("Longitude",longitude) END,
  CASE WHEN latitude IS NULL THEN NULL
        ELSE XMLElement("Latitude",latitude) END)
FROM city
WHERE longitude IS NOT null;
```

[Filename: SQLX/xmlelement2.sql]

- Note: CASE WHEN *cond* THEN *a* ELSE *b* END
is a *functional* construct
(like in “if” in XQuery and <xsl:if> in XSLT)

CONSTRUCTING XML: SEQUENCES OF ELEMENTS

XMLForest: short form for simple elements

```
SELECT XMLElement("Country",
                 XMLForest(name AS Name,
                           code AS car_code,
                           population AS "Population",
                           area AS "Area"))
FROM country
WHERE area > 1000000;
```

[Filename: SQLX/xmlforest.sql]

```
<Country>
  <NAME>Brazil</NAME>                                <!-- note capitalization -->
  <CAR_CODE>BR</CAR_CODE>
  <Population>162661214</Population>
  <Area>8511965</Area>
</Country>
```

⇒ canonical mapping from tuples to XML elements with simple content.

Subqueries

Contents can also be generated by (correlated) Subqueries:

```
SELECT XMLElement("Country",
    XMLAttributes(code AS "car_code"),
    XMLElement("Name",name),
    XMLElement("NoOfCities",
        (SELECT count(*)
         FROM City
         WHERE country=country.code)))
FROM country WHERE area > 1000000;
```

```
SELECT XMLElement("Country",
    XMLAttributes(code AS "car_code"),
    XMLElement("Name",name),
    (SELECT XMLElement("NoOfCities",count(*))
     FROM City
     WHERE country=country.code))
FROM country WHERE area > 1000000;
```

[Filename: SQLX/xmlsubquery.sql]

CONSTRUCTING XML: RESTRUCTURING

- Relational databases have
 - literals (and objects for object-relational tables),
 - tuples,
 - tables.
- XML structures have
 - sequences
 - nesting (generated by XMLElement or by nested subqueries)

How to create

- a sequence from a tuple [note: XMLForest does not create a sequence from a tuple, but from a list that is generated from a tuple], or
- a sequence from a table?

CONSTRUCTING XML: GROUPING INTO A SEQUENCE

Aggregated Lists

- XMLAgg: is a new SQL aggregate function (like `count()` or `sum()`), that does not return a single value but the sequence of nodes of that table
- Note: XMLAgg is not applied to a *sequence* of items, but to a *table* of items!

Simplest case: XMLAgg over a table of XMLType row objects:

```
SELECT XMLElement("Cities",
  (SELECT XMLAgg(name)
   FROM CityXML c))
FROM DUAL;
```

[Filename: SQLX/xmlagg.sql]

Result:

```
<Cities>
  <name>...</name>
  <name>...</name>
  :
</Cities>
```

CONSTRUCTING XML: GROUPING

Aggregated Lists

... another example

- create a sequence of XML nodes from all entries of a relational table:

```
SELECT XMLElement("Continents",
  (SELECT XMLAgg(XMLElement("Continent", XMLAttributes(name AS "Name",
                                                         area AS "Area"))))
    FROM continent))
FROM DUAL;
```

[Filename: SQLX/xmlagg2.sql]

Result:

```
<Continents>
  <Continent Name="Europe" Area="..." />
  <Continent Name="Asia" Area="..." />
  :
</Continents>
```

CONSTRUCTING XML: NESTED GROUPING

Grouping/Aggregation: Nested Lists

- XMLAgg: generate a collection from the tuples inside of a GROUP BY:
In XML, this *list* of items can also be used!

... now we can have the number of cities in a country, together with a list of them:

```
SELECT XMLElement("Country",
    XMLAttributes(country AS car_code),
    XMLElement("NoOfCities", count(*)),
    XMLAgg(XMLElement("city",name) ORDER by population))
FROM city
GROUP BY country;
```

[Filename: SQLX/xmlgroupagg.sql]

Element of the result set:

```
<Country CAR_CODE="D">
  <NoOfCities>85</NoOfCities>
  <city>Erlangen</city> <city>Kaiserslautern</city> ... <city>Berlin</city>
</Country>
```

CONSTRUCTING XML: MAPPING TUPLES INTO SEQUENCES

XMLConcat

- takes a tuple of XML elements and transforms them into a sequence:

```
SELECT XMLElement("City", XMLConcat(name, population, coordinates))  
FROM CityXML  
WHERE country='D';
```

[Filename: SQLX/xmlconcat.sql]

An element of the result set:

```
<City>  
  <name>Berlin</name>  
  <population>...</population>  
  <coordinates><longitude>...</longitude><latitude>...</latitude></coordinates>  
</City>
```

Example: XMLConcat and XMLAgg

- the GROUP BY from Slide 542 can equivalently be expressed by using a (correlated) (Sub)query that returns a tuple for each country (consisting of the number and the aggregation of all cities):

```
SELECT XMLElement("Country",
  XMLAttributes(code AS code),
  XMLElement(name, name),
  (SELECT XMLConcat(
    XMLElement("NoOfCities", count(*)),
    XMLAgg(XMLElement("city",name)))
  FROM City
  WHERE country=code))
FROM country;
```

[Filename: SQLX/xmlconcatagg.sql]

Constructed XML can then be used for filling tables:

FILLING A TABLE ... WITH XML ROW VALUES

```
CREATE TABLE CountryXML OF XMLType;

INSERT INTO CountryXML
(SELECT XMLElement("Country",
  XMLAttributes(code AS "Code",
                population AS "Population"),
  XMLElement("Name",name),
  XMLElement("Area",area),
  (SELECT XMLElement("Capital",
    XMLForest(name AS "Name",
              population AS "Population")))
  FROM city
  WHERE country=country.code
    AND city.name=capital))
FROM country);
```

[Filename: SQLX/fillcountry.sql]

FILLING A TABLE: XML COLUMN VALUES

```
CREATE TABLE CityXML
  (name XMLType,
   province VARCHAR2(35),
   country VARCHAR2(4),
   population XMLType,
   coordinates XMLType);
INSERT INTO CityXML
(SELECT XMLElement("name", name), province, country,
 CASE WHEN population IS NULL THEN NULL
      ELSE XMLElement("population", XMLAttributes(95 as year), population)
      END ,
 CASE WHEN longitude IS NULL THEN NULL
      ELSE XMLElement("coordinates",
                     XMLElement("longitude", longitude),
                     XMLElement("latitude", latitude))
      END
      FROM city);
```

[Filename: SQLX/fillcity.sql]

HANDLING XML DATA FROM WITHIN SQL

- recall: XMLType is defined as an abstract datatype.
- it also has *selectors* that provide an interface for standard XML languages
- signature:
extract: XMLType × XPath_Expression → XMLType ∪ string
extractValue: XMLType × XPath_Expression → string
existsNode: XMLType × XPath_Expression → Boolean
- implementation based on user-defined object types
 - cf. object-relational extensions to SQL
- above operations also available as member methods

```
SELECT extract(value(m), '//city[name="Berlin"]') FROM mondial m;  
SELECT m.extract('//city[name="Berlin"]') FROM mondial m;  
SELECT extractValue(value(m), '//country[@car_code="D"]/population')  
FROM mondial m;  
SELECT m.extractValue('//country[@car_code="D"]/population')  
FROM mondial m; -- buggy (since version 9 ... and still in 11)!!!
```

SELECT: “Extract” Function

`extract(XMLType_instance, XPath_string)`

`XMLType_instance.extract(XPath_string)`

- First argument: selects an attribute with value of type “XMLType” in the current row (use `value(.)` function)
- Second argument: applies `XPath_string` to it
- Result: value of type XMLType or any other SQL type (multi-valued results are concatenated)

XML Row Values

Value of the row is of type XMLType: apply methods directly

```
SELECT extract(value(c), '/Country/@Code'),  
       extract(value(c), '/Country/Capital/Name')  
FROM CountryXML c;
```

```
SELECT c.extract('/Country/@Code'),  
       c.extract('/Country/Capital/Name')  
FROM CountryXML c;
```

ASIDE: SHORT OVERVIEW OF XPATH

(use the SQLX section in different lectures)

- Navigation as in Unix: `/step/step/step`
`/mondial/country/name`
- capitalization is relevant!
- result: a sequence of XML nodes (not only values, but also trees):
`/mondial/country`
- steps to deeper descendants: `/mondial//city/name` , `//city/name`
(latter includes `/mondial/country/city` and `/mondial/country/province/city`)
- attributes: `.../@attributname`: `/mondial/country/@area`
- access to text contents: `/mondial/country/name/text()`
- evaluation of conditions during navigation:
`/mondial/country[@code='D']/@area`
`/mondial/country[name/text()='Germany']/@area`
- Comparisons automatically use only the text contents:
`/mondial/country[name='Germany']/@area`

SELECT: "Extract" Function (Cont'd)

XML Column Values

Recall:

```
CREATE TABLE CityXML (name XMLType, province VARCHAR2(35),
    country VARCHAR2(4), population XMLType, coordinates XMLType);
```

CityXML.population is an XMLType column object:

```
SELECT extract(population,'/') FROM CityXML;
SELECT c.population.extract('/') FROM CityXML c;
SELECT name, extractValue(population,'/population/@YEAR'),
    extractValue(population,'/population')
FROM CityXML;
SELECT name, c.population.extract('/population/@YEAR').getNumberVal(),
    c.population.extract('/population/text()').getNumberVal()
FROM CityXML c
ORDER BY 3;
```

- exact capitalization in XPath argument!
- extractValue currently not implemented as member method (bug)
- use getNumberVal() and getStringVal() functions

SUBQUERIES TO XMLTYPE IN THE WHERE CLAUSE

... for selecting and comparing values, use also extract():

```
SELECT name
FROM CityXML c
WHERE c.population.extract('/population/text()')
      .getNumberVal() > 1000000;
```

```
SELECT c.extract('/Country/Name/text()')
FROM CountryXML c
WHERE c.extract('/Country/@Population')
      .getNumberVal() > 1000000;
```

- Note: comparison takes place on the SQL level (WHERE)
(→ join functionality when variables are used).
- Note: if the XPath expression returns a sequence of results, these are concatenated already during the evaluation of the extract() function ...
- ... thus, one has to use another way.

WHERE: “ExistsNode” Function

`existsNode(XMLType_instance, XPath_string)`

- Checks if item is of XMLType_instance, and XPath_string has a nonempty result set:
- note: the value for the comparison must be given in the XPath *string* – no joins on the SQL level possible.
- result: 1, if a node exists, 0 otherwise.

```
SELECT name, extractValue(population, '/population')  
FROM CityXML  
WHERE existsNode(population, '/population[text()>1000000]') = 1;
```

```
SELECT name, extractValue(population, '/population')  
FROM CityXML c  
WHERE c.population.existsNode('/population[text()>1000000]') = 1;
```

UPDATING XML DATA

- the complete XMLType value is *changed*, not updated
- updateXML(...) as a (transformation) function!
Note: the statement “SELECT updateXML(...) FROM ...” does not update the DB, but returns the value that would result from the update.

updateXML(XMLType_instance, XPath_string, new_value)

- first argument: SQL – selects an (SQL-)attribute of the current tuple (must result in an XMLType object),
- $2n$ th argument: selects the node(s) to be modified by the value of the ...
- $2n + 1$ th argument: new value,
- result: updates instance of type XMLType.
- Note: the expression “SELECT updateXML(...) FROM ...” does not change the database but returns only the value that *would* result from the update.

Updating XML Data (Cont'd)

```
SELECT updateXML(c.population,  
                'Population/text()','1000000',  
                'Population/@YEAR','2004')  
FROM CityXML c WHERE name='Gottingen';  
  
SELECT updateXML(value(c),  
                '/Country/Name/text()','Fidschi')  
FROM CountryXML c  
WHERE extractValue(value(c),'Country/Name')='Fiji';
```

[Filename: SQLX/updatexml.sql]

Updating XML Data (Cont'd)

This function is then used in the SQL SET-Statement:

```
UPDATE CityXML c
SET c.population -- an XMLType element
    = updateXML(c.population,
                '/population/text()', '1000000',
                '/population/@YEAR', '2004')
WHERE name='Gottingen'

UPDATE CountryXML c
SET value(c) = updateXML(value(c),
                          '/Country/Name/text()', 'Fidschi')
WHERE existsNode(value(c), '/Country[Name="Fiji"]') = 1
```

```

CREATE OR REPLACE FUNCTION xslexample RETURN CLOB IS
  xmldoc      CLOB;
  xsldoc      CLOB;
  myParser    dbms_xmlparser.Parser;
  indomdoc    dbms_xmldom.domdocument;
  xsltdomdoc  dbms_xmldom.domdocument;
  xsl         dbms_xslprocessor.stylesheet;
  outdomdocf  dbms_xmldom.domdocumentfragment;
  outnode     dbms_xmldom.domnode;
  proc        dbms_xslprocessor.processor;
  html        CLOB DEFAULT 'BLA';  -- must be initialized;
BEGIN
  -- Get the XML document as CLOB
  SELECT value(m).getClobVal() INTO xmldoc FROM mondial m;
  -- Get the XSL Stylesheet as CLOB
  SELECT s.stylesheet.getClobVal() INTO xsldoc
  FROM stylesheets s WHERE name='mondial-simple.xsl';

  -- Get the new xml parser instance
  myParser := dbms_xmlparser.newParser;
  -- Parse the XML document and get its DOM
  dbms_xmlparser.parseClob(myParser, xmldoc);
  indomdoc := dbms_xmlparser.getDocument(myParser);

  -- Parse the XSL document and get its DOM
  dbms_xmlparser.parseClob(myParser, xsldoc);
  xsltdomdoc := dbms_xmlparser.getDocument(myParser);

  xsl := dbms_xslprocessor.newstylesheet(xsltdomdoc, '');
  -- Get the new xsl processor instance
  proc := dbms_xslprocessor.newProcessor;

  -- Apply stylesheet to DOM document
  outdomdocf := dbms_xslprocessor.processxsl(proc, xsl, indomdoc);
  outnode     := dbms_xmldom.makenode(outdomdocf);

  -- Write the transformed output to the CLOB
  dbms_xmldom.writetoCLOB(outnode, html);
  -- Return the transformed output
  return(html);
END;
/
SELECT xslexample FROM dual;

```

[Filename: SQLX/xslexample.sql]

12.3 XQuery Support in SQLX

SQL function XMLQuery()

- SQL function `xmlquery('query' [passing vars clause] returning content)`
- XQuery function `ora:view(tablename)` turns tables into sequences of XML elements:
 - relational tables: Every row is turned into a ROW element as shown on Slide 514,
 - object table of XMLType: sequence of the XML elements in the object table, comparable to XQuery's `let`
- the result is the sequence of nodes as returned by the XQuery statement (of type "XML content").

```
SELECT
xmlquery(
'for $c in ora:view("countryXML")/Country
where $c/Capital[Population > 1000000]
return $c/Name'
returning content)
from dual;
```

```
SELECT
XMLElement("result",
xmlquery(
'for $c in ora:view("countryXML")/Country
where $c/Capital[Population > 1000000]
return $c/Name'
returning content)) from dual;
```

Passing XML parameters to XMLQuery()

Instances of XMLType can be selected in the SQL environment and passed to XMLQuery:

- context node
- variables

```
SELECT
XMLElement("result",
xmlquery(
  'for $c in ora:view("countryXML")/Country
  where $c/Capital[Population > $pop]
  return $c/Name'
  passing
    (SELECT population FROM City WHERE name='Tokyo') as "pop"
  returning content
)) from dual;
```

- comma-separated value-as-varname-list
- without “as ...”: context node
- “varname” cares for capitalization (\$POP and ... as pop would also be correct)

Syntax example

- Select names of all countries
- from the only XML element stored in table mondial (used as context element)
- that have a city that has a higher population than Tokyo (obtained from an SQL query)

```
SELECT
XMLElement("result",
  xmlquery(
    'for $c in //country
     where $c//city[population > $POP]
     return $c/name'
    passing
      (SELECT value(m) from mondial m),
      (SELECT population FROM City WHERE name='Tokyo') as POP
    returning content
  )) from dual;
```

XMLTable(): from XML contents sequences back to relational tables

- XQuery returns a sequence of nodes, which is of XML type “content” that can be put in an element (see above).

Turn the sequence of nodes into a table of rows:

- SQL function `xmltable('query' [passing vars clause] [COLUMNS column def clause])`
- column-def-clause is a comma-separated list of (datatype, column name, XPath expr.),
- default: a single XMLType pseudo-column, named COLUMN_VALUE,
- the result of XMLTable can be used like a relational table.

```
SELECT column_value
FROM
xmltable ('
  for $j in //country
  return $j/name'
  passing
    (SELECT value(m) FROM mondial m));
```

every row of the table is of XMLType and contains a <name>...</name>element

```
SELECT column_value
FROM
XMLTABLE ('
  for $j in //country
  return $j/name'
  PASSING
    (SELECT value(m) FROM mondial m))
WHERE column_value LIKE '%fr%';
```

Note: “like” is applied to the (contents of the) element.

XMLTABLE columns specification

```
SELECT *
FROM XMLTable ('
  for $j in //country
  return $j/name'
  PASSING
    (SELECT value(m) FROM mondial m)
  COLUMNS
    result XMLTYPE PATH '.',
    x VARCHAR2(35) PATH 'text()');
```

returns <name>...</name>elements.

- Additional specification of namespaces (for the paths): see documentation.

```
SELECT *
FROM XMLTable ('
  for $j in //country
  return $j'
  PASSING
    (SELECT value(m) FROM mondial m)
  COLUMNS
    name VARCHAR2(35) PATH 'name',
    area NUMBER PATH '@area',
    population NUMBER PATH 'population');
```

casts automatically to numbers.

Back and Forth: an example

- the result of XMLTable(...) can be used like a relational table:

```
SELECT u.column_value, u.column_value.extract('//Name/text()')
FROM (
SELECT t.column_value
FROM
XMLTABLE ('
  for $j in $X/*
  return $j'
PASSING
(xmlquery(
  'for $c in ora:view("countryXML")/Country
  where $c/Capital[Population > $pop]
  return $c/Name'
PASSING (SELECT population FROM city WHERE name='Berlin') as "pop"
RETURNING content)
) AS X) t) u
WHERE u.column_value.extract('//Name/text()') like '%ic%';
```

- or e.g. in insert: INSERT INTO ... (SELECT * FROM XMLTable(...)).

XQuery in SQLplus

- simple keyword “xquery”,
- returns the result of applying XMLTable (i.e., one row for each result of the xquery statement):

```
xquery
for $c in ora:view("countryXML")/Country
where $c/Capital[Population > 1000000]
return $c/Name
/
```

In contrast to many XML tools, attribute nodes are output as string values:

```
xquery
for $i in ora:view("mondial")/mondial/country
return $i/@car_code
/
```

Namespaces and Function Declarations

- as usual in XQuery:

```
xquery
declare namespace local = "http://localhost:8080/defaultNS";
declare function local:density($area, $pop)
{
    return $pop div $area
};
for $c in ora:view("mondial")/country
return local:density($c/population,$c/@area)
/
```

Functional Restrictions

(Oracle version 11.1.0.7)

- most XQuery/XPath functionality is supported (aggregation, context functions, string functions, path alternatives, ...)
- id(.) and idref(.) is not supported (recall that documents do not contain a DTD reference)
- any/all is not supported

INDEXES

- Indexes on XML data can be defined over any literal fields:

```
CREATE INDEX countrycodeindex
ON countryxml c
(EXTRACTVALUE(value(c), '//Country/@Code'));
CREATE INDEX countrycapnameindex
ON countryxml c
(EXTRACTVALUE(value(c), '//Country/Capital/Name'));
CREATE INDEX mondialcitynameindex
ON mondial m
(EXTRACTVALUE(value(m), '//Country//City/Name'));
```

12.4 XML Storage in Oracle

- CLOB (Character Large Object): Default.
XML is stored in its ASCII representation.
Note: for content management/delivery (e.g., to a Web server or as a Web service that just requires to get an ASCII stream) this is optimal.
Queries: XML is parsed internally and XPath/XQuery is applied.

- Binary XML

```
CREATE TABLE mondialBin OF XMLType
XMLTYPE STORE AS BINARY XML;
INSERT INTO mondialBin VALUES(
system.getxml(
'http://www.dbis.informatik.uni-goettingen.de/Teaching/DBP/XML/mondial.xml'));
```

- object-relational (only possible, if an XML Schema with oracle-specific annotations is preloaded)

STORAGE: PERFORMANCE COMPARISON

- BinaryXML is much faster
this example: 16:1

```
SET PAUSE OFF;  
SET TIMING ON;
```

```
select xmlquery('
  for $i in ora:view("mondial")/mondial
  let $city := $i//city
  let $country := $i/country
  where $city/@country = $country/@car_code
    and $city/@id = $country/@capital
  return $city/name
  'returning content)
from dual;
```

needs first time: 8.34,
then between 7.30 and 7.90

```
select xmlquery('
  for $i in ora:view("mondialbin")/mondial
  let $city := $i//city
  let $country := $i/country
  where $city/@country = $country/@car_code
    and $city/@id = $country/@capital
  return $city/name
  'returning content)
from dual;
```

needs first time: 0.42,
then between 0.28 and 0.30

Storage: Performance Comparison (Cont'd)

- join between two XPathS on a single XML table

```
SELECT * FROM XMLTABLE('
for $i in ora:view("mondial")//country
    $j in ora:view("mondial")//city,
where $i/@car_code = $j/@country
    and $i/@capital = $j/@id =
return $j/name/text()');
```

- without XMLTYPE STORE AS BINARY XML: to do
- with XMLTYPE STORE AS BINARY XML: to do

Storage: Performance Comparison (Cont'd)

```
CREATE TABLE mcity OF XMLType XMLTYPE STORE AS BINARY XML;
CREATE TABLE mcountry OF XMLType XMLTYPE STORE AS BINARY XML;
INSERT INTO mcity (SELECT COLUMN_VALUE FROM XMLTABLE(
    'for $c in ora:view("mondial")//city return $c'));
INSERT INTO mcountry (SELECT COLUMN_VALUE FROM XMLTABLE(
    'for $c in ora:view("mondial")//country return $c'));

CREATE INDEX mcountrycode ON mcountry c
    (EXTRACTVALUE(value(c), '//country/@car_code'));
CREATE INDEX mcountrycap ON mcountry c
    (EXTRACTVALUE(value(c), '//country/@capital'));
CREATE INDEX mcitycountry ON mcity c
    (EXTRACTVALUE(value(c), '//city/@country'));
CREATE INDEX mcityid ON mcity c
    (EXTRACTVALUE(value(c), '//city/@id'));

SELECT * FROM XMLTABLE('
for $i in ora:view("mcountry")/country,
    $j in ora:view("mcity")/city
where $i/@car_code = $j/@country and $i/@capital = $j/@id
return $j/name');    -- /text() -> error/bug
```


Storage: Performance Comparison (Cont'd)

- without XMLTYPE STORE AS BINARY XML: even for restricted size (cities > 200000 inhabitants, countries with area >1000000) 26 minutes.
- with XMLTYPE STORE AS BINARY XML: 3 minutes
- without indexes: first run needs longer (e.g., 26min/20min); then nearly same time as with indexes.

12.5 Background: XMLType as Object Type

(cf. "Practical Training in SQL" course)

Since SQL3 Standard: Object(-relational) types

- user-definable: CREATE TYPE AS OBJECT ... / CREATE TYPE BODY
- stored as *row* or *column objects*
CREATE TABLE cities OF CityObjectType;
- *member methods*
 - programmed in PL/SQL or recently also in Java
 - calls are embedded into SQL: SELECT *object.method(args)*
- reference attributes:
CREATE TABLE COUNTRY (... , capital REF CityType, ...);
SELECT c.capital ...;

⇒ now used for implementing XMLType

- as predefined internal classes/types
- can be used high-level from SQL, or low-level inside PL/SQL

XSLT IN ORACLE: “TRANSFORM” MEMBER METHOD

Member Method of XMLType: *XML-instance.transform(Stylesheet-as-XMLValue)*

as SQL function: *SELECT XMLTransform(XML-instance,Stylesheet-as-XMLValue)*

```
CREATE TABLE stylesheets
```

```
(name VARCHAR2(100),  
  stylesheet XMLTYPE);
```

```
INSERT INTO stylesheets VALUES('mondial-simple.xml',  
  system.getxml('http://www.dbis.informatik.uni-goettingen.de' ||  
  '/Teaching/DBP/XML/mondial-simple.xml'));
```

```
SELECT value(m).transform(s.stylesheet)  
FROM mondial m, stylesheets s  
WHERE s.name = 'mondial-simple.xml';
```

```
SELECT XMLTransform(value(m),s.stylesheet)  
FROM mondial m, stylesheets s  
WHERE s.name = 'mondial-simple.xml';
```

[Filename: SQLX/applystylesheet.sql]

Using built-in DOM, Parser, and XSL Engine

Tools from several packages can be explicitly used inside PL/SQL procedures:

- `dbms_xmlDOM`: implements DOM (usually, XML is transformed into DOM for processing it in detail)
PL/SQL call: `dbms_xmlDOM.dosomething(object,args)`
- `dbms_xmlparser`: parses documents from CLOB or URL, parses DTD from CLOB or URL (and stores the result);
access to the DOM instance/DTD in the parser then by “getdocument” or “getdoctype”
- `dbms_xslprocessor`: `processxsl(different arguments)`;
`clob2file/file2clob` allows for reading/writing;
`selectnodes/selectsinglenode/valueof`: XPath queries

... for details: Oracle Documentation, google ...

12.6 Storing XML Data in Database Systems

- “shredding” and distributing over suitable tables (of object-relational object types) (queries are translated into SQL joins/dereferencing)
 - Schema-based
 - Generic mapping of the graph structure
- storing it as VARCHAR, CLOB (Character Large Object), or as separate file with special functionality
- storing it “natively”/binary/model-based: internal object model

Literature

Klettke/Meyer “XML & Datenbanken” (dpunkt-Verlag), Ch. 8

Schöning: “XML und Datenbanken” (Hanser), Ch. 7,8

Chaudhri/Rashid/Zicari: XML Data Management

12.6.1 Mapping XML → Relational Model

two basic approaches:

- Schema-based: one or more “customized” tables for each element type
(→ similar to relational normalization theory)
 - (possibly) many null values
 - efficient access on data that belongs together
- one generic large table based on the graph structure:
(element-id, name of the property, value/id of the property)
 - no null values
 - although memory-consuming (keys/names that are stored once in (1) are now stored for each occurrence)
 - data that belongs together is split over several tuples

⇒ in both cases, theory and efficiency of relational database systems can be exploited.

SCHEMA-BASED STORAGE

necessary: DTD or XML Schema of the instance.

1. For each element type that has children or attributes, define a table that contains
 - a column that holds the primary key of the parent,
 - a primary key column if the element type has a member that satisfies (1) or (2),
 - for each scalar attribute and child element type with text-only contents that appears at most once, a column that holds the text contents.
2. for each multi-valued attribute or text-only subelement type that occurs more than once for some element type, a separate table is created with the following columns:
 - key of the parent node,
 - the (attribute or text) value(similar to 1:n relationships in the relational model).
 - for mixed content: possible solutions depend on the specific structure
 - special treatment for infrequent properties (to avoid nulls): handling in a separate XMLType column that holds all these properties together.

Schema-Based Storage: Example

For Mondial countries, provinces and cities, the following relations are created:

- country: key(mondial), key, name, code, population, area, ...
- border: ref(country), ref(other country), length
- language: ref(country), language, percentage
- province: ref(country), key, name, population, area
- city: ref(country), ref(province), key, name, longitude, latitude
- city-population: ref(city), year, value

Exercise

- give an excerpt of the instance
- translate some XPath/XQuery queries to SQL
- extended exercise: generate and populate the schema in SQL

Supported: Oracle (with augmented XML Schema), IBM DB2 (with DAD – Data Access Definition), MS SQL Server (extended Data-Reduced XML)

OBJECT-RELATIONAL APPROACH: INTERNAL OBJECT TYPES

- “shredded storage” of XML data is in general not implemented by plain relational tables, but using object-relational technology:
object types, collections, varrays etc ...
- collections/varrays: with value- and path indexes
- XPath expressions are rewritten into these structures

Integration of “legacy” object types

- application-dependent object types (as used in SQL3 in pre-XML times): standard mapping to XML (e.g. for data exchange)

Oracle & XML Schema

... register XMLSchema (must be typed in one single line!)

```
EXEC dbms_xmlschema.registerURI('http://mondial.de/m.xsd',  
    'http://dbis.informatik.uni-goettingen.de/Mondial/mondial.xsd');
```

can be deleted with

```
EXEC dbms_xmlschema.deleteSchema('http://mondial.de/m.xsd',  
    dbms_xmlschema.DELETE_CASCADE_FORCE);
```

- ... now, it knows `http://mondial.de/m.xsd` and created object tables for all root element types:

```
SELECT * from ALL_XML_TABLES;
```

```
CREATE TABLE mondial2 OF XMLType  
XMLTYPE STORE AS OBJECT RELATIONAL  
XMLSCHEMA "http://mondial.de/m.xsd"  
ELEMENT "mondial";
```

```
INSERT INTO mondial2 VALUES(  
system.getxml(  
    'http://www.dbis.informatik.uni-goettingen.de/Teaching/DBP/XML/mondial.xml'));  
SELECT XMLIsValid(value(m)) FROM mondial2 m;
```

GRAPH-STRUCTURE-BASED STORAGE

Without any schema knowledge, the graph structure can be represented in a single large table:

NodeNumber	ParentNode	[SiblingNo if ordered]	Name	Value
------------	------------	------------------------	------	-------

(see next page)

Alternatives

- separate table for elements and attributes (without node number and sibling number)
- separate between no-value, string value and numeric values for storing adequate types.
- previous-sibling and following-sibling columns instead of sibling-no (DOM style)

Querying

- requires recursive queries (PL/SQL; CONNECT BY)
- large joins (using the same large table several times)
- not implemented in any commercial system [according to Schöning 2003]

NodeNumber	ParentNode	SiblingNo	Name	Value
1	doc	1	mondial	
2	1	1	country	
3	2		@code	D
4	2		@membership	ref(<i>eu</i>)
:	:		:	:
41	2		@membership	ref(<i>un</i>)
42	2		@area	356910
43	2		@capital	ref(92)
44	2	1	name	Germany
45	2	2	population	83536115
:	:		:	:
90	2	47	province	
91	90	1	name	Berlin
92	90	2	city	
93	92		@country	ref(2)
94	92	1	name	Berlin
95	92	2	population	
96	95		@year	1995
97	95		text()	3472009
:	:		:	:

12.6.2 “Opaque” Storage

XML documents are stored as a whole as special datatype that can be used as row type or column data type (most commercial DBS; as described above for SQL/XML)

- approaches with text-based storage (CLOBs, files)
- specialized functionality for this datatype (cf. object-relational DBs: member functions)
 - XPath querying, XSLT support
 - validation
 - text search functions
- syntax embedded into SQL
- supported by indexes
 - full text indexes
 - path indexes/ “functional” indexes (user-defined, e.g. over //city/@country)
 - application and refinement of classical algorithms
- optimization of queries below the relational level!

12.6.3 “Native” Storage

Using “original” concepts of the database for storing XML (internal XML or object model) instead of mapping it or “simply” representing it as ASCII.

- often based on existing object-oriented DB-systems with application of concepts from hierarchical and network-DBs
- no document transformation to another data model
- data model/classes based on the notions of “tree”, “element”, “attribute”, “document order”
- navigation
- XPath/XQuery/XSQL APIs

“Native” Storage: Systems and Products

Many early implementations came from the object-oriented area:

- XML-QL, based on the Strudel system
- LoPiX, based on F-Logic
- Lorel-XML, based on Lorel/OEM
- Tamino (Software AG, Darmstadt, founded 1969 (Adabas, hierarchical DB), Tamino 1999, with XQL, first native XML DBMS),
- Excelon (until 1998: ObjectDesign with “ObjectStore”; since 12.2002: acquired by Progress Software Corp.)
- POET (Hamburg, “Persistent Objects and Extended Database Technology”, product 1990, spin-off from BKS 1992, OQL interfaces, SGML Document Repository 1997, Content Management Suite since 1998, merger with Versant 2004)
- Infonyte (GMD IPSI XQL 1998, based on a “Persistent DOM”, 12.2000: spinoff TU Darmstadt/Fraunhofer-IPSI)

GENERIC DATABASE BEHAVIOR FOR XML DATABASES

Everything that has been developed and discussed for relational databases is also relevant for XML:

- physical storage + storage management
- optimization, evaluation algorithms
- multiuser operation, transactions (ACID), safety, access control
- ECA-rules, triggers

The algorithms and theoretical foundations are very similar.

Often, relational (or hierarchical) DB technology is actually used inside.

COEXISTENCE OF XML AND RELATIONAL DATA

- generating XML (views, data exchange packets, ...) from stored relational data
- relational (and object-relational) techniques used for efficiently storing data-centric XML
- storing text-oriented data in RDB with specialized “native” datatypes
- XPath is also accepted by SQL/XML
- additional XML processing functionality by packages and object types
- XQuery is still not the “winner” for data-oriented applications!
- is it the winner for document-oriented applications?

<http://www.w3.org/TR/xquery-full-text/>

Chapter 13

Miscellaneous

Shortcomings in this lecture

- XML: namespaces (used e.g. for XSL, XLink – but in this lecture really interesting only in combination with RDF)
- XML: processing instructions
- theory: XML Data Model, XPath/XQuery Formal Semantics (currently in an alignment process with XML Schema)

Application areas:

- Web Services
- Multimedia applications with XML
- Semantic Web

XML APPLICATIONS (= LANGUAGES)

- XHTML
- MathML; see e.g. <http://www.w3.org/Math/XSL/>
- SMIL (Synchronized Multimedia Integration Language): description of multimedia presentations
- MPEG-7: meta-metalevel description of audiovisual contents (i.e., used to describe description languages)
- Web Services (kind of XML-style CORBA successor): WSDL (WS Description Language), WSFL (Web Services Flow Language; Workflows), SOAP (Simple Object Access Protocol – messaging format, not only for Web Services) UDDI (Universal Description, Discovery and Integration); ebXML, MS-BizTalk
- Health Care: HL-7; clinical data exchange
- BIOML (polymer structures), GML (geographic ML), CML (chemical ML)

EPILOGUE

What should have been taught?

- knowledge for practical use of XML
- XML is more than only angle brackets
- the XML world provides examples for many basic concepts of computer science and their combination,
- illustrating how concepts in computer science evolve, and
- ... an idea of developments in the near future:
the DBIS group is part of it with
<http://dbis.informatik.uni-goettingen.de/reverse/>

13.1 Overview of some Books ...

- There is not a single book that gives a good introduction the everything about XML
- several books on specific, advanced topics
- german books: comments in german.
- recommended sections are marked with →.
- very recommended sections are marked with ⇒.

Note: the selection of books is a bit randomly. There are also other good ones.

“XML Family of Specifications: A Practical Guide”; Kenneth Sall; Addison-Wesley, May 2002

- historical overview: “professional” focus, no mention of previous research topics
- XML, DTD, SAX; DOM, JDOM; CSS, XSL, XSLT; XLink, XPointer, XPath, XML Schema, RDF, ...
- comprehensive, but often superficial
- 2002 - a bit outdated.

“XPath, XLink, XPointer and XML”; Erik Wilde, David Lowe; Addison-Wesley, July 2002

- ⇒ I recommend the book for its excellent overview of concepts and ideas around XML and the Web (Sections 1-5)
- → Sect.6,7: XLink, XPointer
- Sect.8: Usage, Sect.9: Future

“XSLT Programmers Reference”, 2nd Edition; Michael Kay, Wrox Press, June 2003

- ⇒ *The Book on XSLT.*

“XQuery”, Wolfgang Lehner and Harald Schöning, dpunkt, 2004

- ⇒ *The (german) Book on XQuery.*

“XML & Datenbanken – Konzepte Sprachen und Systeme”; Meike Klettke Holger Meyer; dpunkt-Verlag, 2003

- Kap.3,4 (XML), Kap.10 (Anfragesprachen)
- Kap 2,5,6,7: Allgemeines zu XML und DB
- → Kap.8: XML-Datenbanken, Speicherungstechniken
- → Kap.9: Indexstrukturen
- → Kap.11: XML-Datenbanken: Systeme

“XML und Datenbanken”; Harald Schöning, Hanser-Verlag, 2004

- Kap.1: XML, DTD, DOM, SAX, XSL: oberflächlich
- Kap.2,4: XML Schema, Entwurf
- Kap.3: Allgemeines zum Einsatz von XML
- → Kap.5: XML in Datenbanksystemen: allgemeine Betrachtungen
- Kap.6: XPath/XQuery: oberflächlich
- ⇒ Kap.7,8: XML und Datenbanksysteme, allgemein sowie Produkte (Oracle, IBM, MS SQL, Tamino)

“XML und Datenmodellierung”; Rainer Eckstein, Silke Eckstein; dpunkt-Verlag, 2004

- grobe Einführung XML, DTD, XPath; (kein XQuery)
- ⇒ Kap.4: XML Schema
- (Kap 5: DTD/XML Schema und UML)
- ⇒ Kap.6,7: RDF, RFDS, (OWL)