# XPath-logic and XPathLog: A logic-programming style XML data manipulation language

WOLFGANG MAY

*Institut für Informatik, Universität Göttingen, Göttingen, Germany*
(*e-mail:* may@informatik.uni-goettingen.de)

## Abstract

We define XPathLog as a Datalog-style extension of XPath. XPathLog provides a clear, declarative language for querying and manipulating XML whose perspectives are especially in XML data integration. In our characterization, the formal semantics is defined wrt. an *edge-labeled graph-based* model, which covers the XML data model. We give a complete, logic-based characterization of XML data and the main language concept for XML, XPath. XPath-Logic extends the XPath language with variable bindings and embeds it into first-order logic. XPathLog is then the Horn fragment of XPath-Logic, providing a Datalog-style, rule-based language for querying and manipulating XML data. The model-theoretic semantics of XPath-Logic serves as the base of XPathLog as a logic-programming language, whereas also an equivalent answer-set semantics for evaluating XPathLog queries is given. In contrast to other approaches, the XPath syntax and semantics is also used for a declarative specification how the database should be *updated*: when used in rule heads, XPath filters are interpreted as specifications of elements and properties which should be added to the database.

*KEYWORDS*: XML, XPath, Logic Programming, Information Integration.

## 1 Introduction

Logic-based languages have proven useful in many areas since they allow for small, declarative, and extendible programs. For the database area, Datalog has been investigated for querying and rule-based data manipulation. Extending the Datalog idea, more complex logic-based frameworks like F-Logic (Kifer and Lausen, 1989; Kifer *et al.*, 1995), or the languages of the TSIMMIS project (Garcia-Molina *et al.*, 1997; Abiteboul *et al.*, 1997) have been successfully applied for knowledge representation and data integration. The experiences with a powerful language like F-Logic were the motivation to have a similar "native" language for the XML world that is much simpler than F-Logic, and that is based on the standard XPath language. As a result, we present XPathLog as an XPath-based Datalog-style language for querying and manipulating XML data. By extending XPath with variable bindings and providing a constructive semantics for XPath in rule reads, a declarative XML data manipulation language is obtained. Since both XPath and rule-based

programming by using variable bindings are well-known, intuitive concepts, the
"effect" of the language is easy to understand on an intuitive basis. Additionally, the
well-known logic programming semantics provide concise *global* semantics of such
programs which coincide with the intuitive ideas. Queries and rules for manipulating
and restructuring the internal XML database can be expressed much easier than e.g.
in XQuery (XQuery, 2001) (where update functionality is still in a prototypical state).

*Semistructured Data and XML.* XML has been designed and accepted as *the* frame-
work for semi-structured data where it plays the same role as the relational model
for classical databases. The XML data model applies both to *documents* and
to *databases*: The SGML language was originally defined for *documents* in the
publishing area. On the other hand, the interest in research on *semistructured* data
in the 1990s[1] (e.g. F-Logic (Kifer and Lausen, 1989; Kifer *et al.*, 1995), GraphLog
(Consens and Mendelzon, 1990), UnQL (Buneman *et al.*, 1996; Buneman *et al.*, 2000),
Tsimmis (Garcia-Molina *et al.*, 1997; Abiteboul *et al.*, 1997) with the *OEM* data
model and the *MSL*, *WSL*, and *Lorel* languages, Strudel/StruQL (Fernandez *et al.*,
1997; Fernandez *et al.*, 1998), and *YAT/YATL* (Cluet *et al.*, 1999)) was motivated
by the *database* community, searching for a data model for *data integration* and a
data format for *electronic data interchange*. Here, also the combination of document-
oriented aspects with database aspects was an important motivation to go beyond
classical data models which then resulted in the design of XML.

The XML data model is a hierarchical model which defines an ordered tree with
attributes that can easily be interpreted as a document. The natural relationships
in documents are either (i) substructures, or (ii) references to other parts of the
document (where the term *reference* here means simply a cross-reference in a
document). The nested elements define a document structure whose leaves are the
text contents. Elements (i.e. the structuring components) are annotated by attributes
which do *not* belong to the document contents. Inside the tree, cross-references
(`IDREF` attributes) are allowed.

In contrast, for databases, a hierarchical structure is in general not intuitive. Here,
several kinds of relationships have to be represented, between substructures and pure
references. When using XML for a database-like application, these relationships
have to be represented by reference attributes. On the other hand, order is often not
relevant in databases.

*Mainstream XML Languages.* Specialized languages have been defined for XML
querying, e.g. XQL (Robie, 1999), XML-QL (Deutsch *et al.*, 1999), then XPath
(XPath, 1999) developed from the experiences with XQL and XSL Patterns (and
XPointer) as an addressing language. XQuery (XQuery, 2001) extends XPath
with SQL-like constructs. XSLT (XSLT, 1999) is an XPath-based language for
transforming XML data. A proposal for extending XQuery with update constructs
(XUpdate) has been published in Tatarinov *et al.* (2001); a more detailed proposal
is described in Lehti (2001).

---

[1] We list the approaches in the temporal order of their presentation.

*Other Approaches to XML.* XML-GL (Ceri *et al.*, 1999; Comai *et al.*, 2001) continued the idea of GraphLog for XML. Elog (Baumgartner *et al.*, 2001a) is based on Datalog and classical first-order logic, flattening XML into predicates. It is used as an internal language in Lixto (Baumgartner *et al.*, 2001b). Bry and Schaffert (2002) present the Xcerpt language which regards XML trees as terms, similar to UnQL.

### *1.1 Comparison of design concepts*

Amongst the existing languages for handling semistructured data and XML, there are several facets for distinguishing them in terms of the concepts they use. A more detailed comparison with individual languages can be found in section 6.

*Data model.* Semistructured data can be regarded as a general graph (OEM, UnQL, Strudel, GraphLog, and F-Logic) or as a tree (YATL and XML). Moreover, node-labeled graphs/trees (XML) or edge-labeled graphs (as in Strudel, UnQL, GraphLog, and F-Logic) can be distinguished; for OEM both representations can be found. It is easy to represent a node-labeled instance in a labeled model, whereas the other way requires node replication. Also, ordered (e.g. XML) and unordered (e.g. in OEM, F-Logic, and UnQL) tree/graph models are distinguished.

*Access mechanism.* Generally, there are two approaches for selecting items in a semistructured data tree or graph:

- by matching patterns and templates (GraphLog, MSL/WSL, UnQL, YATL, and later for XML-QL and XML-GL). In UnQL and Xcerpt, (bi)simulation between semistructured data trees is used. If a simulation of a match pattern with variables by the underlying database is found, the appropriate variable bindings are returned and used for generating an answer tree.
- navigational access, like in object-oriented database languages (OQL), as done in Lorel, StruQL, and F-Logic, and later in XPath and also in our XPathLog approach. UnQL provides both patterns and a navigational syntax.

*Functionality.* There are different approaches to either *generating* an answer by instantiating a generating pattern in the rule head according to the variable bindings (UnQL, StruQL, and later XML-QL, XQuery and XML-GL), or *manipulating* the underlying structure by adding information to the underlying database (GraphLog, F-Logic and XUpdate).

Note that this distinction did not exist when considering classical rule-based languages, e.g. Datalog for predicate logic. The facts derived in the rule head were added to the database – either extensionally, or intensionally as view definitions – without directly interfering with the already stored facts. Other rules of the program could easily use both the original data and the derived data. For XML, a semantics where rules generate separate structures is easy to define. In contrast, a semantics where the rule heads interfere with the database contents has to take into account that the evaluation of rules may violate the tree structure.

*Underlying framework.* Some of the languages are based on a kind of model-theoretic semantics: UnQL and Xcerpt directly operate on tree-term structures, employing

and defining mechanisms like tree matching, term unification and (bi)simulation unification. Elog is based on Datalog and classical first-order logic, flattening XML data into predicates. F-Logic defines F-Structures that extend classical first-order logic, and then applies logic programming mechanisms to such models. For the other languages (Tsimmis, Strudel, XML-QL, XQuery, XUpdate), the semantics is directly defined in terms of data structures.

*Rule-based vs. logic programming.* All languages discussed above are *rule-based* and *declarative*, generating variable bindings by a matching/selection part in the "rule body" and then using these bindings in the "rule head" for generating output or updating the database. This rule-based nature is more or less explicit: F-Logic, MSL/WSL (Tsimmis), Elog, and Xcerpt use the ":-" Prolog syntax, whereas UnQL, Lorel, StruQL, XML-QL and XQuery/XUpdate cover their rule-like structure in an SQL-like clause syntax. These clausal languages allow for a straightforward extension with update constructs (as it has been done for Lorel and proposed with XUpdate for XQuery). GraphLog and XML-GL use a graphical representation.

The first, "pure" group separates strictly between the selection part in the rule body and generation/update part in the rule head, whereas UnQL, StruQL, XML-QL and XQuery allow for nested selection-generation parts in the rule bodies.

The global semantics of these languages is influenced by their functionality, distinguishing between query/transformation and query/update languages: UnQL, Xcerpt, XML-QL, and XQuery *generate* (output) structures in their head which are not fed back into the input or internal database.

Only MSL/WSL, Elog, and F-Logic allow to for *additions to the database* or *view definitions* (depending whether bottom-up or top-down semantics is considered), and to use the derived facts in the selection/matching part of other rules. StruQL, XML-QL, and XQuery overcome this restriction by nesting selection-generation parts in the rule bodies. The traverse construct of UnQL (applying a subquery by structural recursion to arbitrary depth) also comes near to local view definitions. Note that these languages require regular path expressions to compute the transitive closure of a binary relation (see Fernandez *et al.* (1997)). We consider the ability to compute a transitive closure as an important feature for a language for handling semistructured data (especially, for a "logic-programming" language, since that is what makes the distinction between Datalog and the relational algebra/calculus).

The difference between *rule-based transformation languages* and *logic programming* languages is mirrored by the fact that the semantics of UnQL, Xcerpt, XML-QL, and XQuery is completely given by he semantics of their rules (qualifying them as *rule-based* languages). In contrast, the *global* semantics of F-Logic and Elog also requires the notions of the $T_P$ operator and of minimal or well-founded models (qualifying them as *logic programming* languages). As a consequence, they require both a *model-theoretic semantics*, and an *answer semantics* for queries.

*Design principles for XPathLog.* XPathLog follows a logic-based approach which has been motivated by the experiences with F-Logic: XML instances are mapped to a semantical structure for interpreting *XPath-Logic* formulas. XPath-Logic is based on (i) first-order logic, and (ii) XPath expressions extended with variable bindings. The

Horn fragment of XPath-Logic, called *XPathLog*, provides a declarative, Datalog-style logic-programming language for manipulation of XML documents.

Regarding the above design principles, XPathLog is positioned as follows:

- XPathLog is completely XPath-based (i.e. navigational access). Matching and generation/update part are strictly distinguished.
  An extended XPath syntax is used for querying (rule bodies) *and* generating/manipulating data (rule heads). The rule body serves for generating variable bindings which are then used in the head for *extending* the current XML database, thereby defining an update semantics for XPath expressions.
- XPathLog uses an *edge-labeled graph* model, which is advantageous when defining several tree views of the internal database. The data model is partly ordered like in XML: subelements are ordered, attributes are unordered.
- XPathLog is a logic-programming language according to the above characterization. It works on an abstract semantical model which represents an *XML database* supporting multiple overlapping XML trees. These *X-Structures* together with the logic, called *XPath-Logic*, provide for a logical characterization of XML data. XPathLog combines the intuitive "local" semantics of addressing XML data by XPath with the appeal of the "global" semantics of logic programming. As an update language, it is based on a bottom-up semantics.
- In contrast to XML-QL, XQuery, and XSLT, the language does not use additional constructs whose semantics has to be defined separately: the only semantic prerequisite is the bottom-up evaluation strategy of Datalog or any other logic programming language.

In this paper, we describe the data model, its logical foundation, the internal semantics of queries, rules, and programs of XPathLog as a true logic programming language for XML. Some aspects have been published in May (2002); May and Behrends (2001) and with the LoPiX prototype in May (2001c). Here, we focus on the theoretical issues of modeling XML and the semantics of a language for queries and basic, elementary updates. A full report on XPathLog can be found in May (2001a).

A possible application area for XPathLog is, for instance, the integration of XML data from several sources, as done in the case study by May (2001b). Here, the power of the combination of XPath expressions with additional variable bindings allows for short and concise declarative and flexible rules. Both, queries and rules for manipulating and restructuring the internal XML database can be expressed much easier than, for example, in XQuery (where update functionality is still in a prototypical state).

*Structure of the paper*. Section 2 defines X-Structures as semantical structures which represent XML documents and presents XPath-Logic. The answer semantics of XPathLog as an XML query language is investigated in section 3. Section 4 defines the semantics of rule heads for generating and modifying XML data, and the semantics and evaluation of XPathLog programs. The implementation in the LoPiX (**Lo**gic **P**rogramming **i**n **X**ML) system and a case study are described in section 5.

An analysis, a general discussion of related work, and the conclusion can be found in section 6. Additional proofs can be found in Appendix Appendix A.

## 2 XPath-Logic: The model-theoretic framework

XPath-Logic and its Horn fragment, XPathLog, extend XPath (XPath, 1999) with Datalog-style variable bindings. XPath-Logic provides the model-theoretic framework for defining a global, logic-programming style semantics for XPathLog.

XPath (XPath, 1999) is the common language for addressing node sets in XML documents. It is based on navigation through the XML tree by *path expressions* of the form *root*/*axisStep*/.../*axisStep* where *root* specifies a starting point of the expression (the root of a document, or a variable that is bound to a node in an XML instance). Every *axisStep* is of the form *axis*::*nodetest* [*qualifier*]*. The *axes* define navigation directions in an XML tree: Given an element *e*, the child axis contains all its children and the attribute axis contains all its attributes. Analogously, parent, ancestor, descendant, preceding-sibling and following-sibling axes are defined. They enumerate the respective nodes by traversing the document tree starting from *e*.

First, along the chosen axis, all elements which satisfy the *nodetest* (which specifies the nodetype or an elementtype which nodes should be considered) are selected; the resulting list is called the *context*. Then, the given *qualifier(s)* (also called *filters*) are applied to each of the nodes (as the *context node*) for finer selection. Inside qualifiers, *relative location paths* are allowed that implicitly start at the context node. Starting with this (local) result set, the next step expression is applied (for details, see (XPath, 1999) or (XQFS, 2001)). The most frequently used axes are abbreviated as *path*/nodetest for *path*/child::*nodetest*, *path*/@*nodetest* for *path*/attribute::*nodetest*, and *path*//*nodetest* for *path*/descendant-or-self::*/ child::*nodetest*.

*Example 1* (*XML, XPath, Result Sets*)
Consider the of the Mondial database (May, 2001e) for illustrations; the DTD is given as follows:

```
<!ELEMENT mondial (country+, organization+, . . . )>
<!ELEMENT country (name, population, encompassed+, border*, city+, . . . )>
    <!ATTLIST country  car_code ID #REQUIRED  capital IDREFS #REQUIRED
                       memberships IDREFS #IMPLIED>
<!ELEMENT name (#PCDATA)>
<!ELEMENT encompassed EMPTY>
    <!ATTLIST encompassed  continent CDATA #REQUIRED
                           percentage CDATA #REQUIRED>
<!ELEMENT border EMPTY>
    <!ATTLIST border  country IDREF #REQUIRED   length CDATA #REQUIRED>
<!ELEMENT city (name, population*)> <!ATTLIST city country IDREF #REQUIRED>
<!ELEMENT population (#PCDATA)> <!ATTLIST population year CDATA #IMPLIED>
<!ELEMENT organization (name, abbrev, members+)>
    <!ATTLIST organization  id ID #REQUIRED   headq IDREF #IMPLIED>
<!ELEMENT abbrev (#PCDATA)>
```

```
<!ELEMENT members EMPTY>
    <!ATTLIST members  type CDATA #REQUIRED  country IDREFS #REQUIRED>
```

An excerpt of the instance is given below (and depicted as a graph can be found in Figure 1 when X-Structures are considered).

```
<country car_code="B" capital="cty-Brussels" memberships="org-eu org-nato . . . ">
    <name>Belgium</name>  <population>10170241</population>
    <encompassed continent="Europe" percentage="100"/>
    <border country="NL" length="450"/>  <border country="D" length="167"/>
        ⋮
    <city id="cty-Brussels" country="B">  <name>Brussels</name>
        <population year="95">951580</population>
    </city>
        ⋮
</country>
<country car_code="D" capital="cty-Berlin" memberships="org-eu org-nato . . . ">
        ⋮
</country>
<organization id="org-eu" headq="cty-Brussels">
    <name>European Union</name>  <abbrev>EU</abbrev>
    <members type="member" country="GR F E A D I B L . . . "/>
    <members type="membership applicant" country="AL CZ . . . "/>
</organization>
<organization id="org-nato" headq="cty-Brussels" . . . >
        ⋮
</organization>
```

The XPath expression

  //country[name]/city[population/text()>100000 and @zipcode]/name/text()

returns all names of cities such that the city belongs (i.e. is a subelement) to a country where a name subelement exists, the city's population is higher than 100,000, and its zipcode is known.

XPath is only an addressing mechanism, not a full query language. It provides the base for most XML query languages, which extend it with their special constructs (e.g. functional style in XSLT, and SQL/OQL style (e.g. joins) in XQuery). In the case of XPath-Logic and XPathLog, the extension feature are Prolog/Datalog style variable bindings, joins, and rules.

*Remark 1 (Relationship to W3C Documents)*
We restrict the considerations to the core concepts of XPath as an addressing and navigation formalism for XML data, i.e. stepwise navigation along the axes and step qualifiers/filters. For the XPath syntax and non-formal semantics, we always refer to the *W3C XPath 2.0 Working Draft* (XPath, 1999). Note that the syntax and semantics of the core concepts of XPath is the same as in XQL (Robie, 1999), XPointer, early drafts of XPath, XPath 1.0 (XPath, 1999), although both the presentation and the naming have been changed several times.

A formal semantics of XPath has been given as a *denotational* semantics in Wadler (1999) that already covers these central notions of XPath. Later, it has been

re-formulated first in the *W3C XML Query Algebra* (XMQ-A, 2001) and then in the *W3C Query Formal Semantics* (XQFS, 2001), where a description in terms of type inference rules and value inference rules is given. Note that since the early implementations (e.g. xt (Clark, 1998)), the actual semantics of XSL Patterns/XPath as its "behavior" has not changed. For comparing our approach with the formal semantics of XPath, we refer to Wadler (1999), which gives a short and concise definition of the central concepts that is best suited as a reference.

### 2.1 Syntax of XPath-Logic

Inspired by the derivation of F-Logic from first-order logic as a logic for dealing with structures containing complex objects, XPath-Logic is defined for expressing properties of XML structures. The main difference between XPath-Logic and first-order logic is that XPath-Logic has an additional type of atomic formulas: *reference expressions* which turn out to be a special kind of predicates with a built-in semantics. The "basic" components of the language are XPath-Logic *PathExpressions* which are syntactically derived from XPath's *PathExpressions* by extending Path with Prolog/Datalog style variable bindings.

*Definition 1* (*XPath-Logic: Syntax*)
The set of basic formulas of an XPath-Logic language is defined as follows:

- every language contains an infinite set Var of variables.
- a specific XPath-Logic language is given by its *signature* $\Sigma$ of element names, attribute names, function names, constant symbols, and predicate names.
- *XPath-Logic reference expressions* over the above names extend the XPath *path expressions*: The syntax of *AxisSteps*, **axis**::*name*[*stepQualifier*]*, may be extended to bind the selected nodes to variables by "-> Var":

```
Step ::= Axis "::" NodeTest StepQualifiers
       | Axis "::" NodeTest StepQualifiers "->" Var StepQualifiers
       | Axis "::" Var StepQualifiers
       | Axis "::" Var StepQualifiers "->" Var StepQualifiers
```

  For an XPath-Logic reference expression, the *underlying XPath expression* is obtained by removing the inserted variable binding constructs.
- An *XPath-Logic predicate* is a predicate over reference expressions.
- *terms* and *atomic formulas* are defined analogously to first-order logic.
- XPath-Logic *compound formulas* are built over predicates and reference expressions, using $\land$, $\lor$, $\neg$, $\exists$, and $\forall$.
- XPath-Logic allows to have formulas in step qualifiers.

Note that XPath-Logic does not use the explicit dereference operator $\Rightarrow$ from XPath 2.0; instead implicit dereferencing of attributes in paths is supported.[2]
The goal of this paper is to introduce the Horn fragment from XPath-Logic, called XPathLog as a Datalog-style XML query and manipulation language. The following

---

[2] XPath-Logic has been designed before XPath 2.0 replaced XPath's id(.) function by the dereferencing operator. Furthermore, we use a data model that directly incorporates references.

example gives some XPathLog queries that review the basic XPath constructs, and illustrate the use of the additional variable binding syntax.

*Example 2* (*XPathLog: Introductory Queries*)
The following examples are evaluated against the MONDIAL database.

**Pure XPath expressions:** pure XPath expressions (i.e. without variables) are interpreted as *existential queries* which return true if the result set is non-empty:

   ?- //country[name/text() = "Belgium"]//city/name/text().
   true

since the country element which has a name subelement with the text contents "Belgium" contains at least one city descendant with a name subelement with non-empty text contents.

**Output Result Set:** The query "?- *xpath*→N" for any XPath expression *xpath* binds $N$ to all nodes belonging to the result set of *xpath*:

   ?- //country[name/text() = "Belgium"]//city/name/text()→N.
   N/"Brussels"
   N/"Antwerp"
   ⋮

respectively, for a result set consisting of elements, logical ids are returned:

   ?- //country[name/text() = "Belgium"]//city→C.
   C/*brussels*
   C/*antwerp*
   ⋮

**Additional Variables:** XPathLog allows to bind all nodes which are traversed by an expression: The following expression returns all tuples $(N_1, C, N_2)$ such that the city with name $N_2$ belongs to the country with name $N_1$ and car code $C$:

   ?- //country[name/text()→N1 and @car_code→C]//city/name/text()→N2.
   N2/"Brussels" C/"B" N1/"Belgium"
   N2/"Antwerp" C/"B" N1/"Belgium"
      ⋮
   N2/"Berlin" C/"D" N1/"Germany"
      ⋮

**Dereferencing IDREF Attributes:** For every organization, give the name of city where the headquarter is located and all names and types of members:

   ?- //organization[name/text()→N and abbrev/text()→A and
                  @headq/name/text()→SN]
                  /members[@type→MT]/@country/name/text()→MN.

One element of the result set is, for example,
   N/"..."   A/"EU"   SN/"Brussels"   MT/"member"   MN/"Belgium"

**Schema Querying:** The use of variables at name positions allows for schema querying, e.g. to give all names of subelements of elements of type city:

```
?- //city/ SubElName .
SubElName/name
SubElName/population
    ⋮
```

**Navigation Variables:** Search for all things that have the name "Monaco". More explicitly, give the element type of all elements that have a name subelement with the text contents "Monaco":

```
?- // Type →X[name/text()→"Monaco"].
    Type/country X/country-monaco
    Type/city    X/city-monaco
```

Closed XPath-Logic formulas can e.g. be used for expressing integrity constraints.

*Example 3* (*Integrity Constraints*)

There are some application-specific integrity constraints on the MONDIAL database:

**Range restrictions:** The text contents of population elements and the value of area attributes must be a non-negative number:

$$\forall\ X: ((//population/text() \rightarrow X\ \text{or}\ //@area \rightarrow X) \rightsquigarrow X \geqslant 0).$$

The sum of percentages of ethnic groups in a country is at most 100%:

$$\forall\ C: (//country \rightarrow C \rightsquigarrow \text{sum}\{N\ [C];\ C/ethnicgroups/@percentage \rightarrow N\} \leqslant 100).$$

**Bidirectional relationships:** The membership of countries in organizations is represented bidirectionally:

$$\forall\ C,O: (//country \rightarrow C[@memberships \rightarrow O] \leftrightarrow$$
$$\exists\ T: //organization \rightarrow O/members[@type \rightarrow T\ \text{and}\ @country \rightarrow C]).$$

**Other conditions:** The country attribute of border subelements of country elements must reference a country which is encompassed by the same continent:

$$\forall\ C,C2: (//country \rightarrow C/border[@country \rightarrow C2] \rightsquigarrow$$
$$(//country \rightarrow C2\ \text{and}\ \exists\ Cont: (C/encompassed/@continent \rightarrow Cont\ \text{and}$$
$$C2/encompassed/@continent \rightarrow Cont))).$$

### 2.2 XML instances as semantical structures

Next, we need a basis for a model-theoretic semantics of XPath-Logic. The *information* that is carried by an XML instance is *abstractly* defined in the *XML Information Set* specification (XMLInf, 1999). It can be represented in different ways (e.g. as the *human-readable* ASCII-based notation, or by using the DOM (DOM-W3C, 1998) that provides an abstract datatype for *implementations*. There are approaches that regard XML trees as database items where the languages operate on (UnQL, Xcerpt). In our approach, the atomic items are the edges of a *graph* (than can be an XML tree, but that can also represent overlapping tree views on an internal graph-like database), called *XTreeGraph*. In contrast to the DOM model and the XML Query Data Model (XMQ-D, 2001) which use a node-labeled tree (i.e. the element and attribute names are associated with the nodes), the XTreeGraph is an *edge-labeled* model. Using an edge-labeled model proves useful for data manipulation and integration (see (May and Behrends, 2001)). Recall that XML-QL (Deutsch *et al.*,

1999) also uses an *edge-labeled* graph, which especially defines the same handling of text contents as ours; influenced by the experiences with the STRUDEL/STRUQL (Fernandez *et al.*, 1998) project for data integration.

Formally, the XTreeGraph is represented by an *X-Structure* (that interprets a signature consisting of element and attribute names, similar to a first order structure). The advantage of that approach is that it allows for *manipulating* an internal database by adding edges to the graph. Thus, XPathLog is not only a query language, but also a manipulation language. Its rule heads do not necessarily *construct* new XML trees/terms, but can update the X-Structure. As a prerequisite for mapping XML instances to X-Structures, some notation for handling lists is needed:

*Notation 1 (Lists)*
Throughout this work, the following usual notation is used:

- For two sets $A$ and $B$, the set of mappings from $A$ to $B$ is denoted by $B^A$.
- A list over a domain $D$ is a mapping from $\mathbb{N}$ to $D$. Thus, the set of lists over $D$ is denoted by $D^{\mathbb{N}}$.
- the empty list is denoted by $\varepsilon$; a unary list containing only the element $x$ is denoted by $(x)$; list concatenation as an operator is denoted by $\circ$.
- $\mathsf{set}(expr_1(x_1,\ldots,x_n) \mid expr_2(x_1,\ldots,x_n))$ stands for

$$\{expr_1(x_1,\ldots,x_n) \mid expr_2(x_1,\ldots,x_n)\}$$

  (i.e. the set of all $expr_1(x_1,\ldots,x_n)$ such that the condition $expr_2(x_1,\ldots,x_n)$ holds). In the following, sets are sometimes used as lists exploiting the fact that a set can be seen as a list by an arbitrary enumeration.
- In a similar way, a list can be constructed by enumerating its elements. For a list $\ell = (i_1, i_2, \ldots)$, $\mathsf{list}_{i \in \ell}(expr_1(i) \mid expr_2(i))$ is the list of all $expr_1(i_j)$ where $expr_2(i_j)$ holds. Similar to list, $\mathsf{concat}_{i \in I}(expr_1(i) \mid expr_2(i))$ does the same if $expr_1(i)$ is already a list.
- For a finite list $\ell = (x_1, \ldots, x_n)$, $\mathsf{reverse}(\ell) = (x_n, \ldots, x_1)$.
- For a list $\ell$, $\ell[i, j]$ denotes the sublist that consists of the $i$th to $j$th elements,
- For a list $\ell$ of pairs i.e. $\ell = ((x_1, y_1), (x_2, y_2), \ldots)$, $\ell \downarrow_1$ denotes the projection of the list to the first component of the list elements, i.e. $\ell \downarrow_1 := (x_1, x_2, \ldots)$.

*X-Structures.* When representing XML instances as X-Structures, (i) their elements/subelement structure, and (ii) the elements' attributes have to be represented. The universe consists of the *element nodes* of the XML instance and the *literals* used as attribute values and text contents. *Element nodes* have properties, defined by (i) subelements (which are ordered) and (ii) attributes (which are unordered). Multivalued attributes (`NMTOKENS` and `IDREFS`) are silently split, and reference attributes are silently resolved. Additionally, X-Structures support named constants and predicates as known from first-order logic.

Each XML instance is represented as a structure with a universe $\mathscr{U}$ over a signature $\Sigma = (\Sigma_N, \Sigma_F, \Sigma_C, \Sigma_P)$ which consists of

- $\Sigma_N$: element names and attribute names,
- $\Sigma_F$: names of XML-built-in functions,

- $\Sigma_C$: constant symbols, denoting elements in the XML instance (e.g. germany, interpreted as the element addressed by /mondial/country[name= "Germany"]).
- $\Sigma_P$: predicates (with arity).
- Additionally, a basic set of literal constants is assumed.

An X-Structure contains only the basic facts about the XML tree, i.e. the child and attribute relationships (similar to the DOM). Note that our approach which associates the order with the children of elements, differs from, for instance, the DOM and XML-QL approaches where a *global* order of all elements is assumed.

*Definition 2* (*X-Structure*)
An X-Structure over a given signature $\Sigma$ is a tuple $\mathscr{X} = (\mathscr{V}, \mathscr{L}, \mathscr{N}, \mathscr{I}, \mathscr{E}, \mathscr{A})$ where the universe $\mathscr{U}$ consists of three sets $\mathscr{V}$, $\mathscr{L}$, and $\mathscr{N}$: $\mathscr{V}$ is a set of nodes (from the graph point of view, vertices), identified by internal names, $\mathscr{L}$ is a set of literals (integers, floats, strings), $\mathscr{N}$ is the set of names (e.g. as occurring in node tests). Names may be further distinguished into $\mathscr{N}_{\mathscr{E}}$, containing the element names (and a special element text() for handling text children), and $\mathscr{N}_{\mathscr{A}}$ containing the attribute names.

- $\mathscr{I}$ is a (partial) mapping, which interprets the signature: $\mathscr{I}_{\mathscr{E}} : \Sigma_N \to \mathscr{N}_{\mathscr{E}}$ and $\mathscr{I}_{\mathscr{A}} : \Sigma_N \to \mathscr{N}_{\mathscr{A}}$ interpret the names in $\Sigma_N$ by element and attribute names. $\mathscr{I}_C : \Sigma_C \to \mathscr{V}$ interprets the constant symbols in $\Sigma_C$ by nodes in $\mathscr{V}$, and $\mathscr{I}_F : \mathscr{V} \times \Sigma_F \times (\mathscr{V} \cup \mathscr{L} \cup \mathscr{N})^* \to \mathscr{V} \cup \mathscr{L} \cup \mathscr{N}$ represents the interpretation of built-in functions (as defined in XPQOF (2001)). Finally, $\mathscr{I}_P : \Sigma_P \times (\mathscr{V} \cup \mathscr{L} \cup \mathscr{N})^* \to \{true, false\}$ represents the interpretation of predicates.
- $\mathscr{E}$ is a (partial) mapping $\mathscr{E} : \mathscr{V} \times \mathbb{N} \times \mathscr{N}_{\mathscr{E}} \to \mathscr{V} \cup \mathscr{L}$ (subelement relationship and text contents; from the graph point of view, an ordered set of edges).
- $\mathscr{A}$ is a (partial) mapping $\mathscr{A} : \mathscr{V} \times \mathscr{N}_{\mathscr{A}} \to 2^{\mathscr{V}} \cup 2^{\mathscr{L}}$ (attribute values). *XML attribute nodes* do not belong to $\mathscr{V}$, but their literal values belong to $\mathscr{L}$. For reference attributes (IDREF), the "results" are not the ID-strings, but the target nodes in $\mathscr{V}$ themselves.

Note that $\mathscr{E}$ and $\mathscr{A}$ are not direct interpretations of $\Sigma$, but mappings that "interprete" $\mathscr{N}$. $\Sigma$ is mapped to $\mathscr{N}$ before being interpreted by $\mathscr{I}$, making attribute and element names full citizens of the language (e.g. as in F-Logic).

There is a canonical mapping from the set of XML instances to the set of X-Structures. The canonical X-Structure to an XML instance is a single XML tree (cf. Figure 1), covering the DOM model.

*Example 4*
Figure 1 shows the X-Structure of the running example given in Example 1.

The elements of $\mathscr{V}$ (representing the element nodes) do not carry information in themselves, they are only of interest as anonymous entities (similar to object ids) which have certain properties that are given by $\mathscr{E}$, $\mathscr{A}$, and $\mathscr{I}$. In the following, mnemonic ids (e.g. *germany*) are used for elements of $\mathscr{V}$. Also, $\Sigma_N$ is identified with $\mathscr{N}_{\mathscr{E}}$ and $\mathscr{N}_{\mathscr{A}}$, omitting $\mathscr{I}_{\mathscr{E}}$ and $\mathscr{I}_{\mathscr{A}}$.

In full generality, an X-Structure can also contain subelement edges and reference edges which are not conforming with the XML tree model, but which are crucial
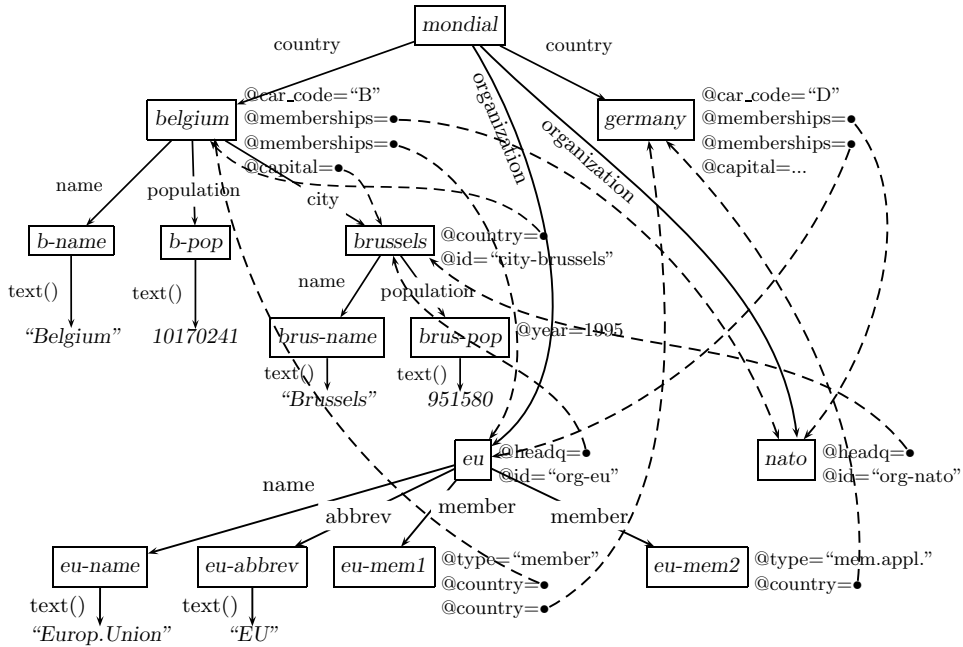
Fig. 1. Example X-Structure.

for data integration: An element may be a subelement of several other elements (as we show in section 4.2 and Figure 3), even with different names of the subelement relationship ("overlapping trees" – thus, the term *XTreeGraph* (May and Behrends, 2001) for the abstract data model).

In the following, X-Structures serve for defining a semantics for XPath-Logic, using the same terms as for XPath.

*Definition 3 (Basic Result Sets: Axes)*
For every node $x$ in an X-Structure $\mathscr{X}$ and every *axis* $a$ as defined in XPath,

$$\mathscr{A}_{\mathscr{X}}(a,x) \in ((\mathscr{V} \cup \mathscr{L}) \times \mathscr{N})^{\mathbb{N}}$$

is the list of pairs (*value, name*) generated by axis $a$ with $x$ as context node (do not confuse $\mathscr{A}_{\mathscr{X}}$ with $\mathscr{A}$ which denotes the interpretation of attributes in $\mathscr{X}$).

$$\mathscr{A}_{\mathscr{X}}(\text{child}, x) \quad := \quad \text{list}_{i \in \mathbb{N}}((y, name) \mid \mathscr{E}(x, i, name) = y)$$
$$\mathscr{A}_{\mathscr{X}}(\text{attribute}, x) \quad := \quad \text{list}((y, name) \mid y \in \mathscr{A}(x, name)) \quad \text{by some enumeration.}$$

For the other axes, $\mathscr{A}_{\mathscr{X}}(a, x)$ is derived according to the XPath specification:

$$\mathscr{A}_{\mathscr{X}}(\text{parent}, x) \quad := \quad \text{set}((p, \mathscr{I}(p, \text{name}, ()))) \mid x \in \mathscr{A}_{\mathscr{X}}(\text{child}, p) \downarrow_1)$$
$$\mathscr{A}_{\mathscr{X}}(\text{preceding-sibling}, x) \quad :=$$
$$\quad \text{concat}_{p \in \mathscr{A}_{\mathscr{X}}(\text{parent}, x) \downarrow_1}(\text{reverse}(\mathscr{A}_{\mathscr{X}}(\text{child}, p)[1, i{-}1]) \mid x = \mathscr{A}_{\mathscr{X}}(\text{child}, p)[i])$$
$$\mathscr{A}_{\mathscr{X}}(\text{following-sibling}, x) \quad :=$$
$$\quad \text{concat}_{p \in \mathscr{A}_{\mathscr{X}}(\text{parent}, x) \downarrow_1}(\mathscr{A}_{\mathscr{X}}(\text{child}, p)[i{+}1, \text{last}()] \mid x = \mathscr{A}_{\mathscr{X}}(\text{child}, p)[i])$$
$$\mathscr{A}_{\mathscr{X}}(\text{ancestor}, x) \quad := \quad \text{concat}_{(p,n) \in \mathscr{A}_{\mathscr{X}}(\text{parent}, x)}(((p, n)) \circ \mathscr{A}_{\mathscr{X}}(\text{ancestor}, p))$$
$$\mathscr{A}_{\mathscr{X}}(\text{descendant}, x) \quad := \quad \text{concat}_{(c,n) \in \mathscr{A}_{\mathscr{X}}(\text{child}, x)}((c, n) \circ \mathscr{A}_{\mathscr{X}}(\text{descendant}, c))$$

*Remark 2*

Recall that in the node-labeled XML/XPath data model, the semantics of expressions is always a list or a set of labeled nodes. In contrast, $\mathscr{A}_{\mathscr{X}}$ does not return a list of (labeled) nodes, but a list of pairs (*node/literal*, *name*), according to the edge-labeled data model underlying our approach.

## 2.3 Semantics

The semantics of XPath-Logic is defined similar to that of first-order logic by induction over the structure of expressions and formulas. The main task here is to define the semantics of reference expressions, handling navigation, order, and filtering. A reference expression simultaneously acts as a term (it has a result (list) and can be compared to terms) and as a predicate (when used in a step qualifier).

The basic result lists are provided by $\mathscr{A}_{\mathscr{X}}(axis, v)$ for every node $v$ of $\mathscr{X}$; recall that $\mathscr{A}_{\mathscr{X}}(\mathsf{attribute}, x)$ contains literals in case of non-reference attributes, and element nodes in case of reference attributes.

### 2.3.1 Semantics of expressions

As for first-order logic, a *variable assignment* $\beta : \mathsf{Var} \to \mathscr{U}$ maps variables to elements of the universe $\mathscr{U}$ (nodes, literals, and names) of the underlying X-Structure. For a variable assignment $\beta$, a variable $x$, and $d \in \mathscr{U}$, the *modified* variable assignment $\beta_x^d$ is identical with $\beta$ except that it assigns $d$ to the variable $x$:

$$\beta_x^d : \mathsf{Var} \to \mathscr{U} : \begin{cases} y \mapsto \beta(y) & \text{if } y \neq x, \\ x \mapsto d & \text{otherwise.} \end{cases}$$

For $\beta$ as above, and a variable $x$, $\beta \setminus \{x\}$ denotes $\beta$ without the mapping for $x$.

Expressions are decomposed into their *axis steps*. Every step consists of choosing an axis, preselecting nodes by a *node test*, and filtering the result by (i) "normal" predicates and (ii) XPath *context functions* (e.g. position() and last()) which use the order of the intermediate result list for selecting a certain element by its index.

*Definition 4* (*Semantics of XPath-Logic expressions*)
The semantics is defined by operators $\mathscr{S}$ and $\mathscr{Q}$ which are derived from the formal semantics given in Wadler (1999).

- $\mathscr{S}_{\mathscr{X}} : \text{Reference\_Expressions} \to (\mathscr{V} \cup \mathscr{L} \cup \mathscr{N})^{\mathbb{N}}$, and
    $(\text{Axes} \times \text{Reference\_Expressions} \times \mathscr{V} \times \text{Var\_Assignments}) \to (\mathscr{V} \cup \mathscr{L} \cup \mathscr{N})^{\mathbb{N}}$
  evaluates reference expressions wrt. an axis, a context node, and a variable assignment and returns a result list. In the second case, we use *any* to denote that the actual value of the node does not matter, and we use $\mathscr{S}^{any}$ to denote that the actual value of *axis* does not matter.
- $\mathscr{Q}_{\mathscr{X}} : (\text{Predicate\_Expressions} \times \mathscr{V} \times \text{Var\_Assignments}) \to \text{Boolean}$
  evaluates step qualifiers wrt. a context node and a variable assignment.

*Reference expressions are evaluated by $\mathscr{S}$:*

1. For closed expressions, $\mathscr{S}_{\mathscr{X}}(refExpr) = \mathscr{S}_{\mathscr{X}}^{any}(refExpr, any, \emptyset)$.

2. Reference expressions are translated into path expressions wrt. a start node:
   - rooted paths: $\mathscr{S}_{\mathscr{X}}^{any}(/p, any, \beta) = \mathscr{S}_{\mathscr{X}}^{any}(p, root, \beta)$ where *root* is as follows:
       * the unique root node if only one XML document is currently stored,
       * the root node that has been used in the outer expression, if $/p$ occurs in an expression of the form $path[/p]$.
   - rooted paths in other documents:
     $$\mathscr{S}_{\mathscr{X}}^{any}(\mathsf{document}(\text{``http://...''})/p, any, \beta) = \mathscr{S}_{\mathscr{X}}^{any}(p, root, \beta)$$
     where *root* is the root node of the document stored at http://....
   - entry points specified by a constant $c$: $\mathscr{S}_{\mathscr{X}}^{any}(c/p, any, \beta) = \mathscr{S}_{\mathscr{X}}^{any}(p, \mathscr{I}_C(c), \beta)$ (this is mainly of interest when multiple documents are used and constants are associated with their roots or some nodes, see section 4.2).
   - entry points specified by variables $v \in \mathsf{Var}$: $\mathscr{S}_{\mathscr{X}}^{any}(v/p, any, \beta) = \mathscr{S}_{\mathscr{X}}^{any}(p, \beta(v), \beta)$

3. Axis steps: $\mathscr{S}_{\mathscr{X}}^{any}(axis :: pattern, x, \beta) = \mathscr{S}_{\mathscr{X}}^{axis}(pattern, x, \beta)$
   where *pattern* is of the form *nodetest remainder* where *remainder* is a sequence of step qualifiers and variable bindings. These are evaluated left to right, always applying the rightmost "operation" (step qualifier or variable) to the result of the left part:

4. Node test: 
   $$\begin{aligned}
   \mathscr{S}_{\mathscr{X}}^{a}(name, x, \beta) &= \mathsf{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,x)}(v \mid n = name) \\
   \mathscr{S}_{\mathscr{X}}^{a}(\mathsf{node}(), x, \beta) &= \mathsf{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,x)}(v \mid v \in \mathscr{V}) \\
   \mathscr{S}_{\mathscr{X}}^{a}(\mathsf{text}(), x, \beta) &= \mathsf{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,x)}(v \mid v \in \mathscr{L}) \\
   \mathscr{S}_{\mathscr{X}}^{a}(N, x, \beta) &= \mathsf{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,x)}(v \mid n = \beta(N))
   \end{aligned}$$

5. Step with variable binding:
   $$\mathscr{S}_{\mathscr{X}}^{a}(pattern \to V, x, \beta) = \begin{cases} (\beta(V)) & \text{if } (\beta(V)) \in \mathscr{S}_{\mathscr{X}}^{a}(pattern, x, \beta) \\ \varepsilon & \text{otherwise.} \end{cases}$$

6. Step qualifiers:
   $$\mathscr{S}_{\mathscr{X}}^{a}(pattern[stepQ], x, \beta) = \mathsf{list}_{y \in \mathscr{S}_{\mathscr{X}}^{a}(pattern, x, \beta)}(y \mid \mathcal{Q}_{\mathscr{X}}(stepQ, y, \beta_{Pos, Size}^{k, n}))$$
   where $L_1 := \mathscr{S}_{\mathscr{X}}^{a}(pattern, x, \beta)$ and $n := size(L_1)$, and for every $y$, let $j$ the index of $y$ in $L_1$, $k := j$ if $a$ is a forward axis, and $k := n + 1 - j$ if $a$ is a backward axis (cf. (Wadler, 1999)). *Pos* and *Size* are only used if the step qualifier contains a context function.

7. Path: $\mathscr{S}_{\mathscr{X}}^{a}(p_1/p_2, x, \beta) = \mathsf{concat}_{y \in \mathscr{S}_{\mathscr{X}}^{a}(p_1, x, \beta)}(\mathscr{S}_{\mathscr{X}}^{any}(p_2, y, \beta))$

*Step Qualifiers are evaluated by $\mathcal{Q}$:*

8. Reference expressions have an existential semantics in step qualifiers:
   $$\mathcal{Q}_{\mathscr{X}}(refExpr, y, \beta) :\Leftrightarrow \mathscr{S}_{\mathscr{X}}^{any}(refExpr, y, \beta) \neq \emptyset$$

9. Predicates (including comparison predicates): The semantics of predicates in XPath is *element-oriented*: $p(refExpr, term)$ evaluates to *true* if at least one pair

taken from the result sets of *refExpr* and *term* satisfies the predicate $p$ (either defined in $\mathscr{I}_P$, or a built-in predicate of XPath (XPQOF, 2001)):

$$Q_{\mathscr{X}}(pred(expr_1,\ldots,expr_n),y,\beta) \quad :\Leftrightarrow$$
$$\text{there are } x_1 \in \mathscr{S}_{\mathscr{X}}^{any}(expr_1,y,\beta),\ldots,x_n \in \mathscr{S}_{\mathscr{X}}^{any}(expr_n,y,\beta)$$
$$\text{such that } (x_1,\ldots,x_n) \in \mathscr{I}_P(pred)$$

10. Boolean Connectives and Quantification are defined as usual.

*Evaluation of terms.*

11. Constants $c \in \Sigma_C$: $\mathscr{S}_{\mathscr{X}}^a(c,x,\beta) = \mathscr{I}_C(c)$. For literals, $\mathscr{S}_{\mathscr{X}}^a(lit,x,\beta) = lit$.
12. Variables: $\mathscr{S}_{\mathscr{X}}^a(var,x,\beta) = \beta(var)$.
13. Functions and arithmetics are also defined element-wise:
$$\mathscr{S}_{\mathscr{X}}^a(f(expr_1,\ldots,expr_n),x,\beta) =$$
$$\{\mathscr{I}_F(\beta(x),f,x_1,\ldots,x_n) \mid x_1 \in \mathscr{S}_{\mathscr{X}}^{any}(expr_1,x,\beta),\ldots,x_n \in \mathscr{S}_{\mathscr{X}}^{any}(expr_n,x,\beta)\}$$
14. Context-related functions use the extension of variable bindings by pseudo-variables $Size$ and $Pos$ in rule (6):
$$\mathscr{S}_{\mathscr{X}}^{any}(\mathsf{position}(),x,\beta) = \beta(Pos) \qquad \text{and} \qquad \mathscr{S}_{\mathscr{X}}^{any}(\mathsf{last}(),x,\beta) = \beta(Size) .$$

The following theorem states the equivalence of our semantics with that given in Wadler (1999), which is in turn equivalent to the one defined by the W3C for XPath in XQFS (2001).

*Theorem 1 (Correctness of $\mathscr{S}$ and $\mathcal{Q}$ wrt. $XPath$)*
For XPath reference expressions without splitting `NMTOKENS` attributes, the semantics coincides with the one given in Wadler (1999) (which already covers all core constructs of XPath as an addressing formalism): For every XPath expression *expr*,

$$\mathscr{S}_{\mathscr{X}}(expr) = \mathscr{S}[[expr]](x)$$

(for arbitrary $x$) where $\mathscr{S}[[expr]]$ is as defined in (Wadler, 1999).

Note that $\mathscr{S}[[expr]]$ defines only a result *set* that is implicitly ordered wrt. document order. Our semantics coincides with the document order as long as no dereferencing is used.

The proof uses the following Lemma which contains the structural induction (for proofs, see Appendix A).

*Lemma 2 (Correctness of $\mathscr{S}$ and $\mathcal{Q}$ wrt. $XPath$: Structural Induction)*
XPath-Logic reference expressions correspond to XPath as follows:

1. For absolute expressions (i.e. $expr = /expr'$, and no free variables):
$$\mathscr{S}_{\mathscr{X}}^{any}(expr,any,\emptyset) = \mathscr{S}[[expr]](x) \text{ for arbitrary } x.$$
2. For expressions, for all $\beta$: $\mathscr{S}_{\mathscr{X}}^{any}(expr,v,\beta) = \mathscr{S}[[expr]](v)$.
3. For step qualifiers, for all $\beta$: $Q_{\mathscr{X}}(stepQ,v,\beta_{Pos,Size}^{k,n}) \Leftrightarrow Q[[stepQ]](v,k,n)$.
4. For arithmetic expressions and built-in functions, for all $\beta$:
$$\mathscr{S}_{\mathscr{X}}^{any}\left(expr,v,\beta_{Pos,Size}^{k,n}\right) = \mathscr{E}[[expr]](v,k,n).$$

where $\mathcal{Q}[[expr]]$, $\mathcal{S}[[expr]]$, and $\mathcal{E}[[expr]]$ are as defined in Wadler (1999). Since XPath expressions are variable-free, $\beta$ is empty except handling the pseudo variables *Size* and *Pos* (which are often also empty).

The above behavior deviates from XPath for special kinds of attributes: When navigating along reference attributes, the result is not in document order, but in the same order as the referencing elements were. Additionally, NMTOKENS that are considered as atomic in XPath, are split in XPath-Logic.

### 2.3.2 Semantics of formulas

*Definition 5* (*Semantics of XPath-Logic Formulas*)
Formulas are interpreted according to the usual first-order semantics

$$\models\ \subseteq\ (\text{X-Structures} \times \text{Var\_Assignments} \times \text{Formulas})$$

15. Reference Expressions: The semantics of reference expressions corresponds to a *predicate* in first-order logic, defining a purely existential semantics:

$$(\mathcal{X}, \beta) \models refExpr \ :\Leftrightarrow\ (\mathcal{S}_{\mathcal{X}}(refExpr, \beta)) \neq \emptyset$$

16. predicates and boolean connectives: same as in first-order logic.

The above definitions associate a *truth value semantics* with XPath-Logic formulas. The $\models$ relationcan be used for expressing integrity constraints on XML documents (see Example 3) and even sets of documents, and for reasoning on X-Structures. In contrast, when defining XPathLog as a *data manipulation language* in the next section, a completely different formalization of the semantics is given: there, as for Datalog queries, the *answer substitutions* for a formula containing free variables have to be computed.

### 2.4 Annotated literals

The XML data model distinguishes between elements and their text contents. Nevertheless, in several situations, elements containing text contents are expected to act as numbers or strings:

*Example 5* (*Annotated Literals*)
Consider again the MONDIAL XML instance. The XPath queries
  //country[population > 5000000]/name/text()      and
  //country[population/text() > 5000000]/name/text()
are equivalent and return "Belgium" in their result set. In the first query, the *element* <population>10170241</population>  is implicitly casted into its literal value.

What happens here is not evident in XML/DTD environments. A corresponding XML Schema instance would show that the *complexType* population is derived from a simple type for integers. The idea here is that an element with a text contents adds structure to a simple type by allowing subelements and attributes. Thus, text

elements with attributes behave as *annotated literals*:

- comparisons, arithmetics, and (optionally) output use the literal value,
- navigation expressions use the element node, and
- in variable bindings, the variable is bound to the element, but it acts as described above when the variable is used, for instance, in a comparison.

## 3 XPathLog: The Horn fragment of XPath-Logic

Similar to the case of Datalog which is the function-free Horn fragment of first-order predicate logic, XPathLog is a logic programming language based on XPath-Logic. The evaluation of a query ?- $L_1, \ldots, L_n$ results in a set of variable bindings (of the free variables of the query) to elements of the universe. The semantics of XPathLog programs (i.e. the semantics of the evaluation of a set of XPathLog rules as a logic program) is then defined in section 4 by combining the answer semantics with the model-theoretic semantics defined in the preceding section.

*Definition 6 (XPathLog)*
Atoms are the basic components of XPathLog rules:

- an *XPathLog atom* is either an XPath-Logic *reference expression* which does not contain quantifiers or disjunction in step qualifiers, or a predicate expression over such expressions.
- an XPathLog atom is *definite* if it uses only the child, sibling, and attribute axes and the atom does not contain negation, disjunction, function applications, and context functions. These atoms are allowed in rule heads (see section 4.2). The excluded features would cause ambiguities what update is intended, e.g. "insert $x$ as a descendant" does not specify where the element should actually be inserted.

Similar to Datalog, an *XPathLog literal* is an atom or a negated atom and an *XPathLog query* is a list ?- $L_1, \ldots, L_n$ of literals (in general, containing free variables). An *XPathLog rule* is a formula of the form $A_1, \ldots, A_k \leftarrow L_1, \ldots, L_n$ where $L_i$ are literals and $A_i$ are definite atoms. $L_1, \ldots, L_n$ is the *body* of the rule, evaluated as a conjunction. $A_1, \ldots, A_k$ is the *head* of the rule, which may contain free variables that must also occur free in the body. In contrast to usual logic programming, we allow for lists of atoms in the rule head which are interpreted as conjunctions.

### 3.1 Queries in XPathLog

The semantics $\mathscr{SB}$ of XPathLog queries associates a result set and a set of *answer substitutions* with every XPathLog query by extending the above definition of $\mathscr{S}$. The semantics provides the formal base for the implementation of an *algebraic evaluation* of XPathLog queries in LoPiX (cf. section 5).

#### 3.1.1 Answers data model

Whereas in Datalog, the answer to a query ?- $L_1, \ldots, L_n$ is a set of variable bindings, the semantics of XPath-Logic reference expressions is defined wrt. an X-Structure

$\mathscr{X}$ as an *annotated result list*, i.e. the semantics of an expression is

(i) a result list (corresponding to the result list of the underlying XPath expression, i.e. without the additional variable bindings), and
(ii) with every element of the result list, a list of variable bindings is associated.

The result list (i) is the same as defined by $\mathscr{S}$ in Definition 4, equivalent to the one defined for XPath expressions in Wadler (1999), and by the W3C for XPath (XQFS, 2001). Whereas from the XPath point of view for "addressing" nodes, only the result list is relevant, XPathLog queries are mapped to a set of variable bindings based on the associated bindings lists.

*Example 6 (Semantics)*
First, the semantics is illustrated by an example. Let $\mathscr{X}$ be the XML structure given in Example 4, and
$$expr := //\mathsf{organization} \rightarrow \mathsf{O}$$
$$\qquad\qquad [\mathsf{member/@country[@car\_code} \rightarrow \mathsf{C\ and\ name/text()} \rightarrow \mathsf{N]}$$
$$\qquad\qquad ]/\mathsf{abbrev/text()} \rightarrow \mathsf{A}.$$
The underlying XPath expression is
$$//\mathsf{organization[member/@country[@car\_code\ and\ name/text()]]/abbrev/text()}\ .$$
with the result list ("UN","EU",...). With each of the results, a list of bindings for the variables O, C, N, and A is associated, yielding the annotated result list
$$\mathscr{SB}_{\mathscr{X}}(expr) = \mathsf{list}((\text{"UN"}, \{(\mathsf{O}/un, \mathsf{A}/\text{"UN"}, \mathsf{C}/\text{"AL"}, \mathsf{N}/\text{"Albania"}),$$
$$(\mathsf{O}/un, \mathsf{A}/\text{"UN"}, \mathsf{C}/\text{"GR"}, \mathsf{N}/\text{"Greece"}),$$
$$\vdots \qquad\qquad \}),$$
$$(\text{"EU"}, \{(\mathsf{O}/eu, \mathsf{A}/\text{"EU"}, \mathsf{C}/\text{"D"}, \mathsf{N}/\text{"Germany"}),$$
$$(\mathsf{O}/eu, \mathsf{A}/\text{"EU"}, \mathsf{C}/\text{"F"}, \mathsf{N}/\text{"France"}),$$
$$\vdots \qquad\qquad \}),$$
$$\vdots \qquad\qquad\qquad\qquad )$$

*Definition 7 (Semantics)*
The domain of *sets of* variable bindings for $V_1, \ldots, V_n$ (i.e. the domain of the second component of our semantics – i.e. the possible answer sets for a query whose free variables are $V_1, \ldots, V_n$) is

$$\mathrm{Var\_Bindings}_{V_1,\ldots,V_n} := \left(2^{((\mathscr{V} \cup \mathscr{L} \cup \mathscr{N})^n)}\right)^{\{V_1,\ldots,V_n\}}.$$

Thus, in the general case for a general set Var of variables where $n$ is unknown,

$$\mathrm{Var\_Bindings} := \bigcup_{n \in \mathbb{N}_0} \left(2^{((\mathscr{V} \cup \mathscr{L} \cup \mathscr{N})^n)}\right)^{(\mathsf{Var}^n)}$$

is the set of sets of variable assignments. For an empty set of variables, $\{true\}$ is the only element in Var_Bindings; in contrast, $\emptyset$ means that there is no variable binding which satisfies a given requirement. We use $\beta$ for denoting an individual variable binding, and $\xi \in \mathrm{Var\_Bindings}$ for denoting a set of variable bindings.

$$\mathrm{AnnotatedResults} := ((\mathscr{V} \cup \mathscr{L}) \times \mathrm{Var\_Bindings})$$

is the set of annotated results (i.e. an annotated result is a pair $(v, \xi)$ where $v$ is a node or a literal and $\xi$ is a set of variable bindings (for the set of variables occurring free in a certain formula)).

*Definition 8 (Operators on Annotated Result Lists)*
From an annotated result list $\theta$, the result list is obtained as $\mathrm{Res}(\theta)$:

$$\mathrm{Res} \ : \ \mathrm{AnnotatedResults}^{\mathbb{N}} \to (\mathscr{V} \cup \mathscr{L} \cup \mathscr{N})^{\mathbb{N}}$$
$$((x_1, \xi_1), \ldots, (x_n, \xi_n)) \mapsto (x_1, \ldots, x_n)$$

For an annotated result list $\theta$ and a given $x \in \mathrm{Res}(\theta)$ contained in the result list, the set of variable bindings associated with $x$ is obtained by $\mathrm{Bdgs}(\theta, x)$:

$$\mathrm{Bdgs} : \ \mathrm{AnnotatedResults}^{\mathbb{N}} \times (\mathscr{V} \cup \mathscr{L} \cup \mathscr{N}) \to \mathrm{Var\_Bindings}$$
$$(((x_1, \xi_1), \ldots, (x, \xi), \ldots, (x_n, \xi_n)), x) \mapsto \xi \ \ (\text{let } \mathrm{Bdgs}(\theta, x) = \emptyset \text{ if } x \notin \mathrm{Res}(\theta))$$

Note that the joins ($\bowtie$) used in this section are always purely relational joins that are applied to the bindings component.

*Example 7 (Semantics (Cont'd))*
Continuing Example 6, $\mathrm{Res}(\mathscr{SB}_{\mathscr{X}}(expr)) = (\text{"UN"}, \text{"EU"}, \ldots)$ is the result list of the underlying XPath expression, and

$$\mathrm{Bdgs}(\mathscr{SB}_{\mathscr{X}}(expr), \text{"EU"}) = \{(\mathrm{O}/eu, \ \mathrm{A}/\text{"EU"}, \ \mathrm{C}/\text{"D"}, \ \mathrm{N}/\text{"Germany"}),$$
$$(\mathrm{O}/eu, \ \mathrm{A}/\text{"EU"}, \ \mathrm{C}/\text{"F"}, \ \mathrm{N}/\text{"France"}),$$
$$\vdots \hspace{5cm} \}$$

yields the variable bindings that are associated with the result value "EU".

### 3.1.2 Safety

The semantics definition evaluates formulas and expressions wrt. a given set of variable bindings which, for example, results from evaluating other subexpressions of the same query. This approach allows for a more efficient evaluation of joins (*sideways information passing strategy*), and is especially needed for evaluating *negated* expressions (by defining negation as a relational "minus" operator). Negated expressions which contain free variables are intended to *exclude* some bindings from a given set of potential results. Thus, for variables occurring in the scope of a negation, the input answer set to the negation must already provide potential bindings. This leads to a *safety* requirement similar to Datalog.

*Definition 9 (Safe Queries)*
First, safety of variables is decided for each individual ocurrence. A variable occurrence $V$ is *safe* wrt. the query if at least one of the following holds:

- if the occurrence is in a literal $L$, and it is not inside the scope of a negation and not in a comparison predicate other than equality (e.g. $X < 3$ is unsafe).
- if the occurrence is in a literal $L_i$ inside a step qualifier *pattern*$[L_1$ and $\ldots$ and $L_n]$ and $V$ has a safe occurrence in *pattern* or in some $L_j$ such that $j < i$.
- if the occurrence is in a literal $L_i$ of the query $?–L_1$ and $\ldots$ and $L_n$, and $V$ has a safe occurrence in some $L_j$ such that $j < i$.

A query $?\text{-} \ \mathsf{L}_1, \ldots, \mathsf{L}_n$ is *safe* if all variable occurrences in the query are safe.

### 3.1.3 Semantics of expressions

In the following, the semantics of safe queries is defined. The basic (non-annotated) result lists are again provided by $\mathscr{A}_{\mathscr{X}}(axis, v)$ for every node $v$ of $\mathscr{X}$.

**Definition 10** (*Answer Semantics of XPath-Logic Expressions*)
The semantics is defined by operators $\mathscr{SB}$ and $\mathscr{QB}$ derived from $\mathscr{S}$ and $\mathscr{Q}$ as defined in Definition 4; the $\mathscr{B}$ stands for the extension with variable bindings:

- $\mathscr{SB}_{\mathscr{X}}$ : (Reference_Expressions) $\rightarrow$ AnnotatedResults$^{\mathbb{N}}$ , and
  (Axes $\times \mathscr{V} \times$ Reference_Expressions $\times$ Var_Bindings)
  $\rightarrow$ AnnotatedResults$^{\mathbb{N}}$

  evaluates reference expressions wrt. an axis, an (optional) context node and a given set of variable bindings and returns an annotated result list.

- $\mathscr{QB}_{\mathscr{X}}$ : (Predicate_Expressions $\times \mathscr{V} \times$ Var_Bindings) $\rightarrow$ Var_Bindings

  evaluates step qualifiers wrt. a context node to sets of variable bindings.

Expressions are evaluated by $\mathscr{SB}$:

1. If no input bindings are given, $\mathscr{SB}_{\mathscr{X}}(refExpr) = \mathscr{SB}_{\mathscr{X}}^{any}(refExpr, any, \emptyset)$
2. Reference expressions are translated into path expressions wrt. a start node:
   - entry points: rooted path: $\mathscr{SB}_{\mathscr{X}}^{any}(/p, any, Bdgs) = \mathscr{SB}_{\mathscr{X}}^{any}(p, root, Bdgs)$
     where *root* is the current root as in Definition 4 (2) for the same case.
   - entry points: constants $c \in \Sigma_C$: $\mathscr{SB}_{\mathscr{X}}^{any}(c/p, any, Bdgs) = \mathscr{SB}_{\mathscr{X}}^{any}(p, c, Bdgs)$
   - rooted paths in other documents:
     $\mathscr{SB}_{\mathscr{X}}^{any}(\text{document}(\text{``http://} \ldots \text{''})/p, any, Bdgs) = \mathscr{S}_{\mathscr{X}}^{any}(p, root, \beta)$
     where *root* is the root node of the document stored at http://. . . .
   - entry points: variables $V \in \text{Var}$:
     $$\mathscr{SB}_{\mathscr{X}}^{any}(V/p, any, Bdgs) = \text{concat}_{x \in \mathscr{V}_{\text{active}}}\left(\mathscr{SB}_{\mathscr{X}}^{any}(p, x, Bdgs \bowtie \{V/x\})\right)$$

     where $\mathscr{V}_{\text{active}}$ is the set of element nodes in the current database.
     Remark: Here, the input bindings are used for optimization: if every $\beta \in Bdgs$ provides already bindings for the variable $V$, the *sideways information passing strategy* directly effects the join $\{V/x\} \bowtie Bdgs$, restricting the possible values for $V$ which in fact results in

     $$\mathscr{SB}_{\mathscr{X}}^{any}(V/p, any, Bdgs) = \text{concat}_{\beta \in Bdgs, x = \beta(V)}\left(\mathscr{SB}_{\mathscr{X}}^{any}(p, x, Bdgs \bowtie \{V/x\})\right)$$

     Thus, the propagation of bindings is not only necessary for handling negation but also provides a relevant optimization for positive literals.
     Note that in the recursive call $\mathscr{SB}_{\mathscr{X}}^{any}(p, x, Bdgs \bowtie \{V/x\})$, the propagated bindings are already augmented with the bindings for $V$.

3. Axis step: $\mathscr{SB}_{\mathscr{X}}^{any}(axis :: pattern, x, Bdgs) = \mathscr{SB}_{\mathscr{X}}^{axis}(pattern, x, Bdgs)$
   where *pattern* is of the form *nodetest remainder* where *remainder* is a sequence of step qualifiers and variable bindings. These are evaluated from left to right, always applying the rightmost "operation" (qualifier or variable) to the result of the left part:

4. Node test:

$$\mathscr{SB}^a_{\mathcal{X}}(name, x, Bdgs) = \mathsf{list}_{(v,n)\in\mathscr{A}_{\mathcal{X}}(a,x),\ n=name}(v, \{true\} \bowtie Bdgs)$$
$$\mathscr{SB}^a_{\mathcal{X}}(\mathsf{node}(), x, Bdgs) = \mathsf{list}_{(v,n)\in\mathscr{A}_{\mathcal{X}}(a,x),\ v\in\mathscr{V}}(v, \{true\} \bowtie Bdgs)$$
$$\mathscr{SB}^a_{\mathcal{X}}(\mathsf{text}(), x, Bdgs) = \mathsf{list}_{(v,n)\in\mathscr{A}_{\mathcal{X}}(a,x),\ v\in\mathscr{L}}(v, \{true\} \bowtie Bdgs)$$
$$\mathscr{SB}^a_{\mathcal{X}}(N, x, Bdgs) = \mathsf{list}_{(v,n)\in\mathscr{A}_{\mathcal{X}}(a,x)}(v, \{N/v\} \bowtie Bdgs)$$

5. Step with variable binding:

$$\mathscr{SB}^a_{\mathcal{X}}(pattern \to V, x, Bdgs) = \mathsf{list}_{(y,\xi)\in\mathscr{SB}^a_{\mathcal{X}}(pattern,x,Bdgs)}(y, \xi \bowtie \{V/y\})$$

6. Step qualifiers:

$$\mathscr{SB}^a_{\mathcal{X}}(pattern[stepQ], x, Bdgs) =$$
$$\mathsf{list}_{(y,\xi)\in\mathscr{SB}^a_{\mathcal{X}}(pattern,x,Bdgs),\ \mathcal{QB}_{\mathcal{X}}(stepQ,y,\xi')\neq\emptyset}(y,\ \mathcal{QB}_{\mathcal{X}}(stepQ, y, \xi') \setminus \{Pos, Size\})$$

If the step qualifier does not contain *context functions*, then $\xi' := \xi$, otherwise let $L := \mathscr{SB}^a_{\mathcal{X}}(pattern, x, Bdgs)$, and then for every $(y, \xi)$ in $L$, $\xi'$ is obtained as follows, extending $\xi$ with bindings of the pseudo variables *Size* and *Pos*:

- start with $\xi' = \emptyset$,
- for every $\beta \in \xi$, the list $L' = \mathsf{list}_{(y,\xi)\in L\ \text{s.t.}\ \beta\in\xi}(y)$ contains all nodes which are selected *for the variable assignment $\beta$*.
- let $size := size(L')$, and for every $y$, let $j$ the index of $x_1$ in $L'$, $pos := j$ if $a$ is a forward axis, and $pos := size+1-j$ if $a$ is a backward axis.
- add $\beta^{Size,Pos}_{size,pos}$ to $\xi'$.

7. Path: $\mathscr{SB}^a_{\mathcal{X}}(p_1/p_2, x, Bdgs) = \mathsf{concat}_{(y,\xi)\in\mathscr{SB}^{any}_{\mathcal{X}}(p_1,x,Bdgs)}\mathscr{SB}^a_{\mathcal{X}}(p_2, y, \xi)$

*Step Qualifiers are evaluated by $\mathcal{QB}$:*

8. Reference expressions (existential semantics) in step qualifiers:

$$\mathcal{QB}_{\mathcal{X}}(refExpr, x, Bdgs) = \bigcup_{(y,\xi)\in\mathscr{SB}^{any}_{\mathcal{X}}(refExpr,x,Bdgs)} \xi$$

9. The built-in equality predicate "$=$" is not only a comparison if both sides are bound, but also serves as an assignment if the left-hand side is a variable $V \in \mathsf{Var}$ which is *not* bound in *Bdgs*:

$$\mathcal{QB}_{\mathcal{X}}(V = expr, x, Bdgs) = \bigcup_{(y,\xi)\in\mathscr{SB}^{any}_{\mathcal{X}}(expr,x,Bdgs)} \xi \bowtie \{V/y\}$$

All other built-in comparisons require all variables to be bound:

$$\mathcal{QB}_{\mathcal{X}}(expr_1\ \textit{op}\ expr_2, x, Bdgs) = \bigcup_{(x_i,\xi_i)\in\mathscr{SB}_{\mathcal{X}}(expr_i,x,Bdgs),\ x_1\ \textsf{op}\ x_2} \xi_1 \bowtie \xi_2.$$

10. Predicates except built-in comparisons:

$$\mathcal{QB}_{\mathcal{X}}(pred(expr_1,\dots,expr_n), x, Bdgs) = \bigcup_{\substack{(x_i,\xi_i)\in\mathscr{SB}_{\mathcal{X}}(expr_i,x,Bdgs) \\ (x_1,\dots,x_n)\in\mathscr{I}(pred)}} \xi_1 \bowtie \dots \bowtie \xi_n$$

11. Negated expressions which do *not contain any free variable*:

$$\mathcal{QB}_{\mathcal{X}}(\text{not } A, x, Bdgs) = \begin{cases} Bdgs & \text{if } \mathcal{QB}_{\mathcal{X}}(A, x, \emptyset) = \emptyset, \\ \emptyset & \text{otherwise, i.e. if } \mathcal{QB}_{\mathcal{X}}(A, x, \emptyset) = \{true\}. \end{cases}$$

12. For negated expressions which contain free variables, negation is interpreted as the "minus" operator (as known e.g. from the relational algebra) wrt. the given input bindings. Thus, all variables which occur free in $A$ must be safe, i.e. every input variable binding has to provide a value for them.

For two variable bindings $\beta_1, \beta_2$, we write $\beta_1 \leqslant \beta_2$ if all variable bindings in $\beta_1$ occur also in $\beta_2$. Intuitively, in this case, if $\beta_1$ is "abandoned", $\beta_2$ should also be abandoned.

$$\mathcal{QB}_{\mathcal{X}}(\text{not } expr, x, Bdgs) = $$
$$Bdgs - \{\beta \in Bdgs \mid \text{ there is a } \beta' \in \mathcal{QB}_{\mathcal{X}}(expr, x, Bdgs) \text{ s.t. } \beta \leqslant \beta'\}$$

13. Conjunction:

$$\mathcal{QB}_{\mathcal{X}}(expr_1 \text{ and } expr_2, x, Bdgs) = $$
$$\mathcal{QB}_{\mathcal{X}}(expr_1, x, Bdgs) \bowtie \mathcal{QB}_{\mathcal{X}}(expr_2, x, \mathcal{QB}_{\mathcal{X}}(expr_1, x, Bdgs))$$

Here, in case of negated conjuncts in the step qualifier, the safety of variables has to be considered. The above definition assumes that by a left-to-right evaluation of conjuncts, the evaluation is safe.

*Evaluation of terms*

14. Constants: for literals, $\mathcal{SB}_{\mathcal{X}}^{any}(lit, x, Bdgs) = (lit, Bdgs)$. For constants $c \in \Sigma_C$, $\mathcal{SB}_{\mathcal{X}}^{any}(c, x, Bdgs) = (\mathcal{I}_C(c), Bdgs)$.

15. Variables: the variable occurrence must be safe, then: $\mathcal{SB}_{\mathcal{X}}^{any}(var, x, Bdgs) = \text{list}_{\beta \in Bdgs}(\beta(var), \beta)$.

16. Function terms and arithmetics:

$$\mathcal{SB}_{\mathcal{X}}^{any}(f(arg_1, \ldots, arg_n)), x, Bdgs) = $$
$$\text{list}_{(x_i, \xi_i) \in \mathcal{SB}_{\mathcal{X}}^{any}(arg_1, x, Bdgs), \ldots,}(f(x_1, \ldots, x_n), \xi_1 \bowtie \ldots \bowtie \xi_n)$$

where $f(x_1, \ldots, x_n)$ results from the built-in evaluation of $f$.

17. Context-related functions use the extension of variable bindings by pseudo-variables $Size$ and $Pos$ in rule (6):

$$\mathcal{SB}_{\mathcal{X}}^{any}(\text{position}(), x, Bdgs) = \text{list}_{\beta \in Bdgs}(\beta(Pos), \{\beta' \in Bdgs \mid \beta(Pos) = \beta'(Pos)\})$$
$$\mathcal{SB}_{\mathcal{X}}^{any}(\text{last}(), x, Bdgs) = \text{list}_{\beta \in Bdgs}(\beta(Pos), \{\beta' \in Bdgs \mid \beta(Size) = \beta'(Size)\})$$

The above semantics is an algebraic characterization of the logical semantics of XPath-Logic expressions which has been defined in section 4:

**Theorem 3** (*Correctness of $\mathcal{SB}$ and $\mathcal{QB}$*)
For every (in general, containing free variables) XPathLog expression *expr*,

$$\text{Res}(\mathcal{SB}_{\mathcal{X}}(expr)) = \bigcup_{\beta \in (\mathcal{V} \cup \mathcal{L} \cup \mathcal{N})^{\text{free}(expr)}} \mathcal{S}_{\mathcal{X}}(expr, \beta).$$

More detailed, for all $x \in \mathcal{V} \cup \mathcal{L} \cup \mathcal{N}$,

$$(x \in \mathrm{Res}(\mathscr{SB}_{\mathcal{X}}(expr)) \text{ and } \beta \in \mathrm{Bdgs}(\mathscr{SB}_{\mathcal{X}}(expr), x)) \;\Leftrightarrow\; x \in \mathscr{S}_{\mathcal{X}}(expr, \beta).$$

Again, the theorem uses a lemma which encapsulates the structural induction.

*Lemma 4* (*Correctness of $\mathscr{SB}$ and $\mathcal{QB}$: Structural Induction*)
The correctness of the answers semantics of XPathLog expressions mirrors the generation of answer sets by the evaluation: The input set *Bdgs* may contain bindings for the free variables of an expression. If for some variable *var*, no binding is given, the result extends *Bdgs* with bindings of *var*. If bindings are given for *var*, this specifies a constraint on the answers to be returned (expressed by joins).

- For every absolute expression *expr*, (i.e. $expr = /expr'$) and every set *Bdgs* of variable bindings,

$$(x \in \mathrm{Res}(\mathscr{SB}_{\mathcal{X}}(expr, Bdgs)) \text{ and } \beta \in \mathrm{Bdgs}(\mathscr{SB}_{\mathcal{X}}(expr, Bdgs), x)) \;\Leftrightarrow\;$$
$$(x \in \mathscr{S}_{\mathcal{X}}(expr, \beta) \text{ and } \beta \text{ completes some } \beta' \in Bdgs \text{ with } \mathsf{free}(expr)).$$

- For every expression *expr*, node *v*, and every set *Bdgs* of variable bindings,

$$(x \in \mathrm{Res}(\mathscr{SB}_{\mathcal{X}}(expr, v, Bdgs)) \text{ and } \beta \in \mathrm{Bdgs}(\mathscr{SB}_{\mathcal{X}}(expr, v, Bdgs), x)) \;\Leftrightarrow\;$$
$$(x \in \mathscr{S}_{\mathcal{X}}(expr, v, \beta) \text{ and } \beta \text{ completes some } \beta' \in Bdgs \text{ with } \mathsf{free}(expr)).$$

- for every step qualifier *stepQ*, node *v*, and every set *Bdgs* of variable bindings,

$$\beta \in \mathcal{QB}_{\mathcal{X}}(stepQ, v, Bdgs) \;\Leftrightarrow\;$$
$$\mathcal{Q}_{\mathcal{X}}(stepQ, v, \beta) \text{ and } \beta \text{ completes some } \beta' \in Bdgs \text{ with } \mathsf{free}(stepQ)).$$

The proof can be found in Appendix A.

### 3.1.4 Semantics of queries

According to Definition 6, XPathLog queries are conjunctions of XPathLog literals. In the following, the evaluation of safe queries is defined. The definition of safety guarantees that a left-to-right evaluation of the body is well-defined (i.e. all variable evaluations in Definition 10(15) are safe). Definition 10(13) already applied left-to-right propagation when evaluating step qualifiers.

*Definition 11*
The evaluation $\mathcal{QB}$ is extended to atoms by

$$\mathcal{QB}_{\mathcal{X}} \;:\; (\mathrm{Atoms} \times \mathrm{Var\_Bindings}) \to \mathrm{Var\_Bindings}$$
$$(A, Bdgs) \mapsto \mathcal{QB}_{\mathcal{X}}(A, root, Bdgs)$$

For safe queries ?- *atom* consisting of only one atom, $\mathcal{QB}$ yields the answer bindings:

$$(\mathcal{X}, \beta) \models atom \;\;:\Leftrightarrow\; \beta \in \mathcal{QB}_{\mathcal{X}}(atom, \emptyset)).$$

*Definition 12* (*Evaluation of Negated Literals*)
The evaluation of negated literals *L* is defined wrt. a set of input bindings which must cover the free variables in *L*, similar to negation in step qualifiers in Definition 10(12):

$$\mathcal{QB}_{\mathcal{X}}(\mathsf{not}\ A, Bdgs) :=$$
$$Bdgs - \{\beta \in Bdgs \mid \text{ there is a } \beta' \in \mathcal{QB}_{\mathcal{X}}(A, Bdgs) \text{ s.t. } \beta \leqslant \beta'\}.$$

*Definition 13* (*Evaluation of Queries*)

The evaluation of a safe query ?- $L_1, \ldots, L_n$ is defined similar to the evaluation of conjunctive step qualifiers in Definition 10(13):

$$
\begin{aligned}
\mathcal{QB}_{\mathcal{X}} \; &: \text{Conj\_Literals} \to \text{Var\_Bindings} \\
\mathcal{QB}_{\mathcal{X}}(L_1 \wedge \ldots \wedge L_i) \; &:= \; \mathcal{QB}_{\mathcal{X}}(L_1 \wedge \ldots \wedge L_{i-1}) \bowtie \\
& \qquad \mathcal{QB}_{\mathcal{X}}(L_i, (\mathcal{QB}_{\mathcal{X}}(L_1 \wedge \ldots \wedge L_{i-1}))|_{\mathsf{free}(L_i)})
\end{aligned}
$$

Given an X-Structure $\mathcal{X}$, the answer to a query ?- $\mathsf{L}_1, \ldots, \mathsf{L}_n$ is the set

$$
\mathsf{answers}_{\mathcal{X}}(\mathsf{L}_1, \ldots, \mathsf{L}_n) \; := \; \mathcal{QB}_{\mathcal{X}}(\mathsf{L}_1 \wedge \ldots \wedge \mathsf{L}_n) \quad \text{of variable bindings.}
$$

*Theorem 5* (*Correctness: Evaluation of Queries*)

For all safe XPathLog queries $Q$, $\quad \beta \in \mathcal{QB}_{\mathcal{X}}(Q) \; \Leftrightarrow \; (\mathcal{X}, \beta) \models Q$.

Note that the semantics of formulas is not based on a Herbrand structure consisting of ground atoms (as "usual" Herbrand semantics are), but directly on the interpretations $\mathscr{A}_{\mathcal{X}}$ of the axes in the X-Structure, and on an interpretation of predicate symbols that can be represented as a finite set of tuples over $\mathscr{V} \cup \mathscr{L} \cup \mathscr{N}$.

## 4 XPathLog programs

In logic programming, rules are used for a declarative specification: if the body of a clause evaluates to *true* for some assignment of its variables, the truth of the head atom for the same variable assignment can be inferred. Depending on the intention, this semantics can be used for (top-down) checking if something is derivable from a given set of facts, or (bottom-up) extending a given set of facts by additional, derived knowledge. In this work, we mainly investigate the bottom-up strategy, regarding XPathLog as an *update language for XML databases*: the evaluation of the body wrt. a given structure yields variable bindings which are propagated to the rule head where facts are added to the model.

Positive XPathLog programs (i.e. the rules contain only positive literals; also step qualifiers may only contain positive expressions) are evaluated bottom-up by a $T_P$-like operator over the X-Structure, providing a minimal model semantics. The formal definition of a $T_P$ operator will be given in Definition 16 for XPathLog programs after explaining the semantics of insertions and updates.

### 4.1 Atomization

In this section, an alternative semantics of conjunctions of *definite XPathLog atoms* is defined which provides the base for the *constructive* semantics of reference expressions in *rule heads*. The semantics is defined by resolving reference expressions syntactically into their constituting atomic steps in the same way as in F-Logic (Frohn *et al.*, 1994). A similar strategy for resolving expressions into atomic steps is followed by several approaches which store XML data in relational databases (Deutsch *et al.*, 2000; Shanmugasundaram *et al.*; Florescu and Kossmann, 1999), by flattening the XML instance to one or more universal relations.

*Definition 14* (*Atomization of Formulas*)

The function    atomize : XPathLogAtoms → $2^{\text{XPathLogAtoms}}$ resolves a definite XPathLog atom into atoms of the form *node*[*axis*::*nodetest*→*result*] and predicates over variables and constants. It will be used in Definition 16 for specifying the semantics of rule heads. atomize is defined by structural induction corresponding to the induction steps when defining $\mathscr{S}_{\mathscr{X}}$. In the following, *path* stands for a *path expression* (or a variable), and *name* for a name (or a variable).

- the entry case:  atomize(/*remainder*) := atomize(root/*remainder*)
- Paths are resolved into steps and step qualifiers are isolated (since context functions are not allowed in definite atoms, it can be assumed that there is at most one step qualifier, optionally preceded by a variable assignment):

$$\begin{aligned}
&\text{atomize}(path/axis :: nodetest \rightarrow var[stepQualifier] \ /remainder) := \\
&\qquad \text{atomize}(path[axis :: nodetest \rightarrow var]) \cup \\
&\qquad \text{atomize}(var[stepQualifier]) \cup \text{atomize}(var/remainder), \\
&\text{atomize}(path/axis :: nodetest[stepQualifier] \ /remainder) := \\
&\qquad \text{atomize}(path[axis :: nodetest \rightarrow \_X]) \cup \\
&\qquad \text{atomize}(\_X[stepQualifier]) \cup \text{atomize}(\_X/remainder)
\end{aligned}$$

  where $\_X$ is a new don't care variable.

- Conjunctions in step qualifiers are separated:

$$\begin{aligned}
&\text{atomize}(var[pred_1 \text{and} \ldots \text{and} \ pred_n]) := \\
&\qquad \text{atomize}(var[pred_1]) \cup \ldots \cup \text{atomize}(var[pred_n])
\end{aligned}$$

- Predicates in step qualifiers:

$$\begin{aligned}
\text{atomize}(var[pred(expr_1,\ldots,expr_n)]) := {}&\text{atomize}(equality(var,expr_1,\_X_1)) \cup \ldots \\
&\text{atomize}(equality(var,expr_n,\_X_n)) \cup \\
&\{pred(\_X_1,\ldots,\_X_n)\}
\end{aligned}$$

  where *equality*(*var*,*expr*,*X*) is defined as follows (if $expr_i$ is a constant, it is not replaced by a variable):

  * *equality*(*var*,*expr*,*X*) = "*expr* → *X*" if *expr* is of the form //*remainder*,
  * *equality*(*var*,*expr*,*X*) = "*var*/*expr* → *X*" if *expr* is of the form *axis* :: *nodetest remainder*.

- Predicate atoms are handled in the same way. Note that here all arguments are absolute expressions (rooted, or starting at a constant, or at a variable).

*Example 8* (*Atomization*)

?- //organization→O[name/text()→ON and
                @headq = members/@country[name/text()→CN]/@capital].

is atomized into

?- root[descendant::organization→O],    O[name→_ON], _ON[text()→ON],
   O[@headq→_S], O[members→_M], _M[@country→_C], _C[@country→_Cap],
   _S = _Cap, _C[child::name→_CN], _CN[text()→CN].

**Theorem 6** (*Correctness of* atomize)
The above semantics is equivalent to the one presented in Definition 10 for all
definite XPathLog atoms $A$ and every X-Structure $\mathscr{X}$, i.e.

$$\mathsf{answers}_{\mathscr{X}}(A) = \mathsf{answers}_{\mathscr{X}}(\mathsf{atomize}(A))$$

Again, the theorem uses a lemma which encapsulates the structural induction, using
the logical semantics for showing the correctness of atomize.

**Lemma 7** (*Correctness of* atomize: *Structural Induction*)
For every X-Structure $\mathscr{X}$ and every definite XPath-Logic atom $A$,

- for every variable assignment $\beta$ of $\mathsf{free}(A)$ such that $(\mathscr{X}, \beta) \models A$, there
  exists a variable assignment $\beta' \supseteq \beta$ of $\mathsf{free}(\mathsf{atomize}(A))$ such that $(\mathscr{X}, \beta') \models$
  $\mathsf{atomize}(A)$, and
- for every variable assignment $\beta'$ of $\mathsf{free}(\mathsf{atomize}(A))$ such that $(\mathscr{X}, \beta') \models$
  $\mathsf{atomize}(A)$, $(\mathscr{X}, \beta'|_{\mathsf{free}(A)}) \models A$.

The proof can be found in Appendix A.

### 4.2 Left hand side

Using *logical expressions* for specifying an update is perhaps the most important
difference to approaches like XSLT, XML-QL, or XQuery where the structure to be
*generated* is always specified by XML patterns, or to the update proposal for XML
described in Tatarinov *et al.* (2001). In contrast, in XPathLog, existing nodes are
communicated via variables to the head, where they are *modified* when appearing
at host position of atoms. The semantics of the left hand side of XPathLog rules –
which is a list of *definite* XPathLog atoms – is now investigated based on the
atomization of expressions. When used in the head, the "/" operator and the
"[...]" construct specify which properties should be added (thus, "[...]" does not
act as a step qualifier, but as a *constructor*). When using the child or attribute axis
for updates, the host of the expression gives the element to be updated or extended; 
when a sibling axis is used, effectively the parent of the host is extended with a new
subelement.

Note that the (pure) XPathLog language does *not* allow to delete or replace
existing elements or attributes[3] – modifications are always monotonic in the sense
that existing "things" remain.

*Generation or extension of attributes.* A ground-instantiated atom of the form $n[@$
$a{\rightarrow}v]$ specifies that the attribute $@a$ of the node $n$ should be set or extended with
$v$. If $v$ is not a literal value but a node, a reference to $v$ is stored.

**Example 9** (*Adding Attributes*)
We add the data code to Switzerland, and make it a member of the European

---

[3] suitable extensions, e.g. of the form delete(*elem,prop,val*) can be defined. Such extensions which would
turn XPathLog into a rule-based imperative language are not investigated in this work.
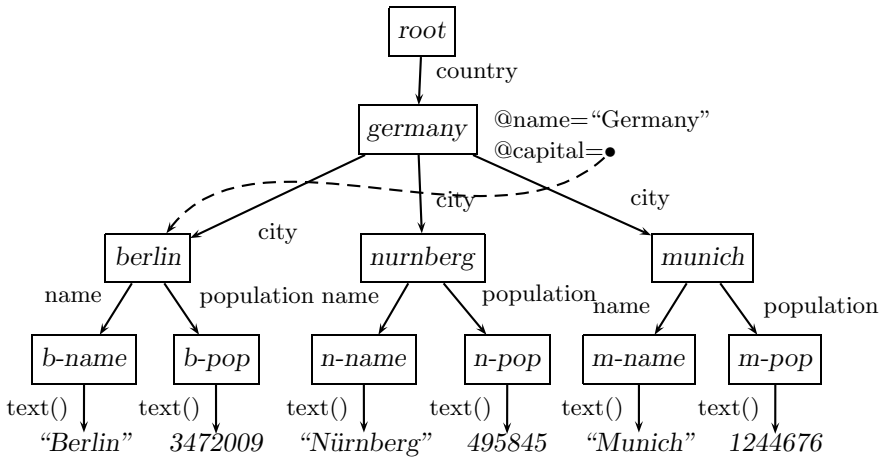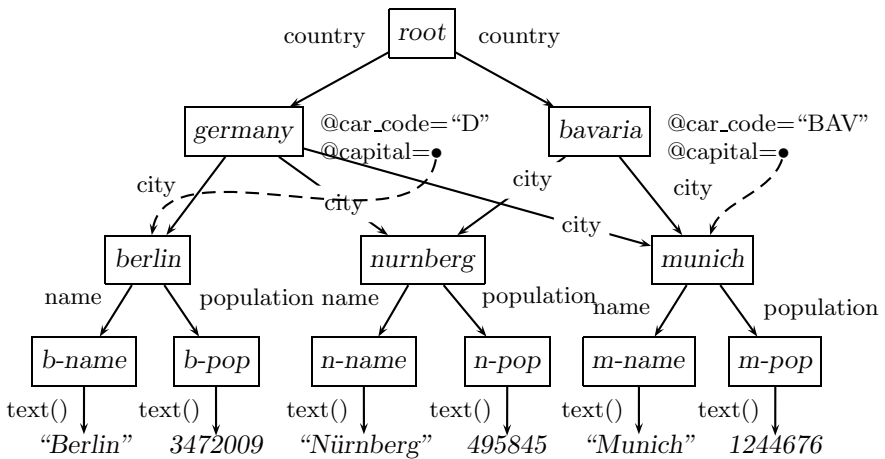
Fig. 2. Linking – before.



Fig. 3. Linking – after.

Union:

  C[@datacode→"ch"], C[@memberships→O] :-
      //country→C[@car_code="CH"], //organization→O[abbrev/text()→"EU"].

results in

‹country | datacode="ch" | car_code="CH" industry="machinery chemicals watches"
          memberships="org-efta org-un | org-eu | ... ">     ...   ‹/country›

*Creation of elements.* Elements can be created as *free* elements by atoms of the form
/*name*[...] (meaning "some element of type *name*" – this is interpreted to create an
element which is not a subelement of any other element), or as subelements.

*Example 10 (Creating Elements)*
We create a new (free) country element with some properties (cf. Figures 2 and 3):

  /country[@car_code→"BAV" and @capital→X and city→X and city→Y] :-
      //city→X[name/text()="Munich"],   //city→Y[name/text()="Nurnberg"].

The two city elements are *linked* as subelements. This operation has no equivalent in the "classical" XML model: these elements are now children of *two* country elements. Thus, changing the elements effects both trees. *Linking* is a crucial feature for efficient restructuring and integration of data (May and Behrends, 2001).

*Insertion of subelements.* Existing elements can be assigned as subelements to other elements: a ground instantiated atom $n$[child :: $s \rightarrow m$] makes $m$ a subelement of type $s$ of $n$. In this case, $m$ is *linked* as $n/s$ at the end of $n$'s children list.

*Example 11 (Inserting Subelements)*
The following two rules are equivalent to the above ones:
  /country[@car_code→"BAV"].
  C[@capital→X and city→X and city→Y] :- //country→C[@car_code="BAV"],
     //city→X[name/text()→"Munich"],   //city→Y[name/text()→"Nurnberg"].

Here, the first rule creates a free element, whereas the second rule uses the variable binding of $C$ to this element for inserting subelements and attributes.

In the above case, the position of the new subelement is not specified. If the atom is of the form $h$[child($i$)::$s$→$v$] or $h$[following/preceding-sibling($j$)::$s$→$v$], this means that the new element to be inserted should be made the $i$th subelement of $h$ or $j$th following/preceding sibling of $h$, respectively.

*Generation of elements by path expressions.* Additionally, subelements can be created by *reference expressions* in the rule head which create nested elements that satisfy the given reference expression. The atomization introduces local variables that occur *only in the head* of the rule. Here, we follow the semantics of PathLog (Frohn *et al.*, 1994) which is implemented in FLORID (Ludäscher *et al.*, 1998) for object creation. After the atomization, the resulting atoms are processed in an order such that the local variables are bound to the nodes/objects which are generated.

*Example 12 (Inserting Text Children)*
Bavaria gets a text subelement name:
  C/name[text()→"Bavaria"] :- //country→C[@car_code="BAV"].
Here, the atomized version of the rule is
  C[name→_N], _N[text()→"Bavaria"] :-
        root[descendant::country→C], C[@car_code="BAV"].
The body produces the variable binding  C/*bavaria*. When the head is evaluated, first, the fact *bavaria*[child::name→$x_1$] is inserted, adding an (empty) name subelement $x_1$ to *bavaria* and binding the local variable _N to $x_1$. Then, the second atom is evaluated, generating the text contents to $x_1$.

*Once-for-each-binding.* In contrast to classical logic programming where it does not matter if a fact is "inserted" into the database several times (e.g. once in every $T_P$ round), here subelements must be created exactly once for each instantiation of a rule. We define a revised $T_P$-operator in Definition 16.

*Using Navigation Variables for Restructuring.* For data restructuring and integration, the intuitiveness and declarativeness of a language gains much from variables ranging not only over data, but also over schema concepts (e.g. as in SchemaSQL (Lakshmanan *et al.*, 1996)). Such features have already been used for HTML-based Web data integration with F-Logic (Ludäscher *et al.*, 1998).

Extending the XPath wildcard concept, XPathLog also allows to have variables at name position. Thus, it allows for schema querying, and also for generating new structures dependent on the data contents of the original one.

*Example 13 (Restructuring, Name Variables)*
Consider a data source which provides data about waters according to the DTD

&lt;!ELEMENT terra (water+,...)&gt;
&lt;!ELEMENT water (...)&gt;   &lt;!ATTLIST water name CDATA #REQUIRED . . . &gt;

which contains, e.g. the following elements:

&lt;water type="river" name="Mississippi"&gt; ... &lt;/water&gt;
&lt;water type="sea" name="North Sea"&gt; ... &lt;/water&gt;.

This tree should be converted into the target DTD

&lt;!ELEMENT geo ((river|lake|sea)*)&gt;
&lt;!ELEMENT river (. . . )&gt;        &lt;!ATTLIST river name CDATA #REQUIRED . . . &gt;
   *(analogously for lakes and seas)*

The first rule,   result/T[@name→N] :- //water[@type→T and @name→N]
creates &lt;river name="Mississippi"/&gt;  and  &lt;sea name="North Sea"/&gt;.
Attributes and contents are then transformed by separate rules which copy properties by using variables at element name and attribute name position:

X[@A→V]:- //water[@type→T and @name→N and @A→V], //T→X[@name→N].
X[S→V]   :- //water[@type→T and @name→N and S→V], //T→X[@name→N].

### 4.3 Global semantics of positive XPathLog programs

An XPathLog program is a declarative specification how to manipulate an XML database, starting with one or more input documents. The semantics of XPathLog programs is defined by bottom-up evaluation based on a $T_P$ operator. Thus, the semantics coincides with the usual understanding of a stepwise process.

For implementing the once-for-each-binding approach, the $T_P$ operator has to be extended with bookkeeping about the instances of inserted rule heads. Additionally, the insertion of subelements adds some nonmonotonicity: adding an atom n[child(i)::e→v] adds a new subelement at the $i$th position, making the original $i$th child/sibling the $i+1$st, etc. In case of multiple extensions to the same element, the positions are determined wrt. the original structure.

*Definition 15 (Extension of X-structures)*
Given an X-Structure $\mathcal{X}$ and a set $\mathcal{I}$ of ground-instantiated atoms as obtained from atomize to be inserted, the new X-Structure $\mathcal{X}' = \mathcal{X} \prec \mathcal{I}$ is obtained as follows:

- initialize $\mathcal{A}_{\mathcal{X}'}(\text{child}, x) := \mathcal{A}_{\mathcal{X}}(\text{child}, x)$, $\mathcal{A}_{\mathcal{X}'}(\text{attribute}, x) := \mathcal{A}_{\mathcal{X}}(\text{attribute}, x)$,
    $\text{preds}(\mathcal{X}') := \text{preds}(\mathcal{X}) \cup \{p \mid p \in \mathcal{I} \text{ is a predicate atom}\}$
  for all node identifiers $x$.

- for all elements of $\mathscr{A}_{\mathscr{X}}(\text{child}, x)$, let $\alpha(\mathscr{A}_{\mathscr{X}}(\text{child}, x)[i]) := \mathscr{A}_{\mathscr{X}'}(\text{child}, x)[i]$ ($\alpha$ maps the indexing from the old list to the new one).
- for all atoms $x[\text{child}(i) :: e \rightarrow y] \in \mathscr{I}$, insert $(y, e)$ into $\mathscr{A}_{\mathscr{X}'}(\text{child}, x)$ immediately after $\alpha(\mathscr{A}_{\mathscr{X}}(\text{child}, x)[i])$.
- for all atoms $x[\text{child} :: e \rightarrow y] \in \mathscr{I}$, append $(y, e)$ at the end of $\mathscr{A}_{\mathscr{X}'}(\text{child}, x)$.
- analogously for sibling axes.
- for all atoms $x[@a \rightarrow y] \in \mathscr{I}$, append $(y, a)$ to $\mathscr{A}_{\mathscr{X}'}(\text{attribute}, x)$.

*Proposition 8* (*Extension of X-Structures*)
The extension operation is correct: $\mathscr{X} \prec \mathscr{I} \models \mathscr{I}$, i.e. when querying the inserted atoms, the query evaluates to true.

With the correctness of atomize, the insertion of rule heads performs correctly:

*Corollary 9* (*Correctness of Insertions*)
For inserting the ground-instantiated head of a rule, it is correct to insert the atomized head: For all ground XPathLog atoms A, $\mathscr{X} \prec \text{atomize}(A) \models A$.

*Definition 16* (*$TX_P$-Operator for XPath-Logic Programs*)
The $TX$-operator works on pairs $(\mathscr{X}, Dic)$ where $\mathscr{X}$ is an X-Structure, and $Dic$ is a dictionary which associates to every rule a set $\xi$ of bindings which have been instantiated in the current iteration:

$$(\mathscr{X}, Dic) + (\{(r_1, \beta_1), \ldots, (r_n, \beta_n)\}) \quad := \quad (\mathscr{X} \prec \{\beta_i(\text{atomize}(head(r_i))) \mid 1 \leq i \leq n\},$$
$$Dic.insert(\{(r_1, \beta_1), \ldots, (r_n, \beta_n)\})),$$
$$(\mathscr{X}, Dic) \downarrow_1 \quad := \quad \mathscr{X}.$$

For an XPathLog program $P$ and an X-Structure $\mathscr{X}$,

$$TX_P(\mathscr{X}, B) := (\mathscr{X}, B) + \{(r, \beta) \mid r = (h \leftarrow b) \in P \text{ and } \mathscr{X} \models \beta(b), \text{ and } (r, \beta) \notin B\},$$
$$TX_P^0(\mathscr{X}) \quad := (\mathscr{X}, \emptyset),$$
$$TX_P^{i+1}(\mathscr{X}) \quad := TX_P(TX_P^i(\mathscr{X})),$$
$$TX_P^\omega(\mathscr{X}) \quad := \begin{cases} (\lim_{i \to \infty} TX_P^i(\mathscr{X})) \downarrow_1 & \text{if } TX_P^0(\mathscr{X}), TX_P^1(\mathscr{X}), \ldots \text{ converges,} \\ \bot & \text{otherwise.} \end{cases}$$

*Remark 3*
Note that for pure Datalog programs $P$ (i.e. only predicates over first-order terms), the evaluation wrt. $TX_P$ does not change the semantics, i.e. $TX_P^\omega(\mathscr{X}) = T_P^\omega(\mathscr{X})$.

*Proposition 10* (*Properties of the $TX_P$ operator*)
The $TX_P$ operator extends the well-known $T_P$ operator. For all positive XPathLog programs $P$, the following holds:

- without considering context functions, the $TX_P$ operator is monotonous (which guarantees that a minimal fixpoint $TX_P^\omega(\mathscr{X})$ exists),
- $TX_P^\omega(\mathscr{X}) \models P$,
- $TX_P$ is order-preserving: for all XPathLog reference expressions *expr* which do not use negation or context functions, $\mathscr{S}_{\mathscr{X}}(expr)$ is a sublist of $\mathscr{S}_{TX_P(\mathscr{X})}(expr)$,
- for all atoms $A$ that do not contain aggregations or function applications, if $A$ holds in $\mathscr{X}$, then it also holds after application of $TX_P$: $\mathscr{X} \models A \Rightarrow TX_P(\mathscr{X}) \models A$.

**Proof** Both properties follow immediately from the definition. The child and attribute axes are extended solely by appending and inserting new "facts".

### 4.4 Semantics of general XPathLog programs

For logic programs which use negation (or similar nonmonotonic features, such as aggregation), there is no *minimal model semantics*. Instead, their semantics is defined wrt. perfect models, well-founded models, or stable models. For practical use – especially when considering bottom-up evaluation – the notion of *perfect models* and *stratification* (Przymusinski, 1988) provides a solution to the problems raised by negation and other nonmonotonic features. Stratification expresses the intuitive notion of process which executes as a sequence of steps.

Note that already not all Datalog programs are stratifiable. For logics over complex structures such as F-Logic, a reasonable notion of stratification can be defined based on the names occurring at property position – as long as variables are not allowed at the property position. With variables allowed at property position, it has been showed for F-Logic in Frohn (1998) that programs are in general not stratifiable. Since (i) even without variables at property position, there are many programs which are not syntactically stratifiable, and (ii) variables at the property position prove to be very useful for data integration (cf. Example 13), syntax-based stratification is not suitable for our approach. Since the intention of XPathLog programs is in general to implement a stepwise process by bottom-up evaluation, often there is a natural, *user-defined* stratification. User defined stratification is supported in the LoPiX system May (2001d) (cf. section 5). The semantics is computed in the same way as for positive programs by iterating the $TX_P$ operator for each stratum.

*Language extensions.* In addition to the pure language as described above, XPathLog supports several extensions. A detailed description of, for example, aggregation (as known, e.g. from SQL) and a class hierarchy and signatures (taken from F-Logic), and data-driven Web access can be found in May (2001a).

## 5 Implementation and application

### 5.1 Implementation: The LoPiX system

XPathLog has been implemented in the LoPiX system May (2001d) which extends the pure XPathLog language with a Web-aware environment and additional functionality for data integration. LoPiX has been developed using major components from the Florid system (FLORID, 1998; Ludäscher *et al.*, 1998), an implementation (in C++) of F-Logic. Due to the similarities between the F-Logic data model and the XML data model in general, and XPathLogic's multi-overlapping-tree model in particular, the Florid modules provided a solid base for an XPathLog implementation. Especially the functionality of the complete module for the evaluation of a deductive language over a data model with complex objects could be reused. The system architecture of LoPiX is depicted in Figure 4.
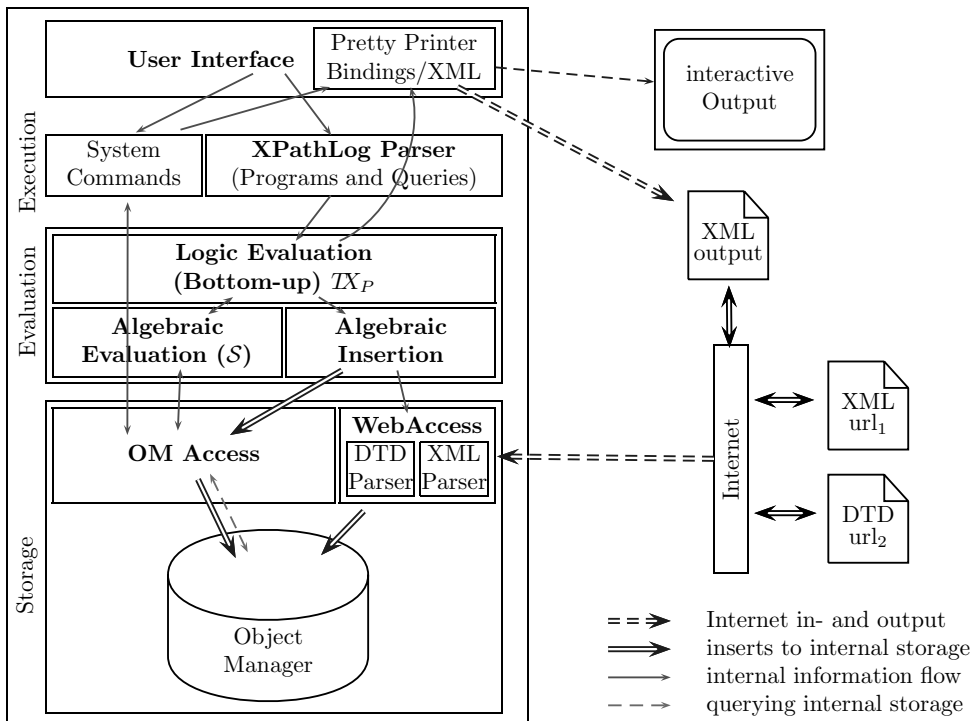
Fig. 4. Architecture of the LoPiX system.

*Storage.* The (extensional) database, is stored in the *ObjectManager*. Here, two variants have been developed: the first one uses a proprietary integrated, frame-based model (from FLORID) that is equipped with indexes for optimized access, whereas the second one is based on a standard DOM implementation.

The *ObjectManagerAccess* encapsulates the storage by implementing the abstract XTreeGraph data model based on the contents of the *ObjectManager*. This abstraction level also adds intensional properties including derived axes, transitivity of class hierarchy, downwards closure of signatures, inheritance, object fusion, synonyms, built-in functionality for data conversion, string handling including matching regular expressions, arithmetics, and aggregation operators.

The *WebAccess* functionality is closely intertwined with the *OMAccess* module: XML sources are mapped to trees in the internal database. Additionally, a method for mapping a DTD to XPathLog signature atoms is provided.

*Evaluation.* The central *Evaluation* module (*LogicEvaluation, AlgebraicEvaluation*, and *AlgebraicInsert*) is taken nearly unchanged from FLORID and provides in fact a *generic* implementation of a deductive language over a data model with complex objects. *LogicEvaluation* implements a seminaive bottom-up evaluation of rules. *AlgebraicEvaluation* translates rule bodies and heads into the underlying object algebra and evaluates the generated algebraic expressions using the query interface of *OMAccess*. The object algebra implements the semantics of XPathLog

queries described in section 3.1, generating sets of tuples of variable bindings. *AlgebraicInsert* instantiates the rule heads with the generated variable bindings and adds the corresponding facts into the database using again the *OMAccess* interface, implementing the $TX_P$-semantics defined in section 4. The evaluation of algebraic expressions does not materialize any intermediate result, but is purely based on nested iterators.

*Execution and UserInterface.* The execution module provides the infrastructure for the system, consisting of a *Parser* (lex/yacc-based) and a *SystemCommands* module that implements (partially) non-logical commands for controlling the evaluation process. The *UserInterface* module allows to use LoPiX from the command shell by invoking system commands and stating interactive queries. The *PrettyPrinter* outputs answers in the variable bindings format known from Datalog; additionally, the result of queries that bind only a single variable can be output as a result set in XML ASCII representation. Additionally, *result views*, i.e. the projections of trees rooted in a given node to a given signature can be exported.

### *5.2  Case study:* **Mondial**

XPathLog/LoPiX has successfully been applied in the MONDIAL case study (May, 2001e; May 2001b). There, the practicability of the approach for data integration is illustrated by integrating a geographical database from the XML representations of its sources (which have been created by FLORID wrappers in May (1999)).

**The CIA World Factbook:** The CIA World Factbook Country Listing (cia:, `http://www.odci.gov/cia/publications/pubs.html`) provides political, economic, and social and some geographical information about the countries. A separate part of the CIA World Factbook provides information about political and economical organizations (orgs:). Here, the data sources overlap by the membership relation: with every organization, the member countries are stored in orgs by name (using the same names as in the cia part).

**Global Statistics: Cities and Provinces:** The *Global Statistics* data (gs:, `http://www.stats.demon.nl`) provides information (grouped by countries) about administrative divisions (area and population, sometimes capital) and main cities (population with year, and province). Whereas the country names are the same as in CIA, the names of cities, that are for example capitals of countries or where the headquarter of a political organization is located, may differ.

The case-study showed that XPathLog allows for an effective, and elegant programming of the integration process. The nature of an XPathLog program as a list of rules allows for grouping rules which together handle a certain task. The programs are modular which also allows for adapting them to potential changes in the source structure and ontology.

For data integration in general, not only "simple" updates are desired, but also specialized operations on tree fragments. The result is constructed using subtrees, elements, and literals of the input sources by the integration operations that extend
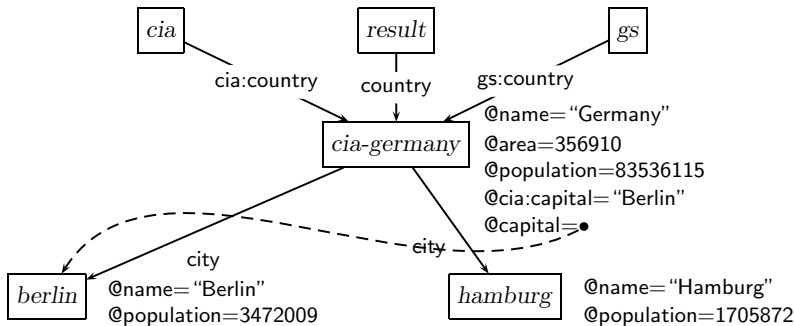
(a) Element fusion – before



Fig. 5. (b) Element fusion – after.

the basic XPathLog in LoPiX. These operations heavily depend on the use of the XTreeGraph data model (May and Behrends, 2001).

*Fusing elements and subtrees. Fusing* elements that represent the same real-world entity from different data sources into a unified element is an important task in information integration. The result is still an element of *both* source trees, and collects the attributes and subelements of both original elements.

*Example 14 (Object Fusion)*
Consider two data sources as shown below and in Figure 5(a). Both describe countries, where cia contains information about name, area, population, and capital, and gs contains information about cities.

```
<!ELEMENT cia (country+)>
<!ELEMENT country (border*)>
    <!ATTLIST country  name CDATA #REQUIRED      car_code ID #REQUIRED
                       area CDATA #IMPLIED       population CDATA #IMPLIED
                       capital CDATA #REQUIRED>
<!ELEMENT border (#PCDATA)>    <!ATTLIST border country IDREF #REQUIRED>
<!ELEMENT gs (country+)>
<!ELEMENT country (city+)>  <!ATTLIST country  name CDATA #REQUIRED>
<!ELEMENT city EMPTY>
    <!ATTLIST city name CDATA #REQUIRED pop CDATA #REQUIRED>
```

*Excerpts of the instances:*

```
<cia>                                  <gs>
  <country car_code='D'  capital='Berlin'    <country name='Germany'>
        name='Germany'  area='356910'      <city name="Berlin" pop="3472009"/>
        population='83536115'>          <city name="Hamburg" pop="1705872"/>
    <border country='F'>451</border>        :
    <border country='A'>784</border>      </country>
      :                                  :
  </country>                           </gs>
    :
</cia>
```

An obvious and typical integration step is to unify the countries in the cia tree with the countries in the gs tree. In XPathLog, this is done by the rule

C1 = C2 :- cia/cia:country→C1[@cia:name→N], gs/gs:country→C2[@gs:name→N].

The example is continued below – Figure 5(b) depicts the final result.

*Synonyms.* Names are also subject of operations, for example, the integrated database uses a unified terminology that differs from the source terminologies. Instead of generating new relationships between nodes, target terminology is introduced by synonyms for already existing relationships.

*Example 15 (Integration: Synonyms)*
Especially, synonyms are an efficient means for taking a whole property from a source tree (and namespace) to the result tree: Consider the situation obtained in Example 14 where the following synonyms are defined:

cia:name = name.        gs:city = city.        gs:text() = text().
cia:area = area.        gs:name = name.
cia:population = population.    gs:pop = population.

*Adding Links.* The integrated database often contains additional links (by subelement or reference attribute relationships) between elements that originally belong to different sources.

*Example 16 (Integration: Additional Links)*
The integration is completed by linking the country subtrees to a result tree and adding the capital reference attributes, here, using *germany*[@cia:capital="Berlin"] and *berlin*[name="Berlin"]. The resulting tree fragment is given in Figure 5(b). In XPathLog, this is done by the rules

result[country→C] :- cia[cia:country→C].
C[@capital→City] :-
    result/country→C[@cia:capital→Name and city→City[@name=Name]].

*Projection.* When the integration and restructuring process is completed, projections are used to define *result views* of the internal database. A result view is an XML tree, e.g. specified by a root node and a DTD.

The complete case study in May (2001b) describes the process of data integration, data cleaning, data restructuring and distinguishing result tree views. The program is easily extendible by additional rules for adding another data source.

## 6 Analysis, related work, and conclusion

### 6.1 Comparison with other XML languages

*XPathLog vs. requirements.* In Fernandez *et al.* (1999), XQL, XML-QL, and the languages YATL Cluet *et al.* (1999) and *Lorel* Abiteboul *et al.* (1997); Goldman *et al.* (1999) are compared and essential features of an XML query language have been identified. XPathLog relates to their requirements as follows:

- existence of some kind of pattern clause, step qualifier clause, and constructor clause: pattern and step qualifier clause are the same as in XPath, extended with variable bindings. The path patterns are superior to XML patterns (e.g. as used in XML-QL) since they allow for dereferencing and navigation along different axes. The constructor clause uses the same XPath-based syntax.
- constructs for imposing nesting and order: nested elements in the result tree are generated by subsequent rules which stepwise generate the result. Grouping (via stepwise generation) and order (via child($i$)::name) is supported.
- combining data from different sources is supported.
- tag variables or regular path expressions: tag variables are supported, regular path expressions are not included in the basic XPathLog language (also not in XPath). They are definable as derived relations.
- alternatives are expressible using a separate rule for each alternative.
- checking for absence of information: existence or non-existence of properties can be tested using negation, e.g. //country[not @indep_date].
- external functions: aggregation, string functions and some data conversion is built-in; the set of functions is extensible.
- navigation along references: implicit dereferencing is supported.

*Semistructured data languages.* We have already mentioned the use of logic programming style languages in pre-XML projects on semistructured data in section 1.

*GraphLog* (Consens and Mendelzon, 1990) and F-Logic/FLORID (Kifer and Lausen, 1989; Kifer *et al.*, 1995) presented logic-programming languages over graph data models that cover the semistructured data model, but did not yet use that notion.

In GraphLog, graphical queries are defined as patterns that are matched with an underlying graph database. The matched vertices are bound to variables that are then used for generating an output instance or for adding edges to the input graph in the rule head. In the graphical representation, the "rule head" is represented as a distinguished edge in the graphical pattern (to be added to the input graph). The language can be seen as a graphical representation of Datalog over binary relations. Thus, according to our criteria stated in section 1.1, GraphLog qualifies as a logic-programming language. GraphLog excludes recursive rules, but allows for *closure literals* that represent the closure of a binary predicate; thus the expressiveness of the language is the same as for stratified linear Datalog.

F-Logic (Kifer and Lausen, 1989; Kifer *et al.*, 1995) is a *deductive object-oriented database language* that can be seen as an early concept of a *semistructured, self-describing* data model. F-Logic defines a data model, a logic, and a database query

and programming language (similar to the relationship between the X-Structures, XPath-Logic and XPathLog). The experiences with F-Logic as a formal framework and as a language for data extraction and integration from the Web (Ludäscher *et al*., 1998; May, 1999) provided the background for the design of XPath-Logic and XPathLog as a crossbreed between XPath and F-Logic, combining the experiences with F-Logic as a successful (but "proprietary") language for data integration with the standards of XML and XPath was a well-grounded evolution step. Especially, the power of the graph-based F-Logic data model compared with the restricted tree model of XML made up a central requirement in the design of XPathLog, leading to the XTreeGraph data model for *virtual trees in a graph database*. Another important aspect taken from F-Logic is to have names as first-order citizens of the language for a seamless incorporation of metadata information. Due to these similarities, it was possible to base the implementation of XPathLog in the LoPiX system on the F-Logic system FLORID.

The *OEM (Object Exchange Model)* of the TSIMMIS project (Garcia-Molina *et al*., 1997; Abiteboul *et al*., 1997) was the first data model that was dedicated explicitly to the notion of *semistructured data*. OEM is a graph based model, for which node-labeled and edge-labeled presentations have been given. With *WSL* and *MSL (Wrapper/Mediator Specification Language)*, Datalog-style programming languages have been presented. The *Lorel* language (McHugh *et al*., 1997) is similar to OQL, combining navigational access (extended with regular path expressions) with clauses. Lorel supports SQL-like, procedural update constructs. *Lorel* has been migrated to XML in Goldman *et al*. (1999). In contrast to the XPathLog/LoPiX migration, LOREL does not support the XML axes.

*UnQL* (Buneman *et al*., 1996, 2000) operates on rooted, edge-labeled graphs. It embeds *graph schemata* that are matched as patterns with the underlying database, combined with navigational access into SQL-like clauses. UnQL's semantics is based on *structural recursion* – similar to the later XSL.

*Strudel/StruQL* (Fernandez *et al*., 1997, 1998) also uses an edge-labeled graph model. Its syntax embeds query patterns that are matched with the underlying database into SQL-like clauses. StruQL rules specify what new elementary structures are created, and what links between them are created. The Strudel project has been continued for XML with XML-QL.

The *YATL* language of the *YAT* system (Cluet *et al*., 1999) is a pre-XML proposal, already using SGML and DTDs. Its trees provide a unified model for relational, object-oriented (ODMG), and semistructured/document data (SGML). The YATL language follows a rule-based design for complex objects in the style of MSL or F-Logic; it supports regular path expressions and tree algebraic operations. In Christophides *et al*. (2000), the YAT system is turned into an XML system for data integration, which still does not use any XML/XPath language constructs. After mapping an XML instance to a YAT tree, there is no notion of attributes. Dereferencing is not explicitly supported, and it has no notion of the XML axes (similar to the same issue for XML-QL).

*XML languages.* XML-QL and XQuery embed XML patterns and XPath expressions, respectively, into SQL-style clauses. Expressions can be nested.

XML-QL (Deutsch *et al.*, 1999) uses XML patterns in the head (CONSTRUCT) and body (WHERE) clause. In that aspect, it is the XML-pattern-counterpart to the XPath-based XPathLog. The XML-QL patterns for selecting elements do not support the XML *axes* except the child axis, and indirectly the descendant by regular path expressions. XML-QL does not support updates; a potential combination of XML patterns and updates is not obvious.

XQuery (XQuery, 2001) embeds XPath expressions in SQL-style FOR – LET – WHERE – RETURN clauses, where the RETURN clause specifies the result as an XML pattern. A proposal for specifying *updates* in XQuery has been published in Tatarinov *et al.* (2001). A more detailed proposal is described in Lehti (2001) and implemented in (Software AG, 2001).

XML-GL (Ceri *et al.*, 1999; Comai *et al.*, 2001) continued the idea of GraphLog for XML. In contrast to GraphLog, the rule body and the rule head are represented by separate graphs, called *extract-match-construct-clip-queries*. The rule heads *generate* separate XML structures. Recursion is excluded. The MIX *(Mediation in XML)* system (Baru *et al.*, 1999) uses the XMAS *(XML Matching and Structuring)* language, derived from XML-QL for data integration; a graphical user interface similar to XML-GL is provided. *XDuce* (Hosoya and Pierce, 2000) is a functional-style tree transformation language which uses *regular expression pattern matching* of (originally, SGML) DTDs for formulating queries against XML instances.

*Xcerpt* (Bry and Schaffert, 2002) is a *pattern-based* language for querying and transforming XML data. It follows a clean, rule-based design where the query (matching) part in the body is separated from the generation part in the rule head. XML instances are regarded as terms that are matched by a term pattern in the rule body, generating variable bindings. The semantics and the implementation is given by *simulation unification* that computes answer substitutions for the variables in the match pattern against the underlying XML term (similar to UnQL). Then, the term in the rule head is instantiated with these variable bindings. Since rule heads have only a *generating* semantics, but not an update semantics, Xcerpt can only be used for querying and transforming XML data, but not for updating/extending an existing internal XML database. It has a rule-based semantics, but there is no global logic programming semantics for the evaluation of programs.

*Elog* (Baumgartner *et al.*, 2001a) is a logic programming language for XML which is used as *internal* language for XML data extraction in the Lixto project (Baumgartner *et al.*, 2001b). It is based on flattening XML data into Datalog with specialized Web Access predicates.

Table 1 gives a comparison of some of the above-mentioned languages. The "paradigm" column indicates the underlying semantics of the languages: the semantics of SQL-like languages is best given as an algebraic semantics that specifies the type and value of expression, allowing for nested expressions. For rule-based languages, a denotational specification of the outcome of the right-hand side (query) and of the result of the left-hand side is required. Logic programming languages

Table 1. *Comparison*

| **SSD** | GraphLog | WSL/MSL | Lorel | UnQL | StruQL | F-Logic |
|---|---|---|---|---|---|---|
| DataModel | graph | graph/atoms | graph | graph | graph | graph |
| Access | patterns | patterns | pat/nav | term unif. | patterns | navigation |
| Views | y | y | y | y | y | y |
| Interfering | | | | | | |
| additions | y | n | y | n | n | y |
| Paradigm | LP | rules | SQL | SQL | SQL | LP |
| **XML** | XML-QL | XQuery | XML-GL | Xcerpt | Elog | XPathLog |
| DataModel | XML | XML | XML tree | XML tree | atoms | XTreeGraph |
| Access | patterns | navigation | patterns | term unif. | (atoms) | navigation |
| Views | y | y | y | y | y | y |
| Interf. add. | n | (XUpdate) | n | n | (+) | y |
| standard- | | | | | | |
| based | no | is standard | no | no | (no) | yes: Xpath |
| Paradigm | SQL | SQL | rules | rules | LP | LP |

require both a model-theoretic semantics (to specify the outcome of rule heads, and for the *global* semantics), and an answer semantics for the querying part.

## 6.2 Contributions

We have described XPath-Logic as a logic-based framework for handling XML data, together with an extended XML data model that is suitable for XML querying, manipulation, and integration. XPathLog combines the intuitive "local" semantics of addressing XML data by XPath with the appeal of the "global" logic programming semantics: it is completely XPath-based, i.e. both the rule bodies and the rule heads use an extended XPath syntax, thereby defining a *constructive* semantics for XPath expressions. Although the syntactic difference between XPath and XPathLog is small, the extension adds much to the language by turning it into a data manipulation language. The close relationship with XPath ensures that its declarative semantics is well understood from the XML perspective. Since both XPath and rule-based programming by using variable bindings are well-known, intuitive concepts, the "effect" of the language is easy to understand on an intuitive basis, making programming easy. The logic programming background provides a strong theoretical foundation of the language concept.

The data model and the language are implemented in the LoPiX system. Its practicability has been demonstrated by the MONDIAL case study.

## Appendix A Proofs

**Proof** of Theorem 1 and Lemma 2: The proof is done by structural induction. The enumeration is the same as in Definition 4. Below, $\beta$ is an assignment of the pseudo variables *Size* and *Pos* (often even empty). We write $\overset{**}{=}$ for "equals by definition

in (Wadler, 1999)". The individual items of the theorem are referred to below by
IH1,…, IH4 (induction hypotheses).

1. For closed, absolute expressions (i.e. without free variables),
$$\mathscr{S}_{\mathscr{X}}(/expr) \stackrel{\text{Def}}{=} \mathscr{S}_{\mathscr{X}}^{any}(/expr, any, \emptyset) \stackrel{\text{IH1}}{=} \mathscr{S}[[/expr]](x) \quad \text{for arbitrary } x.$$

2. Reference expressions ((Wadler, 1999): only absolute expressions):
$$\mathscr{S}_{\mathscr{X}}^{any}(/p, any, \beta) \stackrel{\text{Def}}{=} \mathscr{S}_{\mathscr{X}}^{any}(p, root, \beta) \stackrel{\text{IH2}}{=} \mathscr{S}^{any}[[/expr]](root).$$

3. Axis step:
$$\mathscr{S}_{\mathscr{X}}^{any}(axis :: pattern, x, \beta) \stackrel{\text{Def}}{=} \mathscr{S}_{\mathscr{X}}^{axis}(pattern, x, \beta) \stackrel{\text{IH2}}{=} \mathscr{S}^{axis}[[/pattern]](x).$$

4. The node test is the base case which is directly mapped to the axes:
$$\mathscr{S}_{\mathscr{X}}^{a}(name, x, \beta) \stackrel{\text{Def}}{=} \text{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,x)}(v \mid n = name)$$
which is characterized in Wadler (1999) ($\mathscr{A}[[a]]$ enumerates the axes, $\mathscr{P}(a)$ gives the axes' principal nodetype) by
$$\{x_1 \mid x_1 \in \mathscr{A}[[a]]x, \ nodetype(x_1) = \mathscr{P}(a), \ name(x_1) = name\}$$
which is the definition of $\mathscr{S}^a[[name]](x)$. Note that dereferencing `IDREF(S)` and splitting `NMTOKENS` has been excluded, thus, the result list is still in document order. Similar (note that `node()` is not defined in Wadler (1999), we extend the definition according to the XPath specification)
$$\mathscr{S}_{\mathscr{X}}^{a}(\text{node}(), x, \beta) \stackrel{\text{Def}}{=} \text{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,x)}(v \mid v \in \mathscr{V})$$
$$= \{x_1 \mid x_1 \in \mathscr{A}[[a]]x, \ nodetype(x_1) = element\} = \mathscr{S}^a[[\text{node}()]](x)$$
$$\mathscr{S}_{\mathscr{X}}^{a}(\text{text}(), x, \beta) \stackrel{\text{Def}}{=} \text{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,x)}(v \mid v \in \mathscr{V})$$
$$= \{x_1 \mid x_1 \in \mathscr{A}[[a]]x, \ nodetype(x_1) = Text\} = \mathscr{S}^a[[\text{text}()]](x).$$

5. Step with variable binding: obvious

6. Step qualifiers:
$$\mathscr{S}_{\mathscr{X}}^{a}(pattern[stepQ], x, \beta) \stackrel{\text{Def}}{=} \text{list}_{y \in \mathscr{S}_{\mathscr{X}}^{a}(pattern, x, \beta)}\left(y \mid Q_{\mathscr{X}}\left(stepQ, y, \beta_{Pos,Size}^{k,n}\right)\right)$$
where $L_1 \stackrel{\text{Def}}{=} \mathscr{S}_{\mathscr{X}}^{a}(pattern, x, \beta)$ which equals $\mathscr{S}^a[[pattern]](x, k, n)$ by induction hypothesis IH3 and $n := size(L_1)$, and for every $y$, let $j$ the index of $y$ in $L_1$ (which equals $size(\{x_1 \mid x_1 \in L_1, x_1 \leqslant_{doc} y\})$), $k := j$ if $a$ is a forward axis, and $k := n+1-j$ if $a$ is a backward axis. This is the same as defined for $\mathscr{S}^a[[pattern[stepQualifier]]](x)$ and, by induction hypothesis IH3, the same as
$$\stackrel{\text{IH3}}{=} \text{list}_{y \in \mathscr{S}_{\mathscr{X}}^{a}(pattern, x, \beta)}(y \mid Q[[stepQ]](, y, k, n)) \stackrel{**}{=} \mathscr{S}^a[[pattern[stepQ]]](x).$$

7. Path:
$$\mathscr{S}_{\mathscr{X}}^{a}(p_1/p_2, x, \beta) \stackrel{\text{Def}}{=} \text{concat}_{y \in \mathscr{S}_{\mathscr{X}}^{a}(p_1, x, \beta)}\left(\mathscr{S}_{\mathscr{X}}^{any}(p_2, y, \beta)\right)$$
$$\stackrel{\text{IH2}}{=} \text{concat}_{y \in \mathscr{S}^a[[p_1]](x)}(\mathscr{S}^a[[p_2]](y)) \stackrel{**}{=} \mathscr{S}^a[[p_1/p_2]](x).$$

8. Reference expressions (existential semantics) in step qualifiers:
$$Q_{\mathscr{X}}(refExpr, x, \beta) \stackrel{\text{Def}}{\Leftrightarrow} \mathscr{S}_{\mathscr{X}}^{any}(refExpr, x, \beta) \neq \emptyset$$
$$\stackrel{\text{IH3}}{\Leftrightarrow} \mathscr{S}^{child}[[refExpr]](x) \neq \emptyset \stackrel{**}{\Leftrightarrow} Q[[refExpr]](x, k, n).$$

(for all $k, n$ since these are not used in *refExpr*).

9. Predicates: Wadler (1999) knows only the "=" comparison. The definition is although not complete: e.g. for step qualifiers of the form [a/b/c = "foo"]

which are allowed in XPath, there is no semantics defined. We extend the semantics according to the XPath specification, applying either $\mathscr{S}$ or $\mathscr{E}$.

$\mathcal{Q}_{\mathcal{X}}(pred(expr_1,\ldots,expr_n),x,\beta)$
$\overset{\text{Def}}{\Leftrightarrow}$　there are $x_1 \in \mathscr{S}_{\mathcal{X}}^{any}(expr_1,x,\beta),\ldots,x_n \in \mathscr{S}_{\mathcal{X}}^{any}(expr_n,x,\beta)$
　　　such that $(x_1,\ldots,x_n) \in \mathscr{I}_P(pred)$

$\overset{\text{IH2/4}}{\Leftrightarrow}$　there are $x_1 \in \mathscr{S}^{child}[[expr_1]](x)$ or $x_1 \in \mathscr{E}[[expr_1]](x,\beta(Pos),\beta(Size)),\ldots,$
　　　　$x_n \in \mathscr{S}^{child}[[expr_n]](x)$ or $x_n \in \mathscr{E}[[expr_n]](x,\beta(Pos),\beta(Size))$

　　such that $(x_1,\ldots,x_n) \in \mathscr{I}(pred)$.

10. – 13. Boolean connectives and quantification, constants, and variables: obvious. Functions are not defined in Wadler (1999), but the extension is obvious.

14. Context-related functions use the extension of variable bindings by pseudo-variables *Size* and *Pos* in rule (6):

$\mathscr{S}_{\mathcal{X}}^{any}(\mathsf{position}(),x,\beta) \quad \overset{\text{Def}}{\Leftrightarrow} \quad \beta(Pos) \quad \overset{**}{=} \quad \mathscr{E}[[\mathsf{position}()]](x,\beta(Pos),\beta(Size))$
$\mathscr{S}_{\mathcal{X}}^{any}(\mathsf{last}(),x,\beta) \quad \overset{\text{Def}}{\Leftrightarrow} \quad \beta(Size) \quad \overset{**}{=} \quad \mathscr{E}[[\mathsf{last}()]](x,\beta(Pos),\beta(Size)).$

**Proof** of Lemma 4:

> **Note:** A bit sloppy, we write $(x,\beta) \in \mathscr{SB}_{\mathcal{X}}(expr)$ for "$x \in \text{Res}(\mathscr{SB}_{\mathcal{X}}(expr))$ and $\beta \in \text{Bdgs}(\mathscr{SB}_{\mathcal{X}}(expr),x)$".

1. For closed expressions, $x \in \text{Res}(\mathscr{SB}_{\mathcal{X}}(refExpr)) \quad \overset{\text{Def}}{\Leftrightarrow}$

$x \in \text{Res}(\mathscr{SB}_{\mathcal{X}}(refExpr,\emptyset)) \quad \overset{\text{IH}}{\Leftrightarrow} \quad x \in \mathscr{S}_{\mathcal{X}}(refExpr,\emptyset) \quad \overset{\text{Def.4}}{\Leftrightarrow} \quad x \in \mathscr{S}_{\mathcal{X}}(refExpr).$

2. Reference expressions are translated into path expressions wrt. a start node:
   - entry points: rooted path

   $(x,\beta) \in \mathscr{SB}_{\mathcal{X}}(/p,Bdgs) \quad \overset{\text{Def}}{\Leftrightarrow} \quad (x,\beta) \in \mathscr{SB}_{\mathcal{X}}^{any}(p,root,Bdgs)$
   　　$\overset{\text{IH}}{\Leftrightarrow} \quad x \in \mathscr{S}_{\mathcal{X}}^{any}(p,root,\beta)$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(/p))$
   　　$\overset{\text{Def.4}}{\Leftrightarrow} \quad x \in \mathscr{S}_{\mathcal{X}}^{any}(/p,\beta)$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(/p))$.

   - entry points: constants $c \in \mathscr{V}$ analogously (set $c$ instead of *root* above).
   - entry points: variables $V \in \mathsf{Var}$. By definition,

   $(x,\beta) \in \mathscr{SB}_{\mathcal{X}}(V/p,Bdgs) \quad \overset{\text{Def}}{\Leftrightarrow}$
   　　$(x,\beta) \in \mathsf{concat}_{x \in \mathscr{A}_{\mathcal{X}}(descendants,root)\downarrow_1}\left(\mathscr{SB}_{\mathcal{X}}^{any}(p,x,Bdgs \bowtie \{V/x\})\right)$

   which is exactly the case if there is an $x \in \mathscr{A}_{\mathcal{X}}(\mathsf{descendants},root)\downarrow_1$ such that $(x,\beta) \in (\mathscr{SB}_{\mathcal{X}}^{any}(p,x,Bdgs \bowtie \{V/x\}))$.
   By induction hypothesis, this is equivalent with

   $x \in \mathscr{S}_{\mathcal{X}}^{any}(p,x,\beta)$ and $\beta$ completes some $\beta' \in Bdgs \bowtie \{V/x\}$ with $\mathsf{free}(p)$

   which is exactly the case if $x = \beta(V)$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(V/p)$. By Def. 4, this again is equivalent with

   $x \in \mathscr{S}_{\mathcal{X}}^{any}(V/p,\beta)$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(V/p)$.

3. Axis step: $(x, \beta) \in \mathscr{SB}_{\mathscr{X}}^{any}(axis :: pattern, z, Bdgs)$

$$\overset{\text{Def}}{\Leftrightarrow} \quad (x, \beta) \in \mathscr{SB}_{\mathscr{X}}^{axis}(pattern, z, Bdgs)$$

$$\overset{\text{IH}}{\Leftrightarrow} \quad x \in \mathscr{S}_{\mathscr{X}}^{axis}(pattern, z, \beta)$$

$$\text{and } \beta \text{ completes some } \beta' \in Bdgs \text{ with } \mathsf{free}(pattern)$$

$$\overset{\text{Def.4}}{\Leftrightarrow} \quad x \in \mathscr{S}_{\mathscr{X}}^{any}(axis :: pattern, z, \beta)$$

$$\text{and } \beta \text{ completes some } \beta' \in Bdgs \text{ with } \mathsf{free}(axis :: pattern).$$

4. Node test: $(x, \beta) \in \mathscr{SB}_{\mathscr{X}}^{a}(name, z, Bdgs) \overset{\text{Def}}{\Leftrightarrow}$

$$(x, \beta) \in \mathsf{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,z), \ n=name}(v, \{true\} \bowtie Bdgs)$$

which is exactly the case if $x \in \mathsf{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,z), \ n=name}(v)$ and $\beta \in Bdgs$ which, by Def. 4 is equivalent with $x \in \mathscr{S}_{\mathscr{X}}^{a}(name, z, \beta)$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(name) = \emptyset$. Analogously for $\mathsf{node}()$ and $\mathsf{text}()$.

Variables at nodetest position:

$$(x, \beta) \in \mathscr{SB}_{\mathscr{X}}^{a}(N, z, Bdgs) \quad \overset{\text{Def}}{\Leftrightarrow} \quad (x, \beta) \in \mathsf{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,z)}(v, \{N/n\} \bowtie Bdgs)$$

which is exactly the case if $x \in \mathsf{list}_{(v,n) \in \mathscr{A}_{\mathscr{X}}(a,z)}(v)$ and $\beta \in \{N/n\} \bowtie Bdgs$ which, by Def. 4 is equivalent with $x \in \mathscr{S}_{\mathscr{X}}^{a}(N, z, \beta)$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(N) = \{N\}$.

5. Step with variable binding:

$(x, \beta) \in \mathscr{SB}_{\mathscr{X}}^{a}(pattern \rightarrow V, z, Bdgs)$
$\overset{\text{Def}}{\Leftrightarrow} (x, \beta) \in \mathsf{list}_{(y,\xi) \in \mathscr{SB}_{\mathscr{X}}^{a}(pattern,z,Bdgs)}(y, \xi \bowtie \{V/y\})$
$\Leftrightarrow$ there is a $\beta''$ s.t. $(x, \beta'') \in \mathscr{SB}_{\mathscr{X}}^{a}(pattern, z, Bdgs)$ and $\beta = \beta'' \bowtie \{V/x\}$.

By induction hypothesis, this is exactly the case if there is a $\beta''$ such that $x \in \mathscr{S}_{\mathscr{X}}^{a}(pattern, z, \beta'')$ and $\beta''$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(pattern)$, and $\beta = \beta'' \bowtie \{V/x\}$. Exactly then, since $x = \beta(V)$, by Definition 4, $x \in \mathscr{S}_{\mathscr{X}}^{a}(pattern \rightarrow V, z, \beta)$ and $\beta$ completes $\beta'$ with $\mathsf{free}(pattern \rightarrow V) = \mathsf{free}(pattern) \cup \{V\}$.

6. Step Qualifier(s): $(x, \beta) \in \mathscr{SB}_{\mathscr{X}}^{a}(pattern[stepQualifier], z, Bdgs)$

$$\overset{\text{Def}}{\Leftrightarrow} (x, \beta) \in \mathsf{list}_{\substack{(y,\xi) \in \mathscr{SB}_{\mathscr{X}}^{a}(pattern,z,Bdgs), \\ \mathcal{QB}_{\mathscr{X}}(stepQualifier,y,\xi') \neq \emptyset}}(y, \mathcal{QB}_{\mathscr{X}}(stepQualifier, y, \xi') \setminus \{Pos, Size\})$$

for $\xi$ as defined in Definition 10(6). This is exactly the case if (i) there is a $\beta''$ s.t. $\beta'' \in \mathcal{QB}_{\mathscr{X}}(stepQualifier, x, \xi')$ and $\beta = \beta'' \setminus \{Pos, Size\}$, and (ii) $(x, \xi) \in \mathscr{SB}_{\mathscr{X}}^{a}(pattern, z, Bdgs)$ i.e. $\xi$ is the corresponding set of variable bindings, and (iii) $\mathcal{QB}_{\mathscr{X}}(stepQualifier, x, \xi') \neq \emptyset$.

The first item is by induction hypothesis equivalent to $\mathcal{Q}_{\mathscr{X}}(stepQualifier, x, \beta'')$ and $\beta''$ completes some $\beta' \in \xi'$ with $\mathsf{free}(stepQualifier)$ (*).

The third item is redundant here (it avoids the addition of elements with empty bindings list to the result). Since $\beta''$ completes some $\beta' \in \xi'$ with $\mathsf{free}(stepQualifier)$, we know that $\gamma := \beta' \setminus \{Pos, Size\}$ is an element of $\xi$. Specializing the second item to $\gamma$ yields $(x, \gamma) \in \mathscr{SB}_{\mathscr{X}}^{a}(pattern, z, Bdgs)$.

By induction hypothesis, $x \in \mathscr{S}_{\mathscr{X}}^{a}(pattern, z, \gamma)$ (**) and $\gamma$ completes some $\gamma' \in Bdgs$ with $\mathsf{free}(pattern)$. Above, we derived $\gamma = \beta' \setminus \{Pos, Size\}$. Using (*),

since $\beta''$ is a completion of $\beta'$ with $\mathsf{free}(stepQualifier)$, completing $\gamma' \in Bdgs$ first to $\gamma$ (binding $\mathsf{free}(pattern)$), then to $\beta'$ (binding $Size$ and $Pos$), then to $\beta''$ (binding $\mathsf{free}(stepQualifier)$), we have $\mathcal{Q}_{\mathcal{X}}(stepQualifier, y, \beta'')$.

From (\*\*), since $\beta''$ completes $\gamma$, $x \in \mathcal{S}_{\mathcal{X}}^a(pattern, z, \beta'')$ thus by Def. 4, the desired result $x \in \mathcal{S}_{\mathcal{X}}^a(pattern[stepQualifier], z, Bdgs)$ for $\beta''$ which completes $\gamma' \in Bdgs$ with $\mathsf{free}(pattern[stepQualifier])$.

The argumentation showed the "$\Rightarrow$" direction (which is the more difficult direction since $\gamma$ must be guessed). "$\Leftarrow$" uses the same relationships and variable bindings.

7. Path: $(x, \beta) \in \mathcal{SB}_{\mathcal{X}}^a(p_1/p_2, z, Bdgs)$

$\overset{\text{Def}}{\Leftrightarrow}$ $\quad (x, \beta) \in \mathsf{concat}_{(y,\xi) \in \mathcal{SB}_{\mathcal{X}}^{any}(p_1, z, Bdgs)} \mathcal{SB}_{\mathcal{X}}^a(p_2, y, \xi)$

$\Leftrightarrow$ $\quad$ there is an $(y, \xi) \in \mathcal{SB}_{\mathcal{X}}^{any}(p_1, z, Bdgs)$ s.t. $(x, \beta) \in \mathcal{SB}_{\mathcal{X}}^a(p_2, y, \xi)$

$\overset{\text{IH}}{\Leftrightarrow}$ $\quad$ there is a $\gamma \in \xi$ s.t. there is a $\gamma'$ s.t. $x \in \mathcal{S}_{\mathcal{X}}^a(p_2, y, \gamma')$ and $\gamma'$ completes $\gamma$ with $\mathsf{free}(p_2)$.

For this $\gamma$, $(y, \gamma) \in \mathcal{SB}_{\mathcal{X}}^{any}(p_1, z, Bdgs)$ and by induction hypothesis again $y \in \mathcal{S}_{\mathcal{X}}^a(p_1, z, \gamma)$ and $\gamma$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(p_1)$. Thus, also $x \in \mathcal{S}_{\mathcal{X}}^a(p_2, y, \gamma')$ and $y \in \mathcal{S}_{\mathcal{X}}^a(p_1, z, \gamma')$ and by Def. 4, $x \in \mathcal{S}_{\mathcal{X}}^a(p_1/p_2, z, \gamma')$. $\gamma'$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(p_1) \cup \mathsf{free}(p_2)$.

8. Reference expressions (existential semantics) in step qualifiers:

$\beta \in \mathcal{QB}_{\mathcal{X}}(refExpr, z, Bdgs) \quad \overset{\text{Def}}{\Leftrightarrow} \quad \beta \in \bigcup_{(y,\xi) \in \mathcal{SB}_{\mathcal{X}}^{any}(refExpr, z, Bdgs)} \xi$

$\Leftrightarrow$ $\quad$ there is a $y$ s.t. $(y, \beta) \in \mathcal{SB}_{\mathcal{X}}^{any}(refExpr, z, Bdgs)$

$\overset{\text{IH}}{\Leftrightarrow}$ $\quad y \in \mathcal{S}_{\mathcal{X}}^{any}(refExpr, z, \beta)$
$\quad$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(refExpr)$

$\overset{\text{Def.4}}{\Leftrightarrow}$ $\quad \mathcal{Q}_{\mathcal{X}}(refExpr, z, \beta)$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(refExpr)$.

9. Built-in equality predicate "$=$": similar to predicates and variable assigments by $\rightarrow V$. All other predicates: $\beta \in \mathcal{QB}_{\mathcal{X}}(pred(arg_1, \ldots, arg_n), z, Bdgs)$

$\overset{\text{Def}}{\Leftrightarrow}$ $\quad \beta \in \bigcup_{(x_i, \xi_i) \in \mathcal{SB}_{\mathcal{X}}^{any}(arg_1, z, Bdgs),\ (x_1, \ldots, x_n) \in \mathcal{I}(pred)} \xi_1 \bowtie \ldots \bowtie \xi_n$

$\Leftrightarrow$ $\quad$ there are $(x_1, \xi_1), \ldots, (x_n, \xi_n)$ s.t. $(x_i, \xi_i) \in \mathcal{SB}_{\mathcal{X}}^{any}(arg_i, z, Bdgs)$
$\quad$ and $(x_1, \ldots, x_n) \in \mathcal{I}(pred)$ and $\beta \in \xi_1 \bowtie \ldots \bowtie \xi_n$

$\Leftrightarrow$ $\quad$ (take the right $\beta_i \in \xi_i$)
$\quad$ there are $(x_1, \beta_1), \ldots, (x_n, \beta_n)$ s.t. $(x_i, \beta_i) \in \mathcal{SB}_{\mathcal{X}}^{any}(arg_i, z, Bdgs)$
$\quad$ and $(x_1, \ldots, x_n) \in \mathcal{I}(pred)$ and $\beta = \beta_1 \bowtie \ldots \bowtie \beta_n$

$\overset{\text{IH}}{\Leftrightarrow}$ $\quad$ there are $(x_1, \beta_1), \ldots, (x_n, \beta_n)$ s.t. $x_i \in \mathcal{S}_{\mathcal{X}}^{any}(arg_i, z, \beta_i)$
$\quad$ and $\beta_i$ extends some $\beta_i' \in Bdgs$ with $\mathsf{free}(arg_i)$
$\quad$ and $(x_1, \ldots, x_n) \in \mathcal{I}(pred)$ and $\beta = \beta_1 \bowtie \ldots \bowtie \beta_n$

$\Leftrightarrow$ $\quad$ (the join guarantees that $\beta' := \beta_1' = \ldots = \beta_n'$ holds)
$\quad$ there are $x_1, \ldots, x_n$ s.t. $x_i \in \mathcal{S}_{\mathcal{X}}^{any}(arg_i, z, \beta_i)$
$\quad$ and $\beta$ extends some $\beta' \in Bdgs$ with $\mathsf{free}(arg_1) \cup \ldots \cup \mathsf{free}(arg_n)$

$\overset{\text{Def.4}}{\Leftrightarrow}$ $\quad \mathcal{Q}_{\mathcal{X}}(pred(arg_1, \ldots, arg_n), z, \beta)$
$\quad$ and $\beta$ completes some $\beta' \in Bdgs$ with $\mathsf{free}(pred(arg_1, \ldots, arg_n))$.

10. Negated expressions which do *not contain any free variable*: trivial.
    For negated expressions which contain free variables: note that all variables in
    free(not *expr*) are required to be bound by *Bdgs* (safety).

    $$\beta \in \mathcal{QB}_{\mathcal{X}}(\text{not } expr, z, Bdgs)$$
    $$\overset{\text{Def}}{\Leftrightarrow} \quad \beta \in Bdgs \text{ and there is no } \beta' \in \mathcal{QB}_{\mathcal{X}}(expr, z, Bdgs) \text{ s.t. } \beta \leqslant \beta'$$
    $$\overset{\text{IH}}{\Leftrightarrow} \quad \beta \in Bdgs \text{ and there is no } \beta'' \text{ such that}$$
    $$\mathcal{Q}_{\mathcal{X}}(expr, z, \beta'') \text{ and } \beta'' \text{ extends } \beta' \text{ with free}(expr) \text{ and } \beta \leqslant \beta'$$
    $$\overset{\text{Safety}}{\Leftrightarrow} \quad \beta \in Bdgs \text{ and not } \mathcal{Q}_{\mathcal{X}}(expr, z, \beta)$$
    $$\overset{\text{Def.4}}{\Leftrightarrow} \quad \beta \in Bdgs \text{ and } \mathcal{Q}_{\mathcal{X}}(\text{ not } expr, z, \beta).$$

    Conjunction: $\beta \in \mathcal{QB}_{\mathcal{X}}(expr_1 \text{ and } expr_2, z, Bdgs)$

    $$\overset{\text{Def}}{\Leftrightarrow} \beta \in \mathcal{QB}_{\mathcal{X}}(expr_1, z, Bdgs) \bowtie \mathcal{QB}_{\mathcal{X}}(expr_2, z, \mathcal{QB}_{\mathcal{X}}(expr_1, z, Bdgs))$$
    $$\overset{\text{IH}}{\Leftrightarrow} \text{ there are } \gamma_1 \in \mathcal{QB}_{\mathcal{X}}(expr_1, z, Bdgs)$$
    $$\text{and } \gamma_2 \in \mathcal{QB}_{\mathcal{X}}(expr_1, z, \mathcal{QB}_{\mathcal{X}}(expr_1, z, Bdgs)) \text{ s.t. } \mathcal{Q}_{\mathcal{X}}(expr_1, z, \gamma_1)$$
    $$\text{and } \gamma_1 \text{ completes some } \beta' \in Bdgs \text{ with free}(expr_1) \text{ and}$$
    $$\mathcal{Q}_{\mathcal{X}}(expr_2, z, \gamma_2) \text{ and } \gamma_1 \text{ completes some } \gamma'' \in \mathcal{QB}_{\mathcal{X}}(expr_1, z, Bdgs)$$
    $$\text{with free}(expr_2) \text{ and } \beta = \gamma_1 \bowtie \gamma_2.$$

    $$\Leftrightarrow \quad (\text{join condition: } \gamma_1 = \gamma'' \leqslant \gamma_2) \ \mathcal{Q}_{\mathcal{X}}(expr_1, z, \gamma_2) \text{ and } \mathcal{Q}_{\mathcal{X}}(expr_2, z, \gamma_2)$$
    $$\text{and } \gamma_2 \text{ completes some } \beta' \in Bdgs \text{ with free}(expr_1) \cup \text{free}(expr_2)$$
    $$\Leftrightarrow \quad \mathcal{Q}_{\mathcal{X}}(expr_1 \text{ and } expr_2, z, \gamma_2)$$
    $$\text{and } \gamma_2 \text{ completes some } \beta' \in Bdgs \text{ with free}(expr_1) \cup \text{free}(expr_2).$$

11. – 13.: trivial. (safety for variables; functions similar to predicates).

14. Context-related functions use the extension of variable bindings by pseudo-
    variables *Size* and *Pos* in rule (6):

    $$(x, \beta) \in \mathcal{SB}_{\mathcal{X}}^{any}(\text{position}(), z, Bdgs)$$
    $$\overset{\text{Def}}{\Leftrightarrow} \quad (x, \beta) \in \text{list}_{\beta \in Bdgs}(\beta(Pos), \{\beta' \in Bdgs \mid \beta(Pos) = \beta'(Pos)\})$$
    $$\Leftrightarrow \quad \beta(Pos) = x \text{ for some } \beta \in Bdgs$$
    $$\Leftrightarrow \quad x \in \mathcal{S}_{\mathcal{X}}^{any}(\text{position}(), z, \beta) \text{ for some } \beta \in Bdgs$$
    $$\Leftrightarrow \quad x \in \mathcal{S}_{\mathcal{X}}^{any}(\text{position}(), z, \beta)$$
    $$\text{and } \beta \text{ completes some } \beta' \in Bdgs \text{ by free}(\text{position}()) \text{ (which is empty).}$$

    Analogously for last().

**Proof** of Lemma 7: Structural induction.

- entry case (using $\beta = \beta'$): $(\mathcal{X}, \beta) \models /p \overset{\text{Def.5}}{\Leftrightarrow} (\mathcal{S}_{\mathcal{X}}(/p, \beta)) \neq \emptyset$

  $$\overset{\text{Def.4}}{\Leftrightarrow} (\mathcal{S}_{\mathcal{X}}(p, root, \beta)) \neq \emptyset \overset{\text{Def.4}}{\Leftrightarrow} (\mathcal{S}_{\mathcal{X}}(root/p, \beta)) \neq \emptyset \overset{\text{Def.5}}{\Leftrightarrow} (\mathcal{X}, \beta) \models root/p$$
  $$\overset{\text{IH}}{\Leftrightarrow} (\mathcal{X}, \beta) \models \text{atomize}(root/p) \overset{\text{Def}}{\Leftrightarrow} (\mathcal{X}, \beta) \models \text{atomize}(/p).$$

- Paths are resolved into steps and step qualifiers are isolated (the case where
  a don't care variable is introduced is shown; w.l.o.g., *path* is an absolute path

expression)

$(\mathscr{X}, \beta) \models path/axis :: nodetest[stepQualifier] \ /remainder$

$\Leftrightarrow \mathscr{S}_\mathscr{X}(path/axis :: nodetest[stepQualifier] \ /remainder, \beta) \neq \emptyset$

$\Leftrightarrow \mathscr{S}_\mathscr{X}(path/axis :: nodetest[stepQualifier] \ /remainder, root, \beta) \neq \emptyset$

$\Leftrightarrow \mathsf{concat}_{y \in \mathscr{S}_\mathscr{X}^a(path/axis::nodetest[stepQualifier], root, \beta)}(\mathscr{S}_\mathscr{X}^{any}(remainder, y, \beta)) \neq \emptyset$

$\Leftrightarrow$ there is a node $v \in \mathscr{S}_\mathscr{X}^a(path/axis :: nodetest[stepQualifier], root, \beta)$
   s.t. $\mathscr{S}_\mathscr{X}^{any}(remainder, v, \beta) \neq \emptyset$

$\Leftrightarrow$ there is a node $v \in \mathsf{list}_{y \in \mathscr{S}_\mathscr{X}^a(path/axis::nodetest, x, \beta)}(y \mid \mathcal{Q}_\mathscr{X}(stepQualifier, y, \beta))$
   s.t. $\mathscr{S}_\mathscr{X}^{any}(remainder, v, \beta) \neq \emptyset$

$\Leftrightarrow$ there is a node $v$ s.t. $v \in \mathscr{S}_\mathscr{X}^a(path/axis :: nodetest, x, \beta)$
   and $\mathcal{Q}_\mathscr{X}(stepQualifier, v, \beta)$ and $\mathscr{S}_\mathscr{X}^{any}(remainder, v, \beta) \neq \emptyset$

$\Leftrightarrow$ there is a node $v$ s.t. $v \in \mathscr{S}_\mathscr{X}^a(path/axis :: nodetest \to \_X, x, \beta_{\_X}^v)$
   and $\mathcal{Q}_\mathscr{X}(V[stepQualifier], v, \beta_{\_X}^v)$ and $\mathscr{S}_\mathscr{X}^{any}(V/remainder, v, \beta_{\_X}^v) \neq \emptyset$

$\Leftrightarrow$ there is a node $v$ s.t. $v \in \mathscr{S}_\mathscr{X}^a(path[axis :: nodetest \to \_X], x, \beta_{\_X}^v)$
   and $\mathcal{Q}_\mathscr{X}(V[stepQualifier], v, \beta_{\_X}^v)$ and $\mathscr{S}_\mathscr{X}^{any}(V/remainder, v, \beta_{\_X}^v) \neq \emptyset$

$\Leftrightarrow$ there is a node $v$ s.t. $\mathscr{S}_\mathscr{X}^a(path[axis :: nodetest \to \_X], x, \beta_{\_X}^v) \neq \emptyset$
   and $\mathcal{Q}_\mathscr{X}(V[stepQualifier], v, \beta_{\_X}^v)$ and $\mathscr{S}_\mathscr{X}^{any}(V/remainder, x, \beta_{\_X}^v) \neq \emptyset$

$\Leftrightarrow$ there is a node $v$ s.t. $\mathcal{Q}_\mathscr{X}(path[axis :: nodetest \to \_X], \beta_{\_X}^v)$
   and $\mathcal{Q}_\mathscr{X}(V[stepQualifier], \beta_{\_X}^v)$ and $\mathcal{Q}_\mathscr{X}(V/remainder, \beta_{\_X}^v)$


$\overset{IH}{\Leftrightarrow}$ there is a node $v$ s.t. $\mathcal{Q}_\mathscr{X}(\mathsf{atomize}(path[axis :: nodetest \to \_X]), \beta_{\_X}^v)$
   and $\mathcal{Q}_\mathscr{X}(\mathsf{atomize}(V[stepQualifier]), \beta_{\_X}^v)$ and $\mathcal{Q}_\mathscr{X}(\mathsf{atomize}(V/remainder), \beta_{\_X}^v)$

$\Leftrightarrow$ there is a node $v$ s.t. $\mathcal{Q}_\mathscr{X}(\mathsf{atomize}(\ldots), \beta_{\_X}^v)$.


- Conjunctions in step qualifiers: obvious.

- Predicates in step qualifiers: W.l.o.g., consider a unary predicate with a relative argument expression:

$(\mathscr{X}, \beta) \models V[pred(expr)]$

$\Leftrightarrow \mathscr{S}_\mathscr{X}(V[pred(expr)], \beta(V), \beta) \neq \emptyset$

$\Leftrightarrow \mathsf{list}_{y \in \mathscr{S}_\mathscr{X}^a(V, \beta(V), \beta)}(y \mid \mathcal{Q}_\mathscr{X}(pred(expr), y, \beta)) \neq \emptyset$

$\Leftrightarrow (\beta(V)$ is the only element in $\mathscr{S}_\mathscr{X}^a(V, \beta(V), \beta))$ s.t. $\mathcal{Q}_\mathscr{X}(pred(expr), \beta(V), \beta)$

$\Leftrightarrow$ there is an $x \in \mathscr{S}_\mathscr{X}(expr, \beta(V), \beta)$ such that $pred(x) \in \mathscr{X}$

$\Leftrightarrow$ there is an $x$ s.t. $x \in \mathscr{S}_\mathscr{X}(V/expr \to \_X, root, \beta_{\_X}^x)$ and $(\mathscr{X}, \beta_{\_X}^x) \models pred(\_X)$

$\Leftrightarrow$ there is an $x$ s.t. $(\mathscr{X}, \beta_{\_X}^x) \models V/expr \to \_X$ and $(\mathscr{X}, \beta_{\_X}^x) \models pred(\_X)$

$\overset{IH}{\Leftrightarrow}$ there is an $x$ s.t. $(\mathscr{X}, \beta_{\_X}^x) \models \mathsf{atomize}(V/expr \to \_X)$
   and $(\mathscr{X}, \beta_{\_X}^x) \models pred(\_X)$

$\Leftrightarrow$ there is an $x$ s.t. $(\mathscr{X}, \beta_{\_X}^x) \models \mathsf{atomize}(V[pred(expr)])$.

- Predicate atoms: analogous.

## Acknowledgments

Most of this work has been done when I was a member of the database group at Freiburg University. I want to thank my former colleagues during that time: Lule Ahmedi, Matthias Ihle, Georg Lausen, Pedro Marrón, Martin Weber, and Fang Wei.

## References

ABITEBOUL, S., QUASS, D., MCHUGH, J., WIDOM, J. AND WIENER, J. 1997. The Lorel query language for semistructured data. *Journal on Digital Libraries (JODL) 1*, 1, 68–88.

BARU, C., GUPTA, A., LUDÄSCHER, B., MARCIANO, R., PAPAKONSTANTINOU, Y., VELIKHOV, P. AND CHU, V. 1999. XML-based information mediation with MIX. *ACM Intl. Conference on Management of Data (SIGMOD)*, pp. 597–599.

BAUMGARTNER, R., FLESCA, S. AND GOTTLOB, G. 2001a. The Elog web extraction language. *Intl. Conference on Logic Programming and Automated Reasoning (LPNMR): Lecture Notes in Computer Science 2250*, pp. 548–560. Springer.

BAUMGARTNER, R., FLESCA, S. AND GOTTLOB, G. 2001b. Visual web information extraction with Lixto. *Intl. Conference on Very Large Data Bases (VLDB)*, pp. 119–128.

BRY, F. AND SCHAFFERT, S. 2002. Towards a declarative query and transformation language for XML and semistructured data: Simulation unification. *Intl. Conference on Logic Programming (ICLP)*, pp. 255–270.

BUNEMAN, P., DAVIDSON, S., HILLEBRANDT, G. AND SUCIU, D. 1996. A query language and optimization techniques for unstructured data. *ACM Intl. Conference on Management of Data (SIGMOD)*, pp. 505–516. Montreal, Canada.

BUNEMAN, P., FERNANDEZ, M. AND SUCIU, D. 2000. UnQL: A query language and algebra for semistructured data based on structural recursion. *VLDB Journal 9*, 76–110.

CERI, S., COMAI, S., DAMIANI, E., FRATERNALI, P., PARABOSCHI, S. AND TANCA, L. 1999. XML-GL: a graphical language for querying and restructuring XML documents. *Proceedings 8th International World Wide Web Conference (WWW 8)*, pp. 1171–1187.

CHRISTOPHIDES, V., CLUET, S. AND SIMÉON, J. 2000. On wrapping query languages and efficient XML integration. *ACM Intl. Conference on Management of Data (SIGMOD)*, pp. 141–152.

CLARK, J. 1998. XT: an implementation of XSL Transformations. `http://www.jclark.com/xml/xt.html`.

CLUET, S., DELOBEL, C., SIMÉON, J. AND SMAGA, K. 1999. Your mediators need data conversion. *ACM Intl. Conference on Management of Data (SIGMOD)*, pp. 177–188.

COMAI, S., DAMIANI, E. AND FRATERNALI, P. 2001. Computing graphical queries over XML data. *ACM Transactions on Information Systems (TOIS) 19*, 4, 371–430.

CONSENS, M. AND MENDELZON, A. 1990. GraphLog: a visual formalism for real life recursion. *ACM Symposium on Principles of Database Systems (PODS)*, pp. 404–416.

DEUTSCH, A., FERNANDEZ, M., FLORESCU, D., LEVY, A. AND SUCIU, D. 1999. XML-QL: A Query Language for XML. *8th. WWW Conference*. W3C. World Wide Web Consortium Technical Report, NOTE-xml-ql-19980819, `www.w3.org/TR/NOTE-xml-ql`.

DEUTSCH, A., FERNANDEZ, M. AND SUCIU, D. 2000. Storing semistructured data with STORED. *ACM Intl. Conference on Management of Data (SIGMOD)*, pp. 431–442.

DOM-W3C. 1998. Document object model (DOM). `http://www.w3.org/DOM/`.

FERNANDEZ, M., FLORESCU, D., LEVY, A. AND SUCIU, D. 1997. A query language for a web-site management system. *SIGMOD Record 26*, 3, 4–11.

FERNANDEZ, M., SIMÉON, J. AND WADLER, P. 1999. XML query languages: Experiences and exemplars. draft manuscript, communication to the XML Query W3C Working Group. `http://www-db.research.bell-labs.com/user/simeon/xquery.ps`.

FERNANDEZ, M. F., FLORESCU, D., KANG, J., LEVY, A. Y. AND SUCIU, D. 1998. Catching the boat with Strudel: Experiences with a web-site management system. *ACM Intl. Conference on Management of Data (SIGMOD)*, pp. 414–425.

FLORESCU, D. AND KOSSMANN, D. 1999. A performance evaluation of alternative mapping schemes for storing XML data in a relational database. Technical Report 3684, INRIA.

FLORID. 1998. FLORID homepage. `http://www.informatik.uni-freiburg.de/˜dbis/florid/`.

FROHN, J. 1998. Magic-Set Transformation in deduktiven, objektorientierten Datenbanksprachen. PhD thesis, Institut für Informatik, Universität Freiburg.

FROHN, J., LAUSEN, G. AND UPHOFF, H. 1994. Access to objects by path expressions and rules. *Intl. Conference on Very Large Data Bases (VLDB)*, pp. 273–284.

GARCIA-MOLINA, H., PAPAKONSTANTINOU, Y., QUASS, D., RAJARAMAN, A., SAGIV, Y., ULLMAN, J., VASSALOS, V. AND WIDOM, J. 1997. The TSIMMIS approach to mediation: Data models and languages. *Journal of Intelligent Information Systems 8*, 2, 117–132.

GOLDMAN, R., MCHUGH, J. AND WIDOM, J. 1999. From semistructured data to XML: Migrating the Lore data model and query language. *WebDB 1999*, pp. 25–30.

HOSOYA, H. AND PIERCE, B. C. 2000. Xduce: A typed XML processing language. *WebDB 2000*, pp. 111–116.

KIFER, M. AND LAUSEN, G. 1989. F-Logic: A higher-order language for reasoning about objects, inheritance and scheme. *ACM Intl. Conference on Management of Data (SIGMOD)*, pp. 134–146.

KIFER, M., LAUSEN, G. AND WU, J. 1995. Logical foundations of object-oriented and frame-based languages. *Journal of the ACM 42*, 4, 741–843.

LAKSHMANAN, L. V. S., SADRI, F. AND SUBRAMANIAN, I. N. 1996. SchemaSQL – a language for interoperability in relational multi-database systems. *Intl. Conference on Very Large Data Bases (VLDB)*, pp. 239–250.

LEHTI, P. 2001. Design and implementation of a data manipulation processor for an XML query language. MS thesis, Technische Universität Darmstadt.

LUDÄSCHER, B., HIMMERÖDER, R., LAUSEN, G., MAY, W. AND SCHLEPPHORST, C. 1998. Managing semistructured data with FLORID: A deductive object-oriented perspective. *Information Systems 23*, 8, 589–612.

MAY, W. 1999. Information extraction and integration with FLORID: The MONDIAL case study. Technical Report 131, Universität Freiburg, Institut für Informatik. Available from `http://www.dbis.informatik.uni-goettingen.de/Mondial/`.

MAY, W. 2001a. Habilitation thesis. PhD thesis, Universität Freiburg. Available from `http://www.dbis.informatik.uni-goettingen.de/lopix/`.

MAY, W. 2001b. Information integration in XML: The MONDIAL case study. Technical Report. Available from `http://www.dbis.informatik.uni-goettingen.de/lopix/lopix-mondial.html`.

MAY, W. 2001c. LoPiX: A system for XML data integration and manipulation. *Intl. Conference on Very Large Data Bases (VLDB), Demonstration Track*, pp. 707–708.

MAY, W. 2001d. The LOPIX system. `http://www.dbis.informatik.uni-goettingen.de/lopix/`.

MAY, W. 2001e. The MONDIAL database. `http://www.dbis.informatik.uni-goettingen.de/Mondial/`.

MAY, W. 2002. A rule-based querying and updating language for XML. *Workshop on Databases and Programming Languages (DBPL 2001): Lecture Notes in Computer Science 2397*, pp. 165–181.

MAY, W. AND BEHRENDS, E. 2001. On an XML data model for data integration. In *Intl. Workshop on Foundations of Models and Languages for Data and Objects (FMLDO 2001)*.

MCHUGH, J., ABITEBOUL, S., GOLDMAN, R., QUASS, D. AND WIDOM, J. 1997. Lore: A database management system for semistructured data. *SIGMOD Record 26*, 3, 54–66.

PRZYMUSINSKI, T. C. 1988. On the declarative semantics of deductive databases and logic programs. In: J. Minker, Ed. *Foundations of Deductive Databases and Logic Programming*, pp. 191–216. Morgan Kaufmann.

ROBIE, J. 1999. XQL (XML Query Language). `http://www.metalab.unc.edu/xql/xql-proposal.html`.

SHANMUGASUNDARAM, J., GANG, H., TUFTE, K., ZHANG, C., WITT, D. J. D. AND NAUGHTON, J. Relational databases for querying XML documents: Limitations and opportunities. *Intl. Conference on Very Large Data Bases (VLDB)*, pp. 302–314.

SOFTWARE AG. 2001. Quip: An xquery implementation. `http://www.softwareag.com/developer/quip/`.

TATARINOV, I., IVES, Z. G., HALEVY, A., AND WELD, D. 2001. Updating xml. *ACM Intl. Conference on Management of Data (SIGMOD)*, pp. 133–154.

WADLER, P. 1999. Two semantics for XPath. `http://www.cs.bell-labs.com/who/wadler/topics/xml.html`.

XMLInf. 1999. XML information set. `http://www.w3.org/TR/XML-infoset`.

XMQ-A. 2001. XML Query Algebra. `http://www.w3.org/TR/query-algebra`.

XMQ-D. 2001. XML Query Data Model. `http://www.w3.org/TR/query-datamodel`.

XPath. 1999. XML Path Language (XPath) version 1.0: 1999. `http://www.w3.org/TR/xpath`.

XPQOF. 2001. XQuery 1.0 and XPath 2.0 Functions and Operators. `http://www.w3.org/TR/xquery-operators`.

XQFS. 2001. XQuery 1.0 Formal Semantics. `http://www.w3.org/TR/query-semantics`.

XQuery. 2001. XQuery: A Query Language for XML. `http://www.w3.org/TR/xquery`.

XSLT. 1999. XSL Transformations (XSLT). `http://www.w3.org/TR/xslt`.