

# XPath on Left and Right Sides of Rules: Toward Compact XML Tree Rewriting through Node Patterns

Jean-Yves Vion-Dury

► **To cite this version:**

Jean-Yves Vion-Dury. XPath on Left and Right Sides of Rules: Toward Compact XML Tree Rewriting through Node Patterns. Proceedings of the 2003 ACM Symposium on Document Engineering, DocEng2003, Nov 2003, Grenoble, France. pp.19-25, 2003, <10.1145/958220.958225>. <inria-00423398>

**HAL Id: inria-00423398**

**<https://hal.inria.fr/inria-00423398>**

Submitted on 9 Oct 2009

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# XPath on Left and Right Sides of Rules: Toward Compact XML Tree Rewriting through Node Patterns

Jean-Yves Vion-Dury\*

WAM project[17]

INRIA Rhône-Alpes  
655 Avenue de l'Europe  
Montbonnot, France

Jean-Yves.Vion-Dury@{inria.fr,xrce.xerox.com}

## ABSTRACT

*XPath [3, 5] is a powerful and quite successful language able to perform complex node selection in trees through compact specifications. As such, it plays a growing role in many areas ranging from schema specifications, designation and transformation languages to XML query languages. Moreover, researchers have proposed elegant and tractable formal semantics [8, 9, 10, 14], fostering various works on mathematical properties and theoretical tools [10, 13, 12, 14]. We propose here a novel way to consider XPath, not only for selecting nodes, but also for tree rewriting using rules. In the rule semantics we explore, XPath expressions (noted  $p, p'$ ) are used both on the left and on the right side (i.e. rules have the form  $p \rightarrow p'$ ). We believe that this proposal opens new perspectives toward building highly concise XML transformation languages on widely accepted basis.*

## Categories and Subject Descriptors

D.3.1 [Programming Languages]: Formal Definitions and Theory—*document transformation, XPath, tree rewriting*

## General Terms

Language Theory Design

## 1. INTRODUCTION

Rewriting is a deeply studied branch of theoretical computer science, and influenced an impressive amount of research work, ranging from the general computation theory up to practical applications. We assume in this paper that the reader owns some basic knowledge on rewriting theory

\*Visiting researcher from Xerox Research Centre Europe.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

DocEng '03, November 20–22, 2003, Grenoble, France.  
Copyright 2003 ACM 1-58113-724-9/03/0011 ...\$5.00.

and also on XPath (see [1, 3] for an exhaustive introduction) As a computation model, it has some good properties that attracted the attention of scientists looking for structure transformation languages : decomposition in small basic units (rules), structural matching/filtering, modularity, and static analysis through unification (see for instance termination proofs based on critical peaks analysis, itself based on rule's terms unification [1]).

Rewriting has been applied to string and tree transformation in compiler construction, and also to more complex data structures such as (labelled) graph transformation. However, it has been poorly applied to XML transformation, partly because the richness of the XML tree model require quite complex tree pattern matching specifications<sup>1</sup>. In particular, conventional tree patterns requires specifying not only the interesting sub-structures, but also the ones which have to be ignored during the matching.

XPath in contrast is a node selection language that brings a great simplicity by focusing on paths instead of structures, making powerful hypothesis and metaphors such as implicit context node and navigation. It is today a widely known specification, adopted by many programmers and also recognized as being mathematically tractable [8, 10, 13, 12, 14].

This paper is based on the idea of reusing XPath in systems of rules like  $p \rightarrow p'$ , where  $p$  and  $p'$  are XPath like expressions. In order to build such systems, one has to address two problems:

1. how should XPath be extended in order to embed the filtering variables required for memorizing and reusing sub-structures ?
2. how should an XPath expression be interpreted on the right hand side ( $p'$ ) with respect to tree construction and modification?

Wolfgang May [16] proposed to extend XPath toward a logic based approach, and to describe syntactic and semantic modifications that enable filtering through unification variables. The proposed approach somehow solves the two problems above, but requires a deep modification in the orig-

<sup>1</sup>and also partly because many XML transformations require complex rule application schemes

inal syntax and semantics. Moreover, the latter is deeply related to logic and prolog-like execution model, which makes harder to understand the relationship with the pure XPath semantics and the structure of the computation. In addition, this language addresses only monotonic transformations of trees (nodes and links are never deleted along transformations).

Section 2 presents a formal definition of the syntax and semantics of XPath, inspired from [14] and slightly extended toward handling positional information and node value tests. Then we introduce the general notion of path containment, and give a mathematical definition of this property.

The notion of pattern and rule semantics is defined in section 3, and related to XPath semantics and containment through a translation function  $\Xi$  that transforms an XPath pattern  $\bar{p}$  into an XPath expression  $p$ . Thanks to this approach, the pattern matching operation is nothing else than finding an environment in which  $p$  will select the context node itself. A small and simple set of operations for tree modification is proposed in order to define the semantics of rule application, i.e. the meaning of the right hand side.

We illustrate the expressiveness of the approach in section 4, with realistic examples from MathML [6].

We conclude through a short synthesis and some perspectives on our future work on the topic.

## 2. XPATH: A FORMAL DEFINITION AND SOME PROPERTIES

The purpose of this section is to describe formally a language derived from a fragment of XPath 2.0. We believe that this fragment is expressive enough for most “real world” applications. However, if it deserves today our research purpose, it is planned to extend it step by step as long as we are able to maintain its mathematical properties.

### 2.1 Syntax

The syntax is defined in two stages: (i) the core language, and (ii) the syntactic sugar, which just rewrites into expressions of the core. The following definition is a slight extension of [14]

$$p ::= \wedge \mid \perp \mid p \mid p \mid a::N \mid (p) \mid p/p \mid p[q] \mid \$v \mid \text{for } \$v \text{ in } p \text{ return } p$$

Note that according to XPath 2.0 specification, and to proposal [10], but not in accordance with the proposal of [8], the axes  $a::N$  cannot contain other paths (e.g.  $a::(p_1|p_2)$ ). A node test  $N$ , is either an unqualified name  $n$ , the wildcard symbol  $*$  or a test function among **text()**, **node()**, **comment()**, **processing-instruction()**, **element()**; when  $N$  takes the  $ns:n$  form,  $ns$  is considered as a namespace prefix in accordance to the specification, and processed accordingly. More precisely, it is rewritten into a namespace attribute (see rule  $r_{5f}$  in figure 1).

The important extension we propose with respect to qualifiers is the node set inclusion constraint  $p_1 \sqsubseteq p_2$ , which brings extra expressive power and interesting possibilities for containment inference (see [14]). The inclusion also brings additional facilities for discriminating embedded structures and for extending the standard node set equality test in an

.	→	<i>self</i> :: <b>node</b> ()	( $r_{5a}$ )
$p_1//p_2$	→	$p_1/\text{desc-or-self}::*/p_2$	( $r_{5c}$ )
$N$	→	<i>child</i> :: $N$	( $r_{5d}$ )
$@N$	→	<i>attribute</i> :: $N$	( $r_{5e}$ )
$ns:n$	→	$n[\text{namespace}::ns]$	( $r_{5f}$ )
$p[?v]$	→	$p[. == \$v]$	( $r_{5h}$ )
$p[i]$	→	$p[\text{position}() = i]$	( $r_{5i}$ )
$p[-i]$	→	$p[\text{position}() = -i]$	( $r_{5i}$ )

Figure 1: Syntactic sugars ( $\mathcal{N}_5$ )

elegant way. For instance

$$q ::= \text{true} \mid \text{false} \mid (q) \mid \text{not } q \mid q \text{ and } q \mid q \text{ or } q \mid p \sqsubseteq p$$

However, we propose an extension of qualifier syntax so that minimal positional information can be expressed ( $i$  is the lexical representation for strictly positive integers), and also information on linearized content ( $s$  is the lexical representation for strings).

$$q ::= \text{position}() = i \mid \text{position}() = -i \mid p = s$$

So, expressions like  $a[2]/b[@att = \text{”test”}]^2$  are legal. A (partial) reflexive ordering of  $N$  elements is defined in figure 2.

$$\begin{array}{l} n \leq * \leq \text{node}() \\ \text{processing-instruction}(), \text{text}() \leq \text{node}() \\ \text{element}(), \text{comment}() \leq \text{node}() \end{array}$$

Figure 2: Partial ordering of node tests

The  $a$  symbol denotes *axes*, ranging over the whole set defined in the W3C specification

$$\begin{array}{l} a \in \{ \text{self}, \text{attribute}, \text{namespace} \} \\ \cup \{ \text{child}, \text{parent} \} \\ \cup \{ \text{descendant}, \text{ancestor} \} \\ \cup \{ \text{descendant-or-self}, \text{ancestor-or-self} \} \\ \cup \{ \text{following}, \text{preceding}, \} \\ \cup \{ \text{preceding-sibling}, \text{following-sibling} \} \end{array}$$

We define the syntactic sugar (see figure 1 for their translation into core expressions)

$$p ::= p//p \mid . \mid N \mid @N \mid ns:N \mid p[?v] \mid p[i] \mid p[-i]$$

We have to comment right now on the  $p[?v]$  notation, which will be our mechanism to introduce filtering variables in XPath expressions (see coming section 3.1). More precisely, it rewrites into  $p[. == \$v]$ , where  $\$v$  is a free variable. Following the semantics of the  $==$  node set operator (see section 2.2), it means that the matching requires the existence of a substitution assigning a suitable node value to the variable. It must thus satisfy set inclusion constraints.

We propose syntactic sugars for qualifier (see figure 3):

$$q ::= p \mid p == p \mid p \neq p \mid p \sqsubseteq p \mid p \sqsupset p \mid p \sqsupseteq p$$

<sup>2</sup>equivalent to  $a[\text{position}() = 2]/b[\text{attribute}::att = \text{”test”}]$  after syntactic sugar expansion

$p$	$\rightarrow$	<b>not</b> ( $p \sqsubseteq \perp$ )	( $r_{6a}$ )
$p_1 == p_2$	$\rightarrow$	( $p_1 \sqsubseteq p_2$ ) <b>and</b> ( $p_2 \sqsubseteq p_1$ )	( $r_{6b}$ )
$p_1 \sqsubset p_2$	$\rightarrow$	( $p_1 \sqsubseteq p_2$ ) <b>and not</b> ( $p_2 \sqsubseteq p_1$ )	( $r_{6c}$ )
$p_1 \sqsupset p_2$	$\rightarrow$	$p_2 \sqsubset p_1$	( $r_{6d}$ )
$p_1 \sqsupseteq p_2$	$\rightarrow$	$p_2 \sqsubseteq p_1$	( $r_{6e}$ )
$p_1 \neq p_2$	$\rightarrow$	<b>not</b> ( $p_1 == p_2$ )	( $r_{6f}$ )

**Figure 3: Syntactic sugars (qualifiers) ( $\mathcal{N}_6$ )**

We consider that the / and | operators are fully associative, and that the precedence ordering is (from the tightest to the loosest):

$$a::N < p[q] < p/p < \text{for } \$v \text{ in } p \text{ return } p < p|p$$

so that for instance

$$\text{child}::n[q]/\text{child}::n_2[q'] \mid \text{descendant}::*$$

is syntactically understood as

$$(\text{child}::n[q]) / (\text{child}::n_2[q']) \mid (\text{descendant}::*)$$

The precedence of boolean operators is the standard one (**not** < **and** < **or**).

## 2.2 Semantics

We reuse the denotational definition of [14], also inspired from [8, 10], extended with a node set parameter  $S$  and with an execution context  $\phi$  that uniquely associates variable names to tree nodes.

### 2.2.1 Document model.

The XPath semantics relies on a document-as-a-tree model. A linear (“flat”) document is seen as a well-formed tree after successful parsing. A tree is modeled as a set of “typed” nodes (*element*, *text*, *comment*, *processing-instruction*, *attribute*, *namespace*, and *root*); the type of a node can be checked by a corresponding unary predicate).

A well-formed tree contains only one root node, which has no parent, no attribute and no namespace but may have any other kind of nodes as children. Moreover, only elements can have children. Nodes, identified by  $x$  and  $x_i$  in the sequel, are fully connected in order to form a tree<sup>3</sup>. This structural property relies on the *parent/child* relation  $\rightarrow$  that characterizes the tree, and also its transitive closure  $\rightarrow^+$ . Moreover, a strict ordering, the document ordering  $\ll$ , is defined on every node  $x$  of a tree  $t$ , and reflect the order of opening tags occurrence in the linear document. The ordering relation  $x_1 \ll x_2$  is **true** if the opening tag of  $x_1$  appears before the opening tag of  $x_2$ , **false** otherwise.

### 2.2.2 Semantics of selection.

The selection is defined relatively to a context node  $x$ , and the execution context  $\phi$ . If  $\phi = \{\dots, v = x, \dots\}$  then  $\phi(v) = x$ . The function  $\mathcal{S}$  is inductively defined, i.e. uses itself for its own definition. The induction is even double, since it relies on the  $\mathcal{Q}$  function (defines qualifiers, presented in the next paragraph) which itself uses  $\mathcal{S}$ . This is however quite common in denotational semantics, and just reflects

<sup>3</sup>Every node  $x$  of the tree is reachable from the root ; mathematically,  $\wedge \rightarrow^+ x$

faithfully the syntactic structure of the XPath language.

$$\mathcal{S} : \text{Pattern} \longrightarrow \text{Env} \longrightarrow \text{Node} \longrightarrow \text{Set}(\text{Node}) \longrightarrow \text{Set}(\text{Node})$$

$$\begin{aligned} \mathcal{S}[\perp]_x^{\phi, S} &= \emptyset \\ \mathcal{S}[(p)]_x^{\phi, S} &= \mathcal{S}[p]_x^{\phi, S} \\ \mathcal{S}[p_1 p_2]_x^{\phi, S} &= \mathcal{S}[p_1]_x^{\phi, S} \cup \mathcal{S}[p_2]_x^{\phi, S} \\ \mathcal{S}[p_1 / p_2]_x^{\phi, S} &= \{x_2 \mid x_1 \in \mathcal{S}[p_1]_x^{\phi, S} \wedge x_2 \in \mathcal{S}[p_2]_{x_1}^{\phi, S}\} \\ \mathcal{S}[p[q]]_x^{\phi, S} &= \text{let } S_1 = \mathcal{S}[p]_x^{\phi, S} \text{ in} \\ &\quad \{x_1 \mid x_1 \in s_1 \wedge \mathcal{Q}[q]_{x_1}^{\phi, S_1}\} \\ \mathcal{S}[\$v]_x^{\phi, S} &= \{\phi(v)\} \\ \mathcal{S}[a::N]_x^{\phi, S} &= \{x_1 \mid f_a(x) \wedge \mathcal{T}_a(x_1, N)\} \end{aligned}$$

The **for** expression is the only one that introduces a new (variable, node value) pair into the context

$$\begin{aligned} \mathcal{S}[\text{for } \$v \text{ in } p_1 \text{ return } p_2]_x^{\phi, S} \\ = \{x_2 \mid x_1 \in \mathcal{S}[p_1]_x^{\phi, S} \wedge x_2 \in \mathcal{S}[p_2]_{x_1}^{\phi, v=x_1}\} \end{aligned}$$

The notation  $f_a$  corresponds to the functions defined in the table below.

axis ( $a$ )	$f_a(x)$
<i>self</i>	$\{x\}$
<i>child</i>	$\{x_1 \mid x \rightarrow x_1\}$
<i>parent</i>	$\{x_1 \mid x_1 \rightarrow x\}$
<i>descendant</i>	$\{x_1 \mid x \rightarrow^+ x_1\}$
<i>ancestor</i>	$\{x_1 \mid x_1 \rightarrow^+ x\}$
<i>descendant-or-self</i>	$\{x\} \cup \{x_1 \mid x \rightarrow^+ x_1\}$
<i>ancestor-or-self</i>	$\{x\} \cup \{x_1 \mid x_1 \rightarrow^+ x\}$
<i>following-sibling</i>	$\text{sibling}(x) \cap \text{following}(x)$
<i>preceding-sibling</i>	$\text{sibling}(x) \cap \text{preceding}(x)$
<i>preceding</i>	$\{x_1 \mid x_1 \ll x\}$
<i>following</i>	$\{x_1 \mid x \ll x_1\}$
<i>attribute</i>	$\{x_1 \mid x \rightarrow x_1 \wedge \text{attribute}(x_1)\}$
<i>namespace</i>	$\{x_1 \mid x \rightarrow x_1 \wedge \text{namespace}(x_1)\}$
<i>sibling</i>	$\{x_2 \mid x_1 \rightarrow x \wedge x_1 \rightarrow x_2\}$

The  $\mathcal{T}$  function performs a node test, according to the table below

N	$a$	$\mathcal{T}_a(N, x)$
$n$		name(x)=n
*	<i>attribute</i>	attribute(x)
*	<i>namespace</i>	namespace(x)
*	other	element(x)
<b>text()</b>		text(x)
<b>comment()</b>		comment(x)
<b>processing-instruction()</b>		pi(x)
<b>element()</b>		element(x)
<b>node()</b>		<b>true</b>

### 2.2.3 Semantics of qualifiers.

The originality here is in the inclusion test (last line), directly expressed through a set-theoretic inclusion of selection sets.

$$\begin{aligned} \mathcal{Q} : \text{Qualifier} \longrightarrow \text{Env} \longrightarrow \text{Node} \longrightarrow \text{Boolean} \\ \mathcal{Q}[\text{true}]_x^{\phi, S} &= \text{true} \\ \mathcal{Q}[\text{false}]_x^{\phi, S} &= \text{false} \\ \mathcal{Q}[q_1 \text{ and } q_2]_x^{\phi, S} &= \mathcal{Q}[q_1]_x^{\phi, S} \wedge \mathcal{Q}[q_2]_x^{\phi, S} \\ \mathcal{Q}[q_1 \text{ or } q_2]_x^{\phi, S} &= \mathcal{Q}[q_1]_x^{\phi, S} \vee \mathcal{Q}[q_2]_x^{\phi, S} \\ \mathcal{Q}[(q)]_x^{\phi, S} &= \mathcal{Q}[q]_x^{\phi, S} \\ \mathcal{Q}[\text{not } q]_x^{\phi, S} &= \neg \mathcal{Q}[q]_x^{\phi, S} \\ \mathcal{Q}[p_1 \sqsubseteq p_2]_x^{\phi, S} &= \mathcal{S}[p_1]_x^{\phi, S} \subseteq \mathcal{S}[p_2]_x^{\phi, S} \end{aligned}$$

The former definitions correspond exactly to those of [14], up to the function signature. We now define the positional

and linearization based extensions

$$\begin{aligned} \mathcal{Q}[\mathbf{position}()] &= i \llbracket_x^{\phi, S} &= \text{let } S_1 = \{x_1 \in S \mid x_1 \ll x\} \\ & & \text{in } (\text{size}(S_1) = i - 1) \\ \mathcal{Q}[\mathbf{position}()] &= -i \llbracket_x^{\phi, S} &= \text{let } S_1 = \{x_1 \in S \mid x_1 \ll x\} \\ & & \text{in } (\text{size}(S_1) = \text{size}(S) - i) \\ \mathcal{Q}[p = s] &= s \llbracket_x^{\phi, S} &= \text{let } S_1 = \mathcal{S}[p] \llbracket_x^{\phi, S} \text{ in} \\ & & (S_1 \neq \emptyset) \wedge (\mathbf{value-of}(S_1) = s) \end{aligned}$$

Note also that string value tests ( $p = s$ ) require non empty path selection. The **value-of** function returns the concatenation of all text leaves descendant of each node in the selection, and moreover, in document order.

### 2.3 Containment of XPath expressions

The containment relation  $p \leq p'$  expresses that for all context node  $x$ , the node set resulting from the evaluation of  $p$  is included in the node set selected by  $p'$ . More precisely, this property is defined by

$$p \leq p' \text{ iff } \forall x \quad \mathcal{S}[p] \llbracket_x^{\emptyset, \emptyset} \subseteq \mathcal{S}[p'] \llbracket_x^{\emptyset, \emptyset}$$

The containment inference (see [14]) requires actually a context  $\gamma$  that conveys information about variables (it may be considered as empty, and just ignored as above).

For instance (example taken from [14]), one can prove that  $\$v/b \leq (a|c)/b$  in a context  $\gamma = \{v:c|a\}$

$$\frac{\frac{\gamma \vdash c|a \leq a|c \quad [t_1]}{\gamma \vdash \$v \leq a|c} \quad [d_{3a}] \quad \frac{}{\gamma \vdash b \leq b} \quad [c_2]}{\gamma \vdash \$v/b \leq (a|c)/b} \quad [d_2]$$

In order to express the containment property in a quite general way, the static environment  $\gamma$  and the execution environment  $\phi$  must be considered in conformance ( $\gamma \Vdash \phi$ )

PROPERTY 1. *Conformance of an execution environment*

$$\gamma \Vdash \phi \text{ iff } \forall v \in \gamma, \exists x, S \quad \phi(v) \subseteq \mathcal{S}[\gamma(v)] \llbracket_x^{\phi, S}$$

Thus a containment assertion  $\gamma \vdash p \leq p'$  is expressed through the following general property

PROPERTY 2. *Containment.*

$$\gamma \vdash p \leq p' \text{ iff } \forall x, \phi, S \quad \gamma \Vdash \phi \implies \mathcal{S}[p] \llbracket_x^{\phi, S} \subseteq \mathcal{S}[p'] \llbracket_x^{\phi, S}$$

It turns out to be a quite exciting and fundamental problem, and also difficult from the computational point of view (see [13, 12] for model-based approaches, and our later work [14] for a syntactic-based approach). We propose to use this relation in order to define the semantics of rules, and also to establish later some useful structural properties.

## 3. PATTERNS AND RULES

Rules have the form  $\vec{p} \rightarrow \vec{p}'$ , where  $\vec{p}$  denotes an XPath pattern, i.e. an expression that contains free variables that must be instantiated in order the pattern to match. The matching function  $\mathcal{M}$  returns a context  $\phi'$  that defines a value to each free variable of the pattern.

A single rule is evaluated in a context including the so-called “contextual node” (which defines the notion of relative location in the tree, just as in the standard XPath semantics), and an execution environment  $\phi$  that uniquely defines the value of bound variables.

The left hand side is evaluated in the current node context, and if a matching is found, the tree is transformed by executing a sequence of operations, such as defined by an interpretation function  $\Phi$ .

$$\mathcal{R} \quad : \quad \text{Rule} \rightarrow \text{Env} \rightarrow \text{Node} \rightarrow \text{Bool}$$

$$\mathcal{R}[\vec{p} \rightarrow \vec{p}'] \llbracket_x^{\phi} = \begin{cases} \text{let } \phi' = \mathcal{M}[\vec{p}] \llbracket_x^{\phi} \text{ in} \\ \text{if } \phi' \neq \emptyset \text{ then } (\mathbf{do } \Phi(\phi', p, p'); \mathbf{true}) \\ \text{else } \mathbf{false} \end{cases}$$

We now define the matching, the tree rewriting operations and their construction by function  $\Phi$

### 3.1 Matching and filtering

Wadler first defined XPath matching formally [8] as a boolean function, true if the selection is not empty. Actually we need more than matching, since rewriting tree operations require node identification (*filtering*). Indeed, we have to designate the nodes we want to remove or to bind to other nodes.

Our idea is to mark the various nodes involved in the selection path with a designation tag, i.e. a  $p[?v]$  qualifier, which is a sugar for  $p[. = \$v]$ ,  $v$  being an “existential variable”, not defined in the current  $\phi$  environment. For instance, if the following expression

EXAMPLE 1.

$$\vec{P} = \text{table/tr}[?a]/\text{td}[?b][2]$$

matches in environment  $\phi$  and context node  $x$ , a new environment (a substitution)  $\phi'$  will be established. In that case, the corresponding *tr* and *td* element identifiers  $x_1, x_2$  will be designated by variables  $\$a$  and  $\$b$ :

$$\phi' = \{a = x_1, b = x_2\}$$

We now propose to define the matching function through a translation  $\Xi$  of  $p$  terms into a set of containment constraints. A matching is found when a solution to the system is found. The terms of this equation have a selection semantics in conformance with the definition of section 2.2. The small example 1 above is translated into

$$\Xi(\vec{P}) = \begin{cases} \$a \leq \text{table/tr} \\ \$b \leq \$a/\text{td}[\mathbf{position}() = 2] \end{cases}$$

Another way to present this translation is to generate a unique equivalent XPath term:

$$\Xi(\vec{P}) = .[\$a \sqsubseteq \text{table/tr}][\$b \sqsubseteq \$a/\text{td}[\mathbf{position}() = 2]]$$

which should select the context node in a suitable environment  $\phi'$ . This translation<sup>4</sup> allows us to define the matching semantics in a quite elegant way, since it relates directly this operation to the fundamental XPath selection mechanism.

DEFINITION 1. *Matching of a rule pattern.*

$$\mathcal{M} \quad : \quad \text{Env} \rightarrow \text{Pattern} \rightarrow \text{Node} \rightarrow \text{Env}$$

$$\mathcal{M}[\vec{p}] \llbracket_x^{\phi} = \begin{cases} \text{if } \exists \phi' \mid \mathcal{S}[\Xi(\vec{p})] \llbracket_x^{\phi'} = \{x\} \\ \text{then } \phi' \text{ else } \emptyset \end{cases}$$

<sup>4</sup>We consider that the definition of  $\Xi$  doesn't present technical difficulties, and will not be detailed here.

Many solutions may exist, and moreover, they can be totally ordered through  $\ll$ , the document order. In order to apply a matching rule, it makes sense to choose the first solution as the default one, since it eliminates non-determinism.

### 3.2 Operations on the document tree

We propose now a restricted set of primitive operations used for updating the tree. We define informally the semantics of these operations in the figure 4, through an “object-oriented” style,  $x$  being the node on which the methods are invoked (we could map these definitions into DOM [4] operations, but it would be probably less compact). Note that

operation	arguments	meaning
$x_i = \text{element}()$		create element $x_i$
$x_i = \text{attribute}()$		create attribute $x_i$
$x_i = \text{text}(s)$	content	create text $x_i$
$x_i = \text{comment}(s)$	content	create comment $x_i$
$x_i = \text{pi}(t,i)$	target,inst.	create processing instr. $x_i$
$x.\text{unbind}()$		detach $x$ from his parent
$x.\text{bind}(x_1)$	the parent	let $x$ be child of $x_1$
$x.\text{name}(s)$	the name	set element/attr. name
$x.\text{before}(x_1)$	a node	set $x$ 's position before $x_1$
$x.\text{after}(x_1)$	a node	set $x$ 's position after $x_1$

Figure 4: Basic operations on tree

there is no explicit deletion of nodes. We consider that tree fragments which are let unbound to the tree after a rule application will be freed from memory by an independant garbage collector.

### 3.3 Computation of rewriting operations

Rewriting operations deduced from a given rule are expressed as a sequence of operations described in figure 4. We outline below the  $\Phi$  function principles (*lhs* and *rhs* respectively stand for *left hand side* and *right hand side*).

DEFINITION 2. *The  $\Phi$  generation function: intentional definition*

1. tree nodes designated by variables in the lhs, and not used in the rhs are deleted.
2. tree nodes designated by variable in the rhs, and not used in the lhs are created<sup>5</sup>.
3. nodes defined in the rhs are created unless they are reused from the lhs<sup>5</sup>.
4. nodes of the right hand side are reorganized in such a way that they will finally match the rhs expression.

We tried to capture the essence of our rule semantics through the definition above. However, a more detailed definition is proposed hereafter:

DEFINITION 3. *The  $\Phi$  generation function: a more comprehensive definition*

1. nodes designated by variables in the lhs, and not used in the rhs are deleted (the link to the parent node is deleted through the unbind operation)
2. nodes designated by variable in the rhs, and not used in the lhs are created<sup>5</sup>
3. nodes defined in the rhs and not designated by any variable are created<sup>5</sup>
4. nodes designated both in lhs and rhs, and whose connexity is modified by the new arrangement are let unbound

<sup>5</sup>in accordance with the node type specified by the XPath expression.

5. newly created nodes (or nodes covered by item 4 are bound to the tree according to connexity information of the right hand expression
6. nodes are named in accordance with the rhs
7. position of nodes are set in accordance with position information contained in the rhs

As an illustration, we provide the reader with the translation of operations from right hand sides of examples 3 and 2 (see section 4). These operations are to be executed in the context  $\phi'$  built from matching the left hand side, and with the context node  $x$ .

```

 $\phi'(\$a).\text{unbind}()$ 
 $\phi'(\$b).\text{unbind}()$ 
 $x_1 = \text{element}()$ 
 $x_1.\text{name}(\text{"apply"})$ 
 $x_1.\text{bind}(x)$ 
 $x_2 = \text{element}()$ 
 $x_2.\text{name}(\text{"root"})$ 
 $x_2.\text{bind}(x_1)$ 
 $x_3 = \text{element}()$ 
 $x_3.\text{name}(\text{"degree"})$ 
 $x_4 = \text{text}(\text{"2"})$ 
 $x_4.\text{bind}(x_3)$ 
 $x_3.\text{bind}(x_1)$ 
 $\phi'(\$c).\text{unbind}()$ 
 $\phi'(\$c).\text{bind}(x_1)$ 

```

```

 $x_1 = \text{element}()$ 
 $x_1.\text{name}(\text{"degree"})$ 
 $x_1.\text{bind}(\phi'(\$a))$ 
 $x_2 = \text{text}(\text{"2"})$ 
 $x_2.\text{bind}(x_1)$ 
 $x_1.\text{after}(\phi'(\$b))$ 

```

Operations of example 2

```

 $x_3.\text{after}(x_2)$ 
 $\phi'(\$c).\text{after}(x_3)$ 

```

Operations of example 3

The reader may note that many issues remain open (we will review some of them in the conclusion), and are part of future work. We would like to design an exhaustive translation, simple enough to enable formal treatment and characterization.

### 3.4 Well-formedness of rules

We first propose to force syntactical restrictions on left and right terms: **not**  $p$  qualifiers and  $p|p$  expressions must not contain free variables. This respectively addresses the problem of (i) introducing variables that possibly cannot be defined after matching (e.g. **not**  $p[?a]$ ) and (ii) non-determinism in rule application (we avoid such cases at this point).

Beyond these basic preliminar restrictions, other rule specification could lead to fundamental problems. We propose to detail now all such possibilities. Our underlying method is based on two lines: (i) analyzing hypothesis related to tree connexity and node type constraints and (ii) analyzing illegal operations that may be generated. We summarize these problems below :

1. illegal node type conversions. It can be illustrated by the following ill-formed rules:

$$\text{table}[?a] \rightarrow \text{comment}()[?a]$$

which would require changing the type of a node, thus conduct to a violation of our tree data model, since a *comment* node could get elements as children. Moreover, this situation simply cannot be translated into our basic tree operations, since we do not propose the corresponding primitives

2. illegal node binding. This comes with rules like

$$table[?a][@border[?b]] \rightarrow @border[?b]/table[?a]$$

3. invalid XPath expression. Terms are considered invalid when they always evaluate to the empty set  $\emptyset$ , such as

$$table[\mathbf{not} \ tr]/tr \rightarrow .$$

4. multiply defined filter variables. In the following example, no solution can be found, since nodes designated by  $b$  are always distinct.

$$table[?a]/tr[?b]/td[?b] \rightarrow table[?a]$$

The first point requires a straightforward type analysis (it consists in annotating the variables with their node type, and checking the type invariance). An attractive solution to the second and third point is using a normalization process, such as described in [14]. This approach allows reducing any ill-formed XPath expressions into  $\perp$ , the always void XPath term<sup>6</sup>. Moreover, it reduces the syntactic complexity and eases the definition of subsequent processing. The fourth case requires a quite simple static analysis verifying that no filtering variables are introduced more than once. Note also that such a case might be detected through the normalization process we described above, provided the containment inference system is complete, which is not yet known ([14]; the completeness property is mathematically expressed, but not proved).

#### 4. REWRITING EXAMPLES

Although the paper does not focus on rule application models, one can consider standard approaches such as leftmost-innermost strategies, where the tree is explored node by node, and rules are tried in order. In any case, the transformation terminates when no rule can be applied to any node. We found it interesting to evaluate the expressiveness of our approach against “real world” examples, i.e MathML normalization and transformation ([6, 7]). Normalization is useful in order to eliminate assumptions about default values. For instance the specification says that the *degree* argument of *root* function is 2 if not specified. That is, in content semantics,

$$root(a) = root(2, a)$$

and also in presentation semantics

$$\sqrt{a} = \sqrt[2]{a}$$

The transformation is structurally described as

$$\begin{array}{l} \langle apply \rangle \\ \langle root \rangle \\ \langle ci \rangle a \langle /ci \rangle \\ \langle /apply \rangle \end{array} \rightarrow \begin{array}{l} \langle apply \rangle \\ \langle root \rangle \\ \langle degree \rangle 2 \langle /degree \rangle \\ \langle ci \rangle a \langle /ci \rangle \\ \langle /apply \rangle \end{array}$$

This is expressed as

EXAMPLE 2.

$$\begin{array}{l} apply[?a][\mathbf{not} \ degree]/root[?b] \\ \rightarrow apply[?a]/root[?b]/f-s::*[1][self::degree][. = "2"] \end{array}$$

<sup>6</sup>Still it must be slightly extended to consider positional extensions

Note that the right hand specifies that a new *degree* element must be created just after the original *root* element. The following solution, syntactically simpler, is computationally more expensive, as *apply* and *root* element are first deleted and then created

EXAMPLE 3.

$$\begin{array}{l} apply[?a][\mathbf{not} \ degree]/root[?b]/f-s::*[?c] \\ \rightarrow apply/root/f-s::degree[. = "2"]/f-s::*[?c] \end{array}$$

MathML is split in two parts: one is oriented toward the semantics of mathematical expressions (Content MathML, or C-MathML), and the other is dedicated to the presentation of these expressions (P-MathML, see [7] for a tutorial introduction). This fascinating duality opens interesting transformation problems that we believe can be elegantly addressed through the present proposal. One of them is about transforming C-MathML into P-MathML already addressed through XLST in [15]. For instance, the functional expression

$$power(add(a, 3), 2)$$

can be represented as

$$(a + 3)^2$$

The corresponding MathML transformation is expressed as

$$\begin{array}{l} \langle apply \rangle \\ \langle power \rangle \\ \langle apply \rangle \\ \langle plus \rangle \\ \langle ci \rangle a \langle /ci \rangle \\ \langle cn \rangle 3 \langle /cn \rangle \\ \langle /apply \rangle \\ \langle cn \rangle 2 \langle /cn \rangle \\ \langle /apply \rangle \end{array} \rightarrow \begin{array}{l} \langle msup \rangle \\ \langle mfenced \rangle \\ \langle mi \rangle a \langle /mi \rangle \\ \langle mo \rangle + \langle /mo \rangle \\ \langle mn \rangle 3 \langle /mn \rangle \\ \langle /mfenced \rangle \\ \langle mn \rangle 2 \langle /mn \rangle \\ \langle /msup \rangle \end{array}$$

A solution is proposed through the following rule set

$$\begin{array}{l} (r_1) \quad apply[?a]/power/f-s::*[?b]/f-s::*[?c] \\ \rightarrow msup/*[?b]/f-s::*[?c] \\ (r_2) \quad apply[?a][add]/*[?b]/f-s::*[?c] \\ \rightarrow mfenced/*[?b][1]/f-s::mo[. = "+" ]/f-s::*[?c] \\ (r_3) \quad ci[?a] \rightarrow mi[?a] \\ (r_4) \quad cn[?a][@type = "integer"]/\mathbf{text}()[?b] \\ \rightarrow mn/\mathbf{text}()[?b] \end{array}$$

Note that  $r_4$  just handles integers, and we would have to add more rules to handle all kind of numbers that are allowed by the specification. For instance, the following one addresses *complex-cartesian* numbers.

$$\begin{array}{l} (r_5) \quad cn[?a][@type = "complex-cartesian"]/ \\ \mathbf{text}()[?b]/f-s::sep/f-s::\mathbf{text}()[?c] \\ \rightarrow mn/mrow/mi[\mathbf{text}()[?b]]/ \\ f-s::mo[. = "-"]/f-s::mi[\mathbf{text}()[?c]] \end{array}$$

#### 5. CONCLUSION AND FUTURE WORK

We believe that the two main difficulties described in the introduction are addressed by our proposal:

1. An elegant and simple extension to XPath has been proposed in order to identify useful nodes when applying a pattern to a tree.
2. A natural interpretation of the right hand side is proposed ; moreover, it allows deep reorganization of the tree, such as subtree deletion, node renaming, node creation and positional changes.

However, many issues remain unclear with respect to the generation of transformation operations. Some specifications may generate multiple interpretations, and the  $\Phi$  algorithm will have to decide which is the best one. Let us consider the following rule

$$a/b[?v] \rightarrow a/*[?v]$$

According to our approach, any legal renaming of the  $\$v$  node is legal. Which one should be considered ? From the maximal efficiency point of view, no renaming should be applied. But in a broader context were we would like to perform static type checking analysis (let us imagine we transform an instance of a particular document schema), a good renaming operation would be the one that preserves the validity of the instance (e.g. a  $c$  if the content model of  $a$  nodes defines only  $b$  or  $c$  children).

Let us consider another difficult case, were the generated node must satisfy a complex constraint:

$$a \rightarrow b[?c][* \text{ and } * \sqsubseteq \$c/d]$$

After matching, this rule has to generate a node  $b$  child of the context node, and at least a node  $d$  child of  $b$ . So it looks like that in some cases, the tree operations must satisfy complex constraints, and it is yet unclear to the author how to deal with those issues in a simple way. Part of our future work will be thus to clarify this point and to consider existing literature on this topic, e.g. [2].

Another important issue is to relate pattern unification to pattern matching. For instance, the following two patterns should be considered as equivalent, up to a variable substitution

$$a[?v]/b/c[?w] \sim a[?u]/b/c[?r]$$

One way to explore is to extend XPath containment to pattern containment.

PROPOSITION 1. *pattern containment*

$$\vec{p}_1 \leq \vec{p}_2 \quad \text{iff} \quad \exists \gamma \text{ such that } \gamma \vdash \Xi(\vec{p}_1) \leq \Xi(\vec{p}_2)$$

PROPOSITION 2. *pattern equivalence*

$$\vec{p}_1 \sim \vec{p}_2 \quad \text{iff} \quad (\vec{p}_1 \leq \vec{p}_2 \wedge \vec{p}_2 \leq \vec{p}_1)$$

For the example above, a substitution could be  $\gamma = \{v = \$u, w = \$r, u = \$v, r = \$w\}$

This work on rewriting will be the building block of a new XML transformation language called  $\Omega$ , currently studied in the WAM project [17]. The  $\Omega$  approach tries to associate such transformation rules (a bit more complex, of form  $l, r \rightarrow r'$  where  $l$  is matching against an input document, and  $r \rightarrow r'$  rewrites the output document) with explicit strategies. The idea is to control the rule application and allow complex transformation models while maintaining good properties with respect to type checking and static analysis. We also explore through  $\Omega$  the possibility of defining

non-deterministic transformation strategies, were the output document might be constructed through backtracking and some kind of constraint solving.

## 6. ACKNOWLEDGMENTS

We would like to thanks the many people participating to this reflexion around XPath and supporting the ongoing work on  $\Omega$  transformation language: Vincent Quint and all his team from the WAM [17] INRIA project (especially Nabil Layaïda and Pierre Genevès); Jérôme Euzenat, Lionel Villard and Jean-Marc Andreoli.

## 7. REFERENCES

- [1] N. Dershowitz, J.-P. Jouanaud "Handbook of Theoretical Computer Science", Chapter on Rewriting, pp. 243-320 Elsevier, 1990.
- [2] H. Common, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, S. Tison and M. Tommasi "Tree Automata Techniques and Applications", Technical Report, [www.grappa.univ-lille3.fr/tata](http://www.grappa.univ-lille3.fr/tata), October, 1999.
- [3] XML Path Language (XPath) 1.0, W3C recommendation, 16 November, 1999.
- [4] Document Object Model (DOM) 1.0, W3C recommendation, October, 1998.
- [5] XML Path Language (XPath) 2.0, W3C working draft, 15 November, 2002.
- [6] Mathematical Markup Language (MathML) 2.0, W3C Recommendation, 21 February, 2001.
- [7] Michael Kohlhase, *New MathML tutorial*, W3C working group, [www.w3.org/Math/Documents/mathml-tutorial.pdf](http://www.w3.org/Math/Documents/mathml-tutorial.pdf), 11 March, 2003.
- [8] Phil Wadler, "A formal semantics of patterns in XSLT", Markup Technologies, 1999.
- [9] Phil Wadler, "Two semantics for XPath", [www.cs.bell-labs.com/who/wadler/topics/xml.html](http://www.cs.bell-labs.com/who/wadler/topics/xml.html), 1999.
- [10] Dan Olteanu and Holger Meuss and Tim Furche and François Bry, "Symmetry in XPath", technical report, October 2001, Computer Science Institute, Munich, Germany.
- [11] Alin Deutsch and Val Tannen, "Containment of Regular Path Expressions under Integrity Constraints", Knowledge Representation Meets Databases, 2001.
- [12] Frank Neven and Thomas Schwentick, "XPath containment in the presence of disjunction, DTDs, and variables", International Conference on Database Theory, 2003.
- [13] Jerome Miklau and Dan Suciu, "Containment and Equivalence for an XPath Fragment (Extended Abstract)", Symposium on Principles of Databases Systems, 2002.
- [14] Jean-Yves Vion-Dury and Nabil Layaïda, "Containment of XPath Expressions: an Inference and Rewriting based approach", Extreme Markup Languages, Montréal August, 2003. <http://wam.inrialpes.fr/publications/2003/xtrem2003/xtrem2003.pdf>
- [15] Emmanuel Piétriga, *MathMLc2p*, XSLT based translation of content MathML into presentation MathML, April, 2000. <http://opera.inrialpes.fr/people/Emmanuel.Pietriga/mathml2cp.html>
- [16] Wolfgang May "XPath-Logic and PathLog: A Logic-Based Approach for Declarative XML Data Manipulation", Technical Report No. 149, Institut für Informatik, Freiburg, Germany Feb, 2001.
- [17] WAM "Web Adaptation and Multimedia", Research project from INRIA Rhône-Alpes, <http://wam.inrialpes.fr>.