

# A formal model of pattern matching in XSL

Philip Wadler  
Bell Labs  
Lucent Technologies  
`wadler@research.bell-labs.com`

January 8, 1999

## 1 Introduction

This note presents a formal model of the semantics of pattern matching in XSL [XSL].

The formal semantics brings to light issues that can be hard to spot in an english language description. For instance, here is how [XSL] defines match:

The result of the *MatchExpr* is true if, for any node in the document that contains the context of the *MatchExpr*, the result of evaluating the *SelectExpr* with that node as context contains the context of the *MatchExpr*. Otherwise the result is false.

As James Clark noted [Clark], this sentence is ambiguous: does ‘contains’ qualify ‘dcoument’ or ‘node in the document’? (The intention is the former.) The formal specification given here presents the same information in just one line, and avoids the ambiguity.

The semantics is formalized using sets. When a selection pattern is processed using `for-each`, the selected nodes are presented in document order.

We use the usual notation for sets, including set comprehensions. See any standard textbook for an introduction, such as [Halmos]. We write  $\text{Set}(A)$  for the type of all finite sets with elements of type  $A$ , and  $\text{Set}_1(A)$  for the subtype of  $\text{Set}(A)$  where the sets have at most one element. We write

$$\textit{first}, \textit{last} : \text{Set}(\text{Node}) \rightarrow \text{Set}_1(\text{Node})$$

for the functions that take a set and return the singleton set containing just its first or last node in document order, or return the empty set if the given set is empty.

We represent a relation on nodes by a function  $r : \text{Node} \rightarrow \text{Set}(\text{Node})$ . We define the transitive closure  $r^+$  and the reflexive and transitive closure  $r^*$  in the usual way.

$$\begin{aligned} r^+(x) &= \{x_2 \mid x_1 \in r(x), x_2 \in r^*(x_1)\} \\ r^*(x) &= \{x\} \cup r^+(x) \end{aligned}$$

If  $\text{subnodes} : \text{Node} \rightarrow \text{Set}(\text{Node})$  returns the immdiediate descendants of a node, then  $\text{subnodes}^+(x)$  returns all proper descendants of node  $x$ , and  $\text{subnodes}^*(x)$  returns all descendants of node  $x$  including  $x$  itself.

Throughout we use the following conventions for variables.

$n$	: Name
$p$	: Pattern
$q$	: Qualifier
$s$	: String
$x$	: Node

## 2 Access to nodes

We take as primitive functions to determine the type of a node and to access its components, as specified in [XSL, Section 2.4].

$$\begin{aligned}
isRoot &: Node \rightarrow Boolean \\
isElement &: Node \rightarrow Boolean \\
isAttribute &: Node \rightarrow Boolean \\
isText &: Node \rightarrow Boolean \\
isComment &: Node \rightarrow Boolean \\
isPI &: Node \rightarrow Boolean \\
\\
parent &: Node \rightarrow Set_1(Node) \\
children &: Node \rightarrow Set(Node) \\
attributes &: Node \rightarrow Set(Node) \\
name &: Node \rightarrow Name \\
value &: Node \rightarrow String
\end{aligned}$$

The function *parent* returns a singleton set containing the parent of the node, or returns the empty set if the node is a root. The function *children* returns the empty set if the node is a leaf.

Each node belongs to a document, and each document has a root node.

$$\begin{aligned}
document &: Node \rightarrow Document \\
root &: Document \rightarrow Node
\end{aligned}$$

We represent identifiers by strings. The function *id* takes a string and a document and returns the singleton set containing the unique element within the document which has the string as its ‘id’ attribute, or returns the empty set if no such element exists. The function *split* takes a string into the set of space-separated identifiers within it.

$$\begin{aligned}
id &: Document \times String \rightarrow Set_1(Node) \\
split &: String \rightarrow Set(String)
\end{aligned}$$

The subnodes of a node are all the immediate descendants, whether children or attributes.

$$\begin{aligned}
subnodes &: Node \rightarrow Set(Node) \\
subnodes(x) &= \text{if } isElement(x) \text{ then } children(x) \cup attributes(x) \text{ else } children(x)
\end{aligned}$$

What follows refers only to subnodes, not attributes or children. For all nodes  $x$  we have  $x \in subnodes^*(root(document(x)))$ . Further, for all nodes  $x$  and  $x_1$  we have  $\neg isRoot(x) \wedge parent(x) = x_1$  if and only if  $x \in subnodes(x_1)$ .

The function *siblingElements* returns all siblings of a node that are elements, and the function *siblingsOfType* returns all siblings of an element that have the same tag as the element.

$$\begin{aligned}
siblingElements &: Node \rightarrow Set(Node) \\
siblingElements(x) &= \{x_2 \mid x_1 \in parent(x), x_2 \in subnodes(x_1), isElement(x_2)\} \\
\\
siblingsOfType &: Node \rightarrow Set(Node) \\
siblingsOfType(x) &= \{x_1 \mid x_1 \in siblingElements(x), isElement(x), name(x) = name(x_1)\}
\end{aligned}$$

## 3 Semantics of patterns

Figure 1 shows the semantics of patterns, corresponding to [XSL, Section 2.6]. The term  $\mathcal{M}[p]x$  is true if pattern  $p$  matches node  $x$ , the term  $\mathcal{S}[p]x$  denotes the set of nodes selected by pattern  $p$  from current node  $x$ , and the term  $\mathcal{Q}[q]x$  is true if qualifier  $q$  holds for current node  $x$ .

## References

- [Clark] James Clark, Message to XSL mailing list, 1 January 1999.
- [Halmos] Paul R. Halmos, *Naive Set Theory*, Springer-Verlag, 1970.
- [XSL] Extensible Stylesheet Language (XSL), Version 1.0, World Wide Web Consortium Working Draft 16 December 1998. <http://www.w3.org/TR/WD-xsl>

---

$\mathcal{M}$	$: \text{Pattern} \times \text{Node} \rightarrow \text{Boolean}$
$\mathcal{M}[p]x$	$= \{ x_2 \mid x_1 \in \text{subnodes}^*(\text{root}(\text{document}(x))), x_2 \in \mathcal{S}[p](x_1) \}$
$\mathcal{S}$	$: \text{Pattern} \times \text{Node} \rightarrow \text{Set}(\text{Node})$
$\mathcal{S}[p_1 \mid p_2]x$	$= \mathcal{S}[p_1]x \cup \mathcal{S}[p_2]x$
$\mathcal{S}[\![/p]\!]x$	$= \mathcal{S}[\!/\!p\!]\!(\text{root}(\text{document}(x)))$
$\mathcal{S}[\![/\!/p]\!]x$	$= \{ x_2 \mid x_1 \in \text{subnodes}^*(\text{root}(\text{document}(x))), x_2 \in \mathcal{S}[p]x_1 \}$
$\mathcal{S}[\![p_1 / p_2]\!]x$	$= \{ x_2 \mid x_1 \in \mathcal{S}[p_1]x, x_2 \in \mathcal{S}[p_2]x_1 \}$
$\mathcal{S}[\![p_1 /\!/p_2]\!]x$	$= \{ x_3 \mid x_1 \in \mathcal{S}[p_1]x, x_2 \in \text{subnodes}^*(x_1), x_3 \in \mathcal{S}[p_2]x_2 \}$
$\mathcal{S}[\![p[q]]\!]x$	$= \{ x_1 \mid x_1 \in \mathcal{S}[p]x, \mathcal{Q}[\![q]\!]x_1 \}$
$\mathcal{S}[\![n]\!]x$	$= \{ x_1 \mid x_1 \in \text{subnodes}(x), \text{isElement}(x_1), \text{name}(x_1) = n \}$
$\mathcal{S}[\![*\]\!]x$	$= \{ x_1 \mid x_1 \in \text{subnodes}(x), \text{isElement}(x_1) \}$
$\mathcal{S}[\![@n]\!]x$	$= \{ x_1 \mid x_1 \in \text{subnodes}(x), \text{isAttribute}(x_1), \text{name}(x_1) = n \}$
$\mathcal{S}[\![@*\]\!]x$	$= \{ x_1 \mid x_1 \in \text{subnodes}(x), \text{isAttribute}(x_1) \}$
$\mathcal{S}[\![\text{text}()\]\!]x$	$= \{ x_1 \mid x_1 \in \text{subnodes}(x), \text{isText}(x_1) \}$
$\mathcal{S}[\![\text{comment}()\]\!]x$	$= \{ x_1 \mid x_1 \in \text{subnodes}(x), \text{isComment}(x_1) \}$
$\mathcal{S}[\![\text{pi}(n)]\!]x$	$= \{ x_1 \mid x_1 \in \text{subnodes}(x), \text{isPI}(x_1), \text{name}(x_1) = n \}$
$\mathcal{S}[\![\text{pi}()]\!]x$	$= \{ x_1 \mid x_1 \in \text{subnodes}(x), \text{isPI}(x_1) \}$
$\mathcal{S}[\![\text{id}(p)]\!]x$	$= \{ x_2 \mid x_1 \in \mathcal{S}[p]x, s \in \text{split}(\text{value}(x_1)), x_2 \in \text{id}(\text{document}(x), s) \}$
$\mathcal{S}[\![\text{id}(s)]\!]x$	$= \{ x_1 \mid s_1 \in \text{split}(s), x_1 \in \text{id}(\text{document}(x), s_1) \}$
$\mathcal{S}[\![\text{ancestor}(p)]\!]x$	$= \text{last}(\{ x_1 \mid x_1 \in \text{parent}^+(x), \mathcal{M}[p]x_1 \})$
$\mathcal{S}[\![\text{ancestor-or-self}(p)]\!]x$	$= \text{last}(\{ x_1 \mid x_1 \in \text{parent}^*(x), \mathcal{M}[p]x_1 \})$
$\mathcal{S}[\![.\]\!]x$	$= \{ x \}$
$\mathcal{S}[\![..\]\!]x$	$= \text{parent}(x)$
$\mathcal{Q}$	$: \text{Qualifier} \times \text{Node} \rightarrow \text{Boolean}$
$\mathcal{Q}[\![q_1 \text{ and } q_2]\!]x$	$= \mathcal{Q}[\![q_1]\!]x \wedge \mathcal{Q}[\![q_2]\!]x$
$\mathcal{Q}[\![q_1 \text{ or } q_2]\!]x$	$= \mathcal{Q}[\![q_1]\!]x \vee \mathcal{Q}[\![q_2]\!]x$
$\mathcal{Q}[\![\text{not}(q)]\!]x$	$= \neg \mathcal{Q}[\![q]\!]x$
$\mathcal{Q}[\![\text{first-of-type}()]\!]x$	$= x \in \text{first}(\text{siblingsOfType}(x))$
$\mathcal{Q}[\![\text{last-of-type}()]\!]x$	$= x \in \text{last}(\text{siblingsOfType}(x))$
$\mathcal{Q}[\![\text{first-of-any}()]\!]x$	$= x \in \text{first}(\text{ siblingElements}(x))$
$\mathcal{Q}[\![\text{last-of-any}()]\!]x$	$= x \in \text{last}(\text{ siblingElements}(x))$
$\mathcal{Q}[\![p=s]\!]x$	$= \{ x_1 \mid x_1 \in \mathcal{S}[p]x, \text{value}(x_1) = s \} \neq \emptyset$
$\mathcal{Q}[\![p]\!]x$	$= \mathcal{S}[p]x \neq \emptyset$

---

Figure 1: Semantics of patterns