

# A Pushdown Machine for Recursive XML Processing\*

Keisuke Nakano<sup>1</sup> and Shin-Cheng Mu<sup>2</sup>

<sup>1</sup> Department of Mathematical Informatics,  
University of Tokyo, Japan  
ksk@mist.i.u-tokyo.ac.jp

<sup>2</sup> Institute of Information Science,  
Academia Sinica, Taiwan  
scm@iis.sinica.edu.tw

**Abstract.** XML transformations are most naturally defined as recursive functions on trees. Their direct implementation, however, causes inefficient memory usage because the input XML tree is completely built in memory before being processed. In contrast, programs in stream processing style minimise memory usage since it may output part of the result while reading the input. Stream processing programs, however, are harder to write because the programmer is left with the burden to maintain the state. In this paper, we propose a model for XML stream processing and show that all programs written in a particular style of recursive functions on XML trees, the *macro forest transducer*, can be automatically translated to our stream processors. The stream processor is declarative in style, but can be implemented efficiently by a pushdown machine. We thus get the best of both worlds — program clarity, and efficiency in execution.

## 1 Introduction

Since an XML document has a tree-like structure, it is natural to define XML transformations as recursive functions over trees. Several XML-oriented languages, such as XSLT [34], *fst* [3], XDuce [11] and CDuce [2], allow the programmer to define mutual recursive functions over forests.

As an example, consider the program in Figure 1. Let  $\sigma(f_1)f_2$  denote a forest where the head is a  $\sigma$ -labeled tree whose children constitute the forest  $f_1$ , and the tail is a sibling forest  $f_2$ . The empty forest is denoted by  $\epsilon$  and is usually omitted when enclosed in other trees. The function *Main* in Figure 1 scans through the input tree and reverses the order of all subtrees under nodes labelled **r** by calling the function *Rev*. For example, the input tree  $\mathbf{a}(\mathbf{r}(\mathbf{b}(\mathbf{c}(\mathbf{d}(\mathbf{e}(\mathbf{f}))))))$  is transformed into  $\mathbf{a}(\mathbf{r}(\mathbf{e}(\mathbf{b}(\mathbf{d}(\mathbf{c}(\mathbf{f}))))))$ .

A naive way to execute functions defined in this style is to load the entire forest into memory, so that we have convenient access to the children and siblings for each node. The input stream of tokens, also called *XML events*, is parsed to build the corresponding forest. The forest is then transformed by the function, before the resulting forest is unparsed to an XML stream. Loading the entire tree into memory is not preferable when we have to process large input. However, many XML transformation languages such as XSLT, *fst*, XDuce and CDuce are actually implemented this way.

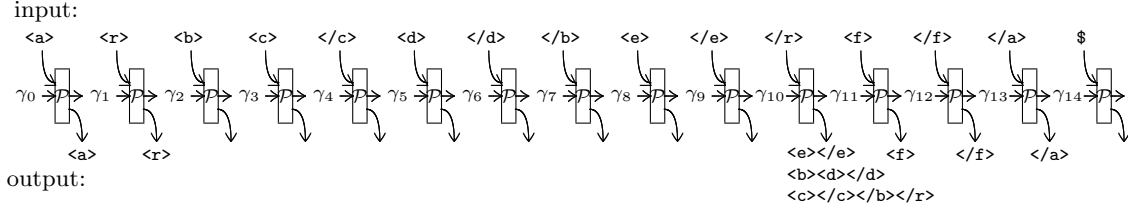
To optimise space usage, we may ask the programmer to switch to another programming style facilitated by APIs such as SAX [32]. The stream processor reads XML events one by one, and the programmer defines respectively what to do when it encounters a start tag  $\langle\sigma\rangle$ , an end tag  $\langle/\sigma\rangle$ , or end of stream  $\$$ . Figure 2 shows what a stream processor would do to accomplish the same task, given the input  $\langle\mathbf{a}\rangle\langle\mathbf{r}\rangle\langle\mathbf{b}\rangle\langle\mathbf{c}\rangle\langle/\mathbf{c}\rangle\langle\mathbf{d}\rangle\langle/\mathbf{d}\rangle\langle\mathbf{e}\rangle\langle/\mathbf{e}\rangle\langle\mathbf{f}\rangle\langle/\mathbf{f}\rangle\langle/\mathbf{a}\rangle$ . The function  $\mathcal{P}$  takes an input XML event and an environment<sup>3</sup>, and returns an updated environment, possibly outputting some XML events. The initial environment in the running example is  $\gamma_0$ . Reading the

\* This work is partially supported by the *Comprehensive Development of e-Society Foundation Software* program of the Ministry of Education, Culture, Sports, Science and Technology, Japan.

<sup>3</sup> A ‘state’ storing information needed to carry out the computation. We use the name ‘environment’ to avoid confusion with mft states.

$Main(\epsilon) = \epsilon$	$Rev(\epsilon, y) = y$
$Main(\mathbf{r}\langle x_1 \rangle x_2) = \mathbf{r}\langle rev(x_1, \epsilon) \rangle (Main(x_2))$	$Rev(\sigma\langle x_1 \rangle x_2, y) = Rev(x_2, \sigma\langle Rev x_1 \epsilon \rangle y)$
$Main(\sigma\langle x_1 \rangle x_2) = \sigma\langle main x_1 \rangle (Main(x_2))$ if $\sigma \neq \mathbf{r}$	

**Fig. 1.** A functional program reversing the subtrees under nodes labelled  $\mathbf{r}$ .



**Fig. 2.** A flow of stream processing

first event  $\langle a \rangle$  yields a new environment  $\gamma_1$  and an output event  $\langle a \rangle$ . The next event  $\langle r \rangle$  is also copied to the output. After that, no output event will be produced for a while, because there is no way for the processor to know what to output before the closing tag  $\langle /r \rangle$  is read. Between  $\langle r \rangle$  and  $\langle /r \rangle$ , the computer reads the input and stores partially computed result in the environment. While stream processing saves memory usage, it is much harder to program in this style. Each problem seems to employ a different data structure to represent the environment, and burden is on the programmer to explicitly maintain it.

Can we write a recursive function on forests and have it automatically transformed to a program in the stream processing style, thereby achieve both clarity and memory efficiency? In this paper, we present a model for an XML stream processor, and shows how to automatically derive XML stream processors from a very expressive class of recursive functions on forests.

Some readers may wonder, can we not use a non-strict programming language, compose a function on forests with a parser constructing the forest on demand, and simply let lazy evaluation do the job? In the latter part of this paper, we will talk about the difficulties of relying on lazy evaluation alone, and the necessities of a more specific model for XML streaming.

We have made two main contributions. Firstly, we propose a model for XML stream processing which is declarative in nature but has an efficient implementation. The environment can be represented uniformly by a partially evaluated stream, called a *temporary expression*. Reading a token triggers a rewriting on the expression. It can be argued that our stream processor is of a higher level than finite automaton based approaches. Yet, for the class of stream processors we are concerned with, there is an efficient method, basing on a pushdown machine, to implement the rewriting.

Secondly, we present a method to derive a stream processor from any function definable in terms of the *macro forest transducer* (mft), proposed by Perst and Seidl [26]. The derivation, which can be seen as a special case of program fusion [30], works by fusing the mft with an XML parser recast as a *top-down tree transducer* (tdtt). The fusion is similar Engelfriet and Vogler's method of composing a (finitary) tdtt and a *macro tree transducer* [6], but we give a proof that the method works for our tdtt, which has an infinite number of states.

The XML stream processor in this paper is implemented in Objective Caml [25] and compared with a few existing tools for forest-style XML transformation.

*Outline* In the next section, we introduce a simplified model of XML documents and mft's. Section 3 gives a formal model of XML stream processors and its derivation from an mft. Section 4 proposes an efficient implementation of the stream processor based on a pushdown machine. In Section 5 we show the benchmark result for a simple example. Section 6 discusses capabilities of our framework. In Section 7 we review some related work and conclude the paper. Appendix A gives a proof of the correctness of our derivation.

## 2 XML and the Macro Forest Transducer

This section formally defines our model of XML and the mft. For simplicity, we deal with a simplified model of XML with only element nodes, and assume that the input XML is well-formed. It is easy to extend our approach to capture other aspects of XML, and it is always possible to perform a well-formedness check before hand.

### 2.1 Forests and XML Streams

Let  $\Sigma$  be an alphabet. A  $\Sigma$ -forest (also called a  $\Sigma$ -hedge [18]), is defined by

$$f ::= \sigma \langle f \rangle f \mid \epsilon,$$

where  $\sigma \in \Sigma$  and  $\epsilon$  denotes the empty forest. We denote by  $\mathcal{F}_\Sigma$  the set of  $\Sigma$ -forests. A  $\Sigma$ -forest  $\mathbf{a} \langle \mathbf{b} \langle \epsilon \rangle \mathbf{c} \langle \epsilon \rangle \epsilon \rangle$  with  $\Sigma = \{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$  represents the XML fragment  $\langle \mathbf{a} \rangle \langle \mathbf{b} \rangle \langle \mathbf{c} \rangle \langle \mathbf{c} \rangle \langle \mathbf{a} \rangle$ . The concatenation of two forests  $f_1, f_2 \in \mathcal{F}_\Sigma$  is written  $f_1 f_2$ . The symbol  $\epsilon$ , being the unit of concatenation, is often omitted.

The  $\Sigma$ -events, written  $\Sigma_{\langle \rangle}$ , is defined by  $\Sigma_{\langle \rangle} = \{\langle \sigma \rangle \mid \sigma \in \Sigma\} \cup \{\langle / \sigma \rangle \mid \sigma \in \Sigma\}$ . An XML stream is a sequence of  $\Sigma$ -events. We denote by  $\Sigma_{\langle \rangle}^\circ$  the set of well-formed sequences of  $\Sigma$ -events and denote by  $\epsilon$  the empty sequence. The symbol  $\$$  denotes the end of an (input) XML stream, which is also regarded as an event. We write  $\Sigma_{\langle \rangle}^\circ \$$  for  $\Sigma_{\langle \rangle}^\circ \cup \{\$\}$ .

Let  $\Sigma$  be an alphabet. The *streaming* of a forest is the function  $[\_ ] : \mathcal{F}_\Sigma \rightarrow \Sigma_{\langle \rangle}^\circ \$$  defined by  $[\sigma \langle f_1 \rangle f_2] = \langle \sigma \rangle [f_1] \langle / \sigma \rangle [f_2]$  and  $[\epsilon] = \epsilon$ . For example,  $[\mathbf{a} \langle \mathbf{b} \langle \mathbf{c} \rangle \rangle] = \langle \mathbf{a} \rangle \langle \mathbf{b} \rangle \langle \mathbf{c} \rangle \langle / \mathbf{c} \rangle \langle / \mathbf{b} \rangle \langle / \mathbf{a} \rangle$ .

### 2.2 Macro Forest Transducers

The *macro forest transducer* (mft) was proposed by Perst and Seidl [26] as an extension to the *macro tree transducer* (mtt) [6] by taking forest concatenation as a basic operator. Functional programmers can think of an mft as a recursive function mapping a forest (and possibly some accumulating parameters) to a forest, with certain restriction on their shapes — the pattern on the forest extracts only the label, the children and the sibling of the first tree; the accumulating parameters cannot be pattern-matched; each function call is passed either the children or the sibling.

We do not propose using the mft as a programming language, but as an intermediate language. A natural question is how expressive it is under these restrictions. It was shown that mft is in fact rather expressive [17]. More discussions will be given in Section 6, where we also show how to convert certain classes of functional programs to mft.

In the convention of mft, a function is called a *state* and its arity its *rank*. Let us write  $\mathbb{N}$  and  $\mathbb{N}^+$  for the set of non-negative integers including and excluding 0, respectively.

**Definition 1.** A *macro forest transducer* (mft) is a tuple  $M = (Q, \Sigma, \Delta, in, R)$ , where

- $Q$  is a finite set of ranked states, the rank of a state given by a function  $rank : Q \rightarrow \mathbb{N}^+$ ,
- $\Sigma$  and  $\Delta$  are alphabets with  $Q \cap (\Sigma \cup \Delta) = \emptyset$ , called the *input alphabet* and the *output alphabet*, respectively,
- $in \in Q$  is the initial ranked state,
- $R$  is a set of rules partitioned by  $R = \bigcup_{q \in Q} R_q$ . For each  $q \in Q$ ,  $R_q$  consists of rules of the form  $q(pat, y_1, \dots, y_n) \rightarrow rhs$ , where  $n = rank(q) - 1$  and
  - $pat$  is either  $\epsilon$  or  $\sigma \langle x_1 \rangle x_2$  for some  $\sigma \in \Sigma$ ,
  - $rhs$  ranges over expressions defined by

$$rhs ::= q'(x_i, rhs, \dots, rhs) \mid \epsilon \mid \delta \langle rhs \rangle \mid y_j \mid rhs \ rhs$$

with  $q' \in Q$ ,  $\delta \in \Delta$ ,  $i = 1, 2$  and  $j = 1, \dots, n$ . Additionally, no variable  $x_i$  occurs in  $rhs$  when  $pat = \epsilon$ .

Perst and Seidl's mft, designed for type checking, can be non-deterministic. Since our focus is on program transformation, our mft's are deterministic and total. That is, for each  $q$  and  $\sigma$  there is exactly one such rule  $q(\sigma\langle x_1 \rangle x_2, \dots) \rightarrow rhs$ . We will denote its right-hand side by  $rhs^{q,\sigma}$ . Similarly  $rhs^{q,\epsilon}$  stands for the right-hand side  $rhs$  of the unique rule  $q(\epsilon, \dots) \rightarrow rhs$ . If a rule for state  $q$  and pattern  $p$  is missing, we assume that there is an implicit rule  $q(p, \dots) \rightarrow \epsilon$ . The semantics of mft's is defined by translating every state into a function [26].

**Definition 2.** Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft. The semantics of a states  $q \in Q$  is given by the function  $\llbracket q \rrbracket : \mathcal{F}_\Sigma \times (\mathcal{F}_\Delta)^n \rightarrow \mathcal{F}_\Delta$  where  $n = rank(q) - 1$ . Each  $\llbracket q \rrbracket$  is defined by:

- $\llbracket q \rrbracket(\sigma\langle \omega_1 \rangle \omega_2, \varphi_1, \dots, \varphi_n) = \llbracket rhs^{q,\sigma} \rrbracket_\rho$  where  $\rho(x_i) = \omega_i$  for  $i = 1, 2$  and  $\rho(y_j) = \varphi_j$  for  $j = 1, \dots, n$ ,
- $\llbracket q \rrbracket(\epsilon, \varphi_1, \dots, \varphi_n) = \llbracket rhs^{q,\epsilon} \rrbracket_\rho$  where  $\rho(y_j) = \varphi_j$  for  $j = 1, \dots, n$ ,

where  $\llbracket \_ \rrbracket_\rho$  evaluates the right-hand side with respect to the environment  $\rho$ :

$$\begin{aligned} \llbracket q'(x_i, rhs_1, \dots, rhs_{n'}) \rrbracket_\rho &= \llbracket q' \rrbracket(\rho(x_i), \llbracket rhs_1 \rrbracket_\rho, \dots, \llbracket rhs_{n'} \rrbracket_\rho), \\ \llbracket \epsilon \rrbracket_\rho &= \epsilon, & \llbracket \delta \langle rhs \rangle \rrbracket_\rho &= \delta \langle \llbracket rhs \rrbracket_\rho \rangle, \\ \llbracket y_j \rrbracket_\rho &= \rho(y_j), & \llbracket rhs \ rhs' \rrbracket_\rho &= \llbracket rhs \rrbracket_\rho \llbracket rhs' \rrbracket_\rho. \end{aligned}$$

**Definition 3.** The transformation induced by an mft  $M = (Q, \Sigma, \Delta, in, R)$  is the function  $\tau_M : \mathcal{F}_\Sigma \rightarrow \mathcal{F}_\Delta$  defined by  $\tau_M(f) = \llbracket in \rrbracket(f, \epsilon, \dots, \epsilon)$ .

*Example 1.* Let  $Q = \{Main, Rev\}$ ,  $\Sigma = \Delta$  be some alphabet containing  $\mathbf{r}$  and  $R$  be the rules in Figure 1 (with the  $=$  sign replaced by  $\rightarrow$ ), then  $M_{rev} = (Q, \Sigma, \Delta, Main, R)$  is an mft.

*Example 2.* The mft  $M_{htm} = (Q, \Sigma, \Delta, Main, R)$  defined below reads an XML document consisting of a title and several paragraphs with some keywords. The output is a (simplified) HTML document where the `para` tag is converted to `p` and `key` tag to `em`. Furthermore, before the `ps` tag we dump the list of keywords we collect so far. Text data is denoted by a node with no children.

$Q = \{Main, Title, InArticle, Key2Em, AllKeys, Copy\}$ ,

$\Sigma = \Delta =$  (some alphabet containing English alphabets and the XML/HTML tags below),

$R = \{ Main(\mathbf{article}\langle x_1 \rangle x_2) \rightarrow \mathbf{html}\langle \mathbf{head}\langle Title(x_1) \rangle \mathbf{body}\langle InArticle(x_1, \epsilon) \rangle \epsilon \rangle,$

$Title(\mathbf{title}\langle x_1 \rangle x_2) \rightarrow \mathbf{title}\langle Copy(x_1) \rangle,$

$InArticle(\mathbf{title}\langle x_1 \rangle x_2, y_1) \rightarrow \mathbf{h1}\langle Copy(x_1) \rangle InArticle(x_2, y_1),$

$InArticle(\mathbf{para}\langle x_1 \rangle x_2, y_1) \rightarrow \mathbf{p}\langle Key2Em(x_1) \rangle InArticle(x_2, y_1 AllKeys(x_1)),$

$InArticle(\mathbf{ps}\langle x_1 \rangle x_2, y_1) \rightarrow \mathbf{h2}\langle \mathbf{Index}\langle \rangle \rangle \mathbf{ul}\langle y_1 \rangle \mathbf{h2}\langle \mathbf{Postscript}\langle \rangle \rangle Copy(x_1),$

$Key2Em(\mathbf{key}\langle x_1 \rangle x_2) \rightarrow \mathbf{em}\langle Copy(x_1) \rangle Key2Em(x_2),$

$Key2Em(\sigma\langle x_1 \rangle x_2) \rightarrow \sigma\langle Key2Em(x_1) \rangle Key2Em(x_2) \quad (\sigma \neq \mathbf{key}), \quad Key2Em(\epsilon) \rightarrow \epsilon,$

$AllKeys(\mathbf{key}\langle x_1 \rangle x_2) \rightarrow \mathbf{li}\langle Copy(x_1) \rangle AllKeys(x_2),$

$AllKeys(\sigma\langle x_1 \rangle x_2) \rightarrow AllKeys(x_1) AllKeys(x_2) \quad (\sigma \neq \mathbf{key}), \quad AllKeys(\epsilon) \rightarrow \epsilon,$

$Copy(\sigma\langle x_1 \rangle x_2) \rightarrow \sigma\langle Copy(x_1) \rangle Copy(x_2) \quad (\sigma \in \Sigma), \quad Copy(\epsilon) \rightarrow \epsilon \}$

### 3 XML Stream Processors and its Derivation

A *temporary expression* is a partially computed stream of XML events. An XML stream processor defines how to rewrite a temporary expression upon reading each input event.

**Definition 4.** An *XML stream processor* (xsp) is a tuple  $S = (Q, \Sigma, \Delta, in, R)$ , where

- $Q$  is a (possibly infinite) set of ranked states, the rank for each state given by  $rank : Q \rightarrow \mathbb{N}$ ,
- $\Sigma$  and  $\Delta$  are (finite) alphabets with  $Q \cap (\Sigma \cup \Delta) = \emptyset$ , called the *input alphabet* and the *output alphabet*, respectively,

- $in \in Q$  is the initial state,
- $R = \{q(y_1, \dots, y_n) \xrightarrow{\chi} rhs \mid q \in Q, \chi \in \Sigma_{\langle \rangle} \}$  is a set of rules, where  $n = rank(q)$  and  $rhs$  ranges over expressions defined by

$$rhs ::= q'(rhs, \dots, rhs) \mid \varepsilon \mid \langle \delta \rangle rhs \langle /\delta \rangle \mid y_j \mid rhs \ rhs$$

where  $q' \in Q$ ,  $\delta \in \Delta$  and  $j = 1, \dots, n$ . Additionally, the pattern  $q'(\dots)$  does not occur in  $rhs$  for any  $q' \in Q$  when  $\chi = \$$ .

### 3.1 Semantics of XML Stream Processors

The semantics of an xsp is defined by translating every rule of the xsp into a transition for temporary expressions.

**Definition 5.** Let  $S = (Q, \Sigma, \Delta, in, R)$  be an xsp. A temporary expression for  $S$ , denoted by  $Tmp_S$ , is defined by  $E ::= \varepsilon \mid \langle \delta \rangle E \mid \langle /\delta \rangle E \mid q(E, \dots, E)E$ .

**Definition 6.** Let  $S = (Q, \Sigma, \Delta, in, R)$  be an xsp and  $s \in \Sigma_{\langle \rangle}^{\circ}$ . The transition over  $Tmp_S$  for an input  $\Sigma$ -event is a function  $\langle \_, \_ \rangle : Tmp_S \times \Sigma_{\langle \rangle} \rightarrow Tmp_S$  defined by

- $\langle \varepsilon, \chi \rangle = \varepsilon$ ,
- $\langle \langle \delta \rangle e, \chi \rangle = \langle \delta \rangle \langle e, \chi \rangle$  where  $\delta \in \Delta$ ,
- $\langle \langle /\delta \rangle e, \chi \rangle = \langle /\delta \rangle \langle e, \chi \rangle$  where  $\delta \in \Delta$ ,
- $\langle q(e_1, \dots, e_n)e, \chi \rangle = (rhs[y_j := \langle e_j, \chi \rangle]_{j=1, \dots, n}) \langle e, \chi \rangle$  where  $(q(y_1, \dots, y_n) \xrightarrow{\chi} rhs) \in R$  with  $q \in Q$  and  $\chi \in \Sigma_{\langle \rangle} \}$ .

The initial temporary expression is  $in(\varepsilon, \dots, \varepsilon)$ . An xsp reads the input stream of events one by one, and updates the temporary expression with the transition  $\langle \_, \_ \rangle$ . The end of the stream is marked by  $\$$ . Let  $\chi_1 \chi_2 \dots \chi_k$  be the input stream with each  $\chi_j \in \Sigma_{\langle \rangle}$ . The final expression is

$$\langle \langle \dots \langle \langle in(\varepsilon, \dots, \varepsilon), \chi_1 \rangle, \chi_2 \rangle, \dots, \chi_k \rangle, \$ \rangle. \quad (1)$$

Note that the final temporary expression is always in  $\Delta_{\langle \rangle}^{\circ}$  since, by Definition 4, the right-hand side of a  $(q, \$)$ -rule does not contain any unevaluated state  $q'(\dots)$ .

**Definition 7.** The transformation induced by an xsp  $S = (Q, \Sigma, \Delta, in, R)$  is the function  $\tau_S : \Sigma_{\langle \rangle}^{\circ} \rightarrow \Delta_{\langle \rangle}^{\circ}$  defined by  $\tau_S(s) = \theta_S(in(\varepsilon, \dots, \varepsilon), s\$)$  where, for  $e \in Tmp_S$ ,

$$\theta_S(e, \varepsilon) = e, \quad \theta_S(e, \chi s) = \theta_S(\langle e, \chi \rangle, s).$$

We will see some examples in the next section.

The induced transformation defines declaratively what the output stream is, given the input stream. The very reason we program in the stream processing style, however, is to be able to print out a prefix of the output stream while reading the input. That is, we would like to ‘squeeze’ some part of the result from after each event read. This will be described in Section 4.3.

### 3.2 Deriving Stream Processors From Macro Forest Transducers

Given an mft  $M = (Q, \Sigma, \Delta, in, R)$ , an input stream  $x$ , and a function  $Parse :: \Sigma_{\langle \rangle}^{\circ} \rightarrow \mathcal{F}_{\Sigma}$  parsing a stream of events into a forest, the expression  $\llbracket in \rrbracket (Parse(x), \varepsilon, \dots, \varepsilon)$  yields a  $\Delta_{\langle \rangle}^{\circ}$  stream. If we can fuse the three functions,  $\llbracket \_ \rrbracket$ ,  $\llbracket in \rrbracket$ , and  $Parse$  into one, we may have a stream processor. Fusing  $\llbracket \_ \rrbracket$  and  $\llbracket in \rrbracket$  is a relatively easy task. The interesting step is fusing them with the parser. An XML parser can be written as a *top-down tree transducer* (tdtt) with an (countably-)infinite number of states

$$\begin{aligned} Parse[1] (\langle \sigma \rangle s) &= \sigma \langle Parse[1] s \rangle (Parse[2] s) \\ Parse[i] (\langle \sigma \rangle s) &= Parse[i+1] s \quad (i > 1) \\ Parse[1] (\langle /\sigma \rangle s) &= \varepsilon \\ Parse[i] (\langle /\sigma \rangle s) &= Parse[i-1] s \quad (i > 1) \\ Parse[i] (\$) &= \varepsilon \end{aligned}$$

for every  $\sigma \in \Sigma$ . Note that we do not need a forest transducer for parsing. The forest is constructed without using forest concatenation. Therefore, although it returns a forest, *Parse* is still technically a tree transducer where the forest is represented by a binary tree. Functional programmers may alert that the first clause seems to imply that  $s$  is traversed twice. However, it is a common style of tdt specifications, and multiple traversals is in fact avoided in the implementation, to be discussed in Section 4. We will also talk about a more typical way to specify the parser, and its effects, in Section 6.1.

Some previous work [23, 21, 24] talked about fusing a tree transducer for parsing with a transformation, but not one as expressive as an mft. More details are given in Section 7. Engelfriet and Vogler [6] described how to fuse a *finitary* tdt and a *macro tree transducer* (mtt). Their method, however, does not apply directly to our application because *Parse* has a infinite number of states. Our derivation from an mft to an xsp, to be presented in this section, is basically Engelfriet and Vogler's transducer fusion extended to mft's and specialised to one particular infinitary tdt, *Parse*. The readers are not required to have knowledge of their method. We will give a new proof of the fusion in Appendix A.

The rationale behind the derivation is as follows. Consider an mft  $M = (Q, \Sigma, \Delta, in, R)$ . For every state  $q \in Q$ , we introduce in the derived xsp a set of states  $\{q[i] \mid i \in \mathbb{N}^+\}$ . Imagine that we are building forests as we read the input stream of events. With each start tag, the forest construction descends by one level. The state  $q[1]$  performs the task that the state  $q$  in the mft is supposed to do. The number 1 indicates that the current forest will be its input. The states  $q[i]$  for  $i > 1$ , on the other hand, represent 'suspended' states which will take effect  $i - 1$  levels *above* the forest currently being built. The number  $i$  denotes the number of end tags expected. When an end tag is read, the number decrease by one, until the number reaches 1 and the state gets activated. When a start tag is read, the number shall increase by one because there is one more start tag to be matched.

**Definition 8.** Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft. We define an xsp  $\mathcal{SP}(M) = (Q', \Sigma, \Delta, in', R')$  where

- $Q' = \{q[i] \mid q \in Q, i \in \mathbb{N}^+\}$  where  $rank(q[i]) = rank(q) - 1$  for every  $q \in Q$  and  $i \in \mathbb{N}^+$ ,
- $in' = in[1] \in Q'$ ,
- $R'$  contains rules introduced by the following three cases:
  - xsp-(1).** for all  $q \in Q$  and  $\sigma \in \Sigma$ , we introduce

$$q[1](y_1, \dots, y_n) \xrightarrow{\langle \sigma \rangle} \mathcal{A}(rhs^{q, \sigma}),$$

- xsp-(2).** for all  $q \in Q$ ,  $\sigma \in \Sigma$  and  $i \in \mathbb{N}^+$ , we introduce:

$$q[1](y_1, \dots, y_n) \xrightarrow{\langle / \sigma \rangle} \mathcal{A}(rhs^{q, \epsilon}),$$

$$q[i](y_1, \dots, y_n) \xrightarrow{\$} \mathcal{A}(rhs^{q, \epsilon}),$$

- xsp-(3).** for all  $q \in Q$ ,  $\sigma \in \Sigma$  and  $i > 1$ , we introduce:

$$q[i](y_1, \dots, y_n) \xrightarrow{\langle \sigma \rangle} q[i+1](y_1, \dots, y_n),$$

$$q[i](y_1, \dots, y_n) \xrightarrow{\langle / \sigma \rangle} q[i-1](y_1, \dots, y_n).$$

The translation  $\mathcal{A}$  is defined by:

$$\mathcal{A}(q(x_i, rhs_1, \dots, rhs_n)) = q[i](\mathcal{A}(rhs_1), \dots, \mathcal{A}(rhs_n)),$$

$$\mathcal{A}(\epsilon) = \epsilon,$$

$$\mathcal{A}(y_j) = y_j,$$

$$\mathcal{A}(\delta \langle rhs \rangle) = \langle \delta \rangle \mathcal{A}(rhs) \langle / \delta \rangle,$$

$$\mathcal{A}(rhs \ rhs') = \mathcal{A}(rhs) \ \mathcal{A}(rhs'),$$

where  $q \in Q$ ,  $n = rank(q)$ ,  $\delta \in \Delta$ ,  $i \in \{1, 2\}$  and  $j \in \{1, \dots, n\}$ .

Note that among the three cases of rule introduction, **xsp-(1)** covers the situation when the state and the input symbols are  $(q[1], \langle \sigma \rangle)$ ; **xsp-(2)** covers  $(q[1], \langle / \sigma \rangle)$  and  $(q[i], \$)$  for  $i \in \mathbb{N}^+$ ; and **xsp-(3)** covers  $(q[i], \langle \sigma \rangle)$  and  $(q[i], \langle / \sigma \rangle)$  for  $i > 1$ . Therefore, the derived xsp  $\mathcal{SP}(M)$  is total if  $M$  is. For the examples below, we define a predicate testing whether the state and the input symbols is in the **xsp-(2)** case:

$$\epsilon_{\Sigma}(i, \chi) = (i = 1 \wedge \chi = \langle / \sigma \rangle) \vee \chi = \$, \text{ with } \sigma \in \Sigma.$$

The correctness of the derivation is stated by the following theorem whose proof is given in Appendix A.

**Theorem 1.** Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft. Then  $\tau_{\mathcal{SP}(M)}(\lfloor f \rfloor) = \lfloor \tau_M(f) \rfloor$  for every  $f \in \mathcal{F}_{\Sigma}$ .

*Example 3.* Apply the derivation to EXAMPLE 1, we get  $\mathcal{SP}(M_{rev}) = (Q', \Sigma, \Sigma, Main[1], R')$ , where  $Q' = \{q[i] \mid q \in \{Main, Rev\}, i \in \mathbb{N}^+\}$  and the set  $R'$  of rules is:

$$\begin{aligned} R' = \{ & Main[1]() \xrightarrow{\langle \mathbf{r} \rangle} \langle \mathbf{r} \rangle Rev[1](\varepsilon) \langle / \mathbf{r} \rangle Main[2](), & Rev[1](y_1) \xrightarrow{\langle \sigma \rangle} Rev[2](\langle \sigma \rangle Rev[1](\varepsilon) \langle / \sigma \rangle y_1) \\ & Main[1]() \xrightarrow{\langle \sigma \rangle} \langle \sigma \rangle Main[1]() \langle / \sigma \rangle Main[2](), & (\sigma \in \Sigma), \\ & (\sigma \neq \mathbf{r}), & \\ & Main[i]() \xrightarrow{\langle \sigma \rangle} Main[i+1]() \quad (\sigma \in \Sigma, i > 1), & Rev[i](y_1) \xrightarrow{\langle \sigma \rangle} Rev[i+1](y_1) \quad (\sigma \in \Sigma, i > 1), \\ & Main[i]() \xrightarrow{\langle / \sigma \rangle} Main[i-1]() \quad (\sigma \in \Sigma, i > 1), & Rev[i](y_1) \xrightarrow{\langle / \sigma \rangle} Rev[i-1](y_1) \quad (\sigma \in \Sigma, i > 1), \\ & Main[i]() \xrightarrow{\chi} \varepsilon \quad (\text{if } \epsilon_{\Sigma}(i, \chi)), & Rev[i](y_1) \xrightarrow{\chi} y_1 \quad (\text{if } \epsilon_{\Sigma}(i, \chi)) \}. \end{aligned}$$

Figure 3 shows a sample run when the input is  $\langle \mathbf{a} \rangle \langle \mathbf{r} \rangle \langle \mathbf{b} \rangle \langle \mathbf{c} \rangle \langle / \mathbf{c} \rangle \langle \mathbf{d} \rangle \langle / \mathbf{d} \rangle \langle \mathbf{e} \rangle \langle / \mathbf{e} \rangle \langle \mathbf{r} \rangle \langle \mathbf{f} \rangle \langle / \mathbf{f} \rangle \langle \mathbf{a} \rangle \$$ .

*Example 4.* The xsp derived from EXAMPLE 2 is  $\mathcal{SP}(M_{htm}) = (Q', \Sigma, \Delta, Main[1], R')$ , where  $Q' = \{q[i] \mid q \in Q, i \in \mathbb{N}^+\}$  and  $R'$  is

$$\begin{aligned} R' = \{ & Main[1]() \xrightarrow{\langle \mathbf{article} \rangle} \langle \mathbf{html} \rangle \langle \mathbf{head} \rangle Title[1]() \langle / \mathbf{head} \rangle \langle \mathbf{body} \rangle InArticle[1](\varepsilon) \langle / \mathbf{body} \rangle \langle / \mathbf{html} \rangle, \\ & Title[1]() \xrightarrow{\langle \mathbf{title} \rangle} \langle \mathbf{title} \rangle Copy[1]() \langle / \mathbf{title} \rangle, \\ & InArticle[1](y_1) \xrightarrow{\langle \mathbf{title} \rangle} \langle \mathbf{h1} \rangle Copy[1]() \langle / \mathbf{h1} \rangle InArticle[2](y_1), \\ & InArticle[1](y_1) \xrightarrow{\langle \mathbf{para} \rangle} \langle \mathbf{p} \rangle Key2Em[1]() \langle / \mathbf{p} \rangle InArticle[2](y_1 AllKeys[1]()), \\ & InArticle[1](y_1) \xrightarrow{\langle \mathbf{ps} \rangle} \langle \mathbf{h2} \rangle Index \langle / \mathbf{h2} \rangle \langle \mathbf{ul} \rangle y_1 \langle / \mathbf{ul} \rangle \langle \mathbf{h2} \rangle Postscript \langle / \mathbf{h2} \rangle Copy[1](), \\ & Key2Em[1]() \xrightarrow{\langle \mathbf{key} \rangle} \langle \mathbf{em} \rangle Copy[1]() \langle / \mathbf{em} \rangle Key2Em[2](), \\ & Key2Em[1]() \xrightarrow{\langle \sigma \rangle} \langle \sigma \rangle Key2Em[1]() \langle / \sigma \rangle Key2Em[2]() \quad (\sigma \neq \mathbf{key}), \\ & AllKeys[1]() \xrightarrow{\langle \mathbf{key} \rangle} \langle \mathbf{li} \rangle Copy[1]() \langle / \mathbf{li} \rangle AllKeys[2](), \\ & AllKeys[1]() \xrightarrow{\langle \sigma \rangle} AllKeys[1]() AllKeys[2]() \quad (\sigma \neq \mathbf{key}), \\ & Copy[1]() \xrightarrow{\langle \sigma \rangle} \langle \sigma \rangle Copy[1]() \langle / \sigma \rangle Copy[2]() \quad (\sigma \in \Sigma), \\ & q[i]() \xrightarrow{\langle \sigma \rangle} q[i+1]() \quad (\sigma \in \Sigma, i > 1, q \neq InArticle), \\ & q[i]() \xrightarrow{\langle / \sigma \rangle} q[i-1]() \quad (\sigma \in \Sigma, i > 1, q \neq InArticle), \\ & q[i]() \xrightarrow{\chi} \varepsilon \quad (\text{if } \epsilon_{\Sigma}(i, \chi)), \\ & InArticle[i](y_1) \xrightarrow{\langle \sigma \rangle} InArticle[i+1](y_1) \quad (\sigma \in \Sigma, i > 1), \\ & InArticle[i](y_1) \xrightarrow{\langle / \sigma \rangle} InArticle[i-1](y_1) \quad (\sigma \in \Sigma, i > 1), \\ & InArticle[i](y_1) \xrightarrow{\chi} \varepsilon \quad ((\chi, i) \in \Sigma_{\varepsilon}) \}. \end{aligned}$$

$Main[1]() \xrightarrow{\langle a \rangle} \langle a \rangle Main[1]() \langle /a \rangle Main[2]()$   
 $\xrightarrow{\langle r \rangle} \langle a \rangle \langle r \rangle Rev[1](\varepsilon) \langle /r \rangle Main[2]() \langle /a \rangle Main[3]()$   
 $\xrightarrow{\langle b \rangle} \langle a \rangle \langle r \rangle Rev[2](\langle b \rangle Rev[1](\varepsilon) \langle /b \rangle) \langle /r \rangle Main[3]() \langle /a \rangle Main[4]()$   
 $\xrightarrow{\langle c \rangle} \langle a \rangle \langle r \rangle Rev[3](\langle b \rangle Rev[2](\langle c \rangle Rev[1](\varepsilon) \langle /c \rangle) \langle /b \rangle) \langle /r \rangle Main[4]() \langle /a \rangle Main[5]()$   
 $\xrightarrow{\langle /c \rangle} \langle a \rangle \langle r \rangle Rev[2](\langle b \rangle Rev[1](\langle c \rangle \langle /c \rangle) \langle /b \rangle) \langle /r \rangle Main[3]() \langle /a \rangle Main[4]()$   
 $\xrightarrow{\langle d \rangle} \langle a \rangle \langle r \rangle Rev[3](\langle b \rangle Rev[2](\langle d \rangle Rev[1](\varepsilon) \langle /d \rangle \langle c \rangle \langle /c \rangle) \langle /b \rangle) \langle /r \rangle Main[4]() \langle /a \rangle Main[5]()$   
 $\xrightarrow{\langle /d \rangle} \langle a \rangle \langle r \rangle Rev[2](\langle b \rangle Rev[1](\langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle) \langle /b \rangle) \langle /r \rangle Main[3]() \langle /a \rangle Main[4]()$   
 $\xrightarrow{\langle /b \rangle} \langle a \rangle \langle r \rangle Rev[1](\langle b \rangle \langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle) \langle /r \rangle Main[2]() \langle /a \rangle Main[3]()$   
 $\xrightarrow{\langle e \rangle} \langle a \rangle \langle r \rangle Rev[2](\langle e \rangle Rev[1](\varepsilon) \langle /e \rangle \langle b \rangle \langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle) \langle /r \rangle Main[3]() \langle /a \rangle Main[4]()$   
 $\xrightarrow{\langle /e \rangle} \langle a \rangle \langle r \rangle Rev[1](\langle e \rangle \langle /e \rangle \langle b \rangle \langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle) \langle /r \rangle Main[2]() \langle /a \rangle Main[3]()$   
 $\xrightarrow{\langle /r \rangle} \langle a \rangle \langle r \rangle \langle e \rangle \langle /e \rangle \langle b \rangle \langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle \langle /b \rangle \langle /r \rangle Main[1]() \langle /a \rangle Main[2]()$   
 $\xrightarrow{\langle f \rangle} \langle a \rangle \langle r \rangle \langle e \rangle \langle /e \rangle \langle b \rangle \langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle \langle /b \rangle \langle /r \rangle \langle f \rangle Main[1]() \langle /f \rangle Main[2]() \langle /a \rangle Main[3]()$   
 $\xrightarrow{\langle /f \rangle} \langle a \rangle \langle r \rangle \langle e \rangle \langle /e \rangle \langle b \rangle \langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle \langle /b \rangle \langle /r \rangle \langle f \rangle \langle /f \rangle Main[1]() \langle /a \rangle Main[2]()$   
 $\xrightarrow{\langle /a \rangle} \langle a \rangle \langle r \rangle \langle e \rangle \langle /e \rangle \langle b \rangle \langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle \langle /b \rangle \langle /r \rangle \langle f \rangle \langle /f \rangle \langle /a \rangle Main[1]()$   
 $\xrightarrow{\$} \langle a \rangle \langle r \rangle \langle e \rangle \langle /e \rangle \langle b \rangle \langle d \rangle \langle /d \rangle \langle c \rangle \langle /c \rangle \langle /b \rangle \langle /r \rangle \langle f \rangle \langle /f \rangle \langle /a \rangle$

Fig. 3. Stream processing induced by  $\mathcal{SP}(M_{rev})$

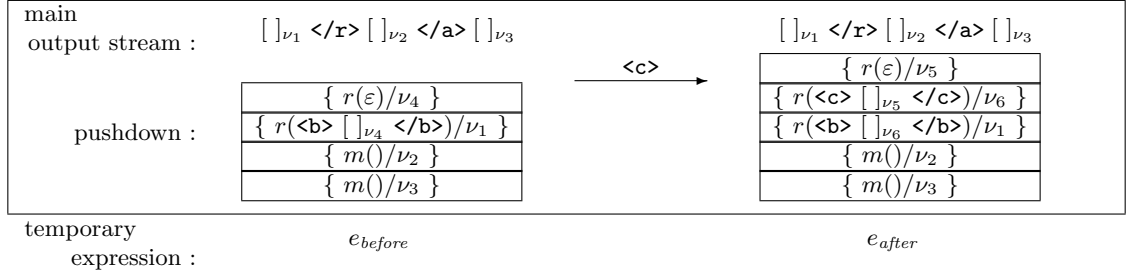


Fig. 4. Pushdown representation for temporary expressions and its updating

## 4 Pushdown XML Stream Processor

The semantics given in Section 3 implies a direct implementation of xsp performing term rewriting each time an event is read. However, an xsp derived from an mft follows a more regular evaluation pattern: most rules are introduced by **xsp-(3)**, which merely increment or decrement the indexes of the states, and term expansion happens only for states indexed 1, triggered by rules in **xsp-(1)**. Such behaviour resembles a stack. In this section, we present an efficient implementation of the xsp's derived from mft's, using a pushdown machine based representation of the temporary expressions.

### 4.1 Summary of Behavior

Let us look at an example first. Consider the sample run of the xsp  $\mathcal{SP}(M_{rev})$  in Figure 3, when event  $\langle c \rangle$  is read. We abbreviate  $Rev$  to  $r$  and  $Main$  to  $m$ . The prefix  $\langle a \rangle \langle r \rangle$  has been 'squeezed' to the output. We need only to keep a suffix in of the temporary expression that remains in memory:

$$e_{before} = r[2](\langle b \rangle r[1](\varepsilon) \langle /b \rangle) \langle /r \rangle m[3]() \langle /a \rangle m[4]()$$

After  $\langle c \rangle$  is read, the expression gets updated to

$$e_{after} = r[3](\langle b \rangle r[2](\langle c \rangle r[1](\varepsilon) \langle /c \rangle) \langle /b \rangle) \langle /r \rangle m[4]() \langle /a \rangle m[5]() .$$

We shall present a data structure such that the update can be performed efficiently.

We represent a temporary expression by a pair of a *main output stream* and a *pushdown*, as shown in Figure 4. The left and right parts in the figure correspond to temporary expressions  $e_{before}$  and  $e_{after}$ , respectively. Consider  $e_{before}$ . Separating the evaluated and unevaluated segments, it can be partitioned into five parts:  $r[2](\dots)$ ,  $\langle /r \rangle$ ,  $m[3]()$ ,  $\langle /a \rangle$  and  $m[4]()$ . If we abstract away the unevaluated parts and replace them with *holes*  $[ ]_{\nu_i}$  using a physical address  $\nu_i$ , we obtain the main output stream  $[ ]_{\nu_1} \langle /r \rangle [ ]_{\nu_2} \langle /a \rangle [ ]_{\nu_3}$ .

The pushdown is a stack of sets, each set consisting of *state frames*. A state frame is a pair of a state  $q(\dots)$  and a hole address  $\nu$ , denoted by  $q(\dots)/\nu$ . The state may have a number of arguments, represented by a sequence in a way similar to the main output stream.

In the pushdown representation, every state  $q[i]$  appears in the  $i$ -th set from the top. Therefore the index  $i$  need not be stored in the representation. Since all states in  $e_{before}$  have distinct indexes, the pushdown contains only singleton sets, which need not be true in general.

Only the states with index 1 gets expanded. In our representation, that means we only need to update the top of the pushdown. Upon reading  $\langle c \rangle$ , the rule of  $r[1]$  that gets triggered is

$$r[1](\varepsilon) \xrightarrow{\langle c \rangle} r[2](\langle \sigma \rangle r[1](\varepsilon) \langle / \sigma \rangle \varepsilon) .$$

That corresponds to popping the set  $\{r(\varepsilon)/\nu_4\}$  (representing  $r[1](\varepsilon)$  in  $e_{before}$ ), and pushing the two sets  $\{r(\varepsilon)/\nu_5\}$  (representing  $r[1](\varepsilon)$  in  $e_{after}$ ) and  $\{r(\langle c \rangle [ ]_{\nu_5} \langle /c \rangle)/\nu_6\}$  (representing  $r[2](\langle \sigma \rangle \dots \langle / \sigma \rangle)$  in  $e_{after}$ ). Now that  $\nu_4$  is expanded, all occurrences of  $[ ]_{\nu_4}$  in the pushdown should be filled with  $[ ]_{\nu_6}$ . Since two items are pushed, all other sets in the pushdown descend for one level. This corresponds to updating all states  $q[i]$  ( $i > 1$ ) to  $q[i + 1]$  at the same time.

In this example, the main output stream remains the same after the updating. In general, the main output stream will be altered when it contains an address referred by the top set of the pushdown, or when its prefix is ready to be squeezed.

## 4.2 Pushdown Representation and its Updating

Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft. An *output stream*  $s$  for  $M$  is defined by

$$m ::= \varepsilon \mid \langle \delta \rangle m \mid \langle / \delta \rangle m \mid [ ]_{\nu} m ,$$

where  $\delta \in \Delta$ , and  $[ ]_{\nu}$  is a hole whose physical address is  $\nu$ . We denote the set of output streams by  $\mathcal{S}_M$ . A *state frame* has the form  $q(m_1, \dots, m_n)/\nu$  where  $\nu$  is a hole address,  $q \in Q$ ,  $n = rank(q)$ , and  $m_i \in \mathcal{S}_M$  ( $i = 1, \dots, n$ ).

A *pushdown* is a mapping from a positive number, representing the depth, to a set of state frames. Furthermore, each hole address  $\nu$  occurs on the right-hand side of  $/$  in a pushdown at most once. The empty pushdown is denoted by  $\emptyset$ . Given a set of state frames  $\Psi$ , we denote by  $\{1 \mapsto \Psi, \dots\}$  a pushdown  $p$  such that  $p(1) = \Psi$ . Two pushdowns  $p_1$  and  $p_2$  can be merged by  $p_1 \oplus p_2 = \{d \mapsto p_1(d) \uplus p_2(d)\}_{d \in \mathbb{N}^+}$ <sup>4</sup>.

Now we define formally how a temporary expression is represented.

**Definition 9.** Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft. A *pushdown representation*  $pd(e)$  for  $e \in \mathcal{Temp}_{\mathcal{SP}(M)}$  is a pair  $\langle m, p \rangle$  of a *main output stream*  $m$  and a *pushdown*  $p$  defined by

$$\begin{aligned} pd(\varepsilon) &= \langle \varepsilon, \emptyset \rangle, \\ pd(q[i](e_1, \dots, e_n) e) &= \\ &\quad \langle [ ]_{\nu} m, \{i \mapsto \{q(m_1, \dots, m_n)/\nu\}\} \oplus p_1 \oplus \dots \oplus p_n \oplus p \rangle, \\ pd(\langle \delta \rangle e) &= \langle \langle \delta \rangle m, p \rangle, \\ pd(\langle / \delta \rangle e) &= \langle \langle / \delta \rangle m, p \rangle, \end{aligned}$$

<sup>4</sup> For any  $d \in \mathbb{N}^+$ , two sets  $p_1(d)$  and  $p_2(d)$  are always disjoint because of the uniqueness of hole addresses.

where  $\langle m, p \rangle = pd(e)$ ,  $\langle m_i, p_i \rangle = pd(e_i)$  and  $\nu$  is a fresh address. We denote the set of pushdown representations for temporary expressions in  $Tmp_{SP(M)}$  by  $Pdr_M$ .

From a pushdown representation, we can recover the temporary expression by filling every hole according to the corresponding state frame in the pushdown. The function  $pd^{-1}$  is defined by

$$\begin{aligned} pd^{-1}(\langle \varepsilon, p \rangle) &= \varepsilon \\ pd^{-1}(\langle \langle \delta \rangle m, p \rangle) &= \langle \delta \rangle pd^{-1}(\langle m, p \rangle) \\ pd^{-1}(\langle \langle / \delta \rangle m, p \rangle) &= \langle / \delta \rangle pd^{-1}(\langle m, p \rangle) \\ pd^{-1}(\langle [ ]_\nu m, p \rangle) &= (p\% \nu) pd^{-1}(\langle m, p \rangle) \end{aligned}$$

where  $p\% \nu$  accesses the corresponding state frame in the pushdown and expands the arguments of the state. When a state frame  $q(m_1, \dots, m_n)/\nu$  is an element of the set  $p(i)$ , we have  $p\% \nu = q[i](pd^{-1}(\langle m_1, p \rangle), \dots, pd^{-1}(\langle m_n, p \rangle))$ . For any temporary expression  $e$ , we have

$$e = pd^{-1}(pd(e)) \quad (2)$$

whose proof is omitted.

We define several operations to manipulate the pushdown representation. An *application* for a hole  $[ ]_\nu$  in an output stream  $m$  with another output stream  $u$  is denoted by  $m@_\nu u$ , i.e., when  $m = m_1 [ ]_\nu m_2$ , we have  $m@_\nu u = m_1 u m_2$ . The hole application can be extended to a set of state frames and a pushdown in the same way, denoted by  $\Psi@_\nu u$  and  $p@_\nu u$ . Let  $p$  be a pushdown and  $\Psi$  a set of state frames. The pushdown obtained by pushing  $\Psi$  on the top of  $p$  is denoted by  $p \ll \Psi = \{1 \mapsto \Psi\} \cup \{d \mapsto p(d-1)\}_{d>1}$ . The dual operation popping the top of  $p$  is denoted by  $\triangleright p = \{d \mapsto p(d+1)\}_{d \in \mathbb{N}^+}$ .

The hole application operation can be efficiently implemented in the sense that the execution time is independent of the size of main output streams and pushdowns. because we know the physical address of the hole. Experimental implementation introduced in Section 5 uses doubly-linked cyclic lists to represent output streams, so we can implement hole application, concatenation and squeeze efficiently.

### 4.3 Pushdown Machines for Macro Forest Transducers

For a given mft  $M$ , we introduce a pushdown machine in stream processing style which simulates the behavior of the xsp  $SP(M)$ .

Since the semantics of an xsp is specified by a transition  $\langle \_, \_ \rangle$  on temporary expressions, we construct the pushdown machine as a transition on pushdown representations. In the following definition, the function  $pd^\circ$  extends  $pd$  by one extra case,  $pd^\circ(y_i) = y_i$  for  $i \in \mathbb{N}^+$ . Therefore  $pd^\circ$  can be applied to the right-hand side of rules in an xsp.

**Definition 10.** Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft. The *pushdown machine for  $M$* , denoted by  $\mathcal{PD}(M)$ , is a function  $\langle \_, \_ \rangle : Pdr_M \times \Sigma_{\langle \rangle \S} \rightarrow Pdr_M$ . For a pushdown representation  $\langle m, p \rangle \in Pdr_M$  and input event  $\chi \in \Sigma_{\langle \rangle \S}$ , a new pushdown representation  $\langle \langle m, p \rangle, \chi \rangle$  is given as follows:

- $\langle \langle m, p \rangle, \langle \sigma \rangle \rangle = \Phi_\sigma(m, (\triangleright p) \ll \emptyset \ll \emptyset, p(1))$ , where the function  $\Phi_\sigma$  is defined by

$$\begin{aligned} \Phi_\sigma(m, p, \emptyset) &= \langle m, p \rangle \\ \Phi_\sigma(m, p, \{q(m_1, \dots, m_n)/\nu\} \uplus \Psi) &= \\ &\Phi_\sigma(m@_\nu m', p@_\nu m' \oplus p', \Psi@_\nu m') \end{aligned}$$

with  $\langle m', p' \rangle = pd^\circ(\mathcal{A}(rhs^{q,\sigma}))[y_j := m_j]_{j=1, \dots, n}$ .

- $\langle \langle m, p \rangle, \langle / \sigma \rangle \rangle = \Phi_\epsilon(m, \triangleright p, p(1))$ , where the function  $\Phi_\epsilon$  is defined by

$$\begin{aligned} \Phi_\epsilon(m, p, \emptyset) &= \langle m, p \rangle \\ \Phi_\epsilon(m, p, \{q(m_1, \dots, m_n)/\nu\} \uplus \Psi) &= \\ &\Phi_\epsilon(m@_\nu m', p@_\nu m' \oplus p', \Psi@_\nu m') \end{aligned}$$

with  $\langle m', p' \rangle = pd^\circ(\mathcal{A}(rhs^{q,\epsilon}))[y_j := m_j]_{j=1, \dots, n}$ .

–  $\langle\langle m, p \rangle, \$ \rangle = \Phi_\epsilon(m, \emptyset, \bigcup_{d \in \mathbb{N}^+} p(d))$ , where  $\Phi_\epsilon$  is as given in the previous case of  $\chi = \langle / \sigma \rangle$ .

For an mft  $M = (Q, \Sigma, \Delta, in, R)$ , the initial pushdown representation of  $\mathcal{PD}(M)$  is

$$\langle [ ]_{\nu_0}, \{1 \mapsto \{in(\epsilon, \dots, \epsilon)/\nu_0\}\} \rangle$$

with address  $\nu_0$ . It corresponds to the initial state of an xsp  $\mathcal{SP}(M)$ , that is,  $in[1](\epsilon, \dots, \epsilon)$ . For a pushdown machine  $P = \mathcal{PD}(M)$ , the transformation  $\tau_P : \Sigma_\diamond^\circ \rightarrow \Delta_\diamond^\circ$  induced by  $P$  is defined in a way similar to  $\tau_{\mathcal{SP}(M)}$ :

$$\tau_P(s) = \zeta_P(\langle [ ]_{\nu_0}, \{1 \mapsto \{in(\epsilon, \dots, \epsilon)/\nu_0\}\} \rangle, s\$)$$

where

$$\zeta_P(\langle m, p \rangle, \epsilon) = m \qquad \zeta_P(\langle m, p \rangle, \chi s) = \zeta_P(\langle\langle m, p \rangle, \chi \rangle, s)$$

for a pushdown representation  $\langle m, p \rangle$ .

For an mft  $M$ , the behaviour of  $\mathcal{PD}(M)$  on pushdown representations mirrors that of  $\mathcal{SP}(M)$  on temporary expressions. Consider the case when a start tag  $\langle \sigma \rangle$  is read. In the xsp  $\mathcal{SP}(M)$ , every state  $q[i]$  ( $i > 1$ ) is rewritten into  $q[i + 1]$ . In the pushdown machine  $\mathcal{PD}(M)$ , the corresponding state frame  $q(\dots)/\nu$  in the  $i$ -th set of the pushdown descends by one level because we perform one pop and two pushes on the pushdown. In  $\mathcal{SP}(M)$ , every state  $q[1]$  is rewritten by  $\mathcal{A}(rhs^{q,\sigma})$ . In the pushdown machine  $\mathcal{PD}(M)$ , for each corresponding state frame  $q(\dots)/\nu$  in the top set of the pushdown, the hole  $[ ]_\nu$  is filled according to  $\mathcal{A}(rhs^{q,\sigma})$ . Since a computation of  $pd^\circ(\mathcal{A}(rhs^{q,\sigma}))$  is invoked, the state  $q[1]$  in  $\mathcal{A}(rhs^{q,\sigma})$  is put as an element of the top set of the pushdown and  $q[2]$  in  $\mathcal{A}(rhs^{q,\sigma})$  is put as an element of the second set from the top.

Consider the case when an end tag  $\langle / \sigma \rangle$  is read. In the xsp  $\mathcal{SP}(M)$ , every state  $q[i]$  ( $i > 1$ ) is rewritten to  $q[i - 1]$ . In the pushdown machine  $\mathcal{PD}(M)$ , the corresponding state frame  $q(\dots)/\nu$  in the  $i$ -th set of the pushdown ascends by one level after popping the pushdown. In the xsp  $\mathcal{SP}(M)$ , every state  $q[1]$  is replaced according to  $\mathcal{A}(rhs^{q,\epsilon})$ . In the pushdown machine  $\mathcal{PD}(M)$ , for the corresponding state frame  $q(\dots)/\nu$  in the top set of the pushdown, the hole  $[ ]_\nu$  is filled according to  $\mathcal{A}(rhs^{q,\epsilon})$ .

After reading the end of stream, the pushdown must be empty since  $\mathcal{A}(rhs^{q,\epsilon})$  contains no pattern  $q[i](\dots)$  and all state frames in the previous pushdown is consumed by  $\Phi_\epsilon(s, \emptyset, \bigcup_{d \in \mathbb{N}^+} p(d))$ . Thereby the final main output stream does not contain any hole.

Since every transition on pushdown representations corresponds to a transition on temporary expressions, we can see that  $\tau_{\mathcal{PD}(M)}(s) = \tau_{\mathcal{SP}(M)}(s)$  for every mft  $M$  and every input stream  $s$ . From THEOREM 1, we have

$$\tau_{\mathcal{PD}(M)}(\lfloor f \rfloor) = \lfloor \tau_M(f) \rfloor$$

for every input forest  $f$  for  $M$ , which shows the equivalence of the original mft and the derived pushdown machine.

The above definition of  $\tau_P$  for a pushdown machine  $P$  can be made more efficient by *squeezing*, that is, printing out the prefix, up to the first hole, of the main output stream, since they are not going to change by subsequent transitions. We define the following function  $sqz$  for output streams:

$$\begin{aligned} sqz(\epsilon) &= (\epsilon, \epsilon) & sqz([ ]_\nu m) &= (\epsilon, [ ]_\nu m) \\ sqz(\langle \delta \rangle m) &= (\langle \delta \rangle m', m'') & sqz(\langle / \delta \rangle m) &= (\langle / \delta \rangle m', m'') \end{aligned}$$

where  $(m', m'') = sqz(m)$ . The function  $sqz$  splits the leading output events and the rest of the output stream, i.e., if  $(m', m'') = sqz(m)$ , then  $m' m'' = m$ .

We can then redefine  $\zeta_P$  with  $sqz$  as follows:

$$\zeta_P(\langle m, p \rangle, \epsilon) = m \qquad \zeta_P(\langle m, p \rangle, \chi s) = m' \zeta_P(\langle\langle m'', p \rangle, \chi \rangle, s)$$

where  $(m', m'') = sqz(\langle\langle m, p \rangle, \chi \rangle)$ .

input size	1MB	4MB	16MB	64MB	256MB
pushdown xsp	0.70sec / 1.14MB	2.71sec / 1.14MB	10.8sec / 1.14MB	43.4sec / 1.14MB	197sec / 1.14MB
direct impl. mft	0.77sec / 3.72MB	3.08sec / 11.5MB	12.9sec / 43.5MB	61.1sec / 171MB	466sec / 609MB
xsltproc	0.77sec / 8.75MB	3.41sec / 33.2MB	19.3sec / 129MB	160sec / 511MB	n/a
saxon	2.48sec / 23.1MB	4.92sec / 39.9MB	14.1sec / 112MB	49.8sec / 403MB	n/a

(execution time / max. memory usage)

**Table 1.** Performance comparison for varying input size

*Remark.* In the actual implementation, we have to take care of dealing with mft rules whose right-hand side does not contain exactly one occurrence of parameter variable  $y_j$  for each  $j$ .

Consider the case that there are some variables which does not occur in the right-hand side. For example, let an mft have a rule  $q(pat, y_1) \rightarrow \epsilon$ . Then the corresponding xsp has a rule  $q[1](y_1) \xrightarrow{\chi} \epsilon$ . When the top set of the pushdown contains  $q(m_1)/\nu$ , the occurrence of  $[ ]_\nu$  will be filled with  $\epsilon$ . In order to avoid ineffectual updating after that, all hole addresses contained in  $m_1$  should be discarded as long as the hole does not occur in the other position.

Consider the case that there are some variables which occurs more than once in the right-hand side. If we use doubly-linked cyclic lists to represent main output streams for efficient updating, we should be carefully deal with them. For example, let an mft have a rule  $q(pat, y_1) \rightarrow y_1 y_1$ . Then the corresponding xsp has a rule  $q[1](y_1) \xrightarrow{\chi} y_1 y_1$ . In this case, a hole  $[ ]_\nu$  may occur twice. When the hole  $[ ]_\nu$  is required to be filled, we cannot replace both occurrence of  $\nu$  with the same doubly-linked list. To solve this problem, we mark the state frame so as to remember it appears twice.

## 5 Benchmarking Results

For benchmarking, we use the random sample generator XMark [33] to produce sample XML documents of various sizes: 1MB, 4MB, 16MB, 64MB and 256MB. Every sample document contains a sequence of `item` nodes, each having a list of children about a dozen lines long. The benchmarking task would be to reverse the order of the subtrees under `item`.

The pushdown machine automatically derived from the mft specification of  $M_{rev}$ , shown as the entry **pushdown xsp** in Table 1, is implemented in Objective Caml, with extensions to handle text nodes. The entry **direct impl. mft** is the program in Figure 1 implemented as mutual recursive functions in Objective Caml. The entry **xsltproc** is one of the fastest XSLT processors bundled with `libxslt` [31], written in C, while **saxon** [13] is one of the fastest XSLT processors in Java. All entries apart from **pushdown xsp** build the entire forest in memory before the transformation. The experiments were conducted on a 1.33 GHz PowerBook G4 with 768 MB of memory. Table 1 compares the total execution time and maximum memory size in seconds and megabytes.

As we expected, **pushdown xsp** uses a much smaller heap comparing too all other entries. That it also outperforms the two XSLT processors may be due to the overhead of the latter maintaining full-fledged XML data, including e.g., namespace URI, number of children, etc. For a fairer comparison, we added the entry **direct impl. mft**. The entry **pushdown xsp** is slightly faster than **direct impl. mft** because it incurs less garbage collection and xsp saves the overhead of building the trees. Therefore we expect that the xsp approach will also deliver competitive speed even after being scaled to handle full XML.

## 6 Discussion

This section discusses capabilities of our framework. Firstly we show the difference between lazy evaluation mechanism and ours. Next possibilities of our framework for existing XML transformation languages are discussed. Finally we describe some limitation of stream processing.

## 6.1 Comparison with Lazy Evaluation

Many of our readers wondered: “Is this all necessary? Can we not just use lazy evaluation?” Consider the program *unparse* (*trans* (*parse* *input*)) in a non-strict functional language, where the function *parse* builds the tree lazily upon the demand of the forest-to-forest transformation *trans*. When the program is run by a lazy evaluator, do we get the desired space behaviour?

To answer the question, we run a number of experiments in the functional language Haskell. We start with choosing a parser to be fused with. The *tdtt*-based parser in Section 3.2 is not a good choice because sharing the input stream *s* among functional causes a space leak. Another possibility is to thread the input stream through recursive calls to the parser, and have the parser return the tree as well as an unprocessed tail of the stream, such as:

$$\begin{aligned} \text{parse}(\langle a \rangle s) = & \text{let } (ts, s') = \text{parse } s \\ & (us, s'') = \text{parse } s' \\ & \text{in } (a \langle ts \rangle us, s'') \end{aligned}$$

The intention is that the input stream can be freed right after being used.

Unfortunately, the space behaviour of the parser is compiler-dependent. It runs fine when compiled with NHC98 [22]. With GHC [7], however, the entire input stream resides in memory throughout the execution of the program. The reason is that *s''* may be viewed as *snd* (*parse* (*snd* (*parse* *s*))), so while the recursive calls are being executed, the system still keeps a pointer to *s* until *s''* is finally demanded by the environment. This problem was addressed by Wadler [29], who proposed a fix in the garbage collector to plug the space leak. The fix is actually implemented in both NHC and GHC, but is fragile in presence of otherwise-valuable optimisations that GHC performs [12].

One may argue that the problem is specific to this particular parser. EXAMPLE 2, however, shows a problem more intrinsic to the nature of lazy evaluation. The thunk that evaluates to the list of keywords appears very late and remain unevaluated until it is finally output. This is in fact what we expect of lazy evaluation – if the programmer applies the function *take* to the generated stream, the thunk need not be evaluated at all. However, the thunk contains a reference to the beginning of the input stream, which means that the entire input stream will reside in memory.

For every *xsp* example we have, we were able to eventually come up with a Haskell equivalent that uses a minimal amount of memory, by inserting strictness annotations after careful tuning and heap profiling, which was no easy task.

Put it in a wider context, we recall Wadler’s claim [29] that we need a parallel evaluator to avoid certain classes of space leaks. This clarifies the relationship between lazy evaluation and *xsp*. Our *xsp* implementation, which evaluates all the states *q*[1] indexed 1, can actually be seen as a parallel evaluator specialised for XML processing.

## 6.2 Streaming for Existing XML Transformation Languages

Is it possibilities to apply our framework to existing transformation languages? The problem reduces to whether we can convert functions defined in languages such as XSLT [34], *fxt* [3], XDuce [11], or CDuce [2], probably with some restrictions imposed, into *mft*’s.

The XML transformation language TL [17] by Maneth, Berlea, Perst and Seidl may give us some hints about what the restrictions could be. In TL programmers also define a collection of mutual recursive functions. TL is like *mft*, but supports pattern matching by *monadic second-order logic* (MSO) formulae. Each rule of TL has the form of  $q(\phi, y_1, \dots, y_n) \rightarrow rhs$ , where  $\phi$  is an MSO formula. When *q* is called, the nodes satisfying  $\phi$  is passed as it argument. Maneth et al. showed that most practical TL programs use only MSO formulae that does not select ancestor nodes, and such programs can be represented in a single deterministic *mft*. It implies that a simple XSLT program which contains only forward XPath expressions can be expressed as an *mft* and implemented by our pushdown machine.

XDuce and CDuce support regular expression pattern [10]. They do not bind variables to the parent of the current node. It is easier to see how *tail-capturing* programs can be translated to *mft*’s. For example, the following XDuce program converting an address book to a telephone list:

```

fun mkTelList (val e as (Name,Addr,Tel?)* =
  match e with
    name[val n], addr[val a], tel[val t], val rest
    -> name[n], tel[t], mkTelList (rest)
  | name[val n], addr[val a], val rest -> mkTelList (rest)
  | () -> ()

```

can be defined by an mft with rules:

$$\begin{aligned}
MkTelList(\mathbf{name}\langle x_1 \rangle x_2) &\rightarrow Name(x_2, \mathbf{name}\langle Val(x_1) \rangle) \\
MkTelList(\epsilon) &\rightarrow \epsilon \\
Name(\mathbf{addr}\langle x_1 \rangle x_2, y_1) &\rightarrow NameAddr(x_2, y_1) \\
NameAddr(\mathbf{tel}\langle x_1 \rangle x_2, y_1) &\rightarrow y_1 \mathbf{tel}\langle Val(x_1) \rangle MkTelList(x_2) \\
NameAddr(\mathbf{name}\langle x_1 \rangle x_2, y_1) &\rightarrow Name(x_2, \mathbf{name}\langle Val(x_1) \rangle) \\
NameAddr(\epsilon, y_1) &\rightarrow y_1
\end{aligned}$$

Here we extend mft's to handle text data, and  $Val$  is the identity function for text. This mft is total if inputs are restricted to the type  $(Name, Addr, Tel?)*$  specified by the original XDuce program. Hosoya and Pierce [10] talked about how to convert non-tail-capturing patterns into tail-capturing equivalents. It will be among our future work to see how this approach works in general.

The mft can be extended to handle primitive datatype other than forests. For example, we can extend the right-hand side with booleans, boolean operators, and conditional branches as follows:

$$rhs ::= \dots \mid true \mid false \mid if(rhs, rhs, rhs)$$

The corresponding extension we need in the xsp is some extra rules [21]:

$$if(true, e_1, e_2) \rightarrow e_1 \qquad if(false, e_1, e_2) \rightarrow e_2$$

By applying these rules in each squeezing phase, we achieve an XML stream processing for the extended mft. Additionally we need to keep the position of  $if$  as a state frame which always locates in the top set of a pushdown. Once we have booleans and conditionals in mft, we can express many transformations, including the `invite/visit` iteration with XPath expressions in XTISP [19], the first author's previous work.

### 6.3 Limitation

There is a class of *inherently memory inefficient* transformations [24]. For example, if we replace the first two rules of  $M_{rev}$  with a single rule:  $Main(\sigma\langle x_1 \rangle x_2) \rightarrow Rev(x_2, \sigma\langle Rev(x_1, \epsilon) \rangle)$ . The mft (call it  $M_{frev}$ ) reverses the subtrees at every nesting level for all nodes. Although we can still derive an xsp from it, the resulting xsp is actually a bit slower than the naive load-all implementation, because it cannot output any result until reading the end of the input stream. Every SAX-like stream processing program has the same problem: this kind of transformation is just not suitable for stream processing.

As a trial experiment, the table below compares  $M_{frev}$  and a program **direct impl. mft** which simply loads the tree and performs the full reverse. The result shows that our implementation does not carry too much overhead even for an inherently inefficient transformation  $M_{frev}$ .

input size	4MB	64MB
pushdown xsp	3.26sec / 12.3MB	61.5sec / 188MB
direct impl. mft	3.08sec / 11.5MB	60.3sec / 170MB

(execution time / max. memory usage)

## 7 Conclusion and Related Work

We have presented a method to automatically derive an XML stream processor from a program expressed as the macro forest transducer — mutual recursive functions on XML forests with certain constraints. The XML stream processor has an efficient implementation based on a pushdown machine. The framework presented in this paper will be the core of the next release of XTiSP [19]. We believe that the mft is expressive enough that we can transform most practical programs written in existing XML processing languages [11, 2, 34] to mft, in order to streamlise them. That will be one of our future work too.

While plenty of work has been devoted to the automatic derivation of XML stream processors from declarative programs, most of them deals with query languages, such as XPath [1, 5, 8, 9] and a subset of XQuery [16]. They are not expressive enough to describe some transformation we would like to have, such as the structure-preserving transformation renaming all the labels **a** to **b**, which can be expressed naturally in recursive functional style.

The key idea of our framework was presented in the first author's previous work [21, 24], based on the composition of (stack-)attributed tree transducers (att) [20]. All programs definable in XTiSP, an XML transformation language designed by the first author [21, 19], can be translated into att's, which can then be composed with an XML parser in a way similar our derivation of xsp's. However, it is well known that att's are less expressive than mft's [6, 26]. Our result in this paper is therefore much more powerful than before. Moreover, without a formal model of stream processors, some part of the implementation of XTiSP is rather ad-hoc and could have been made more efficient. The formalisation in the present helps to produce an implementation that is both correct and efficient.

Kodama, Suenaga, Kobayashi and Yonezawa [15, 28] proposed a translation from tree processing programs. Their tree processing language deals with binary trees which can be easily parsed without end tags. This is fundamentally different from XML processing, since we do not even need a pushdown stack to parse such trees.

Kiselyov [14] gave an XML parser with a general folding function `foldts` over rose trees. An XML transformation is defined by three actions `fup`, `fdown` and `fhere` that specify how to accumulate the seed value. This programming style is not user-friendly and many function closures are stored during the processing. Furthermore, his framework does not mention whether the processor can output a part of the result when reading a single XML event, e.g., a start tag `<a>`.

STX [4] is a template-based XML transformation language that operates on stream of SAX [32] events. While the programmers can define the XML transformation program as well as XSLT [34], they have to explicitly write when and how to store the temporary information like stream processing style.

TransformX presented by Scherzinger and Kemper [27] provides a framework for syntax-directed transformations of XML streams. We can obtain XML stream processors by defining a kind of attribute grammar on the regular tree of the type schema for inputs. Even in their framework, however, we must still keep in mind which information should be buffered before and after reading each subtree in the input.

## References

1. M. Altinel and M. J. Franklin. Efficient filtering of XML documents for selective dissemination of information. *International Journal on Very Large Data Bases*, pages 53–64, 2000.
2. V. Benzaken, G. Castagna, and A. Frisch. CDuce: an XML-centric general-purpose language. In *Proceedings of the 8th International Conference of Functional Programming*, pages 51–63, 2003.
3. A. Berlea and H. Seidl. fxt – a transformation language for XML documents. *Journal of Computing and Information Technology*, 10(1):19–35, 2002.
4. P. Cimprich, O. Becker, C. Nentwich, M. K. H. Jiroušek, P. Brown, M. Batsis, T. Kaiser, P. Hlavnička, N. Matsakis, C. Dolph, and N. Wiechmann. Streaming transformations for XML (STX) version 1.0. <http://stx.sourceforge.net/documents/>.

5. Y. Diao and M. J. Franklin. High-performance XML filtering: An overview of YFilter. In *IEEE Data Engineering Bulletin*, volume 26(1), pages 41–48, 2003.
6. J. Engelfriet and H. Vogler. Macro tree transducers. *Journal of Computer and System Sciences*, 31(1):71–146, 1985.
7. The Glasgow Haskell Compiler. <http://www.haskell.org/ghc/>.
8. T. J. Green, A. Gupta, G. Miklau, M. Onizuka, and D. Suciu. Processing XML streams with deterministic automata and stream indexes. *ACM Transactions on Database Systems*, 29(4):752–788, 2004.
9. A. K. Gupta and D. Suciu. Stream processing of XPath queries with predicates. In *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data*, pages 419–430, 2003.
10. H. Hosoya and B. C. Pierce. Regular expression pattern matching for XML. *Journal of Functional Programming*, 13(6):961–1004, November 2003.
11. H. Hosoya and B. C. Pierce. XDuce: A statically typed XML processing language. *ACM Transactions on Internet Technology*, 3(2):117–148, 2003.
12. S. P. Jones. Space usage. Glasgow Haskell Users Mailing List, <http://www.haskell.org/pipermail/glasgow-haskell-users/2004-August/007023.html>, 17th August 2004.
13. M. Kay. SAXON: The XSLT and XQuery processor. <http://saxon.sourceforge.net/>.
14. O. Kiselyov. A better XML parser through functional programming. In *4th International Symposium on Practical Aspects of Declarative Languages*, volume 2257 of *Lecture Notes in Computer Science*, pages 209–224, 2002.
15. K. Kodama, K. Suenaga, N. Kobayashi, and A. Yonezawa. Translation of tree-processing programs into stream-processing programs based on ordered linear type. In *The 2nd ASIAN Symposium on Programming Languages and Systems*, volume 3302 of *Lecture Notes in Computer Science*, pages 41–56, 2004.
16. B. Ludäscher, P. Mukhopadhyay, and Y. Papakonstantinou. A transducer-based XML query processor. In *Proceedings of 28th International Conference on Very Large Data Bases*, pages 227–238, 2002.
17. S. Maneth, A. Berlea, T. Perst, and H. Seidl. XML type checking with macro tree transducers. In *Proceedings of 24th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, pages 283–294, 2005.
18. M. Murata. Extended path expressions of XML. In *Proceedings of the 20th ACM Symp. on Principles of Database Systems*, pages 153–166, 2001.
19. K. Nakano. XTISP: XML transformation language intended for stream processing. <http://xtisp.org/>.
20. K. Nakano. Composing stack-attributed transducers. Technical Report METR-2004-01, Department of Mathematical Informatics, University of Tokyo, 2004.
21. K. Nakano. An implementation scheme for XML transformation languages through derivation of stream processors. In *The 2nd ASIAN Symposium on Programming Languages and Systems*, volume 3302 of *Lecture Notes in Computer Science*, pages 74–90, 2004.
22. The nhc98 compiler. <http://www.haskell.org/nhc98/>.
23. S. Nishimura. Fusion with stacks and accumulating parameters. In *the 2004 ACM SIGPLAN Workshop on Partial Evaluation and Semantics-based Program Manipulation*, pages 101–112, 2004.
24. S. Nishimura and K. Nakano. XML stream transformer generation through program composition and dependency analysis. *Science of Computer Programming*, 54:257–290, 2005.
25. Objective Caml. <http://caml.inria.fr/ocaml/>.
26. T. Perst and H. Seidl. Macro forest transducers. *Information Processing Letters*, 89:141–149, 2004.
27. S. Scherzinger and A. Kemper. Syntax-directed transformations of XML streams. In *The workshop on Programming Language Technologies for XML*, pages 75–86, 2005.
28. K. Suenaga, N. Kobayashi, and A. Yonezawa. Extension of type-based approach to generation of stream processing programs by automatic insertion of buffering primitives. In *International workshop on Logic-based Program Synthesis and Transformation*, 2005. To appear.
29. P. Wadler. Fixing a space leak with a garbage collector. *Software Practice and Experience*, 17(9):595–608, September 1987.
30. P. Wadler. Deforestation: Transforming programs to eliminate trees. In *Proceedings of the European Symposium on Programming*, volume 300 of *Lecture Notes in Computer Science*, pages 344–358, 1988.
31. libxslt: the XSLT C library for Gnome. <http://xmlsoft.org/XSLT/>.
32. SAX: the simple API for XML. <http://www.saxproject.org/>.
33. XMark: an XML benchmark project. <http://www.xml-benchmark.org/>.
34. XSL transformations (XSLT). <http://www.w3c.org/TR/xslt/>.

## A Proof of Theorem 1

We introduce another definition for each of  $\lfloor \_ \rfloor$  and  $\theta_S$  to prove THEOREM 1.

An *event generator*  $gen$  is another function of streaming  $\lfloor \_ \rfloor$ . Let  $f$  be a  $\Sigma$ -forest with an alphabet  $\Sigma$ . The function  $gen$  can output one by one the next XML event from the XML stream corresponding to  $f$  by updating a *following forest list* (for short, FFL). An FFL ranges over a set given by a syntax:

$$L ::= (F, \#) \mid (F, \sigma) :: L$$

where  $F$  ranges over a set of sub-forest of  $f$ ,  $\sigma \in \Sigma$  and  $\#$  represents an additional root symbol for  $f$ . The  $i$ -th forest in an FFL  $l$  is accessed by  $l[i]$  where  $((f', \sigma) :: l)[1] = f'$ ,  $(f', \#)[1] = f'$ ,  $((f', \sigma) :: l)[i] = l[i-1]$  and  $(f', \#)[i] = \epsilon$  for  $i > 1$ .

A event generator  $gen$  takes an FFL and returns a new FFL following an XML event, which is defined by

$$\begin{aligned} gen((\sigma \langle f_1 \rangle f_2, \sigma') :: l) &= \langle \sigma \rangle gen((f_1, \sigma) :: (f_2, \sigma') :: l) \\ gen((\sigma \langle f_1 \rangle f_2, \#)) &= \langle \sigma \rangle gen((f_1, \sigma) :: (f_2, \#)) \\ gen((\epsilon, \sigma) :: l) &= \langle / \sigma \rangle gen(l) \\ gen((\epsilon, \#)) &= \$ \end{aligned}$$

The function  $gen$  always terminates because every step reduces a total number of nodes of forests in the FFL. We write  $G(l)$  for a set of FFL's occurring as an argument of  $gen$  in the computation of  $gen(l)$ , i.e., if  $l = (\epsilon, \#)$  then  $G(l) = \{(\epsilon, \#)\}$  and otherwise  $G(l) = \{l\} \cup G(l')$  with  $gen(l) = \chi gen(l')$  with  $\chi \in \Sigma_{\diamond}$ . The set  $G((f, \#))$  for a  $f \in \mathcal{F}_{\Sigma}$  is finite because of the termination of  $gen$ . The following lemma shows that  $gen$  simulates a streaming  $\lfloor f \rfloor$  for a forest  $f$  when the initial FFL is  $(f, \#)$ .

**Lemma 1.** *Let  $f$  be a  $\Sigma$ -forest with an alphabet  $\Sigma$ . Then we have*

$$gen((f, \#)) = \lfloor f \rfloor \$ \tag{3}$$

*Proof.* We use an equation

$$gen((f', \sigma) :: l) = \lfloor f' \rfloor \langle / \sigma \rangle gen(l) \tag{4}$$

for a sub-forest  $f'$  of  $f$ ,  $\sigma \in \Sigma$  and an FFL  $l$  to show (3). The equation (4) can be proved by induction on the structure of  $f'$ . We prove the equation (3) by induction on the structure  $f$ . If  $f = \epsilon$ , then (3) holds by the definitions of  $gen$  and  $\lfloor \_ \rfloor$ . If  $f = \sigma \langle f_1 \rangle f_2$ , then we have

$$\begin{aligned} gen((\sigma \langle f_1 \rangle f_2, \#)) &= \langle \sigma \rangle gen((f_1, \sigma) :: (f_2, \#)) \\ &= \langle \sigma \rangle \lfloor f_1 \rfloor \langle / \sigma \rangle gen((f_2, \#)) \\ &= \langle \sigma \rangle \lfloor f_1 \rfloor \langle / \sigma \rangle \lfloor f_2 \rfloor \$ \\ &= \lfloor \sigma \langle f_1 \rangle f_2 \rfloor \$ \end{aligned}$$

from the definitions of  $gen$  and  $\lfloor \_ \rfloor$ , the equation (4) and the induction hypothesis. Therefore (3) holds for every forest  $f$ .  $\square$

Now we introduce another definition for  $\theta_S$  in DEFINITION 7. Let  $S = (Q, \Sigma, \Delta, in, R)$  be an xsp and  $f$  be a  $\Sigma$ -forest. We define a function  $\Theta_S$  as well as  $\theta_S$  by

$$\Theta_S(e, (\epsilon, \#)) = \langle e, \$ \rangle \qquad \Theta_S(e, l) = \Theta_S(\langle e, \chi \rangle, l')$$

for  $e \in Tmp_S$  and an FFL  $l$  where  $gen(l) = \chi gen(l')$  with  $\chi \in \Sigma_{\diamond}$ . The function  $\Theta_S$  takes an FFL as the second argument while  $\theta_S$  takes an event  $\chi$ . The following lemma shows that the function  $\Theta_S$  can simulate a transformation  $\tau_S$  for an input stream  $\lfloor f \rfloor$ .

**Lemma 2.** Let  $S = (Q, \Sigma, \Delta, in, R)$  be an xsp. Then

$$\tau_S(\lfloor f \rfloor) = \Theta_S(in(\varepsilon, \dots, \varepsilon), (f, \#)) \quad (5)$$

for every  $f \in \mathcal{F}_\Sigma$ .

*Proof.* From the definition of  $\tau_S$ , what we have to show is

$$\theta_S(in(\varepsilon, \dots, \varepsilon), \lfloor f \rfloor \mathbb{S}) = \Theta_S(in(\varepsilon, \dots, \varepsilon), (f, \#)). \quad (6)$$

We show an equation

$$\theta_S(e, gen(l)) = \Theta_S(e, l) \quad (7)$$

for  $e \in \text{Tmp}_S$  and  $l \in G((f, \#))$ , which is more general than (6) owing to LEMMA 1. The equation (7) is proved by induction on the cardinality  $\sharp G(l)$  of  $G(l)$ . If  $\sharp G(l) = 1$ , i.e.,  $l = (\varepsilon, \#)$ , then both sides are the same, that is  $\langle e, \mathbb{S} \rangle$ . If  $\sharp G(l) > 1$ , then  $gen(l) = \chi gen(l')$  with  $\chi \in \Sigma_{\langle \rangle}$ . We have  $\sharp G(l') = \sharp G(l) - 1$  since  $l \notin G(l')$  from the definition of  $G$ . Therefore we obtain

$$\begin{aligned} \theta_S(e, gen(l)) &= \theta_S(e, \chi gen(l')) \\ &= \theta_S(\langle e, \chi \rangle, gen(l')) \\ &= \Theta_S(\langle e, \chi \rangle, l') \\ &= \Theta_S(e, l) \end{aligned}$$

from the induction hypothesis, the definitions of  $\theta_S$  and  $\Theta_S$ . Hence the equation (7) holds for  $e \in \text{Tmp}_S$  and  $l \in G((f, \#))$ .  $\square$

Next we define the function  $\mathcal{I} : \text{Tmp}_{\mathcal{SP}(M)} \rightarrow \mathcal{F}_\Delta$  for an mft  $M = (Q, \Sigma, \Delta, in, R)$ . The function  $\mathcal{I}$  translates a temporary expression into a corresponding output forest and has an *invariant property* that  $\mathcal{I}(e, l) = \mathcal{I}(e', l')$  if  $\Theta_S(e, l)$  is computed by  $\Theta_S(e', l')$ , which will be shown as LEMMA 4. Since all temporary expressions for  $\mathcal{SP}(M)$  are *well-defined* from the definition of  $\mathcal{SP}(M)$  in the sense that they range over

$$E ::= q(E, \dots, E) \mid \varepsilon \mid \langle \delta \rangle E \langle /\delta \rangle \mid E E,$$

we define the function  $\mathcal{I}$  by

$$\begin{aligned} \mathcal{I}(q[i](e_1, \dots, e_n), l) &= \llbracket q \rrbracket(l[i], \mathcal{I}(e_1, l), \dots, \mathcal{I}(e_n, l)), \\ \mathcal{I}(\varepsilon, l) &= \varepsilon, & \mathcal{I}(\langle \delta \rangle e \langle /\delta \rangle, l) &= \delta(\mathcal{I}(e, l)), \\ \mathcal{I}(e e', l) &= \mathcal{I}(e, l) \mathcal{I}(e', l). \end{aligned}$$

To prove the invariant property of  $\mathcal{I}$ , we show the following lemma which correlates the right-hand sides of rules of  $M$  and  $\mathcal{SP}(M)$ .

**Lemma 3.** Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft,  $rhs$  be a right-hand side of a rule in  $R_q$  with  $q \in Q$  and  $\text{rank}(q) = n + 1$ ,  $l$  be an FFL for  $f \in \mathcal{F}_\Sigma$  and  $e_j$  with  $j = 1, \dots, n$  be a (well-formed) temporary expression for  $\mathcal{SP}(M)$ . For a binding  $\rho$  given by  $\rho(x_i) = l[i]$  ( $i = 1, 2$ ) and  $\rho(y_j) = \mathcal{I}(e_j, l)$  ( $j = 1, \dots, n$ ), we have

$$\llbracket rhs \rrbracket_\rho = \mathcal{I}(\mathcal{A}(rhs)[y_j := e_j]_{j=1, \dots, n}, l) \quad (8)$$

where  $\mathcal{A}$  is as given in DEFINITION 8. In particular, when  $rhs$  does not contain any pattern  $q(x_i, \dots)$ , we have

$$\llbracket \llbracket rhs \rrbracket_\rho \rrbracket = \mathcal{A}(rhs)[y_j := \lfloor \mathcal{I}(e_j, l) \rfloor]_{j=1, \dots, n}. \quad (9)$$

*Proof.* The statement is proved by induction on the structure of  $rhs$ . In the rest of this proof, we write just  $\rho'(\xi)$  for  $\xi[y_j := e_j]_{j=1,\dots,n}$ . Thus we should prove the equation  $\llbracket rhs \rrbracket_\rho = \mathcal{I}(\rho'(\mathcal{A}(rhs)), l)$ . If  $rhs = \epsilon$ , then both sides of (8) are the same, that is  $\epsilon$ . If  $rhs = \delta\langle rhs' \rangle$ , then the right-hand side of (8) is  $\llbracket \delta\langle rhs' \rangle \rrbracket_\rho = \delta\langle \llbracket rhs' \rrbracket_\rho \rangle$ . From the definitions of  $\mathcal{A}$  and  $\mathcal{I}$ , the left-hand side of (8) is

$$\begin{aligned} & \mathcal{I}(\rho'(\mathcal{A}(\delta\langle rhs' \rangle)), l) \\ &= \mathcal{I}(\langle \delta \rangle \rho'(\mathcal{A}(rhs')) \langle /\delta \rangle, l) \\ &= \delta\langle \mathcal{I}(\rho'(\mathcal{A}(rhs')), l) \rangle \end{aligned}$$

which is equal to  $\delta\langle \llbracket rhs' \rrbracket_\rho \rangle$  from the induction hypothesis. If  $rhs = y_j$ , then both sides of (8) are the same, that is  $\mathcal{I}(e_j, l)$ . If  $rhs = rhs_1 rhs_2$ , then both sides of (8) are the same from the definitions of  $\mathcal{A}$  and  $\mathcal{I}$  and the induction hypothesis. If  $rhs = q(x_i, rhs_1, \dots, rhs_{n'})$ , then the right-hand side of (8) is

$$\begin{aligned} & \llbracket q(x_i, rhs_1, \dots, rhs_{n'}) \rrbracket_\rho \\ &= \llbracket q \rrbracket(l[i], \llbracket rhs_1 \rrbracket_\rho, \dots, \llbracket rhs_{n'} \rrbracket_\rho) \\ &= \llbracket q \rrbracket(l[i], \mathcal{I}(\rho'(\mathcal{A}(rhs_1)), l), \dots, \mathcal{I}(\rho'(\mathcal{A}(rhs_{n'})), l)) \end{aligned}$$

from the definition of  $\llbracket \_ \rrbracket_\rho$  and the induction hypothesis. The left-hand side of (8) is

$$\begin{aligned} & \mathcal{I}(\rho'(\mathcal{A}(q(x_i, rhs_1, \dots, rhs_{n'}))), l) \\ &= \mathcal{I}(\rho'(q[i](\mathcal{A}(rhs_1), \dots, \mathcal{A}(rhs_{n'}))), l) \\ &= \mathcal{I}(q[i](\rho'(\mathcal{A}(rhs_1)), \dots, \rho'(\mathcal{A}(rhs_{n'}))), l) \\ &= \llbracket q \rrbracket(l[i], \mathcal{I}(\rho'(\mathcal{A}(rhs_1)), l), \dots, \mathcal{I}(\rho'(\mathcal{A}(rhs_{n'})), l)) \end{aligned}$$

from the definitions of  $\mathcal{A}$  and  $\mathcal{I}$ . Hence we have the equation (8).

In the particular case that  $rhs$  does not contain any pattern  $q(x_i, \dots)$ , the equation (9) is shown by the induction of the structure of  $rhs$ . When  $rhs = \epsilon$ , the both sides are the same, that is  $\epsilon$ . When  $rhs = \delta\langle rhs' \rangle$ , the left-hand side of (9) is  $\langle \delta \rangle \llbracket \llbracket rhs' \rrbracket_\rho \rangle \langle /\delta \rangle$  and the right-hand side is  $\langle \delta \rangle \rho'(\mathcal{A}(rhs')) \langle /\delta \rangle$ . These are equal from the induction hypothesis. When  $rhs = y_j$ , the both sides are the same, that is  $\llbracket e_j \rrbracket$ . When  $rhs = rhs_1 rhs_2$ , the left-hand side of (9) is  $\llbracket \llbracket rhs_1 \rrbracket_\rho \rrbracket \llbracket \llbracket rhs_2 \rrbracket_\rho \rrbracket$  and the right-hand side is  $\rho'(\mathcal{A}(rhs_1))\rho'(\mathcal{A}(rhs_2))$ . These are equal from the induction hypothesis. Therefore (9) holds.  $\square$

Then we show the following lemma with respect to the invariant property of  $\mathcal{I}$ .

**Lemma 4.** *Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft,  $e$  be a (well-defined) temporary expression for an xsp  $\mathcal{SP}(M)$  and  $l$  be an FFL for a forest  $f \in \mathcal{F}_\Sigma$ . For  $gen(l) = \chi gen(l')$ , we have*

$$\mathcal{I}(e, l) = \mathcal{I}(\langle e, \chi \rangle, l'). \quad (10)$$

*Proof.* We prove the statement by induction on the cardinality  $\sharp G(l)$  of  $G(l)$ . We start with the case of  $\sharp G(l) = 2$  since the equation  $gen(l) = \chi gen(l')$  is assumed.

In the case of  $\sharp G(l) = 2$ , we have  $l = (\epsilon, \sigma) :: (\epsilon, \#)$  with  $gen(l) = \langle / \sigma \rangle gen((\epsilon, \#))$ . We have to show an equation

$$\mathcal{I}(e, (\epsilon, \sigma) :: (\epsilon, \#)) = \mathcal{I}(\langle e, \langle / \sigma \rangle \rangle, (\epsilon, \#)). \quad (11)$$

It is proved by induction on the structure of the (well-defined) temporary expression  $e$ . We only show the case of  $e = q[i](e_1, \dots, e_n)$  with  $q \in Q$  and  $i \in \mathbb{N}^+$  that is the most complicated one in the induction. The other cases can be shown from the definitions of  $\mathcal{I}$  and  $\langle \_ \rangle, \llbracket \_ \rrbracket$ .

If  $e = q[1](e_1, \dots, e_n)$  with  $q \in Q$ , then we have

$$\begin{aligned} & \mathcal{I}(q[1](e_1, \dots, e_n), l) \\ &= \llbracket q \rrbracket(l[1], \mathcal{I}(e_1, l), \dots, \mathcal{I}(e_n, l)) \\ &= \llbracket q \rrbracket(\epsilon, \mathcal{I}(\langle e_1, \langle / \sigma \rangle \rangle, l'), \dots, \mathcal{I}(\langle e_n, \langle / \sigma \rangle \rangle, l')) \\ &= \llbracket rhs^{q, \epsilon} \rrbracket_\rho \\ &= \mathcal{I}(\mathcal{A}(rhs^{q, \epsilon})[y_j := \langle e_j, \langle / \sigma \rangle \rangle]_{j=1,\dots,n}, l') \\ &= \mathcal{I}(\langle q[1](e_1, \dots, e_n), \langle / \sigma \rangle \rangle, l') \end{aligned}$$

from the definitions of  $\mathcal{I}$ ,  $\llbracket q \rrbracket$  and  $\langle \_ , \_ \rangle$ ,  $l = (\epsilon, \sigma) :: (\epsilon, \#)$ , the induction hypothesis and LEMMA 3 where  $\rho$  is a binding defined by  $\rho(y_j) = \mathcal{I}(\langle e_j, \langle \sigma \rangle \rangle, l')$  for  $j = 1, \dots, n$ . Hence the equation (11) holds.

If  $e = q[i](e_1, \dots, e_n)$  with  $q \in Q$  and  $i > 1$ , then we have

$$\begin{aligned} & \mathcal{I}(q[i](e_1, \dots, e_n), l) \\ &= \llbracket q \rrbracket(\epsilon, \mathcal{I}(\langle e_1, \langle \sigma \rangle \rangle, l'), \dots, \mathcal{I}(\langle e_n, \langle \sigma \rangle \rangle, l')) \end{aligned}$$

in a way similar to the case of  $e = q[1](e_1, \dots, e_n)$ . The right-hand side of (11) is

$$\begin{aligned} & \mathcal{I}(\langle q[i](e_1, \dots, e_n), \langle \sigma \rangle \rangle, l') \\ &= \mathcal{I}(q[i-1](\langle e_1, \langle \sigma \rangle \rangle, \dots, \langle e_n, \langle \sigma \rangle \rangle), l') \\ &= \llbracket q \rrbracket(l'[i-1], \mathcal{I}(\langle e_1, \langle \sigma \rangle \rangle, l'), \dots, \mathcal{I}(\langle e_n, \langle \sigma \rangle \rangle, l')) \end{aligned}$$

from the definitions  $\langle \_ , \_ \rangle$  and  $\mathcal{I}$ , a  $(q, \langle \sigma \rangle)$ -rule of  $\mathcal{SP}(M)$ . The equation (11) holds from  $l'[i-1] = \epsilon$ .

Next we show the case of  $\#G(l) > 2$  in (10). Assume that  $gen(l) = \chi gen(l')$ . The equation (10) is proved by induction on the structure of the (well-defined) temporary expression  $e$ . Again we only show the case of  $e = q[i](e_1, \dots, e_n)$  with  $q \in Q$  and  $i \in \mathbb{N}^+$  since the other cases are easily shown from the definitions of  $\mathcal{I}$  and  $\langle \_ , \_ \rangle$ .

If  $e = q[1](e_1, \dots, e_n)$  with  $q \in Q$ , then the left-hand side of (10) is equal to  $\llbracket q \rrbracket(l[1], \mathcal{I}(e_1, l), \dots, \mathcal{I}(e_n, l))$  by the definition of  $\mathcal{I}$ . There are three possibilities for the form of  $l$ ,  $(\sigma \langle f_1 \rangle f_2, \sigma') :: l''$ ,  $(\sigma \langle f_1 \rangle f_2, \#)$  and  $(\epsilon, \sigma) :: l''$ . In the first two cases, we have  $l[1] = \sigma \langle f_1 \rangle f_2$  and  $gen(l) = \langle \sigma \rangle gen(l')$  with  $l'[i] = f_i$  ( $i = 1, 2$ ). The left-hand side of (10) is

$$\begin{aligned} & \llbracket q \rrbracket(l[1], \mathcal{I}(e_1, l), \dots, \mathcal{I}(e_n, l)) \\ &= \llbracket q \rrbracket(\sigma \langle f_1 \rangle f_2, \mathcal{I}(\langle e_1, \langle \sigma \rangle \rangle, l'), \dots, \mathcal{I}(\langle e_n, \langle \sigma \rangle \rangle, l')) \\ &= \llbracket rhs^{q, \sigma} \rrbracket_\rho \\ &= \mathcal{I}(\mathcal{A}(rhs^{q, \sigma})[y_j := \langle e_j, \langle \sigma \rangle \rangle]_{j=1, \dots, n}, l') \\ &= \mathcal{I}(\langle q[1](e_1, \dots, e_n), \langle \sigma \rangle \rangle, l') \end{aligned}$$

from the definitions of  $\llbracket \_ \rrbracket$  and  $\langle \_ , \_ \rangle$ , the induction hypothesis and LEMMA 3 where  $\rho$  is a binding defined by  $\rho(x_i) = f_i$  for  $i = 1, 2$  and  $\rho(y_j) = \mathcal{I}(\langle e_j, \langle \sigma \rangle \rangle, l')$  for  $j = 1, \dots, n$ . Hence we have the equation (10) in this case of  $l$ . In the case of  $l = (\epsilon, \sigma) :: l''$ , the equation (10) can be shown in a way similar to the case of  $G(l) = 2$  and  $e = q[1](e_1, \dots, e_n)$  because  $l[1] = \epsilon$ .

If  $e = q[i](e_1, \dots, e_n)$  with  $q \in Q$  and  $i > 1$ , then the left-hand side of (10) is equal to  $\llbracket q \rrbracket(l[i], \mathcal{I}(e_1, l), \dots, \mathcal{I}(e_n, l))$  by the definition of  $\mathcal{I}$ . There are three possibilities for the form of  $l$ ,  $(\sigma \langle f_1 \rangle f_2, \sigma') :: l''$ ,  $(\sigma \langle f_1 \rangle f_2, \#)$  and  $(\epsilon, \sigma) :: l''$ . In the first two cases, we have  $gen(l) = \langle \sigma \rangle gen(l')$  with  $l[i] = l'[i+1]$  and  $G(l) = G(l') + 1$ . Then we obtain

$$\begin{aligned} & \llbracket q \rrbracket(l[i], \mathcal{I}(e_1, l), \dots, \mathcal{I}(e_n, l)) \\ &= \llbracket q \rrbracket(l'[i+1], \mathcal{I}(\langle e_1, \langle \sigma \rangle \rangle, l'), \dots, \mathcal{I}(\langle e_n, \langle \sigma \rangle \rangle, l')) \\ &= \mathcal{I}(q[i+1](\langle e_1, \langle \sigma \rangle \rangle, \dots, \langle e_n, \langle \sigma \rangle \rangle), l') \\ &= \mathcal{I}(\langle q[i](e_1, \dots, e_n), \langle \sigma \rangle \rangle, l') \end{aligned}$$

from the induction hypotheses, the definition of  $\langle \_ , \_ \rangle$  and a rule in an xsp  $\mathcal{SP}(M)$ . Hence the equation (10) holds in this case of  $l$ . In the case of  $l = (\epsilon, \sigma) :: l''$ , we have  $gen(l) = \langle \sigma \rangle gen(l')$  with  $l[i] = l'[i-1]$  and  $G(l) = G(l') - 1$ . Then the equation (10) is shown in a way similar to the previous case.  $\square$

The following lemma shows the relation of  $\Theta_S$  and  $\mathcal{I}$ .

**Lemma 5.** *Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft and  $l \in G((f, \#))$  be an FFL for  $f \in \mathcal{F}_\Sigma$ . For any temporary expression  $e \in \mathcal{SPM}$ , we have*

$$\llbracket \mathcal{I}(e, l) \rrbracket = \Theta_{\mathcal{SP}(M)}(e, l). \quad (12)$$

*Proof.* For the computation chain of  $\Theta_S(e_0, (f, \#)) = \dots = \Theta_S(e, (\epsilon, \#))$ , we have  $\mathcal{I}(e_0, (f, \#)) = \dots = \mathcal{I}(e, (\epsilon, \#))$  due to LEMMA 4. Hence it is enough to show

$$\lfloor \mathcal{I}(e, (\epsilon, \#)) \rfloor = \Theta_S(e, (\epsilon, \#)) (= \langle e, \$ \rangle) \quad (13)$$

It is proved by induction on  $e$ . We only show the case of  $e = q[i](e_1, \dots, e_n)$  with  $q \in Q$  and  $i \in \mathbb{N}^+$  because the other cases can be shown from the definitions of  $\Theta_{\mathcal{SP}(M)}$ ,  $\lfloor \_ \rfloor$  and  $\mathcal{I}$ . Since  $((\epsilon, \#))[i] = \epsilon$  for any  $i \in \mathbb{N}^+$ , the left-hand side of (13) is

$$\begin{aligned} & \lfloor \mathcal{I}(q[i](e_1, \dots, e_n), (\epsilon, \#)) \rfloor \\ &= \lfloor \lfloor q \rfloor(\epsilon, \mathcal{I}(e_1, (\epsilon, \#)), \dots, \mathcal{I}(e_n, (\epsilon, \#))) \rfloor \\ &= \lfloor \lfloor rhs^{q, \epsilon} \rfloor_\rho \rfloor \end{aligned}$$

with  $\rho(y_j) = \mathcal{I}(e_j, (\epsilon, \#))$  ( $j = 1, \dots, n$ ) from the definitions of  $\mathcal{I}$  and  $\lfloor \_ \rfloor$ . The right-hand side of (13) is

$$\begin{aligned} & \langle q[i](e_1, \dots, e_n), \$ \rangle \\ &= \mathcal{A}(rhs^{q, \epsilon})[y_j := \langle e_j, \$ \rangle]_{j=1, \dots, n} \\ &= \mathcal{A}(rhs^{q, \epsilon})[y_j := \lfloor \mathcal{I}(e_j, (\epsilon, \#)) \rfloor]_{j=1, \dots, n} \end{aligned}$$

from the definitions of  $\langle \_ \rangle$  and the induction hypothesis. Then these are equal from (9) in LEMMA 3 since no pattern  $q'(x_i, \dots)$  occurs in  $rhs^{q, \epsilon}$ .  $\square$

Now we prove THEOREM 1. Let  $M = (Q, \Sigma, \Delta, in, R)$  be an mft and  $S = \mathcal{SP}(M)$  be an xsp. For a forest  $f \in \mathcal{F}_\Sigma$ , we have

$$\begin{aligned} \tau_S(\lfloor f \rfloor) &= \Theta_S(in[1](\epsilon, \dots, \epsilon), (f, \#)) \\ &= \lfloor \mathcal{I}(in[1](\epsilon, \dots, \epsilon), (f, \#)) \rfloor \\ &= \lfloor \lfloor in \rfloor(f, \epsilon, \dots, \epsilon) \rfloor \\ &= \lfloor \tau_M(f) \rfloor \end{aligned}$$

from LEMMA 2, LEMMA 5 and the definitions of  $\mathcal{I}$  and  $\tau_M$ . Therefore THEOREM 1 has been proved.