

# Verification of Tree-Processing Programs via Higher-Order Model Checking

Hiroshi Unno, Naoshi Tabuchi, and Naoki Kobayashi

Tohoku University

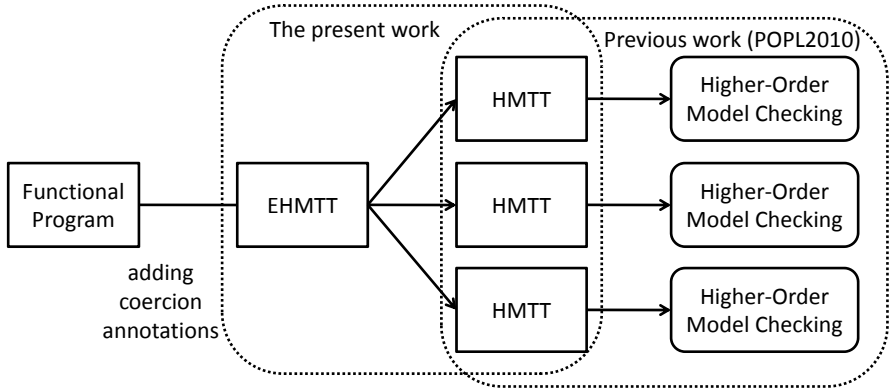
**Abstract.** We propose a new method to verify that a higher-order, tree-processing functional program conforms to an input/output specification. Our method reduces the verification problem to multiple verification problems for higher-order multi-tree transducers, which are then transformed into higher-order recursion schemes and model-checked. Unlike previous methods, our new method can deal with arbitrary higher-order functional programs manipulating algebraic data structures, as long as certain invariants on intermediate data structures are provided by a programmer. We have proved the soundness of the method and implemented a prototype verifier.

## 1 Introduction

The model checking of higher-order recursion schemes [20], or higher-order model checking for short, has been extensively studied recently. Ong [20] has shown the decidability of higher-order model checking. Kobayashi [13,12] then developed a practical model checking algorithm and applied it to program verification. The present work is an extension of that line of work, trying to apply higher-order model checking to verification of a wider range of higher-order programs.

From a programming language point of view, recursion schemes are terms of the simply-typed  $\lambda$ -calculus with recursion and tree constructors (but not destructors). One can also encode finite data domains (such as booleans) by using Church encoding. Based on this observation, Kobayashi [13] applied model checking to resource usage verification of simply-typed functional programs with recursion, booleans, and resource primitives. The limitation of this approach was that programs manipulating infinite data domains such as lists and trees could not be handled. To relax this limitation, in our previous work [15], we have introduced higher-order multi-parameter tree transducers (HMTTs) as an extension of recursion schemes with tree destructors. HMTTs are a kind of tree transducers that take (possibly infinite) input trees, which can be destructed, and outputs a (possibly infinite) tree. However, there still remains a gap between HMTTs and ordinary functional programs that use recursive data structures since HMTTs do not support intermediate data structures: an HMTT cannot destruct trees constructed by the HMTT itself.

In this paper, we propose a verification method for an extension of HMTTs called EHMTTs. In essence, EHMTTs are higher-order, simply-typed functional



**Fig. 1.** Overall Structure of EHMTT Verification Method

programs with recursion and tree primitives. Unlike our previous HMTTs [15], there is no fundamental restriction on tree constructors/destructors, except that special annotations (called coercion) are required for destructing trees constructed in a program. Our method can check whether the output trees generated by a given EHMTT conform to a given output specification whenever the input trees conform to given input specifications. We can apply our method to verification of ordinary functional programs that manipulate algebraic data structures by encoding them as trees and adding annotations to the programs.

The overall structure of our method is shown in Figure 1. A given EHMTT verification problem is reduced to multiple HMTT verification problems, which are then solved by an HMTT verification method presented in our previous work [15]. The HMTT verification method further reduces the HMTT verification problems to model checking problems of recursion schemes, which are finally solved by Kobayashi’s higher-order model checker TRECS [12]. In this paper, we have formalized the reduction from an EHMTT verification problem to HMTT verification problems, and proved the soundness of the reduction. Our verification method is not complete, however, since the verification problem is undecidable in general. We have implemented a prototype verifier and verified functional programs that manipulate XML and user-defined recursive data structures.

The rest of the paper is organized as follows. Section 2 presents some preliminary definitions and notations. In Section 3, we introduce EHMTTs. Section 4 formalizes our verification method for EHMTTs. Section 5 reports on the experimental results. We compare our method with related work in Section 6, and conclude the paper with some remarks on future work in Section 7.

## 2 Preliminaries

We write  $\text{dom}(f)$  for the domain of a map  $f$ , and  $f\{x \mapsto v\}$  for the map  $f'$  such that  $\text{dom}(f') = \text{dom}(f) \cup \{x\}$ ,  $f'(x) = v$  and  $f'(y) = f(y)$  for  $y \in \text{dom}(f) \setminus \{x\}$ .

We write  $X^*$  for the set of sequences of elements of  $X$ . We write  $\epsilon$  for the empty sequence, and  $v_1 \cdots v_n$  for the sequence consisting of  $v_1, \dots, v_n$ . We write  $s_1 \cdot s_2$  for the concatenation of sequences  $s_1$  and  $s_2$ . A sequence  $v_1 \cdots v_n$  is often abbreviated to  $\tilde{v}$ .

A *ranked alphabet*  $\Sigma$  is a map from a finite set of symbols to non-negative integers. For each symbol  $a \in \text{dom}(\Sigma)$ ,  $\Sigma(a)$  denotes the arity of  $a$ . We write  $A_\Sigma$  to denote the largest arity of the symbols in  $\text{dom}(\Sigma)$ . A  $\Sigma$ -*labeled ranked tree*  $T$  is a map from a subset of  $\{1, \dots, A_\Sigma\}^*$  to  $\text{dom}(\Sigma)$  such that:

- $\text{dom}(T)$  is prefix-closed, i.e. if  $\pi \cdot i \in \text{dom}(T)$ , then  $\pi \in \text{dom}(T)$ ; and
- if  $T(\pi) = a$ , then  $\{i \mid \pi \cdot i \in \text{dom}(T)\} = \{1, \dots, \Sigma(a)\}$ .

### 3 Extended HMTTs

In this section, we introduce extended HMTTs (EHMTTs). From a programming language point of view, an EHMTT is a simply-typed, call-by-name, higher-order functional program that takes possibly infinite trees as input and outputs a possibly infinite tree. The main differences from ordinary functional programs is that trees are classified into input and output trees. Input trees can only be destructed, and output trees can only be constructed in a program, as in other tree transducers. Special annotations (**coerce** <sup>$L$</sup> ( $\cdot$ ) introduced below) are however provided to convert output trees to input trees, so that, unlike ordinary tree transducers, trees constructed in a program can be destructed again in the same program. Thus, the class of EHMTTs is actually Turing complete.

We fix below a ranked alphabet  $\Sigma$ . We call elements of  $\text{dom}(\Sigma)$  *terminal symbols*, and use the meta-variable  $a$  for them.

**Definition 1 (EHMTT).** An EHMTT  $\mathcal{P}$  is a pair  $(D, S)$  where  $D$  is a set of function definitions of the form  $\{F_1 \tilde{x}_1 = t_1, \dots, F_n \tilde{x}_n = t_n\}$ , and  $S$  is a function name. Here,  $t$  ranges over the set of terms, given by:

$$t ::= a \mid x \mid F \mid t_1 \ t_2 \mid \text{case } t \text{ of } \{a_i \tilde{y}_i \Rightarrow t_i\}_{i=1}^n \mid \text{coerce}^L(t) \mid \text{gen}^L$$

Here,  $L$  denotes a set of trees. An EHMTT  $(D, S)$  is well-sorted under  $\mathcal{K}$  if  $S : \mathbf{i} \rightarrow \dots \rightarrow \mathbf{i} \rightarrow \mathbf{o} \in \mathcal{K}$  and  $\vdash D : \mathcal{K}$  is derivable by using the sort assignment rules in Figure 2. An HMTT is an EHMTT that does not contain **coerce** <sup>$L$</sup> ( $t$ ).

In the figure, the sorts  $\mathbf{i}$  and  $\mathbf{o}$  describe input and output trees respectively. The sort  $\kappa_1 \rightarrow \kappa_2$  denotes functions that take a tree or tree function of sort  $\kappa_1$  and return a tree or tree function of sort  $\kappa_2$ .  $\tilde{\kappa} \rightarrow \mathbf{o}$  and  $\tilde{x} : \tilde{\kappa}$  are shorthand forms of  $\kappa_1 \rightarrow \dots \rightarrow \kappa_k \rightarrow \mathbf{o}$  and  $x_1 : \kappa_1, \dots, x_k : \kappa_k$  respectively. We consider only well-sorted EHMTTs below.

The term **case**  $t$  of  $\{a_i \tilde{y}_i \Rightarrow t_i\}_{i=1}^n$  is reduced to  $[\tilde{u}_i/\tilde{y}_i]t_i$  if  $t$  evaluates to  $a_i \tilde{u}_i$ . If  $t$  does not match any pattern, the term evaluates to a special terminal symbol **fail**.<sup>1</sup> The term **coerce** <sup>$L$</sup> ( $t$ ) asserts that the tree generated by  $t$  belongs

<sup>1</sup> Thus, verification of the absence of pattern match errors can be encoded as a problem of checking that the tree generated by EHMTT does not contain **fail**, which is an instance of EHMTT verification problems considered below.

*Syntax of Sorts:*

$$\kappa ::= \mathbf{i} \mid \mathbf{o} \mid \kappa_1 \rightarrow \kappa_2$$

*Sort Assignment Rules:*

$$\begin{array}{c}
\mathcal{K} \vdash F : \mathcal{K}(F) \quad (\text{T-FUN}) \\
\mathcal{K} \vdash a : \underbrace{\mathbf{o} \rightarrow \dots \rightarrow \mathbf{o}}_{\Sigma(a)} \rightarrow \mathbf{o} \quad (\text{T-CON}) \\
\mathcal{K}, x : \kappa \vdash x : \kappa \quad (\text{T-VAR}) \\
\mathcal{K} \vdash t_1 : \kappa_1 \rightarrow \kappa_2 \\
\frac{\mathcal{K} \vdash t_2 : \kappa_1}{\mathcal{K} \vdash t_1 \ t_2 : \kappa_2} \quad (\text{T-APP}) \\
\mathcal{K} \vdash t : \mathbf{i} \quad \mathcal{K}, \tilde{y}_i : \tilde{\mathbf{i}} \vdash t_i : \mathbf{o} \\
\text{(for all } i = 1, \dots, N) \\
\frac{}{\mathcal{K} \vdash \mathbf{case} \ t \ \mathbf{of} \ \{a_i \ \tilde{y}_i \Rightarrow t_i\}_{i=1}^n : \mathbf{o}} \quad (\text{T-CASE}) \\
\mathcal{K} \vdash \mathbf{gen}^L : \mathbf{i} \quad (\text{T-GEN}) \\
\frac{\mathcal{K} \vdash t : \mathbf{o}}{\mathcal{K} \vdash \mathbf{coerce}^L(t) : \mathbf{i}} \quad (\text{T-COERCE}) \\
\mathcal{K} = \{F_1 : \tilde{\kappa}_1 \rightarrow \mathbf{o}, \dots, F_n : \tilde{\kappa}_n \rightarrow \mathbf{o}\} \\
\mathcal{K}, \tilde{x}_i : \tilde{\kappa}_i \vdash t_i : \mathbf{o} \text{ (for each } i) \\
\frac{}{\vdash \{F_1 \ \tilde{x}_1 = t_1, \dots, F_n \ \tilde{x}_n = t_n\} : \mathcal{K}} \quad (\text{T-DEF})
\end{array}$$

*Operational Semantics*

$$\begin{array}{c}
t \text{ (extended terms)} ::= \dots \mid \underline{a} \mid \mathbf{o2i}(t) \mid \mathbf{assert}^L(t) \\
E \text{ (evaluation contexts)} ::= [] \mid a \ t_1 \dots t_{j-1} \ E \ t_{j+1} \dots t_{\Sigma(a)} \\
\mid \mathbf{case} \ E \ \mathbf{of} \ \{a_i \ \tilde{y}_i \Rightarrow t_i\}_{i=1}^n \mid \mathbf{o2i}(E) \mid \mathbf{assert}^L(E) \\
\frac{F \ \tilde{x} = t \in D}{E[F \ \tilde{t}] \longrightarrow_{\mathcal{P}} E[\tilde{t}/\tilde{x}]t} \quad (\text{E-APP}) \\
E[\mathbf{case} \ \underline{a_i} \ \tilde{t} \ \mathbf{of} \ \{a_i \ \tilde{y}_i \Rightarrow t_i\}_{i=1}^n] \longrightarrow_{\mathcal{P}} E[\tilde{t}/\tilde{y}_i]t_i \quad (\text{E-CASE}) \\
\frac{a \notin \{a_1, \dots, a_n\}}{E[\mathbf{case} \ \underline{a} \ \tilde{t} \ \mathbf{of} \ \{a_i \ \tilde{y}_i \Rightarrow t_i\}_{i=1}^n] \longrightarrow_{\mathcal{P}} E[\mathbf{fail}]} \quad (\text{E-CASE-FAIL}) \\
\frac{a \ L_1 \dots L_n \subseteq L}{E[\mathbf{gen}^L] \longrightarrow_{\mathcal{P}} E[\underline{a} \ \mathbf{gen}^{L_1} \dots \mathbf{gen}^{L_n}]} \quad (\text{E-GEN}) \\
E[\mathbf{coerce}^L(t)] \longrightarrow_{\mathcal{P}} \mathbf{assert}^L(t) \quad (\text{E-COERCE-ASSERT}) \\
E[\mathbf{coerce}^L(t)] \longrightarrow_{\mathcal{P}} E[\mathbf{o2i}(t)] \quad (\text{E-COERCE-INPUT}) \\
E[\mathbf{o2i}(a_i \ \tilde{t})] \longrightarrow_{\mathcal{P}} E[\underline{a_i} \ \mathbf{o2i}(\tilde{t})] \quad (\text{E-INPUT}) \\
\frac{t^\perp \notin L^\perp}{\mathbf{assert}^L(t) \longrightarrow_{\mathcal{P}} \mathbf{Error}} \quad (\text{E-ASSERT-ERROR})
\end{array}$$

**Fig. 2.** Sort Assignment Rules and Call-by-Name Operational Semantics

to a set  $L$  of trees, and converts the tree to an input tree.  $\mathbf{gen}^L$  generates an element of  $L$  non-deterministically.

*Example 1.* Consider the EHMTT  $\mathcal{P}_{rev} = (D, Reverse)$ , where  $D$  consists of:

$$\begin{aligned}
 Reverse\ x &= \mathbf{case}\ x\ \mathbf{of}\ e \Rightarrow e \\
 &\quad | \ a\ x' \Rightarrow Append\ (\mathbf{coerce}^{b^*a^*e}(Reverse\ x'))\ (a\ e) \\
 &\quad | \ b\ x' \Rightarrow Append\ (\mathbf{coerce}^{b^*e}(Reverse\ x'))\ (b\ e). \\
 Append\ x\ y &= \mathbf{case}\ x\ \mathbf{of}\ e \Rightarrow y \\
 &\quad | \ a\ x' \Rightarrow a\ (Append\ x'\ y) \\
 &\quad | \ b\ x' \Rightarrow b\ (Append\ x'\ y).
 \end{aligned}$$

$Prev$  takes a tree of the form  $a^m(b^n(e))$  as input, and outputs a tree  $b^n(a^m(e))$ . The coercions  $\mathbf{coerce}^{b^*a^*e}(\cdot)$  and  $\mathbf{coerce}^{b^*e}(\cdot)$  assert that their arguments belong to  $\{b^m(a^n(e)) \mid m, n \geq 0\}$  and  $\{b^m(e) \mid m \geq 0\}$  respectively, and convert them to input trees. Note that  $Reverse$  and  $Append$  have sorts  $\mathbf{i} \rightarrow \mathbf{o}$  and  $\mathbf{i} \rightarrow \mathbf{o} \rightarrow \mathbf{o}$  respectively, so that  $Reverse\ x'$  returns an output tree.  $\square$

Figure 2 shows the formal semantics of the language. In the semantics, the set of terms are extended as follows. An underlined symbol  $\underline{a}$  denotes an input tree constructor (which, by the restriction of EHMTT, occurs only at run-time, not in source programs).  $\mathbf{o2i}(t)$  and  $\mathbf{assert}^L(t)$  are used to define the semantics of  $\mathbf{coerce}^L(t)$ : the former converts an output tree to an input tree, and the latter asserts that the tree generated by  $t$  belong to  $L$ . In the rule E-ASSERT-ERROR,  $t^\perp$  is a finite  $(\Sigma \cup \{\perp \mapsto 0\})$ -labeled ranked tree, defined by:

$$t^\perp = \begin{cases} a\ t_1^\perp \cdots t_n^\perp & (\text{if } t = a\ t_1 \cdots t_n) \\ \perp & (\text{otherwise}) \end{cases}$$

$L^\perp$  is the set  $\{T \mid T \preceq T' \in L\}$ , where  $T \preceq T'$  means that  $T$  is obtained from  $T'$  by replacing some nodes of  $T'$  with  $\perp$ .

*Remark 1.* Note that EHMTT is call-by-name. This is because our verification method is based on the model checking of higher-order recursion schemes, whose semantics is call-by-name. To deal with call-by-value programs, it suffices to apply CPS transformation before applying our verification method. The reasons why we allow infinite trees as inputs and outputs for EHMTTs are as follows. First, we would like to verify programs that manipulate not only finite but also infinite data structures (such as streams). Secondly, we would like to model a program that contains non-deterministic branches (which is typically obtained by abstracting branching information of a user program) as an EHMTT that generates a single tree describing all the possible outputs of the program. In that case, even if a program manipulates only finite data structures, the output of the EHMTT can be an infinite tree.

The goal of our verification is to check that a given EHMTT conforms to a given specification on input and output. As EHMTTs manipulate infinite trees, we use

top-down tree automata called *trivial automata* (which are Büchi tree automata with a trivial acceptance condition) as specifications (as well as for annotations  $L$  in  $\mathbf{coerce}^L(\cdot)$  and  $\mathbf{gen}^L$ ).

**Definition 2 (trivial automaton).** A trivial automaton  $\mathcal{M}$  is a quadruple  $(\Sigma, Q, \Delta, q_0)$ , where:

- $\Sigma$  is a ranked alphabet.
- $Q$  is a finite set of states.
- $\Delta$  is a finite subset of  $Q \times \text{dom}(\Sigma) \times Q^*$  called a transition relation such that if  $(q, a, \tilde{q}) \in \Delta$ , then the length of the sequence  $\tilde{q}$  is  $\Sigma(a)$ .
- $q_0$  is a state called an initial state.

A  $\Sigma$ -labeled ranked tree  $T$  is accepted by  $\mathcal{M}$  if there is a  $Q$ -labeled tree  $R$  such that:

- $\text{dom}(T) = \text{dom}(R)$ .
- For any  $\pi \in \text{dom}(R)$ ,  $(R(\pi), T(\pi), R(\pi \cdot 1) \cdots R(\pi \cdot \Sigma(T(\pi)))) \in \Delta$ .
- $R(\epsilon) = q_0$ .

We write  $\mathcal{L}(\mathcal{M})$  for the set of  $\Sigma$ -labeled ranked trees accepted by  $\mathcal{M}$ .

When restricted to finite trees, the class of languages recognized by trivial automata is equivalent to the class of regular tree languages.

*Example 2.* Recall Example 1. A trivial automaton for accepting  $\mathbf{b}^* \mathbf{a}^* \mathbf{e}$  is defined by  $(\Sigma, \{q_0, q_1\}, \Delta, q_0)$ , where:

$$\Sigma = \{\mathbf{a} \mapsto 1, \mathbf{b} \mapsto 1, \mathbf{e} \mapsto 0\}$$

$$\Delta = \{(q_0, \mathbf{b}, q_0), (q_0, \mathbf{a}, q_1), (q_0, \mathbf{e}, \epsilon), (q_1, \mathbf{a}, q_1), (q_1, \mathbf{e}, \epsilon)\}$$

□

We now formalize our verification problem:

**Definition 3.** Given an EHMTT  $\mathcal{P} = (D, S)$  and trivial automata  $\mathcal{M}_1, \dots, \mathcal{M}_k$ ,  $\mathcal{M} = (\Sigma, Q, \Delta, q_0)$ , we write  $\models (\mathcal{P}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{M})$  if for all  $T_1 \in \mathcal{L}(\mathcal{M}_1), \dots, T_k \in \mathcal{L}(\mathcal{M}_k)$ ,

1.  $S \ T_1 \cdots T_k \xrightarrow{*}_{\mathcal{P}} t$  implies  $t^\perp \in \mathcal{L}(\mathcal{M}^\perp)$ , and
2.  $S \ T_1 \cdots T_k \not\xrightarrow{*}_{\mathcal{P}} \mathbf{Error}$ .

Here,  $\mathcal{M}^\perp$  is the trivial automaton  $(\Sigma \cup \{\perp \mapsto 0\}, Q, \Delta \cup \{(q, \perp, \epsilon) \mid q \in Q\}, q_0)$ . An EHMTT verification problem  $(\mathcal{P}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{M})$  is the problem to check that  $\models (\mathcal{P}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{M})$ .

The first condition of an EHMTT verification problem says that given input trees that conform to the input specification, the EHMTT generates a valid tree.<sup>2</sup> The second condition means that the EHMTT in fact never causes a coercion error.

<sup>2</sup> Because of the presence of  $\perp$ , only safety properties are guaranteed; there is no guarantee that the EHMTT eventually generates a tree that belongs to  $\mathcal{L}(\mathcal{M})$ .

In [15], we have presented a (sound but incomplete) method for the restricted case (which we call *HMTT verification problems*) where  $\mathcal{P}$  is an HMTT (i.e., for the case where  $\mathcal{P}$  does not contain  $\mathbf{coerce}^L(\cdot)$ ). In the next section, we reduce an EHMTT verification problem to HMTT verification problems.

## 4 Verification Method for EHMTTs

We now present a method for reducing an EHMTT verification problem to HMTT verification problems (which can then be solved by the previous method [15]). The idea is to reduce each of the two conditions in Definition 3 to HMTT verification problems.

Let  $(\mathcal{P}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{M})$  be a given EHMTT verification problem, and suppose that  $\mathcal{P}$  contains  $m$  occurrences of coercions:  $\mathbf{coerce}^{L_1}(\cdot), \dots, \mathbf{coerce}^{L_m}(\cdot)$ . We construct HMTTs  $\mathcal{P}^{\mathcal{A}}, \mathcal{P}^{\mathcal{B}_1}, \dots, \mathcal{P}^{\mathcal{B}_m}$  such that:

- $\mathcal{P}^{\mathcal{A}}$  approximates the output of  $\mathcal{P}$ , by assuming that coercions never fail.
- $\mathcal{P}^{\mathcal{B}_i}$  approximates all the possible arguments of  $\mathbf{coerce}^{L_i}(\cdot)$ .

Then, the verification problem  $(\mathcal{P}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{M})$  can be reduced to  $m + 1$  HMTT verification problems:  $(\mathcal{P}^{\mathcal{A}}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{M})$ ,  $(\mathcal{P}^{\mathcal{B}_1}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{B}(L_1))$ ,  $\dots$ ,  $(\mathcal{P}^{\mathcal{B}_m}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{B}(L_m))$  (where  $\mathcal{B}(L)$  is an automaton for accepting trees representing subsets of  $L$ ; see Section 4.2 below). Sections 4.1 and 4.2 below show the constructions of  $\mathcal{P}^{\mathcal{A}}$  and  $\mathcal{P}^{\mathcal{B}_i}$  respectively.

### 4.1 Construction of $\mathcal{P}^{\mathcal{A}}$

Let  $\mathcal{P}^{\mathcal{A}}$  be the HMTT obtained by just replacing every occurrence of  $\mathbf{coerce}^{L_i}(\cdot)$  in  $\mathcal{P}$  with  $\mathbf{gen}^{L_i}$ . Then,  $\mathcal{P}^{\mathcal{A}}$  approximates the output of  $\mathcal{P}$ , assuming that no coercion error occurs.

**Theorem 1.** *Let  $\mathcal{P} = (D, S)$  be an EHMTT such that  $\models (\mathcal{P}^{\mathcal{A}}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{M})$  holds. For any  $T_1 \in \mathcal{L}(\mathcal{M}_1), \dots, T_k \in \mathcal{L}(\mathcal{M}_k)$ , if  $S \ T_1 \cdots T_k \not\vdash_{\mathcal{P}}^* \mathbf{Error}$  and  $S \ T_1 \cdots T_k \longrightarrow_{\mathcal{P}}^* t$ , then  $t^\perp \in \mathcal{L}(\mathcal{M}^\perp)$ .*

A proof is given in the full version of this paper [23].

*Example 3.* Recall  $\mathcal{P}_{rev}$  in Example 1.  $\mathcal{P}_{rev}^{\mathcal{A}}$  is  $(D, \text{Reverse}^{\mathcal{A}})$  where  $D$  is given by:

$$\begin{aligned}
 \text{Reverse}^{\mathcal{A}} \ x^{\mathcal{A}} &= \mathbf{case} \ x^{\mathcal{A}} \ \mathbf{of} \ e \Rightarrow e \\
 &\quad | \ a \ x'^{\mathcal{A}} \Rightarrow \text{Append}^{\mathcal{A}} \ \mathbf{gen}^{b^* a^* e} \ (a \ e) \\
 &\quad | \ b \ x'^{\mathcal{A}} \Rightarrow \text{Append}^{\mathcal{A}} \ \mathbf{gen}^{b^* e} \ (b \ e) \\
 \text{Append}^{\mathcal{A}} \ x^{\mathcal{A}} \ y^{\mathcal{A}} &= \mathbf{case} \ x^{\mathcal{A}} \ \mathbf{of} \ e \Rightarrow y^{\mathcal{A}} \\
 &\quad | \ a \ x'^{\mathcal{A}} \Rightarrow a \ (\text{Append}^{\mathcal{A}} \ x'^{\mathcal{A}} \ y^{\mathcal{A}}) \\
 &\quad | \ b \ x'^{\mathcal{A}} \Rightarrow b \ (\text{Append}^{\mathcal{A}} \ x'^{\mathcal{A}} \ y^{\mathcal{A}})
 \end{aligned}$$

## 4.2 Construction of $\mathcal{P}^{\mathcal{B}_i}$

The construction of  $\mathcal{P}^{\mathcal{B}_i}$  is more involved, for the following reasons.

1. Given an input,  $\mathcal{P}$  may invoke  $\mathbf{coerce}^{L_i}(\cdot)$  more than once. For example, given  $\mathbf{b}(\mathbf{b}(\mathbf{e}))$  as input,  $\mathcal{P}_{rev}$  in Example 1 invoke  $\mathbf{coerce}^{\mathbf{b}^*\mathbf{e}}(\cdot)$  twice, with different parameters  $\mathbf{e}$  and  $\mathbf{b}(\mathbf{e})$ . Thus,  $\mathcal{P}^{\mathcal{B}_i}$  should approximate the *set* of trees that are passed to  $\mathbf{coerce}^{L_i}(\cdot)$ .
2. How a function invokes  $\mathbf{coerce}^{L_i}(\cdot)$  may depend on its arguments. For example, consider a higher-order function  $F$  defined by  $F\ g\ x = g(x)$ . Obviously, how  $\mathbf{coerce}^{L_i}(\cdot)$  is invoked during evaluation of  $F\ t_1\ t_2$  depends on  $t_1$  and  $t_2$ .

To address the first issue, we represent a (possibly infinite) set of trees by a single (possibly infinite) tree. We use special terminal symbols  $\mathbf{br}$  and  $\mathbf{emp}$ , which represent the set union and an empty set respectively. For example, the set  $\{\mathbf{e}, \mathbf{b}(\mathbf{e})\}$  is represented by  $\mathbf{br}\ \mathbf{e}\ (\mathbf{b}(\mathbf{e}))$ .  $\mathcal{P}^{\mathcal{B}_i}$  outputs such a tree representation of (an over-approximation of) the set of trees passed to  $\mathbf{coerce}^{L_i}(\cdot)$ .

To address the second issue, we duplicate each parameter  $x$  of a function into  $x^{\mathcal{A}}$  and  $x^{\mathcal{B}}$ . The parameter  $x^{\mathcal{A}}$  is used to compute (an approximation of) the original value of  $x$ , while the parameter  $x^{\mathcal{B}}$  computes (an approximation of) the set of trees passed to  $\mathbf{coerce}^{L_i}(\cdot)$  during evaluation of  $x$ . For example,  $F\ g\ x = g(x)$  is transformed to:  $F^{\mathcal{B}}\ g^{\mathcal{A}}\ g^{\mathcal{B}}\ x^{\mathcal{A}}\ x^{\mathcal{B}} = g^{\mathcal{B}}\ x^{\mathcal{A}}\ x^{\mathcal{B}}$ . Here,  $F^{\mathcal{B}}$  computes an approximation of the set of trees passed to  $\mathbf{coerce}^{L_i}(\cdot)$  by calling  $g^{\mathcal{B}}$  with duplicated parameters  $x^{\mathcal{A}}$  and  $x^{\mathcal{B}}$ .

We give below more concrete examples to explain the construction of  $\mathcal{P}^{\mathcal{B}_i}$ .

*Example 4.* Recall  $\mathcal{P}_{rev}$  in Example 1. For the first coercion  $\mathbf{coerce}^{L_1}(\text{Reverse } x')$  (where  $L_1 = \mathbf{b}^*\mathbf{a}^*\mathbf{e}$ ), we construct the following HMTT  $\mathcal{P}_{rev}^{\mathcal{B}_1}$ :

$$\begin{aligned}
 \text{Reverse}^{\mathcal{B}}\ x^{\mathcal{A}}\ x^{\mathcal{B}} = & \\
 & \mathbf{br}\ x^{\mathcal{B}}\ (\mathbf{case}\ x^{\mathcal{A}}\ \mathbf{of}\ \mathbf{e} \Rightarrow \mathbf{emp} \\
 & \quad | \mathbf{a}\ x'^{\mathcal{A}} \Rightarrow \text{Append}^{\mathcal{B}}\ \mathbf{gen}^{\mathbf{b}^*\mathbf{a}^*\mathbf{e}}\ (\mathbf{br}\ (\text{Reverse}^{\mathcal{A}}\ x'^{\mathcal{A}}) \\
 & \quad \quad (\text{Reverse}^{\mathcal{B}}\ x'^{\mathcal{A}}\ \mathbf{emp})) \\
 & \quad \quad (\mathbf{a}\ \mathbf{e})\ \mathbf{emp} \\
 & \quad | \mathbf{b}\ x'^{\mathcal{A}} \Rightarrow \text{Append}^{\mathcal{B}}\ \mathbf{gen}^{\mathbf{b}^*\mathbf{e}}\ (\text{Reverse}^{\mathcal{B}}\ x'^{\mathcal{A}}\ \mathbf{emp}) \\
 & \quad \quad (\mathbf{b}\ \mathbf{e})\ \mathbf{emp}) \\
 \text{Append}^{\mathcal{B}}\ x^{\mathcal{A}}\ x^{\mathcal{B}}\ y^{\mathcal{A}}\ y^{\mathcal{B}} = & \\
 & \mathbf{br}\ x^{\mathcal{B}}\ (\mathbf{case}\ x^{\mathcal{A}}\ \mathbf{of}\ \mathbf{e} \Rightarrow y^{\mathcal{B}} \\
 & \quad | \mathbf{a}\ x'^{\mathcal{A}} \Rightarrow \text{Append}^{\mathcal{B}}\ x'^{\mathcal{A}}\ \mathbf{emp}\ y^{\mathcal{A}}\ y^{\mathcal{B}} \\
 & \quad | \mathbf{b}\ x'^{\mathcal{A}} \Rightarrow \text{Append}^{\mathcal{B}}\ x'^{\mathcal{A}}\ \mathbf{emp}\ y^{\mathcal{A}}\ y^{\mathcal{B}})
 \end{aligned}$$

As mentioned above, the parameters of *Reverse* and *Append* have been duplicated. When *Reverse*  $t$  is called in  $\mathcal{P}_{rev}$ , there are two cases where  $\mathbf{coerce}^{L_1}(\cdot)$



may be called: the case where  $t$  contains  $\mathbf{coerce}^{L_1}(\cdot)$  and it is called when  $t$  is evaluated by the case statement (note that our language is call-by-name); and the case where  $\mathbf{coerce}^{L_1}(\cdot)$  is called in a case branch. In the body of the definition of  $\mathbf{Reverse}^B$ , the part  $x^B$  approximates the set of trees passed to  $\mathbf{coerce}^{L_1}(\cdot)$  in the former case, and the part  $\mathbf{case } x^A \mathbf{ of } \dots$  approximates the set of trees for the latter case.

In the clause for  $\mathbf{a}(x'^A)$ ,  $\mathbf{Append}^B$  is used to compute an approximation of the set of trees passed to  $\mathbf{coerce}^{L_1}(\cdot)$ . The first and third parameters approximate the values of the original parameters of  $\mathbf{Append}$ . The second parameter ( $\mathbf{br } (Reverse^A x'^A) (Reverse^B x'^A \mathbf{ emp})$ ) approximates the set of trees passed to  $\mathbf{coerce}^{L_1}(\cdot)$  during the computation of  $\mathbf{coerce}^{L_1}(Reverse x')$ . Here, there are two cases where coercion can occur: (i) the value of  $Reverse x'$  is computed and passed to  $\mathbf{coerce}^{L_1}(\cdot)$  and (ii)  $\mathbf{coerce}^{L_1}(\cdot)$  is invoked during the computation of  $Reverse x'$ . The parts  $(Reverse^A x'^A)$  and  $(Reverse^B x'^A \mathbf{ emp})$  cover the former and the latter cases respectively. In the latter, the second parameter of  $(Reverse^B$  is an empty set, as the trees passed to  $\mathbf{coerce}^{L_1}(\cdot)$  during the computation of  $x'$  are already covered by  $x^B$ .  $\square$

The reduction works similarly for EHMTTs with higher-order functions.

*Example 5.* Let us consider a higher-order version of the list reverse program:

$$\begin{aligned} Reverse\ x &= Reverseh\ Append\ x \\ Reverseh\ f\ x &= \mathbf{case } x \mathbf{ of } e \Rightarrow e \\ &\quad | a\ x' \Rightarrow f\ (\mathbf{coerce}^{b^*a^*e}(Reverseh\ f\ x'))\ (a\ e) \\ &\quad | b\ x' \Rightarrow f\ (\mathbf{coerce}^{b^*e}(Reverseh\ f\ x'))\ (b\ e) \end{aligned}$$

We get the following HMTT for the first coercion  $\mathbf{coerce}^{b^*a^*e}(Reverseh\ f\ x')$ :

$$\begin{aligned} Reverse^B\ x^A\ x^B &= Reverseh^B\ Append^A\ Append^B\ x^A\ x^B \\ Reverseh^B\ f^A\ f^B\ x^A\ x^B &= \\ \mathbf{br } x^B\ (\mathbf{case } x^A \mathbf{ of } e \Rightarrow \mathbf{emp} & \\ \quad | a\ x'^A \Rightarrow f^B\ \mathbf{gen}^{b^*a^*e}\ (\mathbf{br } (Reverseh^A\ f^A\ x'^A) & \\ &\quad (Reverseh^B\ f^A\ f^B\ x'^A\ \mathbf{emp})) \\ &\quad (a\ e)\ \mathbf{emp} \\ \quad | b\ x'^A \Rightarrow f^B\ \mathbf{gen}^{b^*e}\ (Reverseh^B\ f^A\ f^B\ x'^A\ \mathbf{emp}) & \\ &\quad (b\ e)\ \mathbf{emp}) \end{aligned}$$

Here,  $\mathbf{Append}^A$  is the one obtained in Example 3. Note that  $Reverseh^B$  requires an additional argument  $f^B$ , which generates all the trees passed to the coercion by  $f$ .  $\square$

Formally, given an EHMTT  $\mathcal{P} = (D, S)$ ,  $\mathcal{P}^{\mathcal{B}_i}$  is  $(\mathcal{B}_i(D), S)$  where:

$$\begin{aligned}
\mathcal{B}_i(D) = & \{S \ x_1 \cdots x_k = S^{\mathcal{B}_i} \ x_1 \ \mathbf{emp} \cdots x_k \ \mathbf{emp}\} \cup \\
& \{a^{\mathcal{B}_i} \ x_1^{\mathcal{A}} \ x_1^{\mathcal{B}_i} \cdots x_{\Sigma(a)}^{\mathcal{A}} \ x_{\Sigma(a)}^{\mathcal{B}_i} = \mathbf{br} \ x_1^{\mathcal{B}_i} \cdots x_{\Sigma(a)}^{\mathcal{B}_i} \mid a \in \text{dom}(\Sigma)\} \cup \\
& \{F^{\mathcal{A}} \ x_1^{\mathcal{A}} \cdots x_n^{\mathcal{A}} = \mathcal{A}(t) \mid F \ x_1 \cdots x_n = t \in D\} \cup \\
& \{F^{\mathcal{B}_i} \ x_1^{\mathcal{A}} \ x_1^{\mathcal{B}_i} \cdots x_n^{\mathcal{A}} \ x_n^{\mathcal{B}_i} = \mathcal{B}_i(t) \mid F \ x_1 \cdots x_n = t \in D\} \\
\\
\mathcal{B}_i(a) = & a^{\mathcal{B}_i} \quad \mathcal{B}_i(x) = x^{\mathcal{B}_i} \quad \mathcal{B}_i(F) = F^{\mathcal{B}_i} \\
\mathcal{B}_i(t_1 \ t_2) = & \mathcal{B}_i(t_1) \ \mathcal{A}(t_2) \ \mathcal{B}_i(t_2) \\
\mathcal{B}_i(\mathbf{case} \ t \ \mathbf{of} \ \{a_j \ \tilde{y}_j \Rightarrow t_j\}_{j=1}^n) = & \\
& \mathbf{br} \ \mathcal{B}_i(t) \ (\mathbf{case} \ \mathcal{A}(t) \ \mathbf{of} \ \{a_j \ \tilde{y}_j^{\mathcal{A}} \Rightarrow [\widetilde{\mathbf{emp}}/\tilde{y}_j^{\mathcal{B}_i}] \mathcal{B}_i(t_j)\}_{j=1}^n) \\
\mathcal{B}_i(\mathbf{gen}^L) = & \mathbf{emp} \\
\mathcal{B}_i(\mathbf{coerce}^{L_j}(t)) = & \begin{cases} \mathbf{br} \ \mathcal{A}(t) \ \mathcal{B}_i(t) & (\text{if } i = j) \\ \mathcal{B}_i(t) & (\text{otherwise}) \end{cases}
\end{aligned}$$

Here,  $\mathcal{A}(t)$  is the term obtained by replacing every coercion  $\mathbf{coerce}^L(\cdot)$ , variable  $x$ , and function name  $F$  in  $t$  with  $\mathbf{gen}^L$ ,  $x^{\mathcal{A}}$ , and  $F^{\mathcal{A}}$  respectively.  $\mathbf{br} \ t_1 \cdots t_n$  stands for  $\mathbf{br} \ t_1 \ (\mathbf{br} \ t_2 \ (\mathbf{br} \cdots (\mathbf{br} \ t_{n-1} \ t_n)))$  if  $n \geq 2$ ,  $t_1$  if  $n = 1$ , and  $\mathbf{emp}$  if  $n = 0$ . For each terminal  $a \in \text{dom}(\Sigma)$ , we obtain the new function  $a^{\mathcal{B}_i}$  that generates all the trees passed to the  $i$ -th coercion by the actual arguments of  $a$ .

Given a trivial automaton  $\mathcal{M}(L_i) = (\Sigma, Q, \Delta, q_0)$  for accepting  $L_i$ , the output specification for  $\mathcal{P}^{\mathcal{B}_i}$  is the trivial automaton  $\mathcal{B}(L_i) = (\Sigma', Q, \Delta', q_0)$ , where:

$$\begin{aligned}
\Sigma' &= \Sigma \cup \{\mathbf{br} \mapsto 2, \mathbf{emp} \mapsto 0\} \\
\Delta' &= \Delta \cup \{(q, \mathbf{br}, q \cdot q), (q, \mathbf{emp}, \epsilon) \mid q \in Q\}
\end{aligned}$$

The following theorem states the correctness of the construction of  $\mathcal{P}^{\mathcal{B}_i}$ . See the full version of this paper for the proof [23].

**Theorem 2.** *Let  $\mathcal{P} = (D, S)$  be an EHMTT and suppose that the coercions in  $\mathcal{P}$  are  $\mathbf{coerce}^{L_1}(\cdot), \dots, \mathbf{coerce}^{L_m}(\cdot)$ . If  $\models (\mathcal{P}^{\mathcal{B}_i}, \mathcal{M}_1, \dots, \mathcal{M}_k, \mathcal{B}(L_i))$  holds for each  $i \in \{1, \dots, m\}$ , for any  $T_1 \in \mathcal{L}(\mathcal{M}_1), \dots, T_k \in \mathcal{L}(\mathcal{M}_k)$ ,  $S \ T_1 \cdots T_k \not\vdash_{\mathcal{P}}^* \mathbf{Error}$ .*

Our reduction from EHMTT to HMTT verification problems is incomplete, however, i.e. there is a case that an EHMTT satisfies a given specification, but the generated HMTTs do not satisfy the required properties. There are two main reasons for this.

- Coercion annotations may not be good enough. For example, if coercions are annotated with the empty language  $\emptyset$ , the derived HMTTs obviously do not satisfy the property. Actually, there may be no good way to annotate coercions. For example, consider the EHMTT  $S \ x = \text{Zip}(\mathbf{coerce}^L(\text{Unzip} \ x))$ , where  $\text{Unzip}$  takes an input  $\mathbf{s}^n \mathbf{z}$  that encodes a natural number  $n$  and returns an output tree  $\mathbf{pair}(\mathbf{s}^n \mathbf{z})(\mathbf{s}^n \mathbf{z})$ , and  $\text{Zip}$  takes an input tree of the form  $\mathbf{pair}(\mathbf{s}^{n_1} \mathbf{z})(\mathbf{s}^{n_2} \mathbf{z})$  and outputs  $\mathbf{fail}$  if and only if  $n_1 \neq n_2$ . To verify that  $S \ x$  never outputs  $\mathbf{fail}$  for any  $x \in \{\mathbf{s}^n \mathbf{z} \mid n \geq 0\}$ , we need the coercion annotation  $L = \{\mathbf{pair}(\mathbf{s}^n \mathbf{z})(\mathbf{s}^n \mathbf{z}) \mid n \geq 0\}$ , which cannot be expressed by

a trivial automaton or a regular language. As another example, consider the following variant of a reverse function:

$$\begin{array}{l} \text{Reverse } x = \text{case } x \text{ of } e \Rightarrow e \\ \quad | \text{cons } z \ x' \Rightarrow \text{Append } (\text{coerce}^L(\text{Reverse } x')) \ (\text{cons } z \ e) \\ \text{Append } x \ y = \dots \end{array}$$

Here, we have used a list-like representation of sequences of **a**, **b**. In this case, the appropriate annotation depends on the value of  $z$  ( $L$  should be **b\*a\*e** if  $z$  is **a** while **b\*e** if  $z$  is **b**), which cannot be expressed in our language. One way to avoid this problem is to duplicate a part of the code so that appropriate annotations can be inserted.

- The output specification  $\mathcal{B}(L_i)$  for  $\mathcal{P}^{\mathcal{B}_i}$  is too restrictive. When the automaton for accepting  $L_i$  is non-deterministic,  $\mathcal{B}(L_i)$  does not accept all the tree representations of subsets of  $L_i$ . This problem can easily be remedied, however, by using a more elaborate construction of  $\mathcal{B}(L_i)$ , hence not a fundamental limitation.

Our overall method is also incomplete because of the incompleteness of the HMTT verification method [15] (unsurprisingly, as the HMTT verification problem is undecidable in general).

## 5 Experiments

We have implemented the reduction method from an EHMTT verification problem to HMTT verification problems presented in Section 4. For solving the HMTT verification problems, we adopted an HMTT verification method in [15] and Kobayashi’s higher-order model checker TRECS [12].

Table 1 shows the results of preliminary experiments. The column “O” shows the order of each EHMTT which is the largest order of the sorts of the functions. The order of a sort is defined by:

$$order(\mathbf{i}) = order(\mathbf{o}) = 0 \quad order(\kappa_1 \rightarrow \kappa_2) = \max(order(\kappa_1) + 1, order(\kappa_2))$$

The column “C” shows the number of coercions in each EHMTT. The columns “R” and “S” are the number of rules and the size of each EHMTT respectively. The size of an EHMTT is measured by the number of symbols occurring in the right-hand side of the rewriting rules. “Sum<sub>R</sub>” and “Sum<sub>S</sub>” respectively are the sum of the numbers of the rules and the sum of the sizes of all HMTTs derived from each EHMTT. “Q<sub>I</sub>” and “Q<sub>O</sub>” respectively show the numbers of the states of trivial automata for the input and output specifications. The column “T<sub>Red</sub>” shows the elapsed time, in milliseconds, of reduction from an EHMTT verification problem to HMTT verification problems. The column “Y/N” indicates whether each EHMTT was proved correct (Y) or rejected (N). The column “T<sub>MC</sub>” shows the total running time of the higher-order model checker TRECS to solve all the HMTT verification problems derived from each EHMTT.

**Table 1.** Experimental Results

Programs	O	C	R	S	Sum <sub>R</sub>	Sum <sub>S</sub>	Q <sub>I</sub>	Q <sub>O</sub>	T <sub>Red</sub>	Y/N	T <sub>MC</sub>
<b>Reverse</b>	1	2	3	32	23	222	4	2	1	Y	4
<b>Isort</b>	1	1	4	29	16	115	3	2	1	Y	3
<b>Msort</b>	2	4	8	131	88	1,731	3	2	2	Y	224
<b>HomRep-Rev</b>	4	1	12	90	43	362	6	2	1	Y	31
<b>Split</b>	2	1	6	126	33	572	23	9	3	Y	132
<b>Bib2Html</b>	2	1	13	493	126	2,303	59	50	52	Y	52
<b>XMarkQ1</b>	2	1	12	454	118	2,136	99	23	29	Y	168
<b>XMarkQ2</b>	1	2	9	461	207	3,797	99	4	77	Y	92
<b>Gapid-Html</b>	3	1	17	374	75	1,642	16	7	2	Y	112
<b>JWIG-guess</b>	2	1	6	465	98	2,331	64	50	588	Y	50
<b>JWIG-cal</b>	1	2	12	475	222	4,045	60	50	72	Y	73
<b>MinCaml-K</b>	2	8	19	605	563	16,117	5	3	5	Y	647
<b>Split'</b>	2	1	6	126	33	572	23	9	3	N	27
<b>JWIG-guess'</b>	2	1	6	465	98	2,331	64	50	586	N	49
<b>JWIG-cal'</b>	1	2	12	475	222	4,045	60	50	2	N	55

The program **Reverse** is the same as the one presented in the paper. **Isort** performs insertion sort on the lists encoded as linear trees over  $\Sigma = \{a \mapsto 1, b \mapsto 1, e \mapsto 0\}$ . **Msort** performs merge sort instead of insertion sort on the same linear trees. **HomRep-Rev** takes a word homomorphism  $h$  over linear trees  $(a + b)^*e$ , a number  $n$  and a word  $w \in (a + b)^*e$ , and produces the reverse of the image  $h^n(w)$ . We let  $h = \{a \mapsto bb, b \mapsto a\}$  and verified that if  $n$  is an even number and  $w \in a^*b^*e$ , then the reversed image is in  $b^*a^*e$ . The program **Split** presented in Figure 3 is taken from sample programs of CDuce [2], a higher-order XML-oriented functional language. **Split** takes a list of persons, and splits it into two lists of men and women. **Bib2Html** also simulates a CDuce program that transforms a list of bibliography into an XHTML. **XMarkQ1** and **XMarkQ2** taken from Q1 and Q2 of XMark benchmark suite [21] simulate simple XQuery queries. The program **Gapid-Html** is a composition of Tozawa's high-level tree transducers [22]. It takes a document of the following DTD:

```

type Doc      = doc[Preface, (Div|P|Note)*]
type Preface  = preface[Header, P*]
type Header   = header[A*]
type P        = p[A*]
type Div      = div[(Div|P|Note|A)*]
type Note     = note[(P|A)*]
type A        = a[A*]

```

and another tree as inputs. It checks whether the children of each node of the document are empty, and if so, replaces the empty children with a “hole”. The program then inserts the given tree into the holes. The program finally transforms the result to an XHTML. The programs **JWIG-guess** and **JWIG-cal** are

taken from sample programs of JWIG (<http://www.brics.dk/JWIG/>), a programming language for interactive Web services. A main feature of JWIG is document templates. For example, the following document template represents an HTML document with a hole named  $x$ :

```
<html>
  <head><title> ... </title></head>
  <body><[x]></body>
</html>
```

We can instantiate the template by substituting another document or template for  $x$ . In EHMTTs, the template can be encoded as the following rule for a function  $T$  with an argument  $x$ :

$$T\ x \rightarrow \text{html}(\text{head}(\text{title}(\text{text leaf leaf})\ \text{leaf})\ (\text{body } x\ \text{leaf}))\ \text{leaf}$$

The program **JWIG-guess** is a number guessing game. The program **JWIG-cal** is a web-based calendar service. **MinCaml-K** simulates the K-normalization routine of the MinCaml compiler (<http://min-caml.sourceforge.net/index-e.html>). Finally, **Split'**, **JWIG-guess'** and **JWIG-cal'** are respectively the same as **Split**, **JWIG-guess** and **JWIG-cal** except that they involve wrong coercions which lead to an **Error** to see that programs with wrong coercions are rejected. These programs (except for **Reverse**, **Isort**, **Msort** and **HomRep-Rev**) are manually translated from original source codes to EHMTTs.

All the valid programs have been proved correct by our verification method despite its incompleteness, while wrong programs are correctly rejected. The number of coercions (thus the number of annotations required by our method) is much smaller than the number of rules (functions) in all cases. Though these numbers depend on the particular encoding, this result witnesses that our verification method usually requires fewer annotations than existing verification methods [9,2,4], which require type annotations for every function definition. Further comparison with the existing methods on this point is given in Section 6.

All the programs were proved correct within 1 second. From this, we can expect that our method can verify non-trivial programs reasonably fast despite the high time complexity ( $n$ -EXPTIME complete, where  $n$  is the order) of higher-order model checking.

## 6 Related Work

As shown in Section 1, our verification method is based on recent advances on higher-order model checking [11,1,20,14]. Ong [20] has proven the decidability of the model checking problem for recursion schemes, and Kobayashi has developed and implemented a type-based model checking algorithm [14].

As we have shown in Section 5, our EHMTT verification method can be applied to verification of functional programs that manipulate various data structures such as strings, lists, trees, XML, and user-defined recursive data structures

```

Split x = case x of person g n c =>
  case c of children cs => case g of gender gend =>
    Let gend n (coerce qPair (MakePair cs nil nil))
  Let gend n x = case gend of
    m => Make man (Copy n) x
  | f => Make woman (Copy n) x
  MakePair ps ms fs = case ps of nil => pair ms fs
  | cons p sib => case p of person g n c => case g of gender gend =>
    case gend of
      m => MakePair sib (cons (person (gender m) (Copy n) (Copy c)) ms) fs
    | f => MakePair sib ms (cons (person (gender f) (Copy n) (Copy c)) fs)
  Make tag name sdpair = case sdpair of pair s d =>
    tag name (sons (RevMap Split s nil))
    (daughters (RevMap Split d nil))
  RevMap f l ac =
    case l of nil => ac | cons x xs => RevMap f xs (cons (f x) ac)

```

Fig. 3. Split

by encoding them as trees and adding coercion annotations to the programs. We compare our method with existing verification methods below.

Refinement types [7,4] can be used for verification of functional programs that manipulate user-defined recursive data structures. The original refinement type system [7] uses a naïve least fixed-point algorithm to infer the most precise refinement types of functions, and does not seem to scale for higher-order functions. Another refinement type system proposed by Davies [4] requires users to write type annotations for each function. For example, for  $\mathcal{P}_{rev}$  in Example 1, *Reverse* and *Append* need to be annotated with the following intersection types:

$$\begin{aligned}
Reverse &: (a^*b^*e \rightarrow b^*a^*e) \wedge (b^*e \rightarrow b^*e) \\
Append &: (b^*a^*e \rightarrow a^*e \rightarrow b^*a^*e) \wedge (b^*e \rightarrow b^*e \rightarrow b^*e)
\end{aligned}$$

In contrast, our method requires only the annotations  $\mathbf{coerce}^{b^*a^*e}(Reverse\ x')$  and  $\mathbf{coerce}^{b^*e}(Reverse\ x')$  in the definition of *Reverse*. As in this case, we expect that coercion annotations required in our approach tends to be simpler than refinement type declarations. Because of the limitation of our approach discussed at the end of Section 4, however, it may be useful to combine both approaches.

Several research groups have proposed typed XML processing languages [9,2]. Their type systems can be used for verification of XML processing programs. As in the Davies's refinement type system, these type systems require type annotations for each function for type checking. Thus, our method can be used as an alternative for a verification purpose. Meanwhile their type systems support advanced programming features such as parametric polymorphism [8] and regular expression pattern matching [10]. Extensions of our method with these features are left for future work. While the type checking of the XML processing languages are incomplete, extensive work has been done on complete type checking of various tree transducers [18,17]. They are not Turing-complete, however, and

thus less expressive than our EHMTTs. As shown in [16], ordinary macro and high-level tree transducers [5,6] are subsumed by linear HMTTs, for which our EHMTT verification method is sound and complete.

String analysis [3,19] can verify programs that manipulate strings by approximating a string-processing program as a regular or a context-free grammar. In contrast, our method is more precise since we can naturally model programs as EHMTTs, which are strictly more expressive than context-free grammars.

Our approach of reducing EHMTT verification to simpler verification problems based on coercion annotations is a reminiscent of program verification techniques for imperative languages based on verification condition generation from loop invariants: coercion annotations are invariants, and generated HMTT verification problems can be considered verification conditions. The main differences are that our target is a higher-order functional language and that not all recursions (or loops) need to be annotated with invariants.

## 7 Conclusion

We have proposed a verification method for tree-processing programs based on reduction to higher-order model checking, and shown its effectiveness through experiments. We plan to investigate techniques for inferring coercion annotations, which would enable fully automatic verification of tree-processing programs. If the derived HMTTs are ordinary macro or high-level tree transducers [5,6], annotations can indeed be inferred by the inverse inference technique. For general EHMTTs, we plan to apply techniques of machine learning. Addressing the limitations discussed at the end of Section 4 is also left for future work.

## Acknowledgment

We would like to thank anonymous referees for useful comments.

## References

1. Aehlig, K., de Miranda, J.G., Ong, C.H.L.: The monadic second order theory of trees given by arbitrary level-two recursion schemes is decidable. In: Urzyczyn, P. (ed.) TLCA 2005. LNCS, vol. 3461, pp. 39–54. Springer, Heidelberg (2005)
2. Benzaken, V., Castagna, G., Frisch, A.: CDuce: an XML-centric general-purpose language. In: ICFP 2003, pp. 51–63. ACM, New York (2003)
3. Christensen, A.S., Møller, A., Schwartzbach, M.I.: Precise analysis of string expressions. In: Cousot, R. (ed.) SAS 2003. LNCS, vol. 2694, pp. 1–18. Springer, Heidelberg (2003)
4. Davies, R.: Practical refinement-type checking. Ph.D. thesis, Carnegie Mellon University, chair-Pfenning, Frank (2005)
5. Engelfriet, J., Vogler, H.: Macro tree transducers. *Journal of Computer and System Sciences* 31(1), 71–146 (1985)
6. Engelfriet, J., Vogler, H.: High level tree transducers and iterated pushdown tree transducers. *Acta Informatica* 26(1/2), 131–192 (1988)

7. Freeman, T., Pfenning, F.: Refinement types for ML. In: PLDI 1991, pp. 268–277. ACM, New York (1991)
8. Hosoya, H., Frisch, A., Castagna, G.: Parametric polymorphism for XML. *ACM Transactions on Programming Languages and Systems* 32(1), 1–56 (2009)
9. Hosoya, H., Pierce, B.C.: XDuce: A statically typed XML processing language. *ACM Transactions on Internet Technology* 3(2), 117–148 (2003)
10. Hosoya, H., Vouillon, J., Pierce, B.C.: Regular expression types for XML. In: ICFP 2000, pp. 11–22. ACM, New York (2000)
11. Knapik, T., Niwinski, D., Urzyczyn, P.: Higher-order pushdown trees are easy. In: Nielsen, M., Engberg, U. (eds.) FOSSACS 2002. LNCS, vol. 2303, pp. 205–222. Springer, Heidelberg (2002)
12. Kobayashi, N.: Model-checking higher-order functions. In: PPDP 2009, pp. 25–36. ACM, New York (2009)
13. Kobayashi, N.: Types and higher-order recursion schemes for verification of higher-order programs. In: POPL 2009, pp. 416–428. ACM, New York (2009)
14. Kobayashi, N., Ong, C.-H.L.: A type system equivalent to the modal mu-calculus model checking of higher-order recursion schemes. In: LICS 2009, pp. 179–188. IEEE, Los Alamitos (2009)
15. Kobayashi, N., Tabuchi, N., Unno, H.: Higher-order multi-parameter tree transducers and recursion schemes for program verification. In: POPL 2010, pp. 495–508. ACM, New York (2010)
16. Kobayashi, N., Tabuchi, N., Unno, H.: Higher-order multi-parameter tree transducers and recursion schemes for program verification. An extended version (2010), <http://www.kb.ecei.tohoku.ac.jp/~koba/papers/hmtt.pdf>
17. Maneth, S., Berlea, A., Perst, T., Seidl, H.: XML type checking with macro tree transducers. In: PODS 2005, pp. 283–294. ACM, New York (2005)
18. Milo, T., Suciu, D., Vianu, V.: Typechecking for XML transformers. *Journal of Computer and System Sciences* 66(1), 66–97 (2003)
19. Minamide, Y.: Static approximation of dynamically generated web pages. In: WWW 2005, pp. 432–441. ACM, New York (2005)
20. Ong, C.-H.L.: On model-checking trees generated by higher-order recursion schemes. In: LICS 2006, pp. 81–90. IEEE, Los Alamitos (2006)
21. Schmidt, A., Waas, F., Kersten, M., Carey, M.J., Manolescu, I., Busse, R.: XMark: a benchmark for XML data management. In: VLDB 2002, pp. 974–985. VLDB Endowment (2002)
22. Tozawa, A.: XML type checking using high-level tree transducer. In: Hagiya, M., Wadler, P. (eds.) FLOPS 2006. LNCS, vol. 3945, pp. 81–96. Springer, Heidelberg (2006)
23. Unno, H., Tabuchi, N., Kobayashi, N.: Verification of tree-processing programs via higher-order model checking. An extended version (2010), <http://www.kb.ecei.tohoku.ac.jp/~uhiro/papers/aplas2010.pdf>