

Computing Simulations over Tree Automata (Efficient Techniques for Reducing Tree Automata)

Parosh A. Abdulla¹, Ahmed Bouajjani², Lukáš Holík³, Lisa Kaati¹, and Tomáš Vojnar³

¹ University of Uppsala, Sweden
{parosh,lisa.kaati}@it.uu.se

² LIAFA, University Paris 7, France
abou@liafa.jussieu.fr

³ FIT, Brno University of Technology, Czech Republic
{holik,vojnar}@fit.vutbr.cz

Abstract. We address the problem of computing simulation relations over tree automata. In particular, we consider downward and upward simulations on tree automata, which are, loosely speaking, analogous to forward and backward relations over word automata. We provide simple and efficient algorithms for computing these relations based on a reduction to the problem of computing simulations on labelled transition systems. Furthermore, we show that downward and upward relations can be combined to get relations compatible with the tree language equivalence, which can subsequently be used for an efficient size reduction of nondeterministic tree automata. This is of a very high interest, for instance, for symbolic verification methods such as regular model checking, which use tree automata to represent infinite sets of reachable configurations. We provide experimental results showing the efficiency of our algorithms on examples of tree automata taken from regular model checking computations.

1 Introduction

Tree automata are widely used for modelling and reasoning about various kinds of structured objects such as syntactical trees, structured documents, configurations of complex systems, algebraic term representations of data or computations, etc. (see [9]). For instance, in the framework of regular model checking, tree automata are used to represent and manipulate sets of configurations of infinite-state systems such as parameterized networks of processes with a tree-like topology, or programs with dynamic linked data-structures [7,3,5,6].

In the above context, checking language equivalence and reducing automata wrt. the language equivalence is a fundamental issue, and performing these operations efficiently is crucial for all practical applications of tree automata. Computing a minimal canonical tree automaton is, of course, possible, but it requires determinisation, which may lead to an exponential blow-up in the size of the automaton. Therefore, even if the resulting automaton can be small, we may not be able to compute it in practice due to the very expensive determinisation step, which is, indeed, a major bottleneck when using canonical tree automata.

A reasonable and pragmatic approach is to consider a notion of equivalence that is stronger than language equivalence, but which can be checked efficiently, using a polynomial algorithm. Here, a natural trade-off between the strength of the considered

equivalence and the cost of its computation arises. In the case of word automata, an equivalence which is widely considered as a good trade-off in this sense is simulation equivalence. It can be checked in polynomial time, and efficient algorithms have been designed for this purpose (see, e.g., [10,14]). These algorithms make the computation of simulation equivalence quite affordable even in comparison with the one of bisimulation, which is cheaper [13], but which is also stronger, and therefore leads in general to less significant reductions in the sizes of the automata.

In this paper, we study notions of entailment and equivalence between tree automata, which are suitable in the sense discussed above, and we also provide efficient algorithms for their computation.

We start by considering a basic notion of tree simulation, called *downward simulation*, corresponding to a natural extension of the usual notion of simulation defined on *or*-structures to *and-or* structures. This relation can be shown to be compatible with the tree language equivalence.

The second notion of simulation that we consider, called *upward simulation*, corresponds intuitively to a generalisation of the notion of backward simulation to *and-or* structures. The definition of an upward simulation is parametrised by a downward simulation: Roughly speaking, two states q and q' are upward similar if whenever one of them, say q , considered within some vector (q_1, \dots, q_n) at position i , has an upward transition to some state s , then q' appears at position i of some vector (q'_1, \dots, q'_n) that has also an upward transition to a state s' , which is upward similar to s , and moreover, for each position $j \neq i$, q_j is downward similar to q'_j .

Upward simulation is not compatible with the tree language equivalence. It is rather compatible with the so-called context language equivalence, where a context of a state q is a tree with a hole on the leaf level such that if we plug a tree in the tree language of q into this hole, we obtain a tree recognised by the automaton. However, we show an interesting fact that when we restrict ourselves to upward relations compatible with the set of final states of automata, the downward and upward simulation equivalences can be *combined* in such a way that they give rise to a new equivalence relation which is compatible with the tree language equivalence. This combination is not trivial. It is based on the idea that two states q_1 and q_2 may have different tree languages and different context languages, but for every t in the tree language of one of them, say q_1 , and every C in the context language of the other, here q_2 , the tree $C[t]$ (where t is plugged into C) is recognised by the automaton. The combined relation is coarser than (or, in the worst case, as coarse as) the downward simulation and according to our practical experiments, it usually leads to significantly better reductions of the automata.

In this way, we obtain two candidates for simulation-based equivalences for use in automata reduction. Then, we consider the issue of designing efficient algorithms for computing these relations. A deep examination of downward and upward simulation equivalences shows that they can be computed using essentially the same algorithmic pattern. Actually, we prove that, surprisingly, computing downward and upward tree simulations can be reduced in each case to computing simulations on standard labelled transition systems. These reductions provide a simple and elegant way of solving in a uniform way the problem of computing tree simulations by reduction to computing simulations in the word case. The best known algorithm for solving the latter problem,

published recently in [14], considers simulation relations defined on Kripke structures. The use of this algorithm requires its adaptation to labelled transition systems. We provide such an adaptation and we provide also a proof for this algorithm which can be seen as an alternative, more direct, proof of the algorithm of [14]. The combination of our reductions with the labelled transition systems-based simulation algorithm leads to efficient algorithms for our equivalence relations on tree automata, whose precise complexities are also analysed in the paper.

We have implemented our algorithms and performed experiments on automata computed in the context of regular tree model checking (corresponding to representations of the set of reachable configurations of parametrised systems). The experiments show that, indeed, the relations proposed in this paper provide significant reductions of these automata and that they perform better than (existing) bisimulation-based reductions [11].

Related work. As far as we know, this is the first work which addresses the issue of computing simulation relations for tree automata. The downward and upward simulation relations considered in this work have been introduced first in [4] where they have been used for proving soundness of some acceleration techniques used in the context of regular tree model checking. However, the problem of computing these relations has not been addressed in that paper. A form of combining downward and upward relations has also been defined in [4]. However, the combinations considered in that paper require some restrictions which are computationally difficult to check and that are not considered in this work. Bisimulations on tree automata have been considered in [2,11]. The notion of a backward bisimulation used in [11] corresponds to what can be called a downward bisimulation in our terminology.

Outline. The rest of the paper is organised as follows. In the next section, we give some preliminaries on tree automata, labelled transition systems, and simulation relations. Section 3 describes an algorithm for checking simulation on labelled transition systems. In Section 4 resp. Section 5, we translate downward resp. upward simulation on tree automata into corresponding simulations on labelled transition systems. Section 6 gives methods for reducing tree automata based on equivalences derived from downward and upward simulation. In Section 7, we report some experimental results. Finally, we give conclusions and directions for future research in Section 8.

Remark. For space reasons, all proofs are deferred to [1].

2 Preliminaries

In this section, we introduce some preliminaries on trees, tree automata, and labelled transition systems (LTS). In particular, we recall two simulation relations defined on tree automata in [4], and the classical (word) simulation relation defined on LTS. Finally, we will describe an encoding which we use in our algorithms to describe pre-order relations, e.g., simulation relations.

For an equivalence relation \equiv defined on a set Q , we call each equivalence class of \equiv a *block*, and use Q/\equiv to denote the set of blocks in \equiv .

Trees. A *ranked alphabet* Σ is a set of symbols together with a function $\text{Rank} : \Sigma \rightarrow \mathbb{N}$. For $f \in \Sigma$, the value $\text{Rank}(f)$ is said to be the *rank* of f . For any $n \geq 0$, we denote by Σ_n the set of all symbols of rank n from Σ . Let ε denote the empty sequence. A *tree* t over an alphabet Σ is a partial mapping $t : \mathbb{N}^* \rightarrow \Sigma$ that satisfies the following conditions:

- $\text{dom}(t)$ is a finite, prefix-closed subset of \mathbb{N}^* , and
- for each $p \in \text{dom}(t)$, if $\text{Rank}(t(p)) = n > 0$, then $\{i \mid pi \in \text{dom}(t)\} = \{1, \dots, n\}$.

Each sequence $p \in \text{dom}(t)$ is called a *node* of t . For a node p , we define the i^{th} *child* of p to be the node pi , and we define the i^{th} *subtree* of p to be the tree t' such that $t'(p') = t(pip')$ for all $p' \in \mathbb{N}^*$. A *leaf* of t is a node p which does not have any children, i.e., there is no $i \in \mathbb{N}$ with $pi \in \text{dom}(t)$. We denote by $T(\Sigma)$ the set of all trees over the alphabet Σ .

Tree Automata. A (finite, non-deterministic, bottom-up) *tree automaton* (TA) is a 4-tuple $A = (Q, \Sigma, \Delta, F)$ where Q is a finite set of states, $F \subseteq Q$ is a set of final states, Σ is a ranked alphabet, and Δ is a set of transition rules. Each transition rule is a triple of the form $((q_1, \dots, q_n), f, q)$ where $q_1, \dots, q_n, q \in Q$, $f \in \Sigma$, and $\text{Rank}(f) = n$. We use $(q_1, \dots, q_n) \xrightarrow{f} q$ to denote that $((q_1, \dots, q_n), f, q) \in \Delta$. In the special case where $n = 0$, we speak about the so-called *leaf rules*, which we sometimes abbreviate as $\xrightarrow{f} q$. We use $\text{Lhs}(A)$ to denote the set of *left-hand sides* of rules, i.e., the set of tuples of the form (q_1, \dots, q_n) where $(q_1, \dots, q_n) \xrightarrow{f} q$ for some f and q . Finally, we denote by $\text{Rank}(A)$ the smallest $n \in \mathbb{N}$ such that $n \geq m$ for each $m \in \mathbb{N}$ where $(q_1, \dots, q_m) \in \text{Lhs}(A)$ for some $q_i \in Q$, $1 \leq i \leq m$.

A *run* of A over a tree $t \in T(\Sigma)$ is a mapping $\pi : \text{dom}(t) \rightarrow Q$ such that for each node $p \in \text{dom}(t)$ where $q = \pi(p)$, we have that if $q_i = \pi(pi)$ for $1 \leq i \leq n$, then Δ has a rule $(q_1, \dots, q_n) \xrightarrow{t(p)} q$. We write $t \xrightarrow{\pi} q$ to denote that π is a run of A over t such that $\pi(\varepsilon) = q$. We use $t \Longrightarrow q$ to denote that $t \xrightarrow{\pi} q$ for some run π . The *language* of a state $q \in Q$ is defined by $L(q) = \{t \mid t \Longrightarrow q\}$, while the *language* of A is defined by $L(A) = \bigcup_{q \in F} L(q)$.

Labelled Transition Systems. A (finite) *labelled transition system* (LTS) is a tuple $T = (S, \mathcal{L}, \rightarrow)$ where S is a finite set of states, \mathcal{L} is a finite set of labels, and $\rightarrow \subseteq S \times \mathcal{L} \times S$ is a transition relation.

Given an LTS $T = (S, \mathcal{L}, \rightarrow)$, a label $a \in \mathcal{L}$, and two states $q, r \in S$, we denote by $q \xrightarrow{a} r$ the fact that $(q, a, r) \in \rightarrow$. We define the set of *a-predecessors* of a state r as $\text{pre}_a(r) = \{q \in S \mid q \xrightarrow{a} r\}$. Given $X, Y \subseteq S$, we denote $\text{pre}_a(X)$ the set $\bigcup_{s \in X} \text{pre}_a(s)$, we write $q \xrightarrow{a} X$ iff $q \in \text{pre}_a(X)$, and $Y \xrightarrow{a} X$ iff $Y \cap \text{pre}_a(X) \neq \emptyset$.

Simulations. For a tree automaton $A = (Q, \Sigma, \Delta, F)$, a *downward simulation* D is a binary relation on Q such that if $(q, r) \in D$ and $(q_1, \dots, q_n) \xrightarrow{f} q$, then there are r_1, \dots, r_n such that $(r_1, \dots, r_n) \xrightarrow{f} r$ and $(q_i, r_i) \in D$ for each i such that $1 \leq i \leq n$. It is easy to show [4] that any downward simulation can be closed under reflexivity and transitivity. Moreover, there is a unique maximal downward simulation over a given tree automaton, which we denote as \preceq_{down} in the sequel.

Given a TA $A = (Q, \Sigma, \Delta, F)$ and a downward simulation D , an *upward simulation* U induced by D is a binary relation on Q such that if $(q, r) \in U$ and $(q_1, \dots, q_n) \xrightarrow{f} q'$ with $q_i = q$, $1 \leq i \leq n$, then there are r_1, \dots, r_n, r' such that $(r_1, \dots, r_n) \xrightarrow{f} r'$ where $r_i = r$, $(q', r') \in U$, and $(q_j, r_j) \in D$ for each j such that $1 \leq j \neq i \leq n$. In [4], it is shown that any upward simulation can be closed under reflexivity and transitivity. Moreover, there is a unique maximal upward simulation with respect to a fixed downward simulation over a given tree automaton, which we denote as \preceq_{up} in the sequel.

Given an *initial* pre-order $I \subseteq Q \times Q$, it can be shown that there are unique maximal downward as well as upward simulations included in I on the given TA, which we denote \preceq_x^I in the sequel, for $x \in \{down, up\}$. Further, we use \cong_x to denote the equivalence relation $\preceq_x \cap \preceq_x^{-1}$ on Q for $x \in \{down, up\}$. Likewise, we define the equivalence relations \cong_x^I for an initial pre-order I on Q and $x \in \{down, up\}$.

For an LTS $T = (S, \mathcal{L}, \rightarrow)$, a (*word*) *simulation* is a binary relation R on S such that if $(q, r) \in R$ and $q \xrightarrow{a} q'$, then there is an r' with $r \xrightarrow{a} r'$ and $(q', r') \in R$. In a very similar way as for simulations on trees, it can be shown that any given simulation on an LTS can be closed under reflexivity and transitivity and that there is a unique maximal simulation on the given LTS, which will we denote by \preceq . Moreover, given an *initial* pre-order $I \subseteq S \times S$, it can be shown that there is a unique maximal simulation included in I on the given LTS, which we denote \preceq^I in the sequel. We use \cong to denote the equivalence relation $\preceq \cap \preceq^{-1}$ on S and consequently \cong^I to denote $\preceq^I \cap (\preceq^I)^{-1}$.

Encoding. Let S be a set. A *partition-relation pair* over S is a pair (P, Rel) where (1) $P \subseteq 2^S$ is a partition of S (i.e., $S = \cup_{B \in P} B$, and for all $B, C \in P$, if $B \neq C$, then $B \cap C = \emptyset$), and (2) $Rel \subseteq P \times P$. We say that a partition-relation pair (P, Rel) over S *induces* (or *defines*) the relation $\delta = \bigcup_{(B,C) \in Rel} B \times C$.

Let \preceq be a pre-order defined on a set S , and let \equiv be the equivalence $\preceq \cap \preceq^{-1}$ defined by \preceq . The pre-order \preceq can be represented—which we will use in our algorithms below—by a partition-relation pair (P, Rel) over S such that $(B, C) \in Rel$ iff $s_1 \preceq s_2$ for all $s_1 \in B$ and $s_2 \in C$. In this representation, if the partition P is as coarse as possible (i.e., such that $s_1, s_2 \in B$ iff $s_1 \equiv s_2$), then, intuitively, the elements of P are blocks of \equiv , while Rel reflects the partial order on P corresponding to \preceq .

3 Computing Simulations on Labelled Transition Systems

We now introduce an algorithm to compute the (unique) maximal simulation relation \preceq^I on an LTS for a given initial pre-order I on states. Our algorithm is a re-formulation of the algorithm proposed in [14] for computing simulations over *Kripke structures*.

3.1 An Algorithm for Computing Simulations on LTS

For the rest of this section, we assume that we are given an LTS $T = (S, \mathcal{L}, \rightarrow)$ and an initial pre-order $I \subseteq S \times S$. We will use Algorithm 1 to compute the maximum simulation $\preceq^I \subseteq S \times S$ included in I . In the algorithm, we use the following notation. Given $\rho \subseteq S \times S$ and an element $q \in S$, we denote $\rho(q)$ the set $\{r \in S \mid (q, r) \in \rho\}$.

The algorithm performs a number of iterations computing a sequence of relations, each induced by a partition-relation pair (P, Rel) . During each iteration, the states belonging to a block $B' \in P$ are those which are currently assumed as capable of simulating those from any B with $(B, B') \in Rel$. The algorithm starts with an initial partition-relation pair (P_{init}, Rel_{init}) that induces the initial pre-order I on S . The partition-relation pair is then gradually refined by splitting blocks of the partition P and by restricting the relation Rel on P . When the algorithm terminates, the final partition-relation pair (P_{sim}, Rel_{sim}) induces the required pre-order \preceq^I .

The refinement performed during the iterations consists of splitting the blocks in P and then updating the relation Rel accordingly. For this purpose, the algorithm maintains a set $Remove_a(B)$ for each $a \in \mathcal{L}$ and $B \in P$. Such a set contains states that do not have an a -transition going into states that are in B nor to states of any block B' with $(B, B') \in Rel$. Clearly, the states in $Remove_a(B)$ cannot simulate states that have an a -transition going into $\bigcup_{(B, B') \in Rel} B'$. Therefore, for any $Remove_a(B) \neq \emptyset$, we can split each block $C \in P$ to $C \cap Remove_a(B)$ and $C \setminus Remove_a(B)$. This is done using the function *Split* on line 6.

After performing the *Split* operation, we update the relation Rel and the *Remove* sets. This is carried out in two steps. First, we compute an approximation of the next values of Rel and *Remove*. More precisely, after a split, all Rel relations between the original “parent” blocks of states are inherited to their “children” resulting from the split (line 8)—the notation $parent_{p_{prev}}(C)$ refers to the parent block from which C arose within the split. On line 10, the remove sets are then inherited from parent blocks to their children. To perform the second step, we observe that the inheritance of the original relation Rel on parent blocks to the children blocks is not consistent with the split we have just performed. Therefore, on line 14, we subsequently prune Rel such that blocks C that have an a -transition going into B states cannot be considered as simulated by blocks D which do not have an a -transition going into $\bigcup_{(B, B') \in Rel} B'$ —notice that due to the split that we have performed, the D blocks are now included in *Remove*. This pruning can then cause a necessity of further refinements as the states that have some b -transition into a D block (that was freshly found not to simulate C), but not to C nor any block that is still viewed as capable of simulating C , have to stop simulating states that can go into $\bigcup_{(C, C') \in Rel} C'$. Therefore, such states are added into $Remove_b(C)$ on line 17.

3.2 Correctness and Complexity of the Algorithm

In the rest of the section, we assume that Algorithm 1 is applied on an LTS $T = (S, \mathcal{L}, \rightarrow)$ with an initial partition-relation pair (P_{init}, Rel_{init}) . The correctness of the algorithm is formalised in Theorem 1.

Theorem 1. *Suppose that I is the pre-order induced by (P_{init}, Rel_{init}) . Then, Algorithm 1 terminates and the final partition-relation pair (P_{sim}, Rel_{sim}) computed by it induces the simulation relation \preceq^I , and, moreover, $P_{sim} = S / \simeq^I$.*

A similar correctness result is proved in [14] for the algorithm on Kripke structures, using notions from the theory of abstract interpretation. In [1], we provide an alternative, more direct proof, which is, however, beyond the space limitations of this paper. Therefore, we will only mention the key idea behind the termination argument. In particular, the key point is that if we take any block B from P_{init} and any $a \in \mathcal{L}$, if B or any

Algorithm 1. Computing simulations on states of an LTS

Input: An LTS $T = (S, \mathcal{L}, \rightarrow)$, an initial partition-relation pair (P_{init}, Rel_{init}) on S inducing a pre-order $I \subseteq S \times S$.

Data: A partition-relation pair (P, Rel) on S , and for each $B \in P$ and $a \in \mathcal{L}$, a set $Remove_a(B) \subseteq S$.

Output: The partition-relation pair (P_{sim}, Rel_{sim}) inducing the maximal simulation on T contained in I .

```

/* initialisation */
1 (P, Rel) ← (Pinit, Relinit);
2 forall a ∈ ℒ, B ∈ P do Removea(B) ← S \ prea(∪ Rel(B));

/* computation */
3 while ∃a ∈ ℒ. ∃B ∈ P. Removea(B) ≠ ∅ do
4   Remove ← Removea(B); Removea(B) ← ∅;
5   Pprev ← P; Bprev ← B; Relprev ← Rel;
6   P ← Split(P, Remove);
7   forall C ∈ P do
8     Rel(C) ← {D ∈ P | D ⊆ ∪ Relprev(parentPprev(C))};
9     forall b ∈ ℒ do
10      Removeb(C) ← Removeb(parentPprev(C))
11   forall C ∈ P. C  $\xrightarrow{a}$  Bprev do
12     forall D ∈ P. D ⊆ Remove do
13       if (C, D) ∈ Rel then
14         Rel ← Rel \ {(C, D)};
15         forall b ∈ ℒ do
16           forall r ∈ preb(D) \ preb(∪ Rel(C)) do
17             Removeb(C) ← Removeb(C) ∪ {r}
18 (Psim, Relsim) ← (P, Rel);

```

of its children B' , which arises by splitting, is repeatedly selected to be processed by the while loop on line 3, then the $Remove_a(B)$ (or $Remove_a(B')$) sets can never contain a single state $s \in S$ at an iteration i of the while loop as well as on a later iteration j , $j > i$. Therefore, as the number of possible partitions as well as the number of states is finite, the algorithm must terminate.

The complexity of the algorithm is equal to that of the original algorithm from [14], up to the new factor \mathcal{L} that is not present in [14] (or, equivalently, $|\mathcal{L}| = 1$ in [14]). The complexity is stated in Theorem 2.

Theorem 2. *Algorithm 1 has time complexity $O(|\mathcal{L}| \cdot |P_{sim}| \cdot |S| + |P_{sim}| \cdot |\rightarrow|)$ and space complexity $O(|\mathcal{L}| \cdot |P_{sim}| \cdot |S|)$.*

A proof of Theorem 2, based on a similar reasoning as in [14], can be found in [1]. Here, let us just mention that the result expects the input LTS and the initial partition-relation pair be encoded in suitable data structures. This fact is important for the complexity analyses presented later on as they build on using Algorithm 1.

In particular, the input LTS is represented as a list of records about its states—we call this representation as the *state-list* representation of the LTS. The record about

each state $s \in S$ contains a list of nonempty $pre_a(s)$ sets¹, each of them encoded as a list of its members. The partition P_{init} (and later any of its refinements) is encoded as a doubly-linked list (DLL) of blocks. Each block is represented as a DLL of (pointers to) states of the block. The relation Rel_{init} (and later any of its refinements) is encoded as a Boolean matrix $P_{init} \times P_{init}$.

4 Computing Downward Simulation

In this section, we describe algorithms for computing downward simulation on tree automata. Our approach consists of two parts: (1) we translate the maximal downward simulation problem over tree automata into a corresponding maximal simulation problem over LTSs (i.e., basically word automata), and (2) we compute the maximal word simulation on the obtained LTS using Algorithm 1. Below, we describe how the translation is carried out.

We translate the downward simulation problem on a TA $A = (Q, \Sigma, \Delta, F)$ to the simulation problem on a derived LTS A^\bullet . Each state and each left hand side of a rule in A is represented by one state in A^\bullet , while each rule in A is simulated by a set of rules in A^\bullet . Formally, we define $A^\bullet = (Q^\bullet, \Sigma^\bullet, \Delta^\bullet)$ as follows:

- The set Q^\bullet contains a state q^\bullet for each state $q \in Q$, and it also contains a state $(q_1, \dots, q_n)^\bullet$ for each $(q_1, \dots, q_n) \in Lhs(A)$.
- The set Σ^\bullet contains each symbol $a \in \Sigma$ and each index $i \in \{1, 2, \dots, n\}$ where n is the maximal rank of any symbol in Σ .
- For each transition rule $(q_1, \dots, q_n) \xrightarrow{f} q$ of A , the set Δ^\bullet contains both the transition $q^\bullet \xrightarrow{f} (q_1, \dots, q_n)^\bullet$ and transitions $(q_1, \dots, q_n)^\bullet \xrightarrow{i} q_i^\bullet$ for each $i: 1 \leq i \leq n$.
- The sets Q^\bullet , Σ^\bullet , and Δ^\bullet do not contain any other elements.

The following theorem shows correctness of the translation.

Theorem 3. *For all $q, r \in Q$, we have $q^\bullet \preceq r^\bullet$ iff $q \preceq_{down} r$.*

Due to Theorem 3, we can compute the simulation relation \preceq_{down} on Q by constructing the LTS A^\bullet and running Algorithm 1 on it with the initial partition-relation pair being simply $(P^\bullet, Rel^\bullet) = (\{Q^\bullet\}, \{(Q^\bullet, Q^\bullet)\})^2$.

4.1 Complexity of Computing the Downward Simulation

The complexity naturally consists of the price of compiling a given TA $A = (Q, \Sigma, \Delta, F)$ into its corresponding LTS A^\bullet , the price of building the initial partition-relation pair (P^\bullet, Rel^\bullet) , and the price of running Algorithm 1 on A^\bullet and (P^\bullet, Rel^\bullet) .

We assume the automata not to have unreachable states and to have at most one (final) state that is not used in the left-hand side of any transition rule—general automata

¹ We use a list rather than an array having an entry for each $a \in \mathcal{L}$ in order to avoid a need to iterate over alphabet symbols for which there is no transition.

² We initially consider all states of the LTS A^\bullet equal, and hence they form a single class of P^\bullet , which is related to itself in Rel^\bullet .

can be easily pre-processed to satisfy this requirement. Further, we assume the input automaton A to be encoded as a list of states $q \in Q$ and a list of the left-hand sides $l = (q_1, \dots, q_n) \in Lhs(A)$. Each left-hand side l is encoded by an array of (pointers to) the states q_1, \dots, q_n , plus a list containing a pointer to the so-called f -list for each $f \in \Sigma$ such that there is an f transition from l in Δ . Each f -list is then a list of (pointers to) all the states $q \in Q$ such that $l \xrightarrow{f} q$. We call this representation the *lhs-list* automata encoding. Then, the complexity of preparing the input for computing the downward simulation on A via Algorithm 1 is given by the following lemma.

Lemma 1. *For a TA $A = (Q, \Sigma, \Delta, F)$, the LTS A^\bullet and the partition-relation pair (P^\bullet, Rel^\bullet) can be derived in time and space $O(Rank(A) \cdot |Q| + |\Delta| + (Rank(A) + |\Sigma|) \cdot |Lhs(A)|)$.*

In order to instantiate the complexity of running Algorithm 1 for A^\bullet and (P^\bullet, Rel^\bullet) , we first introduce some auxiliary notions. First, we extend \preceq_{down} to the set $Lhs(A)$ such that $(q_1, \dots, q_n) \preceq_{down} (r_1, \dots, r_n)$ iff $q_i \preceq_{down} r_i$ for each $i : 1 \leq i \leq n$. We notice that $P_{sim} = Q^\bullet / \cong$. From an easy generalisation of Theorem 3 to apply not only for states from Q , but also the left-hand sides of transition rules from $Lhs(A)$, i.e., from the fact that $\forall l_1, l_2 \in Lhs(A). l_1 \preceq_{down} l_2 \Leftrightarrow l_1^\bullet \preceq l_2^\bullet$, we have that $|Q^\bullet / \cong| = |Q / \cong_{down}| + |Lhs(A) / \cong_{down}|$.

Lemma 2. *Given a tree automaton $A = (Q, \Sigma, \Delta, F)$, Algorithm 1 computes the simulation \preceq on the LTS A^\bullet for the initial partition-relation pair (P^\bullet, Rel^\bullet) with the time complexity $O((|\Sigma| + Rank(A)) \cdot |Lhs(A)| \cdot |Lhs(A) / \cong_{down}| + |\Delta| \cdot |Lhs(A) / \cong_{down}|)$ and the space complexity $O((|\Sigma| + Rank(A)) \cdot |Lhs(A)| \cdot |Lhs(A) / \cong_{down}|)$.*

The complexity of computing the downward simulation for a tree automaton A via the LTS A^\bullet can now be obtained by simply adding the complexities of computing A^\bullet and (P^\bullet, Rel^\bullet) and of running Algorithm 1 on them.

Theorem 4. *Given a tree automaton A , the downward simulation on A can be computed in time $O((|\Sigma| + Rank(A)) \cdot |Lhs(A)| \cdot |Lhs(A) / \cong_{down}| + |\Delta| \cdot |Lhs(A) / \cong_{down}|)$ and space $O((|\Sigma| + Rank(A)) \cdot |Lhs(A)| \cdot |Lhs(A) / \cong_{down}| + |\Delta|)$.³*

Moreover, under the standard assumption that the maximal rank and size of the alphabet are constants, we get the time complexity $O(|\Delta| \cdot |Lhs(A) / \cong_{down}|)$ and the space complexity $O(|Lhs(A)| \cdot |Lhs(A) / \cong_{down}| + |\Delta|)$.

5 Computing Upward Simulation

In a similar manner to the downward simulation, we translate the upward simulation problem on a tree automaton $A = (Q, \Sigma, \Delta, F)$ to the simulation problem on an LTS A^\odot . To define the translation from the upward simulation, we first make the following definition. An *environment* is a tuple of the form $((q_1, \dots, q_{i-1}, \square, q_{i+1}, \dots, q_n), f, q)$ obtained

³ Note that in the special case of $Rank(A) = 1$ (corresponding to a word automaton viewed as a tree automaton), we have $|Lhs(A)| = |Q|$, which leads to the same complexity as Algorithm 1 has when applied directly on word automata.

by removing a state q_i , $1 \leq i \leq n$, from the i^{th} position of the left hand side of a rule $((q_1, \dots, q_{i-1}, q_i, q_{i+1}, \dots, q_n), f, q)$, and by replacing it by a special symbol $\square \notin Q$ (called a *hole* below). Like for transition rules, we write $(q_1, \dots, \square, \dots, q_n) \xrightarrow{f} q$ provided $((q_1, \dots, q_{i-1}, q_i, q_{i+1}, \dots, q_n), f, q) \in \Delta$ for some $q_i \in Q$. Sometimes, we also write the environment as $(q_1, \dots, \square_i, \dots, q_n) \xrightarrow{f} q$ to emphasise that the hole is at position i . We denote the set of all environments of A by $Env(A)$.

The derivation of A° differs from A^\bullet in two aspects: (1) we encode environments (rather than left-hand sides of rules) as states in A° , and (2) we use a non-trivial initial partition on the states of A° , taking into account the downward simulation on Q . Formally, we define $A^\circ = (Q^\circ, \Sigma^\circ, \Delta^\circ)$ as follows:

- The set Q° contains a state q° for each state $q \in Q$, and it also contains a state $((q_1, \dots, \square_i, \dots, q_n) \xrightarrow{f} q)^\circ$ for each environment $(q_1, \dots, \square_i, \dots, q_n) \xrightarrow{f} q$.
- The set Σ° contains each symbol $a \in \Sigma$ and also a special symbol $\lambda \notin \Sigma$.
- For each transition rule $(q_1, \dots, q_n) \xrightarrow{f} q$ of A , the set Δ° contains both the transitions $q_i^\circ \xrightarrow{\lambda} ((q_1, \dots, \square_i, \dots, q_n) \xrightarrow{f} q)^\circ$ for each $i \in \{1, \dots, n\}$ and the transition $((q_1, \dots, \square_i, \dots, q_n) \xrightarrow{f} q)^\circ \xrightarrow{f} q^\circ$.
- The sets Q° , Σ° , and Δ° do not contain any other elements.

We define I to be the smallest binary relation on Q° containing all pairs of states of the automaton A , i.e., all pairs (q_1°, q_2°) for each $q_1, q_2 \in Q$, as well as all pairs of environments $((q_1, \dots, \square_i, \dots, q_n) \xrightarrow{f} q)^\circ, ((r_1, \dots, \square_i, \dots, r_n) \xrightarrow{f} r)^\circ$ such that $(q_j, r_j) \in D$ for each $j: 1 \leq j \neq i \leq n$.

The following theorem shows correctness of the translation.

Theorem 5. *For all $q, r \in Q$, we have $q \preceq_{up} r$ iff $q^\circ \preceq^I r^\circ$.*

The relation I is clearly a pre-order and so the relation $\iota = I \cap I^{-1}$ is an equivalence. Due to Theorem 5, we can compute the simulation relation \preceq_{up} on Q by constructing the LTS A° and running Algorithm 1 on it with the initial partition-relation pair (P°, Rel°) inducing I , i.e., $P^\circ = Q^\circ / \iota$ and $Rel^\circ = \{(B, C) \in P^\circ \times P^\circ \mid B \times C \subseteq I\}$.

5.1 Complexity of Computing the Upward Simulation

Once the downward simulation \preceq_{down} on a given TA $A = (Q, \Sigma, \Delta, F)$ is computed, the complexity of computing the simulation \preceq_{up} naturally consists of the price of compiling A into its corresponding LTS A° , the price of building the initial partition-relation pair (P°, Rel°) , and the price of running Algorithm 1 on A° and (P°, Rel°) .

We assume the automaton A to be encoded in the same way as in the case of computing the downward simulation. Compared to preparing the input for computing the downward simulation, the main obstacle in the case of the upward simulation is the need to compute the partition P_e° of the set of environments $Env(A)$ wrt. I , which is a subset of the partition P° (formally, $P_e^\circ = P^\circ \cap 2^{Env(A)}$). If the computation of P_e° is done naively (i.e., based on comparing each environment with every other environment), it can introduce a factor of $|Env(A)|^2$ into the overall complexity of the procedure. This

would dominate the complexity of computing the simulation on A° where, as we will see, $|Env(A)|$ is only multiplied by $|Env(A)/\cong_{up}|$.

Fortunately, this complexity blowup can be to a large degree avoided by exploiting the partition $Lhs(A)/\cong_{down}$ computed within deriving the downward simulation as shown in detail in [1]. Here, we give just the basic ideas.

For each $1 \leq i \leq Rank(A)$, we define an i -weakened version D_i of the downward simulation on left-hand sides of A such that $((q_1, \dots, q_n), (r_1, \dots, r_m)) \in D_i \iff n = m \geq i \wedge (\forall 1 \leq j \leq n. j \neq i \implies q_j \preceq_{down} r_j)$. Clearly, each D_i is a pre-order, and we can define the equivalence relations $\approx_i = D_i \cap D_i^{-1}$. Now, a crucial observation is that there exists a simple correspondence between P_e° and $Lhs(A)/\approx_i$. Namely, we have that $L \in Lhs(A)/\approx_i$ iff for each $f \in \Sigma$, there is a block $E_{L,f} \in P_e^\circ$ such that $E_{L,f} = \{(q_1, \dots, \square_i, \dots, q_n) \xrightarrow{f} q \mid \exists q_i, q \in Q. (q_1, \dots, q_i, \dots, q_n) \in L \wedge (q_1, \dots, q_i, \dots, q_n) \xrightarrow{f} q\}$.

The idea of computing P_e° is now to first compute blocks of $Lhs(A)/\approx_i$ and then to derive from them the P_e° blocks. The key advantage here is that the computation of the \approx_i -blocks can be done on blocks of $Lhs(A)/\cong_{down}$ instead of directly on elements of $Lhs(A)$. This is because, for each i , blocks of $Lhs(A)/\cong_{down}$ are sub-blocks of blocks of $Lhs(A)/\approx_i$. Moreover, for any blocks K, L of $Lhs(A)/\cong_{down}$, the test of $K \times L \subseteq D_i$ can simply be done by testing whether $(k, l) \in D_i$ for any two representatives $k \in K, l \in L$. Therefore, all \approx_i -blocks can be computed in time proportional to $|Lhs(A)/\cong_{down}|^2$.

From each block $L \in Lhs(A)/\approx_i$, one block $E_{L,f}$ of P_e° is generated for each symbol $f \in \Sigma$. The $E_{L,f}$ blocks are obtained in such a way that for each left-hand side $l \in L$, we generate all the environments which arise by replacing the i^{th} state of l by \square , adding f , and adding a right-hand side state $q \in Q$ which together with l form a transition $l \xrightarrow{f} q$ of A . This can be done efficiently using the lhs-list encoding of A . An additional factor $|\Delta| \cdot \log |Env(A)|$ is, however, introduced due to a need of not having duplicates among the computed environments, which could result from transitions that differ just in the states that are replaced by \square when constructing an environment. The factor $\log |Env(A)|$ comes from testing a set membership over the computed environments to check whether we have already computed them before or not.

Moreover, it can be shown that Rel° can be computed in time $|P^\circ|^2$. The complexity of constructing A° and (P°, Rel°) is then summarised in the below lemma.

Lemma 3. *Given a tree automaton $A = (Q, \Sigma, \Delta, F)$, the downward simulation \preceq_{down} on A , and the partition $Lhs(A)/\cong_{down}$, the LTS A° and the partition-relation pair (P°, Rel°) can be derived in time $O(|\Sigma| \cdot |Q| + Rank(A) \cdot (|Lhs(A)| + |Lhs(A)/\cong_{down}|^2) + Rank(A)^2 \cdot |\Delta| \cdot \log |Env(A)| + |P^\circ|^2)$ and in space $O(|\Sigma| \cdot |Q| + |Env(A)| + Rank(A) \cdot |Lhs(A)| + |Lhs(A)/\cong_{down}|^2 + |P^\circ|^2)$.*

In order to instantiate the complexity of running Algorithm 1 for A° and (P°, Rel°) , we again first introduce some auxiliary notions. Namely, we extend \preceq_{up} to the set $Env(A)$ such that $(q_1, \dots, \square_i, \dots, q_n) \xrightarrow{f} q \preceq_{up} (r_1, \dots, \square_j, \dots, r_m) \xrightarrow{f} r \iff m = n \wedge i = j \wedge q \preceq_{up} r \wedge (\forall k \in \{1, \dots, n\}. k \neq i \implies q_k \preceq_{down} r_k)$. We notice that $P_{sim} = Q^\circ / \cong^I$. From an easy generalisation of Theorem 5 to apply not only for states from Q , but also environments from $Env(A)$, i.e., from the fact that $\forall e_1, e_2 \in Env(A). e_1 \preceq_{up} e_2 \iff e_1^\circ \preceq^I e_2^\circ$, we have that $|Q^\circ / \cong^I| = |Q / \cong_{up}| + |Lhs(A) / \cong_{up}|$.

Lemma 4. *Given a tree automaton $A = (Q, \Sigma, \Delta, F)$, the upward simulation \preceq_{up} on A can be computed by running Algorithm 1 on the LTS A° and the partition-relation pair (P°, Rel°) in time $O(Rank(A) \cdot |\Delta| \cdot |Env(A)/\cong_{up}| + |\Sigma| \cdot |Env(A)| \cdot |Env(A)/\cong_{up}|)$ and space $O(|\Sigma| \cdot |Env(A)| \cdot |Env(A)/\cong_{up}|)$.*

The complexity of computing upward simulation on a TA A can now be obtained by simply adding the price of computing downward simulation, the price of computing A° and (P°, Rel°) , and the price of running Algorithm 1 on A° and (P°, Rel°) .

Theorem 6. *Given a tree automaton $A = (Q, \Sigma, \Delta, F)$, let $T_{down}(A)$ and $S_{down}(A)$ denote the time and space complexity of computing the downward simulation \preceq_{down} on A . Then, the upward simulation \preceq_{up} on A can be computed in time*

$$O((|\Sigma| \cdot |Env(A)| + Rank(A) \cdot |\Delta|) \cdot |Env(A)/\cong_{up}| + Rank(A)^2 \cdot |\Delta| \cdot \log |Env(A)| + T_{down}(A))$$

and in space $O(|\Sigma| \cdot |Env(A)| \cdot |Env(A)/\cong_{up}| + S_{down}(A))$.⁴

Finally, from the standard assumption that the maximal rank and the alphabet size are constants and from observing that $|Env(A)| \leq Rank(A) \cdot |\Delta| \leq Rank(A) \cdot |\Sigma| \cdot |Q|^{Rank(A)+1}$, we get the time complexity $O(|\Delta| \cdot (|Env(A)/\cong_{up}| + \log |Q|) + T_{down}(A))$ and the space complexity $O(|Env(A)| \cdot |Env(A)/\cong_{up}| + S_{down}(A))$.

6 Reducing Tree Automata

In this section, we describe how to reduce tree automata while preserving the language of the automaton. The idea is to identify suitable equivalence relations on states of tree automata, and then collapse the sets of states which form equivalence classes. We will consider two reduction methods: one which uses downward simulation, and one which is defined in terms of both downward and upward simulation. The choice of the equivalence relation is a trade-off between the amount of reduction achieved and the cost of computing the relation. The second mentioned equivalence is heavier to compute as it requires that both downward and upward simulation are computed and then suitably composed. However, it is at least as coarse as—and often significantly coarser than—the downward simulation equivalence, and hence can give much better reductions as witnessed even in our experiments.

Consider a tree automaton $A = (Q, \Sigma, \Delta, F)$ and an equivalence relation \equiv on Q . The *abstract tree automaton* derived from A and \equiv is $A(\equiv) = (Q(\equiv), \Sigma, \Delta(\equiv), F(\equiv))$ where:

- $Q(\equiv)$ is the set of blocks in \equiv . In other words, we collapse all states which belong to the same block into one abstract state.
- $(B_1, \dots, B_n) \xrightarrow{f} B$ iff $(q_1, \dots, q_n) \xrightarrow{f} q$ for some $q_1 \in B_1, \dots, q_n \in B_n, q \in B$. This is, there is a transition in the abstract automaton iff there is a transition between states in the corresponding blocks.
- $F(\equiv)$ contains a block B iff $B \cap F \neq \emptyset$. Intuitively, a block is accepting if it contains at least one state which is accepting.

⁴ Note that in the special case of $Rank(A) = 1$ (corresponding to a word automaton viewed as a tree automaton), we have $|Env(A)| \leq |\Sigma| \cdot |Q|$, which leads to almost the same complexity (up to the logarithmic component) as Algorithm 1 has when applied directly on word automata.

6.1 Downward Simulation Equivalence

Given a tree automaton $A = (Q, \Sigma, \Delta, F)$, we consider the abstract automaton $A \langle \cong_{\text{down}} \rangle$ constructed by collapsing states of A which are equivalent with respect to \cong_{down} . We show that the two automata accept the same language, i.e., $L(A) = L(A \langle \cong_{\text{down}} \rangle)$. Observe that the inclusion $L(A) \subseteq L(A \langle \cong_{\text{down}} \rangle)$ is straightforward. We can prove the inclusion in the other direction as follows. Using a simple induction on trees, one can show that downward simulation implies language inclusion. In other words, for states $q, r \in Q$, if $q \preceq_{\text{down}} r$, then $L(q) \subseteq L(r)$. This implies that for any $B \in Q \langle \cong_{\text{down}} \rangle$, it is the case that $L(B) \subseteq L(r)$ for any $r \in B$. Now suppose that $t \in L(A \langle \cong_{\text{down}} \rangle)$. It follows that $t \in L(B)$ for some $B \in F \langle \cong_{\text{down}} \rangle$. Since $B \in F \langle \cong_{\text{down}} \rangle$, there is some $r \in B$ with $r \in F$. It follows that $t \in L(r)$, and hence $t \in L(A)$. This gives the following Theorem.

Theorem 7. $L(A) = L(A \langle \cong_{\text{down}} \rangle)$ for each tree automaton A .

In fact, $A \langle \cong_{\text{down}} \rangle$ is the minimal automaton which is equivalent to A with respect to downward simulation and which accepts the same language as A .

6.2 Composed Equivalence

Consider a tree automaton $A = (Q, \Sigma, \Delta, F)$. Let I_F be a partitioning of Q such that $(q, r) \in I_F$ iff $q \in F \implies r \in F$. Consider a reflexive and transitive downward simulation D , and a reflexive and transitive upward simulation U induced by D . Assume that $U \subseteq I_F$. We will reduce A with respect to relations of the form \equiv_R which preserve language equivalence, but which may be much coarser than downward simulations. Here, each \equiv_R is an equivalence relation $R \cap R^{-1}$ defined by a pre-order R satisfying certain properties. More precisely, we use $D \oplus U$ to denote the set of relations on Q such that for each $R \in (D \oplus U)$, the relation R satisfies the following two properties: (i) R is transitive and (ii) $D \subseteq R \subseteq (D \circ U^{-1})$. For a state $r \in Q$ and a set $B \subseteq Q$ of states, we write $(B, r) \in D$ to denote that there is a $q \in B$ with $(q, r) \in D$. We define $(B, r) \in U$ analogously. We will now consider the abstract automaton $A \langle \equiv_R \rangle$ where the states of A are collapsed according to \equiv_R . We will relate the languages of A and $A \langle \equiv_R \rangle$.

To do that, we first define the notion of a *context*. Intuitively, a context is a tree with “holes” instead of leaves. Formally, we consider a special symbol $\circ \notin \Sigma$ with rank 0. A *context* over Σ is a tree c over $\Sigma \cup \{\circ\}$ such that for all leaves $p \in c$, we have $c(p) = \circ$. For a context c with leaves p_1, \dots, p_n , and trees t_1, \dots, t_n , we define $c[t_1, \dots, t_n]$ to be the tree t , where

- $\text{dom}(t) = \text{dom}(c) \cup \{p_1 \cdot p' \mid p' \in \text{dom}(t_1)\} \cup \dots \cup \{p_n \cdot p' \mid p' \in \text{dom}(t_n)\}$,
- for each $p = p_i \cdot p'$, we have $t(p) = t_i(p')$, and
- for each $p \in \text{dom}(c) \setminus \{p_1, \dots, p_n\}$, we have $t(p) = c(p)$.

In other words, $c[t_1, \dots, t_n]$ is the result of appending the trees t_1, \dots, t_k to the holes of c . We extend the notion of runs to contexts. Let c be a context with leaves p_1, \dots, p_n . A *run* π of A on c from (q_1, \dots, q_n) is defined in a similar manner to a run on a tree except that for a leaf p_i , we have $\pi(p_i) = q_i$, $1 \leq i \leq n$. In other words, each leaf labelled with \circ is annotated by one q_i . We use $c[q_1, \dots, q_n] \xrightarrow{\pi} q$ to denote that π is a run of A on c from (q_1, \dots, q_n) such that $\pi(\varepsilon) = q$. The notation $c[q_1, \dots, q_n] \implies q$ is explained in a similar manner to runs on trees.

Using the notion of a context, we can relate runs of A with those of the abstract automaton $A\langle\equiv_R\rangle$. More precisely, we can show that for blocks $B_1, \dots, B_n, B \in Q\langle\equiv_R\rangle$ and a context c , if $c[B_1, \dots, B_n] \Longrightarrow B$, then there exist states $r_1, \dots, r_n, r \in Q$ such that $(B_1, r_1) \in D, \dots, (B_n, r_n) \in D, (B, r) \in U$, and $c[r_1, \dots, r_n] \Longrightarrow r$. In other words, each run in $A\langle\equiv_R\rangle$ can be simulated by a run in A which starts from larger states (with respect to downward simulation) and which ends up at a larger state (with respect to upward simulation). This leads to the following lemma.

Lemma 5. *If $t \Longrightarrow B$, then $t \Longrightarrow w$ for some w with $(B, w) \in U$. Moreover, if $B \in F\langle\equiv_R\rangle$, then also $w \in F$.*

In other words, each tree t which leads to a block B in $A\langle\equiv_R\rangle$ will also lead to a state in A which is larger than (some state in) the block B with respect to upward simulation. Moreover, if t can be accepted at B in $A\langle\equiv_R\rangle$ (meaning that B contains a final state of A , i.e., $B \cap F \neq \emptyset$), then it can be accepted at w in A (i.e., $w \in F$) too.

Notice that Lemma 5 holds for any downward and upward simulations satisfying the properties mentioned in the definition of \oplus . We now instantiate the lemma for the maximal downward and upward simulation to obtain the main result. We take D and U to be \preceq_{down} and $\preceq_{up}^{I_F}$, respectively, and we let \preceq_{comp} be any relation from the set of relations $(\preceq_{down} \oplus \preceq_{up}^{I_F})$. We let \cong_{comp} be the corresponding equivalence.

Theorem 8. $L(A\langle\cong_{comp}\rangle) = L(A)$ for each tree automaton A .

Proof. The inclusion $L(A\langle\cong_{comp}\rangle) \supseteq L(A)$ is trivial. Let $t \in L(A\langle\cong_{comp}\rangle)$, i.e., $t \Longrightarrow B$ for some block B where $B \cap F \neq \emptyset$. Lemma 5 implies that $t \Longrightarrow w$ such that $w \in F$. \square

Note that it is clearly the case that $\cong_{down} \subseteq \cong_{comp}$. Moreover, note that a relation $\preceq_{comp} \in (\preceq_{down} \oplus \preceq_{up}^{I_F})$ can be obtained, e.g., by a simple (random) pruning of the relation $\preceq_{down} \circ (\preceq_{up}^{I_F})^{-1}$ based on iteratively removing links not being in \preceq_{down} and at the same time breaking transitivity of the so-far computed composed relation. Such a way of computing \preceq_{comp} does not guarantee that one obtains a relation of the greatest cardinality possible among relations from $\preceq_{down} \oplus \preceq_{up}^{I_F}$, but, on the other hand, it is cheap (in the worst case, cubic in the number of states). Moreover, our experiments show that even this simple way of computing the composed relation can give us a relation \cong_{comp} that is much coarser (and yields significantly better reductions) than \cong_{down} .

Remark. Our definition of a context coincides with the one of [8] where all leaves are holes. On the other hand, a context in [9] and [3] is a tree with a *single* hole. Considering single-hole contexts, one can define the *language of contexts* $L_c(q)$ of a state q to be the set of contexts on which there is an accepting run if the hole is replaced by q . Then, for all states q and r , it is the case that $q \preceq_{up} r$ implies $L_c(q) \subseteq L_c(r)$.

7 Experiments with Reducing Tree Automata

We have implemented our algorithms in a prototype tool written in Java. We have run the prototype on a number of tree automata that arise in the framework of *tree regular model checking*. Tree regular model checking is the name of a family of techniques for analysing infinite-state systems in which states are represented by trees, (infinite) sets of states by

Table 1. Reduction of the number of states and rules using different reduction algorithms

Protocol	original		\cong_{down}		\cong_{comp}		backward bisimulation	
	states	rules	states	rules	states	rules	states	rules
percolate	10	72	7	45	7	45	10	72
	20	578	17	392	14	346	20	578
	28	862	13	272	13	272	15	341
arbiter	15	324	10	248	7	188	11	252
	41	313	28	273	19	220	33	285
	109	1248	67	1048	55	950	83	1116
leader	17	153	11	115	6	47	16	152
	25	384	16	235	6	59	23	382
	33	876	10	100	7	67	27	754

finite tree automata, and transitions by tree transducers. Most of the algorithms in the framework rely crucially on efficient automata reduction methods since the size of the generated automata often explodes, making computations infeasible without reduction. The (nondeterministic) tree automata that we have considered arose during verification of the *Percolate* protocol, the *Arbiter* protocol, and the *Leader* election protocol [4].

Our experimental evaluation was carried out on an AMD Athlon 64 X2 2.19GHz PC with 2.0 GB RAM. The time for minimising the tree automata varied from a few seconds up to few minutes. Table 1 shows the number of states and rules of the various considered tree automata before and after computing \cong_{down} , \cong_{comp} , and the backward bisimulation from [11]. Backward bisimulation is the bisimulation counterpart of downward simulation. The composed simulation equivalence \cong_{comp} was computed in the simple way based on the random pruning of the relation $\preceq_{down} \circ (\preceq_{up}^{I_F})^{-1}$ as mentioned at the end of Section 6.2. As Table 1 shows, \cong_{comp} achieves the best reduction (often reducing to less than one-third of the size of the original automaton). As expected, both \cong_{down} and \cong_{comp} give better reductions than backward bisimulation in all test cases.

8 Conclusions and Future Work

We have presented methods for reducing tree automata under language equivalence. For this purpose, we have considered two kinds of simulation relations on the states of tree automata, namely downward and upward simulation. We give procedures for efficient translation of both kinds of relations into simulations defined on labelled transition systems. Furthermore, we define a new, language-preserving equivalence on tree automata, derived from compositions of downward and upward simulation, which (according to our experiments) usually gives a much better reduction on the size of automata than downward simulation.

There are several interesting directions for future work. First, we would like to implement the proposed algorithms in a more efficient way, perhaps over automata encoded in a symbolic way using BDDs like in MONA [12], in order to be able to experiment with bigger automata. Further, for instance, we can define *upward* and *downward bisimulation* for tree automata in an analogous way to the case of simulation. It is straightforward to show that the encoding we use in this paper can also be used to translate

bisimulation problems on tree automata into corresponding ones for LTSs. Although reducing according to a bisimulation does not give the same reduction as for a simulation, it is relevant since it generates more efficient algorithms. Also, we plan to investigate coarser relations for better reductions of tree automata by refining the ideas behind the definition of the composed relation introduced in Section 6. We believe that it is possible to define a refinement scheme allowing one to define an increasing family of such relations between downward simulation equivalence and tree language equivalence. Finally, we plan to consider extending our reduction techniques to the class of unranked trees which are used in applications such as reasoning about structured documents or about configurations of dynamic concurrent processes.

Acknowledgement. The work was supported by the ANR-06-SETI-001 French project AVERISS, the Czech Grant Agency (projects 102/07/0322 and 102/05/H050), the Czech-French Barrande project 2-06-27, and the Czech Ministry of Education by the project MSM 0021630528 *Security-Oriented Research in Information Technology*.

References

1. Abdulla, P., Bouajjani, A., Holík, L., Kaati, L., Vojnar, T.: Computing Simulations over Tree Automata. Technical report, FIT-TR-2007-001, FIT, Brno University of Technology, Czech Republic (2007)
2. Abdulla, P., Högberg, J., Kaati, L.: Bisimulation Minimization of Tree Automata. In: H. Ibarra, O., Yen, H.-C. (eds.) CIAA 2006. LNCS, vol. 4094, Springer, Heidelberg (2006)
3. Abdulla, P., Jonsson, B., Mahata, P., d'Orso, J.: Regular Tree Model Checking. In: Brinksma, E., Larsen, K.G. (eds.) CAV 2002. LNCS, vol. 2404, Springer, Heidelberg (2002)
4. Abdulla, P., Legay, A., d'Orso, J., Rezine, A.: Tree Regular Model Checking: A Simulation-based Approach. The Journal of Logic and Algebraic Programming 69(1-2), 93–121 (2006)
5. Bouajjani, A., Habermehl, P., Rogalewicz, A., Vojnar, T.: Abstract Regular Tree Model Checking. In: ENTCS, vol. 149(1), pp. 37–48. Elsevier, Amsterdam (2006)
6. Bouajjani, A., Habermehl, P., Rogalewicz, A., Vojnar, T.: Abstract Regular Tree Model Checking of Complex Dynamic Data Structures. In: Yi, K. (ed.) SAS 2006. LNCS, vol. 4134, Springer, Heidelberg (2006)
7. Bouajjani, A., Touili, T.: Extrapolating Tree Transformations. In: Brinksma, E., Larsen, K.G. (eds.) CAV 2002. LNCS, vol. 2404, Springer, Heidelberg (2002)
8. Bouajjani, A., Touili, T.: Reachability Analysis of Process Rewrite Systems. In: Pandya, P.K., Radhakrishnan, J. (eds.) FSTTCS 2003. LNCS, vol. 2914, Springer, Heidelberg (2003)
9. Comon, H., Dauchet, M., Gilleron, R., Jacquemard, F., Lugiez, D., Tison, S., Tommasi, M.: Tree Automata Techniques and Applications (1997), Available on: <http://www.grappa.univ-lille3.fr/tata>
10. Henzinger, M., Henzinger, T., Kopke, P.: Computing simulations on finite and infinite graphs. In: Proc. of FOCS 1995, IEEE, Los Alamitos (1995)
11. Maletti, A., Högberg, J., May, J.: Backward and forward bisimulation minimisation of tree automata. In: Holub, J., Žďárek, J. (eds.) CIAA 2007. LNCS, vol. 4783, pp. 109–121. Springer, Heidelberg (2007)
12. Klarlund, N., Møller, A.: MONA Version 1.4 User Manual, BRICS, Department of Computer Science, University of Aarhus, Denmark (2001)
13. Paige, R., Tarjan, R.: Three Partition Refinement Algorithms. SIAM Journal on Computing 16, 973–989 (1987)
14. Ranzato, F., Tapparo, F.: A New Efficient Simulation Equivalence Algorithm. In: Proc. of LICS 2007, IEEE CS, Los Alamitos (2007)