

Improving the Representation of Infinite Trees to Deal with Sets of Trees

Laurent Mauborgne

LIENS – DMI, École Normale Supérieure, 45 rue d’Ulm, 75 230 Paris cedex 05, France
Tel: +33 (0) 1 44 32 20 66; Email: Laurent.Mauborgne@ens.fr
WWW home page: <http://www.dmi.ens.fr/~mauborgn/>

Abstract. In order to deal efficiently with infinite regular trees (or other pointed graph structures), we give new algorithms to store such structures. The trees are stored in such a way that their representation is unique and shares as much as possible. This maximal sharing allows substantial memory gain and speed up. For example, equality testing becomes constant time. The algorithms are incremental, and as such allow good reactive behavior. These new algorithms are then applied to the representation of sets of trees. The expressive power of this new representation is exactly what is needed by set-based analysis.

1 Introduction

When applying set-based analysis techniques for practical applications, one is surprised to see that the representation of the sets of trees is not very efficient. Even when we use tree automata, we cannot overcome this problem without performing a minimization of the whole automaton at each step. We propose a new way of dealing with this kind of structure to get a representation that is as small as possible during the computation.

After analysis of the problem, it appears that the underlying structure we want to optimize can be described mathematically as regular infinite trees. Because tree structures appear everywhere in computer science where a hierarchy occurs, we found it interesting to present the algorithms in an independent way. In this way, our technique appears as an extension of an efficient solution to store finite trees.

The representation we extend uses just the minimum amount of memory by sharing equivalent subtrees. This saves a lot of space. It is used, for example, with sets of words represented as a tree to share common prefixes. It is possible to share the subtrees incrementally, and at the same time to give a unique representation to different versions of the same trees. Such a technique allows constant time equality testing and a great speed up for many other algorithms manipulating trees. It has been the source of the success of Binary Decision Diagrams (BDDs) [2], which are considered one of the best representations for boolean functions so far.

But as soon as a loop occurs somewhere in the data, finite tree techniques are no longer adequate. The main contribution of this article is to extend the good

results of unique sharing representation from finite trees to infinite trees. These techniques are applied to the representation of sets of trees in set-based analysis, but they can also be applied directly to the representation and manipulation of finite automata, or infinite boolean functions [14].

After a recollection of the classic results over finite trees in section 2, we present the solutions for the most difficult problems with infinite trees in the section 3 on cycles. The general problem is then treated in section 4, with a full example. Complexity issues and algorithms to manipulate infinite trees are discussed in section 5. The application to sets of trees implies the description of a new encoding to keep the uniqueness of the representation. This new contribution is described in section 6.

2 Classic Representation of Trees

2.1 Trees and Graphs

As we deal with the computer representation of data structures, we must give a clear meaning to the word representation, and in particular clearly distinguish between what is represented and what is the representation. For this reason, we will give a mathematical definition of what is a tree, and another one for the way it is usually stored in a computer.

Let \mathbb{N}^* be the set of words over \mathbb{N} , ε denoting the empty word. We note \prec the prefix ordering on words and $u.v$ the concatenation of the words u and v . Let F be a finite set of labels.

Definition 1. A tree t labeled by F is a function of $pos(t) \rightarrow F$ such that $pos(t) \subset \mathbb{N}^*$ and $\forall p \in \mathbb{N}^*, \forall i \in \mathbb{N}, p.i \in pos(t) \Rightarrow (p \in pos(t) \text{ and } \forall j < i, p.j \in pos(t))$

Let $p \in pos(t)$. The subtree of t in p , written $t_{[p]}$ is defined by: $pos(t_{[p]}) \stackrel{\text{def}}{=} \{q \in \mathbb{N}^* \mid p.q \in pos(t)\}$, and $t_{[p]}(q) \stackrel{\text{def}}{=} t(p.q)$. A tree is uniquely determined by the label of its root, $t(\varepsilon)$, and by the children of the root, the different $t_{[i]}, i \in \mathbb{N}$.

In the sequel, a generic tree will be denoted $\begin{matrix} f \\ \swarrow \searrow \\ t_0 \dots t_{n-1} \end{matrix}$, where f is the label of the root, and $(t_i)_{i < n}$ are the children of the root.

When representing a tree in a computer, we usually use one computer location for each position p in $pos(t)$, where we store the label $t(p)$ and the location of the different children (the $p.i$'s in $pos(p)$) of this position. Such a representation is well modeled by a graph, where each node of the graph corresponds to a computer location. We do not give the most general definition of graphs, but the definition that is useful in this article to represent trees.

Definition 2. A graph G labeled by F is composed of two sets, the node set, G^N , and the edge set, $G^E \subset G^N \times G^N \times \mathbb{N}$, and every node of the graph is associated with a label in F .

We define the notion of path in a graph: let $p \in \mathbb{N}^*$, p is a *path* of the node N if and only if $p = \varepsilon$ or $p = i.q$ and there is an $M \in G^N$ such that $(N, M, i) \in G^E$ and q is a path of M . If O is the only node at the end of the path, we write $N.p = O$. We define $\mathcal{G}(N)$ as the graph defined by the modes which can be reached from N . We will often identify a node N and the graph $\mathcal{G}(N)$.

Definition 3. A node N represents a tree t if and only if the set of paths of N is $pos(t)$, and $\forall p \in pos(t)$, $N.p$ is well defined, and its label is $t(p)$.

A *finite tree* t is a tree such that $pos(t)$ is finite. There is always a possible representation by a finite graph for finite trees. In the most common use, one node corresponds to each path of the finite tree.

A *regular tree* t is a tree such that the number of distinct subtrees of t is finite. Such a tree can be infinite, but it can still be represented by a finite graph [6], see Fig. 1 for an example.

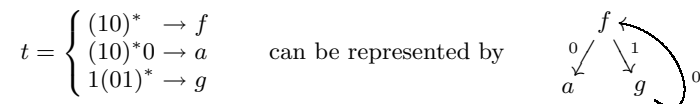


Fig. 1. An infinite regular tree

2.2 Best Representation

The naive representation, which consists in using any graph representing the tree [6], is very easy to deal with and quite widely used for small problems. But we can do far better if we observe that some nodes can represent different paths of the tree, as long as the subtrees at these paths are the same. This is called sharing the subtrees (see e.g. [1]). In fact, the best we can do is to have exactly one node for each distinct subtree. This is what we call the best representation of a tree. In the case of finite trees, this can save a lot of space, and even time by memoizing [15], and in the case of infinite regular trees, we avoid the possibility of unbounded representation for a given tree.

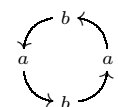
When dealing with many trees, we can do even better: considering the entire computer memory as one graph, we can optimize the representation for all the trees, and have in effect exactly one memory location for each distinct tree we need to store. An immediate consequence is that we just have to compare the location of the roots (the node representing the trees) to compare entire trees. Such a technique is used e.g. in BDDs [2] to achieve impressive speed up and memory gain.

The technique to obtain the best representation of the trees uses a dictionary mechanism linking keys to nodes of the graph, usually a hash table. The keys

are built incrementally: if the keys for the $(t_i)_{i < n}$ are known and linked to the nodes $(N_i)_{i < n}$, then the key for $\underset{t_0 \dots t_{n-1}}{\vee}^f$ is $(f, (N_i)_{i < n})$. Each time a key is not present in the dictionary, it is associated with a new node N , with edges to the N_i 's. If we come to a tree whose key is already in the dictionary, we use the corresponding node. As the trees are always built from leaves to root, we have indeed a best representation for the trees.

3 Dealing with Cycles

When representing infinite trees, though, we cannot go from the leaves to the root, so we cannot start the key mechanism which leads to the best representation. The difficulty lies in the infinite paths of the tree, that is the cycles of the graph representing the tree. Whereas in finite trees there is no need to see beyond the immediate children of a given node, when dealing with cycles, we can have reasons to look further, in order to detect the two causes of cycle unfolding: cycle growth and root unfolding. For example, consider the cycle $a \xrightarrow{b} b$.

 is an example of cycle growth, and $a \rightarrow b \xrightarrow{a} a$ is an example of root

unfolding. In this very simple example, it is easy to reduce root unfolding by looking at the key of the root, but it is much more difficult if the root itself is still in another cycle. In order to concentrate on the real difficulties, we suppose in this section that we deal with strongly connected graphs, that is graphs such that there is a path between any pair of nodes.

3.1 Cycle Growth and Tree Keys

We give \equiv_{tree} as the equivalence between nodes representing the same tree. The goal of cycle growth reduction is to find an equivalent graph with the minimum number of nodes. In such a graph, whatever the nodes N and M , $N \equiv_{\text{tree}} M \Rightarrow N = M$. Such a problem is called a partitioning problem. It has been solved in time $n \log(n)$ by Hopcroft [10] for finite automata, and in the general case by [4]. We call **share**(N) the algorithm that takes a node N and modifies the associated graph so that it has the fewest possible nodes (Fig. 2).

Cycle growth reduction corresponds to the state of the art in automata representation. But we want to go further: we need that the representation be unique whatever the different versions of the same tree. To perform this, we give a key which distinguishes between non isomorphic graphs. This key is associated to a given node N of the graph. It is a finite tree which corresponds to the graph as long as we do not loop, but as soon as we loop, the label of the node is replaced by its access path from N . It is described as **treeKey**(N). See Fig. 3 for an example. The isomorphism between graphs is not the same thing as \equiv_{tree} . In general it can differentiate two graphs which represent the same tree. The

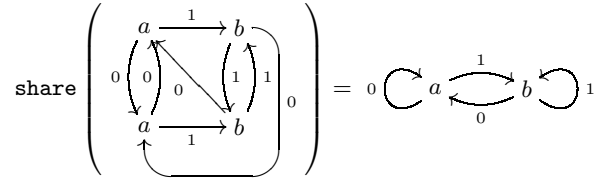


Fig. 2. Application of the **share** algorithm.

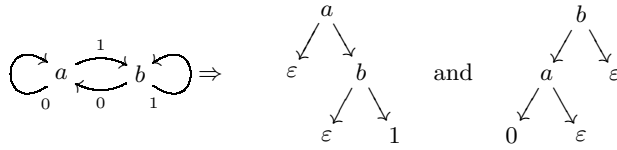


Fig. 3. A graph, followed by the tree keys of its two nodes

interesting point is that it is indeed the same relation on graphs with a minimal number of nodes.

Proposition 1. *Whatever M and N , such that $\mathcal{G}(M)$ and $\mathcal{G}(N)$ are graphs with minimal number of nodes, $\mathbf{treeKey}(M) = \mathbf{treeKey}(N) \Leftrightarrow M \equiv_{\text{tree}} N$.*

Proof. The difficult point is $M \equiv_{\text{tree}} N \Rightarrow \mathbf{treeKey}(M) = \mathbf{treeKey}(N)$. Suppose there are M and N such that $\mathcal{G}(M)$ and $\mathcal{G}(N)$ are graphs with minimal number of nodes, $M \equiv_{\text{tree}} N$ and $\mathbf{treeKey}(M) \neq \mathbf{treeKey}(N)$. Let $t_M = \mathbf{treeKey}(M)$ and $t_N = \mathbf{treeKey}(N)$. Because $t_M \neq t_N$, there is a path p such that $t_M(p) \neq t_N(p)$. But if $t_M(p)$ is a label of the graph, $t_M(p)$ is the label of $M.p$, and the same holds for N . Because $M \equiv_{\text{tree}} N$, $M.p$ and $N.p$ have the same label, so at least one of $t_M(p)$ or $t_N(p)$ is not a label of the graphs (and so is in \mathbb{N}^*), say $t_M(p)$. It means there is a $q \prec p$ such that $M.q \equiv_{\text{tree}} M.p$. So $N.q \equiv_{\text{tree}} N.p$, but by minimality of the number of nodes of $\mathcal{G}(N)$, $N.q$ and $N.p$ must be the same node, and so $t_N(p) = q = t_M(p)$. \square

Because we can find an equivalent graph with minimal number of nodes for strongly connected graphs, we have a valid key mechanism for any strongly connected graph: we first apply **share**, then **treeKey**.

3.2 Root Unfolding and Partial Keys

With just **share** and **treeKey** (applied to every node), we can have a unique representation that shares common subtrees. But as we need to start the whole process from the beginning for each little modification in the trees, such a process would be quite slow. Moreover, it is much better to apply the **share** algorithm on

the smallest possible graphs. As it is not a linear algorithm, we have better results if we can split the graph and apply the algorithm to each separate subgraph only.

The finite parts of the tree can always be treated in the classic way, while the loops will need a special treatment. In order to decompose the graph and mark those parts of the graph which have been definitely treated, we introduce partial keys. A partial key looks like a node key for a finite tree, a label followed by a vector of nodes, except that for some parts of the vector, there is no node (see Sect 4.3 for an example). A partial key k has a name: $\text{name}(k) \in F$ and is a partial function from \mathbb{N} to nodes. A graph labeled by partial keys is such that for every node N in the graph, if k is the partial key for N , the edges in the graph correspond to those integers for which the partial key is not defined. For example, if a node is labeled by f of arity 3, we can have a partial key which is not defined on 0 and 1 (we write a \bullet), and on 2 its value is the node number 4. We write $(f, \bullet \bullet \square_4)$ for this partial key. The only edges that can leave from such a node would be labeled by 0 and 1. The idea is that what is in the partial keys is uniquely represented. In our example, the node number 4, \square_4 , is a unique representation of some tree. Later on during the computation, it is possible that we have a unique representation for the first component, say with node \square_2 , and the partial key becomes $(f, \square_2 \bullet \square_4)$. When a partial key is full (defined everywhere), then the node should be a unique representation.

This new graphs have new equivalence relation, \equiv_{pk} which is implied by \equiv_{tree} . This new equivalence relation corresponds to \equiv_{tree} after the expansion of the partial keys into the graph.

But now, with those partial keys, we can have a strongly connected graph such that, by root unfolding, one of its nodes is equivalent to a node in a partial key. Figure 4 shows a case of root unfolding, which can be as big as we want, even after cycle growth reduction¹. So, we must look for such a node, even before

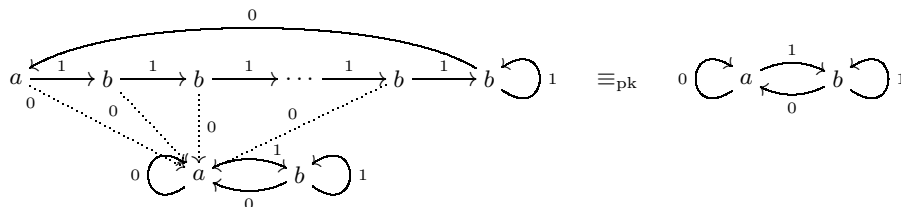


Fig. 4. Root unfolding of a cycle

applying the **share** algorithm.

The name of the algorithm performing this task is **shareWithDone**(N). It returns N if and only if no other node in the partial keys is equivalent to N .

¹ In this figure, dotted lines correspond to nodes stored in partial keys.

Otherwise, it returns the node in the partial keys that is equivalent to N . This algorithm uses some properties of the graph to reduce the complexity of the computation. Let G be the graph associated with N . As always in this section, we suppose that G is strongly connected. We call H the graph already computed and that is reachable from the partial keys of G . The algorithm determines whether a node of G is equivalent to a node of H . If it is the case, then there is root unfolding. If not, there is no root unfolding. We show that it is enough to verify this property for one node to treat the entire graph G because G is strongly connected. Suppose N is equivalent to M in H . Then, whatever the legible path p , $N.p$ is equivalent to $M.p$. Because H has been treated already, any $M.p$ is in H , and because G is strongly connected, any node of G is a $N.p$.

There is a kind of reciprocal property that is exploited too: for some subsets of H^N , if no node of the subset is equivalent to a particular node of G , then they are not equivalent to any node of G . A subset of H^N is said to be closed if and only if, for every legible path p , for every node N in the subset, $N.p$ is in the subset.

Proposition 2. $\forall H' \subset H^N$ such that H' is closed, if $\exists N \in G^N$ such that $\forall M \in H', N \not\equiv_{pk} M$, then this holds for every $N \in G^N$.

Proof. Let H' be such a subset and N a node of G . If N is not equivalent to any node in H' , then, suppose there is a $M \in G^N$ and a $O \in H'$ such that M is equivalent to O . As G is strongly connected, there is a p such that $M.p = N$. So, N would be equivalent to $O.p$, which is in H' . This proves that no element of G^N is equivalent to any element of H' . \square

Because of these properties, we can use the following algorithm for **shareWithDone**: we just compare every nodes of G with the nodes that are reachable from their partial keys and not already encountered. This comparison can be quite efficient by exploiting the fact that the nodes in the partial keys are unique representations of trees, although we have a quadratic worst case complexity.

We will show in the next section, that by applying first **shareWithDone**, then **share** and then **treeKey**, we can indeed represent uniquely (and with the least possible number of nodes) any strongly connected graph, in an incremental process.

4 The Best Representation for Infinite Trees

4.1 Informal Presentation

In order to show how we can produce the best representation for an infinite tree, we solve the following problem: considering a graph representing a tree t , return an equivalent graph with a minimal number of nodes. To achieve this in an incremental way, we use two dictionary mechanisms and a decomposition of the graph. First, we apply the classic algorithm, using the dictionary D , on the finite subtrees of the tree. When a finite subtree is entirely treated, it is

incorporated in the graph through partial keys. Second, when there is no more finite subtree, there is a subtree represented by a strongly connected graph. The dictionary D_G stores the tree keys of such graphs, and after `shareWithDone` and if necessary, `share`, we can decide whether another equivalent graph has already been encountered, and if not, use new nodes. When the strongly connected graph is treated, it is considered as just a node, and so we can iterate on our algorithm until we give the representation of the root.

4.2 The Algorithm

We suppose given a dictionary D which maps full keys to nodes corresponding to a unique representation of the associated tree, and a dictionary D_G which maps tree keys (in fact keys of these finite trees) to nodes corresponding to a unique representation of the associated strongly connected graph.

The algorithm uses local dictionaries too, which we assume to be empty when the process starts on a tree. The dictionary `encountered` contains the nodes of the original representation already encountered (so that we do not loop). The set `returnNodes` is used to detect the roots of the loops.

A node is considered “treated” when it is in the dictionary D (and so it represents uniquely a tree). To decide whether a node is “treated”, we just have to look at its key: it is “treated” if the key is full.

`representation(t)`

```

Step 1 if  $t \in \text{encountered}$  then
        if encountered( $t$ ) is not treated add it in returnNodes
        return encountered( $t$ )
Step 2  $N$  is a new node labeled by the empty partial key  $k$  of name
        the label of  $t$ 
Step 3 for each child  $t_i$  of  $t$  do
        3a  $N_i \leftarrow \text{representation}(t_i)$ 
        3b if  $N_i$  is treated, then add it to  $k$ 
            else  $N.i \leftarrow N_i$ 
Step 4 if  $k$  is full then
        if  $k \in D$  return  $D(k)$ 
        else add  $k \rightarrow N$  to  $D$  and return  $N$ 
Step 5 remove  $N$  from returnNodes
Step 6 if returnNodes =  $\emptyset$  then return representCycle( $N$ )
Step 7 return  $N$ 

```

`representCycle(N)`

```

Step 1 if shareWithDone( $N$ )  $\neq N$  then return shareWithDone( $N$ )
Step 2 share( $N$ )
Step 3 if treeKey( $N$ )  $\in D_G$  then return  $D_G(\text{treeKey}(N))$ 
Step 4 for each node  $M$  in the graph defined by  $N$  do
        4a add treeKey( $M$ )  $\rightarrow M$  to  $D_G$ 
        4b add the children of  $M$  to its partial key  $m$ 

```


4c add $m \rightarrow M$ to D
 Step 5 **return** N

4.3 Example

We present the algorithm to represent regular trees on an example, the graph of Fig. 5, where each node is assigned a number. We will write t_i for the tree

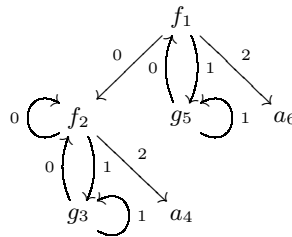


Fig. 5. Example

represented by the node number i , \square_i .

`representation`(t_1) calls `representation` for t_2 , t_5 and t_6 . The call to `representation` on t_2 will return the node \square_2 . It will also store various nodes in D , and in particular $(a) \rightarrow \square_4$. The call on t_5 will just return an untreated node \square_5 , with nothing added in the dictionaries. The call on t_6 will recognize on step 4 that a is in D and so it will return \square_4 .

Thus, at step 5, `returnNodes` = $\{\square_1\}$ becomes empty, and we call `re-`

`presentCycle` with the graph² $(f, \square_2 \bullet \square_4)$. A call to `shareWithDone` returns $(g, \bullet \bullet)$.

the node \square_2 . So the return value of `representation` on t_1 is \square_2 , the node

labeled by f in the graph $(f, \square_2 \bullet \square_4)$. Moreover, the dictionaries will be:

$$D = \{(a) \rightarrow \square_4, (g, \square_3 \square_2) \rightarrow \square_3, (f, \square_2 \square_3 \square_4) \rightarrow \square_2\}$$

² Remember that $(f, \square_2 \bullet \square_4)$ is the partial key which is not defined on its second component.

$$D_G = \left\{ \begin{array}{c} \begin{array}{c} (f, \bullet \bullet \square_4) \\ \swarrow \quad \downarrow \\ \varepsilon \quad (g, \bullet \bullet) \\ \swarrow \quad \searrow \\ \varepsilon \quad 1 \end{array} \quad \rightarrow \square_2, \quad \begin{array}{c} (g, \bullet \bullet) \\ \downarrow \quad \swarrow \\ (f, \bullet \bullet \square_4) \quad \varepsilon \\ \swarrow \quad \searrow \\ 1 \quad \varepsilon \end{array} \rightarrow \square_3 \end{array} \right\}$$

4.4 Proof of the Algorithm

The algorithm returns the node of a graph. We must prove that this graph represents the same tree as the original graph, and that it is a graph of maximal sharing.

First, notice that the algorithm terminates, because of the dictionary **encountered** which implies that each node of the original graph is treated only once.

The correctness of the algorithm is derived from the fact that we return the same graph as the original, except when we recognize that an equivalent node had already been encountered (through the node keys or the tree keys), in which case we replace one node by the other. It is the case step 4 of **representation**, and steps 1, 2 and 3 of **representCycle**

The fact that the resulting graph has the minimal number of nodes lies in the use of the dictionaries D and D_G to ensure that we never duplicate any node. The dictionary D contains the node keys of every node encountered, and the dictionary D_G contains the tree key of every node of every strongly connected graph with minimal number of nodes we encounter. We can prove that each time we definitely introduce new nodes, there is no duplication. Definitive introduction is performed in two points: step 4 of **representation**, and step 4 of **representCycle**.

Step 4 of **representation**, we know that the key k is not in D . Moreover, each one of the N_i composing the key is unique because nodes in partial keys

have already been treated. So if a tree $\begin{array}{c} f \\ \swarrow \quad \searrow \\ t_0 \quad \dots \quad t_{n-1} \end{array}$ had already been encountered,

the key $(f, (N_i)_{i < n})$ would already have been encountered.

Step 4 of **representCycle**, we know that the key $\text{treeKey}(\text{share}(N))$ has never been encountered before. Because such a key is valid for strongly connected graphs, it means that no other node M such that $M \equiv_{\text{tree}} N$ have been encountered before. But the problem is that we have a partial key semantics on these graphs, and $\equiv_{\text{tree}} \subset \equiv_{\text{pk}}$, so we could have $M \not\equiv_{\text{tree}} N$ but $M \equiv_{\text{pk}} N$ in effect representing the same tree. Because $M \not\equiv_{\text{tree}} N$, there is a path p such that $M.p$ and $N.p$ do not have the same label, k_M and k_N . But as N and M represent the same tree, k_M and k_N must have the same name, so their only possible difference is in the partial function. It means there is an i such that one of the keys is defined on i and not the other key (if both of them were defined on i , their value would be the same on i , as the nodes in partial keys are unique representations). By construction, the nodes M and N are in strongly connected graphs. So if one of the keys is not defined on i , there is a q such that

$M.piq = M$ or $N.piq = N$. If t is the tree represented by both nodes, it means that $t_{[piq]} = t$. Suppose k_M is defined on i , then there is a node reachable from $k_M(i)$ which represents the same tree as M , and as such it would have been found by `shareWithDone`. So the graph defined by M would never have gone beyond the step 1 of `representCycle`. It means that another representative is stored for the cycle (we go on like this until we find one which is equivalent to N , which means that the test step 3 could not have been false). If k_N is defined on i , by the same argument, we could not have been beyond the step 1, and so no new node is created.

If no node equivalent to N has been encountered, it is the same for every other node M in the graph represented by N . It is due to the strong connectivity of the graph which implies that if M has already been encountered, N has already been encountered.

5 Complexity Issues

Algorithms on shared trees can be more difficult than standard algorithms on trees, because we must keep the uniqueness of the representation, and for efficiency, we must do it incrementally. Comparing complexities of algorithms on the two representations (the naive and the sharing ones) is difficult, though. The complexity is measured with respect to the size of the inputs of the algorithms, which can be reduced to the number of nodes of the inputs in our case. In the case of shared regular trees, the number of nodes is exactly the number of distinct subtrees of the tree, but when the tree is not shared, the number of nodes can be of any value greater than the number of distinct subtrees. In the sequel, we denote by n this number of nodes, but we must keep in mind that this n can be much bigger in the case of non-shared trees.

The basic property of shared trees is the uniqueness of the representation. Thus, testing tree equality is really immediate: we just compare the memory location of the root. In the classic case, the best method uses a partitioning algorithm. Another case where we can avoid such a computation with shared trees is testing if a tree is a subtree of another one. In the shared case, we just have to compare the root of the first tree with all the nodes of the second one. Not only is it linear, but the second tree is very likely to have very less nodes in the shared case than in the classic representation.

When building finite trees, we need only one operation, which we call root

construction: we give a label f and the nodes $(N_i)_{i < n}$, and we build $\begin{matrix} f \\ \swarrow \searrow \\ N_0 \dots N_{n-1} \end{matrix}$.

Such an operation is constant time in the naive representation and in the sharing representation for finite trees (assuming hashing is constant time [12, 3]). It is indeed also constant time for infinite trees, but this operation does not suffice to build any regular tree. We need also some loop building mechanism. We call this second operation recursive construction. Considering a tree t and a label x , it consists in replacing every edge going to x by an edge to the root, and then apply `representCycle` to maintain the uniqueness of the representation. Concerning

the complexity of this algorithm, it seems that the prevailing operation is the final (and unique) call to `share`, which is applied on the smallest possible subgraph, but in the worst case, the quadratic complexity of `shareWithDone` will take precedence.

Many other operations can be adapted to shared trees while preserving the uniqueness of the representation by derivation from the `representation` algorithm. But due to lack of space, we let the reader write their own adaptations.

	sharing representation	naive representation
testing $t_1 = t_2$	$\mathcal{O}(1)$	$\mathcal{O}((n_1 + n_2) \log(n_1 + n_2))$
testing t_1 subtree of t_2	$\mathcal{O}(n_2)$	$\mathcal{O}((n_1 + n_2) \log(n_1 + n_2))$
building $t_{[p]}$	$\mathcal{O}(p)$	$\mathcal{O}(p)$
root construction	$\mathcal{O}(1)$	$\mathcal{O}(1)$
recursive construction	$\mathcal{O}(n^2)$	$\mathcal{O}(n)$

Fig. 6. Summary of worst case time complexities

The summary suggests that if we are to perform equality testing, it can be beneficial to perform sharing during the calculus. What we show here are worst case complexity, though, and the difficult cases are quite pathological, and thanks to some simple optimizations, they are quite rare. The situation is very similar to the complexity of operations on BDDs [2] compared to the operations on boolean formulas. The size of the formula representing a given boolean function is unbounded, but the basic operations, like conjunctions, are linear in the size of one of the formulas whereas they are quadratic for the BDDs. Nevertheless, in practice BDDs are far more efficient.

6 Application: Set-Based Analysis

We propose to use these techniques to improve the representations of sets of trees. The expressive power of this improved representation is exactly what is needed in set-based analysis [9], where sets of trees are approximated by ignoring the dependencies between variables (an idea which was already present in [16, 11]).

6.1 Tree Automata and Graphs

Because the cartesian approximation eliminates any dependencies between children of a tree, we can use deterministic top-down tree automata in set-based analysis. The idea we use here is that deterministic top-down tree automata can be seen as graphs, where the only properties that matter are path properties, and so it can be represented efficiently as a regular infinite tree.

A deterministic top-down tree automaton [17, 8] is a tuple (Q, I, δ, F) where Q is a finite set of states, $I \in Q$ is the initial state, $F \subset Q$ is a set of final

states, and $\delta : A \times Q \rightarrow Q \times \dots \times Q$ is the transition function which takes a label in A and a state, and returns a sequence of states (as many as the arity of the label). The corresponding graph G is such that $G^N = Q$, $G^E = \{(q, q', a_i) \mid \delta(a, q) = (\dots, q', \dots)$ and q' in i^{th} position $\}$. This connection means that we can represent the sets used in set-based analysis without any variable name in the representation, and in a shared way.

6.2 Tree Skeletons

In order to represent the sets of set-based analysis as trees, we use a new label to represent the anonymous states of the tree automata. This label, which we call a choice label corresponds to a possible union in the interpretation of the infinite tree. We denote this label \bigcirc . We call the infinite trees with this extra label a tree skeleton. The set of trees represented by a tree skeleton is defined³ by:

$$\text{Set} \left(\begin{array}{c} f \\ \swarrow \searrow \\ t_0 \dots t_{n-1} \end{array} \right) \stackrel{\text{def}}{=} \left\{ \begin{array}{c} f \\ \swarrow \searrow \\ u_0 \dots u_{n-1} \end{array} \mid \forall i < n, u_i \in \text{Set}(t_i) \right\}$$

$$\text{Set} \left(\begin{array}{c} \bigcirc \\ \swarrow \searrow \\ t_0 \dots t_{n-1} \end{array} \right) \stackrel{\text{def}}{=} \bigcup_{i < n} \text{Set}(t_i)$$

In order to have a unique representation of the sets of trees (and so keep the constant time equality testing and memoizing properties), we make some restrictions on what infinite trees are considered valid tree skeletons. First we eliminate unnecessary choices: if a choice node has only one child, it is replaced by its child. If a choice node is the child of a choice node, it is replaced by its children. We perform the cartesian approximation: if two children of a choice node have the same label, they are merged (replaced by their cartesian upper approximation). Finally, the children of a choice node are ordered according to their labels. See the summary of figure 7.

Any deterministic top-down tree automaton can be represented by a valid tree skeleton. Consider an automaton (Q, I, δ, F) . We first build the infinite tree labeled by Q and A , such that the root is labeled by I , the children of a given state q are the different a such that $\delta(q, a)$ is defined, and the children of such a a are the $\delta(q, a)$. This tree is regular because there is at most one subtree labeled by a given $q \in Q$, and at most $|Q|$ subtrees labeled by a given $a \in A$. The second step consists in removing every label of arity 0 which does not come from a state in F , and in replacing every state by \bigcirc . Then we derive the valid tree skeleton.

6.3 Using Tree Skeletons in Analysis

Manipulation of tree skeletons uses basic algorithms on shared infinite regular trees. Once we can keep the maximal sharing property, it is easy to keep track of

³ Set is defined as the least fixpoint of this set of equations. The ordering is the pointwise ordering of the inclusion of the images. If we wanted to include infinite trees (as in [5]), we would take the greatest fixpoint.

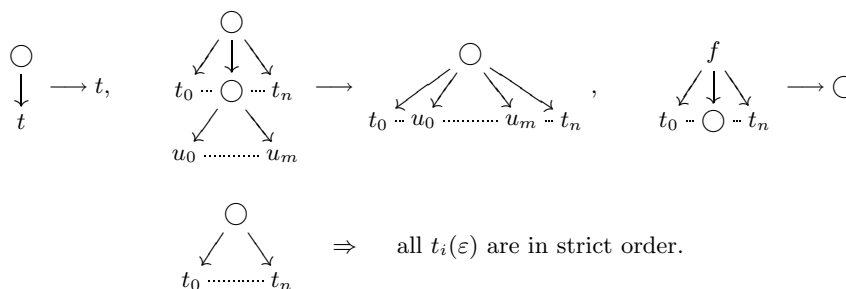


Fig. 7. Rules to obtain a valid tree skeleton

the other rules for tree skeletons. Then tree skeletons can be used everywhere we consider a set of trees in the analysis. It can replace some of the tree automata of [7] (if we keep the original restrictions of set-based analysis), or the tree grammars of [13], as the approximation on union corresponds indeed to cartesian approximation.

In practice, you can try to use the toolbox under development at the following address: <http://www.di.ens.fr/~mauborgn/skeleton.tar.gz>.

7 Conclusion

While trying to improve the representation of sets of trees in set-based analysis, we presented generic algorithms to manipulate efficiently any structure encoded as infinite regular trees. These algorithms allow a very compact representation of such structures and a constant time equality testing. One of their advantages is their incrementality which allows their use on dynamic structures. The complexity analysis cannot describe the potential benefit of this new representation, but it suggests the same gain as for Binary Decision Diagrams which use similar techniques.

We also described a new way of representing sets of trees using infinite regular trees. This new representation is sharing, incremental and unique. Current work includes the integration of the representation in an actual analyzer to show experimentally its benefits.

Acknowledgments

Many thanks are due to the anonymous referees for their very useful comments.

References

- [1] AHO, A. V., HOPCROFT, J. E., AND ULLMAN, J. D. *Data Structures and Algorithms*. Addison-Wesley, 1983.

- [2] BRYANT, R. E. Graph based algorithms for boolean function manipulation. *IEEE Transactions on Computers C-35* (August 1986), 677–691.
- [3] CAI, J., AND PAIGE, R. Using multiset discrimination to solve language processing without hashing. *Theoretical Computer Science* (1994). also, U. of Copenhagen Tech. Report, DIKU-TR Num. D-209, 94/16, URL <ftp://ftp.diku.dk/diku/semantics/papers/D-209.ps.Z>.
- [4] CARDON, A., AND CROCHEMORE, M. Partitioning a graph in $O(|A|\log_2|V|)$. *Theoretical Computer Science 19* (1982), 85–98.
- [5] CHARATONIK, W., AND PODELSKI, A. Co-definite set constraints. In *9th International Conference on Rewriting Techniques and Applications* (March-April 1998), T. Nipkow, Ed., vol. 1379 of *Lecture Notes in Computer Science*, Springer-Verlag, pp. 211–225.
- [6] COLMERAUER, A. PROLOG and infinite trees. In *Logic Programming* (1982), K. L. Clark and S.-A. Tärnlund, Eds., vol. 16 of *APIC Studies in Data Processing*, Academic Press, pp. 231–251.
- [7] DEVIENNE, P., TALBOT, J., AND TISON, S. Solving classes of set constraints with tree automata. In *3th International Conference on Principles and Practice of Constraint Programming* (October 1997), G. Smolka, Ed., vol. 1330 of *Lecture Notes in Computer Science*, Springer-Verlag, pp. 62–76.
- [8] GÉCSEGE, F., AND STEINBY, M. *Tree Automata*. Akadémia Kiadó, 1984.
- [9] HEINTZE, N. *Set Based Program Analysis*. PhD thesis, School of Computer Science, Carnegie Mellon University, October 1992.
- [10] HOPCROFT, J. An $n \log n$ algorithm for minimizing states in a finite automaton. In *Theory of machines and computations* (1971), Z. Kohavi and A. Paz, Eds., Academic Press, pp. 189–196.
- [11] JONES, N. D., AND MUCHNICK, S. S. Flow analysis and optimization of LISP-like structures. In *6th POPL* (January 1979), ACM Press, pp. 244–256.
- [12] KNUTH, D. E. *Sorting and Searching*, vol. 3 of *The Art of Computer Programming*. Addison-Wesley, 1973.
- [13] LIU, Y. A. Dependence analysis for recursive data. In *IEEE International Conference on Computer Languages* (May 1998), pp. 206–215.
- [14] MAUBORGNE, L. Binary decision graphs. In *Static Analysis Symposium (SAS'99)* (1999), A. Cortesi and G. Filé, Eds., vol. 1694 of *Lecture Notes in Computer Science*, Springer-Verlag, pp. 101–116.
- [15] MICHIE, D. “memo” functions and machine learning. *Nature 218* (April 1968), 19–22.
- [16] REYNOLDS, J. Automatic computation of data set definitions. In *Information Processing '68* (1969), Elsevier Science Publisher, pp. 456–461.
- [17] THATCHER, J. W., AND WRIGHT, J. B. Generalized finite automata with an application to a decision problem of second-order logic. *Mathematical Systems Theory 2* (1968), 57–82.