

# On Sharing, Memoization, and Polynomial Time\* (Long Version)

Martin Avanzini<sup>†</sup>      Ugo Dal Lago<sup>‡</sup>

November 19, 2018

We study how the adoption of an evaluation mechanism with sharing and memoization impacts the class of functions which can be computed in polynomial time. We first show how a natural cost model in which lookup for an already computed value has no cost is indeed invariant. As a corollary, we then prove that the most general notion of ramified recurrence is sound for polynomial time, this way settling an open problem in implicit computational complexity.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Related Work . . . . .	3
<b>2</b>	<b>The Need for Sharing and Memoisation</b>	<b>3</b>
<b>3</b>	<b>Preliminaries</b>	<b>4</b>
<b>4</b>	<b>Memoization and Sharing, Formally</b>	<b>9</b>
<b>5</b>	<b>GRSR is Sound for Polynomial Time</b>	<b>18</b>
<b>6</b>	<b>Conclusion</b>	<b>19</b>

---

\*This work was partially supported by FWF project number J 3563.

<sup>†</sup>INRIA Sophia Antipolis

<sup>‡</sup>University of Bologna & INRIA Sophia Antipolis

# 1 Introduction

Traditionally, complexity classes are defined by giving bounds on the amount of resources algorithms are allowed to use while solving problems. This, in principle, leaves open the task of understanding the *structure* of complexity classes. As an example, a given class of functions is not necessarily closed under composition or, more interestingly, under various forms of recursion. When the class under consideration is not too large, say close enough to the class of *polytime computable functions*, closure under recursion does not hold: iterating over an efficiently computable function is not necessarily efficiently computable, e.g. when the iterated function grows more than linearly. In other words, characterizing complexity classes by purely recursion-theoretical means is non-trivial.

In the past twenty years, this challenge has been successfully tackled, by giving *restricted* forms of recursion for which not only certain complexity classes are closed, but which *precisely* generate the class. This has been proved for classes like PTIME, PSPACE, the polynomial hierarchy PH, or even smaller ones like NC (more information about related work is in Section 1.1). A particularly fruitful direction has been the one initiated by Bellantoni and Cook, and independently by Leivant, which consists in restricting the primitive recursive scheme by making it *predicative*, thus forbidding those nested recursive definitions which lead outside the classes cited above. Once this is settled, one can tune the obtained scheme by either adding features (e.g. parameter substitutions) or further restricting the scheme (e.g. by way of linearization).

Something a bit disappointing in this field is that the expressive power of the simplest (and most general) form of predicative recurrence, namely *simultaneous* recurrence on *generic algebras* is unknown. If algebras are restricted to be *string* algebras, or if recursion is not simultaneous, soundness for polynomial time computation is known to hold [15, 20]. The two soundness results are obtained by quite different means, however: in presence of trees, one is forced to handle *sharing* [15] of common sub-expressions, while simultaneous definitions by recursion requires a form of *memoization* [20].

In this paper, we show that sharing and memoization can indeed be reconciled, and we exploit both to give a new invariant time cost model for the evaluation of rewrite systems. That paves the way towards a polytime soundness for simultaneous predicative recursion on generic algebras, thus solving the open problem we were mentioning. More precisely, with the present paper we make the following contributions:

1. We define a simple functional programming language. The domain of the defined functions is a free algebra formed from constructors. Hence we can deal with functions over strings, lists, but also trees (see Section 3). We then extend the underlying rewriting based semantics with *memoization*, i.e. intermediate results are automatically tabulated to avoid expensive re-computation (Section 4). As standard for functional programming languages such as Haskell or OCaml, data is stored in a *heap*, facilitating *sharing* of common sub-expression. To measure the *runtime* of such programs, we employ a novel cost model, called *memoized runtime complexity*, where each function application counts one time unit, but lookups of tabulated calls do not have to be accounted.
2. Our *invariance theorem* (see Theorem 4.17) relates, within a polynomial overhead, the memoized runtime complexity of programs to the cost of implementing the defined functions on a classical model of computation, e.g. *Turing* or *random access machines*. The invariance theorem thus confirms that our cost model truthfully represents the computational complexity of the defined function.
3. We extend upon Leivant's notion of *ramified recursive functions* [19] by allowing definitions

by *generalised ramified simultaneous recurrence* (*GRSR* for short). We show that the resulting class of functions, defined over arbitrary free algebras have, when implemented as programs, polynomial memoized runtime complexity (see Theorem 5.3). By our invariance theorem, the function algebra is sound for polynomial time, and consequently GSR characterizes the class of polytime computable functions.

## 1.1 Related Work

That predicative recursion *on strings* is sound for polynomial time, even in presence of simultaneous recursive definitions, is known for a long time [8]. Variations of predicative recursion have been later considered and proved to characterize classes like PH [9], PSPACE [22], EXPTIME [3] or NC [11]. Predicative recursion on trees has been claimed to be sound for polynomial time in the original paper by Leivant [19], the long version of which only deals with strings [20]. After fifteen years, the non-simultaneous case has been settled by the second author with Martini and Zorzi [15]; their proof, however, relies on an ad-hoc, infinitary, notion of graph rewriting. Recently, ramification has been studied in the context of a simply-typed  $\lambda$ -calculus in an unpublished manuscript [16]; the authors claim that a form of ramified recurrence on trees captures polynomial time; this, again, does not take simultaneous recursion into account.

The formalism presented here is partly inspired by the work of Hoffmann [18], where sharing and memoization is shown to work well together in the realm of term graph rewriting. The proposed machinery, although powerful, is unnecessarily complicated for our purposes. Speaking in Hoffmann's terms, our results require a form of full memoization, which *is* definable in Hoffmann's system. However, most crucially for our concerns, it is unclear how the overall system incorporating full memoization can be implemented efficiently, if at all.

## 2 The Need for Sharing and Memoisation

This Section is an informal, example-driven, introduction to ramified recursive definitions and their complexity. Our objective is to convince the reader that those definitions do *not* give rise to polynomial time computations if naively evaluated, and that sharing and memoization are *both* necessary to avoid exponential blowups.

In Leivant's system [20], functions and variables are equipped with a *tier*. Composition must preserve tiers and, crucially, in a function defined by primitive recursion the tier of the recurrence parameter must be higher than the tier of the recursive call. This form of *ramification* of functions effectively tames primitive recursion, resulting in a characterisation of the class of *polytime computable functions*.

Of course, ramification also controls the growth rate of functions. However, as soon as we switch from strings to a domain where tree structures are definable, this control is apparently lost. For illustration, consider the following definition.

$$\mathbf{tree}(\mathbf{0}) = \mathbf{L} \quad \mathbf{tree}(\mathbf{S}(n)) = \mathbf{br}(\mathbf{tree}(n)) \quad \mathbf{br}(t) = \mathbf{B}(t, t).$$

The function  $\mathbf{tree}$  is defined by primitive recursion, essentially from basic functions. As a consequence, it is easily seen to be ramified in the sense of Leivant. Even though the number of recursive steps is linear in the input, the result of  $\mathbf{tree}(\mathbf{S}^n(\mathbf{0}))$  is the complete binary tree of height  $n$ . As thus the length of the output is exponential in the one of its input, there is, at least apparently, little hope to prove  $\mathbf{tree}$  a polytime function. The way out is

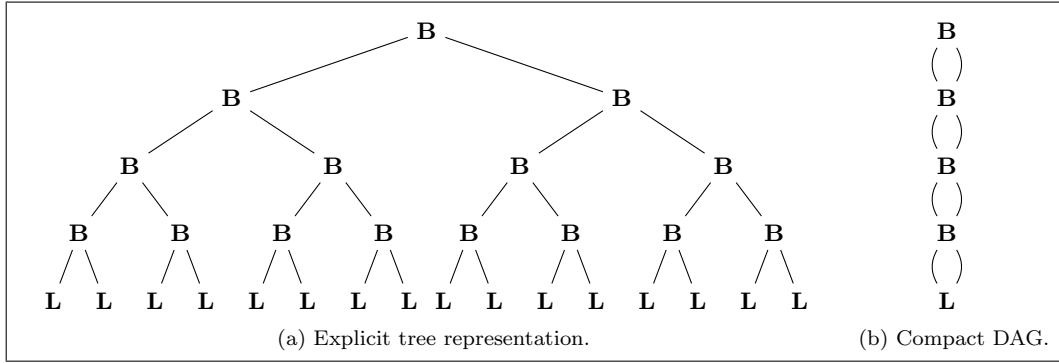


Figure 1: Complete Binary Tree of Height Four, as Computed by  $\text{tree}(\mathbf{S}^4(\mathbf{0}))$ .

*sharing*: the complete binary tree of height  $n$  can be compactly represented as a *directed acyclic graph* (DAG for short) of linear size (see Figure 1). Indeed, using the compact DAG representation it is easy to see that the function  $\text{tree}$  is computable in polynomial time. This is the starting point of [15], in which general ramified recurrence is proved sound for polynomial time. A crucial observation here is that *not only* the output's size, but also the total amount of work can be kept under control, thanks to the fact that evaluating a primitive recursive definition on a compactly represented input can be done by constructing an isomorphic DAG of recursive calls.

This does not scale up to *simultaneous* ramified recurrence. The following example computes the genealogical tree associated with *Fibonacci's rabbit problem* for  $n \in \mathbb{N}$  generations. Rabbits come in pairs. After one generation, each *baby* rabbit pair ( $\mathbf{N}$ ) matures. In each generation, an *adult* rabbit pair ( $\mathbf{M}$ ) bears one pair of babies.

$$\begin{array}{lll} \text{rabbits}(\mathbf{0}) = \mathbf{N}_L & \mathbf{a}(\mathbf{0}) = \mathbf{M}_L & \mathbf{b}(\mathbf{0}) = \mathbf{N}_L \\ \text{rabbits}(\mathbf{S}(n)) = \mathbf{b}(n) & \mathbf{a}(\mathbf{S}(n)) = \mathbf{M}(\mathbf{a}(n), \mathbf{b}(n)) & \mathbf{b}(\mathbf{S}(n)) = \mathbf{N}(\mathbf{a}(n)) . \end{array}$$

The function  $\text{rabbits}$  is obtained by case analysis from the functions  $\mathbf{a}$  and  $\mathbf{b}$ , which are defined by *simultaneous* primitive recursion: the former recursively calls itself *and* the latter, while the latter makes a recursive call to the former. The output of  $\text{rabbits}(\mathbf{S}^n(\mathbf{0}))$  is tightly related to the sequence of Fibonacci numbers: the number of nodes at depth  $i$  is given by the  $i^{\text{th}}$  Fibonacci number. Hence the output tree has exponential size in  $n$  but, again, can be represented compactly (see Figure 2). This does not suffice for our purposes, however. In presence of simultaneous definitions, indeed, avoiding re-computation of previously computed values becomes more difficult, the trick described above does not work, and the key idea towards that is the use of *memoization*.

What we prove in this paper is precisely that sharing and memoization can indeed be made to work together, and that they together allow to prove polytime soundness for all ramified recursive functions, also in presence of tree algebras and simultaneous definitions.

### 3 Preliminaries

**General Ramified Simultaneous Recurrence** Let  $\mathbb{A}$  denote a finite (*untyped*) signatures  $\mathcal{F}$  of constructors  $\mathbf{c}_1, \dots, \mathbf{c}_k$ , each equipped with an arity  $\text{ar}(\mathbf{c}_i)$ . In the following, the set of

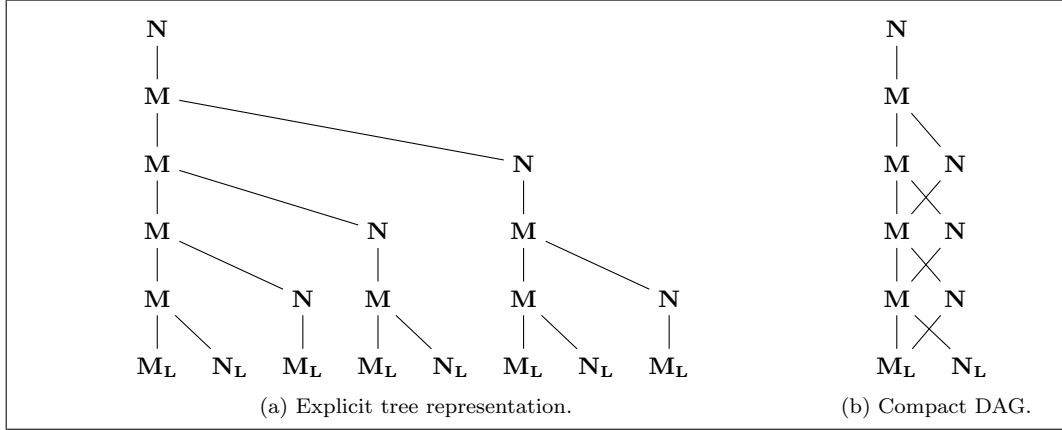


Figure 2: Genealogical Rabbit Tree up to the Sixth Generation, as Computed by  $\text{rabbits}(\mathbf{S}^6(\mathbf{0}))$ .

terms  $\mathcal{T}(\mathbb{A})$  is also denoted by  $\mathbb{A}$  if this does not create ambiguities. We are interested in total functions from  $\mathbb{A}^n = \underbrace{\mathbb{A} \times \dots \times \mathbb{A}}_{n \text{ times}}$  to  $\mathbb{A}$ .

**Definition 3.1.** *The following are so-called basic functions:*

- For each constructor  $\mathbf{c}$ , the constructor function  $\mathbf{f}_{\mathbf{c}} : \mathbb{A}^{\text{ar}(\mathbf{c})} \rightarrow \mathbb{A}$  for  $\mathbf{c}$ , defined as follows:  
 $\mathbf{f}_{\mathbf{c}}(x_1, \dots, x_{\text{ar}(\mathbf{c})}) = \mathbf{c}(x_1, \dots, x_{\text{ar}(\mathbf{c})})$
- For each  $1 \leq n \leq m$ , the  $(m, n)$ -projection function  $\Pi_n^m : \mathbb{A}^m \rightarrow \mathbb{A}$  defined as follows:  
 $\Pi_n^m(x_1, \dots, x_m) = x_n$ .

**Definition 3.2.**

- Given a function  $\mathbf{f} : \mathbb{A}^n \rightarrow \mathbb{A}$  and  $n$  functions  $\mathbf{g}_1, \dots, \mathbf{g}_n$ , all of them from  $\mathbb{A}^m$  to  $\mathbb{A}$ , the composition  $\mathbf{h} = \mathbf{f} \circ (\mathbf{g}_1, \dots, \mathbf{g}_n)$  is a function from  $\mathbb{A}^m$  to  $\mathbb{A}$  defined as follows:  
 $\mathbf{h}(\vec{x}) = \mathbf{f}(\mathbf{g}_1(\vec{x}), \dots, \mathbf{g}_n(\vec{x}))$ .
- Suppose given the functions  $\mathbf{f}_i$  where  $1 \leq i \leq k$  such that for some  $m$ ,  $\mathbf{f}_i : \mathbb{A}^{\text{ar}(\mathbf{c}_i)} \times \mathbb{A}^n \rightarrow \mathbb{A}$ . Then the function  $\mathbf{g} = \text{case}(\{\mathbf{f}_i\}_{1 \leq i \leq k})$  defined by case distinction from  $\{\mathbf{f}_i\}_{1 \leq i \leq k}$  is a function from  $\mathbb{A} \times \mathbb{A}^n$  to  $\mathbb{A}$  defined as follows:  $\mathbf{g}(\mathbf{c}_i(\vec{x}), \vec{y}) = \mathbf{f}_i(\vec{x}, \vec{y})$ .
- Suppose given the functions  $\mathbf{f}_i^j$ , where  $1 \leq i \leq k$  and  $1 \leq j \leq n$ , such that for some  $m$ ,  $\mathbf{f}_i^j : \mathbb{A}^{\text{ar}(\mathbf{c}_i)} \times \mathbb{A}^{n \cdot \text{ar}(\mathbf{c}_i)} \times \mathbb{A}^m \rightarrow \mathbb{A}$ . The functions  $\{\mathbf{g}_j\}_{1 \leq j \leq n} = \text{simrec}(\{\mathbf{f}_i^j\}_{1 \leq i \leq k, 1 \leq j \leq n})$  defined by simultaneous primitive recursion from  $\{\mathbf{f}_i^j\}_{1 \leq i \leq k, 1 \leq j \leq n}$  are all functions from  $\mathbb{A} \times \mathbb{A}^m$  to  $\mathbb{A}$  such that for  $\vec{x} = x_1, \dots, x_{\text{ar}(\mathbf{c}_i)}$ ,

$$\mathbf{g}_j(\mathbf{c}_i(\vec{x}), \vec{y}) = \mathbf{f}_i^j(\vec{x}, \mathbf{g}_1(x_1, \vec{y}), \dots, \mathbf{g}_1(x_{\text{ar}(\mathbf{c}_i)}, \vec{y}), \dots, \mathbf{g}_n(x_1, \vec{y}), \dots, \mathbf{g}_n(x_{\text{ar}(\mathbf{c}_i)}, \vec{y}), \vec{y}) .$$

We denote by  $\text{SIMREC}(\mathbb{A})$  the class of *simultaneous recursive functions over  $\mathbb{A}$* , defined as the smallest class containing the basic functions of Definition 3.1 and that is closed under the schemes of Definition 3.2.

*Tiering*, the central notion underlying Leivant’s definition of *ramified recurrence*, consists in attributing *tiers* to inputs and outputs of some functions among the ones constructed as above, with the goal of isolating the polytime computable ones. Roughly speaking, the role of tiers is to single out “a copy” of the signature by a level: this level permits to control the

$\frac{\mathbf{f}_c \triangleright \mathbb{A}_n^{\text{ar}(c)} \rightarrow \mathbb{A}_n}{\mathbf{f}_c \triangleright \mathbb{A}_n^{\text{ar}(c)} \rightarrow \mathbb{A}_n}$	$\frac{\Pi_m^n \triangleright \mathbb{A}_{p_1} \times \dots \times \mathbb{A}_{p_m} \rightarrow \mathbb{A}_{p_n}}{\Pi_m^n \triangleright \mathbb{A}_{p_1} \times \dots \times \mathbb{A}_{p_m} \rightarrow \mathbb{A}_{p_n}}$	$\frac{\mathbf{f}_i \triangleright \mathbb{A}_p^{\text{ar}(c_i)} \times \mathbf{A} \rightarrow \mathbb{A}_m}{\mathbf{f}_i \triangleright \mathbb{A}_p^{\text{ar}(c_i)} \times \mathbf{A} \rightarrow \mathbb{A}_m}$
$\frac{\mathbf{f}_c \triangleright \mathbb{A}_n^{\text{ar}(c)} \rightarrow \mathbb{A}_n \quad \Pi_m^n \triangleright \mathbb{A}_{p_1} \times \dots \times \mathbb{A}_{p_m} \rightarrow \mathbb{A}_{p_n} \quad \text{case}(\{\mathbf{f}_i\}_{1 \leq i \leq k}) \triangleright \mathbb{A}_p \times \mathbf{A} \rightarrow \mathbb{A}_m}{\mathbf{f}_c \triangleright \mathbb{A}_n^{\text{ar}(c)} \rightarrow \mathbb{A}_n \quad \Pi_m^n \triangleright \mathbb{A}_{p_1} \times \dots \times \mathbb{A}_{p_m} \rightarrow \mathbb{A}_{p_n} \quad \text{case}(\{\mathbf{f}_i\}_{1 \leq i \leq k}) \triangleright \mathbb{A}_p \times \mathbf{A} \rightarrow \mathbb{A}_m}$		
$\mathbf{f} \triangleright \mathbb{A}_{p_1} \times \dots \times \mathbb{A}_{p_n} \rightarrow \mathbb{A}_m$	$\mathbf{g}_i \triangleright \mathbf{A} \rightarrow \mathbb{A}_{p_i}$	$\frac{\mathbf{f}_i^j \triangleright \mathbb{A}_p^{\text{ar}(c_i)} \times \mathbb{A}_m^{n \cdot \text{ar}(c_i)} \times \mathbf{A} \rightarrow \mathbb{A}_m \quad p > m}{\mathbf{f}_i^j \triangleright \mathbb{A}_p^{\text{ar}(c_i)} \times \mathbb{A}_m^{n \cdot \text{ar}(c_i)} \times \mathbf{A} \rightarrow \mathbb{A}_m \quad p > m}$
$\frac{\mathbf{f} \circ (\mathbf{g}_1, \dots, \mathbf{g}_n) \triangleright \mathbf{A} \rightarrow \mathbb{A}_m \quad \text{simrec}(\{\mathbf{f}_i^j\}_{1 \leq i \leq k, 1 \leq j \leq n}) \triangleright \mathbb{A}_p \times \mathbf{A} \rightarrow \mathbb{A}_m}{\mathbf{f} \circ (\mathbf{g}_1, \dots, \mathbf{g}_n) \triangleright \mathbf{A} \rightarrow \mathbb{A}_m \quad \text{simrec}(\{\mathbf{f}_i^j\}_{1 \leq i \leq k, 1 \leq j \leq n}) \triangleright \mathbb{A}_p \times \mathbf{A} \rightarrow \mathbb{A}_m}$		

Figure 3: Tiering as a Formal System.

recursion nesting. Tiering can be given as a formal system, in which judgments have the form  $\mathbf{f} \triangleright \mathbb{A}_{p_1} \times \dots \times \mathbb{A}_{p_{\text{ar}(\mathbf{f})}} \rightarrow \mathbb{A}_m$  for  $p_1, \dots, p_{\text{ar}(\mathbf{f})}, m$  natural numbers and  $\mathbf{f} \in \text{SIMREC}(\mathbb{A})$ . The system is defined in Figure 3, where  $\mathbf{A}$  denotes the expression  $\mathbb{A}_{q_1} \times \dots \times \mathbb{A}_{q_k}$  for some  $q_1, \dots, q_k \in \mathbb{N}$ . Notice that composition preserves tiers. Moreover, recursion is allowed only on inputs of tier higher than the tier of the function (in the case  $\mathbf{f} = \text{simrec}(\{\mathbf{f}_i^j\}_{1 \leq i \leq k, 1 \leq j \leq n})$ , we require  $p > m$ ).

**Definition 3.3.** We call a function  $\mathbf{f} \in \text{SIMREC}(\mathbb{A})$  definable by general ramified simultaneous recurrence (GRSR for short) if  $\mathbf{f} \triangleright \mathbb{A}_{p_1} \times \dots \times \mathbb{A}_{p_{\text{ar}(\mathbf{f})}} \rightarrow \mathbb{A}_m$  holds.

**Remark 3.4.** Consider the word algebra  $\mathbb{W} = \{\epsilon, \mathbf{a}, \mathbf{b}\}$  consisting of a constant  $\epsilon$  and two unary constructors  $\mathbf{a}$  and  $\mathbf{b}$ , which is in bijective correspondence to the set of binary words. Then the functions definable by ramified simultaneous recurrence over  $\mathbb{W}$  includes the ramified recursive functions from Leivant [20], and consequently all polytime computable functions.

**Example 3.5.**

1. Consider  $\mathbb{N} := \{\mathbf{0}, \mathbf{S}\}$  with  $\text{ar}(\mathbf{0}) = 0$  and  $\text{ar}(\mathbf{S}) = 1$ , which is in bijective correspondence to the set of natural numbers. We can define addition  $\text{add} : \mathbb{N}_i \times \mathbb{N}_j \rightarrow \mathbb{N}_j$  for  $i > j$ , by

$$\text{add}(\mathbf{0}, y) = \Pi_1^1(y) = y \quad \text{add}(\mathbf{S}(x), y) = (\mathbf{f}_S \circ \Pi_2^3)(x, \text{add}(x, y), y) = \mathbf{S}(\text{add}(x, y)),$$

using general simultaneous ramified recursion, i.e.  $\{\text{add}\} = \text{simrec}(\{\{\Pi_1^1, \mathbf{f}_S \circ \Pi_2^3\}\})$ .

2. Let  $\mathbb{T} := \{\mathbf{N}_L, \mathbf{M}_L, \mathbf{N}, \mathbf{M}\}$ , where  $\text{ar}(\mathbf{N}_L) = \text{ar}(\mathbf{M}_L) = 0$ ,  $\text{ar}(\mathbf{N}) = 1$  and  $\text{ar}(\mathbf{M}) = 2$ . Then we can define the functions  $\text{rabbits} : \mathbb{N}_i \rightarrow \mathbb{T}_j$  for  $i > j$  from Section 2 by composition from the following two functions, defined by simultaneous ramified recurrence.

$$\begin{aligned} \mathbf{a}(\mathbf{0}) &= \mathbf{M}_L & \mathbf{a}(\mathbf{S}(n)) &= (\mathbf{f}_M \circ (\Pi_2^3, \Pi_3^3))(n, \mathbf{a}(n), \mathbf{b}(n)) = \mathbf{M}(\mathbf{a}(n), \mathbf{b}(n)) \\ \mathbf{b}(\mathbf{0}) &= \mathbf{N}_L & \mathbf{b}(\mathbf{S}(n)) &= (\mathbf{f}_N \circ \Pi_3^3)(n, \mathbf{a}(n), \mathbf{b}(n)) = \mathbf{N}(\mathbf{a}(n)). \end{aligned}$$

3. We can define a function  $\#\text{leaves} : \mathbb{T} \rightarrow \mathbb{N}$  by simultaneous primitive recursion which counts the number of leaves in  $\mathbb{T}$ -trees as follows.

$$\begin{aligned} \#\text{leaves}(\mathbf{N}_L) &= \mathbf{S}(\mathbf{0}) & \#\text{leaves}(\mathbf{M}_L) &= \mathbf{S}(\mathbf{0}) \\ \#\text{leaves}(\mathbf{N}(t)) &= \#\text{leaves}(t) & \#\text{leaves}(\mathbf{M}(l, r)) &= \text{add}(\#\text{leaves}(l), \#\text{leaves}(r)). \end{aligned}$$

However, this function cannot be ramified, since  $\text{add}$  in the last equation requires different tiers. Indeed, having a ramified recursive function  $\#\text{leaves} : \mathbb{T}_i \rightarrow \mathbb{N}_1$  (for some  $i > 1$ ) defined as above would allow us to ramify  $\text{fib} = \#\text{leaves} \circ \text{rabbits}$  which on input  $n$  computes the  $n^{\text{th}}$  Fibonacci number, and is thus an exponential function.

$$\boxed{
\begin{array}{c}
\frac{\mathbf{f} \in \mathcal{F} \quad t_i \downarrow v_i \quad \mathbf{f}(v_1, \dots, v_k) \downarrow v}{\mathbf{f}(t_1, \dots, t_k) \downarrow v} \quad \frac{\mathbf{c} \in \mathcal{C} \quad t_i \downarrow v_i}{\mathbf{c}(t_1, \dots, t_k) \downarrow \mathbf{c}(v_1, \dots, v_k)} \\
\frac{\mathbf{f}(p_1, \dots, p_k) \rightarrow r \in \mathcal{R} \quad \forall i. p_i \sigma = v_i \quad r \sigma \downarrow v}{\mathbf{f}(v_1, \dots, v_k) \downarrow v}
\end{array}
}$$

Figure 4: Operational Semantics for Program  $(\mathcal{F}, \mathcal{C}, \mathcal{R})$ .

**Computational Model, Syntax and Semantics** We introduce a simple, *rewriting based*, notion of program for computing functions over term algebras.

**Definition 3.6.** A program  $P$  is given as a triple  $(\mathcal{F}, \mathcal{C}, \mathcal{R})$  consisting of two disjoint signatures  $\mathcal{F}$  and  $\mathcal{C}$  of operation symbols  $f_1, \dots, f_m$  and constructors  $\mathbf{c}_1, \dots, \mathbf{c}_n$  respectively, and a finite set  $\mathcal{R}$  of rules  $l \rightarrow r$  over terms  $l, r \in \mathcal{T}(\mathcal{F} \cup \mathcal{C}, \mathcal{V})$ . For each rule, the left-hand side  $l$  is of the form  $f_i(p_1, \dots, p_k)$  where the patterns  $p_j$  consist only of variables and constructors, and all variables occurring in the right-hand side  $r$  also occur in the left-hand side  $l$ .

We keep the program  $P = (\mathcal{F}, \mathcal{C}, \mathcal{R})$  fixed throughout the following. Moreover, we require that  $P$  is *orthogonal*, that is, the following two requirements are met:

1. *left-linearity*: the left-hand sides  $l$  of each rule  $l \rightarrow r \in \mathcal{R}$  is *linear*; and
2. *non-ambiguity*: there are no two rules with overlapping left-hand sides in  $\mathcal{R}$ .

Orthogonal programs define a class of deterministic first-order functional programs, see e.g. [5]. The domain of the defined functions is the constructor algebra  $\mathcal{T}(\mathcal{C})$ . Correspondingly, elements of  $\mathcal{T}(\mathcal{C})$  are called *values*, which we denote by  $v, u, \dots$ . In Figure 4 we present the operational semantics, realizing standard *call-by-value* evaluation order. The statement  $t \downarrow v$  means that the term  $t$  *reduces* to the value  $v$ . We say that  $P$  computes the function  $\mathbf{f} : \mathcal{T}(\mathcal{C})^k \rightarrow \mathcal{T}(\mathcal{C})$  if there exists an operation  $f \in \mathcal{F}$ , such that  $\mathbf{f}(v_1, \dots, v_k) = v$  if and only if  $f(v_1, \dots, v_k) \downarrow v$  holds for all inputs  $v_i \in \mathcal{T}(\mathcal{C})$ .

**Example 3.7** (Continued from Example 3.5). *The definition of rabbits from Section 2 can be turned into a program  $P_R$  over constructors of  $\mathbb{N}$  and  $\mathbb{T}$ , by orienting the underlying equations from left to right and replacing applications of functions  $\mathbf{f} \in \{\mathbf{rabbits}, \mathbf{a}, \mathbf{b}\}$  with corresponding operation symbols  $f \in \{\mathbf{rabbits}, \mathbf{a}, \mathbf{b}\}$ . For instance, concerning the function  $\mathbf{a}$ , the defining equations are turned into  $\mathbf{a}(\mathbf{0}) \rightarrow \mathbf{M}_L$  and  $\mathbf{a}(\mathbf{S}(n)) \rightarrow \mathbf{M}(\mathbf{a}(n), \mathbf{b}(n))$ .*

**Definition 3.8.** For  $\mathbf{f} \in \text{SIMREC}(\mathbb{A})$ , by  $P_{\mathbf{f}}$  we denote the program  $(\mathcal{F}_{\mathbf{f}}, \mathcal{C}_{\mathbf{f}}, \mathcal{R}_{\mathbf{f}})$  where:

- the set of operations  $\mathcal{F}_{\mathbf{f}}$  contains for each function  $\mathbf{g}$  underlying the definition of  $\mathbf{f}$  a corresponding operation symbol  $g$ .
- the set of constructors  $\mathcal{C}_{\mathbf{f}}$  contains the constructors of  $\mathbb{A}$ ;
- the set of rules  $\mathcal{R}_{\mathbf{f}}$  contain for each equation  $l = r$  defining a function  $\mathbf{g}$  underlying the definition  $\mathbf{f}$  the orientation  $l \rightarrow r$ .

Notice that due to the inductive definition of the class  $\text{SIMREC}(\mathbb{A})$ , the program  $P_{\mathbf{f}}$  is finite. From the shape of the initial functions and operations (Definition 3.1 and Definition 3.2) it is also clear that  $P_{\mathbf{f}}$  is orthogonal.

**Terms and Term Graphs** Furthermore, we fix a set of *variables*  $\mathcal{V}$  disjoint from function symbols. Terms over a signature  $\mathcal{F}$  and  $\mathcal{V}$  are defined as usual, and form a set  $\mathcal{T}(\mathcal{F}, \mathcal{V})$ . A

term  $t$  is called *ground* if it does not contain variables, it is called *linear* if every variable occurs at most once in  $t$ . The ground terms are collected in  $\mathcal{T}(\mathcal{F})$ . The set of *subterms*  $\text{STs}(t)$  of a term  $t$  is defined by  $\text{STs}(t) := \{t\}$  if  $t \in \mathcal{V}$  and  $\text{STs}(t) := \bigcup_{1 \leq i \leq \text{ar}(f)} \text{STs}(t_i) \cup \{t\}$  if  $t = f(t_1, \dots, t_{\text{ar}(f)})$ . A *substitution*, is a finite mapping  $\sigma$  from variables to terms. By  $t\sigma$  we denote the term obtained by replacing in  $t$  all variables  $x \in \text{dom}(\sigma)$  by  $\sigma(x)$ . If  $s = t\sigma$ , we also say that  $s$  is an *instance* of the term  $t$ .

We borrow key concepts from *term graph rewriting* (see e.g. the survey of Plump [24] for an overview) and follow the presentation of Barendregt et al. [7]. A *term graph*  $T$  over a signature  $\mathcal{F}$  is a *directed acyclic graph* whose nodes are labeled by symbols in  $\mathcal{F} \cup \mathcal{V}$ , and where outgoing edges are ordered. Formally,  $T$  is a triple  $(N, \text{suc}, \text{lab})$  consisting of *nodes*  $N$ , a *successors function*  $\text{suc} : N \rightarrow N^*$  and a *labeling function*  $\text{lab} : N \rightarrow \mathcal{F} \cup \mathcal{V}$ . We require that term graphs are *compatible* with  $\mathcal{F}$ , in the sense that for each node  $o \in N$ , if  $\text{lab}_T(o) = f \in \mathcal{F}$  then  $\text{suc}_T(o) = [o_1, \dots, o_{\text{ar}(f)}]$  and otherwise, if  $\text{lab}_T(o) = x \in \mathcal{V}$ ,  $\text{suc}_T(o) = []$ . In the former case, we also write  $T(o) = f(o_1, \dots, o_{\text{ar}(f)})$ , the latter case is denoted by  $T(o) = x$ . We define the *successor relation*  $\rightarrow_T$  on nodes in  $T$  such that  $o \rightarrow_T p$  holds iff  $p$  occurs in  $\text{suc}(o)$ , if  $p$  occurs at the  $i^{\text{th}}$  position we also write  $o \xrightarrow{i}_T p$ . Throughout the following, we consider only *acyclic* term graphs, that is, when  $\rightarrow_T$  is acyclic. Hence the *unfolding*  $[o]_T$  of  $T$  at node  $o$ , defined by  $[o]_T := x$  if  $T(o) = x \in \mathcal{V}$ , and otherwise  $[o]_T := f([o_1]_T, \dots, [o_k]_T)$  where  $T(o) = f(o_1, \dots, o_k)$ , results in a finite term. We called the term graph  $T$  *rooted* if there exists a unique node  $o$ , the *root* of  $T$ , with  $o \rightarrow_T^* p$  for every  $p \in N$ . We denote by  $T|o$  the *sub-graph* of  $T$  rooted at  $o$ . Consider a symbol  $f \in \mathcal{F}$  and nodes  $\{o_1, \dots, o_{\text{ar}(f)}\} \subseteq N$  of  $T$ . The extension  $S$  of  $T$  by a fresh node  $o_f \notin N$  with  $S(o_f) = f(o_1, \dots, o_{\text{ar}(f)})$  is denoted by  $T \uplus \{o_f \mapsto f(o_1, \dots, o_{\text{ar}(f)})\}$ . We write  $f(T|o_1, \dots, T|o_{\text{ar}(f)})$  for the term graph  $S|o_f$ .

For two rooted term graphs  $T = (N_T, \text{suc}_T, \text{lab}_T)$  and  $S = (N_S, \text{suc}_S, \text{lab}_S)$ , a mapping  $m : N_T \rightarrow N_S$  is called *morphic* in  $o \in N_T$  if (i)  $\text{lab}_T(o) = \text{lab}_S(m(o))$  and (ii)  $o \xrightarrow{i}_T p$  implies  $m(o) \xrightarrow{i}_S m(p)$  for all appropriate  $i$ . A *homomorphism* from  $T$  to  $S$  is a mapping  $m : N_T \rightarrow N_S$  that (i) maps the root of  $T$  to the root of  $S$  and that (ii) is morphic in all nodes  $o \in N_T$  not labeled by a variable. We write  $T \geq_m S$  to indicate that  $m$  is, possibly an extension of, a homomorphism from  $T$  to  $S$ .

Every term  $t$  is trivially representable as a *canonical tree*  $\Delta(t)$  unfolding to  $t$ , using a fresh node for each occurrence of a subterm in  $t$ . For  $t$  a linear term, to each variable  $x$  in  $t$  we can associate a *unique node* in  $\Delta(t)$  labeled by  $x$ , which we denote by  $o_x$ . The following proposition relates matching on terms and homomorphisms on trees. It essentially relies on the imposed linearity condition.

**Proposition 3.9** (Matching on Graphs). *Let  $t$  be a linear term,  $T$  be a term graph and let  $o$  be a node of  $T$ .*

1. *If  $\Delta(t) \geq_m T|o$  then there exists a substitution  $\sigma$  such that  $t\sigma = [o]_T$ .*
2. *Vice versa, if  $t\sigma = [o]_T$  holds for some substitution  $\sigma$  then there exists a morphism  $\Delta(t) \geq_m T|o$ .*

Here, the substitution  $\sigma$  and homomorphism  $m$  satisfy  $\sigma(x) = [m(o_x)]_T$  for all variables  $x$  in  $t$ .

*Proof.* The proof is by induction on  $t$ . We first proof the direction from left to right. Assume  $\Delta(t) \geq_m S|o$ . When  $t$  is a variable, the substitution  $\sigma := \{t \mapsto [o]_S\}$  satisfies  $t\sigma = [o]_S$ . Since  $m(\epsilon) = o$ , we conclude the base case. For the inductive step, assume  $t = f(t_1, \dots, t_k)$ .



$$\begin{array}{c}
\frac{f \in \mathcal{F} \quad (C_{i-1}, t_i) \Downarrow_{n_i} (C_i, v_i) \quad (C_k, f(v_1, \dots, v_k)) \Downarrow_n (C_{k+1}, v) \quad m = n + \sum_{i=1}^k n_i}{(C_0, f(t_1, \dots, t_k)) \Downarrow_m (C_{k+1}, v)} \text{ (Split)} \\
\frac{\mathbf{c} \in \mathcal{C} \quad (C_{i-1}, t_i) \Downarrow_{n_i} (C_i, v_i) \quad m = \sum_{i=1}^k n_i}{(C_0, \mathbf{c}(t_1, \dots, t_k)) \Downarrow_m (C_k, \mathbf{c}(v_1, \dots, v_k))} \text{ (Con)} \quad \frac{(f(v_1, \dots, v_k), v) \in C}{(C, f(v_1, \dots, v_k)) \Downarrow_0 (C, v)} \text{ (Read)} \\
\frac{(f(v_1, \dots, v_k), v) \notin C \quad f(p_1, \dots, p_k) \rightarrow r \in \mathcal{R} \quad \forall i. p_i \sigma = v_i \quad (C, r\sigma) \Downarrow_m (D, v)}{(C, f(v_1, \dots, v_k)) \Downarrow_{m+1} (D \cup \{(f(v_1, \dots, v_k), v)\}, v)} \text{ (Update)}
\end{array}$$

Figure 5: Cost Annotated Operational Semantics with Memoization for Program  $(\mathcal{F}, \mathcal{C}, \mathcal{R})$ .

Fix  $i = 1, \dots, k$ , and define  $m_i(p) = m(i \cdot p)$  for each position  $i \cdot p$  in  $t$ . By case analysis on the nodes of  $\Delta(t_i)$  one verifies  $\Delta(t_i) \geq_{m_i} [n_i]_S$ . Thus by induction hypothesis  $t_i \sigma_i = [n_i]_S$  for a substitution  $\sigma_i$ , where without loss of generality  $\sigma_i$  is restricted to variables in  $t_i$ . Define  $\sigma := \bigcup_{i=1}^k \sigma_i$ . Then  $t\sigma = f(t_1 \sigma_1, \dots, t_k \sigma_k) = f([o_1]_S, \dots, [o_k]_S) = [o]_S$ , where the last equality follows as  $m$  is morphic on  $o$ . Moreover, from the shape of  $\sigma_i$  and  $m_i$  it is not difficult to see that by construction the substitution  $\sigma$  and homomorphism  $m$  are related as claimed by the lemma.

Now for the inverse direction, suppose  $t\sigma = [o]_S$ . If  $t$  is a variable and thus  $\Delta(t)$  consists of a single unlabeled node, trivially  $\Delta(t) \geq_m S \downarrow o$  holds for  $m$  the homomorphism which maps the root of  $\Delta(t)$  to  $o$ . Observe that  $\sigma(t) = [o]_S = [m(\epsilon)]_S$ , which concludes the base case. For the inductive step suppose  $t = f(t_1, \dots, t_k)$ , hence  $S(o) = f(o_1, \dots, o_k)$ ,  $t_i \sigma = [o_i]_S$  ( $i = 1, \dots, k$ ) and thus by induction hypothesis  $\Delta(t_i) \geq_{m_i} [o_i]_S$  for homomorphisms  $m_i$ . Define the function  $m$  by  $m(\epsilon) := o$  and  $m(i \cdot p) := m_i(p)$  for all  $i = 1, \dots, k$  and positions  $i \cdot p$  of  $t$ . Observe that  $m$  is defined on all nodes of  $\Delta(t)$ . By definition of  $m$  one finally concludes the lemma, using the induction hypotheses together with the equalities  $\Delta(t)(i \cdot p) = \Delta(t_i)(p)$  for nodes  $i \cdot p$  ( $i = 1, \dots, k$ ) of  $\Delta(t)$ .  $\square$

## 4 Memoization and Sharing, Formally

To incorporate *memoization*, we make use of a *cache*  $C$  which stores results of intermediate functions calls. A *cache*  $C$  is modeled as a set of tuples  $(f(v_1, \dots, v_{\text{ar}(f)}), v)$ , where  $f \in \mathcal{F}$  and  $v_1, \dots, v_{\text{ar}(f)}$  as well as  $v$  are values. Figure 5 collects the *memoizing operational semantics* with respect to the program  $P = (\mathcal{F}, \mathcal{C}, \mathcal{R})$ . Here, a statement  $(C, t) \Downarrow_m (D, v)$  means that starting with a cache  $C$ , the term  $t$  reduces to the value  $v$  with updated cache  $D$ . The natural number  $m$  indicating the *cost* of this reduction. The definition is along the lines of the standard semantics (Figure 4), carrying the cache throughout the reduction of the given term. The last rule of Figure 4 is split into two rules **(Read)** and **(Update)**. The former performs a read from the cache, the latter the reduction in case the corresponding function call is not tabulated, updating the cache with the computed result. Notice that in the semantics, a read is attributed zero cost, whereas an update is accounted with a cost of one. Consequently the cost  $m$  in  $(C, t) \Downarrow_m (D, v)$  refers to the number of non-tabulated function applications.

**Lemma 4.1.** *We have  $(\emptyset, t) \Downarrow_m (C, v)$  for some  $m \in \mathbb{N}$  and cache  $C$  if and only if  $t \downarrow v$ .*

*Proof.* Call a cache  $C$  *proper* if  $(f(v_1, \dots, v_k), v) \in C$  implies  $f(v_1, \dots, v_k) \downarrow v$ . For the direction from right to left, we show the following stronger claim.

**Claim 4.2.** *Suppose  $t \downarrow v$  and let cache  $C_1$  be a proper cache. Then  $(C_1, t) \Downarrow_m (C_2, v)$  for some  $m$ .*

The proof is by induction on the deduction  $\Pi$  of the statement  $t \downarrow v$ .

1. Suppose that the last rule in  $\Pi$  has the form

$$\frac{t_i \downarrow v_i \quad f(p_1, \dots, p_k) \rightarrow r \in \mathcal{R} \quad p_i \sigma = v_i \quad r \sigma \downarrow v}{f(t_1, \dots, t_k) \downarrow v}$$

We consider the more involved case where at least one  $t_i$  is not a value. By induction hypothesis, we obtain proper caches  $D_0, \dots, D_k$  with  $D_0 = C_1$  and  $(D_{i-1}, t_i) \Downarrow_{m_i} (D_{i-1}, v_i)$ . By the rule (**Split**), it suffices to show  $(D_k, f(v_1, \dots, v_k)) \Downarrow_n (C_2, v)$  for  $C_2$  a proper cache. We distinguish two cases. Consider the case  $(f(v_1, \dots, v_k), u) \in D_k$  for some  $u$ . Using that  $f(v_1, \dots, v_k) \downarrow u$  implies  $v = u$  for orthogonal programs, we conclude the case by one application of rule (**Read**). Otherwise, we conclude by rule (**Update**) using the induction hypothesis on  $r \sigma \downarrow v$ . Note that the resulting cache is also in this case proper.

2. The final case follows directly from induction hypothesis, using the rule (**Constructor**).

For the direction from left to right we show the following stronger claim

**Claim 4.3.** *Suppose  $(C_1, t) \Downarrow_m (C_2, v)$  for a proper cache  $C_1$ . Then  $t \downarrow v$ .*

The proof is by induction on the deduction  $\Pi$  of the statement  $(C_1, t) \Downarrow_m (C_2, v)$

1. Suppose first that the last rule in  $\Pi$  is of the form

$$\frac{(C_{i-1}, t_i) \Downarrow_{m_i} (C_i, v_i) \quad (C_k, f(v_1, \dots, v_k)) \Downarrow_n (C_{k+1}, v)}{(C_0, f(t_1, \dots, t_k)) \Downarrow_m (C_{k+1}, v)}$$

By induction hypothesis, we see  $t_i \downarrow v_i$ , using also that the configurations  $c_i$  are all proper by the previous claim. As we also have  $f(v_1, \dots, v_k) \downarrow v$  by induction hypothesis, it follows that some rule  $f(p_1, \dots, p_k) \rightarrow r \in \mathcal{R}$  matches  $f(v_1, \dots, v_k)$ . Putting things together, we conclude by one application (**Function**).

2. The remaining cases where the last rule in  $\Pi$  is (**Constructor**), (**Read**) or (**Update**) follow either from the assumption that  $C_1$  is proper, or from induction hypothesis using that by the previous claim the intermediate and resulting caches are all proper.

□

The lemma confirms that the call-by-value semantics of Section 3 is correctly implemented by the memoizing semantics. To tame the growth rate of values, we define *small-step semantics* corresponding to the memoizing semantics, facilitating sharing of common sub-expressions.

**Small-Step Semantics with Memoization and Sharing** To incorporate sharing, we extend the pairs  $(C, t)$  by a *heap*, and allow *references* to the heap both in terms and in caches. Let  $\text{Loc}$  denote a countably infinite set of *locations*. We overload the notion of *value*, and define *expressions*  $e$  and (*evaluation*) *contexts*  $E$  according to the following grammar:

$$\begin{aligned} v &:= \ell \mid \mathbf{c}(v_1, \dots, v_k); \\ e &:= \ell \mid \langle \mathbf{f}(\ell_1, \dots, \ell_k), e \rangle \mid \mathbf{f}(e_1, \dots, e_k) \mid \mathbf{c}(e_1, \dots, e_k); \\ E &:= \square \mid \langle \mathbf{f}(\ell_1, \dots, \ell_k), E \rangle \mid \mathbf{f}(\ell_1, \dots, \ell_{i-1}, E, e_{i+1}, \dots, e_k) \mid \mathbf{c}(\ell_1, \dots, \ell_{i-1}, E, e_{i+1}, \dots, e_k). \end{aligned}$$

Here,  $\ell_1, \dots, \ell_k, \ell \in \text{Loc}$ ,  $\mathbf{f} \in \mathcal{F}$  and  $\mathbf{c} \in \mathcal{C}$  are  $k$ -ary symbols. An expression is a term including references to values that will be stored on the heap. The additional construct  $\langle \mathbf{f}(\ell_1, \dots, \ell_k), e \rangle$  indicates that the partially evaluated expression  $e$  descends from a call  $\mathbf{f}(v_1, \dots, v_k)$ , with arguments  $v_i$  stored at location  $\ell_i$  on the heap. A context  $E$  is an expression with a unique *hole*, denoted as  $\square$ , where all sub-expression to the left of the hole are references pointing to values. This syntactic restriction is used to implement a *left-to-right, call-by-value* evaluation order. We denote by  $E[e]$  the expression obtained by replacing the hole in  $E$  by  $e$ .

A *configuration* is a triple  $(D, H, e)$  consisting of a *cache*  $D$ , *heap*  $H$  and expression  $e$ . Unlike before, the cache  $D$  consists of pairs of the form  $(\mathbf{f}(\ell_1, \dots, \ell_k), \ell)$  where instead of values, we store references  $\ell_1, \dots, \ell_k, \ell$  pointing to the heap. The heap  $H$  is represented as a (multi-rooted) term graph  $H$  with nodes in  $\text{Loc}$  and constructors  $\mathcal{C}$  as labels. If  $\ell$  is a node of  $H$ , then we say that  $H$  stores at location  $\ell$  the value  $[\ell]_H$  obtained by unfolding  $H$  starting from location  $\ell$ . We keep the heap in a *maximally shared* form, that is,  $H(\ell_a) = \mathbf{c}(\ell_1, \dots, \ell_k) = H(\ell_b)$  implies  $\ell_a = \ell_b$  for two locations  $\ell_a, \ell_b$  of  $H$ . Thus crucially, values are stored once only, by the following lemma.

**Lemma 4.4.** *Let  $H$  be a maximally shared heap with locations  $\ell_1, \ell_2$ . If  $[\ell_1]_H = [\ell_2]_H$  then  $\ell_1 = \ell_2$ .*

The operation  $\text{merge}(H, \mathbf{c}(\ell_1, \dots, \ell_k))$ , defined as follows, is used to extend the heap  $H$  with a constructor  $\mathbf{c}$  whose arguments point to  $\ell_1, \dots, \ell_k$ , retaining maximal sharing. Let  $\ell_f$  be the first location not occurring in the nodes  $N$  of  $H$  (with respect to an arbitrary, but fixed enumeration on  $\text{Loc}$ ). For  $\ell_1, \dots, \ell_k \in N$  we define

$$\text{merge}(H, \mathbf{c}(\ell_1, \dots, \ell_k)) := \begin{cases} (H, \ell) & \text{if } H(\ell) = \mathbf{c}(\ell_1, \dots, \ell_k), \\ (H \cup \{\ell_f \mapsto \mathbf{c}(\ell_1, \dots, \ell_k)\}, \ell_f) & \text{otherwise.} \end{cases}$$

Observe that the first clause is unambiguous on maximally shared heaps.

Figure 6 collects the small step semantics with respect to a program  $P = (\mathcal{F}, \mathcal{C}, \mathcal{R})$ . We use  $\rightarrow_{\text{rsm}}$  to abbreviate the relation  $\rightarrow_{\mathbf{r}} \cup \rightarrow_{\mathbf{s}} \cup \rightarrow_{\mathbf{m}}$  and likewise we abbreviate  $\rightarrow_{\mathbf{R}} \cup \rightarrow_{\text{rsm}}$  by  $\rightarrow_{\text{Rrsm}}$ . Furthermore, we define  $\rightarrow_{\text{R/rsm}} := \rightarrow_{\text{rsm}}^* \cdot \rightarrow_{\mathbf{R}} \cdot \rightarrow_{\text{rsm}}^*$ . Hence the *m-fold composition*  $\rightarrow_{\text{R/rsm}}^m$  corresponds to a  $\rightarrow_{\text{Rrsm}}$ -reduction with precisely  $m$  applications of  $\rightarrow_{\mathbf{R}}$ . Throughout the following, we are interested in reductions over *well-formed* configurations:

**Definition 4.5.** *A configuration  $(D, H, e)$  is well-formed if the following conditions hold.*

1. *The heap  $H$  is maximally shared.*
2. *The cache  $D$  is a function, and compatible with  $e$ , that is, if  $\langle \mathbf{f}(\ell_1, \dots, \ell_k), e' \rangle$  occurs as a sub-expression in  $e$ , then  $(\mathbf{f}(\ell_1, \dots, \ell_k), \ell) \notin D$  for any  $\ell$ .*

$$\begin{array}{c}
\frac{(\mathbf{f}(\ell_1, \dots, \ell_k), \ell) \notin D \quad \mathbf{f}(p_1, \dots, p_k) \rightarrow r \in \mathcal{R} \quad T := \Delta(\mathbf{f}(p_1, \dots, p_k))}{T \gg_m \mathbf{f}(H \setminus \ell_1, \dots, H \setminus \ell_k) \quad \sigma_m := \{x \mapsto m(\ell_x) \mid \ell_x \in \text{Loc}, T(\ell_x) = x \in \mathcal{V}\}} \text{ (apply)} \\
\frac{(\mathbf{f}(\ell_1, \dots, \ell_k), \ell) \in D}{(D, H, E[\mathbf{f}(\ell_1, \dots, \ell_k)]) \rightarrow_r (D, H, E[\langle \mathbf{f}(\ell_1, \dots, \ell_k), r \sigma_m \rangle])} \text{ (read)} \\
\frac{}{(D, H, E[\langle \mathbf{f}(\ell_1, \dots, \ell_k), \ell \rangle]) \rightarrow_s (D \cup \{(\mathbf{f}(\ell_1, \dots, \ell_k), \ell)\}, H, E[\ell])} \text{ (store)} \\
\frac{(H', \ell) = \text{merge}(H, \mathbf{c}(\ell_1, \dots, \ell_k))}{(D, H, E[\mathbf{c}(\ell_1, \dots, \ell_k)]) \rightarrow_m (D, H', E[\ell])} \text{ (merge)}
\end{array}$$

Figure 6: Small Step Semantics with Memoization and Sharing for Program  $(\mathcal{F}, \mathcal{C}, \mathcal{R})$ .

3. *The configuration contains no dangling locations, that is,  $H(\ell)$  is defined for each location  $\ell$  occurring in  $D$  and  $e$ .*

**Lemma 4.6. 1.** *If  $(D, H, E[e])$  is well-formed then so is  $(D, H, e)$ .*

2. *If  $(D_1, H_1, e_1) \rightarrow_{\text{Rrsm}} (D_2, H_2, e_2)$  and  $(D_1, H_1, e_1)$  is well-formed then so is  $(D_2, H_2, e_2)$ .*

*Proof.* It is not difficult to see that Assertion 1 holds. To see that Assertion 2 holds, fix a well-formed configuration  $(D_1, H_1, e_1)$  and suppose  $(D_1, H_1, e_1) \rightarrow_{\text{Rrsm}} (D_2, H_2, e_2)$ . We check that  $(D_2, H_2, e_2)$  is well-formed by case analysis on  $\rightarrow_{\text{Rrsm}}$ .

1. *The heap  $H_2$  is maximally shared:* As only the relation  $\rightarrow_m$  modifies the heap, it suffices to consider the case  $(D_1, H_1, e_1) \rightarrow_m (D_2, H_2, e_2)$ . Then  $(H_2, \ell) = \text{merge}(H_1, \mathbf{c}(\ell_1, \dots, \ell_k))$  for some location  $\ell$ , and the property follows as merge preserves maximal sharing.
2. *The cache  $D_2$  is a function:* It suffices to consider the rules  $\rightarrow_s$ . As immediate consequence of compatibility of  $D_1$  with  $e_1$  it follows that  $D_2$  is a function.
3. *The cache  $D_2$  is compatible with  $e_2$ :* Only the rules  $\rightarrow_r$  and  $\rightarrow_s$  potentially contradict compatibility. In the former case, the side conditions ensure that  $e_2$  and  $D_2$  are compatible, in the latter case compatibility follows trivially from compatibility of  $D_1$  with  $e_1$ .
4. *No dangling references:* Observe that only rule  $\rightarrow_m$  introduces a fresh location. The merge operations guarantees that this location occurs in the heap  $H_2$ .

□

From now on, if not mentioned otherwise we will suppose that configurations are well-formed, tacitly employing Lemma 4.6.

It is now time to show that the model of computation we have just introduced fits our needs, namely that it faithfully simulates big-step semantics as in Figure 5 (itself a correct implementation of call-by-value evaluation from Section 3). This is proven by first showing how big-step semantics can be *simulated* by small-step semantics, later proving that the latter is in fact *deterministic*.

In the following, we denote by  $[e]_H$  the term obtained from  $e$  by following pointers to the heap, ignoring the annotations  $\langle f(\ell_1, \dots, \ell_k), \cdot \rangle$ . Formally, we define

$$[e]_H := \begin{cases} f([e_1]_H, \dots, [e_k]_H) & \text{if } e = f(e_1, \dots, e_k), \\ [e']_H & \text{if } e = \langle f(\ell_1, \dots, \ell_k), e' \rangle. \end{cases}$$

Observe that this definition is well-defined as long as  $H$  contains all locations occurring in  $e$ . Likewise, we set  $[D]_H := \{([e]_H, [\ell]_H) \mid (e, \ell) \in D\}$ .

Our simulation result relies on the following two auxiliary lemmas concerning heaps. The first is based on the observation that in  $\rightarrow_{\text{Rrsm}}$ -reductions, the heap is monotonically increasing.

**Lemma 4.7.** *If  $(D_1, H_1, e_1) \rightarrow_{\text{Rrsm}} (D_2, H_2, e_2)$  then the following properties hold:*

1.  $[\ell]_{H_2} = [\ell]_{H_1}$  for every location  $\ell$  of  $H_1$ ;
2.  $[D_1]_{H_2} = [D_1]_{H_1}$  and  $[e_1]_{H_2} = [e_1]_{H_1}$ .

*Proof.* As for any other step the heap remains untouched, the only non-trivial case is  $(D_1, H_1, E[\mathbf{c}(\ell_1, \dots, \ell_k)]) \rightarrow_{\text{m}} (D_1, H_2, E[\ell])$  with  $(H_2, \ell) = \text{merge}(H_1, \mathbf{c}(\ell_1, \dots, \ell_k))$ . Observe that by definition of  $\text{merge}$ ,  $H_2(\ell) = H_1(\ell)$  for every  $\ell \in N_{H_1}$ . From this Assertion 1 is easy to establish. Assertion 2 follows then by standard inductions on  $D_1$  and  $E$ , respectively.  $\square$

**Lemma 4.8.** *Let  $(D, H, E[v])$  be a configuration for a value  $v$ . Then  $(D, H, E[v]) \rightarrow_{\text{m}}^* (D, H', E[\ell])$  with  $[\ell]_{H'} = [v]_H$ .*

*Proof.* Note that by assumptions,  $e$  consist only of constructors and locations. We proof the lemma by induction on the number of constructor symbols in  $e$ . In the base case,  $e = \ell$  and the lemma trivially holds. For the inductive step, it is not difficult to see that  $e = E'[\mathbf{c}(\ell_1, \dots, \ell_k)]$  for some evaluation context  $E'$ , and hence  $(D, H, E[e]) \rightarrow_{\text{m}} (D, H', E[E'[\ell]])$ , where  $(H', \ell) = \text{merge}(H, \mathbf{c}(\ell_1, \dots, \ell_k))$ . Using that  $[\ell]_{H'} = [\mathbf{c}(\ell_1, \dots, \ell_k)]_{H'}$  by definition of  $\text{merge}$  and Lemma 4.7(2) we conclude  $[E[E'[\ell]]]_{H'} = [E[e]]_H$ . We complete this derivation to the desired form, by induction hypothesis.  $\square$

An *initial configuration* is a configuration of the form  $(\emptyset, H, e)$  with  $H$  a maximally shared heap and  $e = f(v_1, \dots, v_k)$  an expression unfolding to a function call. Notice that the arguments  $v_1, \dots, v_k$  are allowed to contain references to the heap  $H$ .

**Lemma 4.9** (Simulation). *Let  $(\emptyset, H, e)$  be an initial configuration. If  $(\emptyset, [e]_H) \Downarrow_m (C, v)$  holds for  $m \geq 1$  then  $(\emptyset, H, e) \rightarrow_{\text{R/rsm}}^m (D, G, \ell)$  for a location  $\ell$  in  $G$  with  $[\ell]_G = v$ .*

*Proof.* Call a configuration  $(D, H, e)$  *proper* if it is well-formed and  $e$  does not contain a sub-expression  $\langle f(v_1, \dots, v_k), e' \rangle$ . We show the following claim:

**Claim 4.10.** *For every proper configuration  $(D_1, H_1, e_1)$ ,  $([D_1]_{H_1}, [e_1]_{H_1}) \Downarrow_m (C_2, v)$  implies  $(D_1, H_1, e_1) \rightarrow_{\text{rsm}}^* \cdot \rightarrow_{\text{R/rsm}}^m (D_2, H_2, \ell)$  with  $([D_2]_{H_2}, [\ell]_{H_2}) = (C_2, v)$ .*

Observe that  $\rightarrow_{\text{rsm}}^* \cdot \rightarrow_{\text{R/rsm}}^m = \rightarrow_{\text{R/rsm}}^m$  whenever  $m > 0$ . Since an initial configuration is trivially proper, the lemma follows from the claim.

To prove the claim, abbreviate the relation  $\rightarrow_{\text{rsm}}^* \cdot \rightarrow_{\text{R/rsm}}^m$  by  $\rightarrow^m$  for all  $m \in \mathbb{N}$ . Below, we tacitly employ  $\rightarrow^{m_1} \cdot \rightarrow^{m_2} = \rightarrow^{m_1+m_2}$  for all  $m_1, m_2 \in \mathbb{N}$ . The proof is by induction on the deduction  $\Pi$  of the statement  $([D_1]_H, [e]_{H_1}) \Downarrow_m (C, v)$ .

1. Suppose that the last rule in  $\Pi$  has the form

$$\frac{\mathbf{c} \in \mathcal{C} \quad (C_{i-1}, t_i) \Downarrow_{m_i} (C_i, v_i) \quad m = \sum_{i=1}^k m_i}{(C_0, \mathbf{c}(t_1, \dots, t_k)) \Downarrow_m (C_k, \mathbf{c}(v_1, \dots, v_k))}$$

Fix a proper configuration  $(D_0, H_0, e_0)$  unfolding to  $(C_0, \mathbf{c}(t_1, \dots, t_k))$ . Under these assumptions, either  $e_0$  is a location or  $e_0 = \mathbf{c}(e_1, \dots, e_k)$ . The former case is trivial, as then  $t$  is a value and thus  $m = 0$ . Hence suppose  $e_0 = \mathbf{c}(e_1, \dots, e_k)$ . We first show that for all  $i \leq k$ ,

$$(D_0, H_0, \mathbf{c}(e_1, \dots, e_k)) \rightarrow_{\sum_{j=1}^i m_j} (D_i, H_i, \mathbf{c}(\ell_1, \dots, \ell_i, e_{i+1}, \dots, e_k)), \quad (\dagger)$$

for a configuration  $(D_i, H_i, \mathbf{c}(\ell_1, \dots, \ell_i, e_{i+1}, \dots, e_k))$  which unfolds to the configuration  $(C_i, \mathbf{c}(v_1, \dots, v_i, t_{i+1}, \dots, t_k))$ . The proof is by induction on  $i$ , we consider the step from  $i$  to  $i+1$ . Induction hypothesis yields a well-formed configuration  $(D_i, H_i, E[e_{i+1}])$  for  $E = \mathbf{c}(\ell_1, \dots, \ell_i, \square, e_{i+2}, \dots, e_k)$  reachable by a Derivation  $(\dagger)$ . As the configuration  $(D_i, H_i, e_{i+1})$  unfolds to  $(C_i, t_{i+1})$ , the induction hypothesis of the claim on the assumption  $(C_i, t_{i+1}) \Downarrow_{m_{i+1}} (C_{i+1}, v_{i+1})$  yields  $(D_i, H_i, e_{i+1}) \rightarrow^{m_{i+1}} (D_{i+1}, H_{i+1}, \ell_{i+1})$  where the resulting configuration unfolds to  $(C_{i+1}, v_{i+1})$ . As a consequence, its not difficult to see that also  $(D_i, H_i, E[e_{i+1}]) \rightarrow^{m_{i+1}} (D_{i+1}, H_{i+1}, E[\ell_{i+1}])$  holds. Since  $[\ell_{i+1}]_{H_{i+1}} = v_{i+1}$  and  $[E[e_{i+1}]]_{H_i} = \mathbf{c}(v_1, \dots, v_i, t_{i+1}, t_{i+2}, \dots, t_k)$ , using Lemma 4.7(2) on the last equality it is not difficult to see that  $[E[\ell_{i+1}]]_{H_{i+1}} = \mathbf{c}(v_1, \dots, v_i, v_{i+1}, t_{i+2}, \dots, t_k)$ . As we already observed  $[D_{i+1}]_{H_{i+1}} = C_{i+1}$ , we conclude the Reduction  $(\dagger)$ .

In total, we thus obtain a reduction  $(D_0, H_0, \mathbf{c}(e_1, \dots, e_k)) \rightarrow^m (D_k, H_k, \mathbf{c}(\ell_1, \dots, \ell_k))$  where  $m = \sum_{i=1}^k m_i$  and  $(D_k, H_k, \mathbf{c}(\ell_1, \dots, \ell_k))$  is a well-formed, in fact proper, configuration which unfolds to  $(C_k, \mathbf{c}(v_1, \dots, v_k))$ . Employing  $\rightarrow^m \cdot \rightarrow_{\mathbf{m}}^* = \rightarrow^m$  we conclude the case with Lemma 4.8.

2. Suppose that the last rule in  $\Pi$  has the form

$$\frac{(C_{i-1}, t_i) \Downarrow_{m_i} (C_i, v_i) \quad (C_k, \mathbf{f}(v_1, \dots, v_k)) \Downarrow_n (C_{k+1}, v) \quad m = n + \sum_{i=1}^k m_i}{(C_0, \mathbf{f}(t_1, \dots, t_k)) \Downarrow_m (C_{k+1}, v)}$$

Fix a proper configuration  $(D_0, H_0, e_0)$  unfolding to  $(C_0, \mathbf{f}(t_1, \dots, t_k))$ . By induction on  $k$ , exactly as in the previous case, we obtain a proper configuration  $(D_k, H_k, \mathbf{f}(\ell_1, \dots, \ell_k))$  unfolding to  $(C_k, \mathbf{f}(v_1, \dots, v_k))$  with

$$(D_0, H_0, e_0) \rightarrow_{\sum_{i=1}^k m_i} (D_k, H_k, \mathbf{f}(\ell_1, \dots, \ell_k)).$$

The induction hypothesis also yields configuration  $(D_{k+1}, H_{k+1}, \ell)$  unfolding to  $(C_{k+1}, v)$  with  $(D_k, H_k, \mathbf{f}(\ell_1, \dots, \ell_k)) \rightarrow^n (D_{k+1}, H_{k+1}, \ell)$ . Summing up we conclude the case.

3. Suppose that the last rule in  $\Pi$  has the form

$$\frac{(\mathbf{f}(v_1, \dots, v_k), u) \in \mathcal{C}}{(C, \mathbf{f}(v_1, \dots, v_k)) \Downarrow_0 (C, v)}$$

Consider a proper configuration  $(D, H, e)$  that unfolds to  $(C, \mathbf{f}(v_1, \dots, v_k))$ . Then  $e = \mathbf{f}(e_1, \dots, e_k)$ , and using  $k$  applications of Lemma 4.8 we construct a reduction

$$(D, H, \mathbf{f}(e_1, \dots, e_k)) \rightarrow_{\mathbf{m}}^* (D, H_1, \mathbf{f}(\ell_1, e_2, \dots, e_k)) \rightarrow_{\mathbf{m}}^* \dots \rightarrow_{\mathbf{m}}^* (D, H_k, \mathbf{f}(\ell_1, \dots, \ell_k)),$$

with  $(D, H_k, f(\ell_1, \dots, \ell_k))$  unfolding to  $(C, f(v_1, \dots, v_k))$ . Lemma 4.4 and the assumption on  $C = [D]_{H_k}$  implies that there exists a *unique* pair  $(f(\ell_1, \dots, \ell_k), \ell) \in D$  with  $[f(\ell_1, \dots, \ell_k)]_{H_k} = f(v_1, \dots, v_k)$  and  $[\ell]_{H_k} = v$ . Thus overall

$$(D, H, e) = (D, H, f(e_1, \dots, e_k)) \rightarrow_m^* (D, H_k, f(\ell_1, \dots, \ell_k)) \rightarrow_r (D, H_k, \ell),$$

where  $(D, H_k, \ell)$  unfolds to  $(C, v)$ . Using  $\rightarrow_m^* \cdot \rightarrow_r \subseteq \rightarrow^0$  we conclude the case.

4. Finally, suppose that the last rule in  $\Pi$  has the form

$$\frac{(f(v_1, \dots, v_k), v) \notin C \quad f(l_1, \dots, l_k) \rightarrow r \in \mathcal{R} \quad \forall i. l_i \sigma = v_i \quad (C, r\sigma) \Downarrow_m (D, v)}{(C, f(v_1, \dots, v_k)) \Downarrow_{m+1} (D \cup \{(f(v_1, \dots, v_k), v)\}, v)}$$

Fix a proper configuration  $(D, H, e)$  that unfolds to  $(C, f(v_1, \dots, v_k))$ , in particular  $e = f(e_1, \dots, e_k)$ . As above, we see  $(D, H, e) \rightarrow_m^* (D, H_k, f(\ell_1, \dots, \ell_k))$  for a configuration  $(D, H_k, f(\ell_1, \dots, \ell_k))$  also unfolding to  $(C, f(v_1, \dots, v_k))$ . As  $(f(v_1, \dots, v_k), v) \notin C$ , we have  $(f(\ell_1, \dots, \ell_k), \ell) \notin D_k$  for any location  $\ell$ , by Lemma 4.4. Since Proposition 3.9 on the assumption yields  $\Delta(f(l_1, \dots, l_k)) \geq_m f(H \upharpoonright_{l_1}, \dots, H \upharpoonright_{l_k})$  for a matching morphism  $m$ , in total we obtain

$$(D, H, e) \rightarrow_m^* (D, H, f(\ell_1, \dots, \ell_k)) \rightarrow_r (D, H_k, \langle f(\ell_1, \dots, \ell_k), r\sigma_m \rangle).$$

Note that by Proposition 3.9 the substitution  $\sigma$  and induced substitution  $\sigma_m$  satisfy  $\sigma(x) = [\sigma_m(x)]_{H_k}$  for all variables  $x$  in  $t$ . Hence by a standard induction on  $r$ ,  $[r\sigma_m]_{H_k} = r\sigma$  follows. We conclude that  $(D, H_k, \langle f(\ell_1, \dots, \ell_k), r\sigma_m \rangle)$  unfolds to  $(C, r\sigma)$ . Thus the induction hypothesis yields a well-formed configuration  $(D', G, \ell)$  unfolding to  $(C', v)$  with  $(D, H_k, r\sigma_m) \rightarrow^m (D', G, \ell)$ . Thus

$$\begin{aligned} (D, H_k, \langle f(\ell_1, \dots, \ell_k), r\sigma_m \rangle) &\rightarrow^m (D', G, \langle f(\ell_1, \dots, \ell_k), \ell \rangle) \\ &\rightarrow_s (D' \cup \{(f(\ell_1, \dots, \ell_k), \ell)\}, G, \ell). \end{aligned}$$

Using that  $[\ell]_G = v$  and  $[f(\ell_1, \dots, \ell_k)]_{H_k} = f(v_1, \dots, v_k)$ , Lemma 4.7 yields

$$[D' \cup \{(f(\ell_1, \dots, \ell_k), \ell)\}]_G = C' \cup \{(f(v_1, \dots, v_k), v)\}.$$

Putting things together, employing  $\rightarrow_m^* \cdot \rightarrow_r \subseteq \rightarrow^1$  and  $\rightarrow^m \cdot \rightarrow_s = \rightarrow^m$  we conclude  $(D, H, e) \rightarrow^{m+1} (D' \cup \{(f(\ell_1, \dots, \ell_k), \ell)\}, G, \ell)$ , where the resulting configuration unfolds to  $(C' \cup \{(f(v_1, \dots, v_k), v)\}, v)$ . We conclude this final case. □

The next lemma shows that the established simulation is *unique*, that is, there is exactly one derivation  $(\emptyset, H, e) \rightarrow_{\mathbb{R}/\text{rsm}}^m (D, G, \ell)$ . Here, a relation  $\rightarrow$  is called *deterministic on a set*  $A$  if  $b_1 \leftarrow a \rightarrow b_2$  implies  $b_1 = b_2$  for all  $a \in A$ .

**Lemma 4.11** (Determinism).

1. The relations  $\rightarrow_{\mathbb{R}}$ ,  $\rightarrow_r$ ,  $\rightarrow_s$  and  $\rightarrow_m$  are deterministic on well-formed configurations.
2. The relation  $\rightarrow_{\text{Rrsm}}$  is deterministic on well-formed configurations.

*Proof.* For Assertion 1, fix  $\rightarrow_r \in \{\rightarrow_{\mathbf{R}}, \rightarrow_{\mathbf{r}}, \rightarrow_{\mathbf{s}}, \rightarrow_{\mathbf{m}}\}$ . Let  $(D, H, e)$  be a well-formed configuration. We show that any *peak*  $(D_1, H_1, e_1) \xrightarrow{r_1} (D, H, e) \xrightarrow{r_2} (D_2, H_2, e_2)$  is *trivial*, i.e.  $(D_1, H_1, e_1) = (D_2, H_2, e_2)$ . Observe that independent on  $\rightarrow_r$ , the evaluation context  $E$  in the corresponding rule is unique. From this, we conclude the assertion by case analysis on  $\rightarrow_r$ . The non-trivial cases are  $\rightarrow_r = \rightarrow_{\mathbf{R}}$  and  $\rightarrow_r = \rightarrow_{\mathbf{r}}$ . In the former case, we conclude using that rules in  $\mathcal{R}$  are non-overlapping, tacitly employing Proposition 3.9. The latter case we conclude using that  $D$  is well-formed.

Finally, for Assertion 2 consider a peak  $(D_1, H_1, e_1) \xrightarrow{r_1} (D, H, e) \xrightarrow{r_2} (D_2, H_2, e_2) \xrightarrow{r_3}$ ,  $\rightarrow_{r_2} \in \{\rightarrow_{\mathbf{R}}, \rightarrow_{\mathbf{r}}, \rightarrow_{\mathbf{s}}, \rightarrow_{\mathbf{m}}\}$ . We show that this peak is trivial by induction on the expression  $e$ . By the previous assertion, it suffices to consider only the case  $\rightarrow_{r_1} \neq \rightarrow_{r_2}$ .

The base case constitutes of the cases (i)  $e = f(\ell_1, \dots, \ell_k)$ , (ii)  $e = \langle f(\ell_1, \dots, \ell_k), \ell \rangle$  and (iii)  $e = \mathbf{c}(\ell_1, \dots, \ell_k)$ . The only potential peak can occurs in case (i) between relations  $\rightarrow_{\mathbf{R}}$  and  $\rightarrow_{\mathbf{r}}$ . Here, a non-trivial peak is prohibited by the pre-conditions put on  $H$ . For the inductive step, we consider a peak

$$(D_1, H_1, E[e'_1]) \xrightarrow{r_1} (D, H, E[e']) \xrightarrow{r_2} (D_2, H_2, E[e'_2]),$$

where  $e = E[e']$  for a context  $E$ . As we thus have a peak  $(D_1, H_1, e'_1) \xrightarrow{r_1} (D, H, e') \xrightarrow{r_2} (D_2, H_2, e'_2)$ , which by induction hypothesis is trivial, we conclude the assertion.  $\square$

**Theorem 4.12.** *Suppose  $(\emptyset, f(v_1, \dots, v_k)) \Downarrow_m (C, v)$  holds for a reducible term  $f(v_1, \dots, v_k)$ . Then for each initial configuration  $(\emptyset, H, e)$  with  $[e]_H = f(v_1, \dots, v_k)$ , there exists a unique sequence  $(\emptyset, H, e) \xrightarrow{m}_{\mathbf{R}/\mathbf{rsm}} (D, G, \ell)$  for a location  $\ell$  in  $G$  with  $[\ell]_G = v$ .*

*Proof.* As  $f(v_1, \dots, v_k)$  is reducible, it follows that  $m \geq 1$ . Hence the theorem follows from Lemma 4.9 and Lemma 4.11.  $\square$

**Invariance** Theorem 4.12 tells us that a term-based semantics (in which sharing is *not* exploited) can be simulated step-by-step by another, more sophisticated, graph-based semantics. The latter's advantage is that each computation step does not require copying, and thus does not increase the size of the underlying configuration too much. This is the key observation towards *invariance*: the number of reduction step is a sensible cost model from a complexity-theoretic perspective. Precisely this will be proven in the remaining of the section.

Define the *size*  $|e|$  of an expression recursively by  $|\ell| := 1$ ,  $|f(e_1, \dots, e_k)| := 1 + \sum_{i=1}^k |e_i|$  and  $|\langle f(\ell_1, \dots, \ell_k), e \rangle| := 1 + |e|$ . In correspondence we define the *weight*  $\text{wt}(e)$  by ignoring locations, i.e.  $\text{wt}(\ell) := 0$ . Recall that a reduction  $(D_1, H_1, e_1) \xrightarrow{m}_{\mathbf{R}/\mathbf{rsm}} (D_2, H_2, e_2)$  consists of  $m$  applications of  $\rightarrow_{\mathbf{R}}$ , all interleaved by  $\rightarrow_{\mathbf{rsm}}$ -reductions. As a first step, we thus estimate the overall length of the reduction  $(D_1, H_1, e_1) \xrightarrow{m}_{\mathbf{R}/\mathbf{rsm}} (D_2, H_2, e_2)$  in  $m$  and the size of  $e_1$ . Set  $\Delta := \max\{|r| \mid l \rightarrow r \in \mathcal{R}\}$ . The following serves as an intermediate lemma.

**Lemma 4.13.** *The following properties hold:*

1. *If  $(D_1, H_1, e_1) \xrightarrow{\mathbf{rsm}} (D_2, H_2, e_2)$  then  $\text{wt}(e_2) < \text{wt}(e_1)$ .*
2. *If  $(D_1, H_1, e_1) \xrightarrow{\mathbf{R}} (D_2, H_2, e_2)$  then  $\text{wt}(e_2) \leq \text{wt}(e_1) + \Delta$ .*

*Proof.* The first assertion follows by case analysis on  $\rightarrow_{\mathbf{rsm}}$ . For the second, suppose  $(D_1, H_1, e_1) \xrightarrow{\mathbf{R}} (D_2, H_2, e_2)$  where  $e_1 = E[f(\ell_1, \dots, \ell_k)]$  and  $e_2 = E[\langle f(\ell_1, \dots, \ell_k), r\sigma_m \rangle]$



for a rule  $f(\ell_1, \dots, \ell_k) \rightarrow r \in \mathcal{R}$ . Observe that since the substitution  $\sigma_m$  replaces variables by locations,  $\Delta \geq |r| = |r\sigma_m| \geq \text{wt}(r\sigma_m)$  holds. Consequently,

$$\text{wt}(f(\ell_1, \dots, \ell_k)) + \Delta \geq 1 + \text{wt}(r\sigma_m) = \text{wt}(\langle f(\ell_1, \dots, \ell_k), r\sigma_m \rangle).$$

Then the assertion follows by a standard induction on  $E$ .  $\square$

Then essentially an application of the *weight gap principle* [17], a form of *amortized cost analysis*, binds the overall length of an  $\rightarrow_{\mathbb{R}/\text{rsm}}^m$ -reduction suitably.

**Lemma 4.14.** *If  $(D_1, H_1, e_1) \rightarrow_{\mathbb{R}/\text{rsm}}^m (D_2, H_2, e_2)$  then  $(D_1, H_1, e_1) \rightarrow_{\mathbb{R}\text{rsm}}^n (D_2, H_2, e_2)$  for  $n \leq (1 + \Delta) \cdot m + \text{wt}(e)$  and  $\Delta \in \mathbb{N}$  a constant depending only on  $\mathbb{P}$ .*

*Proof.* For a configuration  $c = (D, H, e)$  define  $\text{wt}(c) := \text{wt}(e)$  and let  $\Delta$  be defined as in Lemma 4.13. Consider  $(D_1, H_1, e_1) \rightarrow_{\mathbb{R}/\text{rsm}}^m (D_2, H_2, e_2)$  which can be written as a reduction

$$(D_1, H_1, e_1) = c_0 \xrightarrow{\text{rsm}}^{n_0} d_0 \rightarrow_{\mathbb{R}} c_1 \xrightarrow{\text{rsm}}^{n_1} d_1 \rightarrow_{\mathbb{R}} \dots \xrightarrow{\text{rsm}}^{n_m} d_m, \quad (\ddagger)$$

of length  $n := m + \sum_{i=0}^k n_i$ . Lemma 4.13 yields (i)  $n_i \leq \text{wt}(c_i) - \text{wt}(d_i)$  for all  $0 \leq i \leq m$ ; and (ii)  $\text{wt}(c_{i+1}) - \text{wt}(d_i) \leq \Delta$  for all  $0 \leq i < m$ . Hence overall, the Reduction  $(\ddagger)$  is of length

$$\begin{aligned} n &\leq m + (\text{wt}(c_0) - \text{wt}(d_0)) + \dots + (\text{wt}(c_m) - \text{wt}(d_m)) \\ &= m + \text{wt}(c_0) + (\text{wt}(c_1) - \text{wt}(d_0)) + \dots + (\text{wt}(c_m) - \text{wt}(d_{m-1})) - \text{wt}(d_m) \\ &\leq m + \text{wt}(c_0) + m \cdot \Delta \\ &= (1 + \Delta) \cdot m + \text{wt}(e). \end{aligned}$$

The lemma follows.  $\square$

Define the size of a configuration  $|(D, H, e)|$  as the sum of the sizes of its components. Here, the size  $|D|$  of a cache  $D$  is defined as its cardinality, similar, the size  $|H|$  of a heap is defined as the cardinality of its set of nodes. Notice that a configuration  $(D, H, e)$  can be straight forward encoded within logarithmic space-overhead as a string  $\lceil (D, H, e) \rceil$ , i.e. the length of the string  $\lceil (D, H, e) \rceil$  is bounded by a function in  $O(\log(n) \cdot n)$  in  $|(D, H, e)|$ , using constants to encode symbols and an encoding of locations logarithmic in  $|H|$ . Crucially, a step in the small-step semantics increases the size of a configuration only by a constant.

**Lemma 4.15.** *If  $(D_1, H_1, e_1) \rightarrow_{\mathbb{R}\text{rsm}} (D_2, H_2, e_2)$  then  $|(D_2, H_2, e_2)| \leq |(D_1, H_1, e_1)| + \Delta$ .*

*Proof.* The lemma follows by case analysis on the rule applied in  $(D_1, H_1, e_1) \rightarrow_{\mathbb{R}\text{rsm}} (D_2, H_2, e_2)$ , using  $1 \leq \Delta$ .  $\square$

**Theorem 4.16.** *There exists a polynomial  $p : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$  such that for every initial configuration  $(\emptyset, H_1, e_1)$ , a configuration  $(D_2, H_2, e_2)$  with  $(\emptyset, H_1, e_1) \rightarrow_{\mathbb{R}/\text{rsm}}^m (D_2, H_2, e_2)$  is computable from  $(\emptyset, H_1, e_1)$  in time  $p(|H_1| + |e_1|, m)$ .*

*Proof.* It is tedious, but not difficult to show that the function which implements a step  $c \rightarrow_{\mathbb{R}\text{rsm}} d$ , i.e. which maps  $\lceil c \rceil$  to  $\lceil d \rceil$ , is computable in polynomial time in  $\lceil c \rceil$ , and thus in the size  $|c|$  of the configuration  $c$ . Iterating this function at most  $n := (1 + \Delta) \cdot m + |(\emptyset, H_1, e_1)|$  times on input  $\lceil (\emptyset, H_1, e_1) \rceil$ , yields the desired result  $\lceil (D_2, H_2, e_2) \rceil$  by Lemma 4.14. Since each iteration increases the size of a configuration by at most the constant  $\Delta$  (Lemma 4.15), in particular the size of each intermediate configuration is bounded by a linear function in  $|(\emptyset, H_1, e_1)| = |H_1| + |e_1|$  and  $n$ , the theorem follows.  $\square$

Combining Theorem 4.12 and Theorem 4.16 we thus obtain the following.

**Theorem 4.17.** *There exists a polynomial  $p : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$  such that for  $(\emptyset, f(v_1, \dots, v_k)) \Downarrow_m (C, v)$ , the value  $v$  represented as DAG is computable from  $v_1, \dots, v_k$  in time  $p(\sum_{i=1}^k |v_i|, m)$ .*

Theorem 4.17 thus confirms that the cost  $m$  of a reduction  $(\emptyset, f(v_1, \dots, v_k)) \Downarrow_m (C, v)$  is a suitable cost measure. In other words, the *memoized runtime complexity* of a function  $f$ , relating input size  $n \in \mathbb{N}$  to the maximal cost  $m$  of evaluation  $f$  on arguments  $v_1, \dots, v_k$  of size up to  $n$ , i.e.  $(\emptyset, f(v_1, \dots, v_k)) \Downarrow_m (C, v)$  with  $\sum_{i=1}^k |v_i| \leq n$ , is an *invariant cost model*.

**Example 4.18** (Continued from Example 3.7). *Reconsider the program  $\text{PR}$  and the evaluation of a call  $\text{rabbits}(\mathbf{S}^n(\mathbf{0}))$  which results in the genealogical tree  $v_n$  of height  $n \in \mathbb{N}$  associated with Fibonacci's rabbit problem. Then one can show that  $\text{rabbits}(\mathbf{S}^n(\mathbf{0})) \Downarrow_m v_n$  with  $m \leq 2 \cdot n + 1$ . Crucially here, the two intermediate functions  $\mathbf{a}$  and  $\mathbf{b}$  defined by simultaneous recursion are called only on proper subterms of the input  $\mathbf{S}^n(\mathbf{0})$ , hence in particular the rules defining  $\mathbf{a}$  and  $\mathbf{b}$  respectively, are unfolded at most  $n$  times. As a consequence of the bound on  $m$  and Theorem 4.17 we obtain that the function `rabbits` from the introduction is polytime computable.*

**Remark 4.19.** *Strictly speaking, our DAG representation of a value  $v$ , viz the part of the final heap reachable from a corresponding location  $\ell$ , is not an encoding in the classical, complexity theoretic setting. Different computations resulting in the same value  $v$  can produce different DAG representations of  $v$ , however, these representations differ only in the naming of locations. Even though our encoding can be exponentially compact in comparison to a linear representation without sharing, it is not exponentially more succinct than a reasonable encoding for graphs (e.g. representations as circuits, see Papadimitriou [23]). In such succinct encodings not even equality can be decided in polynomial time. Our form of representation does clearly not fall into this category. In particular, in our setting it can be easily checked in polynomial time that two DAGs represent the same value.*

## 5 GRSR is Sound for Polynomial Time

Sometimes (e.g., in [10]), the first step towards a proof of soundness for ramified recursive systems consists in giving a proper bound precisely relating the size of the result and the size of the inputs. More specifically, if the result has tier  $n$ , then the size of it depends polynomially on the size of the inputs of tier higher than  $n$ , but only *linearly*, and in very restricted way, on the size of inputs of tier  $n$ . Here, a similar result holds, but size is replaced by *minimal shared size*.

The *minimal shared size*  $\|v_1, \dots, v_k\|$  for a *sequence* of elements  $v_1, \dots, v_k \in \mathbb{A}$  is defined as the number of subterms in  $v_1, \dots, v_k$ , i.e. the cardinality of the set  $\bigcup_{1 \leq i \leq k} \text{STs}(v_i)$ . Then  $\|v_1, \dots, v_k\|$  corresponds to the number of locations necessary to store the values  $v_1, \dots, v_k$  on a heap (compare Lemma 4.4). If  $\mathbf{A}$  is the expression  $\mathbb{A}_{n_1} \times \dots \times \mathbb{A}_{n_m}$ ,  $n$  is a natural number, and  $\vec{t}$  is a sequence of  $m$  terms, then  $\|\vec{t}\|_{\mathbf{A}}^{\geq n}$  is defined to be  $\|t_{i_1}, \dots, t_{i_k}\|$  where  $i_1, \dots, i_k$  are precisely those indices such that  $n_{i_1}, \dots, n_{i_k} > n$ . Similarly for  $\|\vec{t}\|_{\mathbf{A}}^{\leq n}$ .

**Proposition 5.1** (Max-Poly). *If  $\mathbf{f} \triangleright \mathbf{A} \rightarrow \mathbb{A}_n$ , then there is a polynomial  $p_{\mathbf{f}} : \mathbb{N} \rightarrow \mathbb{N}$  such that  $\|\mathbf{f}(\vec{v})\| \leq \|\vec{v}\|_{\mathbf{A}}^{\leq n} + p_{\mathbf{f}}(\|\vec{v}\|_{\mathbf{A}}^{\geq n})$ .*

Once we know that ramified recursive definitions are not too fast-growing for the minimal shared size, we know that all terms around do not have a too-big minimal shared size. As a consequence:

**Proposition 5.2.** *If  $\mathbf{f} \triangleright \mathbf{A} \rightarrow \mathbb{A}_n$ , then there is a polynomial  $p_{\mathbf{f}} : \mathbb{N} \rightarrow \mathbb{N}$  such that for every  $v$ ,  $(\emptyset, \mathbf{f}(\vec{v})) \Downarrow_m (C, v)$ , with  $m \leq p_{\mathbf{f}}(\|\vec{v}\|)$ .*

The following, then, is just a corollary of Proposition 5.2 and Invariance (Theorem 4.17).

**Theorem 5.3.** *Let  $\mathbf{f} : \mathbb{A}_{p_1} \times \dots \times \mathbb{A}_{p_k} \rightarrow \mathbb{A}_m$  be a function defined by general ramified simultaneous recursion. There exists then a polynomial  $p_{\mathbf{f}} : \mathbb{N}^k \rightarrow \mathbb{N}$  such that for all inputs  $v_1, \dots, v_k$ , a DAG representation of  $\mathbf{f}(v_1, \dots, v_k)$  is computable in time  $p_{\mathbf{f}}(|v_1|, \dots, |v_k|)$ .*

**Example 5.4** (Continued from Example 4.18). *In Example 3.5 we indicated that the function `rabbits` :  $\mathbb{N} \rightarrow \mathbb{T}$  from Section 2 is definable by GRSR. As a consequence of Theorem 5.3, it is computable in polynomial time, e.g. on a Turing machine. Similar, we can prove the function `tree` from Section 2 polytime computable.*

## 6 Conclusion

In this work we have shown that simultaneous ramified recurrence on generic algebras is sound for polynomial time, resolving a long-lasting open problem in implicit computational complexity theory. We believe that with this work we have reached the *end of a quest*. Slight extensions, e.g. the inclusion of *parameter substitution*, lead outside polynomial time as soon as simultaneous recursion over trees is permissible.

Towards our main result, we introduced the notion of memoized runtime complexity, and we have shown that this cost model is invariant under polynomial time. Crucially, we use a compact DAG representation of values to control duplication, and tabulation to avoid expensive re-computations. To the authors best knowledge, our work is the first where sharing and memoization is reconciled, in the context of implicit computational complexity theory. Both techniques have been extensively employed, however separately. Essentially relying on sharing, the invariance of the unitary cost model in various rewriting based models of computation, e.g. the  $\lambda$ -calculus [1, 2, 14] and term rewrite systems [4, 13] could be proven. Several works (e.g. [6, 12, 21]) rely on memoization, employing a measure close to our notion of memoized runtime complexity. None of these works integrate sharing, instead, inputs are either restricted to strings or dedicated bounds on the size of intermediate values have to be imposed. We are confident that our second result is readily applicable to resolve such restrictions.

## References

- [1] B. Accattoli and U. Dal Lago. On the Invariance of the Unitary Cost Model for Head Reduction. In *Proc. of 23<sup>rd</sup> RTA*, volume 15 of *LIPICs*, pages 22–37. Dagstuhl, 2012.
- [2] B. Accattoli and U. Dal Lago. Beta Reduction is Invariant, Indeed. In *Proc. of 23<sup>rd</sup> CSL*, page 8. ACM, 2014.
- [3] T. Arai and N. Eguchi. A New Function Algebra of EXPTIME Functions by Safe Nested Recursion. *TOCL*, 10(4), 2009.
- [4] M. Avanzini and G. Moser. Closing the Gap Between Runtime Complexity and Polytime Computability. In *Proc. of 21<sup>st</sup> RTA*, volume 6 of *LIPICs*, pages 33–48. Dagstuhl, 2010.

- [5] F. Baader and T. Nipkow. *Term Rewriting and All That*. Cambridge University Press, 1998.
- [6] P. Baillot, U. Dal Lago, and J.-Y. Moyén. On Quasi-interpretations, Blind Abstractions and Implicit Complexity. *MSCS*, 22(4):549–580, 2012.
- [7] H. P. Barendregt, M. v. Eekelen, J. R. W. Glauert, J. R. Kennaway, M. J. Plasmeijer, and M. R. Sleep. Term Graph Rewriting. In *PARLE (2)*, volume 259 of *LNCS*, pages 141–158. Springer, 1987.
- [8] S. Bellantoni. *Predicative Recursion and Computational Complexity*. PhD thesis, University of Toronto, 1992.
- [9] S. Bellantoni. Predicative Recursion and the Polytime Hierarchy. In *Feasible Mathematics II*. Birkhäuser Boston, 1994.
- [10] S. Bellantoni and S. Cook. A new Recursion-Theoretic Characterization of the Polytime Functions. *CC*, 2(2):97–110, 1992.
- [11] G. Bonfante, R. Kahle, J.-Y. Marion, and I. Oitavem. Recursion Schemata for NCK. In *Proc. of 22<sup>nd</sup> CSL*, volume 5213 of *LNCS*, pages 49–63. Springer, 2008.
- [12] G. Bonfante, J.-Y. Marion, and J.-Y. Moyén. Quasi-interpretations: A Way to Control Resources. *TCS*, 412(25):2776–2796, 2011.
- [13] U. Dal Lago and S. Martini. Derivational Complexity is an Invariant Cost Model. In *Revised Selected Papers of 1<sup>st</sup> FOPARA*, volume 6324 of *LNCS*, pages 100–113. Springer, 2009.
- [14] U. Dal Lago and S. Martini. On Constructor Rewrite Systems and the Lambda Calculus. *LMCS*, 8(3):1–27, 2012.
- [15] U. Dal Lago, S. Martini, and M. Zorzi. General Ramified Recurrence is Sound for Polynomial Time. In *Proc. of 1<sup>st</sup> DICE*, volume 23 of *EPTCS*, pages 47–62, 2010.
- [16] N. Danner and J. S. Royer. Ramified Structural Recursion and Corecursion. *CoRR*, abs/1201.4567, 2012.
- [17] N. Hirokawa and G. Moser. Automated Complexity Analysis Based on the Dependency Pair Method. In *Proc. of 4<sup>th</sup> IJCAR*, volume 5195 of *LNAI*, pages 364–380. Springer, 2008.
- [18] B. Hoffmann. Term Rewriting with Sharing and Memoization. In *Proc. of 3<sup>rd</sup> ALP*, volume 632 of *LNCS*, pages 128–142. Springer, 1992.
- [19] D. Leivant. Stratified Functional Programs and Computational Complexity. In *Proc. of 20<sup>th</sup> POPL*, pages 325–333. ACM, 1993.
- [20] D. Leivant. Ramified Recurrence and Computational Complexity I: Word Recurrence and Poly-time. In *Feasible Mathematics II*, volume 13, pages 320–343. Birkhäuser Boston, 1995.
- [21] J.-Y. Marion. Analysing the Implicit Complexity of Programs. *IC*, 183:2–18, 2003.

- [22] Isabel Oitavem. Implicit Characterizations of Pspace. In *PTCS*, pages 170–190, 2001.
- [23] Christos H. Papadimitriou. *Computational Complexity*. Addison Wesley Longman, second edition, 1995.
- [24] D. Plump. Essentials of Term Graph Rewriting. *ENTCS*, 51:277–289, 2001.