

EPiC Series in Computing

Volume 57, 2018, Pages 543-563

LPAR-22. 22nd International Conference on Logic for Programming, Artificial Intelligence and Reasoning



Automatic Space Bound Analysis for Functional Programs with Garbage Collection

Yue Niu and Jan Hoffmann

Carnegie Mellon University, Pittsburgh, PA, United States {yuen,jhoffmann}@cs.cmu.edu

Abstract

This article introduces a novel system for deriving upper bounds on the heap-space requirements of functional programs with garbage collection. The space cost model is based on a perfect garbage collector that immediately deallocates memory cells when they become unreachable. Heap-space bounds are derived using type-based automatic amortized resource analysis (AARA), a template-based technique that efficiently reduces bound inference to linear programming. The first technical contribution of the work is a new operational cost semantics that models a perfect garbage collector. The second technical contribution is an extension of AARA to take into account automatic deallocation. A key observation is that deallocation of a perfect collector can be modeled with destructive pattern matching if data structures are used in a linear way. However, the analysis uses destructive pattern matching to accurately model deallocation even if data is shared. The soundness of the extended AARA with respect to the new cost semantics is proven in two parts via an intermediate linear cost semantics. The analysis and the cost semantics have been implemented as an extension to Resource Aware ML (RaML). An experimental evaluation shows that the system is able to derive tight symbolic heap-space bounds for common algorithms. Often the bounds are asymptotic improvements over bounds that RaML derives without taking into account garbage collection.

1 Introduction

The memory footprint of a program is an important performance metric that determines if a program can be safely executed on a given system. Ideally, developers should describe or approximate the memory footprint of programs as functions of the inputs. However, such memory bounds are often difficult to derive and to prove sound. To assist programmers with deriving memory bounds, the programming language community has developed automatic and semi-automatic analysis techniques [24, 12, 2]. These systems are often special cases of more general resource bound analyses that are based on abstract interpretation [18, 7, 37], recurrence solving [16, 1, 14, 28], type systems [27, 22, 29, 43, 42, 15], program logics [5, 10, 9, 35], proof assistants [33, 11], and term rewriting [6, 34, 17].

This article introduces a novel type system for automatically deriving upper bounds on the heap-space requirements of functional programs with garbage collection (GC). Due to the

G. Barthe, G. Sutcliffe and M. Veanes (eds.), LPAR-22 (EPiC Series in Computing, vol. 57), pp. 543–563

challenges of modeling and predicting garbage collection, most existing techniques for automating and guiding the derivation of bounds on the heap memory requirements assume manual memory management or simply ignore deallocation in the analysis [24, 26, 36, 13, 12, 2]. As a result, the derived bounds are not accurate when the underlying system employs garbage collection. The only exceptions we are aware of are the works by Albert et al. [3, 4], Braberman et al. [8], and Unnikrishnan et al. [40, 39]. They analyze the heap-space usage of programs with GC in two steps. First, they make the deallocation of GC explicit; for example with a static analysis for estimating object lifetimes [4] or with a program translation [39]. Second, they extract and solve recurrence relations to derive a bound. The difference of our work is that our technique is based on a type system, which is proved sound with respect to a formal cost semantics. Advantages of a type-based approach include natural compositionality and the use of type derivations as certificates for resource bounds.

We model the (highwater mark) memory usage based on a perfect garbage collector that immediately deallocates memory cells when they become unreachable. The bounds that are derived with respect to this cost model are not only a good theoretical measure of the heap-space consumption of the program but also have practical relevance. Consider a function $f: A \to B$ and assume we derived a bound $b_f: [A] \to \mathbb{N}$. In an execution of f(a), we can then keep track of the memory usage and start the garbage collector whenever the bound $b_f(a)$ is reached. It is then guaranteed that the evaluation will succeed using $b_f(a)$ heap-memory cells.¹ To improve performance, we could trigger GC more often (to compactify the heap) or allow memory use of more than $b_f(a)$ cells (to amortize the cost of garbage collection).

The first technical contribution of the work is a new operational cost semantics that models a perfect garbage collector. The cost semantics is a big-step (or natural) semantics that keeps track of the reachable memory cells in the style of Spoonhower et al. [38] and Minamide [30]. Operationally, this cost is the highwater mark on the heap usage, or the maximum number of cells used in the mutable store during evaluation. If we traverse the evaluation tree in preorder and view each node as a "step" of the computation, then a cell is used in the current node if it is reachable from the reminder of the computation. Our formalization of reachability is identical to the concept that garbage collectors implement to decide if a cell can be freed during evaluation. For simplicity, we assume that evaluation of the cons node allocates one fresh heap cell and that all other operations do not allocate heap cells. However, the semantics can be instantiated with more realistic cost metrics. A difference to existing formulations of cost semantics with GC [31, 38, 30] is that we update the highwater mark when reachability changes at inner nodes of the derivation of the evaluation judgement instead of at leaves. Moreover, we use a *freelist*, which represents *named* cells available for evaluation. This alternative formulation is equivalent to the existing semantics and mainly motivated by the soundness proof of our type system for bound analysis. However, the cost semantics is a natural approach and different enough from its predecessors [38, 30] to be of interest in its own right.

Our second technical contribution is the type system for deriving bounds on the heap-space for programs with perfect GC. The type system is an extension of type-based automatic amortized resource analysis (AARA) [24, 27, 41, 21, 22, 32]. AARA is a template-based technique that introduces potential functions to efficiently and automatically reduce bound inference to linear programming. Existing type systems based on AARA can derive bounds on the highwater mark of the heap usage for programs with manual deallocation [27], but can only derive a bound on the number of total heap allocations for programs with GC [22]. This is usually a gross over-approximation of the actual memory requirement. Our extension is based on the observation that deallocation of a perfect collector can be modeled with destructive pattern

¹We are not considering memory fragmentation, which can be avoided using a copying collector.

matching (deallocating the matched cell) if data structures are used in a linear way. In the type system, we extend this observation to non-linear programs and use destructive pattern matching to accurately model deallocation even if data is shared.

The third technical contribution is to prove the soundness of the extended AARA with respect to the GC-based cost semantics. The proof is non-trivial and proceeds in two parts: First, we prove the soundness of the type system with respect to a semantics that copies data structures if they are shared. Second, we prove for all programs that our GC semantics uses less memory than this copying semantics. While the proofs are relatively standard, many details—like relating program states of the two semantics in the simulation proof—are quite involved. Briefly, we have to provide and maintain a mapping γ from the heap used in the GC semantics H_{gc} to subsets of the heap used in the copying semantics H_{copy} such that the image of H_{gc} under γ forms a partition on the second heap. The intuition is that given a cell $l \in H_{gc}$, there must be multiple cells $\gamma(l) \in H_{copy}$ that were allocated during sharing, and thus "morally the same" as l.

The analysis and the cost semantics have been implemented as an extension to Resource Aware ML (RaML) [21, 22]. RaML is an implementation of AARA for a subset of OCaml that can derive multivariate polynomial bounds. However, we restrict the technical development in this paper to a simple first-order language with tuples and lists. The proofs and ideas carry over to the more complex case of RaML.² An experimental evaluation shows that the system is able to derive tight symbolic heap-space bounds for common algorithms. Our results suggest that our new analysis provides asymptotic bound improvements to several classes of commonly used functions and programming patterns. We examine the reasons for these improvements and design decisions throughout the system.

2 Setting the Stage

In the technical part of the paper, we focus our attention to a first-order, strictly evaluated functional language. One can think of this language as a simple subset of OCaml or SML. The only recursive data type in the language is the list type. However, our work extends to the expected algebraic data types definable in RaML. Being first order, the language does not allow arbitrary local functional definitions. Instead, all functions are defined at the top level and are mutually recursive by default. The types of these functions form a signature for the program, and the semantics and typing judgments will be indexed by this signature. Thus, the function types of the language can be expressed as arrows between zero-order (base) types. Types are formally defined in Figure 1. Like in all grammars, we provide the abstract (left) and concrete (right) syntax for every type former [19]. A signature $\Sigma : \text{Var} \to \text{FTypes}$ is a map from variables to first-order types. A program P is a Σ indexed map from Var to pairs $(y_f, e_f)_{f \in \Sigma}$, where $\Sigma(y_f) = \tau \to \tau'$, and $\Sigma; y_f : \tau \vdash e_f : \tau'$ (the type system is discussed in Section 4). We write $P : \Sigma$ to mean P is a program with signature Σ .

To simplify the presentation, the expressions of our language (see Figure 1) are in *let normal* form (also A normal form). The one nonstandard construct is **share** x **as** x_1, x_2 **in** e, which we will explain in more detail in the following sections. We introduce two distinct notions of *linearity*, one on the syntactic level, and one on the semantic level. Syntactic linearity is linearity in expression variables, while semantic linearity is linearity in locations (defined below). We say that a semantics is linear if it respects semantic linearity.

In line with previous works on space cost semantics [38, 30], we employ a heap, which

 $^{^{2}}$ An exception are function closures that we discuss in the Section 7.

Niu and Hoffmann

BTypes $ au$::=	E	Exp e ::=	
nat	nat	$\mathtt{var}(x)$	x
unit	unit	$\mathtt{nat}[n]$	\overline{n}
bool	bool	unit	()
$\texttt{prod}(\tau_1;\tau_2)$	$ au_1 imes au_2$	Т	Т
$\texttt{list}(\tau)$	$L(\tau)$	F	F
		$if(x;e_1;e_2)$	$\texttt{if}x\texttt{then}e_1\texttt{else}e_2$
FTypes ρ ::=		$\mathtt{ap}(f;x)$	f(x)
$\texttt{arr}(\tau_1;\tau_2)$	$ au_1 o au_2$	$\mathtt{tpl}(x_1;x_2)$	$\langle x_1, x_2 angle$
		$\mathtt{match}_\mathtt{P}(x_1, x_2.e_1)$	$\texttt{match} \ p \ \{(x_1; x_2) \hookrightarrow e_1\}$
Val v ::=		nil	[]
$\mathtt{val}(n)$	n	$\mathtt{cons}(x_1;x_2)$	$x_1 :: x_2$
val(T)	Т	$\mathtt{match}_{\mathtt{L}}\{l\}(e_1; x, xs.e_2)$	$\texttt{match}l\{\texttt{nil} \hookrightarrow e_1 \mid \texttt{cons}(x;xs) \hookrightarrow e_2\}$
val(F)	F	$\mathtt{let}(e_1;x: au.e_2)$	$\mathtt{let}\; x=e_1\; \mathtt{in}\; e_2$
$\mathtt{val}(\mathtt{Null})$	Null	$\mathtt{share}(x; x_1, x_2.e)$	share x as x_1, x_2 in e
$\mathtt{val}(l)$	l		
$\mathtt{val}(\mathtt{pair}(v_1;v_2))$	$\langle v_1, v_2 \rangle$		

Figure 1: Simple Types, Values, and Expressions

persistently binds locations to values (normalized terms). As usual, we derive the cost of a (terminating) program from the number of heap locations used during execution, which in our case is the *maximum difference* between the sizes of the initial and final freelist. We let Loc be an infinite set of names for addressing the heap. For the rest of the paper, we use the following: Stack $\triangleq \{V \mid V : \mathsf{Var} \to \mathsf{Val}\}$ and $\mathsf{Heap} \triangleq \{H \mid H : \mathsf{Loc} \to \mathsf{Val}\}$ for the set of stacks and heaps respectively.

Reachability Before we define the rules for the cost semantics, we relate the heap locations to values with the 3-place reachability relation reach(H, v, L) on $\text{Heap} \times \text{Val} \times \wp(\text{Loc})$, where \wp is the powermultiset. This is read as "under heap H, the value v reaches the multiset of locations L". Write $L = reach_H(v)$ to indicate this is a functional relation justified by the (valid) mode (+, +, -). We say that the reachable set of v is L.

$$\frac{A = reach_H(v_1) \qquad B = reach_H(v_2)}{A \uplus B = reach_H(\langle v_1, v_2 \rangle)} \qquad \frac{A = reach_H(H(l))}{\{l\} \uplus A = reach_H(l)} \qquad \frac{v \in \mathbb{N} \cup \{\mathsf{T}, \mathsf{F}, \mathsf{Null}\}}{\emptyset = reach_H(v)}$$

In the rules, \uplus is multiset union. L is a multiset because we need to keep track of the *number* of ways a location might be reached in order to prove soundness. However, the cost semantics can be read by truncating any multiset to a set. Furthermore, we will sometimes mix multiset and set operations as the situation calls for. For example, we will write $l \in S$ for a multiset S if $S(l) \geq 1$.

The notion of reachability naturally lifts to expressions and contexts:

$$reach_{H}(V) = \biguplus_{x \in dom(V)} reach_{H}(V(x)) \qquad locs_{V,H}(e) = reach_{H}(V \upharpoonright_{FV(e)})$$

Where $FV : \mathsf{Exp} \to \mathcal{P}(\mathsf{Var})$ denotes the set of free-variables of expressions as usual.

Towards the Garbage Collection Cost Semantics Now we are ready to give a first attempt to modeling the cost semantics for a tracing garbage collector. Before we present our

new semantics, we explain an existing cost semantics we experimented with [30]. Judgements have the form $V, H, R \vdash e \Downarrow^s v, H'$, which can be read as follows. Under stack $V \in \mathsf{Stack}$, heap $H \in \mathsf{Heap}$, and continuation set $R \subseteq \mathsf{Loc}$, e evaluates to v and H' using s heap locations. The idea is that R keeps track of the set of locations necessary to complete the evaluation *after* e is evaluated (hence the name continuation). For example, we have the let rule:

$$\frac{V, H, R \uplus locs_{V,H}(x.e_2) \vdash e_1 \Downarrow^{s_1} v_1, H_1 \qquad V[x \mapsto v_1], H, R \vdash e_2 \Downarrow^{s_2} v_2, H_2}{V, H, R \vdash \mathsf{let}(e_1; x : \tau.e_2) \Downarrow^{\max s_1, s_2} v_2, H_2}$$

Notice that to evaluate e_1 , we have to extend the continuation R with locations in e_2 , which will be used *after* e_1 is evaluated. The total space used is the max of the component, indicating that locations used for e_1 can be reused for e_2 . This is clear when we look at the variable rule:

$$\frac{V(x) = v}{V, H, R \vdash x \Downarrow^{|dom(R \uplus reach_H(v))|} v, H}$$

It states that evaluating a variable x requires the locations reachable from x as well as the continuation set R. While this way of counting heap locations does model a tracing garbage collector, it is not compatible with the existing type systems for amortized analysis. In these systems, such as RaML, the type rules count the heap locations as data is created, i.e. at each data constructor. Thus looking up a variable incurs no cost, since it was accounted for during creation. On the other hand, the cost of indexing a variable in the semantics includes the cost of the entire continuation set, which is potentially unbounded. This mismatch between the dynamics and statics of language prevents us from proving the soundness of the analysis. We give a new cost semantics that is 1) compatible with the type system and 2) also a more concrete model of a garbage collector since costs are realized with explicit locations.

3 Garbage Collection Cost Semantics

In this section, we present our novel cost semantics by combining *freelist semantics* from [25] with the cost semantics for modeling perfect GC [30] that we discussed in the previous section. The resulting semantics, called \mathcal{E}_{gc} , is well suited for proving the soundness of the novel type-based bound analysis.

The garbage collection cost semantics \mathcal{E}_{gc} is defined by a collection of judgement of the form

$$\mathcal{C} \vdash_{P:\Sigma} e \Downarrow v, H', F'$$

Where $C \in \text{Stack} \times \text{Heap} \times \wp(\text{Loc}) \times \mathcal{P}(\text{Loc})$ is a *configuration* usually written with variables V, H, R, F. Because the signature Σ for the mapping of function names to first-order functions does not change during evaluation, we drop the subscript $P : \Sigma$ from $\vdash_{P:\Sigma}$ when the context of evaluation is clear. Given a configuration $\mathcal{C} = (V, H, R, F)$, the evaluation judgment states that under stack V, heap H, continuation (multi)set R, freelist F, and program P with signature Σ , the expression e evaluates to value v, and engenders a new heap H' and freelist F'. In comparison with the attempt from the previous section, the key ingredient we added is the freelist, which serves as the set of available locations. Similar to the predicate *reach*, We call R a (multi)set since the fact that it's a multiset is only useful during the soundness proof. For evaluation, it is convenient to just view R as a set. Define a *computation* as a pair (\mathcal{C}, e) of a configuration \mathcal{C} and an expression e.

The semantics \mathcal{E}_{gc} is designed to model the heap usage of a program running with a tracing counting garbage collector: whenever a heap cell becomes unreachable from the root set, it

Niu and Hoffmann

$$\begin{split} & V_{1} = V \upharpoonright_{FV(e_{1})} R' = R \cup locs_{V,H}(lam(x: \tau.e_{2})) \\ & V_{1},H,R',F \vdash e_{1} \Downarrow v_{1},H_{1},F_{1} V_{2} = (V[x \mapsto v_{1}]) \upharpoonright_{FV(e_{2})} g = \{l \in H \mid l \notin F_{1} \cup R \cup locs_{V_{2},H_{1}}(e_{2})\} \\ & V_{2},H_{1},R,F \cup g \vdash e_{2} \Downarrow v_{2},H_{2},F_{2} \\ & V(x) = T V' = V \upharpoonright_{FV(e_{1})} V',H,R,F \cup g \vdash e_{1} \Downarrow v,H',F' \\ & \frac{g = \{l \in H \mid l \notin F \cup R \cup locs_{V',H}(e_{1})\} V',H,R,F \cup g \vdash e_{1} \Downarrow v,H',F' \\ & V(x) = T V' = V \upharpoonright_{FV(e_{2})} V',H,R,F \cup g \vdash e_{2} \Downarrow v,H',F' \\ & V(x) = F V' = V \upharpoonright_{FV(e_{2})} V',H,R,F \cup g \vdash e_{2} \Downarrow v,H',F' \\ & V(x) = F V' = V \upharpoonright_{FV(e_{2})} V',H,R,F \cup g \vdash e_{2} \Downarrow v,H',F' \\ & \frac{V(x) = v' P(f) = (y_{f},e_{f})}{V,H,R,F \vdash if(x;e_{1};e_{2}) \Downarrow v,H',F'} (F:CondF) \\ & \frac{V' = ([y_{f} \mapsto v']) \upharpoonright_{FV(e_{f})} g = \{l \in H \mid l \notin F \cup R \cup locs_{V,H}(e_{f})\} V',H,R,F \cup g \vdash e_{f} \Downarrow v,H',F' \\ & V(x) = v' P(f) = (y_{f},e_{f}) V',H,R,F \cup g \vdash e_{f} \Downarrow v,H',F' \\ & V(x) = v' P(f) = (y_{f},e_{f}) V',H,R,F \cup g \vdash e_{f} \Downarrow v,H',F' \\ & V(x) = v' P(f) = (v_{f},v_{1}) V',H,R,F \cup g \vdash e_{f} \Downarrow v,H',F' \\ & V(x) = v' V(x) = v (V(x_{1}),V(x_{2})) l \in F H' = H\{l \mapsto v\} \\ & V,H,R,F \vdash nil \Downarrow val(Null),H,F \\ & (F:Nil) V' = V \upharpoonright_{FV(e_{1})} V' = V \upharpoonright_{FV(e_{1})} V',H,R,F \cup g \vdash e_{1} \Downarrow v,H',F' \\ & \frac{V(x) = Null V' = V \upharpoonright_{FV(e_{1})} V' = V \upharpoonright_{FV(e_{1})} V',H,R,F \sqcup g \vdash e_{1} \Downarrow v,H',F' \\ & \frac{V(x) = Null V' = V \upharpoonright_{FV(e_{1})} V',H,R,F \vdash u \Downarrow v,H',F' \\ & \frac{V(x) = Null V' = V \upharpoonright_{FV(e_{1})} V',H,R,F \sqcup g \vdash e_{1} \Downarrow v,H',F' \\ & \frac{V(x) = V V' (y') = V (y',y',y') \lor_{FV(e_{2})} V',H,R,F \sqcup g \vdash e_{2} \Downarrow v,H',F' \\ & \frac{V(x) = V V V' = (V[x_{h} \mapsto w_{h},x_{t} \mapsto v_{h}) \upharpoonright_{FV(e_{2})} V',H,R,F \vdash match_{k}\{e_{1};x_{h},x_{t},e_{2} \Downarrow v,H',F' \\ & \frac{V(x) = v V V' (y') (V' = (V[x_{h} \mapsto w_{h},x_{t} \mapsto v_{h}) \upharpoonright_{FV(e_{2})} V,H,R,F \vdash match_{k}\{e_{1};x_{h},x_{t},e_{2} \lor v,H',F' \\ & \frac{V(x) = v' V' V' (V' = (V[x_{h} \mapsto w_{h},x_{t} \mapsto v_{h}) \lor_{F} V' \\ & V,H,R,F \vdash match_{k}\{e_{1};x_{h},x_{t},e_{2} \lor v,H',F' \\ & \frac{V(x) = v' V' V' (V' = (V[x_{h} \mapsto v',x_{2} \mapsto v']) \upharpoonright_{FV(e_{1})} V' \\ & V,H,R,F \vdash match_{k}\{e_{1};x_{h},x_{t},e_{2} \lor v,H',F' \\ & \frac{V(x) = v' V' V' V' (Y' = V[x_{1} \mapsto v',x_{2$$

Figure 2: Cost Semantics for Perfect Garbage Collection

becomes collected and added to the freelist as available for reallocation. As before, the continuation set R represents the set of locations required to compute the continuation *excluding* the current expression. We define the *root set* as the union of the locations in the continuation set R and the locations in the current expression e.

The inference rules for the semantics are given in Figure 2. For example, the rule F:CondT states that, to evaluate a conditional, look in the stack for the value of the branching boolean. In the case it is true, we proceed to evaluate the first branch. Furthermore, we collect cells in the heap that are not reachable from the root set $(R \cup locs_{V',H}(e_1))$ or already in the current free-list F, and add them (g) to the available cells for evaluating e_1 .

Another example is the rule F:Let for let expressions: to evaluate the expressions $let(e_1; x:\tau.e_2)$, we evaluate the first expression with the corresponding restricted stack V_1 and an expanded continuation set R'. The extra locations come from the free variables of e_2 (not including the bound variable x), which we cannot collect during the evaluation of e_1 . Next, we restrict the extended stack to only free variables of e_2 , and evaluate e_2 with this stack and the original continuation set R. The other rules are similar.

Note that in contrast to the semantics in the previous section, evaluating a variable does not incur any cost. This ensures that we will be able prove the soundness of the type system. Also,

since we don't allow local function definitions, we do not create closures during evaluation. Also note that we restrict the domain of the stack to the appropriate variables during evaluation. This is only to facilitate the proof of the linearity of the copying semantics introduced later, and not necessary for the implementation.

For example, we can implement the **append** and **appTwice** function, which has variable sharing. First, we analyze the heap usage of **append** under \mathcal{E}_{gc} . We case on the first component of the input. In case it's nil, we just return 12, and there are no allocations or deallocations. In case it's cons of x and xs, we need to allocate one heap location for the cons cell binding x and the recursive result, for which we can use the just matched-on cell. Again, the net overhead is zero. Thus, the total space overhead of **append** is zero.

For appTwice, we first share the list 1 as 11 and 12. In the first let, the locations in 12 are added to the continuation set, which prevents the first call to append from destructing 11. Thus size of 11 new locations are allocated from the freelist to construct 11'. The second call has no net increase in heap allocations since 12 can be destructed along the way. The return value is a pair which is stack-allocated and doesn't require a heap allocation. Thus, the total space overhead for appTwice is size of the input list 1.

From this, we see that the minimum size for the initial freelist to successively evaluate a call to **appTwice** is exactly the length of the input. In general, we define the cost of a closed program to be the minimum size of the initial freelist that guarantees successful evaluation, which is equivalent to the cost annotation in the previous cost semantics introduced in Section 2.

4 Automatic Amortized Heap-Space Analysis with GC

Automatic Amortized Resource Analysis (AARA) The idea of AARA [24, 27, 21, 22] is to automate the potential method of amortized analysis using a type system. Types introduce potential functions that map data structures of the given type to non-negative numbers. The type rules ensure that there is always sufficient potential to cover the evaluation cost of the next step and the potential of the next program state.

To illustrate the idea, we informally explain the linear potential method for the functions in Figure 3. We will use the allocation/heap metric which simply counts the number of cons constructor calls during the evaluation.³ With this metric, the cost of evaluating append(11,12) is m, where m is the number of cons constructors in 11, and the resource annotated type of append is $L^1(\text{int}) \times L^0(\text{int}) \xrightarrow{0/0} L^0(\text{int})$. This type says that to type append(11,12), we need 11 to have 1 potential per element, 12 to have 0

```
let rec append (11, 12) =
   match 11 with
   | [] -> 12
   | x::xs -> x::(append (xs, 12))
let appTwice 1 =
    share 1 as 11,12 in
   let 11' = append (11, []) in
   let 12' = append (12, []) in
   (11',12')
```

Figure 3: Functions append and appTwice

per element, and the result will be a list with 0 potential per element. Additionally, the function uses 0 constant potential, and leaves 0 constant potential after evaluating. This translates to a bound which states that the number of allocations append makes is bounded by 1 times size of the first list. For appTwice(1), the cost under the heap metric is 2m, where m is the number of cons constructors in 1. This is because we have to share the input list across two calls of append, which each requires lists with unit potential per element. For example, if $1 : L^2(int)$,

³This is in contrast to the highwater mark for the GC semantics \mathcal{E}_{gc} that is targeted by our new analysis.

Niu and Hoffmann

$$\frac{n \in \mathbb{Z}}{H \models \operatorname{val}(n) \mapsto n : \operatorname{nat}} (V:\operatorname{ConstI}) \qquad \qquad \overline{H \models \operatorname{val}(\operatorname{Null}) \mapsto \operatorname{val}(\operatorname{Null}) : \operatorname{unit}} (V:\operatorname{ConstI}) \\ \frac{A \in \operatorname{BType}}{H \models \operatorname{val}(\operatorname{Null}) \mapsto \operatorname{val}(\operatorname{Null}) : L(A)} (V:\operatorname{Nil}) \qquad \qquad \overline{H \models \operatorname{val}(T) \mapsto \operatorname{val}(T) : \operatorname{bool}} (V:\operatorname{True}) \\ \frac{H \models \operatorname{val}(F) \mapsto \operatorname{val}(F) : \operatorname{bool}}{H \models \operatorname{val}(F) : \operatorname{bool}} (V:\operatorname{False}) \qquad \qquad \frac{H \models v_1 \mapsto a_1 : A_1 \qquad H \models v_2 \mapsto a_2 : A_2}{H \models \langle v_1, v_2 \rangle \mapsto \langle a_1, a_2 \rangle : A_1 \times A_2} (V:\operatorname{Pair}) \\ \frac{l \in \operatorname{Loc} \qquad H(l) = \langle v_h, v_t \rangle}{H \models l \mapsto [a_1, \dots, a_n] : L(A)} \qquad (V:\operatorname{Cons})$$

Figure 4: Mapping Locations to Semantic Values

then 11 and 12 both get 1 potential per element so that $l1 : L^1(int), l2 : L^1(int)$, which covers the cost of the next 2 calls to append, and the resulting pair of lists both have 0 potential per element.

More generally, we can give the following types to append and appTwice:

append:
$$L^p(\text{int}) \times L^q(\text{int}) \xrightarrow{r/r'} L^s(\text{int})$$
, where $p \ge s+1$, $q \ge s$ and $r \ge r'$
appTwice: $L^p(\text{int}) \xrightarrow{q/q'} L^r(\text{int}) \times L^s(\text{int})$, where $p \ge r+s+2$ and $q \ge q'$

Notice that the constant potentials r and q are unconstrained since the functions don't use any potential in the base cases. With AARA, the type system keeps track of this collection of constraints on resource annotations and passes them to an off-the-shelf LP-solver which finds the minimum solution. This is then translated to concrete resource bounds like the ones we derived by hand. It has been shown that this technique can be extended to polynomial potential functions, user-defined data types, and higher-order functions while still relying on linear constraint solving [21, 22].

Linear Potential Functions Before giving the type rules, we need to formalize linear potential as explained above. Since potential is associated with the *structure* of a value and not the particular heap locations, it is helpful to introduce a mapping from heap values to semantic values of a type. First, we give a denotational semantics for (define the structures of) the first-order types:

$$\begin{split} \llbracket \texttt{unit} \rrbracket = \{\texttt{val}(\texttt{Null})\} & \llbracket A_1 \times A_2 \rrbracket = \llbracket A_1 \rrbracket \times \llbracket A_2 \rrbracket \\ \llbracket \texttt{bool} \rrbracket = \{\texttt{val}(\texttt{T}),\texttt{val}(\texttt{F})\} & \texttt{nil} \in \llbracket L(A) \rrbracket \\ \llbracket \texttt{nat} \rrbracket = \mathbb{N} & \texttt{cons}(a; l) \in \llbracket L(A) \rrbracket \text{ if } a \in \llbracket A \rrbracket \text{ and } l \in \llbracket L(A) \rrbracket \end{split}$$

The meaning of each type is the least set such that the above holds. As usual, we write $[a_1, ..., a_n]$ for $cons(a_1; ..., cons(a_n; nil))$.

In Figure 4 we give the judgements relating heap values to semantic values, in the form $H \models v \mapsto a : A$, which can be read as follows: Under heap H, heap value v defines the semantic value $a \in \llbracket A \rrbracket$. Given a stack V, we write $H \models V : \Gamma$ if $dom(V) \subseteq dom(\Gamma)$ and for every $x \mapsto v \in V, H \models V(x) \mapsto a : \Gamma(x)$ for some $a \in \llbracket A \rrbracket$.

We introduce linear potential for structures corresponding to the base types. The definition of linear potential is standard [20]. Below is the grammar for resource-annotated types:

BTypes
$$A ::=$$
 FTypes $\rho ::=$
... $\operatorname{arr}(A_1; A_2; p; q) \quad A_1 \xrightarrow{p/q} A_2$
list^p(A) $L^p(A)$

The intended meaning is that a list of $L^p(A)$ has p units of potential per cons cell, and a function of type $A \xrightarrow{p/q} B$ takes constant potential p to run and q is the constant potential left afterwards.

With linear potential, each component of a structure is associated with a constant amount of potential. Given a structure a in a heap H, where $H \vDash v \mapsto a : A$, we define its potential $\Phi_H(a:A)$ by recursion on A:

$$\begin{split} \Phi_H(v:A) &= 0 & \text{if } A \in \{\texttt{unit},\texttt{bool},\texttt{nat}\} \\ \Phi_H(\langle v_1, v_2 \rangle : A_1 \times A_2) &= \Phi_H(v_1:A_1) + \Phi_H(v_2:A_2) \\ \Phi_H(l:L^p(A)) &= p + \Phi_H(v_h:A) + \Phi_H(v_T:L^p(A)) & \text{if } H(l) = \langle v_h, v_h \rangle \end{split}$$

Write $\Phi_{V,H}(\Gamma)$ for $\sum_{x \in dom(V)} \Phi_H(V(x) : \Gamma(x))$.

Now define $A
ightarrow A_1, A_2, n$ as the sharing relation for resource-annotated types:

$L^p(A) \Upsilon^n L^q(A_1), L^r(A_2)$	if $p = q + r + n$ and $A \Upsilon^n A_1, A_2$
$A \times B \Upsilon^n A_1 \times B_1, A_2 \times B_2$	if $A \uparrow^n A_1, A_2$ and $B \uparrow^n B_1, B_2$
$A \Upsilon^n A, A$	$\text{if }A\in\{\texttt{unit},\texttt{bool},\texttt{nat}\}$

The sharing relation captures the amount of potential needed to copy a type A where each cons node in any structure in [A] has a copying overhead n.

Type Rules The type system FO^{gc} consists of rules of the form $\Sigma; \Gamma \mid \frac{q}{q'} e : A$, read as under signature $\Sigma : \mathsf{Var} \to \mathsf{FTypes}$, typing environment $\Gamma : \mathsf{Var} \to \mathsf{BTypes}$, e has type A starting with q units of constant potential and ending with q' units.

Our type system is based on the one of classic linear AARA [24]. We give a review of the rules in Figure 5. Since we are interested in the number of heap locations, there is an implicit side condition in all rules which ensures all constants are assumed to be nonnegative.

For example, L:Cons states that to add an element to a list with p potential per element, we need p + 1 units of constant potential: p to maintain the potential of the list, and 1 for allocating the cons cell. L:MatL states that matching on a list with type $L^p(A)$, we need to type the nil case with the same constant potentials, and we need to type the cons case with an additional p units of constant potential, since we get the spill of p from the definition of linear potential. As the last example, we look at L:Share, which states that to share a variable x of type A, we need to split the potential between A_1 and A_2 , and type the rest of the expression with the two new variables $x_1 : A_1, x_2 : A_2$.

New Rules The new type system for programs with garbage collection replaces the rules L:MatL and L:Share. The observation is that if we ensure that locations are used linearly, we can use destructive pattern matching to model local garbage collection by returning the

Niu and Hoffmann

$$\frac{\sum(f) = A \frac{q/q'}{\Sigma; x : B \mid \frac{q}{q} x : B}}{\sum; x : A \mid \frac{q}{q'} f(x) : B} (L:Fun) \qquad \frac{\sum(f) = A \frac{q/q'}{\Sigma; x : A \mid \frac{q}{q'} f(x) : B}}{\sum; x : A \mid \frac{q}{q'} f(x) : B} (L:Fun) \qquad \frac{\sum; \Gamma \mid \frac{q}{q'} e : B}{\sum; \Gamma, x : bool \mid \frac{q}{q'} if x then e_t else e_f : B} (L:Cond)$$

$$\frac{\overline{\Sigma}; x : A_1, x_2 : A_2 \mid \frac{q}{q} \langle x_1, x_2 \rangle : A_1 \times A_2}{\sum; x_1 : A_1, x_2 : A_2 \mid \frac{q}{q} \langle x_1, x_2 \rangle : A_1 \times A_2} (L:Pair) \qquad \frac{\Sigma; \Gamma, x_1 : A_1, x_2 : A_2 \mid \frac{q}{q'} e : B}{\sum; \Gamma, x : (A_1, A_2) \mid \frac{q}{q'} match x \{(x_1; x_2) \hookrightarrow e\} : B} (L:MatP)$$

$$\frac{\overline{\Sigma}; \emptyset \mid \frac{q}{q} nil : L^p(A)}{\sum; \emptyset \mid \frac{q}{q'} e_1 : B} \qquad \sum; \Gamma, x_h : A, x_t : L^p(A) \mid \frac{q+p+1}{q} cons(x_h; x_t) : L^p(A) (L:Cons)$$

$$\frac{\sum; \Gamma \mid \frac{q}{q'} e_1 : B}{\sum; \Gamma, x : L^p(A) \mid \frac{q}{q'} match x \{nil \hookrightarrow e_1 \mid cons(x_h; x_t) \hookrightarrow e_2\} : B} (L:MatL)$$

$$\frac{A \ \Upsilon \mid A_1, A_2 \qquad \sum; \Gamma, x_1 : A_1, x_2 : A_2 \mid \frac{q}{q'} e : B}{\sum; \Gamma, x_1 : A_1, x_2 : A_2 \mid \frac{q}{q'} e : B} (L:Share) \qquad \frac{\sum; \Gamma \mid \frac{q}{p} e_1 : A \qquad \sum; \Gamma_2, x : A \mid \frac{p}{q'} e_2 : B}{\sum; \Gamma_1, \Gamma_2 \mid \frac{q}{q'} let(e_1; x : \tau.e_2) : B} (L:Let)$$

Figure 5: Type Rules of Classic AARA [24]

potential associated with the constructor location (notice the extra +1 in the second premise):

$$\frac{\Sigma; \Gamma \left| \frac{q}{q'} e_1 : B \qquad \Sigma; \Gamma, x_h : A, x_t : L^p(A) \left| \frac{q+p+1}{q'} e_2 : B \right.}{\Sigma; \Gamma, x : L^p(A) \left| \frac{q}{q'} \operatorname{match} x \left\{ \operatorname{nil} \hookrightarrow e_1 \mid \operatorname{cons}(x_h; x_t) \hookrightarrow e_2 \right\} : B} (L:MatLD)$$

This is validated by the fact (Lemma 3) that in the auxiliary copying semantics (introduced in later), once a cons-cell is matched on, there can be no live references from the root set to it, and thus we are justified in restituting the potential to type the subexpression e_2 .

However, the rule L:MatLD is not sound for programs with aliasing of data. We address this issue by replacing the rule L:Share with the rule L:ShareCopy:

$$\frac{A \ \Upsilon^1 \ A_1, A_2 \qquad \Sigma; \Gamma, x_1 : A_1, x_2 : A_2 \left| \frac{q}{q'} e : B \right|}{\Sigma; \Gamma, x : A \left| \frac{q}{q'} \right|^q} \text{ share } x \text{ as } x_1, x_2 \text{ in } e : B} (\text{L:ShareCopy})$$

To share a variable of type A, we need to split the potential between two new annotated types A_1 and A_2 as usual. In addition, we have to pay an "overhead" of 1 for every cons node in any structure in $[\![A]\!]$. The idea is that we treat data as if it is actually copied. This is sound w.r.t. the copying semantics because the size of the domain of the reachable set of a value v is exactly the linear potential of v : A with all resource annotations set to 1.

For example, Figure 6 contains derivations for append and appTwice. Here, A is short for int and $\Sigma = [\text{append} \mapsto L^p(A) \times L^p(A) \xrightarrow{q/q} L^p(A)]$ is the program signature (since append is the only function actually invoked in the program). From these derivations, we get the improved space overhead bound to append and appTwice:

append:
$$L^p(\text{int}) \times L^q(\text{int}) \xrightarrow{r/r'} L^s(\text{int})$$
, where $p \ge s, q \ge s$, and $r \ge r'$
appTwice: $L^p(\text{int}) \xrightarrow{q/q'} L^r(\text{int}) \times L^s(\text{int})$, where $p \ge r + s + 1$ and $q \ge q'$

Cost Metrics In previous versions of AARA [27, 21], the typing judgment and cost semantics are parametrized by a *cost metric* $m : \text{res}_\text{const} \to \mathbb{Q}$, which assigns a constant cost to each

Niu and Hoffmann

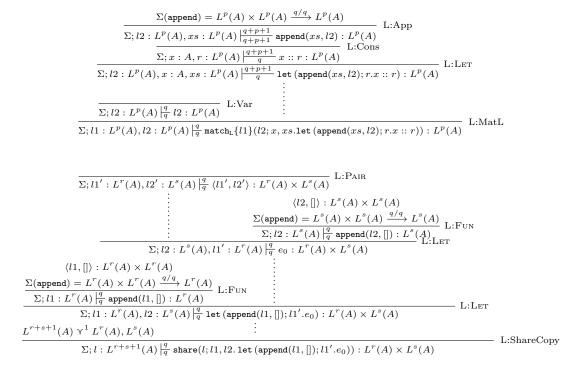


Figure 6: Type derivations for the functions append and appTwice. In the derivation for appTwice we write e_0 for let (append(l2, []); $l2' \cdot \langle l1', l2' \rangle$).

step in the semantics. Recall the heap metric introduced above; formally, this is the function $k \mapsto \mathbb{1}_{k=k^{\text{cons}}}$. We instantiate the previous type system with this metric (which only accounts for heap *allocations*), resulting in a concrete type system RaML^{heap}. We give a full evaluation of the improvements of FO^{gc} over RaML^{heap} in Section 7. Although we defined the constructor to cost 1 heap location (as shown in L:Cons and L:MatLD), it can be any constant as long as the introduction and elimination rules agree on the constant. Thus we can extend the type system to accurately track constructors which vary in size depending on the argument.

Type Inference One of the benefits of AARA is efficient type inference using off-the-shelve LP solvers [24], even for non-linear potential functions [21, 22]. The new rules do not complicate inference and previous techniques still apply. In a nutshell, inference is performed in three steps: First, perform a standard Hindley-Milner type inference for the base types. Then, annotate the type derivation with (yet unknown) variables for the potential annotations and collect linear constraints that are derived from the type rules. Finally, solve the constraints with an LP solver and minimize the potential annotations of the inputs. Details can be found in previous work [24, 22].

5 Soundness of FO^{gc}

We seek to prove the following theorem.

Niu and Hoffmann

Theorem 1 (Soundness). Let $H \vDash V:\Gamma$, $\Sigma; \Gamma \mid \frac{q}{q'} e : B$, and $V, H \vdash^{\mathcal{E}_{oper}} e \Downarrow v, H'$. Then for all configurations W, Y, F, R, if $V, H \sim W, Y$ and $|F| \ge \Phi_{V,H}(\Gamma) + q$, there exists a value w, and a freelist F' such that

$$W, Y, R, F \vdash^{\mathcal{E}_{gc}} e \Downarrow w, Y', F' \quad and \quad v \sim^{H'}_{Y'} w$$

Here, \mathcal{E}_{oper} is a standard big-step semantics, with judgments of the form $V, H \vdash e \Downarrow v, H'$ derived from $\mathcal{E}_{gc}, V, H \sim W, Y$ is context equivalence, and $v \sim_{Y'}^{H'} w$ is value equivalence (these are defined below). The theorem states that, given a terminating expression and a freelist that is sufficiently large (as predicated by the type derivation), a run with \mathcal{E}_{gc} will normalize to an equivalent value.

To facilitate the proof, we define an intermediate semantics \mathcal{E}_{copy} which is semantically linear. The proof has two stages: First, we show \mathcal{E}_{copy} over-approximates \mathcal{E}_{gc} , meaning that any computation that succeeds with \mathcal{E}_{gc} will succeed with an equally-sized or smaller freelist with \mathcal{E}_{gc} . Then we show FO^{gc} is sound with respect to \mathcal{E}_{copy} , and thus by the previous step sound with respect to \mathcal{E}_{gc} .

As mentioned above, we introduce a big step semantics \mathcal{E}_{oper} that does not use freelists or account for garbage collection. We use it to characterize expressions that normalize to values when initialized with a sufficient freelist. This technique has also been employed in earlier work on AARA [25]. In the judgment $V, H \vdash e \Downarrow v, H'$, the "freelist" is the whole ambient set of locations Loc, and thus we never run out of locations during evaluation. This introduces a problem for value and context equivalence: when comparing evaluation results between a run with \mathcal{E}_{copy} and \mathcal{E}_{oper} , the return values might not be syntactically equal. Due to the difference in allocation strategies and the fact that both are nondeterministic, we need a more robust notion of equality for values. Luckily, we can just use the structures from the denotational semantics (defined in Section 4). In both runs, the return value maps to the semantic value [5]. Thus we use structural equality as the basis for value and context equivalence:

Definition 1 (Value Equivalence). Two values v_1, v_2 are equivalent (with the presupposition that they are well-formed w.r.t. heaps H_1, H_2) iff $H_1 \vDash v_1 \mapsto a : A$ and $H_2 \vDash v_2 \mapsto a : A$. Write value equivalence as $v_1 \sim_{H_2}^{H_1} v_2$.

Definition 2 (Context Equivalence). Two contexts $(V_1, H_1), (V_2, H_2)$ are equivalent iff $dom(V_1) = dom(V_2)$ and for all $x \in dom(V_1), V_1(x) \sim_{H_2}^{H_1} V_2(x)$. Write context equivalence as $(V_2, H_2) \sim (V_2, H_2)$

Linear Garbage Collection Cost Semantics To establish the soundness of the type system, we need an intermediary semantics \mathcal{E}_{copy} , which is *semantically linear*. As mentioned in Section 2, this means that locations are treated linearly, that is, no location can be used twice in a program. Variable sharing is achieved via *copying*: the shared value is created by allocating a fresh set of locations from the freelist and copying the locations of the original value one by one. This is also sometimes referred to as deep copying. Let copy(H, L, v, H', v') be a 5-place relation on Heap $\times \mathcal{P}(Loc) \times Val \times Heap \times Val$. Similar to reachability, we write this as H', v = copy(H, L, v) to signify the intended mode for this predicate: (+, +, +, -, -).

$$\frac{v \in \{n, \mathtt{T}, \mathtt{F}, \mathtt{Null}\}}{H, v = copy(H, L, v)} \qquad \qquad \frac{l' \in L \quad H', v = copy(H, L \setminus \{l'\}, H(l))}{H'\{l' \mapsto v\}, l' = copy(H, L, l)}$$

$$\frac{L_1 \sqcup L_2 \subseteq L \quad |L_1| = |dom(reach_H(v_1)|)}{|H_1, v_1' = copy(H, L_1, v_1) \quad H_2, v_2' = copy(H_1, L_2, v_2)}{H_2, \langle v_1', v_2' \rangle = copy(H, L, \langle v_1, v_2 \rangle)}$$

554

Primitives require no cells to copy; a location value is copied recursively; a pair of values is copied sequentially, and the total number of cells required is the size of the reachable set of the value. Now, consider \mathcal{E}_{gc} with the share rule F:Share replaced with the following rule.

$$\begin{array}{c} V(x) = v' \quad L \subseteq F \\ |L| = |dom(reach_H(v'))| \quad H', v'' = copy(H, L, v') \quad V' = (V[x_1 \mapsto v', x_2 \mapsto v'']) \upharpoonright_{FV(e)} \\ F' = F \setminus L \quad g = \{l \in H \mid l \notin F' \cup R \cup locs_{V',H}(e)\} \quad V', H', R, F' \cup g \vdash e \Downarrow v, H'', F'' \\ \hline V, H, R, F \vdash \text{share } x \text{ as } x_1, x_2 \text{ in } e \Downarrow v, H'', F'' \end{array}$$
(E:Share)

To share a variable, we first copy the shared value. The number of cells required is equal to the size of the reachable set from the value. This copying sharing semantics is what justifies the analysis to restitute the potential when matching on a cons node, since even if the node was shared, we had to pay for the cost by copying the node when sharing the original value. Next, we restrict the stack to the appropriate variables. Lastly, any locations not reachable from the current subexpression e are collected. This is for the case when a variable is shared but not used later.

Recall that a *computation* is a pair (\mathcal{C}, e) consisting of a configuration $\mathcal{C} = (V, H, R, F)$ and an expression e. Since the cost semantics can only *preserve* the linearity of a computation, we restrict our attention to computations that are linear initially, and show that \mathcal{E}_{copy} respects the linearity of any initially linear computation. This motivates the following definitions:

Definition 3. (Linear context) Given a context (V, H), let $x, y \in dom(V)$, $x \neq y$, and $r_x = reach_H(V(x))$, $r_y = reach_H(V(y))$. It is *linear* given that $set(r_x)$, $set(r_y)$, and $r_x \cap r_y = \emptyset$.

Where set(S) means S a proper set $(\forall x, S(x) \leq 1)$. Denote this by linearCtxt(V, H). Whenever linearCtxt(V, H) holds, there is at most one path from a variable on the stack V to any location in H. Now we can formalize our intuition for linear computations:

Definition 4 (Linear computation). Given a configuration C = (V, H, R, F) and an expression e, we say the 5-tuple (C, e) is a *computation*; it is a *linear computation* given that dom(V) = FV(e), linearCtxt(V, H), and disjoint $(\{R, F, locs_{V,H}(e)\})$. We write linearComp(V, H, R, F, e) (equivalently linearComp(C, e)) to denote this fact.

Over-Approximation Intuitively, we expect that any terminating compution with \mathcal{E}_{copy} has a corresponding run with \mathcal{E}_{gc} that can be instantiated with an equally-sized or smaller freelist. Although this seems quite straightforward to prove, a complete characterization of the relationship between the space allocations of two runs with each semantics is necessary. To demonstrate the difficulties involved, consider the following proof attempt:

Attempt 1. Let $C_2 = (V, H, R, F)$ be a configuration and (C_2, e) be a linear computation. Given that $C_2 \vdash^{\mathcal{E}_{copp}} e \Downarrow v, H', F'$, for all configurations $C_1 = (W, Y, R, M)$ such that $W, Y \sim V, H$ and |M| = |F|, there exists a triple $(w, Y', M') \in Val \times Heap \times Loc$ such that

$$\mathcal{C}_1 \vdash^{\mathcal{E}_{gc}} e \Downarrow w, Y', M'$$
 and $v \sim^{H'}_{Y'} w$ and $|M'| \ge |F'|$

We proceed with induction on the derivation of the judgment in \mathcal{E}_{copy} . Almost every case goes through, save for E:Let. First, we get $W_1, Y \sim V_1, H$ and we have the following from induction on the first premise:

$$W_1, Y, R', M \vdash^{\mathcal{E}_{gc}} e \Downarrow w_1, Y_1, M_1$$
 and $v_1 \sim^{H_1}_{Y_1} w_1$ and $|M_1| \ge |F_1|$

To instantiate the induction hypothesis on the second premise, we need to show that, among other things, $|M_1 \cup j| \geq |F_1 \cup g|$, where j is the set of collected locations in the \mathcal{E}_{gc} judgment. We cannot show this precisely because g might contain more cells then j due to the linearity of \mathcal{E}_{copy} , thus preventing a piecewise comparison. But of course |j| is always less than |g|, since \mathcal{E}_{gc} doesn't copy to share values! This shows that there is a mismatch between the induction hypothesis and the relationship between the sizes of the respective freelists and the garbage sets. Specifically, we need to know exactly how much larger M_1 is compared to F_1 at any given step.

Having a sense of what is missing, we formulate the criteria which characterize the required equivalence between two configurations, which we call *copy extension*.

Definition 5. A configuration $C_2 = (V_2, H_2, R_2, F_2)$ is a *copy extension* of another configuration $C_1 = (V_1, H_1, R_1, F_1)$ iff

- 1. $V_1, H_1 \sim V_2, H_2$
- 2. There is a proper partition $\gamma : dom(H_1) \setminus F_1 \to \mathcal{P}(dom(H_2) \setminus F_2)$ such that for all $l \in dom(\gamma), |\gamma(l)| = reach_{H_1}(V_1)(l) + R_1(l)$
- 3. For all $l \in dom(\gamma)$, $x \in dom(V_1)$, sequence of directions P which is valid w.r.t. $V_1(x)$, $|reach_{H_2}(V_2(x;P)) \cap \gamma(l)| = reach_{H_1}(V_1(x;P))(l)$.
- 4. For all $l \in dom(\gamma)$, $|\gamma(l) \cap R_2| = R_1(l)$
- 5. $|F_1| = |F_2| + | \oslash (\gamma)|$, where $\oslash (\gamma) = \bigcup_{P \in ec(\gamma)} P \setminus \{rep(P)\}$

Write this as $C_1 \leq C_2$.

The intention is that C_2 is a configuration for an evaluation using \mathcal{E}_{copy} , and C_1 a configuration for \mathcal{E}_{gc} . The first condition is the straightforward context equivalence. The second condition requires the existence of a mapping γ that tells us given a location in $H_1 \setminus F_1$, which locations in $H_2 \setminus F_2$ are shared instances.

For example, consider the expression share x as x_1, x_2 in e and assume the stack is $[\mathbf{x} \mapsto 1]$, and the heap equals $[1 \mapsto \langle 0, \text{Null} \rangle]$, i.e. \mathbf{x} is the list [0]. In an evaluation with \mathcal{E}_{gc} , the stack becomes $[\mathbf{x1} \mapsto 1, \mathbf{x2} \mapsto 1]$, and the heap does not change. With \mathcal{E}_{copy} , we allocate a new location in the heap: $[1 \mapsto \langle 0, \text{Null} \rangle, 2 \mapsto \langle 0, \text{Null} \rangle]$, and the stack changes accordingly: $[\mathbf{x1} \mapsto 1, \mathbf{x2} \mapsto 2]$. Now γ would map 1 to $\{1, 2\}$, since both are shared instances of the former.

Thus, the image of γ is a collection of disjoint subsets whose union is $dom(H_2) \setminus F_2$, and each location in $dom(H_2) \setminus F_2$ belongs to a unique class whose preimage is the unique representative in $dom(H_1) \setminus F_1$. Furthermore, we noticed it is crucial to include the fact that the size of $\gamma(l)$ must be the sum of the number of references from the stack and the continuation set. Furthermore, we also require each subset $\gamma(l)$ (also referred to as class) to be nonempty (this is the *proper* partition condition).

While γ gives us a relation between the active regions of two respective heaps, we still need to know exactly how variables on the stack factor in this relationship. Let $l \in H_1$. Specifically, we need to know that the number of references to l from every *sub*value in V_1 is equal to the size of the corresponding part of the class $\gamma(l)$. First, we need to access subvalues of a value using directions:

Definition 6. Let Dir be the set {L,R,N}, denoting left, right, and next respectively. We define the function $get_H : (1 \oplus Val) \times Dir \rightarrow 1 \oplus Val$ which indexes values via directions:

Let $P \in \mathcal{S}(\mathsf{Dir})$, where $\mathcal{S}(X)$ denotes the set of sequence with elements from X. We define $find_H : (1 \oplus \mathsf{Val}) \times \mathcal{S}(\mathsf{Dir}) \to 1 \oplus \mathsf{Val}$, which extends get_H to sequences of directions:

$$find_H(v, D :: P) = find_H(get_H(v, D), P)$$
$$find_H(v, []) = v$$

Call P valid w.r.t. a value v if $find_H(v, P) = Just(v')$ for some v'. Given a sequence P that is valid w.r.t. V(x), write $V_H(x; P)$ for $from Just(find_H(V(x), P))$ and $reach_H(V(x; P))$ for $reach_H(V_H(x; P))$.

With this, the third condition gives us a more fine grained restriction: for any subvalue in V_1 , the number of references from it to l is equal to the size of the intersection of the reachable set of the corresponding subvalue in V_2 with the appropriate class $\gamma(l)$.

The next condition simply states that the continuation sets respect γ . Lastly, we have that F_1 is greater than F_2 , with the overhead $\oslash(\gamma)$ being exactly the sum $\sum_{l \in \gamma} |\gamma(l)| - 1$. Here $ec(\gamma)$ is the image of γ : $\{\gamma(l) \mid l \in dom(\gamma)\}$. Since each class $\gamma(l)$ is non-empty, we use rep(l) to choose an arbitrary element from the class. Now we can state the key lemma:

Lemma 2. Let $(\mathcal{C}_2, e) = (V, H, R, F, e)$ be a linear computation. Given that $\mathcal{C}_2 \vdash^{\mathcal{E}_{copy}} e \Downarrow v, H', F'$, for all configurations \mathcal{C}_1 such that $\mathcal{C}_1 \preceq \mathcal{C}_2$, there exists a triple $(w, Y', M') \in \mathsf{Val} \times \mathsf{Heap} \times \mathsf{Loc}$ and $\gamma' : dom(Y') \setminus M' \to \mathcal{P}(dom(H') \setminus F')$ s.t.

- 1. $\mathcal{C}_1 \vdash^{\mathcal{E}_{\mathsf{gc}}} e \Downarrow w, Y', M'$
- 2. $v \sim_{Y'}^{H'} w$
- 3. γ' is a proper partition, and for all $l \in dom(\gamma')$, $|\gamma'(l)| = |reach_{Y_1}(w_1)(l)| + S(l)$
- 4. For all $l \in dom(\gamma')$ and $P \in \mathcal{S}(\text{Dir})$ that is valid w.r.t. v, $|reach_{H'}(find_{H'}(v; P)) \cap \gamma'(l)| = reach_{Y'}(find_{Y'}(w; P))(l)$
- 5. For all $l \in dom(\gamma')$, $\gamma'(l) \cap R = \gamma(l) \cap R$
- 6. $|M'| = |F'| + |\oslash(\gamma')|$

Thus we can execute a computation using \mathcal{E}_{gc} given that the computation succeeded in a run with \mathcal{E}_{copy} , which shows that \mathcal{E}_{copy} is an over-approximation of \mathcal{E}_{gc} .

Soundness of FO^{gc} For the second part of the proof, we show FO^{gc} is sound w.r.t. \mathcal{E}_{copy} :

Lemma 3 (Linearity of \mathcal{E}_{copy}). Let $\mathcal{C} = (V, H, R, F)$ be a configuration, $H \models V : \Gamma$, $\Sigma; \Gamma \mid_{q'}^{q} e : B$, and $\mathcal{C} \vdash^{\mathcal{E}_{copy}} e \Downarrow v, H', F'$. Given linearComp (\mathcal{C}, e) , we have set $(reach_{H'}(v))$ and disjoint $(\{R, F', reach_{H'}(v)\})$.

Theorem 4 (Soundness). let $H_o \models V_o : \Gamma, \Sigma; \Gamma \mid_{q'}^q e : B, V_o, H_o \vdash e \Downarrow v_o, H'_o$. Then for all $C \in \mathbb{Q}^+$ and configurations (V, H, R, F) s.t. $V_o, H_o \sim V, H$, linearComp(V, H, R, F, e), and $|F| \ge \Phi_{V,H}(\Gamma) + q + C$, there exists a triple (v, H', F') and a freelist F' s.t. $V, H, R, F \vdash^{\mathcal{E}_{\mathsf{copy}}} e \Downarrow v, H', F', v_o \sim_{H'}^{H'_o} v$, and $|F'| \ge \Phi_{H'}(v : B) + q' + C$

In other words, given a terminating expression (verified by succeeding with the run using \mathcal{E}_{oper}) and given a freelist that is sufficiently large (as predicated by the type derivation), a run with \mathcal{E}_{copy} will normalize to an equivalent value, and the resulting freelist will be sufficiently large (as predicated by the type derivation).

6 Implementation and Evaluation

Implementation We have implemented the novel cost semantics and the type system in Resource Aware ML (RaML). The implementation covers full RaML, including user-defined data types, higher-order functions, and polynomial potential functions. However, there is no destructive match for function closures and analyzing the heap-space usage of closures still amounts to counting allocations only. The main changes that where necessary have been in the rules for sharing and pattern matching as described earlier. We also needed to change some elaboration passes that were no longer cost preserving with the GC cost model.

The garbage collection cost semantics is implemented as an alternative evaluation module inside RaML. As mentioned before, RaML leverages the syntax of OCaml programs. First, we take the OCaml type checked abstract syntax tree and perform a series of transformations. The evaluation modules operate on the resulting RaML syntax tree. In the gc evaluation module, evaluate has the following signature:

Here, the second argument **int** specifies the size of the initial freelist. The result is an option triple of the return value, heap, and freelist; None is returned in case the freelist was not sufficient for the evaluation. Whereas the normal evaluation boxes every value (everything evaluates to a location), the gc module follows the cost semantics and only boxes data constructors. The rationale is that the size for other values can be computed statically and thus stack allocated. One difference between the cost semantics and its implementation is that while in the language presented here list is the only data type, our implementation supports user defined data types. The extension is straightforward except the treatment of the nil constructor, or generally "empty" constructors that have arity zero. For simplicity of presentation, we evaluate all nil constructors to the same null value in the cost semantics. This is natural for lists because all nil constructors are the same, and every list has at most one nil node. However, for custom data types that have more than one kind of empty constructor, it is not possible to map every constructor to the same null value. Thus, the implementation treats all constructors uniformly, so each empty constructor also costs one heap location.

As mentioned before, all functions used in a program are declared in a global mutually recursive block, and we do not account for the constant space overhead for this block in the cost semantics. In order to implement this global function block, we allow closure creation during program evaluation. However, we allocate all closures from a separate freelist into a separate heap. This ensures that data constructors are allocated from the correct freelist and no space overhead is created by allocating closures for function declarations.

Evaluation We evaluated our new analysis on a number of functions. Table 1 contains a representative compilation. It shows the type signature for each function. Table 2 presents the test data that showcase the difference between $\text{RaML}^{\text{heap}}$, the previous RaML type system instantiated with the heap metric (the old analysis which only counts heap allocations), and FO^{gc} , which includes deallocations and copying cost for sharing. For each type system, we show the heap space bound computed by RaML, the number of constraints generated, and the time elapsed during analysis. The last column gives the expression for the exact heap high watermark derived by hand and verified by running the cost semantics.

Except for bfs and dfs, all functions in the table take a *principal* argument of type list. The variables in the table refer to this argument (for example, the type of the principal argument

of quicksort is 'a list). In general, M refers to the number of cons constructors of the principal argument (or the number of *outer* cons nodes in case of nested lists); L refers to the maximum number of cons nodes of the inner lists.

For the sorting functions, aside from mergesort, the new analysis using the gc metric derived asymptotically better bounds when compared to the heap metric. Furthermore, all bounds are *exact* with respect to the cost semantics. In regards to mergesort, the analysis was not able to derive a tight bound due to the limitations of AARA in deriving logarithmic bounds. A particu-

function	type
quicksort	['a ->'a ->bool; 'a list] ->'a list
mergesort	$[['a; 'a] \rightarrow bool; 'a list] \rightarrow 'a list$
ocamlsort	$[['a; 'a] \rightarrow bool; 'a list] \rightarrow 'a list$
selection sort	int list ->int list
eratosthenes	int list ->int list
dfs	[btree; int] ->btree option
bfs	[btree; int] ->btree option
transpose	'a list list ->'a list list
map it	['a ->'b; 'a list list] ->'b list list * 'b list list
pairs	'a list ->('a * 'a) list

Table 1: Signature of Test Functions

larly nice result is that for quicksort, we derive that the space usage is exactly 0, which justifies its use as a zero space-overhead (or "in place") sorting algorithm.

Next, we have have the graph search algorithms operating on a binary tree. Again, the gc metric was able to derive exact space overheads, while the heap metric derived linear bounds for both. For transpose, the gc metric derived an asymptotically better bound, but was not able to derive the exact overhead. We implement matrices as lists-of-lists in row-major order. The transpose function is implemented tail-recursively, with the accumulator starting as the empty list. When "flipping" the first row r of the input and appending this to the accumulator, we need to create |r| many new nil and cons constructors to store the row as a column. While this overhead only occurs once, RaML is unable to infer this from the source code, and thus the cost is repeated over the entire input matrix, resulting in the linear bound (w.r.t. the size of the matrix). This artifact is unrelated to the new extension; it is a limitation due to the implementation of RaML.

The last two functions demonstrate how the gc metric performs when there is variable sharing. map_it maps the input function across each list in the principal argument twice, returning a tuple of nested lists. The gc metric dictates that every outer data constructor in the principal argument needs to be copied, and thus gives the linear bound M + 1. In this case, the bound is exact. The function "pairs" takes a list and outputs all pairs of the input list which are ordered ascending in input position. For example, pairs [1;2;3;4] = [(1,2);(1,3);(1,4);(2,3);(2.4);(3,4)]. For pairs, the gc metric derived a bound that is asymptotically the same as the heap metric, but with better constants. An exact bound could not be derived because the deallocation potential from the pattern match in the definition of pairs is wasted since the matched body could already be typed with zero cost. However, this deallocation is used as usual in the cost semantics. Thus the slack in the bound totals to the size of the input.

7 Conclusion and Future Work

In this article, we introduced a novel operational cost semantics that models a perfect tracing garbage collector and an extension to AARA that is sound with respect to the new semantics. We implemented the new semantics and analysis as modules in RaML and found through experimental testing that the extended AARA was able to derive asymptotically better bounds for several commonly used functions and programming patterns; often, the bounds are optimal

	$\operatorname{RaML}^{heap}$			FO^{gc}			
function	computed bound	constraints	time	computed bound	constraints	time	optimal
quicksort	$1.00 + 3.50M + 1.50M^2$	8515	0.52	0	8519	0.48	0
mergesort	$1.00 - 4.67M + 6.33M^2$	9572	0.64	$-0.50M + 0.50M^2$	9578	0.58	$\log(M)$
ocamlsort	$7.50 + 5.50M + 1.00M^2$	8565	0.51	1.00 + 1.00M	8573	0.50	M + 1
selection sort	$2.00 + 3.00M + 1.00M^2$	639	0.06	0	642	0.05	0
eratosthenes	$1.00 + 1.50M + 0.50M^2$	515	0.06	0	517	0.04	0
dfs	3.00 + 2.00M	5481	0.90	2	5483	0.36	2
bfs	5.00 + 10.00M	24737	4.15	4	24742	1.62	4
transpose	$1.00 + 3.50LM + 0.50LM^2$	10680	0.50	1.00 + 2.00LM	10684	0.50	$\max(0, 2L - 1)$
map it	2.00 + 2.00LM + 4.00M	30699	1.58	1.00M + 1.00	30703	1.57	M + 1
pairs	$1.00 + 1.00M^2$	10214	0.60	$0.50M + 0.50M^2$	10217	0.64	$0.5M^2 - 1.5M + 2$

Table 2: Automatic Bound Analysis with RaML

with respect to the cost semantics.

One direction for future work is using the cost free metric cf to model global garbage collection. In cf, all resource constants, including constructor nodes, are set to 0. A cost-free typing judgment then captures how an expression manipulates the structures in the context into the structure induced by its type. Using this fact, we could express the maximum space usage in the sequential composition $let(e_1; x : \tau.e_2)$ by analyzing e_1 twice—once with the cost-free metric and once with the regular metric—and assign potential to x using the result type in the cost-free typing. In prior work [23], the authors have successfully employed this cost-free metric to analyze parallel programs. Here, the difficulty is showing the simultaneous soundness of both destructive pattern matching and the cost-free composition. Another complication is the choice between local variable sharing and global context sharing. We leave the exploration of this area to future work.

Another direction for future work are function closures. The current treatment in our implementation is unsatisfactory since there is no equivalent to the destructive pattern match for closures. As a result, the GC metric in RaML only accounts for allocation of closures, which is not an improvement over the existing implementation. Ideally, we would like to account for deallocation at function applications and treat closures similar to other data structures in sharing. However, the size of closures cannot be determined easily statically and closures can not capture potential and are currently shared freely in RaML. As a result, the techniques we developed here do not directly carry over to closures.

Finally, we are interested in exploring if our work can be used to improve the efficiency of garbage collection in languages like OCaml. A guaranteed upper bound on the heap space can be used in different ways to control the frequency of the collections and the total memory that is requested from the operating system.

Acknowledgments

This article is based on research supported by the United States Air Force under DARPA AA Contract FA8750-18-C-0092 and DARPA STAC Contract FA8750-15-C-0082, and by the National Science Foundation under SaTC Award 1801369 and SHF Award 1812876. Yue Niu has been supported by Carnegie Mellon's Undergraduate Research Office through a Summer Undergraduate Research Fellowship (SURF). Any opinions, findings, and conclusions contained in this document are those of the authors and do not necessarily reflect the views of the sponsoring organizations.

References

- Elvira Albert, Jesús Correas Fernández, and Guillermo Román-Díez. Non-cumulative Resource Analysis. In Tools and Algorithms for the Construction and Analysis of Systems - 21st International Conference, (TACAS'15), 2015.
- [2] Elvira Albert, Samir Genaim, and Miguel Gomez-Zamalloa. Heap space analysis for java bytecode. In Proceedings of the 6th International Symposium on Memory Management (ISMM'07), 2007.
- [3] Elvira Albert, Samir Genaim, and Miguel Gómez-Zamalloa Gil. Live Heap Space Analysis for Languages with Garbage Collection. In Proceedings of the 2009 International Symposium on Memory Management (ISMM'09), 2009.
- [4] Elvira Albert, Samir Genaim, and Miguel Gómez-Zamalloa. Heap space analysis for garbage collected languages. Science of Computer Programming, 78(9):1427 – 1448, 2013.
- [5] Robert Atkey. Amortised Resource Analysis with Separation Logic. In 19th Euro. Symp. on Prog. (ESOP'10), 2010.
- [6] Martin Avanzini and Georg Moser. A Combination Framework for Complexity. In 24th International Conference on Rewriting Techniques and Applications (RTA'13), 2013.
- [7] Régis Blanc, Thomas A. Henzinger, Thibaud Hottelier, and Laura Kovács. ABC: Algebraic Bound Computation for Loops. In Logic for Prog., AI., and Reasoning - 16th Int. Conf. (LPAR'10), 2010.
- [8] Víctor A. Braberman, Federico Fernández, Diego Garbervetsky, and Sergio Yovine. Parametric prediction of heap memory requirements. In 7th Int. Symp. on Memory Management (ISMM'08), pages 141–150, 2008.
- [9] Quentin Carbonneaux, Jan Hoffmann, Thomas Reps, and Zhong Shao. Automated Resource Analysis with Coq Proof Objects. In 29th International Conference on Computer-Aided Verification (CAV'17), 2017.
- [10] Quentin Carbonneaux, Jan Hoffmann, and Zhong Shao. Compositional Certified Resource Bounds. In 36th Conference on Programming Language Design and Implementation (PLDI'15), 2015. Artifact submitted and approved.
- [11] Arthur Charguéraud and François Pottier. Machine-Checked Verification of the Correctness and Amortized Complexity of an Efficient Union-Find Implementation. In Interactive Theorem Proving - 6th International Conference (ITP'15), 2015.
- [12] Wei-Ngan Chin, Huu Hai Nguyen, Corneliu Popeea, and Shengchao Qin. Analysing Memory Resource Bounds for Low-level Programs. In Proceedings of the 7th International Symposium on Memory Management (ISMM'08), 2008.
- [13] Wei-Ngan Chin, Huu Hai Nguyen, Shengchao Qin, and Martin Rinard. Memory Usage Verification for OO Programs. In Proceedings of the 12th International Conference on Static Analysis (SAS'05), 2005.
- [14] Norman Danner, Daniel R. Licata, and Ramyaa Ramyaa. Denotational Cost Semantics for Functional Languages with Inductive Types. In 29th Int. Conf. on Functional Programming (ICFP'15), 2012.
- [15] Ankush Das, Jan Hoffmann, and Frank Pfenning. Parallel complexity analysis with temporal session types. In 23rd International Conference on Functional Programming (ICFP'18), 2018. Conditionally accepted.
- [16] Antonio Flores-Montoya and Reiner Hähnle. Resource Analysis of Complex Programs with Cost Equations. In Programming Languages and Systems - 12th Asian Symposiu (APLAS'14), 2014.
- [17] Florian Frohn, M. Naaf, Jera Hensel, Marc Brockschmidt, and Jürgen Giesl. Lower Runtime Bounds for Integer Programs. In Automated Reasoning - 8th International Joint Conference (IJ-CAR'16), 2016.
- [18] Sumit Gulwani, Krishna K. Mehra, and Trishul M. Chilimbi. SPEED: Precise and Efficient Static Estimation of Program Computational Complexity. In 36th ACM Symp. on Principles of Prog. Langs. (POPL'09), 2009.

- [19] Robert Harper. Practical Foundations for Programming Languages. Cambridge University Press, 2016.
- [20] Jan Hoffmann. Types with Potential: Polynomial Resource Bounds via Automatic Amortized Analysis. PhD thesis, Ludwig-Maximilians-Universität München, 2011.
- [21] Jan Hoffmann, Klaus Aehlig, and Martin Hofmann. Multivariate Amortized Resource Analysis. In 38th Symposium on Principles of Programming Languages (POPL'11), 2011.
- [22] Jan Hoffmann, Ankush Das, and Shu-Chun Weng. Towards Automatic Resource Bound Analysis for OCaml. In 44th Symposium on Principles of Programming Languages (POPL'17), 2017.
- [23] Jan Hoffmann and Zhong Shao. Automatic static cost analysis for parallel programs. In Proceedings of the 24th European Symposium on Programming on Programming Languages and Systems -Volume 9032, pages 132–157, New York, NY, USA, 2015. Springer-Verlag New York, Inc.
- [24] Martin Hofmann and Steffen Jost. Static Prediction of Heap Space Usage for First-Order Functional Programs. In 30th ACM Symp. on Principles of Prog. Langs. (POPL'03), 2003.
- [25] Martin Hofmann and Steffen Jost. Static prediction of heap space usage for first-order functional programs. In Proceedings of the 30th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '03, pages 185–197, New York, NY, USA, 2003. ACM.
- [26] Martin Hofmann and Steffen Jost. Type-Based Amortised Heap-Space Analysis. In 15th Euro. Symp. on Prog. (ESOP'06), 2006.
- [27] Steffen Jost, Kevin Hammond, Hans-Wolfgang Loidl, and Martin Hofmann. Static Determination of Quantitative Resource Usage for Higher-Order Programs. In 37th ACM Symp. on Principles of Prog. Langs. (POPL'10), 2010.
- [28] Zachary Kincaid, Jason Breck, Ashkan Forouhi Boroujeni, and Thomas Reps. Compositional recurrence analysis revisited. In Conference on Programming Language Design and Implementation (PLDI'17), 2017.
- [29] Ugo Dal Lago and Marco Gaboardi. Linear Dependent Types and Relative Completeness. In 26th IEEE Symp. on Logic in Computer Science (LICS'11), 2011.
- [30] Yasuhiko Minamide. Space-profiling semantics of the call-by-value lambda calculus and the CPS transformation. *Electr. Notes Theor. Comput. Sci.*, 26:105–120, 1999.
- [31] Greg Morrisett, Matthias Felleisen, and Robert Harper. Abstract models of memory management. In Proceedings of the Seventh International Conference on Functional Programming Languages and Computer Architecture (FPCA'95), 1995.
- [32] Van Chan Ngo, Mario Dehesa-Azuara, Matthew Fredrikson, and Jan Hoffmann. Verifying and Synthesizing Constant-Resource Implementations with Types. In 38th IEEE Symposium on Security and Privacy (S&P '17), 2017.
- [33] Tobias Nipkow. Amortized Complexity Verified. In Interactive Theorem Proving 6th International Conference (ITP'15), 2015.
- [34] Lars Noschinski, Fabian Emmes, and Jürgen Giesl. Analyzing Innermost Runtime Complexity of Term Rewriting by Dependency Pairs. J. Autom. Reasoning, 51(1):27–56, 2013.
- [35] Ivan Radiček, Gilles Barthe, Marco Gaboardi, Deepak Garg, and Florian Zuleger. Monadic Refinements for Relational Cost Analysis. Proc. ACM Program. Lang., 2(POPL), 2017.
- [36] Hugo R. Simões, Pedro B. Vasconcelos, Mário Florido, Steffen Jost, and Kevin Hammond. Automatic Amortised Analysis of Dynamic Memory Allocation for Lazy Functional Programs. In 17th Int. Conf. on Funct. Prog. (ICFP'12), 2012.
- [37] Moritz Sinn, Florian Zuleger, and Helmut Veith. A Simple and Scalable Approach to Bound Analysis and Amortized Complexity Analysis. In Computer Aided Verification - 26th Int. Conf. (CAV'14), 2014.
- [38] Daniel Spoonhower, Guy E. Blelloch, Robert Harper, and Phillip B. Gibbons. Space profiling for parallel functional programs. In *Proceedings of the 13th ACM SIGPLAN International Conference* on Functional Programming, ICFP '08, pages 253–264, New York, NY, USA, 2008. ACM.

- [39] Leena Unnikrishnan and Scott D. Stoller. Parametric Heap Usage Analysis for Functional Programs. In Proceedings of the 2009 International Symposium on Memory Management (ISMM '09), 2009.
- [40] Leena Unnikrishnan, Scott D. Stoller, and Yanhong A. Liu. Optimized Live Heap Bound Analysis. In Verification, Model Checking, and Abstract Interpretation, 4th International Conference (VMCAI'03), pages 70–85, 2003.
- [41] Pedro B. Vasconcelos, Steffen Jost, Mário Florido, and Kevin Hammond. Type-Based Allocation Analysis for Co-recursion in Lazy Functional Languages. In 24th European Symposium on Programming (ESOP'15), 2015.
- [42] P. Wang, D. Wang, and A. Chlipala. TiML: A Functional Language for Practical Complexity Analysis with Invariants. In Object-Oriented Prog., Syst., Lang., and Applications (OOPSLA'17), 2017.
- [43] Ezgi Çiçek, Gilles Barthe, Marco Gaboardi, Deepak Garg, and Jan Hoffmann. Relational Cost Analysis. In 44th Symposium on Principles of Programming Languages (POPL'17), 2017.