## **Computer Science**

Abstract Models of Memory Management

Greg Morrisett

Matthias Felleisen January 1995 CMU-CS-95-110 Robert Harper



# Carnegie Mellon

19950317 123

DISTRIBUTION STATEMENT A

Approved for public release; Distribution Unlimited

## Abstract Models of Memory Management

Greg Morrisett

Matthias Felleisen January 1995

CMU-CS-95-110

Robert Harper

School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213



Also published as Fox Memorandum CMU-CS-FOX-95-01

Accesion For			
NTIS CRA&I DTIC TAB Unannounced  Justification			
By			
Availability Codes			
Dist	Avail and <b>/ or</b> Specia <b>l</b>		
A-L			

This research was sponsored by the Defense Advanced Research Projects Agency, CSTO, under the title "The Fox Project: Advanced Development of Systems Software", ARPA Order No. 8313, issued by ESD/AVS under Contract No. F19628-91-C-0168.

The views and conclusions contained in this document are those of the authors and should not be interpreted as representing official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the U.S. Government.

Approved for public release;
Distribution Unlimited

 $\textbf{Keywords:} \ \texttt{garbage} \ \texttt{collection}, \\ \texttt{allocation}, \\ \texttt{type} \ \texttt{theory}, \\ \texttt{polymorphism}$ 

#### Abstract

Most specifications of garbage collectors concentrate on the low-level algorithmic details of how to find and preserve accessible objects. Often, they focus on bit-level manipulations such as "scanning stack frames," "marking objects," "tagging data," etc. While these details are important in some contexts, they often obscure the more fundamental aspects of memory management: what objects are garbage and why?

We develop a series of calculi that are just low-level enough that we can express allocation and garbage collection, yet are sufficiently abstract that we may formally prove the correctness of various memory management strategies. By making the heap of a program syntactically apparent, we can specify memory actions as rewriting rules that allocate values on the heap and automatically dereference pointers to such objects when needed. This formulation permits the specification of garbage collection as a relation that removes portions of the heap without affecting the outcome of the evaluation.

Our high-level approach allows us to compactly specify and prove correct a wide variety of memory management techniques, including standard trace-based garbage collectors (*i.e.*, the family of copying and mark/sweep collectors), generational collection, and type-based tag-free collection. Furthermore, since the definition of garbage is based on the *semantics* of the underlying language instead of the conservative approximation of inaccessibility, we are able to specify and formally prove the idea that type inference can be used to collect some objects that are accessible but never used.

## 1 Memory Safety

Advanced programming languages manage memory allocation and deallocation automatically. Automatic memory managers, or garbage collectors, significantly facilitate the programming process because programmers can rely on the language implementation for the delicate tasks of finding and freeing unneeded objects. Indeed, the presence of a garbage collector ensures *memory safety* in the same way that a type system guarantees *type safety*: no program in an advanced programming language will crash due to dangling pointer problems while pointer allocation, access, and deallocation is transparent. However, in contrast to type systems, memory management strategies and particularly garbage collectors rarely come with a compact formulation and a formal proof of soundness. Indeed, since garbage collectors work on the machine representations of abstract values, the very idea of providing a proof of memory safety sounds unrealistic given the lack of simple models of memory operations.

The recently developed syntactic approaches to the specification of language semantics by Felleisen and Hieb [11] and Mason and Talcott [23, 24] are the first execution models that are intensional enough to permit the specification of memory management actions and yet are sufficiently abstract to permit compact proofs of important properties. Starting from the  $\lambda_v$ -S calculus of Felleisen and Hieb, we design compact specifications of a number of memory management ideas and prove several correctness theorems.

The basic idea underlying the development of our garbage collection calculi is the representation of a program's run-time memory as a global series of syntactic declarations. The program evaluation rules allocate large objects in the global declaration, which represents the heap, and automatically dereference pointers to such objects when needed. As a result, garbage collection can be specified as any relation that removes portions of the current heap without affecting the result of a program's execution.

In Section 2, we present a small functional programming language,  $\lambda gc$ , with a rewriting semantics that makes allocation explicit. In Section 3, we define a semantic notion of garbage collection for  $\lambda gc$  and show that there is no *optimal* collection strategy that is computable. In Section 4, we specify the "free-variable" garbage collection rule which models tracing-based collectors including mark/sweep and copying collectors. We prove that the free-variable rule is correct and provide two "implementations" at the syntactic level: the first corresponds to a copying collector, the second to a generational one. In the context of generational collection, we also illustrate how mutable reference cells render the original generational collector unsound and how the problem can be fixed at an abstract level.

In Section 5, we show how the use of abstract syntax facilitated our work in the first few sections. Specifically, the abstract syntax of  $\lambda$ gc implicitly carries the *shape* (size and pointer information) of allocated objects, which permits specifications of GC algorithms that ignore the problem of traversing the memory graph in the absence of guiding information. To illustrate how to refine  $\lambda$ gc just enough to address this problem, we formalize so-called "tag-free" collection algorithms for explicitly-typed, monomorphic languages such as Pascal and Algol [7, 36, 8]. We show how to reconstruct necessary shape information about values from types during garbage collection. We are able to prove the correctness of the garbage collection algorithm by using a well known type preservation argument. We then sketch how the type-based garbage collection algorithms can be extended to a language with a polymorphic type system in the style of the Girard-Reynold's System F [13, 14, 30]. Our formulation leads to a fairly simple proof of the correctness of Tolmach's garbage collector for the Gallium ML compiler [34].

In Section 6, we justify our semantic definition of garbage by showing that Milner-style type inference can be used to prove that an object is semantically garbage even though the object is still

reachable. While previous authors have sketched this idea [6, 3, 16, 12], we are the first to present a formal proof of this result. The proof is obtained by casting the well-known interpretation of types as logical relations into our framework.

Section 7 discusses related work and Section 8 closes with a summary.

#### Programs:

```
(variables)
                      x, y, z \in
                                                  ::= \cdots \mid -2 \mid -1 \mid 0 \mid 1 \mid 2 \mid \cdots
(integers)
                                 \in
                                       Int
(expressions)
                                 \in
                                       Exp
                                                  ::= x \mid i \mid \langle e_1, e_2 \rangle \mid \pi_1 \mid e \mid \pi_2 \mid e \mid \lambda x.e \mid e_1 \mid e_2
(heap values)
                                 \in
                                       Hval
                                                 ::= i \mid \langle x_1, x_2 \rangle \mid \lambda x.e
(heaps)
                           H
                                                        \{x_1=h_1,\ldots,x_n=h_n\}
                                 \in
                                       Heap
                                                 ::=
                                       Prog
(programs)
                                                 ::= letrec H in e
                                                         letrec H in x
(answers)
                                       Ans
```

#### **Evaluation Contexts and Instruction Expressions:**

```
(contexts) E \in Ctxt ::= [] |\langle E, e \rangle | \langle x, E \rangle | \pi_i E | E e | x E |
(instructions) I \in Instr ::= h | \pi_i x | x y
```

Figure 1: The Syntax of  $\lambda gc$ 

## 2 Modeling Allocation: $\lambda gc$

Syntax: The syntax of  $\lambda gc$  (see Figure 1) is that of a conventional, higher-order, applicative programming language based on the  $\lambda$ -calculus. Following the tradition of functional programming, a  $\lambda gc$  program (P) consists of some mutually recursive definitions (H) and an expression (e). The global definitions are useful for defining mutually recursive procedures, but their primary purpose here is to represent the run-time heap of a program. In general, there can be cycles in a heap so we use letrec instead of let as the binding form. Expressions are either variables (x), integers (i), pairs  $(\langle e_1, e_2 \rangle)$ , projections  $(\pi_i)$ , abstractions  $(\lambda x.e)$ , or applications  $(e_1 e_2)$ .

```
FV(i) = \emptyset
FV(x) = \{x\}
FV(\langle e_1, e_2 \rangle) = FV(e_1) \cup FV(e_2)
FV(\pi_i e) = FV(e)
FV(\lambda x.e) = FV(e) \setminus \{x\}
FV(e_1 e_2) = FV(e_1) \cup FV(e_2)
FV(\{x_1 = h_1, \dots, x_n = h_n\}) = (FV(h_1) \cup \dots \cup FV(h_n)) \setminus \{x_1, \dots, x_n\}
FV(\text{letrec } H \text{ in } e) = (FV(H) \cup FV(e)) \setminus Dom(H)
```

Figure 2: Calculating Free Variables

Formally, the heap is a series of pairs, called *bindings*, consisting of variables and heap values. Heap values are a semantically significant subset of expressions. The order of the bindings is

irrelevant. In addition, each variable must be bound to at most one heap value in a heap. Hence, we treat heaps as sets and, when convenient, as finite functions. We write Dom(H) to denote the bound variables of H and Rng(H) to denote the heap values bound in H. We sometimes refer to variables bound in the heap as locations or pointers.

The language contains two binding constructs:  $\lambda x.e$  binds x in e and letrec H in e binds the variables in Dom(H) to the expressions Rng(H) in both the values in Rng(H) and in e. Following convention, we consider expressions to be equivalent up to a consistent  $\alpha$ -conversion of bound variables. The free variables of expressions, heaps, and programs are determined by the definitions of Figure 2.

Considering programs equivalent modulo  $\alpha$ -conversion and the treatment of heaps as sets instead of sequences hides many of the complexities of memory management. In particular, programs are automatically considered equivalent if the heap is re-arranged and locations are re-named as long as the "graph" of the program is preserved. This abstraction allows us to focus on the issues of determining what bindings in the heap are garbage without specifying how such bindings are represented in a real machine.

Mathematical Notation: We use  $X \uplus X'$  to denote the union of two disjoint sets, X and X'. We use  $H \uplus H'$  to denote the union of two heaps whose domains are disjoint. We use  $\{e_1/x\}e_2$  to denote capture-avoiding substitution of the expression  $e_1$  for the free variable x in the expression  $e_2$ . We use  $X \setminus X'$  to denote  $\{x \in X \mid x \notin X'\}$ .

Semantics: The rewriting semantics for  $\lambda gc$  is an adaptation of the standard reduction function of the  $\lambda_v$ -S calculus [11]. Roughly speaking, this kind of semantics describes an abstract machine whose states are programs and whose instructions are relations between programs. The desired final state of this abstract machine is an answer program (A) whose body is a pointer to some value, such as an integer, in the heap.

Each rewriting step of a program letrec H in e proceeds according to a simple algorithm. If the body of the program, e, is not a variable, it is partitioned into an evaluation context E (an expression with a hole  $[\ ]$  in the place of a sub-expression), which represents the control state, and an instruction expression I, which roughly corresponds to a program counter: e = E[I]. The instruction expression determines the next expression e' and any changes to the heap resulting in a new heap H'. Putting the pieces together yields the next program in the evaluation sequence: letrec H' in E[e']. Each instruction determines one transition rule of the abstract machine. Formally, a rule denotes a relation between programs. A set of rules denotes the union of the respective relations.

We use the following conventions: Let G be a set of program relations, and let P and P' be programs. Then,

 $P \stackrel{G}{\longmapsto} P'$  means P rewrites to P' according to one of the rules in G and  $\stackrel{G}{\longmapsto}^*$  is the reflexive, transitive closure of  $\stackrel{G}{\longmapsto}$ .

G+r is the union of G with the rule r.

A program P is canonical with respect to G iff there is no rule in G and no P' such that  $P \stackrel{G}{\longmapsto} P'$ .

 $P \Downarrow_G P'$  means that  $P \stackrel{\mathsf{G}}{\longmapsto} {}^*P'$  and P' is canonical with respect to G.

 $P \uparrow_{G}$  means that there exists an infinite sequence of programs  $P_{i}$  such that  $P \stackrel{G}{\longmapsto} P_{1} \stackrel{G}{\longmapsto} \cdots$ .

Figure 1 defines the set of evaluation contexts and instruction expressions for  $\lambda gc$ . The definition of evaluation contexts (E) reflect the left-to-right evaluation order of the language. All terms to the left of the path from the root to the hole are variables; the terms on the right are arbitrary. Instruction expressions (I) consist of heap values (h), applications of (pointers to heap-allocated) procedures to (pointers to heap-allocated) values, and projections of (pointers to heap-allocated) tuples.

```
 \begin{array}{ll} (\textbf{alloc}) & \text{letrec } H \text{ in } E[h] \xrightarrow{\text{alloc}} \text{letrec } H \uplus \{x=h\} \text{ in } E[x] \\ (\textbf{proj}) & \text{letrec } H \text{ in } E[\pi_i \ x] \xrightarrow{\pi_i} \text{letrec } H \text{ in } E[x_i] \\ (\textbf{app}) & \text{letrec } H \text{ in } E[x \ y] \xrightarrow{\text{app}} \text{letrec } H \uplus \{z=H(y)\} \text{ in } E[e] \\ \end{array}   \begin{array}{ll} (H(x) = \langle x_1, x_2 \rangle \text{ and } i=1,2) \\ (H(x) = \lambda z.e) \end{array}
```

Figure 3: Rewriting Rules for  $\lambda gc$ 

The evaluation rules for  $\lambda gc$  (see Figure 3) reflect the intentions behind our choice of instruction expressions. The transition rule **alloc** models the allocation of values in the heap by binding the value to a new variable and using this variable in its place in the program. Note that the " $\forall$ " notation carries the implicit requirement that the newly allocated variable x cannot be in the domain of the heap H. The transition **proj** specifies how a projection instruction extracts the appropriate component from a pointer to a heap-allocated pair. Similarly, **app** is a transliteration of the conventional  $\beta$ -value rule into our modified setting. It binds the formal parameter of a heap-allocated procedure to the value of the pointer given as the actual argument, and places the expression part of the procedure into the evaluation context. Multiple applications of the same procedure require  $\alpha$ -conversion to ensure that the formal parameter does not conflict with bindings already in the heap<sup>1</sup>. We use R to abbreviate the union of the rules **alloc**, **proj**, and **app**.

As a simple example of evaluation under R, consider rewriting the program  $P = \text{letrec } \{x = 1\}$  in  $(\lambda y.\pi_1 y)\langle x, x \rangle$ :

```
\begin{array}{cccc} \mathsf{letrec}\ \{x=1\}\ \mathsf{in}\ (\lambda y.\pi_1\ y)\langle x,x\rangle & \stackrel{\mathsf{alloc}}{\longmapsto} & \mathsf{letrec}\ \{x=1,z=\lambda y.\pi_1\ y\}\ \mathsf{in}\ z\ \langle x,x\rangle \\ & \stackrel{\mathsf{alloc}}{\longmapsto} & \mathsf{letrec}\ \{x=1,z=\lambda y.\pi_1\ y,w=\langle x,x\rangle\}\ \mathsf{in}\ z\ w \\ & \stackrel{\mathsf{app}}{\longmapsto} & \mathsf{letrec}\ \{x=1,z=\lambda y.\pi_1\ y,w=\langle x,x\rangle,y=\langle x,x\rangle\}\ \mathsf{in}\ \pi_1\ y \\ & \stackrel{\pi_1}{\longmapsto} & \mathsf{letrec}\ \{x=1,z=\lambda y.\pi_1\ y,w=\langle x,x\rangle,y=\langle x,x\rangle\}\ \mathsf{in}\ x \end{array}
```

We can break P into an evaluation context  $E_0 = [\ ]\langle x,x\rangle$  and instruction  $I_0 = \lambda y.\pi_1 \ y$ . Since  $I_0$  is a heap value, alloc applies. We choose a new variable z, bind the heap value to z, and replace it with z in the hole of  $E_0$ , resulting in the program letrec  $\{x=1,z=\lambda y.\pi_1 \ y\}$  in  $z\ \langle x,x\rangle$ . This program can be broken into evaluation context  $E_1 = z\ [\ ]$  and instruction  $I_1 = \langle x,x\rangle$ .  $I_1$  is also a heap value, so alloc applies, we choose a new variable w, bind the heap value to w, and replace the instruction with w, resulting in the program letrec  $\{x=1,z=\lambda y.\pi_1 \ y,w=\langle x,x\rangle\}$  in  $z\ w$ . This program can be broken into an empty context  $(E_2 = [\ ])$  and instruction  $I_2 = z\ w$ . Since z is bound to a procedure and w is bound in the heap, app applies. The heap value that w is bound to  $(\langle x,x\rangle)$  is bound to the formal argument of the procedure (y) in the heap and the body of the procedure replaces the instruction, resulting in the program letrec  $\{x=1,z=\lambda y.\pi_1\ y,w=\langle x,x\rangle,y=\langle x,x\rangle\}$  in  $\pi_1\ y$ . This program can be broken into an empty evaluation context  $(E_3 = [\ ])$  and instruction  $I_3 = \pi_1\ y$ . Since y is bound to a pair, the **proj** rule applies and the instruction is replaced with the first

<sup>&</sup>lt;sup>1</sup>An alternative rule for application substitutes the actual argument (y) for the formal (z) within e and performs no allocation. This rule is essentially equivalent to app, but the definition above simplifies the proofs of Section 6.

component of the pair, namely x, resulting in the answer program letrec  $\{x = 1, z = \lambda y.\pi_1 \ y, w = \langle x, x \rangle, y = \langle x, x \rangle \}$  in x.

The canonical programs of  $\lambda gc$  are either answers or stuck programs. The latter correspond to machine states that result from the misapplication of primitive program operations or unbound variables.

Definition 2.1 (Stuck Programs) A program is stuck if it is of one of the following forms:

- letrec H in  $E[\pi_i \ x]$   $(x \notin Dom(H) \ or \ H(x) \neq \langle x_1, x_2 \rangle) \ or$
- letrec H in E[x y]  $(x \text{ or } y \notin Dom(H) \text{ or } H(x) \neq \lambda z.e).$

All programs either diverge or evaluate to an answer or a stuck program. Put differently, the evaluation process defines a partial function from  $\lambda gc$  programs to canonical programs [11, 37].

**Theorem 2.2** If P is a closed program, then there exists at most one P' such that  $P \Downarrow_R P'$ .

**Proof** (sketch): The theorem relies on the following lemma which shows that at each evaluation step, the instruction expression is uniquely determined.

Lemma 2.3 (Unique Decomposition) If P = letrec H in e is a  $\lambda gc$  program, then either P is stuck, P is an answer, or else there exists a unique evaluation context E and instruction I such that e = E[I].

**Proof** (sketch): We show by induction on the structure of e that either e is a variable or there exists a unique E and I, such that e = E[I]. Here, we give one case of the induction as an example: suppose  $e = (e_1 \ e_2)$ . There are four cases to consider. If  $e_1$  is not a variable, then by induction, there exists an E', I, such that  $e_1 = E'[I]$ , in which case we take  $E = E' \ e_2$ . If  $e_1$  is a variable, then  $e_1 = x$  for some x and  $H(x) = \lambda z.e'$  else e would be stuck. Suppose  $e_2$  is not a variable. Then by induction,  $e_2 = E'[I]$  and we can take E = x E'. Otherwise,  $e_2 = y$  and we can take E = [I] and  $I = (e_1 \ e_2)$ . Other cases follow similarly.

The rest follows from the lemma in a straightforward manner, given that P is closed.  $\Box$ 

## 3 A Semantic Definition of Garbage

Since the semantics of  $\lambda gc$  makes the allocation of values explicit, including the implicit pointer dereferencing in the language, we can also define what it means to garbage collect a value in the heap and then analyze some basic properties. A binding x=h in the heap of a program is garbage if removing the binding has no "observable" effect on running the program. In our case, we consider only integer results and non-termination to be observable.

**Definition 3.1** (Kleene Equivalence)  $(P_1, G_1) \simeq (P_2, G_2)$  means  $P_1 \Downarrow_{G_1}$  letter  $H_1$  in x where  $H_1(x) = i$  if and only if  $P_2 \Downarrow_{G_2}$  letter  $H_2$  in y and  $H_2(y) = i$ . If  $G_1 = G_2 = R$ , then we simply write  $P_1 \simeq P_2$ .

A binding is *garbage* if removing it results in a program that is Kleene equivalent to the original program:

**Definition 3.2** (Garbage) Given a program  $P = \text{letrec } H \uplus \{x = h\} \text{ in } e, \text{ the binding } x = h \text{ is garbage } with \text{ respect to } P \text{ iff } P \simeq \text{letrec } H \text{ in } e.$ 

A collection of a program is the same program with some garbage bindings removed. An optimal collection of a program is a program with as many garbage bindings removed as possible.

**Definition 3.3 (Collection, Optimal Collection)** The program letrec H in e is a collection of  $P = \text{letrec } H \uplus H'$  in e iff all bindings x = h in H' are garbage with respect to P. The program letrec H in e is an optimal collection of P iff no binding in H is garbage with respect to P.

Unfortunately, there can be no optimal garbage collector because determining whether a binding is garbage or not is undecidable.

**Proposition 3.4** (Garbage Undecidable) Determining if a binding is garbage in an arbitrary closed  $\lambda gc$  program is undecidable.

**Proof:** Since  $\lambda$ gc is essentially an untyped  $\lambda$ -calculus, determining whether an arbitrary closed program P =letrec H in e halts with an answer or not is undecidable. Assume that GC is an algorithm for determining whether or not some binding in the heap is garbage. Then, in particular, GC can tell whether or not  $x = \langle x_1, x_2 \rangle$  is garbage in the program:

$$P' =$$
letrec  $H \uplus \{x_1 = 1, x_2 = 2, x = \langle x_1, x_2 \rangle \}$  in  $(\lambda y.\pi_1 \ x) \ e$ 

where x does not occur in Dom(H). We will show that GC claims that the binding for x is garbage if and only if P terminates with an answer, which proves that GC cannot exist.

Assume GC claims the binding for x is garbage. Then, the program letrec H in  $(\lambda y.\pi_1 \ x)e$  is Kleene equivalent to P', which clearly implies that P' either diverges or gets stuck, because this program will always either diverge or get stuck. To see this, note that evaluation will become stuck when x is dereferenced by the projection operation. Given an infinite sequence of reduction steps for P', we can easily produce an infinite sequence of reduction steps for P, since this sequence must occur while evaluating e. Similarly, given a sequence of reduction steps which takes e0 to a stuck program, we can produce a sequence of reduction steps that causes e1 to become stuck, because e2 must become stuck while evaluating e3, since e3 is still bound in e4. Consequently, if e6 claims e8 is garbage, then e9 either diverges or gets stuck.

Now assume GC claims that the binding for x is not garbage. This implies the existence of a heap value h and variable z such that

$$P' \xrightarrow{\mathbb{R}} \text{*letrec } H' \uplus \{x_1 = 1, x_2 = 2, x = \langle x_1, x_2 \rangle, z = h\} \text{ in } (\lambda y. \pi_1 x) z$$

Otherwise, x would never be involved in a reduction step. Given this finite reduction sequence for P', it follows that  $P \stackrel{\mathbb{R}}{\longmapsto} *$ letrec  $H' \uplus \{z = h\}$  in z In short, the evaluation of P terminates with an answer given the assumption that x is not garbage.

Since the two cases show that GC claims x is not garbage if and only if P terminates with an answer, we have established the proposition.

## 4 Reachability-Based Garbage Collection

Since computing an optimal collection is undecidable, a garbage collection algorithm must conservatively approximate the set of garbage bindings. Most garbage collectors compute the *reachable* set of bindings in a program, given the variables in use in the current instruction expression and control state. All reachable bindings are preserved; the others are eliminated.

In Section 4.1 we present a general specification for reachability-based garbage collection and prove that it only collects true garbage. In Section 4.2 we give a high-level description of an algorithm based on tracing collection and prove that the algorithm satisfies our specification. In Section 4.3 we specialize our algorithm to one that takes advantage of a generational partition and in Section 4.4, we show how generational collection works in the presence of assignment.

#### 4.1 The Free-Variable Rule

Following Felleisen and Hieb [11], reachability in  $\lambda gc$  is formalized by considering free variables. The following "free-variable" GC rule describes bindings as garbage if there are no references to these bindings in the other bindings, nor in the currently evaluating expression:

(fv) letrec 
$$H_1 \uplus H_2$$
 in  $e \stackrel{\text{fv}}{\longmapsto}$  letrec  $H_1$  in  $e \pmod{H_2} \cap FV$ (letrec  $H_1$  in  $e) = \emptyset$ )

The fv rule is correct in that it only removes garbage, and thus computes valid collections. The keys to the proof of correctness of fv are a postponement lemma and a diamond lemma. The statements of these lemmas can be summarized by the diagrams of Figure 4 respectively, where solid arrows denote relations that are assumed to exist and dashed arrows denote relations that can be derived from the assumed relations.

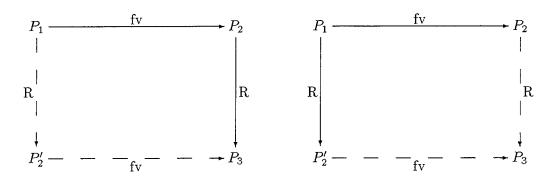


Figure 4: Postponement and Diamond Properties

**Lemma 4.1 (Postponement)** If  $P_1 \xrightarrow{fv} P_2 \xrightarrow{\mathbb{R}} P_3$ , then there exists a  $P_2'$  such that  $P_1 \xrightarrow{\mathbb{R}} P_2' \xrightarrow{fv} P_3$ .

**Proof:** The proof proceeds by case analysis on the elements of R:

alloc: Assume letrec  $H_1 \uplus H_2$  in  $E[h] \xrightarrow{\text{fv}} \text{letrec } H_1 \text{ in } E[h] \xrightarrow{\text{alloc}} \text{letrec } H_1 \uplus \{x = h\} \text{ in } E[x]$ . We only need to show that  $\mathbf{fv}$  applies after performing the allocation:

letrec 
$$H_1 \uplus H_2$$
 in  $E[h] \stackrel{ ext{alloc}}{\longmapsto}$  letrec  $H_1 \uplus H_2 \uplus \{x=h\}$  in  $E[x]$ 

But this is clear since allocation only adds x to the locations of the resulting program, and x cannot be in  $Dom(H_2)$  by the side-condition of the alloc rule. Consequently:

letrec 
$$H_1 \uplus H_2 \uplus \{x=h\}$$
 in  $E[x] \xrightarrow{\text{fv}} \text{letrec } H_1 \uplus \{x=h\}$  in  $E[x]$ 

**proj:** Assume letrec  $H_1 \uplus H_2$  in  $E[\pi_i \ x] \xrightarrow{f^v}$  letrec  $H_1$  in  $E[\pi_i \ x] \xrightarrow{\pi_i}$  letrec  $H_1$  in  $E[x_i]$  Since **proj** applies, x must be bound to some pair  $\langle x_1, x_2 \rangle$  and furthermore, x must be in  $H_1$  else **fv** does not apply. Thus we may swap the steps:

letrec 
$$H_1 \uplus H_2$$
 in  $E[\pi_i \ x] \xrightarrow{\pi_i}$  letrec  $H_1 \uplus H_2$  in  $E[x_i] \xrightarrow{fv}$  letrec  $H_1$  in  $E[x_i]$ 

app: Assume letrec  $H_1 \uplus H_2$  in  $E[x \ y] \xrightarrow{fv}$  letrec  $H_1$  in  $E[x \ y] \xrightarrow{app}$  letrec  $H_1 \uplus \{z = h\}$  in E[e] Since app applies, x must be bound to  $\lambda z.e$ , y must be bound to h, and z cannot occur in the domain of  $H_1$ . Furthermore, both x and y must be bound in  $H_1$ . Performing the application first yields:

letrec 
$$H_1 \uplus H_2$$
 in  $E[x \ y] \xrightarrow{\text{app}} \text{letrec } H_1 \uplus H_2 \uplus \{z = h\} \text{ in } E[e]$ 

The resulting program is well formed only if z is not already bound in  $H_2$ . This can be ensured by  $\alpha$ -converting the program if necessary. Let x' be a variable that occurs in e. Then either x' is z or x' occurs free in  $\lambda z.e$ . In either case, x' is necessarily bound in  $H_1 \uplus \{z = h\}$  in order for the fv rule to apply, thus x' is not bound in  $H_2$ . Consequently, we can take the following step:

letrec 
$$H_1 \uplus H_2 \uplus \{z = h\}$$
 in  $E[e] \xrightarrow{\text{fv}} \text{letrec } H_1 \uplus \{z = h\}$  in  $E[e]$ 

**Lemma 4.2 (Diamond)** If  $P_1 \stackrel{\text{fv}}{\longmapsto} P_2$  and  $P_1 \stackrel{\mathbb{R}}{\longmapsto} P_2'$ , then there exists a  $P_3$  such that  $P_2 \stackrel{\mathbb{R}}{\longmapsto} P_3$  and  $P_2' \stackrel{\text{fv}}{\longmapsto} P_3$ .

**Proof:** Assume  $P_1 = \text{letrec } H_1 \uplus H_2 \text{ in } E[I]$ ,  $P_2 = \text{letrec } H_1 \text{ in } E[I]$ , and  $P_1 \overset{\text{fv}}{\longmapsto} P_2$ . We can easily show by case analysis on the elements of R that if  $P_1 \overset{\text{R}}{\longmapsto} P_2'$  where  $P_2' = \text{letrec } H_1 \uplus H_2 \uplus H_3 \text{ in } E[e]$ , for some  $H_3$  and e, then  $P_2' \overset{\text{fv}}{\longmapsto} P_3$  and  $P_2 \overset{\text{R}}{\longmapsto} P_3$  where  $P_3 = \text{letrec } H_1 \uplus H_3 \text{ in } E[e]$ .

With the Postponement and Diamond Lemmas in hand, it is straightforward to show that **fv** is a correct GC rule.

Theorem 4.3 (Correctness of fv) If  $P \xrightarrow{fv} P'$ , then P' is a collection of P.

**Proof:** Let  $P = \text{letrec } H_1 \uplus H_2 \text{ in } e \text{ and } P' = \text{letrec } H_1 \text{ in } e \text{ such that } P \xrightarrow{\text{fv}} P'$ . We must show P evaluates to an integer value iff P' evaluates to the same integer. Suppose  $P' \Downarrow_R \text{ letrec } H \text{ in } x$  and H(x) = i. By induction on the number of rewriting steps using the Postponement Lemma, we can show that  $P \Downarrow_R \text{ letrec } H \uplus H_2 \text{ in } x \text{ and clearly } (H \uplus H_2)(x) = H(x) = i$ . Now suppose  $P \Downarrow_R \text{ letrec } H \text{ in } x \text{ and } H(x) = i$ . By induction on the number of rewriting steps using the Diamond Lemma, we know that there exists an H' such that  $P' \Downarrow_R \text{ letrec } H' \text{ in } x \text{ and letrec } H \text{ in } x \xrightarrow{\text{fv}} \text{ letrec } H' \text{ in } x$ . Thus, x must be bound in H' and since fv only drops bindings, H'(x) = i.

This theorem shows that a single application of  $\mathbf{fv}$  results in a Kleene equivalent program. A real implementation interleaves garbage collection with evaluation. The following theorem shows that adding  $\mathbf{fv}$  to R preserves evaluation.

Theorem 4.4 For all programs P,  $(P,R) \simeq (P,R+\mathbf{fv})$ .

**Proof:** Clearly any evaluation under R can be simulated by  $R+\mathbf{fv}$  simply by not performing any  $\mathbf{fv}$  steps. Thus, if  $P \Downarrow_R A$  then  $P \Downarrow_{R+\mathbf{fv}} A$ . Now suppose  $P \Downarrow_{R+\mathbf{fv}}$  lettrec  $H_1$  in  $x_1$  and  $H_1(x_1) = i$ . Then there exists a finite rewriting sequence using  $R+\mathbf{fv}$  as follows:

$$P \xrightarrow{\mathbb{R} + \mathbf{fv}} P_1 \xrightarrow{\mathbb{R} + \mathbf{fv}} P_2 \xrightarrow{\mathbb{R} + \mathbf{fv}} \cdots \xrightarrow{\mathbb{R} + \mathbf{fv}} \mathsf{letrec} \ H_1 \ \mathsf{in} \ x_1$$

We can show by induction on the number of rewriting steps in this sequence, using the Postponement Lemma, that all  $\mathbf{fv}$  steps can be performed at the end of the evaluation sequence. This provides us with an alternative evaluation sequence where all the R steps are performed at the beginning:

$$P \overset{\mathbb{R}}{\longmapsto} P_1' \overset{\mathbb{R}}{\longmapsto} P_2' \overset{\mathbb{R}}{\longmapsto} \cdots \overset{\mathbb{R}}{\longmapsto} P_n' \overset{\mathbf{fv}}{\longmapsto} P_{n+1} \overset{\mathbf{fv}}{\longmapsto} P_{n+2} \overset{\mathbf{fv}}{\longmapsto} \cdots \overset{\mathbf{fv}}{\longmapsto} \text{ letrec } H_1 \text{ in } x_1$$

Since **fv** does not affect the expression part of a program and only removes bindings from the heap,  $P'_n = \text{letrec } H_1 \uplus H_2 \text{ in } x \text{ for some } H_2$ . Thus,  $P \Downarrow_R \text{ letrec } H_1 \uplus H_2 \text{ in } x$ . Since  $H_1(x) = i$ ,  $(H_1 \uplus H_2)(x) = i$ .

#### 4.2 The Free-Variable Tracing Algorithm

The free-variable GC rule is a *specification* of a garbage collection algorithm. It assumes some mechanism for partitioning the set of bindings into two disjoint pieces such that one set of bindings is unreachable from the second set of bindings and the body of the program. Real garbage collection algorithms need a deterministic mechanism for generating this partitioning. It is possible to formulate an abstract version of such a mechanism, the free-variable tracing algorithm, by lifting the ideas of mark-sweep and copying collectors to the level of program syntax.

We adopt the terminology of copying collection in the description of the free-variable tracing algorithm. We use two heaps and a set: a "from-heap"  $(H_f)$ , a "scan-set" (S), and a "to-heap"  $(H_t)$ . The from-heap is the set of bindings in the current program and the to-heap will become the set of bindings preserved by the algorithm. The scan-set records the set of variables reachable from the to-heap that have not yet been moved from the from-heap to the to-heap. The scan-set is often referred to as the "frontier."

The body of the algorithm proceeds as follows: A variable x is removed from S such that  $H_f$  has a binding for x. If no such locations are in S, the algorithm terminates. Otherwise, it scans the heap value h to which x is bound in the from-set  $H_f$ , looking for free variables. For each  $y \in FV(h)$ , it checks to see if y has already been forwarded to the to-set  $H_t$ . Only if y is not bound in  $H_t$  does it add the variable to the scan-set S. This ensures that a variable moves at most once from the from-heap to the scan-set.

Formulating the free-variable tracing algorithm as a rewriting system is easy. It requires only one rule that relates triples of from-sets, scan-sets, and to-sets:

$$\langle H_f \uplus \{x = h\}, S \uplus \{x\}, H_t \rangle \implies \langle H_f, S \cup (FV(h) \setminus (Dom(H_t) \uplus \{x\})), H_t \uplus \{x = h\} \rangle$$

Initially the free variables of the evaluation context and instruction expression, which correspond to the "roots" of a computation, are placed in S. Computing the free variables of the context represents the scanning of the "stack" of a conventional implementation while computing the free variables of the instruction expression corresponds to scanning the "registers." The initial tuple is re-written until we reach a state where no variable in the scan-set is bound in the from-heap. At this point, we have forwarded enough bindings to the to-heap. This leads to the following free-variable tracing algorithm rule:

$$\begin{array}{ccc} & \mathsf{letrec}\ H\ \mathsf{in}\ e \xrightarrow{\mathsf{fva}} \mathsf{letrec}\ H'\ \mathsf{in}\ e \\ & \mathsf{if} \\ & \langle H, FV(e), \emptyset \rangle \Longrightarrow^* \langle H'', S, H' \rangle \ \mathsf{and}\ \mathit{Dom}(H'') \cap S = \emptyset \end{array}$$

Clearly, the algorithm always terminates since the size of the from-heap strictly decreases with each step. Furthermore, this new rewriting rule is a subset of the rule **fv**, which implies the correctness of the algorithm.

**Theorem 4.5** If  $P \xrightarrow{\text{fva}} P'$ , then  $P \xrightarrow{\text{fv}} P'$ .

**Proof:** Let P = letrec H in e be a  $\lambda \text{gc}$  program. The first step is to prove the basic invariants of the garbage collection rewriting system. If  $\langle H, FV(e), \emptyset \rangle \Longrightarrow^* \langle H_f, S, H_t \rangle$ , then  $H_f \uplus H_t = H$  and  $FV(\text{letrec } H_t \text{ in } e) = S$ . Now let  $P' = \text{letrec } H_1 \text{ in } e$  and suppose  $P \stackrel{\text{fva}}{\longmapsto} P'$ . Then,

$$\langle H_1 \uplus H_2, FV(e), \emptyset \rangle \Longrightarrow^* \langle H_2, S, H_1 \rangle.$$

and  $Dom(H_2) \cap S = \emptyset$ . By the invariants,  $FV(\text{letrec } H_1 \text{ in } e) = S$ , so  $Dom(H_2) \cap FV(P') = \emptyset$ . Consequently,  $P \stackrel{\text{fv}}{\longmapsto} P'$ .

If we require that a collection algorithm produce a *closed* program, then **fva** is an optimal collection algorithm. That is, if P is a closed program and  $P \stackrel{\text{fva}}{\longmapsto} P'$ , then P' has the fewest bindings needed to keep the program closed without affecting evaluation. Assuming each step in the free-variable tracing algorithm takes time proportional to the size (in symbols) of the heap object forwarded to the to-heap, the time cost of the algorithm is proportional to the amount of data preserved, not the total amount of data in the original heap.

#### 4.3 Generational Garbage Collection

The free-variable tracing algorithm examines all of the reachable bindings in the heap to determine that a set of bindings may be removed. By carefully partitioning the heap into smaller heaps, a garbage collector can scan less than the whole heap and still free significant amounts of memory. A generational partition of a program's heap is a sequence of sub-heaps ordered in such a way that "older" generations never have pointers to "younger" generations.

**Definition 4.6 (Generational Partition)** A generational partition of a heap H is a sequence of heaps  $H_1, H_2, \ldots, H_n$  such that  $H = H_1 \uplus H_2 \uplus \cdots \uplus H_n$  and for all i such that  $1 \le i < n$ ,  $FV(H_i) \cap Dom(H_{i+1} \uplus H_{i+2} \uplus \cdots \uplus H_n) = \emptyset$ . The  $H_i$  are referred to as generations and  $H_i$  is said to be an older generation than  $H_j$  if i < j.

Given a generational partition of a program's heap, a free-variable based garbage collector can eliminate a set of bindings in younger generations without looking at any older generations.

Theorem 4.7 (Generational Collection) Let  $H_1, \ldots, H_i, \ldots, H_n$  be a generational partition of the heap of P = letrec H in e. Suppose  $H_i = (H_i^1 \uplus H_i^2)$ , and  $Dom(H_i^2) \cap FV(\text{letrec } H_i^1 \uplus H_{i+1} \uplus \cdots \uplus H_n \text{ in } e) = \emptyset$ . Then  $P \stackrel{\text{tv}}{\longmapsto} \text{letrec } (H \setminus H_i^2) \text{ in } e$ .

**Proof:** We must show that  $Dom(H_i^2) \cap FV(\text{letrec } (H \setminus H_i^2) \text{ in } e) = \emptyset$ . Since  $H_1, \dots, H_n$  is a generational partition of H, for all  $j, 1 \leq j < i$ ,  $FV(H_j) \cap Dom(H_{j+1} \uplus \dots \uplus H_n) = \emptyset$ . Hence,  $FV(H_1 \uplus \dots \uplus H_{i-1}) \cap Dom(H_i^2) = \emptyset$ . Now,

```
 FV(\operatorname{letrec}\ H\setminus H_i^2\ \operatorname{in}\ e)\cap Dom(H_i^2) \\ = (FV(H\setminus H_i^2)\cup FV(e))\cap Dom(H_i^2) \\ = (FV(H_1\uplus\cdots\uplus H_{i-1})\cup FV(H_i^1\uplus\cdots\uplus H_n)\cup FV(e))\cap Dom(H_i^2) \\ = (FV(H_1\uplus\cdots\uplus H_{i-1})\cap Dom(H_i^2))\cup ((FV(H_i^1\uplus\cdots\uplus H_n)\cup FV(e))\cap Dom(H_i^2)) \\ = \emptyset\cup ((FV(H_i^1\uplus\cdots\uplus H_n)\cup FV(e))\cap Dom(H_i^2)) \\ = FV(\operatorname{letrec}\ H_i^1\uplus\cdots\uplus H_n\ \operatorname{in}\ e)\cap Dom(H_i^2) \\ = \emptyset
```

Generational collection is important for three practical reasons: First, evaluation of closed, pure  $\lambda$ gc programs makes it easy to maintain generational partitions.

Theorem 4.8 (Generational Preservation for R) Let P = letrec H in e be a closed program. If  $H_1, \ldots, H_n$  is a generational partition of H and  $P \stackrel{\mathbb{R}}{\longmapsto} \text{letrec } H \uplus H' \text{ in } e, \text{ then } H_1, \ldots, H_n, H'$  is a generational partition of  $H \uplus H'$ .

**Proof:** The **proj** rule does not modify the heap, so it clearly preserves the generational partition. If the **alloc** rule is applied, then e = E[h] for some E and h and  $H' = \{x = h\}$ . Since P is closed, x cannot occur free in H. Thus for all  $i, 1 \le i \le n$ ,  $FV(H_i) \cap Dom(H_{i+1} \uplus H_{i+2} \uplus \cdots \uplus H_n \uplus H') = \emptyset$ . If the **app** rule is applied, then  $e = E[x \ y]$  for some E, x, and y such that x is bound to  $\lambda z.e'$  and y is bound to some h in H. Thus,  $H' = \{z = h\}$  and since P is closed, z cannot occur free in H. Thus for all  $i, 1 \le i \le n$ ,  $FV(H_i) \cap Dom(H_{i+1} \uplus H_{i+2} \uplus \cdots \uplus H_n \uplus H') = \emptyset$ .

The second reason generational collection is important is that, given a generational partition, we can directly use the free-variable tracing algorithm to generate a collection of a program. The following rule invokes the free-variable algorithm on the program letrec  $H_2$  in e where  $H_1$ ,  $H_2$  is a generational partition of the original program's heap. The resulting heap is glued onto  $H_1$  to produce a collection of the program.

$$(\mathbf{gen}) \qquad \begin{array}{c} \mathsf{letrec}\ H\ \mathsf{in}\ e \stackrel{\mathsf{gen}}{\longmapsto} \mathsf{letrec}\ H_1 \uplus H_2'\ \mathsf{in}\ e \\ \\ \mathsf{if}\ \mathsf{letrec}\ H_2\ \mathsf{in}\ e \stackrel{\mathsf{fva}}{\longmapsto} \mathsf{letrec}\ H_2'\ \mathsf{in}\ e \\ \\ \mathsf{and}\ H_1, H_2\ \mathsf{form}\ \mathsf{a}\ \mathsf{generational}\ \mathsf{partition}\ \mathsf{of}\ H \end{array}$$

The rule's soundness follows directly from the Generational Collection theorem, together with the soundness of the free-variable tracing algorithm.

Theorem 4.9 If 
$$P \xrightarrow{\text{gen}} P'$$
, then  $P \xrightarrow{\text{fv}} P'$ .

The third reason generational collection is important is that empirical evidence shows that "objects tend to die young" [35]. That is, recently allocated bindings are more likely to become garbage in a small number of evaluation steps. Thus, if we place recently allocated bindings in younger generations we can concentrate our collection efforts on these generations, ignoring older generations, and still eliminate most of the garbage.

## 4.4 Assignment and Generational Garbage Collection

The addition of certain language features or implementation techniques such as assignment or lazy evaluation tends to break the Generational Preservation theorem. For instance, consider the addition of an assignment operator (:=) to  $\lambda gc$ :

(expressions) 
$$e \in Exp ::= \cdots \mid e_1 := e_2$$
  
(contexts)  $E \in Ctxt ::= \cdots \mid E := e \mid x := E$   
(instructions)  $I \in Instr ::= \cdots \mid x := y$ 

$$(\mathbf{set}) \qquad \mathsf{letrec} \ H \uplus \{x = h, y = h'\} \ \mathsf{in} \ E[x := y] \ \stackrel{\mathsf{set}}{\longmapsto} \ \mathsf{letrec} \ H \uplus \{x = h', y = h'\} \ \mathsf{in} \ E[x]$$

resulting in a language  $\lambda$ gc-set. Evaluation of  $\lambda$ gc-set programs under  $R + \mathbf{set}$  does not satisfy the Generational Preservation theorem. To see this, consider the program:

letrec 
$$\{x_0 = 0, x_1 = 1, x_2 = \langle x_0, x_0 \rangle, x_3 = \langle x_2, x_2 \rangle \}$$
 in  $\pi_1$   $(\pi_1 (x_1 := x_3))$ 

where the heap can be generationally partitioned into  $H_1 = \{x_0 = 0, x_1 = 1\}$  and  $H_2 = \{x_2 = \langle x_0, x_0 \rangle, x_3 = \langle x_2, x_2 \rangle\}$ . Performing the assignment results in the program:

letrec 
$$\{x_0 = 0, x_1 = \langle x_2, x_2 \rangle, x_2 = \langle x_0, x_0 \rangle, x_3 = \langle x_2, x_2 \rangle \}$$
 in  $\pi_1$   $(\pi_1 \ x_1)$ 

but taking  $H_1 = \{x_0 = 0, x_1 = \langle x_2, x_2 \rangle\}$  and  $H_2 = \{x_2 = \langle x_0, x_0 \rangle, x_3 = \langle x_2, x_2 \rangle\}$  no longer yields a generational partition. If we attempted to do a generational collection on  $H_2$  using this partition, then both  $x_2$  and  $x_3$  would be collected since there is no reference to them in the expression part of the program. This would result in the program:

letrec 
$$\{x_0 = 0, x_1 = \langle x_2, x_2 \rangle\}$$
 in  $\pi_1$   $(\pi_1 \ x_1)$ 

which after one proj step would be stuck:

letrec 
$$\{x_0 = 0, x_1 = \langle x_2, x_2 \rangle\}$$
 in  $\pi_1 x_2$ 

because  $x_2$  is unbound. Yet, if no collection was performed, the computation would not get stuck and would return an answer of 0.

It is possible to maintain a generational partition for  $\lambda$ gc-set during evaluation. We simply keep track of all "older" bindings that are updated and move them from the older generation to a younger generation. The following theorem formalizes this idea for two generations. The generalization to multiple generations is straightforward.

Theorem 4.10 (Generational Preservation for R + set) If  $P = \text{letrec } H_1 \uplus H_2$  in e is a closed  $\lambda$  gc-set program and  $P \overset{\mathbb{R}+\text{set}}{\longrightarrow} {}^*P'$  where  $P' = \text{letrec } H_1 \uplus H_2' \uplus H_3$  in e' and  $Dom(H_2) = Dom(H_2')$ , then  $H_1, (H_2' \uplus H_3)$  is a generational partition of P'.

The theorem states that if we evaluate P for some number of steps resulting in P', where  $H_1 \uplus H_2$  is conceptually the older generation,  $H_2$  is the subset of the older generation that is modified via set to become  $H'_2$ , and  $H_3$  is the younger generation of newly allocated bindings, then  $H_1$  and  $(H'_2 \uplus H_3)$  form a generational partition. There are a variety of techniques for determining which values in the older generation must be moved to the younger generation. We refer the interested reader to Hosking *et al.* [19].

## 5 Garbage Collection via Type Recovery

The delimiters and other tokens of the abstract syntax mark or "tag" heap values with enough information that we can distinguish pairs from functions, pointers from integers, etc.. This allows us to navigate through the memory unambiguously, but placing tags on heap values and stripping them off to perform a computation can impose a heavy overhead on the running time and space requirements of programs [33]. An alternative to tagging is the use of types to determine the shape of an object. If types are determined at compile time and evaluation maintains enough information that the types of reachable objects can always be recovered, then there is no need to tag values.

A number of researchers have made attempts to explore this alternative [7, 8, 3, 15, 26, 34, 2], but none of them presented concise characterizations of the underlying techniques with correctness proofs. In this section, we present the basic idea behind type-recovery based garbage collection. In the following subsection, we give a simple overview of the techniques used to record and reconstruct the types of objects in collectors such as Britton's [8] and Tolmach's [34]. We then introduce  $\lambda$ gc-mono, an explicitly typed, monomorphic variant of  $\lambda$ gc. We show how to adapt the free-variable tracing algorithm to recover types of objects in the heap and to use these types in the

traversal of heap objects instead of abstract syntax. The proof of correctness for this garbage collection algorithm is given by extending the proof of soundness of the type system for  $\lambda$ gc-mono. The last subsection extends the work to an explicitly typed, polymorphic language.

## 5.1 An Overview of Type Recovery

The type of a data structure such as a tuple or integer determines the types of the components of the data structure (if any), but the type of a function does not contain the types of the free variables of the function. If the compiler associates the types of all free variables with the code of a function, it becomes possible to recover the types of all reachable objects in the heap, given the set of reachable objects from the current evaluation context and instruction. This provides a potential for space savings compared to tagging all values, because type information, like code, can be shared among different applications of the same function.

For example, in a low-level model that uses closures (an explicit environment mapping variables to values paired with an expression) instead of substitution, the expression part or "code" of the closure can be shared. The type information needed to find and preserve the free variables of the closure is simply a type environment that lists the types of each value in the closure's environment. This information can also be shared because the types of the free variables are the same for each value environment paired with a given expression.

As noted earlier, determining the free variables of the evaluation context and instruction corresponds to scanning a stack and registers in a conventional implementation. Thus, the compiler must provide a type map for each stack frame. Like the type environment associated with a closure, this information can be recorded at compile time and a pointer to the information can be pushed on the stack as a new stack frame is allocated. The type information is "collected" as stack frames are popped.

In our high-level framework, we use a global heap instead of closures, rely upon  $\alpha$ -conversion to generate unique names, and use evaluation contexts instead of an explicit stack to model evaluation. Consequently, the set of free variables changes as expressions are evaluated and  $\alpha$ -conversion is performed. Thus, the type maps associated with functions and evaluation contexts must also change. Our approach is to represent these type maps implicitly by initially tagging each expression with its type and by allowing the substitution that occurs as part of  $\alpha$ -conversion to replace a variable with a variable, but leave the type tag intact. The type map information associated with a closure or evaluation context can then be extracted by a simple function similar to the FV function of Figure 2. As data structures are allocated, the type information is stripped, leaving the heap essentially tag free and ensuring that we do not use the abstract syntax to navigate the heap. However, the type information remains intact within functions, corresponding to the explicit type maps needed for closures.

#### 5.2 $\lambda$ gc-mono

 $\lambda$ gc-mono is an explicitly typed, monomorphic variant of  $\lambda$ gc. The set of types  $(\tau)$  of  $\lambda$ gc-mono contains the conventional basic types and constructed types for typing a functional programming language like  $\lambda$ gc. The expressions of  $\lambda$ gc-mono are the same as for  $\lambda$ gc, except that each raw expression (u) is paired with some type information (see Figure 5). Heap values are not paired with their type, reflecting the fact that the memory is "almost tag free."

The evaluation contexts (E) consist of a raw context (U) and a type  $(\tau)$ . Raw contexts are filled with instructions (I) which are a subset of the raw expressions (u). The evaluation rules for  $\lambda$ gc-mono, named RM, are variants of the rules for  $\lambda$ gc that largely ignore the type information

```
Types:
                                                              \in Type ::= int | \tau_1 \times \tau_2 | \tau_1 \rightarrow \tau_2
Programs:
                           (expressions)
                                                              TExp
                                                                                  (u:\tau)
                                                                                  x \mid i \mid \langle e_1, e_2 \rangle \mid \pi_i \mid e \mid \lambda x : \tau . e \mid e_1 \mid e_2
                                                              UExp
                           (heap values)
                                                        \in
                                                              Hval
                                                                                  i \mid \langle x_1, x_2 \rangle \mid \lambda x : \tau.e
                           (heaps)
                                                  H
                                                        \in
                                                              Heap
                                                                                  \{x_1=h_1,\ldots,x_n=h_n\}
                                                   P
                                                        \in
                                                                                  letrec H in e
                           (programs)
                                                              Prog
                           (answers)
                                                   Α
                                                        \in
                                                              Ans
                                                                                  letrec H in (x:\tau)
```

#### **Evaluation Contexts and Instructions:**

Figure 5: The Syntax of  $\lambda$ gc-mono

on the sub-expressions of the program's body.

```
 \begin{array}{lll} \text{(alloc-int)} & \text{letrec } H \text{ in } E[i] & \stackrel{\text{alloc}}{\longmapsto} \text{letrec } H \uplus \{x=i\} \text{ in } E[x] \\ \text{(alloc-pair)} & \text{letrec } H \text{ in } E[\langle x_1 : \tau_1, x_2 : \tau_2 \rangle] & \stackrel{\text{alloc}}{\longmapsto} \text{letrec } H \uplus \{x=\langle x_1, x_2 \rangle\} \text{ in } E[x] \\ \text{(alloc-fn)} & \text{letrec } H \text{ in } E[\lambda y : \tau.e] & \stackrel{\text{alloc}}{\longmapsto} \text{letrec } H \uplus \{x=\lambda y : \tau.e\} \text{ in } E[x] \\ \text{(proj)} & \text{letrec } H \text{ in } E[\pi_i \ (x : \tau)] & \stackrel{\pi_i}{\longmapsto} \text{letrec } H \text{ in } E[x_i] \ (H(x)=\langle x_1, x_2 \rangle, i=1,2) \\ \text{(app)} & \text{letrec } H \text{ in } E[(x : \tau_1) \ (y : \tau_2)] & \stackrel{\text{app}}{\Longrightarrow} \text{letrec } H \uplus \{z=H(y)\} \text{ in } E[u] \ (H(x)=\lambda z : \tau'.(u : \tau'')) \\ \end{array}
```

Allocation of integers, tuples and functions strips the type tag off of the heap value before placing it in the heap. Allocation of tuples also removes the tags from the components of the tuple. The removal of type information corresponds to the passage of a value from code to data and does not necessarily reflect a runtime cost. Projection and application are essentially the same as for  $\lambda gc$ . Note that substitution of a result expression for an instruction occurs "in place," and hence the type of the instruction is ascribed to the expression.

The notion of a stuck state is adapted in accordance with the type structure of the language.

**Definition 5.1** ( $\lambda$ gc-mono Stuck Programs) A program is stuck if it is of one of the following forms:

```
• letrec H in E[\pi_i(x:\tau)] (x \notin Dom(H) \text{ or } H(x) \neq \langle x_1, x_2 \rangle);
```

```
• letrec H in E[(x:\tau_1) \ (y:\tau_2)] \quad (x \notin Dom(H) \ or \ H(x) \neq \lambda z:\tau.e \ or \ y \notin Dom(H)).
```

#### 5.3 Static Semantics

The static semantics of  $\lambda gc$ -mono consists of four judgements about program components. The first judgement,  $\Gamma \triangleright e$ , is a binary relation between a type assignment  $\Gamma$ , mapping a finite set of variables to types, and a typed expression e. Figure 6 (Expressions) contains the fairly conventional inference rules for expressions that generate the static semantics. Typing heaps and complete programs requires three additional judgements. The first is  $\Gamma \triangleright h : \tau$  which asserts that the heap value h has type  $\tau$  under the assumptions in  $\Gamma$ . The second is  $\Gamma \triangleright H : \Gamma'$ , which asserts that the variables

Expressions:

$$(\mathbf{var}) \ \frac{\Gamma \uplus \{x : \tau\} \trianglerighteq (x : \tau)}{\Gamma \uplus \{x : \tau\} \trianglerighteq (x : \tau)} \qquad (\mathbf{int}) \ \frac{\Gamma \trianglerighteq (i : \mathbf{int})}{\Gamma \trianglerighteq (i : \mathbf{int})}$$

$$(\mathbf{pair}) \ \frac{\Gamma \trianglerighteq (u_1 : \tau_1) \qquad \Gamma \trianglerighteq (u_2 : \tau_2)}{\Gamma \trianglerighteq (\langle (u_2 : \tau_1), (u_2 : \tau_2) \rangle : \tau_1 \times \tau_2)} \qquad (\mathbf{proj}) \ \frac{\Gamma \trianglerighteq (u : \tau_1 \times \tau_2)}{\Gamma \trianglerighteq (\pi_i \ (u : \tau_1 \times \tau_2) : \tau_i)} \qquad (i = 1, 2)$$

$$(\mathbf{fn}) \ \frac{\Gamma \uplus \{x : \tau_1\} \trianglerighteq (u : \tau_2)}{\Gamma \trianglerighteq (\lambda x : \tau_1. (u : \tau_2) : \tau_1 \to \tau_2)} \qquad (\mathbf{app}) \ \frac{\Gamma \trianglerighteq (u_1 : \tau_1 \to \tau_2) \qquad \Gamma \trianglerighteq (u_2 : \tau_1)}{\Gamma \trianglerighteq ((u_1 : \tau_1 \to \tau_2) \ (u_2 : \tau_1) : \tau_2)}$$

Heaps and Programs:

$$\begin{aligned} & (\mathbf{h}\textbf{-int}) \ \overline{\Gamma \rhd i : \mathsf{int}} \\ & (\mathbf{h}\textbf{-pair}) \ \frac{\Gamma \rhd (x_1 : \tau_1) \quad \Gamma \rhd (x_2 : \tau_2)}{\Gamma \rhd \langle x_1, x_2 \rangle : \tau_1 \times \tau_2} \qquad & (\mathbf{h}\textbf{-fn}) \ \frac{\Gamma \rhd (\lambda x : \tau_1.e : \tau_1 \to \tau_2)}{\Gamma \rhd \lambda x : \tau_1.e : \tau_1 \to \tau_2)} \\ & (\mathbf{heap}) \ \frac{\forall \, x \in Dom(\Gamma') \cdot \Gamma \uplus \Gamma' \rhd H(x) : \Gamma'(x)}{\Gamma \rhd H : \Gamma'} \qquad & (\mathbf{prog}) \ \frac{\emptyset \rhd H : \Gamma \quad \Gamma \rhd e}{\rhd \mathsf{letrec} \ H \ \mathsf{in} \ e} \end{aligned}$$

Figure 6: The Static Semantics of  $\lambda$ gc-mono

given type  $\tau$  in  $\Gamma'$  are bound to heap values in H of the appropriate type, under the assumptions in  $\Gamma$ . The judgement's definition via the **heap** rule, similar to Harper's store typing [18] and the typing for Wright & Felleisen's Reference ML [37], requires "guessing" the types of the values in the heap and verifying these guesses simultaneously due to potential cycles. The third judgement,  $\triangleright P$ , asserts the well-typing of a complete program. Figure 6 (Heaps and Programs) contains the necessary inference rules for all of these judgements.

**Definition 5.2** (Well Formed  $\lambda$ gc-mono Programs) A  $\lambda$ gc-mono program P is well formed if  $\triangleright P$  is derivable.

The calculus  $\lambda gc\text{-mono}$  is type sound in that evaluation of well formed programs cannot get stuck [37]. Our proof of soundness uses a subject-reduction based argument in the style of Wright and Felleisen [37].

**Theorem 5.3 (Type Soundness)** If  $\triangleright P$  then either P is an answer or else there exists some P' such that  $P \stackrel{\mathbb{R}^M}{\longmapsto} P'$  and  $\triangleright P'$ .

**Proof** (sketch): The first two lemmas (below) are needed to prove properties about the allocation rule. The Weakening lemma allows us to throw extra information into a type assignment and still type check an expression; the Allocation lemma allows us to replace an expression with a variable, as long as that location is placed in the type assignment with the proper type.

Lemma 5.4 (Weakening) If  $\Gamma \triangleright e$  and  $x \notin FV(e)$ , then  $\Gamma \uplus \{x:\tau\} \triangleright e$ .

**Lemma 5.5 (Allocation)** If  $\Gamma \triangleright (u : \tau)$  and  $\Gamma \triangleright E[u]$  and  $x \notin Dom(\Gamma)$ , then  $\Gamma \uplus \{x : \tau\} \triangleright E[(x : \tau)]$ .

The Unique Decomposition lemma allows us to deconstruct a program into a unique evaluation context and instruction expression.

**Lemma 5.6** (Unique Decomposition) If P = letrec H in e has no free variables, then either P is stuck, P is an answer, or else there exists unique E and I such that e = E[I].

Lemma 5.7 (Stuck Programs Untypeable) If P is a stuck program, then  $\not\triangleright P$ .

Finally, the key lemma is Type Preservation, which asserts that if a program is well typed and it takes a step, then the resulting program is also well typed.

Lemma 5.8 (Type Preservation) If  $\triangleright P$  and  $P \stackrel{\text{RM}}{\longmapsto} P'$ , then  $\triangleright P'$ .

**Proof** (sketch): Follows from the Unique Decomposition, Weakening, and Allocation lemmas, together with an examination of each of the rules of RM.

Given these lemmas, the proof of type soundness is straightforward. Suppose  $\triangleright P$  and P does not diverge. Then there exists some P' such that  $P \stackrel{\text{RM}}{\longmapsto} {}^*P'$  and P' is canonical with respect to RM. Further suppose that P' = letrec H in e for some H and e. If  $e = (x:\tau)$ , then by Type Preservation,  $\triangleright$  letrec H in  $x:\tau$  and the theorem is satisfied. Otherwise, by the Unique Decomposition lemma, P' must be stuck and thus  $\triangleright P'$  is not derivable, contradicting type preservation.  $\square$ 

#### 5.4 Using Types in Garbage Collection

Specifying a valid garbage collection rule that exploits types is now straightforward. The key property that the rule must preserve is *type preservation*. That is, if  $P \stackrel{gc}{\mapsto} P'$ , and P is well typed, then P' must be well typed. One way to achieve this goal is to make it into a side-condition of the GC rule:

(mono) letrec 
$$H_1 \uplus H_2$$
 in  $e \stackrel{\text{mono}}{\longmapsto}$  letrec  $H_1$  in  $e$  if  $\triangleright$  letrec  $H_1$  in  $e$ 

**Theorem 5.9** For all well formed  $\lambda gc\text{-mono programs }P,\ (P,RM)\simeq (P,RM+\mathbf{mono}).$ 

**Proof:** Since **mono** preserves types, the type soundness proof trivially extends to the dynamic semantics with **mono**. This implies that, since P is well formed, P cannot get stuck under either system. Since bindings are only dropped and not updated by **mono**, the results of evaluating under either system must be the same.

Like the free-variable rule of  $\lambda gc$ , **mono** needs to be refined before it can serve as the basis for an implementation. The refinement is similar to the one from **fv** to **fva** but uses the types embedded in programs to traverse heap values. The basis for the collection algorithm is a function that determines a minimal, with respect to set-inclusion, type assignment  $\Gamma$  for any expression e such that  $\Gamma \triangleright e$ :

```
\begin{array}{rcl} MTA(x:\tau) & = & \{x:\tau\} \\ MTA(i:\tau) & = & \emptyset \\ MTA(\langle e_1,e_2\rangle:\tau_1\times\tau_2) & = & MTA(e_1)\cup MTA(e_2) \\ MTA(\pi_i\ e:\tau) & = & MTA(e) \\ MTA(\lambda x:\tau_1.e:\tau_1\to\tau_2) & = & MTA(e)\setminus\{x:\tau_1\} \\ MTA(e_1\ e_2:\tau) & = & MTA(e_1)\cup MTA(e_2) \end{array}
```

**Lemma 5.10** If  $\Gamma \triangleright e$ , then  $MTA(e) \triangleright e$  and  $MTA(e) \subseteq \Gamma$ .

If P = letrec H in e and P is closed, then MTA(e) determines the types of the locations in the heap H that are immediately reachable from the expression e. A garbage collector can use this type information together with the following Tag function to traverse the reachable heap values based on their types instead of their abstract syntax.

$$\begin{array}{lcl} Tag[\operatorname{int}](i) & = & (i : \operatorname{int}) \\ Tag[\tau_1 \times \tau_2](\langle x_1, x_2 \rangle) & = & (\langle (x_1 : \tau_1), (x_2 : \tau_2) \rangle : \tau_1 \times \tau_2) \\ Tag[\tau_1 \to \tau_2](\lambda x : \tau_1 . e) & = & ((\lambda x : \tau_1 . e) : \tau_1 \to \tau_2) \end{array}$$

Tag is a curried function that takes a type and then takes a heap value of that type, and annotates that heap value with enough information to turn it back into an expression. It is important to note that Tag pattern-matches and operates according to the type argument given and not the abstract syntax of the heap value given. MTA can be used on the resulting expression to find the minimal type assignment for the heap value. This provides us with the free locations and their types for the heap value. The following lemma summarizes the relationship between Tag and MTA:

**Lemma 5.11** If  $\emptyset \triangleright H : \Gamma \uplus \{x : \tau\}$  then there exists an h such that

- 1.  $\{x=h\}\subseteq H$
- 2.  $MTA(Tag[\tau](h)) \triangleright h : \tau$
- 3.  $MTA(Tag[\tau](h)) \subseteq \Gamma \uplus \{x:\tau\}$

Equipped with MTA, we can now redefine the free-variable tracing algorithm so that it uses Tag to traverse heap values. The algorithm is specified in the same manner as  $\mathbf{fva}$ , i.e., as a rewriting system among tuples of the form  $\langle H_f, \Gamma_s, H_t, \Gamma_t \rangle$  where  $H_f$  is the from-heap,  $\Gamma_s$  is a type assignment corresponding to the scan-set,  $H_t$  is the to-heap, and  $\Gamma_t$  is a type assignment for the to-heap:

$$\langle H_f \uplus \{x = h\}, \Gamma_s \uplus \{x : \tau\}, H_t, \Gamma_t \rangle \stackrel{\text{mono}}{\Longrightarrow} \langle H_f, \Gamma_s', H_t \uplus \{x = h\}, \Gamma_t' \rangle$$
where
$$\Gamma_s' = (\Gamma_s \cup MTA(Tag[\tau](h)) \setminus \Gamma_t'$$

$$\Gamma_t' = \Gamma_t \uplus \{x : \tau\}$$

The algorithm is initialized by taking the program heap H as the initial from-heap and the minimal type assignment of the program expression as  $\Gamma_s$ . At each step in the algorithm,  $\Gamma_s$  describes the types of all locations that are immediately reachable from e or  $H_t$ , but have not yet been forwarded to  $H_t$ .  $\Gamma_t$  describes the types of all locations that have been forwarded to  $H_t$ .

When a variable x is found in  $\Gamma_s$  with type  $\tau$  and x is bound in  $H_f$  to the heap value h, the collector forwards the binding x = h to  $H_t$  and adds  $x:\tau$  to  $\Gamma_t$ . It then uses  $Tag[\tau]$  to traverse h, placing the necessary type information on the components so that MTA can determine the heap value's minimal type assignment. This step provides the locations (and their types) that are immediately reachable from h. Finally, the collector adds each of these locations to  $\Gamma_s$  unless they have already been forwarded to  $H_t$ .

Using this algorithm instead of an *a priori* partitioning of the heaps, **mono-a** becomes a high-level specification of a collector whose traversal of heap values uses types instead of tag information:

$$(\textbf{mono-a}) \begin{tabular}{l} \textbf{letrec $H$ in $e$} & \stackrel{\texttt{mono-a}}{\longmapsto} \textbf{letrec $H'$ in $e$} \\ \textbf{if} \\ \langle H, \mathit{MTA}(e), \emptyset, \emptyset \rangle & \stackrel{\texttt{mono}}{\Longrightarrow} {}^*\langle H'', \emptyset, H', \Gamma \rangle \\ \end{tabular}$$

The following definition gives the primary invariants of the algorithm:

<sup>&</sup>lt;sup>2</sup>The garbage collection rewriting system only maintains  $\Gamma_t$  to simplify the presentation and proof; an implementation will not have to construct  $\Gamma_t$ .

Definition 5.12 (mono-a Invariants)  $\langle H_f, \Gamma_s, H_t, \Gamma_t \rangle$  has the mono-a invariants with respect to a program letrec H in e and type assignment  $\Gamma_0$  iff:

- 1.  $H = H_f \uplus H_t$  (each binding is either in the from-heap or to-heap.)
- 2.  $\Gamma_s \uplus \Gamma_t \triangleright e$  (every binding needed for e is in the scan-set or to-heap.)
- 3.  $Dom(\Gamma_s) \subseteq Dom(H_f)$  (the scan-set corresponds to bindings in the from-heap.)
- 4.  $\Gamma_s \triangleright H_t : \Gamma_t$  (the scan-set holds all free variables in the to-heap.)
- 5.  $\Gamma_s \uplus \Gamma_t \subseteq \Gamma_0$  (the scan-set and to-heap agree with  $\Gamma_0$ .)

**Lemma 5.13** If  $\emptyset \triangleright H : \Gamma_0$ , T has the mono-a invariant properties with respect to P and  $\Gamma_0$ , and  $T \stackrel{\text{mono}}{\Longrightarrow} T'$ , then T' has the mono-a invariant properties with respect to P and  $\Gamma_0$ .

**Proof:** Assume  $P = \text{letrec } H \text{ in } e \text{ and } T = \langle H_f \uplus \{x = h\}, \Gamma_s \uplus \{x : \tau\}, H_t, \Gamma_t \rangle \text{ and this rewrites to } T' = \langle H_f, \Gamma_s', H_t \uplus \{x = h\}, \Gamma_t' \rangle \text{ where } \Gamma_t' = \Gamma_t \uplus \{x : \tau\} \text{ and } \Gamma_s' = (\Gamma_s \cup MTA(Tag[\tau](h))) \setminus \Gamma_t'. \text{ We must show:}$ 

- 1.  $H = H_f \uplus H_t \uplus \{x = h\}$ : Follows directly from the first invariant for T.
- 2.  $\Gamma_s' \uplus \Gamma_t' \rhd e$ : Follows directly from the second invariant for T together with weakening.
- 3.  $Dom(\Gamma'_s) \subseteq Dom(H_f)$ : This follows from the third invariant if  $(MTA(Tag[\tau](h)) \setminus \Gamma'_t) \subseteq Dom(H_f)$ . This in turn follows from  $\triangleright H : \Gamma_0, \{x:\tau\} \subseteq \Gamma_0$ , invariant 5 for T, and Lemma 5.11.
- 4.  $\Gamma_s' \triangleright H_t \uplus \{x = h\} : \Gamma_t'$ : This follows from the fourth invariant of T and the **heap** rule of the static semantics if we can show  $(\Gamma_s \cup (MTA(Tag[\tau](h)) \setminus \Gamma_t')) \uplus \Gamma_t' \triangleright h : \tau$ . This in turn follows from  $\{x : \tau\} \subseteq \Gamma_0$  and Lemma 5.11.
- 5.  $\Gamma_s' \uplus \Gamma_t' \subseteq \Gamma_0$ : This follows from the fifth invariant of T if  $MTA(Tag[\tau](h)) \subseteq \Gamma_0$ . This in turn follows from  $\{x:\tau\} \subseteq \Gamma_0$  and Lemma 5.11.

**Lemma 5.14** If  $\emptyset \triangleright H : \Gamma_0$ ,  $\Gamma_0 \triangleright e$ , and  $\langle H, MTA(e), \emptyset, \emptyset \rangle \stackrel{\text{mono}}{\Longrightarrow} {}^*T$ , then T has the mono-a invariant properties with respect to H and  $\Gamma_0$ .

**Proof:** By induction on the length of the rewriting sequence  $\langle H, MTA(e), \emptyset, \emptyset \rangle \stackrel{\text{mono}}{\Longrightarrow} T$ . The base case directly follows from Lemma 5.10. The inductive step directly follows from Lemma 5.13.

The correctness of the algorithmic type-based garbage collection rule can now easily be verified.

**Theorem 5.15** If P is a well formed program and  $P \stackrel{\text{mono-a}}{\longmapsto} P'$ , then  $P \stackrel{\text{mono}}{\longmapsto} P'$ .

**Proof:** We must verify that the **mono-a** garbage collection rule preserves typability in order for **mono** to apply. That is, we must show that if P is a well formed program, and  $P \stackrel{\text{mono-a}}{\longmapsto} P'$ , then  $\triangleright P'$ .

Let  $P = \text{letrec } H \text{ in } e \text{ and suppose } \triangleright P$ . Then, for some  $\Gamma_0$ ,  $\emptyset \triangleright H : \Gamma_0 \text{ and } \Gamma_0 \triangleright e$ . If

$$\langle H, MTA(e), \emptyset, \emptyset \rangle \stackrel{\text{mono}}{\Longrightarrow} {}^* \langle H'', \emptyset, H', \Gamma_t \rangle,$$

then by Lemma 5.14 taking  $\Gamma_s = \emptyset$ , we know that  $\emptyset \triangleright H' : \Gamma_t$  (invariant 4) and  $\Gamma_t \triangleright e$  (invariant 2). Thus, by the **prog** rule of the static semantics,  $\triangleright$  letrec H' in e.

Furthermore, the algorithm never gets stuck, so it always applies:

**Theorem 5.16** If P is a well formed program, then there exists a program P' such that  $P \stackrel{\text{mono-a}}{\longmapsto} P'$ .

**Proof:** If the algorithm takes a step, the size of the from-heap strictly decreases, so we know the collection either terminates or gets stuck. Here we show that the collection cannot get stuck, so it must terminate. Suppose P = letrec H in e and  $\langle H, MTA(e), \emptyset, \emptyset \rangle \stackrel{\text{mono}}{\Longrightarrow} {}^*\langle H_f, \Gamma_s \uplus \{x:\tau\}, H_t, \Gamma_t \rangle$ . By Lemma 5.14, we know that  $x \in Dom(H_f)$ , since the domain of the scan type-assignment is a subset of the from-heap. Consequently, there exists an h such that  $H_f = H'_f \uplus \{x = h\}$  and  $\langle H_f, \Gamma_s \uplus \{x:\tau\}, H_t, \Gamma_t \rangle \stackrel{\text{mono}}{\Longrightarrow} {}^*\langle H'_f, \Gamma'_s, H_t \uplus \{x = h\}, \Gamma'_t \rangle$  with appropriate  $\Gamma'_t$  and  $\Gamma'_s$ .

The mono-a rule does not perform type inference but rather type recovery. That is, the types of values in the heap are recovered without unification. This was accomplished by making the types of subexpressions explicit in the program. We have traded tags on data structures such as pairs, with tags on sub-expressions of functions. However, the tags can be shared among functions that use the same code, differing only in their environment. Furthermore, some space savings can be realized for parametric data structures such as lists. For example, if we have  $x:(\text{int}\times\text{int})$  list, as a subexpression, x may be bound to a large list, but each of the components of the list have the same type and consequently share the same "type tag" in an expression. A true comparison between tag-oriented and type-oriented garbage collection will have to be based on implementations.

#### 5.5 Explicit Polymorphism

Extending type-recovery based garbage collection to an explicitly polymorphic language where types are passed to polymorphic routines at run-time is straightforward. The language of the extended calculus,  $\lambda$ gc-poly, is an adaptation of the Girard-Reynolds polymorphic  $\lambda$ -calculus. The extended language of types contains type variables and quantified types:

(type variables) 
$$t \in Tvar$$
  
(types)  $\tau \in Type ::= \cdots \mid t \mid \forall t.\tau$ 

The marker  $\forall$  binds a type variable t in a type expression  $\tau$ . Types that only differ in their choice of bound type variables are considered to be equivalent.  $FTV(\tau)$  denotes the free type variables of type  $\tau$ .

The syntactic categories for  $\lambda$ gc-poly are the same as for  $\lambda$ gc-mono, except for the following additions:

```
\begin{array}{lll} u & \in & \text{untyped expr} & ::= & \cdots & \mid \Lambda t.e \mid e \mid \tau \\ h & \in & \text{heap values} & ::= & \cdots & \mid \Lambda t.e \end{array}
```

The  $\Lambda$ -expression, or type abstraction, is a construct for introducing and scoping polymorphism; type application  $(e[\tau])$  utilizes type abstractions at concrete types.

The set of evaluation contexts for  $\lambda gc$ -poly is the same as for  $\lambda gc$ -mono, modulo the extended syntax. In particular, there is no evaluation context within  $\Lambda$ -abstractions, because type abstracted expressions are allocated heap values. The one-step rewriting rules are the same as for  $\lambda gc$ -mono with the addition of two new rules, one for allocating  $\Lambda$ -abstractions and one for governing the application of  $\Lambda$ -abstractions:

$$\begin{array}{ll} \textbf{(alloc-t)} & \mathsf{letrec}\ H\ \mathsf{in}\ E\left[\Lambda t.e\right] \stackrel{\mathsf{alloc}}{\longmapsto} \mathsf{letrec}\ H \uplus \{x = \Lambda t.e\}\ \mathsf{in}\ E\left[x\right] \\ \textbf{(type-app)} & \mathsf{letrec}\ H\ \mathsf{in}\ E\left[(x:\tau)\left[\tau'\right]\right] \stackrel{\mathsf{t-app}}{\longmapsto} \mathsf{letrec}\ H\ \mathsf{in}\ E\left[\{\tau'/t\}u\right] \end{array} \qquad (H(x) = \Lambda t.(u:\tau'')) \end{array}$$

As with other heap values, type abstracted expressions are allocated on the heap. Type application  $((x:\tau)[\tau'])$  replaces all free occurrences of the bound type variable in a  $\Lambda$ -expression with a type. A compiler for this language might represent  $\Lambda$ -expressions in the same way that it represents

 $\lambda$ -abstractions and treat type application the same as term application. Since types are passed to type abstractions during evaluation, we say that this semantics is a *type-passing* interpretation. We refer to the extended set of rules as RP.

The static semantics for  $\lambda$ gc-poly defines the following judgements:

- 1.  $\Delta \triangleright \tau$  (type  $\tau$  is well formed.)
- 2.  $\Delta \triangleright \Gamma$  (type assignment  $\Gamma$  is well formed.)
- 3.  $\Delta$ ;  $\Gamma \triangleright e$  (expression e is well formed.)
- 4.  $\Delta$ ;  $\Gamma \triangleright h : \tau$  (heap value h is well formed.)
- 5.  $\Delta$ ;  $\Gamma \triangleright H : \Gamma'$  (heap H is well formed.)
- 6.  $\triangleright P$  (program P is well formed.)

where  $\Delta$  is a set of type variables and  $\Gamma$  is a type assignment. The first judgement asserts that the type  $\tau$  is well formed with respect to  $\Delta$  in that the free type variables of  $\tau$  occur in  $\Delta$ . The second judgement asserts that  $\Gamma$  is well formed with respect to  $\Delta$ , in that all of the types in the range of  $\Gamma$  are well formed with respect to  $\Delta$ . The third judgement asserts that an expression is well formed. The judgement is defined in the same fashion as the expression-level typing judgement for  $\lambda$ gc-mono, except that  $\Gamma$  is required to be well formed with respect to  $\Delta$ , and the resulting type must also be well formed with respect to  $\Delta$ . The following rules are added for  $\Lambda$ -abstractions and type applications respectively:

$$\frac{\Delta \triangleright \Gamma \qquad \Delta \uplus \{t\}; \Gamma \triangleright (u : \tau)}{\Delta; \Gamma \triangleright (\Lambda t. (u : \tau) : \forall t. \tau)} \quad (t \not\in \Delta) \qquad \frac{\Delta; \Gamma \triangleright (u : \forall t. \tau') \qquad \Delta \triangleright \tau}{\Delta; \Gamma \triangleright ((u : \forall t. \tau') \lceil \tau \rceil : \{\tau/t\}\tau')}$$

Similarly, a rule is added to typecheck a Λ-abstraction as a heap value:

$$\frac{\Delta \triangleright \Gamma \qquad \Delta \uplus \{t\}; \Gamma \triangleright (u : \tau)}{\Delta; \Gamma \triangleright \Lambda t. (u : \tau) : \forall t. \tau}$$

The last two judgements assert that heaps and programs are well formed and are defined by the following inference rules:

$$\frac{\forall \, x \in Dom(\Gamma') \, . \, \Delta; \Gamma \uplus \Gamma' \rhd H(x) : \Gamma'(x)}{\Delta; \Gamma \rhd H : \Gamma'} \qquad \qquad \underbrace{\emptyset; \emptyset \rhd H : \Gamma \qquad \Gamma \rhd e}_{\rhd \, \text{letrec $H$ in $e$}}$$

A program is well formed if the heap and expression each type check with an empty  $\Delta$ . Consequently, there can be no free type variables in a well formed program.

As for  $\lambda$ gc-mono, evaluation preserves typing, so if a program is well formed and then takes some number of steps, the resulting program will have no free type variables. Furthermore, the type system for  $\lambda$ gc-poly is sound:

**Theorem 5.17 (Type Soundness)** If  $\triangleright P$  then either P is an answer or else there exists a P' such that  $P \stackrel{\text{RP}}{\longmapsto} P'$  and  $\triangleright P'$ .

The type-based free-variable algorithm can be adapted for  $\lambda$ gc-poly as follows. First, MTA and Tag are extended to work with the new language constructs:

$$MTA(\Lambda t.e) = MTA(e)$$
  
 $MTA(e[\tau]) = MTA(e)$   
 $Tag[\forall t.\tau](\Lambda t.e) = (\Lambda t.e : \forall t.\tau)$ 

The garbage collection algorithm itself remains unchanged:

$$\langle H_f \uplus \{x=h\}, \Gamma_s \uplus \{x:\tau\}, H_t, \Gamma_t \rangle \stackrel{\text{poly}}{\Longrightarrow} \langle H_f, \Gamma_s', H_t \uplus \{x=h\}, \Gamma_t' \rangle$$
 where 
$$\Gamma_s' = (\Gamma_s \cup MTA(Tag[\tau](h)) \setminus \Gamma_t'$$
 
$$\Gamma_t' = \Gamma_t \uplus \{x:\tau\}$$

To prove the algorithm does not get stuck, we must show that we never find a variable in  $\Gamma_s$  that is assigned an unknown type (t). But this follows from the well-typedness of the original program. Consequently, it is straightforward to adapt the proofs of termination and correctness of **mono-a** to show that the rule **poly-a** is correct and always applies.

## 6 Collecting Reachable Garbage Using Type Inference

So far, we have only considered specifications and algorithms for collecting unreachable bindings. In this section, we show that by using type inference during the garbage collection process, some bindings that *are* reachable can still be safely collected. That is, type inference can be used to prove that an object is garbage even though it is reachable.

As a simple example, consider the following  $\lambda$ gc program:

letrec 
$$\{x_1 = 1, x_2 = 2, x_3 = \langle x_2, x_2 \rangle, x_4 = \langle x_1, x_3 \rangle \}$$
 in  $\pi_1 x_4$ 

Every binding in the heap is accessible from the program's expression  $(\pi_1 \ x_4)$ , so the free-variable based collection rules can collect nothing. But clearly the program will never dereference  $x_3$  nor  $x_2$ . The inference-based collection scheme described in this section will allow us to conclude that replacing the binding  $x_3 = \langle x_2, x_2 \rangle$  with  $x_3 = 0$  (or any other binding) will have no observable effect on evaluation. That is, the inference collection scheme shows that

letrec 
$$\{x_1 = 1, x_2 = 2, x_3 = 0, x_4 = \langle x_1, x_3 \rangle \}$$
 in  $\pi_1 \ x_4$ 

is Kleene equivalent to the original program. Now by applying the free-variable rule, we can conclude that the binding  $x_2 = 2$  can be safely collected.

In subsection 6.1, we introduce type inference for  $\lambda$ gc based on the presentation of Mitchell [25, Section 4.5]. In Section 6.2 we show that if type inference can assign a binding an unconstrained type and the rest of the program still has a valid typing, then the binding can effectively be collected. Our proof of this theorem is based on the method of *logical relations*, a standard proof technique from type theory. (See Mitchell [25, Section 3] for an overview of logical relations and various references.)

## 6.1 $\lambda gc$ and Type Inference

We start by considering the original  $\lambda gc$  language as an *implicitly* typed, monomorphic language, where the types of the language are the same as for  $\lambda gc$ -mono except for the addition of type variables:

(types) 
$$\tau \in \mathit{Type} ::= t \mid \mathsf{int} \mid \tau_1 \times \tau_2 \mid \tau_1 \to \tau_2$$

By implicitly typed, we mean that the terms of the language are not decorated with types as in  $\lambda$ gc-mono. We add type variables to the set of types so that each well-formed expression has a principal or most general type (explained below).

Type inference is the process of decorating  $\lambda$ gc programs with types so that the resulting program type checks under the  $\lambda$ gc-mono rules. (Refer to Figures 1 and 3 for the syntax and

dynamic semantics of  $\lambda gc$  and Figure 6 for the typing rules for  $\lambda gc$ -mono.) Alternatively, we may directly specify a set of typing rules for  $\lambda gc$  programs by taking the typing rules for  $\lambda gc$ -mono and erasing the type information from the terms, resulting in the inference system of Figure 7. These rules define judgements of the form  $\Gamma \vdash e : \tau$ , where  $\Gamma$  is a type assignment and e is a  $\lambda gc$  expression. We use " $\vdash$ " instead of " $\triangleright$ " to keep from confusing these typing judgements with the  $\lambda gc$ -mono and  $\lambda gc$ -poly judgements.

Expressions:

$$\begin{array}{c} \text{(var)} \ \overline{\Gamma \uplus \{x : \tau\} \vdash x : \tau} & \text{(int)} \ \overline{\Gamma \vdash i : \text{int}} \\ \\ \text{(tuple)} \ \overline{\Gamma \vdash e_1 : \tau_1} & \Gamma \vdash e_2 : \tau_2 \\ \overline{\Gamma \vdash \langle e_1, e_2 \rangle : \tau_1 \times \tau_2} & \text{(proj)} \ \overline{\Gamma \vdash e : \tau_1 \times \tau_2} \\ \\ \text{(fn)} \ \overline{\Gamma \uplus \{x : \tau_1\} \vdash e : \tau_2} \\ \overline{\Gamma \vdash \lambda x . e : \tau_1 \to \tau_2} & \text{(app)} \ \overline{\Gamma \vdash e_1 : \tau_1 \to \tau_2} & \overline{\Gamma \vdash e_2 : \tau_1} \\ \hline \Gamma \vdash e_1 \ e_2 : \overline{\tau_2} \\ \end{array}$$

Heaps and Programs:

$$(\mathbf{heap}) \ \frac{\forall \ x \in Dom(\Gamma') \ . \ \Gamma \uplus \Gamma' \vdash H(x) : \Gamma'(x)}{\Gamma \vdash H : \Gamma'} \qquad \qquad (\mathbf{prog}) \ \frac{\emptyset \vdash H : \Gamma \qquad \Gamma \vdash e : \tau}{\vdash \mathsf{letrec} \ H \ \mathsf{in} \ e : \tau}$$

Figure 7: Type Inference Rules for  $\lambda gc$ 

A given  $\lambda$ gc expression can have multiple typing derivations according to these rules and consequently multiple typings, but an expression's typings may be ordered so that there is a most general, or *principal* typing and every other typing is an *instance* of this principal typing and is thus derivable. The ordering is defined in terms of a substitution:

**Definition 6.1 (Substitution)** A substitution S is a function from type variables to type expressions.

We write  $S\tau$  to denote the type obtained by replacing each type variable t in  $\tau$  with S(t) and we write  $S\Gamma$  to denote the type assignment  $\{x:S\tau\mid x:\tau\in\Gamma\}$ .

**Definition 6.2 (Instance, Subsumes, Principal Typing)** A typing  $\Gamma' \vdash e : \tau'$  is an instance of the typing  $\Gamma \vdash e : \tau$  if there exists a substitution S such that  $S\Gamma \subseteq \Gamma'$  and  $S\tau = \tau'$ . In such a case, we say that  $\langle \Gamma, \tau \rangle$  subsumes  $\langle \Gamma', \tau' \rangle$ . A typing  $\Gamma \vdash e : \tau$  is principal if for all other typings  $\Gamma' \vdash e : \tau'$ ,  $\langle \Gamma, \tau \rangle$  subsumes  $\langle \Gamma', \tau' \rangle$ .

A consequence of the definition of a principal typing is that any instance of the principal typing is derivable:

**Lemma 6.3** If  $\langle \Gamma, \tau \rangle$  is a principal typing for e, then for all substitutions  $S, S\Gamma \vdash e : S\tau$ .

It is straightforward to see that principal typings for  $\lambda$ gc expressions are unique (up to  $\alpha$ -conversion of type variables). Consequently, we can search for the principal typing of a  $\lambda$ gc program in order to determine if there is any typing derivation of a program. The critical component of a type inference algorithm is an algorithm for unifying types and type assignments:

**Definition 6.4 (Unifier, Most General Unifier)** A unifier of a set of type equations E is a substitution S such that for all  $\tau_1 = \tau_2 \in E$ ,  $S\tau_1$  is syntactically equal to  $S\tau_2$ . A unifier S of E is more general than  $S_1$  if there exists a unifier  $S_2$  such that  $S_1 = S_2 \circ S$ .

**Lemma 6.5** (Unification) [31] Let E be a set of type equations. There is an algorithm Unify(E) that computes a most general unifier of E if one exists and fails otherwise.

If  $\Gamma_1$  and  $\Gamma_2$  are type assignments, then the unification algorithm can be used to compute a most general substitution S such that  $S\Gamma_1 \cup S\Gamma_2$  is well formed, if such a substitution exists. (Recall that a well formed type assignment maps term variables to at most one type.) We simply compute Unify(E) where  $E = \{\tau_1 = \tau_2 \mid x:\tau_1 \in \Gamma_1 \text{ and } x:\tau_2 \in \Gamma_2\}$ . We write  $Unify(\Gamma_1, \Gamma_2)$  for the substitution obtained in this manner.

```
(t fresh)
                            =\langle \{x:t\},t\rangle
            PT(x)
                                     \langle \emptyset, int \rangle
PT(\langle e_1, e_2 \rangle) = \text{let } \langle \Gamma_1, \tau_1 \rangle = PT(e_1)
                                      in let \langle \Gamma_2, \tau_2 \rangle = PT(e_2)
                                          in let S = Unify(\Gamma_1, \Gamma_2)
                                               \operatorname{in}\langle S\Gamma_1 \cup S\Gamma_2, S\tau_1 \times S\tau_2 \rangle
      PT(\pi_i \ e) = \text{let } \langle \Gamma, \tau \rangle = PT(e)
                                      in let S = Unify(\{x:\tau\}, \{x:t_1 \times t_2\})
                                                                                                                                            (t_1, t_2 \text{ fresh})
                                          in \langle S\Gamma, St_i \rangle
     PT(\lambda x.e) = |\text{let } \langle \Gamma, \tau_1 \times \tau_2 \rangle = PT(\langle x, e \rangle)
                                     in \langle \Gamma \setminus \{x:\tau_1\}, \tau_1 \to \tau_2 \rangle
   PT(e_1 \ e_2) = \text{let } \langle \Gamma_1, \tau_1 \rangle = PT(e_1) \text{ in let } \langle \Gamma_2, \tau_2 \rangle = PT(e_2)
                                     in let S = Unify(\Gamma_1 \uplus \{x:\tau_1\}, \Gamma_2 \uplus \{x:\tau_2 \to t\})
                                                                                                                                                (x, t \text{ fresh})
                                          in \langle S\Gamma_1 \cup S\Gamma_2, St \rangle
         PTA(e) = \operatorname{let} \langle \Gamma, \tau \rangle = PT(e) \operatorname{in} \Gamma
        PTT(e)
                                     let \langle \Gamma, \tau \rangle = PT(e) in \tau
```

Figure 8: An Algorithm for Computing the Principal Typing of an Expression

Given a  $\lambda gc$  expression e, the algorithm PT(e) in Figure 8 computes the most general type assignment and type for that expression. The algorithm is similar to the MTA algorithm of Section 5.4, but uses the Unify(-) routine to compute most general substitutions that allow two type assignments to be merged. The variable case returns a type assignment mapping the variable to a fresh type variable and returns that type variable as the type of the expression. The integer case returns an empty type assignment and the integer type. The tuple case computes the principal typings of the components of the tuple, and then computes a unifying substitution of their respective type assignments. This step is necessary because the same variable could occur within two different components, and the types of these two occurrences must unify. The substitution is applied to each component's principal type assignment and the union of the resulting assignments is returned as the principal assignment of the tuple. The type is calculated by applying the substitution to each of the components' principal types and then forming a tuple type. The type of an abstraction  $\lambda x.e$  is calculated by finding the principal typing of the expression  $\langle x,e\rangle$ . The type of the resulting expression contains the domain type of the abstraction paired with the range type. The pairing is needed since x may not occur in e and thus may not occur in the resulting type assignment  $\Gamma$ . The principal type assignment of the abstraction is simply  $\Gamma$  minus the occurrence of  $x:\tau_1$ .

The following lemma establishes the key properties of PT(e):

**Lemma 6.6** If  $PT(e) = \langle \Gamma, \tau \rangle$ , then  $\Gamma \vdash e : \tau$ . Furthermore, if  $\Gamma' \vdash e : \tau'$  for some  $\Gamma'$  and  $\tau'$ , then PT(e) exists and PT(e) subsumes  $\langle \Gamma', \tau' \rangle$ .

The algorithm for determining the principal typing of an expression can be easily extended to find a principal typing for a heap and a program. The former returns a principal pair of type assignments  $\langle \Gamma, \Gamma_H \rangle$  for H, such that  $\Gamma \vdash H : \Gamma_H$  and the latter returns the principal type assignment and type for the program:

$$PT(\{x_1 = h_1, \dots, x_n = h_n\}) = \operatorname{let} \langle \Gamma_1, \tau_1 \rangle = PT(h_1) \cdots \langle \Gamma_n, \tau_n \rangle = PT(e_n)$$

$$\operatorname{in} \operatorname{let} \Gamma = \{x_1 : \tau_1, \dots, x_n : \tau_n\}$$

$$\operatorname{in} \operatorname{let} S = \operatorname{Unify}(\Gamma, \Gamma_1, \dots, \Gamma_n)$$

$$\operatorname{in} \langle (S\Gamma_1 \cup \dots \cup S\Gamma_n) \setminus S\Gamma, S\Gamma \rangle$$

$$PT(\operatorname{letrec} H \operatorname{in} e) = \operatorname{let} \langle \Gamma, \Gamma_H \rangle = PT(H)$$

$$\operatorname{in} \operatorname{let} \langle \Gamma_e, \tau \rangle = PT(e)$$

$$\operatorname{in} \operatorname{let} S = \operatorname{Unify}(\Gamma \uplus \Gamma_H, \Gamma_e)$$

$$\operatorname{in} \langle S\Gamma \cup (S\Gamma_e \setminus S\Gamma_H), S\tau \rangle$$

The following lemma shows that PT(P) computes the principal type for a program and that any typing for the program is an instance of the principal typing.

**Lemma 6.7** If  $PT(P) = \langle \Gamma, \tau \rangle$ , then  $\Gamma \vdash P : \tau$ . Furthermore, if  $\Gamma' \vdash P : \tau'$ , then  $\langle \Gamma, \tau \rangle$  subsumes  $\langle \Gamma', \tau' \rangle$ .

The algorithm is complete in that if there exists some typing for the program, then the algorithm will not fail to find the principal typing of the program.

**Lemma 6.8** If  $\Gamma \vdash P : \tau$ , then PT(P) exists.

Finally, soundness for the static semantics follows by showing that principal typings of a program are preserved by evaluation.

**Theorem 6.9 (Type Soundness)** If  $\vdash P$  then either P is an answer or else there exists a P' such that  $P \stackrel{\mathbb{R}}{\longmapsto} P'$  and  $\vdash P'$ .

#### 6.2 The Inference GC Specification

We will show that if we can find a typing for a program that assigns a location bound in the heap a type variable, then that location may be bound to any value without affecting evaluation. Consequently, any pointers contained in the location's binding do not need to be scanned and traced during garbage collection. The intuition behind the theorem is that a location's type is unconstrained only if the location is not applied nor projected nor used in some manner that would constrain the type. Consequently, we can replace the binding in the heap with any binding we choose. In particular, if the location is bound to a large heap value, we can bind the location to an integer or dummy piece of code without affecting evaluation. This replacement allows us to collect any bindings that used to be reachable through this binding without knowing anything about the shape of the original heap value.

The proof of the theorem relies upon the semantic interpretation of types as logical relations. In our case, the relations are a type-indexed family of binary relations relating programs to programs,

answers to answers, and heaps to heaps. The relations are contrived so that, if two programs are related, then they are Kleene equivalent, so one program converges to an answer iff the other converges to a related answer and related answers at base type (int) yield equal values. Roughly speaking, the relations are logically extended to relate answers of higher type ( $\rightarrow$ ) if, whenever such answers are applied to appropriately related answers, the resulting computations yield related results.

The relations are defined in Figure 9 by induction on types. The definitions are parameterized by an arbitrary relational interpretation of type variables,  $\Theta$ . If t is a type variable, then  $\Theta(t)$  determines some fixed, but arbitrary relation between answer programs. This is consistent with the idea that well-typed programs have an implicit " $\forall$ " quantifier for the type variables in a program. The parameterization of the interpretation of type variables makes it straightforward to extend the definition of the relations to account for predicative polymorphism.

#### Computations:

```
\begin{split} \Theta \models P_1 \sim P_2 : \tau & \text{ iff } \quad \vdash P_1 : \tau \text{ and } \vdash P_2 : \tau \text{ and } \\ & P_1 \Downarrow_R A_1 \text{ implies } P_2 \Downarrow_R A_2 \text{ and } \Theta \models A_1 \approx A_2 : \tau \\ & \text{ and } \\ & P_2 \Downarrow_R A_2 \text{ implies } P_1 \Downarrow_R A_1 \text{ and } \Theta \models A_1 \approx A_2 : \tau \end{split}
```

#### Answers:

#### Heaps:

```
\Theta \models H_1 \approx H_2 : \Gamma \quad \text{iff} \quad \emptyset \vdash H_1 : \Gamma \text{ and } \emptyset \vdash H_2 : \Gamma \text{ and} \\ \forall x \in Dom(\Gamma).\Theta \models \text{letrec } H_1 \text{ in } x \approx \text{letrec } H_2 \text{ in } x : \Gamma(x)
```

Figure 9: Relational Interpretation of Types

Two well-typed programs  $P_1$  and  $P_2$  are related at a type  $\tau$ , written  $\Theta \models P_1 \sim P_2 : \tau$  iff, whenever one of the programs terminates with an answer, then the other program terminates with a related answer at type  $\tau$ .

Two answer programs  $A_1$  and  $A_2$  are related at type  $\tau$ , written  $\Theta \models A_1 \approx A_2 : \tau$ , as follows: If  $\tau$  is a type variable t, then the answers are related iff they are in the relation  $\Theta(t)$ . If  $\tau$  is an integer, then the answers are related iff, when we lookup the answer variables in their respective heaps, they are bound to the same integer value. If  $\tau$  is a pair, then the answers are related iff, whenever we project a component, the resulting programs are related. Finally, if  $\tau$  is an arrow type  $\tau_1 \to \tau_2$ , the answers are related iff, whenever we apply the answer variables to related arguments at type  $\tau_1$ , we get related programs at type  $\tau_2$ . Even though the relations between programs and answers are defined in terms of one another, the relations are well-founded because the size of the

type index always decreases when one relation refers to another.

The definition of the relations ensures that related programs remain related even if more bindings are added to the programs' heaps.

**Lemma 6.10**  $If \Theta \models \mathsf{letrec}\ H_1 \ \mathsf{in}\ e_1 \sim \mathsf{letrec}\ H_2 \ \mathsf{in}\ e_2 : \tau \ and \vdash \mathsf{letrec}\ H_1 \uplus H_1' \ \mathsf{in}\ e : \tau \ and \vdash \mathsf{letrec}\ H_2 \uplus H_2' \ \mathsf{in}\ e : \tau \ and$ 

Since evaluation only adds new bindings and leaves existing bindings intact, it is clear that evaluation preserves the relations. If the language permitted assignment, then this property would not necessarily hold.

For the statement of the following lemma, we need to extend the answer relation to heaps. Two heaps,  $H_1$  and  $H_2$ , are related at a context  $\Gamma$ , written  $\Theta \models H_1 \approx H_2 : \Gamma$ , if they can be  $\alpha$ -converted so that, for all variables x in  $\Gamma$ , the answers letrec  $H_1$  in x and letrec  $H_2$  in x are related at  $\Gamma(x)$ .

The following lemma establishes our primary result: an expression is related to itself in the context of any two related heaps.

**Lemma 6.11** For all  $\Theta$ :  $Tvar \rightarrow (Ans \leftrightarrow Ans)$ , if  $\Gamma \vdash e : \tau$  and  $\Theta \models H_1 \approx H_2 : \Gamma$ , then  $\Theta \models \text{letrec } H_1 \text{ in } e \sim \text{letrec } H_2 \text{ in } e : \tau$ .

**Proof:** By induction on the derivation of  $\Gamma \vdash e : \tau$ .

var:  $(\Gamma \uplus \{x:\tau\} \vdash x:\tau)$  By the definition of  $\Theta \models H_1 \approx H_2:\Gamma$ .

int: ( $\Gamma \vdash i : \text{int}$ ) letrec  $H_j$  in  $i \stackrel{\mathbb{R}}{\longmapsto}$  letrec  $H_j \uplus \{x_j = i\}$  in  $x_j$  for j = 1, 2 and these two answers are related by the definition of  $\Theta \models A_1 \approx A_2 : \text{int}$ .

tuple:  $(\Gamma \vdash \langle e_1, e_2 \rangle : \tau_1 \times \tau_2)$  Suppose letrec  $H_1$  in  $\langle e_1, e_2 \rangle$  does not diverge. Then by examination of the rewriting rules, it is clear that:

$$\begin{array}{ll} \operatorname{letrec} \ H_1 \ \operatorname{in} \ \langle e_1, e_2 \rangle \stackrel{\mathbb{R}}{\longmapsto}^{\star} & \operatorname{letrec} \ H_1 \uplus H_1' \uplus \{x_1 = h_1\} \ \operatorname{in} \ \langle x_1, e_2 \rangle \stackrel{\mathbb{R}}{\longmapsto}^{\star} \\ \operatorname{letrec} \ H_1 \uplus H_1' \uplus \{x_1 = h_1\} \uplus H_1'' \uplus \{x_2 = h_2\} \ \operatorname{in} \ \langle x_1, x_2 \rangle \stackrel{\mathbb{R}}{\longmapsto}^{\star} \\ \operatorname{letrec} \ H_1 \uplus H_1' \uplus \{x_1 = h_1\} \uplus H_1'' \uplus \{x_2 = h_2\} \uplus \{x = \langle x_1, x_2 \rangle\} \ \operatorname{in} \ x \end{array}$$

Thus, letrec  $H_1$  in  $e_1 \stackrel{\mathbb{R}}{\longmapsto} *$  letrec  $H_1 \uplus H_1' \uplus \{x_1 = h_1\}$  in  $x_1$ . By the induction hypothesis applied to the first hypothesis of the **pair** rule, we conclude that letrec  $H_2$  in  $e_1 \stackrel{\mathbb{R}}{\longmapsto} *$  letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\}$  in  $x_1'$  and these two answers are related. Consequently, letrec  $H_2$  in  $\langle e_1, e_2 \rangle \stackrel{\mathbb{R}}{\longmapsto} *$  letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\}$  in  $\langle x_1', e_2 \rangle$ .

From the rewriting sequence above, we can conclude that letrec  $H_1 \uplus H_1' \uplus \{x_1 = h_1\}$  in  $e_2 \overset{\mathbb{R}}{\mapsto} *$  letrec  $H_1 \uplus H_1' \uplus \{x_1 = h_1\} \uplus H_1'' \uplus \{x_2 = h_2\}$  in  $x_2$ . By hypothesis,  $\Theta \models H_1 \approx H_2 : \Gamma$  and thus by Lemma  $6.10 \Theta \models H_1 \uplus H_1' \uplus \{x_1 = h_1\} \approx H_2 \uplus H_2' \uplus \{x_2 = h_2\} : \Gamma$ , so the induction hypothesis applies to the second hypothesis of the **pair** rule and we can conclude that letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\}$  in  $e_2 \overset{\mathbb{R}}{\mapsto} *$  letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\} \uplus H_2'' \uplus \{x_2' = h_2'\}$  in  $x_2'$  and these two answers are related. Consequently, letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\}$  in  $\langle x_1', e_2 \rangle \overset{\mathbb{R}}{\mapsto} *$  letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\} \uplus H_2'' \uplus \{x_2' = h_2'\}$  in  $\langle x_1', x_2' \rangle$ . Since related answers remain related under heap extensions,

$$\Theta \models \mathsf{letrec} \ H_1 \uplus H_1' \uplus H_1'' \uplus \{x_1 = h_1, x_2 = h_2, x = \langle x_1, x_2, \} \rangle \ \mathsf{in} \ x_i \approx \\ \mathsf{letrec} \ H_2 \uplus H_2'' \uplus \{x_1' = h_1', x_2' = h_2', x' = \langle x_1', x_2', \} \rangle \ \mathsf{in} \ x_i' : \tau_i \qquad (i = 1, 2)$$

Consequently,

$$\Theta \models \mathsf{letrec}\ H_1 \uplus H_1' \uplus H_1'' \uplus \{x_1 = h_1, x_2 = h_2, x = \langle x_1, x_2, \} \rangle \ \mathsf{in}\ \pi_i\ x \sim \\ \mathsf{letrec}\ H_2 \uplus H_2'' \uplus H_2'' \uplus \{x_1' = h_1', x_2' = h_2', x' = \langle x_1', x_2', \} \rangle \ \mathsf{in}\ \pi_i\ x' : \tau_i \qquad (i = 1, 2)$$

Thus,  $\Theta \models \text{letrec } H_1 \text{ in } \langle e_1, e_2 \rangle \sim \text{letrec } H_2 \text{ in } \langle e_1, e_2 \rangle : \tau_1 \times \tau_2$ .

proj:  $(\Gamma \vdash \pi_i \ e : \tau_i)$  Suppose letrec  $H_1$  in  $\pi_i \ e$  converges. Then,

letrec 
$$H_1$$
 in  $\pi_i$   $e \stackrel{\mathbb{R}}{\longmapsto} {}^*$ letrec  $H_1 \uplus H_1'$  in  $\pi_i$   $x \stackrel{\mathbb{R}}{\longmapsto}$  letrec  $H_1 \uplus H_1'$  in  $y$ 

for some  $H_1'$ , x, and y. Thus, letrec  $H_1$  in  $e \mapsto^{\mathbb{R}} *$ letrec  $H_1 \uplus H_1'$  in x. By the induction hypothesis, letrec  $H_2$  in  $e \mapsto^{\mathbb{R}} *$ letrec  $H_2 \uplus H_2'$  in x' and  $\Theta \models$  letrec  $H_1 \uplus H_1'$  in  $x \approx$  letrec  $H_2 \uplus H_2'$  in  $x' : \tau_1 \times \tau_2$ . By the definition of this relation, we know that

$$\Theta \models \mathsf{letrec}\ H_1 \uplus H_1' \ \mathsf{in}\ \pi_i\ x \sim \mathsf{letrec}\ H_2 \uplus H_2' \ \mathsf{in}\ \pi_i\ x' : \tau_1 \times \tau_2$$

for i = 1, 2. Thus,  $\Theta \models \text{letrec } H_1 \text{ in } \pi_i \ e \sim \text{letrec } H_2 \text{ in } \pi_i \ e : \tau_1 \times \tau_2$ .

fn:  $(\Gamma \vdash \lambda x.e : \tau_1 \to \tau_2)$  letrec  $H_i$  in  $\lambda x.e \overset{\mathbb{R}}{\longmapsto}$  letrec  $H_i \uplus \{y_i = \lambda x.e\}$  in  $y_i$  for i = 1, 2. We must show that these two answers are related. Suppose:

$$\Theta \models \mathsf{letrec}\ H_1 \uplus \{y_1 = \lambda x.e\} \uplus H_1' \ \mathsf{in}\ z_1 \approx \mathsf{letrec}\ H_1 \uplus \{y_2 = \lambda x.e\} \uplus H_2' \ \mathsf{in}\ z_2 : \tau_1$$

for some  $H'_1$ ,  $H'_2$ ,  $z_1$ , and  $z_2$ . We need to show that:

$$\Theta \models \mathsf{letrec}\ H_1 \uplus \{y_1 = \lambda x.e\} \uplus H_1' \ \mathsf{in}\ y_1\ z_1 \sim \mathsf{letrec}\ H_1 \uplus \{y_2 = \lambda x.e\} \uplus H_2' \ \mathsf{in}\ y_2\ z_2 : \tau_2 \sqcup \tau_2$$

This follows if:

$$\Theta \models \mathsf{letrec}\ H_1 \uplus \{y_1 = \lambda x.e, x = h_1\} \uplus H_1' \ \mathsf{in}\ e \sim \mathsf{letrec}\ H_1 \uplus \{y_2 = \lambda x.e, x = h_2\} \uplus H_2' \ \mathsf{in}\ e : \tau_2 \uplus H_2' \ \mathsf{in}\ \mathsf{in}\$$

where  $h_i = (H_i \uplus \{y_i = \lambda x.e\} \uplus H_i')(z_i)$  for i = 1, 2. By the induction hypothesis, it suffices to show that:

$$\Theta \models H_1 \uplus \{y_1 = \lambda x.e, x = h_1\} \uplus H_1' \approx H_1 \uplus \{y_2 = \lambda x.e, x = h_2\} \uplus H_2' : \Gamma \uplus \{x : \tau_1\}$$

By the lemma's hypothesis, we know that  $\Theta \models H_1 \approx H_2 : \Gamma$  and thus by Lemma 6.10

$$\Theta \models H_1 \uplus \{y_1 = \lambda x.e, x = h_1\} \uplus H_1' \approx H_2 \uplus \{y_2 = \lambda x.e, x = h_2\} \uplus H_2' : \Gamma.$$

Since

 $\Theta \models \text{letrec } H_1 \uplus \{y_1 = \lambda x.e, x = h_1\} \uplus H_1' \text{ in } z_1 \approx \text{letrec } H_1 \uplus \{y_2 = \lambda x.e, x = h_2\} \uplus H_2' \text{ in } z_2 : \tau_1$  and  $h_1$  and  $h_2$  are bound to  $z_1$  and  $z_2$  in these respective heaps, it follows that:

 $\Theta \models \text{letrec } H_1 \uplus \{y_1 = \lambda x.e, x = h_1\} \uplus H_1' \text{ in } x \approx \text{letrec } H_1 \uplus \{y_2 = \lambda x.e, x = h_2\} \uplus H_2' \text{ in } x : \tau_1$ Consequently,

$$\Theta \models H_1 \uplus \{y_1 = \lambda x.e, x = h_1\} \uplus H_1' \approx H_1 \uplus \{y_2 = \lambda x.e, x = h_2\} \uplus H_2' : \Gamma \uplus \{x : \tau_1\}$$

holds and, working backwards we know that

$$\Theta \models \mathsf{letrec}\ H_1 \uplus \{y_1 = \lambda x.e\} \uplus H_1' \ \mathsf{in}\ y_1\ z_1 \sim \mathsf{letrec}\ H_1 \uplus \{y_2 = \lambda x.e\} \uplus H_2' \ \mathsf{in}\ y_2\ z_2 : \tau_2 \sqcup \tau_2$$

and thus  $\Theta \models \text{letrec } H_1 \uplus \{y_1 = \lambda x.e\} \text{ in } y_1 \approx \text{letrec } H_2 \uplus \{y_2 = \lambda x.e\} \text{ in } y_2 : \tau_1 \to \tau_2.$ 

app:  $(\Gamma \vdash e_1 \ e_2 : \tau_2)$  Suppose letrec  $H_1$  in  $e_1 \ e_2$  does not diverge. Then by examination of the rewriting rules, it is clear that:

letrec 
$$H_1$$
 in  $e_1$   $e_2 \overset{\mathbb{R}}{\longmapsto} {}^*$  letrec  $H_1 \uplus H_1' \uplus \{x_1 = h_1\}$  in  $x_1$   $e_2 \overset{\mathbb{R}}{\longmapsto} {}^*$  letrec  $H_1 \uplus H_1' \uplus \{x_1 = h_1\} \uplus H_1'' \uplus \{x_2 = h_2\}$  in  $x_1$   $x_2$ 

Thus, letrec  $H_1$  in  $e_1 \stackrel{\mathbb{R}}{\longmapsto} *$ letrec  $H_1 \uplus H_1' \uplus \{x_1 = h_1\}$  in  $x_1$ . By the induction hypothesis applied to the first hypothesis of the app rule, we conclude that letrec  $H_2$  in  $e_1 \stackrel{\mathbb{R}}{\longmapsto} *$ letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\}$  in  $x_1'$ 

and these two answers are related at  $\tau_1 \to \tau_2$ . Consequently, letrec  $H_2$  in  $e_1$   $e_2 \mapsto^{\mathbb{R}} *$  letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\}$  in  $x_1'$   $e_2$ .

From the rewriting sequence above, we can conclude that letrec  $H_1 \uplus H_1' \uplus \{x_1 = h_1\}$  in  $e_2 \overset{\mathbb{R}}{\longmapsto} *$ letrec  $H_1 \uplus H_1' \uplus \{x_1 = h_1\} \uplus H_1'' \uplus \{x_2 = h_2\}$  in  $x_2$ . By hypothesis,  $\Theta \models H_1 \approx H_2 : \Gamma$  and thus  $\Theta \models H_1 \uplus H_1' \uplus \{x_1 = h_1\} \approx H_2 \uplus H_2' \uplus \{x_2 = h_2\} : \Gamma$ , so the induction hypothesis applies to the second hypothesis of the app rule and we can conclude that letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\}$  in  $e_2 \overset{\mathbb{R}}{\longmapsto} *$ letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\} \uplus H_2'' \uplus \{x_2' = h_2'\}$  in  $x_2'$  and these two answers are related at  $\tau_1$ . Consequently, letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\}$  in  $x_1' e_2 \overset{\mathbb{R}}{\longmapsto} *$ letrec  $H_2 \uplus H_2' \uplus \{x_1' = h_1'\} \uplus H_2'' \uplus \{x_2' = h_2'\}$  in  $x_1' x_2'$ . Since

$$\Theta \models \mathsf{letrec}\ H_1 \uplus H_1' \uplus \{x_1 = h_1\} \uplus H_1'' \uplus \{x_2 = h_2\} \ \mathsf{in}\ x_2 \approx \\ \mathsf{letrec}\ H_2 \uplus H_2' \uplus \{x_1' = h_1'\} \uplus H_2'' \uplus \{x_2' = h_2'\} \ \mathsf{in}\ x_2' : \tau_1,$$

and

$$\Theta \models \mathsf{letrec}\ H_1 \uplus H_1' \uplus \{x_1 = h_1\} \uplus H_1'' \uplus \{x_2 = h_2\} \ \mathsf{in}\ x_1 \approx \\ \mathsf{letrec}\ H_2 \uplus H_2' \uplus \{x_1' = h_1'\} \uplus H_2'' \uplus \{x_2' = h_2'\} \ \mathsf{in}\ x_1' : \tau_1 \to \tau_2$$

by the definition of  $\approx$  at arrow-types, it follows that

$$\Theta \models \mathsf{letrec}\ H_1 \uplus H_1' \uplus \{x_1 = h_1\} \uplus H_1'' \uplus \{x_2 = h_2\} \ \mathsf{in}\ x_1\ x_2 \sim \\ \mathsf{letrec}\ H_2 \uplus H_2' \uplus \{x_1' = h_1'\} \uplus H_2'' \uplus \{x_2' = h_2'\} \ \mathsf{in}\ x_1'\ x_2' : \tau_2$$

Thus,  $\Theta \models \mathsf{letrec}\ H_1 \ \mathsf{in}\ e_1\ e_2 \sim \mathsf{letrec}\ H_2 \ \mathsf{in}\ e_1\ e_2 : \tau_2$ .

Our goal is to show that if we are given a context  $\Gamma$ , expression e, and heap H such that  $\Gamma \vdash e : \tau$  and  $\emptyset \vdash H : \Gamma$ , then letrec H' in e is Kleene equivalent to letrec H in e where H' is defined as follows:

$$H' = \{x = H(x) \mid \Gamma(x) \notin Tvar\} \uplus \{x = 0 \mid \Gamma(x) \in Tvar\}$$

This follows from Lemma 6.11 if we can show that, taking  $\Theta_0(t)$  to be the everywhere-defined relation on answer programs,  $\Theta_0 \models H \approx H' : \Gamma$ . This in turn follows if we can show that  $\Theta_0 \models H \approx H : \Gamma$  (H is related to itself), since  $\Theta_0(t)$  relates every program and H'(x) differs from H(x) only when  $\Gamma(x) = t$ .

Unfortunately, we cannot directly show that a well formed heap is related to itself! The problem is that if we attempt to argue by induction on the derivation of  $\emptyset \vdash H : \Gamma$ , the uses of the heap rule require that we assume what we are trying to prove. The same problem is encountered when using logical relations to reason about conventional calculi with recursion or iteration operators.

If we forbid cycles in the heap, then we can transform the derivation of the heap's well formedness into a derivation that only uses a let-style rule instead of the recursive letrec-style heap rule:

$$\begin{aligned} \text{(let-heap)} \ \frac{\Gamma_1 \uplus \Gamma_2 \vdash h : \tau \qquad \Gamma_1 \vdash H : \Gamma_2}{\Gamma_1 \vdash H \uplus \{x = h\} : \Gamma_2 \uplus \{x : \tau\}} \end{aligned}$$

Consequently, if a heap is cycle free, we may show by induction on the derivation using the **let-heap** rule that the heap is related to itself.

**Lemma 6.12** If  $\Gamma_1 \vdash H : \Gamma_2$ ,  $\Theta \models H_1 \approx H_2 : \Gamma_1$  and H is cycle free, then  $\Theta \models H_1 \uplus H \approx H_2 \uplus H : \Gamma_1 \uplus \Gamma_2$ .

**Proof:** Since H is cycle free, there exists some derivation of  $\Gamma_1 \vdash H$ :  $\Gamma$  using only the **letheap** rule. The proof proceeds by induction on the height of this derivation. If H is empty, then the lemma holds trivially by the assumptions about  $H_1$  and  $H_2$ . Suppose  $\Gamma_1 \vdash H$ :  $\Gamma_2$ 

and  $\Gamma_1 \uplus \Gamma_2 \vdash h : \tau$ . By induction,  $\Theta \models H_1 \uplus H \approx H_2 \uplus H : \Gamma_1 \uplus \Gamma_2$ . By Lemma 6.11,  $\Theta \models \mathsf{letrec}\ H_1 \uplus H \ \mathsf{in}\ h \sim \mathsf{letrec}\ H_2 \uplus H \ \mathsf{in}\ h : \tau$ . Thus,

$$\Theta \models \mathsf{letrec}\ H_1 \uplus H \uplus \{x = h\} \ \mathsf{in}\ x \approx \mathsf{letrec}\ H_2 \uplus H \uplus \{x = h\} \ \mathsf{in}\ x : \tau.$$

Consequently, 
$$\Theta \models H_1 \uplus H \uplus \{x = h\} \approx H_2 \uplus H \uplus \{x = h\} : \Gamma_1 \uplus \Gamma_2 \uplus \{x : \tau\}.$$

Finally, we can state and prove the following Inference GC specification:

Theorem 6.13 (Inference GC) Let  $\Gamma_1 = \{x_1:t_1,\ldots,x_n:t_n\}$  and  $H_1 = \{x_1 = h_1,\ldots,x_n = h_n\}$  and  $H'_1 = \{x_1 = 0,\ldots,x_n = 0\}$ . If

- 1.  $\Gamma_1 \uplus \Gamma_2 \vdash e : \tau \ (\tau \not\in Tvar)$ , and
- 2.  $\Gamma_1 \vdash H_2 : \Gamma_2$ , and
- 3.  $\exists S.\emptyset \vdash H_1 : S\Gamma_1$ , and
- 4. H<sub>2</sub> is cycle free,

then letrec  $H_1 \uplus H_2$  in  $e \simeq$  letrec  $H_1' \uplus H_2$  in e.

**Proof:** Taking  $\Theta_0(t)$  to be the everywhere defined relation,  $\Theta_0 \models H_1 \approx H_1' : \Gamma_1$  holds trivially. By Lemma 6.12, since  $H_2$  is cycle free and  $\Gamma_1 \vdash H_2 : \Gamma_2$ , we know that  $\Theta_0 \models H_1 \uplus H_2 \approx H_1' \uplus H_2 : \Gamma_1 \uplus \Gamma_2$ . Since  $\Gamma_1 \uplus \Gamma_2 \vdash e : \tau$ , we know by Lemma 6.11 that  $\Theta_0 \models \text{letrec } H_1 \uplus H_2 \text{ in } e \sim \text{letrec } H_1' \uplus H_2 \text{ in } e \downarrow \text{letrec } H_1 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_1 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_1 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_1 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_1 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_1 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_1 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_2 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_2 \uplus H_2 \text{ in } e \downarrow \text{letrec } H_2 \uplus H_$ 

It should be possible to extend our argument to cyclic heaps if we can show that every cyclic heap is appropriately approximated by some finite "unwinding" of the heap. One approach is to use a fixed-point semantics for  $\lambda gc$  where types are interpreted as CPOs formed by suitably lifting answer programs. The meaning of an answer program letrec H in x would essentially be the least fixed-point of the chain of meanings of the approximations letrec  $H^0$  in x,letrec  $H^1$  in x,letrec  $H^2$  in  $x,\ldots$  where  $H^i$  denotes the  $i^{th}$  unwinding of the heap H. See Gunter [17, Chapter 4.4] for an example of the treatment of unwinding and approximation in the context of a conventional  $\lambda$ -calculus.

#### 7 Related Work

The literature on garbage collection in sequential programming languages per se contains few papers that attempt to provide a compact characterization of algorithms or correctness proofs. Demers et al. [10] give a model of memory parameterized by an abstract notion of a "points-to" relation. As a result, they can characterize reachability-based algorithms including mark-sweep, copying, generational, "conservative," and other sophisticated forms of garbage collection. However, their model is intentionally divorced from the programming language and cannot take advantage of any semantic properties of evaluation, such as type preservation. Consequently, their framework cannot easily model the type-based collectors of Sections 5 and 6. Nettles [27] provides a concrete specification of a copying garbage collection algorithm using the Larch specification language. Our specification of the free-variable tracing algorithm is essentially a high-level, one-line description of his specification. Like Demers et al., he uses a traditional definition of garbage based on unreachability in the memory graph and purposefully distances the specification from any language evaluation details.

Hudak gives a denotational model that tracks reference counts for a first-order language [20]. He presents an abstraction of the model and gives an algorithm for computing approximations of reference counts statically. Chirimar, Gunter, and Riecke give a framework for proving invariants regarding memory management for a language with a *linear* type system [9]. Their low-level semantics specifies explicit memory management based on reference counting. The goal of the work was to determine whether, using the linear type system, in-place update would be a sound optimization. Both Hudak and Chirimar *et al.* assume a fairly weak approximation of garbage (reference counts).

Tolmach [34] built a type-recovery collector for a variant of SML that passes type information to polymorphic routines during execution, effectively implementing our  $\lambda$ gc-poly language and the corresponding collector of Subsection 5.5. Aditya and Caro gave a type-recovery algorithm for an implementation of Id that uses a technique that appears to be equivalent to type passing [1] and Aditya, Flood, and Hicks extended this work to garbage collection for Id [2].

There have been a variety of papers regarding inference-based collection for monomorphic [7, 36, 8] and polymorphic languages [3, 15, 16, 12]. Appel [3] argued informally that "tag-free" collection is possible for polymorphic languages such as SML by a combination of recording information statically and performing what amounts to type inference during the collection process, though the connections between inference and collection were not made clear. Baker [6] recognized that Milner-style type inference can be used to prove that reachable objects can be safely collected, but did not give a formal account of this result. Goldberg and Gloger [16] recognized that it is not possible to reconstruct the concrete types of all reachable values in an implementation of an ML-style language that does not pass types to polymorphic routines. They gave an informal argument based on traversal of stack frames to show that such values are semantically garbage. Fradet [12] gave another argument based on Reynolds' abstraction/parametricity theorem [29]. Fradet's formulation is closer to ours than Goldberg and Gloger's, since he represented the evaluation "stack" as a source-language term. However, none of these papers gives a complete formulation of the underlying dynamic and static semantics of the language and thus the proofs of correctness are necessarily ad hoc.

Finally, Purushothaman and Seaman [28, 32] and Launchbury [21] have proposed "natural" semantics for call-by-need (lazy) languages where the semantic objects include an explicit heap. This allows sharing and memoization of computations to be expressed in the semantics. More recently, Ariola et al. [5] (see also [4, 22]) have presented a purely syntactic theory of the call-by-need  $\lambda$ -calculus that is largely compatible with our work.

## 8 Summary

Our paper provides a unifying semantic framework for a variety of garbage collection ideas including standard copying and mark-sweep collection, generational collection, tag-free collection, and inference-based collection. By making allocation and the heap explicit, we are able to reason about memory management using traditional  $\lambda$ -calculus techniques. In particular, we are able to make strong connections between garbage collection and type theory. By abstracting away such low-level details as evaluation stacks, registers, and addresses, we are able to formulate complicated collection algorithms in a compact manner and yet give a formal proof of correctness. The framework is reasonably robust, as demonstrated by the breadth of topics covered in this paper, and we expect that the framework can be extended to deal with other work on garbage collection.

## 9 Acknowledgements

Thanks to Edo Biagoni, Andrzej Filinski, Mark Leone, Scott Nettles, Chris Okasaki, and Jeannette Wing for proofreading earlier drafts of this document and for providing valuable insight. This manuscript is submitted for publication with the understanding that the U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes, notwithstanding any copyright notation thereon.

#### References

- [1] S. Aditya and A. Caro. Compiler-directed type reconstruction for polymorphic languages. In *Proceedings* of the Conference on Functional Programming Languages and Computer Architecture, pages 74-82, Copenhagen, Denmark, June 1993.
- [2] S. Aditya, C. Flood, and J. Hicks. Garbage collection for strongly-typed languages using run-time type reconstruction. In *Proceedings of the 1994 ACM Conference on Lisp and Functional Programming*, pages 12-23, June 1994.
- [3] A. W. Appel. Runtime tags aren't necessary. Journal of Lisp and Symbolic Computation, 2:153-162, 1989.
- [4] Z. M. Ariola and M. Felleisen. The call-by-need lambda calculus. Technical Report CIS-TR-94-23, University of Oregon, Oct. 1994.
- [5] Z. M. Ariola, M. Felleisen, J. Maraist, M. Odersky, and P. Wadler. A call-by-need lambda calculus. In Conference Record of the 22nd Annual ACM Symposium on Principles of Programming Languages, San Francisco, CA, Jan. 1995.
- [6] H. Baker. Unify and conquer (garbage, updating, aliasing ...) in functional languages. In *Proceedings* of the 1990 ACM Conference on Lisp and Functional Programming, pages 218-226, 1990.
- [7] P. Branquart and J. Lewi. A scheme for storage allocation and garbage collection for Algol-68. In Algol-68 Implementation. North-Holland Publishing Company, Amsterdam, 1970.
- [8] D. E. Britton. Heap storage management for the programming language Pascal. Master's thesis, University of Arizona, 1975.
- [9] J. Chirimar, C. A. Gunter, and J. G. Riecke. Proving memory management invariants for a language based on linear logic. In *Proceedings of the 1992 ACM Conference on Lisp and Functional Programming*, pages 139-150, June 1992.
- [10] A. Demers, M. Weiser, B. Hayes, H. Boehm, D. Bobrow, and S. Shenker. Combining generational and conservative garbage collection: Framework and implementations. In Conference Record of the 17th Annual ACM Symposium on Principles of Programming Languages, pages 261-269, Jan. 1990.
- [11] M. Felleisen and R. Hieb. The revised report on the syntactic theories of sequential control and state. Technical Report 89-100, Rice University, June 1989. Also appears in: *Theoretical Computer Science*, 102, 1992.
- [12] P. Fradet. Collecting more garbage. In Proceedings of the 1994 ACM Conference on Lisp and Functional Programming, pages 24-33, June 1994.
- [13] J.-Y. Girard. Une extension de l'interpretation de Gödel à l'analyse, et son application a l'elimination des coupures dans l'analyse et la théorie des types. In Proceedings of the Second Scandinavian Logic Symposium, edited by J.E. Fenstad. North-Holland, Amsterdam, pages 63-92, 1971.
- [14] J.-Y. Girard. Interprétation Fonctionnelle et Elimination des Coupures dans l'Arithmétique d'Ordre Supérieur. PhD thesis, Université Paris VII, 1972.

- [15] B. Goldberg. Tag-free garbage collection for strongly typed programming languages. In *Proceedings* of the ACM SIGPLAN '91 Conference on Programming Language Design and Implementation, pages 165-176, June 1991.
- [16] B. Goldberg and M. Gloger. Polymorphic type reconstruction for garbage collection without tags. In *Proceedings of the 1992 ACM Conference on Lisp and Functional Programming*, pages 53-65, June 1992.
- [17] C. A. Gunter. Semantics of Programming Languages. MIT Press, 1992.
- [18] R. Harper. A simplified account of polymorphic references. Technical Report CMU-CS-93-169, School of Computer Science, Carnegie Mellon University, June 1993.
- [19] A. L. Hosking, J. E. B. Moss, and D. Stefanović. A comparative performance evaluation of write barrier implementations. In *Proceedings OOPSLA '92, ACM SIGPLAN Notices*, pages 92-109, Oct. 1992. Published as Proceedings OOPSLA '92, ACM SIGPLAN Notices, volume 27, number 10.
- [20] P. Hudak. A semantic model of reference counting and its abstraction. In *Proceedings of the 1986 ACM Conference on Lisp and Functional Programming*, pages 351-363, 1986.
- [21] J. Launchbury. A natural semantics for lazy evaluation. In Conference Record of the 20th Annual ACM Symposium on Principles of Programming Languages, Charleston, SC, Jan. 1993.
- [22] J. Maraist, M. Odersky, and P. Wadler. The call-by-need lambda calculus (unabridged). Technical Report 28/94, Universität Karlsruhe, Fakultät für Informatik, Oct. 1994.
- [23] I. Mason and C. Talcott. Reasoning about programs with effects. In *Proceedings of Programming Language Implementation and Logic Programming*, LNCS 582. Springer-Verlag, 1990.
- [24] I. Mason and C. Talcott. Equivalences in functional languages with effects. *Journal of Functional Programming*, 2(1), 1991.
- [25] J. C. Mitchell. Type systems for programming languages. Technical Report STAN-CS-89-1277, Department of Computer Science, Stanford University, 1989.
- [26] R. Morrison, A. Dearle, R. C. H. Connor, and A. L. Brown. An ad hoc approach to the implementation of polymorphism. ACM Trans. Prog. Lang. Syst., 13(3):342-371, July 1991.
- [27] S. Nettles. A Larch specification of copying garbage collection. Technical Report CMU-CS-92-219, School of Computer Science, Carnegie Mellon University, Dec. 1992.
- [28] Purushothaman and J. Seaman. An adequate operational semantics of sharing in lazy evaluation. In *Proceedings of the 4th European Symposium on Programming*, LNCS 582. Springer-Verlag, 1992.
- [29] J. Reynolds. Types, abstraction, and parametric polymorphism. In *Proceedings of Information Processing 83*, pages 513-523, 1983.
- [30] J. C. Reynolds. Towards a theory of type structure. In *Proceedings, Colloque sur la Programmation*. Lecture Notes in Computer Science, volume 19, pages 408-425. Springer-Verlag, Berlin, 1974.
- [31] J. A. Robinson. A machine-oriented logic based on the resolution principle. *Journal of the Association for Computing Machinery*, 12, 1965.
- [32] J. M. Seaman. An Operational Semantics of Lazy Evaluation for Analysis. PhD thesis, Pennsylvania State University, 1993.
- [33] P. Steenkiste and J. Hennessey. Tags and type checking in LISP: Hardware and software approaches. In Proceedings of the Second International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-II), pages 50-59, Oct. 1987.
- [34] A. Tolmach. Tag-free garbage collection using explicit type parameters. In *Proceedings of the 1994 ACM Conference on Lisp and Functional Programming*, pages 1-11, June 1994.

- [35] D. Ungar. Generational scavenging: A non-disruptive high performance storage management reclamation algorithm. In ACM SIGPLAN Software Engineering Symposium on Practical Software Development Environments, pages 15–167, Pittsburgh, Pennsylvania, Apr. 1984.
- [36] P. Wodon. Methods of garbage collection for Algol-68. In Algol-68 Implementation. North-Holland Publishing Company, Amsterdam, 1970.
- [37] A. Wright and M. Felleisen. A syntactic approach to type soundness. Technical Report TR91-160, Department of Computer Science, Rice University, Apr. 1991.