# Liveness-Based Garbage Collection

Rahul Asati[1], Amitabha Sanyal[1], Amey Karkare[2], and Alan Mycroft[3]

[1] IIT Bombay, Mumbai 400076, India
{rahulasati,as}@cse.iitb.ac.in,
[2] IIT Kanpur, Kanpur 208016, India
karkare@cse.iitk.ac.in,
[3] Computer Laboratory, University of Cambridge, CB3 0FD, UK
alan.mycroft@cl.cam.ac.uk

**Abstract.** Current garbage collectors leave much heap-allocated data uncollected because they preserve data *reachable* from a root set. However, only *live* data—a subset of reachable data—need be preserved.

Using a first-order functional language we formulate a context-sensitive liveness analysis for structured data and prove it correct. We then use a 0-CFA-like conservative approximation to annotate each allocation and function-call program point with a finite-state automaton—which the garbage collector inspects to curtail reachability during marking. As a result, fewer objects are marked (albeit with a more expensive marker) and then preserved (e.g. by a copy phase).

Experiments confirm the expected performance benefits—increase in garbage reclaimed and a consequent decrease in the number of collections, a decrease in the memory size required to run programs, and reduced overall garbage collection time for a majority of programs.

## 1 Introduction

Most modern programming languages support dynamic allocation of heap data. Static analysis of heap data is much harder than analysis of static and stack data. Garbage collectors, for example, conservatively approximate the liveness of heap objects by their reachability from a set of memory locations called the *root set*. Consequently, many objects that are reachable but not live remain uncollected, causing a larger-than-necessary memory demand. This is confirmed by empirical studies on Haskell [1], Scheme [2] and Java [3] programs.

Here we consider a first-order pure functional language and propose a liveness analysis which annotates various program points with a description of variables and fields whose object references may be dereferenced in the future. The garbage collector then only marks objects pointed to by live references and leaves other, merely reachable, objects to be reclaimed. (Although not strictly necessary, a collector would normally *nullify* dead variables and fields rather than leaving dangling references.) Since there are fewer live objects than reachable objects, more memory is reclaimed. Additionally, since the collector traverses a smaller portion of the heap, the time spent for each collection is also smaller. The work is presented in the context of a stop-the-world non-incremental garbage collector (mark-and-sweep, compacting or copying) for which we

also show a monotonicity result: that our technique can never cause more garbage collections to occur in spite of changing the rather unpredictable execution points at which collections occur. We anticipate that our technique is applicable to more modern collectors (generational, concurrent, parallel), but leave such extensions to future work.

We first define a *fully context-sensitive* (in the sense that its results are unaffected by function inlining) liveness analysis and prove it correct. However, fully context-sensitive methods often do not scale, and this analysis would also require us to determine, at run-time, the internal liveness of a function body at each call. Hence, similarly to the 0-CFA approach, we determine a context-independent *summary* of liveness for each function which safely approximates the context-dependence of all possible calls [4–6]. (Note that an intraprocedural context-insensitive method which assumes no information about function callers would be too imprecise for our needs.) In essence our approach sets up interprocedural data-flow equations for the liveness summaries of functions and shows how these can be solved symbolically as context-free grammars (CFGs). We can then determine a CFG for each program point; these are then safely approximated with finite-state automata which are encoded as tables for each program point. For garbage collection purposes only automata corresponding to *GC points* need to be stored. GC points are program points associated with a call to a user function or to **cons**—see Section 4.

We previously proposed an intraprocedural method for heap liveness analysis for a Java-like language [7] which statically inserted statements nullifying dead references to improve garbage collection; by contrast nullification here occurs dynamically (which can work better with aliasing) when the garbage collector acts on liveness annotations to avoid traversing dead references. A workshop paper [8] outlined the basic 0-CFA-style-summary interprocedural approach to functional-program liveness analysis. The current paper adds the context-sensitive analysis and better formalisation along with experimental results.

**Motivating Example**  Figure 1(a) shows an example program. The label $\pi$ of an expression $e$ denotes a program point. During execution of the program, it represents the instant of time just before the evaluation of $e$. We view the heap as a graph. Nodes in the heap, also called (**cons**) *cells* contain **car** and **cdr** *fields* containing values. Edges in the graph are *references* and emanate from *variables* or fields. Variable and field values may also be atomic values (**nil**, integers etc.) While it is convenient to box these in diagrams, our presented analysis treats them as non-heap values.

Figure 1(b) shows the heap at $\pi$. The edges shown by thick continuous arrows are those which are made live by the program. In addition, assuming that the value of any reachable part of the program result may be explored or printed, the edges marked by thick dashed arrows are also live. A cell is marked and preserved during garbage collection, only if it is reachable from the root set through a path of live edges. All other cells can be reclaimed. We model the liveness properties of the heap as automata and pass these automata to the garbage collector. Thus if a garbage collection happens at $\pi$ with the heap shown in Figure 1(b), only the cells $w$ and (**cdr** $w$), along with (**car** (**cdr** $w$)) and all cells reachable from it, will be marked and preserved.

```
(define (append l1 l2)
  (if (null? l1) l2
      (cons (car l1)
            (append (cdr l1) l2))))

(let z ←(cons (cons 4 (cons 5 nil))
              (cons 6 nil)) in
 (let y ← (cons 3 nil) in
    (let w ← (append y z) in
      π:(car (cdr w)))))
```

(a) Example program.

(b) Memory graph at π. Thick edges denote live links. Traversal stops during garbage collection at edges marked ✕.
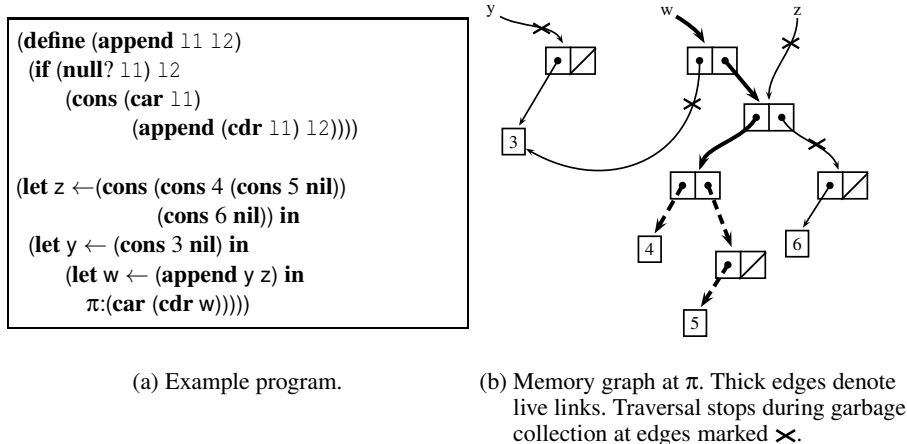
**Fig. 1.** Example Program and its Memory Graph.

**Organisation of the paper** Section 2 gives the syntax and semantics of the language used to illustrate our analysis along with basic concepts and notations. Liveness analysis is described in Section 3 followed by a sketch of a correctness proof relative to a non-standard semantics. Section 4 shows how to encode liveness as finite-state automata. Section 5 reports experimental results and Section 6 proves that a liveness-based collector can never do more garbage collections than a reachability-based collector.

## 2 The target language—syntax and semantics

We let $x$, $y$, $z$ range over variables, $f$ over user-functions and $p$ over primitive functions (**cons**, $+$ etc.). The syntax of our language is shown in Figure 2; it has eager semantics and restricts programs to be in Administrative Normal Form (ANF) [9] where all actual parameters to functions are variables. This restriction does not affect expressibility (and indeed we feel free to ignore it in examples when inessential), but simplifies the analysis formulation. Additionally, as in the three-address instruction form familiar from compiler texts, it forces each temporary to be named and function calls to be serialised (necessary to get an unambiguous definition of liveness). We further require that each variable in a program is distinct, so that no scope shadowing occurs—this simplifies proofs of soundness. In this formulation expressions: either perform a test (**if**), make a computation step (**let**) or return (**return**). The **return** keyword is logically redundant, but we find it clarifies the semantics and analysis.

The body of the program is the expression denoted by $e_{\mathbf{main}}$; for analysis purposes it is convenient to regard $e_{\mathbf{main}}$ as part of a function definition (**define** (main) $e_{\mathbf{main}}$) as in C. We write $π\!:\!e$ to associate the label $π$ (not part of the language syntax) with the program point just before expression $e$.

In spite of the ANF restrictions it is still possible to inline non-recursive functions (a fact we use to prove the safety of liveness analysis). A user-function call (**let** $x \leftarrow (f\, y_1 \ldots y_n)$ **in** $e$) to a function defined (after renaming its formals and locals to be disjoint from existing variables) by (**define** $(f\, z_1 \ldots z_n)\, e_f$) is replaced by a sequence of **let**s of the form $z_i \leftarrow (\mathbf{id}\, y_i)$ followed by the body $e_f$ but with its (**return** $w$) expressions replaced by (**let** $x \leftarrow (\mathbf{id}\, w)$ **in** $e$). (We prefer to use **id** as a form of no-op

3

$$
\begin{array}{lll}
p \in Prog & ::= d_1 \dots d_n\, e_{\textbf{main}} & \text{— } program \\
d \in Fdef & ::= (\textbf{define}\ (f\ x_1\ \dots\ x_n)\ e) & \text{— } function\ definition \\
e \in Expr & ::= \begin{cases} (\textbf{if}\ x\ e_1\ e_2) & \text{— } conditional \\ (\textbf{let}\ x \leftarrow s\ \textbf{in}\ e) & \text{— } let\ binding \\ (\textbf{return}\ x) & \text{— } return\ from\ function \end{cases} \\[4ex]
s \in Stmt & ::= \begin{cases} k & \text{— } constant\ (numeric\ or\ \textbf{nil}) \\ (\textbf{cons}\ x_1\ x_2) & \text{— } constructor \\ (\textbf{car}\ x) \qquad (\textbf{cdr}\ x) & \text{— } selectors \\ (\textbf{null}?\ x) \qquad (+\ x_1\ x_2) & \text{— } tester\ and\ generic\ arithmetic \\ (\textbf{id}\ x) & \text{— } identity\ function\ (for\ inlining) \\ (f\ x_1\ \dots\ x_n) & \text{— } function\ application \end{cases}
\end{array}
$$

**Fig. 2.** The syntax of our language

function rather than introducing the form ($\textbf{let}\ x \leftarrow w\ \textbf{in}\ e$) where the *Stmt* part of **let** is a simple variable.)

**Semantics** We now give an operational semantics for our language. Later, a refinement of the operational semantics, which we call minefield semantics, will serve to prove liveness analysis correct. We give a small-step semantics because, unlike big-step semantics, correctness for non-terminating programs does not need special treatment. We start with the domains used by the semantics:

$$
\begin{array}{lll}
v : Val & = \mathbb{N} + \{\textbf{nil}\} + Loc & \text{– Values} \\
\rho : Env & = Var \rightarrow Val & \text{– Environment} \\
H : Heap & = Loc \rightarrow (Val \times Val) & \text{– Heap}
\end{array}
$$

Here *Loc* is a countable set of locations which hold **cons** cells. A value is either a number, the empty list **nil**, or a location $\ell$. Our liveness analysis does not track numeric values, and thus is neutral as to whether these are boxed or represented as immediates. An environment is a finite mapping from variables to values, and a heap a finite mapping from locations to pairs of values. Finally, S is a stack (using $\bullet$ for push and $[]$ for empty stack) of frames of unfinished function calls. A frame is a triple $(e, x, \rho)$ representing the call site ($\textbf{let}\ x \leftarrow (f\ y_1\ \dots\ y_n)\ \textbf{in}\ e$) being evaluated in environment $\rho$. Frames can also be viewed as continuations, in this view the (ORD-RETURN) rule in the small-step operational semantics (Figure 3) invokes them.

The semantics of statements $s$ are given by the judgement form $\rho, H, s \rightsquigarrow H', v$ and those for expressions $e$ by the form $\rho, S, H, e \rightarrow \rho', S', H', e'$. The start state is $(\{\}, [], \{\}, e_{\textbf{main}})$ and the program terminates successfully with result value $\rho(x)$ on reaching the halt state $(\rho, [], H, (\textbf{return}\ x))$

Notation: we write $\rho[x \mapsto v]$ for the environment which is as $\rho$ but has value $v$ at $x$. We also write $[\vec{x} \mapsto \vec{v}]$ which respectively has values $v_1, \dots, v_n$ at $x_1, \dots, x_n$ and write $[\vec{x} \mapsto \rho(\vec{y})]$ when $v_1, \dots, v_n$ are $\rho(y_1), \dots, \rho(y_n)$.

4

$$\frac{}{\rho, H, k \leadsto H, k} \text{ (ORD-CONST)} \qquad \frac{\ell \notin \mathrm{dom}(H) \text{ is a fresh location}}{\rho, H, (\mathbf{cons}\ x\ y) \leadsto H[\ell \mapsto (\rho(x), \rho(y))], \ell} \text{ (ORD-CONS)}$$

$$\frac{H(\rho(x)) = (v_1, v_2)}{\rho, H, (\mathbf{car}\ x) \leadsto H, v_1} \text{ (ORD-CAR)} \qquad \frac{H(\rho(x)) = (v_1, v_2)}{\rho, H, (\mathbf{cdr}\ x) \leadsto H, v_2} \text{ (ORD-CDR)}$$

$$\frac{}{\rho, H, (\mathbf{id}\ x) \leadsto H, \rho(x)} \text{ (ORD-ID)} \qquad \frac{\rho(x) \in \mathbb{N} \qquad \rho(y) \in \mathbb{N}}{\rho, H, (+\ x\ y) \leadsto H, \rho(x) + \rho(y)} \text{ (ORD-PRIM)}$$

$$\frac{\rho(x) \neq \mathbf{nil}}{\rho, H, (\mathbf{null?}\ x) \leadsto H, 0} \qquad \frac{\rho(x) = \mathbf{nil}}{\rho, H, (\mathbf{null?}\ x) \leadsto H, 1} \text{ (ORD-NULL)}$$

$$\frac{\rho(x) \in \mathbb{N} \setminus \{0\}}{\rho, S, H, (\mathbf{if}\ x\ e_1\ e_2) \longrightarrow \rho, S, H, e_1} \qquad \frac{\rho(x) = 0}{\rho, S, H, (\mathbf{if}\ x\ e_1\ e_2) \longrightarrow \rho, S, H, e_2} \text{ (ORD-IF)}$$

$$\frac{\rho, H, s \leadsto H', v \qquad s \text{ is not } (f\ y_1 \dots y_n)}{\rho, S, H, (\mathbf{let}\ x \leftarrow s\ \mathbf{in}\ e) \longrightarrow \rho[x \mapsto v], S, H', e} \text{ (ORD-LET-NONFN)}$$

$$\frac{s \text{ is } (f\ y_1 \dots y_n) \qquad f \text{ defined as } (\mathbf{define}\ (f\ z_1\ \dots\ z_n)\ e_f)}{\rho, S, H, (\mathbf{let}\ x \leftarrow s\ \mathbf{in}\ e) \longrightarrow [\vec{z} \mapsto \rho(\vec{y})], (\rho, x, e) \bullet S, H, e_f} \text{ (ORD-LET-FNCALL)}$$

$$\frac{}{\rho, (\rho', x', e') \bullet S, H, (\mathbf{return}\ x) \longrightarrow \rho'[x' \mapsto \rho(x)], S, H, e'} \text{ (ORD-RETURN)}$$

**Fig. 3.** The small-step operational semantics

*Stuck states.* Note that certain forms of $e$ do not reduce with $\rightarrow$ (perhaps because $\leadsto$ could not reduce a contained $s$). Some of these we eliminate syntactically, e.g. ensuring all variables and functions are defined and are called with the correct number of parameters. Others include $(\mathbf{cdr}\ \mathbf{nil}), (\mathbf{car}\ 3), (+\ \mathbf{nil}\ 4)$ and $(\mathbf{if}\ \mathbf{nil}\ e_1\ e_2)$. All but the first can be eliminated with a static type system but, treating our program as dynamically typed, we regard all these as stuck states.

## 3 Liveness

In classical liveness analysis a variable is either 'live' (its value may be used in future computation) or 'dead' (definitely not used). Semantically, a variable is dead at a given program point if arbitrary changes to its value have no effect on the computation. Later we will use $\bot$ to represent a value which 'explodes' when it is used in a computation; dead variables can safely have their value replaced with $\bot$. For heap-allocated data we need a richer model of liveness in that both variables and fields of **cons** cells may be dead or live. Using $\mathbf{0}, \mathbf{1}$ to represent access using **car**, **cdr** respectively, liveness of the structure reachable from a variable is a set of *access paths* which we represent as a subset of $\{\mathbf{0}, \mathbf{1}\}^*$, and use conventional grammar notation. Thus the liveness of $x$ being $\{\mathbf{10}, \mathbf{110}\}$ means that future computation can only refer to the second and third members of $x$ considered as a list. Semantically, access paths are prefix-closed, as accessing a field requires accessing all the paths from the variable to the field, and hence the above liveness is properly written $\{\varepsilon, \mathbf{1}, \mathbf{10}, \mathbf{11}, \mathbf{110}\}$. The classical notions of a scalar variable being live or dead correspond to $\{\varepsilon\}$ and $\{\}$.

The overall liveness (also written *liveness environment* for emphasis) at a program point is conceptually a mapping from variables to subsets of $\{0,1\}^*$, but we often abuse notation, for example writing $\{x.01, x.1, y.\varepsilon\}$ instead of the map $[x \mapsto \{\varepsilon, 0, 01, 1\}, y \mapsto \{\varepsilon\}, z \mapsto \{\}]$. Analogously to classical liveness, the liveness at program point $\pi$ in $\pi : e$ is the liveness just before executing $e$.

A complementary notion to liveness is *demand*. The demand for expression $e$ is again an access path—that subset of $\{0,1\}^*$ which the context of $e$ may explore of $e$'s result. So, for example given a demand $\sigma$ and the expression $\pi : (\textbf{return } x)$, the liveness at $\pi$ is exactly $x.\sigma$. The classical analogy of this is in *strong liveness*, where an assignment node $n : x := y + z$ causes $y$ and $z$ to be live on entry to $n$ if (and only if) $x$ is live at exit of $n$—the liveness of $x$ at exit from $n$ becomes the demand on $y + z$. Note that, for an operation like division which may raise an exception, the assignment $n : x := y/z$ makes $y$ and $z$ live regardless of the liveness of $x$.

We use $\sigma$ to range over demands, $\alpha$ to range over access paths and $L$ to range over liveness environments. The notation $\sigma_1\sigma_2$ denotes the set $\{\alpha_1\alpha_2 \mid \alpha_1 \in \sigma_1, \alpha_2 \in \sigma_2\}$. Often we shall abuse notation to juxtapose an edge label and a set of access paths: $0\sigma$ is a shorthand for $\{0\}\sigma$. Finally, we use $\mathsf{LF}$ to range over *demand transformers*; given user function $f$, $\mathsf{LF}_f$ transforms demands on a call to $f$ into demands on its formal parameters: if $f$ is defined by $(\textbf{define } (f\ x_1\ \ldots\ x_n)\ e_f)$ and called with demand $\sigma$, then $\mathsf{LF}_f^i(\sigma)$ is the liveness of $x_i$ at $e_f$.

Note that liveness refers to variables and fields, and not to **cons** cells (i.e. to edges in the memory graph, not to locations themselves). Hence liveness of $\{x.\varepsilon, x.0\}$ means that future computation may refer to the value $\ell$ of variable $x$, and also to the **car** field of location $\ell$. In the absence of other pointers to heap location $\ell$, we are certain that the **cdr** field of $\ell$ will not be referenced and may hence be corrupted arbitrarily. Note therefore, that while $\ell$ cannot be garbage collected, any location $\ell'$ stored in the **cdr** field of $\ell$ would be garbage (again provided there are no other aliases to $\ell$ or $\ell'$).

## 3.1 Liveness Analysis

First recall the classical formulation of liveness (as sets of simple variables) on three-address instructions, $live_{in}(I) = live_{out}(I) \setminus def(I) \cup ref(I)$, and then note that *strong liveness* needs, when $I$ is the instruction $z := x + y$, that $ref(I)$ be refined to $\{x, y\}$ if $z \in live_{out}(I)$ and $\{\}$ otherwise.

Our liveness analysis formulated in Figure 4 is analogous. Firstly, the function *ref*, when given a statement $s$, returns the liveness *generated* by $s$. Because we generalise *strong* liveness, *ref* needs a second parameter, specifying the demand $\sigma$ on the result of $s$, to determine which access paths of its free variables are made live. The cases for $(\textbf{id } x)$ and $(+\ x\ y)$ exemplify this. A demand of $\sigma$ on $(\textbf{car } x)$ is transformed to the demand $0\sigma$ on $x$. In addition, **car** always dereferences its argument (even if its result is never used). This generates the liveness $\{x.\varepsilon\} \cup x.0\sigma$ (note $\sigma$ may be $\{\}$). In the opposite sense, the demand of $0\sigma$ on $(\textbf{cons } x\ y)$ is transformed to the demand $\sigma$ on $x$. Note that **cons** does not, by itself, dereference its arguments. Thirdly, for the case of a user-function call, a third parameter $\mathsf{LF}$ to *ref* expresses how the demand $\sigma$ on the result is transformed into demands on its parameters. Constants generate no liveness.

$$ref(\kappa, \sigma, \mathsf{LF}) = \{\,\}, \text{ for } \kappa \text{ a constant, including } \mathbf{nil}$$

$$ref((\mathbf{cons}\ x\ y), \sigma, \mathsf{LF}) = \{x.\alpha \mid \mathbf{0}\alpha \in \sigma\} \cup \{y.\alpha \mid \mathbf{1}\alpha \in \sigma\}$$

$$ref((\mathbf{car}\ x), \sigma, \mathsf{LF}) = \{x.\varepsilon\} \cup \{x.\mathbf{0}\alpha \mid \alpha \in \sigma\}$$

$$ref((\mathbf{cdr}\ x), \sigma, \mathsf{LF}) = \{x.\varepsilon\} \cup \{x.\mathbf{1}\alpha \mid \alpha \in \sigma\}$$

$$ref((\mathbf{id}\ x), \sigma, \mathsf{LF}) = \{x.\sigma\}$$

$$ref((+\ x\ y), \sigma, \mathsf{LF}) = \{x.\varepsilon, y.\varepsilon\}$$

$$ref((\mathbf{null?}\ x), \sigma, \mathsf{LF}) = \{x.\varepsilon\}$$

$$ref((f\ y_1\ \cdots\ y_n), \sigma, \mathsf{LF}) = \bigcup_{i=1}^{n} y_i.\mathsf{LF}_f^i(\sigma)$$

$$\mathcal{L}((\mathbf{return}\ x), \sigma, \mathsf{LF}) = x.\sigma$$

$$\mathcal{L}((\mathbf{if}\ x\ e_1\ e_2), \sigma, \mathsf{LF}) = \mathcal{L}(e_1, \sigma, \mathsf{LF}) \cup \mathcal{L}(e_2, \sigma, \mathsf{LF}) \cup \{x.\varepsilon\}$$

$$\mathcal{L}((\mathbf{let}\ x \leftarrow s\ \mathbf{in}\ e), \sigma, \mathsf{LF}) = \mathsf{L} \setminus x.\{\mathbf{0}, \mathbf{1}\}^* \cup ref(s, \mathsf{L}(x), \mathsf{LF}), \text{ where } \mathsf{L} = \mathcal{L}(e, \sigma, \mathsf{LF})$$

$$\frac{\mathcal{L}(e_f, \sigma, \mathsf{LF}) = \bigcup_{i=1}^{n} z_i.\mathsf{LF}_f^i(\sigma) \text{ for each } f \text{ and } \sigma}{d_1 \ldots d_k \vdash^l \mathsf{LF}} \quad \text{(LIVE-DEFINE)}$$

$$\text{where } (\mathbf{define}\ (f\ z_1\ \ldots\ z_n)\ e_f) \text{ is a member of } d_1 \ldots d_k$$

**Fig. 4.** Liveness equations and judgement rule

The function $\mathcal{L}$ now gives the (total) liveness of an expression $e$. The cases **return** and **if** are straightforward, but note the liveness $x.\varepsilon$ generated by the latter. The case $(\mathbf{let}\ z \leftarrow s\ \mathbf{in}\ e')$ resembles a three-address instruction: the liveness of $e$ is given by taking the liveness, $\mathsf{L}$, of $e'$, killing any liveness of $z$ and adding any liveness generated by $s$. The main subtlety is how the liveness of $z$ in $\mathsf{L}$ is converted to a demand $\mathsf{L}(z)$ to be placed on $s$ via $ref(s, \mathsf{L}(z), \mathsf{LF})$.

Finally, the judgement form $Prog \vdash^l \mathsf{LF}$ is used to determine $\mathsf{LF}$. Analogously to classical liveness being computed as a solution of dataflow equations, we require, via inference rule (LIVE-DEFINE), $\mathsf{LF}$ to satisfy the fixed-point property that: when we assume $\mathsf{LF}$ to be the demand transformer for the program then the calculated liveness of each function body $\mathcal{L}(e_f, \sigma, \mathsf{LF})$ agrees with the assumed $\mathsf{LF}_f$. As usual, there are often multiple solutions to $\mathsf{LF}$; all are safe (see Section 3.2) but we prefer the least one as giving the least liveness subject to safety—and hence greatest amount of garbage collected.

We make three observations: firstly the rule (LIVE-DEFINE) has a least solution as $\mathcal{L}(\cdot)$ is monotonic in $\sigma$; secondly that (LIVE-DEFINE) resembles the rule for type inference of mutually recursive function definitions, and thirdly the asymmetry of demand and liveness (compared to post- and pre-liveness classically) is due to the functional formulation here.

Section 4 shows how the demand transformer $\mathsf{LF}$ for a program (representing a fully context-sensitive analysis) can be safely approximated, for each function, by a *procedure summary* (unifying the contexts in the style of 0-CFA). The summary consists of a pair of a single demand and, for this demand, the corresponding tuple of demands the function makes on its arguments.

### 3.2 Minefield semantics and correctness

This section gives a modified semantics which checks liveness annotations at run time, and 'explodes' when these are found to be inconsistent with execution behaviour, but otherwise behaves as the standard semantics. We show that such explosions never occur and hence run-time checks can be elided. We first assume an arbitrary demand transformer LF (below we assume $Prog \vdash^l$ LF). We then enrich the abstract machine state $\rho, S, H, e$ to $\rho, S, H, e, \sigma, \Sigma$. Here $\sigma$ is the demand to the currently active function, thus the liveness L at $e$ is $\mathcal{L}(e, \sigma, \text{LF})$,[4] and $\Sigma$ is a stack of demands—one for each function frame pushed in S.

Second, we augment *Val* with a value $\perp$. To model strong liveness $\perp$ may be copied freely, including into a **cons** cell, but explodes when used computationally (in a primitive operation other than a copy). Additionally we define $GC(\text{L}, \Sigma) : (\rho, H, S) \mapsto (\rho', H', S')$ which determines *live-reachability*[5] using $\rho$ and the $\rho$'s in S as the root set and following links in H *only as far as allowed by* L *and* $\Sigma$. Hence $GC(\text{L}, \Sigma)$ replaces live-unreachable values—in $\rho$, in H and the $\rho$'s in S—with $\perp$. For example, if $x.\varepsilon \notin \text{L}$ then $\rho'(x) = \perp$. Only $GC(\cdots)$ introduces $\perp$.

Third, we update the semantics in four ways: (*i*) we arrange that all ($\rightarrow$) transitions on expressions $e$ first use $GC(\cdots)$ to update the state and then continue as before; and (*ii*) whenever the value $\perp$ is used computationally in a reduction ($\rightsquigarrow$), we enter a distinguished stuck state BANG. For example, supposing $\rho(x) = \perp$ then $\rho, H, s \rightsquigarrow$ BANG if $s$ is (**car** $x$), (**cdr** $x$) or ($+ x y$), but not if $s$ is (**id** $x$), (**cons** $x y$) or ($f x y$). Finally (*iii*) we make a similar change to the ($\rightarrow$) reduction for (ORD-IF) and (*iv*) augment the (ORD-LET) rule for primitives to propagate BANG from ($\rightsquigarrow$) to ($\rightarrow$).

The resulting minefield semantics behaves identically (identical heap, identical steps, including possible non-termination) to the standard semantics, except for the sole possibility of the minefield semantics going BANG while the standard semantics continues (either to a halt state, to a stuck state, or reduces forever).

We now prove a result that relates liveness analysis to the semantics.

**Proposition 1.** *Given program P with $P \vdash^l$ LF, then in the minefield semantics $P \rightarrow$ BANG can never occur (cf. 'well-typed programs do not go wrong').*

*Proof outline:* Space does not permit a full proof, but we give the two main steps. We proceed by contradiction and assume there is a program $P$ for which $P \vdash^l$ LF can enter state BANG. The first step is to construct a program $P'$ with identical behaviour, but with no user-function definitions, using inlining. This is possible because a program which goes BANG does so after a finite number of reductions, and hence even recursive functions have only had a finite number of invocations. We hence repeatedly inline-expand user-function calls in $P$ until we obtain a program $P'$ which behaves identically[6] to $P$ in the standard semantics, but executes no user-function calls. Any remaining, non-executed, calls can be replaced with a new primitive with the same demand-to-liveness transfer function—thus making $P'$ a simple expression $e'$. Not only do the program

---

[4] This a simple liveness propagation using $\mathcal{L}(\cdots)$ and $ref(\cdots)$ as LF is assumed given.

[5] Reachability curtailed by liveness information.

[6] Modulo replacement of (ORD-CALL) and (ORD-RETURN) steps with (ORD-ID) steps.

points in reducing $e'$ correspond one-one to states during evaluation of $P$, but also the liveness associated with a point in $e'$ is identical to the liveness at the corresponding state $(\rho, S, H, e, \sigma, \Sigma)$ of $P$ (concatenating the liveness $\mathcal{L}(e, \sigma, \mathsf{LF})$ at $e$ with the liveness, obtained from $\Sigma$, of the call sites in $S$ and after renaming the variables corresponding to the inlining which produced $e'$). This assertion relies on the analysis being fully context-sensitive, and noting that while the change of scope caused by inlining changes variable visibility between $P$ and $e'$ it does not change the liveness—as local-to-a-function **let**-variables whose scope has been prolonged due to inlining are dead in the prolonged scope.

The second step of the proof is to show that $e'$ cannot go BANG. We proceed by induction. Correctness of the **if** and final **return** forms are immediate; the (**let** $z \leftarrow s$ **in** $e$) form requires showing the inference rules ensure that any value referenced via $z$ in $e$ was already live-reachable (via another variable and path) in any enclosing expression (so that $GC(\cdots)$ could not re-write it to $\bot$).

## 4 Computing liveness and its encoding as a table

Section 3 gave a context-sensitive liveness analysis and proved it correct with reference to a *minefield* semantics. For practical use we need to solve the liveness equations *finitely and symbolically*. As expressed mathematically, and given a fixed program, three things are potentially unbounded: (*i*) the number of call strings (and hence arguments $\sigma$ to $\mathsf{LF}$); (*ii*) the length of access paths $\alpha \in \sigma$ and (*iii*) the number of such access paths.

We commonly first solve (*i*) by reducing the number of distinct calling contexts (e.g. to a single unified context in 0-CFA style). However, it turns out that we can solve the equations for $\mathsf{LF}$ symbolically without this reduction, so here we defer this to Section 4.2. (It also allows easier extension to dynamically determined liveness—future work.) We address (*ii*) and (*iii*) by re-interpreting the liveness definitions in Figure 4 symbolically as a grammar rather than a mutually recursive set of equations on sets of access paths.

This requires two ideas. Firstly we have to control the use of functions—they tend to be infinitary and do not occur naturally in CFGs; in particular our $\mathsf{L}$ maps names to access paths, and $\mathsf{LF}$ maps access paths to a tuple of access paths. The former is achieved by using a separate meta-variable (later non-terminal) $\mathsf{L}_i^x$ for each variable $x$ and each program point $\pi_i$ ($\mathsf{L}_i^x$ represents $\mathsf{L}(x)$ at $\pi_i$). Section 4.1 shows how the latter $\mathsf{LF}_f^i$ is also expressible finitely (it is a linear form).

Secondly, there are also two technical issues in re-interpreting Figure 4 as a grammar. One is the use of the set-difference operator $\setminus$ in "$\cdots \setminus x.\{\mathbf{0}, \mathbf{1}\}^*$" reflecting the classical *gen/kill* dataflow formulation. However after separating, as above, liveness environments into per-variable liveness $\mathsf{L}_i^x$ the '$\setminus$' operator reduces to the harmless grammar rule $\mathsf{L}_i^x = \{\}$ (assuming $i$ labels the (**let** $x$ ...) expression). The other is that *ref* for **cons** decomposes strings and thus gives a general grammar not a CFG. Below we show how symbols $\bar{\mathbf{0}}, \bar{\mathbf{1}}$ can give an equivalent CFG.

Finally, Section 4.3 uses a construction due to Mohri and Nederhof [10] to over-approximate context-free grammars with regular grammars; these are more appropriate for run-time use. Hence the overall 'big picture' view is that each GC point is annotated

with a table encoding the DFA for that program point. When garbage collection occurs, each saved return address on the run-time stack identifies the call-site GC point. The DFA annotating each such GC point is then used by the garbage collector to curtail (to access paths accepted by the DFA) its local-variable reachability-based marking.

**GC points.** Given a call site (an expression $\pi_1 : (\textbf{let } x \leftarrow (f \ y_1 \ldots y_n) \textbf{ in } \pi_2 : e))$ its associated GC point is $\pi_2$, as it is liveness at $\pi_2$ that should be encoded in the DFA associated with the call.[7] In the case of a call to **cons** as in $\pi_1 : (\textbf{let } x \leftarrow (\textbf{cons } y_1 \ y_2) \textbf{ in } \pi_2 : e)$, the situation is slightly more complex. We may either treat **cons** as doing a full procedure call (and mark its formal parameters separately during garbage collection, which again leads to its GC point being $\pi_2$), or we may regard **cons** as being inlined, in which case it is vital that liveness of $y_1$ and $y_2$ are represented in the DFA (which is achieved by using $\pi_1$ rather than $\pi_2$ as the GC point). We adopt the latter approach.

**Modifying the cons rule.** The *ref* rule for **cons**, shown in Figure 4, requires us to remove the leading **0** and **1** from access paths in $\sigma$. Mathematically this is fine but causes problems when solving the liveness equations *symbolically* since such decomposition cannot be expressed as a context-free grammar. To handle this, we introduce two new symbols $\bar{\textbf{0}}$ and $\bar{\textbf{1}}$ with the properties:

$$\bar{\textbf{0}}\sigma \triangleq \{\alpha \mid \textbf{0}\alpha \in \sigma\} \qquad \text{and} \qquad \bar{\textbf{1}}\sigma \triangleq \{\alpha \mid \textbf{1}\alpha \in \sigma\}$$

We can now rewrite the **cons** rule as:

$$ref((\textbf{cons } x \ y), \sigma, \textsf{LF}) = x.\bar{\textbf{0}}\sigma \cup y.\bar{\textbf{1}}\sigma$$

We call the liveness equations with this modification $\mathcal{L}'$. The definitions of $\bar{\textbf{0}}$ and $\bar{\textbf{1}}$ induce the following relation $\hookrightarrow$ over sets of access paths:

$$\sigma_1\bar{\textbf{0}}\sigma_2 \hookrightarrow \sigma_1\sigma_2', \text{ where } \sigma_2' = \{\alpha \mid \textbf{0}\alpha \in \sigma_2\}, \quad \text{and}$$
$$\sigma_1\bar{\textbf{1}}\sigma_2 \hookrightarrow \sigma_1\sigma_2', \text{ where } \sigma_2' = \{\alpha \mid \textbf{1}\alpha \in \sigma_2\}$$

The reflexive transitive closure of $\hookrightarrow$ will be denoted as $\overset{*}{\hookrightarrow}$. The following proposition relates $\mathcal{L}$ and $\mathcal{L}'$:

**Proposition 2.** *Assume that a liveness computation based on $\mathcal{L}$ gives the liveness of the variable $x$ at a program point $\pi_i$ as $\sigma$ (symbolically, $\mathsf{L}_i^x = \sigma$). Further, suppose $\mathsf{L}_i^x = \sigma'$ when $\mathcal{L}'$ is used for liveness computation instead of $\mathcal{L}$. Then $\sigma' \overset{*}{\hookrightarrow} \sigma$.*

To see why the proposition is true, consider an analysis based on $\mathcal{L}'$ in which $\sigma$ appears in the context $ref((\textbf{cons } x \ y), \sigma, \textsf{LF})$. Let $\alpha \in \sigma$. The symbol $\bar{\textbf{0}}$ (respectively $\bar{\textbf{1}}$) merely marks a place in $\alpha$ where the original **cons** rule would have erased an immediately following **0** (respectively **1**), or, in absence of such a symbol, would have dropped $\alpha$ itself. Since the application of any rule in $\mathcal{L}'$ merely adds symbols at the beginning of $\alpha$, the markers and other symbols in $\alpha$ are propagated to other dependent parts of program in their same relative positions. Consequently, the erasure carried out at the end of the analysis with $\overset{*}{\hookrightarrow}$ gives the same result as obtained through $\mathcal{L}$.

[7] A subtlety is that at machine code level the assignment to $x$ does not take place until after the call, and so for garbage-collection purposes the DFA need not represent liveness of $x$.

```
(define (append l1 l2)
    π₁:(let test ← (null? l1) in
        π₂:(if test π₃:(return l2)
            π₄:(let tl ← (cdr l1) in
                π₅:(let rec ← (append tl l2) in
                    π₆:(let hd ← (car l1) in
                        π₇:(let ans ← (cons hd rec) in (return ans))))))))


πₘₐᵢₙ: ...
        π₈:(let y ← (append a b) in
            π₉:(let w ← (append y z) in
                π₁₀: (let c ← (cdr w) in
                    π₁₁: (let d ← (car c) in  (return d)))))))
```

**Fig. 5.** An example program. GC points are $\pi_6$, $\pi_7$, $\pi_9$ and $\pi_{10}$.

### 4.1 Generating equations for the demand transformer LF

We shall consider the program in Figure 5 as a running example. Unlike the program in Figure 1, this program is in ANF.

To generate the equations defining $\mathsf{LF}_f$, we follow the rule DEFINE-LIVE. We start with a *symbolic demand* $\sigma$ and determine $\mathsf{L} = \mathcal{L}(e_f, \sigma, \mathsf{LF})$, treating $\mathsf{LF}$ as an uninterpreted function symbol. We then generate equations of the form $\mathsf{LF}_f^i(\sigma) = \mathsf{L}(x_i)$ where $x_i$ is the $i$th formal parameter of $f$ and $\mathsf{L}(x_i)$ is the liveness of $x_i$. For our example program which has a single function **append**, this generates the following equations:

$$\mathsf{LF}_{\mathbf{append}}^1(\sigma) = \{\epsilon\} \cup \mathbf{0\bar{0}}\sigma \cup \mathbf{1}\mathsf{LF}_{\mathbf{append}}^1(\bar{\mathbf{1}}\sigma)$$
$$\mathsf{LF}_{\mathbf{append}}^2(\sigma) = \sigma \cup \mathsf{LF}_{\mathbf{append}}^2(\bar{\mathbf{1}}\sigma)$$

In general, the equations for $\mathsf{LF}$ are recursive since $\mathsf{L}$ may, in turn, be expressed in terms of $\mathsf{LF}$. We assume that $\mathsf{LF}_f$ is expressible in the closed form as:[8]

$$\mathsf{LF}_f^i(\sigma) = \mathsf{I}_f^i \cup \mathsf{D}_f^i \sigma \qquad (1)$$

where $\mathsf{I}_f^i$ and $\mathsf{D}_f^i$ are sets of strings over the alphabet $\{\mathbf{0}, \mathbf{1}, \bar{\mathbf{0}}, \bar{\mathbf{1}}\}$. The reason why $\mathsf{LF}$ has this form is as follows. Recall that $\mathsf{LF}_f^i(\sigma)$ gives the access paths starting from $i$ that have to be dereferenced to produce the sub-structure $\sigma$ of the result of $f$. $\mathsf{I}_f^i$ represents the access paths that would be dereferenced, but do not contribute to the result. This happens, for instance, when the argument is used only within the condition of an **if**. $\mathsf{D}_f^i$, in contrast, represents the paths that are dereferenced to actually produce the result.

---

[8] This is similar to solving the differential equation $ay'' + by' + c = 0$, where we guess that the solution has the form $y = e^{rx}$. Substituting the solution in the equation yields a quadratic equation in $r$, and each solution of $r$ gives rise to a solution of the differential equation (in our setup we can effectively pick the least solution rather than needing linear combinations).

11

To solve for $\mathsf{LF}_f$, we substitute the guessed form into its equations. $\mathsf{LF_{append}}$ gives:

$$\mathsf{I^1_{append}} \cup \mathsf{D^1_{append}}\sigma = \{\epsilon\} \cup \mathbf{0\bar{0}}\sigma \cup \mathbf{1}(\mathsf{I^1_{append}} \cup \mathsf{D^1_{append}}\bar{\mathbf{1}}\sigma)$$

$$\mathsf{I^2_{append}} \cup \mathsf{D^2_{append}}\sigma = \sigma \cup \mathsf{I^2_{append}} \cup \mathsf{D^2_{append}}\bar{\mathbf{1}}\sigma$$

Equating the terms containing $\sigma$ on the two sides of each equation, and doing the same for the terms without $\sigma$, we get equations for $\mathsf{I}^i_f$ and $\mathsf{D}^i_f$ that are independent of $\sigma$.

$$\mathsf{I^1_{append}} = \{\epsilon\} \cup \mathbf{1}\mathsf{I^1_{append}} \qquad\qquad \mathsf{I^2_{append}} = \mathsf{I^2_{append}}$$

$$\mathsf{D^1_{append}} = \{\mathbf{0\bar{0}}\} \cup \mathbf{1}\mathsf{D^1_{append}}\bar{\mathbf{1}} \qquad\qquad \mathsf{D^2_{append}} = \{\epsilon\} \cup \mathsf{D^2_{append}}\bar{\mathbf{1}}$$

Note that these equations can be viewed as CFGs, with all but $\mathsf{D^1_{append}}$ being regular, and that any solution of $\mathsf{I}^i_f$ and $\mathsf{D}^i_f$ yields a solution of $\mathsf{LF}_f$.

### 4.2   Generating liveness equations $\mathsf{L}$ for function bodies

We now calculate a 0-CFA-style summary liveness for each GC point of a program. There are two parts to this. First, for each function $f$, we determine a *summary demand* $\sigma_f$ over-approximating any demand $\sigma$ passed to $f$. Such demands are caused by calls to $f$ occurring at call sites. We introduce the notation $\delta_f(\pi, g)$ for the contribution to $\sigma_f$ caused a call site $\pi$ occurring in function $g$. So, suppose function $g$ contains a call site $\pi$ to $f$, say $\pi : (\mathbf{let}\ x \leftarrow (f\ y_1 \ldots y_n)\ \mathbf{in}\ e)$. Under the assumption that the demand on $g$ is $\sigma_g$, the liveness at $e$ is $\mathsf{L} = \mathcal{L}(e, \sigma_g, \mathsf{LF})$, and the **let** case of Figure 4 tells us this call site contributes $\mathsf{L}(x)$ to the demand $\sigma_f$ placed on $f$; hence $\delta_f(\pi, g)$ is simply $\mathsf{L}(x)$.

Now, supposing the $k$ call sites to function $f$ are $\pi^1$ (in function $g^1$) $\ldots \pi^k$ (in function $g^k$), then the over-approximation requirement on $\sigma_f$ is achieved by taking $\sigma_f = \delta_f(\pi^1, g^1) \cup \cdots \cup \delta_f(\pi^k, g^k)$.

The expression $e_{\mathbf{main}}$ is a special case; we assume it may be called externally with demand $\sigma_{\mathbf{main}} = \{\mathbf{0}, \mathbf{1}\}^*$ (denoted $\sigma_{all}$). This is because any part of its value may be used by the environment—for printing the result, for instance.

For the running example, **append** has calls from **main** at $\pi_9$ and a recursive call at $\pi_5$. So $\sigma_{\mathbf{append}} = \delta_{\mathbf{append}}(\pi_9, \mathbf{main}) \cup \delta_{\mathbf{append}}(\pi_5, \mathbf{append})$. Calculating the $\delta_{\mathbf{append}}(\pi, g)$ for the two call sites, and substituting gives:

$$\sigma_{\mathbf{append}} = (\{\epsilon, \mathbf{1}\} \cup \mathbf{10}\sigma_{all})\ \cup\ \bar{\mathbf{1}}\sigma_{\mathbf{append}}$$

Second, for each function $f$ (possibly **main**) we need the liveness at each contained GC point $\pi$. Given $\sigma_f$ calculated above, this is simply $\mathcal{L}(\pi, \sigma_f, \mathsf{LF})$. For the running example, containing GC points $\pi_6, \pi_7$ in **append** and $\pi_9, \pi_{10}$ in $e_{\mathbf{main}}$, this gives (recall Equation (1) above states $\mathsf{LF}^i_f(\sigma) = \mathsf{I}^i_f \cup \mathsf{D}^i_f.\sigma$):

$$\mathsf{L}^{11}_6 = \{\epsilon\} \cup \mathbf{0\bar{0}}\sigma_{\mathbf{append}} \qquad\qquad \mathsf{L}^{\mathtt{rec}}_6 = \bar{\mathbf{1}}\sigma_{\mathbf{append}}$$

$$\mathsf{L}^{\mathtt{hd}}_7 = \bar{\mathbf{0}}\sigma_{\mathbf{append}} \qquad\qquad\quad \mathsf{L}^{\mathtt{rec}}_7 = \bar{\mathbf{1}}\sigma_{\mathbf{append}}$$

$$\mathsf{L}^{\mathtt{y}}_9 = \mathsf{LF^1_{append}}(\{\epsilon, \mathbf{1}\} \cup \mathbf{10}\sigma_{all}) \qquad \mathsf{L}^{\mathtt{z}}_9 = \mathsf{LF^2_{append}}(\{\epsilon, \mathbf{1}\} \cup \mathbf{10}\sigma_{all})$$

$$\mathsf{L}^{\mathtt{w}}_{10} = \{\epsilon, \mathbf{1}\} \cup \mathbf{10}\sigma_{all}$$

In summary, the equations generated during liveness analysis are:

1. For each function $f$, equations defining $I_f^i$ and $D_f^i$ for use by $\mathsf{LF}_f$.
2. For each function $f$, an equation defining the summary demand $\sigma_f$ on $e_f$.
3. For each function $f$ (including **main** for $e_{\mathbf{main}}$) an equation defining liveness at each GC point of $e_f$.

### 4.3   Solving liveness equations—the grammar interpretation

The liveness equations above (of the form $X = \ldots$) can now be re-interpreted as a context-free grammar (CFG) on the alphabet $\{\mathbf{0}, \mathbf{1}, \bar{\mathbf{0}}, \bar{\mathbf{1}}\}$. We use $\langle X \rangle$ to denote the corresponding non-terminal which then appears in a production $\langle X \rangle \to \ldots$ . We can think of the resulting productions as being associated with several grammars, one for each non-terminal $\langle L_i^x \rangle$ regarded as a start symbol. As an example, the grammar for $\langle L_9^y \rangle$ comprises the following productions:

$$\langle L_9^y \rangle \to \langle I_{\mathbf{append}}^1 \rangle \mid \langle D_{\mathbf{append}}^1 \rangle (\varepsilon \mid \mathbf{1} \mid \mathbf{10}\langle \sigma_{all} \rangle)$$
$$\langle I_{\mathbf{append}}^1 \rangle \to \varepsilon \mid \mathbf{1}\langle I_{\mathbf{append}}^1 \rangle$$
$$\langle D_{\mathbf{append}}^1 \rangle \to \mathbf{0}\bar{\mathbf{0}} \mid \mathbf{1}\langle D_{\mathbf{append}}^1 \rangle \bar{\mathbf{1}}$$
$$\langle \sigma_{all} \rangle \to \varepsilon \mid \mathbf{0}\langle \sigma_{all} \rangle \mid \mathbf{1}\langle \sigma_{all} \rangle$$

Other equations can be converted similarly. The language generated by $\langle L_i^x \rangle$, denoted $\mathscr{L}(\langle L_i^x \rangle)$, is the desired solution of $L_i^x$. However, recall from our earlier discussion that the decision problem that we are interested in during garbage collection is:

Let $x.\alpha$ be a *forward access path*—consisting only of edges $\mathbf{0}$ and $\mathbf{1}$ (but not $\bar{\mathbf{0}}$ or $\bar{\mathbf{1}}$). Let $\mathscr{L}(\langle L_i^x \rangle) \overset{*}{\hookrightarrow} \sigma$, where $\sigma$ consists of forward paths only. Then does $\alpha \in \sigma$?

We could convert the rules defining $\hookrightarrow$ into productions and add them to the grammar. However, this results in an *unrestricted grammar* [11], and the membership problem for such grammars is undecidable. We circumvent the problem by over-approximating the CFG generated by the analysis to *strongly regular* CFGs which have easy translations to non-deterministic finite state automata (NFA). The NFAs are then simplified on the basis of the $\hookrightarrow$ rules to enable checking of membership of forward access paths. The resulting NFAs are finally converted to DFAs for use during garbage collection.

**Approximating CFGs using NFAs.** We use the algorithm by Mohri and Nederhof [10] to approximate a CFG to a *strongly regular* grammar. The transformation has the property that if $\mathsf{L}$ is a non-terminal in the grammar $G$ and $G'$ is the grammar after the Mohri-Nederhof transformation, then $\mathscr{L}_G(\mathsf{L}) \subseteq \mathscr{L}_{G'}(\mathsf{L})$. This is required for the approximation to be safe with respect to liveness.

We exemplify the Mohri-Nederhof transformation on the $\langle L_9^y \rangle$ grammar above. We pick the only production that is affected by the transformation—the production for $D_{\mathbf{append}}^1$. The production for $I_{\mathbf{append}}^1$, while recursive, is already in strongly regular form and is therefore unaffected by the transformation.

$$\langle D_{\mathbf{append}}^1 \rangle \to \mathbf{0}\bar{\mathbf{0}}\langle D_{\mathbf{append}}^1 \rangle' \mid \mathbf{1}\langle D_{\mathbf{append}}^1 \rangle$$
$$\langle D_{\mathbf{append}}^1 \rangle' \to \bar{\mathbf{1}}\langle D_{\mathbf{append}}^1 \rangle' \mid \varepsilon$$

```
Input: NFA $\overline{N}$ with underlying alphabet $\{\mathbf{0}, \mathbf{1}, \bar{\mathbf{0}}, \bar{\mathbf{1}}\}$
Output: NFA $N$ with underlying alphabet $\{\mathbf{0}, \mathbf{1}\}$ such that $\mathscr{L}(\overline{N}) \overset{*}{\hookrightarrow} \mathscr{L}(N)$.
Steps:
    $i \leftarrow 0$
    $N_0 \leftarrow$ Equivalent NFA of $\overline{N}$ without ε-moves [11]
    repeat
        $N'_{i+1} \leftarrow N_i$
        for all states $q$ in $N_i$ such that $q$ has an incoming edge from $q'$ with label $\bar{\mathbf{0}}$ and outgoing
        edge to $q''$ with label $\mathbf{0}$ do
            add an edge in $N'_{i+1}$ from $q'$ to $q''$ with label ε. {bypass $\bar{\mathbf{0}}\mathbf{0}$ using ε}
        end for
        for all states $q$ in $N_i$ such that $q$ has an incoming edge from $q'$ with label $\bar{\mathbf{1}}$ and outgoing
        edge to $q''$ with label $\mathbf{1}$ do
            add an edge in $N'_{i+1}$ from $q'$ to $q''$ with label ε. {bypass $\bar{\mathbf{1}}\mathbf{1}$ using ε}
        end for
        $N_{i+1} \leftarrow$ Equivalent NFA of $N'_{i+1}$ without ε-moves
        $i \leftarrow i + 1$
    until $(N_i = N_{i-1})$
    $N \leftarrow N_i$
```

**Fig. 6.** Algorithm for transforming an NFA to accept forward paths only.

The languages generated for $\langle \mathrm{D}^1_{\mathbf{append}} \rangle$ in the original grammar and the new grammar
are $\mathbf{1}^i\mathbf{0}\bar{\mathbf{0}}\bar{\mathbf{1}}^i$ and $\mathbf{1}^*\mathbf{0}\bar{\mathbf{0}}\bar{\mathbf{1}}^*$, showing a loss of precision.

**Transforming NFAs to accept forward paths:** The strongly regular CFGs obtained
after the Mohri-Nederhof transformation are first converted into NFAs. The algorithm
described in Figure 6 converts an NFA $\overline{N}$ to a NFA $N$ such that $\mathscr{L}(\overline{N}) \overset{*}{\hookrightarrow} \mathscr{L}(N)$, where
$N$ accepts forward paths only. Thus $N$ can be used to check membership of forward
paths.

The algorithm repeatedly introduces ε edges to bypass a pair of consecutive edges
labelled $\bar{\mathbf{0}}\mathbf{0}$ or $\bar{\mathbf{1}}\mathbf{1}$. The process is continued until a fixed point is reached. When the
fixed point is reached, the resulting NFA contains all possible reductions corresponding
to all the paths in the original NFA. The proofs of the termination and correctness of
the algorithm are given in our earlier paper [8].

We illustrate the algorithm in Figure 6 by constructing the automaton for $\langle \mathrm{L}^{\mathrm{y}}_9 \rangle$. Fig-
ure 7 (c) shows the automaton for $\langle \mathrm{L}^{\mathrm{y}}_9 \rangle$ constructed by composing the automata for
$\langle \mathrm{l}^1_{\mathbf{append}} \rangle$, $\langle \mathrm{D}^1_{\mathbf{append}} \rangle$ and $(\varepsilon \mid \mathbf{1} \mid \mathbf{10}\sigma_{all})$. After ε removal we get (d). We add an ε edge
from $q_2$ to $q_3$ bypassing the $\bar{\mathbf{1}}\mathbf{1}$ pair from $q_2$ to $q_2$ and then to $q_3$. This is shown in (e).
The ε edge is removed in (f) and a second ε is added to the automaton bypassing the $\bar{\mathbf{0}}\mathbf{0}$
pair from $q_1$ to $q_2$. Removing this ε edge gives the automaton shown in (g). Restricting
this automaton to forward edges only, we get the final automaton shown in (h). This
automaton recognises $\mathbf{1}^* \mid \mathbf{1}^*\mathbf{0}\sigma_{all}$, showing that the entire list y, including its elements,
is live at $\pi_9$. Also note that the language accepted by the final automaton satisfies the
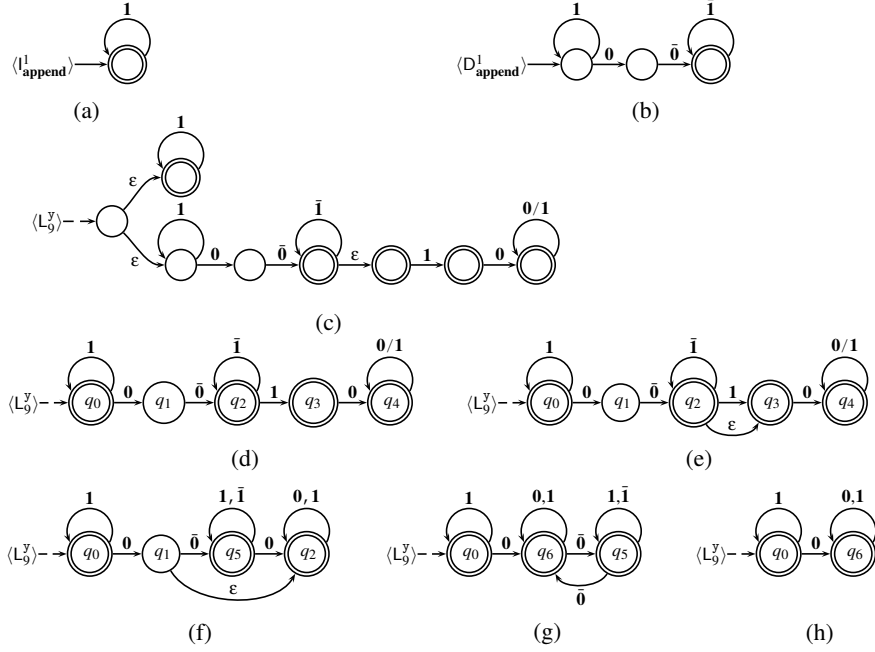prefix-closed property.

**Fig. 7.** Automata for the example program

## 5 Prototype and evaluation

To demonstrate the effectiveness of liveness-based garbage collection, we have built a prototype consisting of an interpreter for our language, a liveness analyser and a copying collector that can optionally use the results of liveness analysis for marking instead of reachability. When the collector uses liveness for marking, we call it a *liveness-based collector* (LGC), else we use the term *reachability-based collector* (RGC). The collector is neither incremental nor generational. As a consequence, any cell that becomes unreachable or dead is assuredly collected in the next round of garbage collection.

When LGC is invoked (by a call to **cons**) the activation records on the stack all correspond to functions suspended at GC points, and by construction at each GC point we have a DFA specifying liveness of each local variable in the activation record. As usual such local variables form the root set for garbage collection.

Let $\mathrm{dfa}_\pi^x$ denote the DFA for the variable and program point pair $(x, \pi)$. We write $\mathrm{initial}(\mathrm{dfa}_\pi^x)$ for the initial state of $\mathrm{dfa}_\pi^x$. Considering a DFA as a table, $\mathrm{dfa}_\pi^x(q, sym)$ returns the next state for the state $q$ and the symbol $sym$, where $sym$ is **0** or **1**. We shall also write $\mathrm{dfa}_\pi^x(q, sym)$? for a predicate indicating whether there is a transition from $q$ on $sym$. The LGC action to chase the root variable $x$ at $\pi$ can be described as follows: If $\mathscr{L}(\mathrm{dfa}_\pi^x)$ is empty, then nothing needs to be done. Otherwise we call $\mathrm{copy}(\mathrm{dfa}_\pi^x, \mathrm{initial}(\mathrm{dfa}_\pi^x), x)$ in Figure 8 and assign the returned pointer to $x$. The function move_to_tospace($x$) copies the value of $x$ in the other semi-space and returns the new

15

```
function copy(dfa, q, x)
    let y ← move_to_tospace(x)
    if x.tag ≠ cons then skip
    else if dfa(q, 0)? then y.car = copy(dfa, dfa(q, 0), x.car)
         if dfa(q, 1)? then y.cdr = copy(dfa, dfa(q, 1), x.cdr)
    return y
```

**Fig. 8.** Function for copying a root set variable.

address. It hides details such as returning the forwarding pointer if the value of $x$ is already copied, and creating the forwarding pointer otherwise.

The graphs in Figure 9[9] show the number of cells in the heap over time for RGC and LGC—here time is measured in terms of the number of **cons** cells allocated. In addition, they also show the number of reachable cells and the number of cells that are actually live (this is statically approximated by our liveness analysis). Since the programs have different memory requirements, we have tuned the size of heap for each program to ensure a reasonable number of collections. An invocation of RGC decreases the number of cells in heap until it touches the curve of reachable cells. An invocation of LGC decreases the number of heap cells to no lower than the curve of live cells.

To construct the reachable and live curves, we record for every cell its creation time (Create_time), its last use time (Use_time), and the earliest time when the cell becomes unreachable and can be garbage collected (Collection_time). For accurate recording of Collection_time, we force frequent invocations of a reachability-based collector in a separate run. A cell is live at time $T$ if Create_time $\leq T \leq$ Use_time. If Create_time $\leq T \leq$ Collection_time, it is reachable.

The benchmark programs are drawn from the no-fib suite and other sources and have been manually converted to ANF. All graphs except fibheap show strictly fewer garbage collector invocations for LGC; fibheap is an exception in that the number of reachable cells first grows steadily until it almost fills the heap. This triggers garbage collections in both LGC and RGC. The number of reachable cells then drops steeply to a low level and remains low resulting in no further garbage collections. The graphs also show the precision of our liveness analysis. For all programs except nperm and lambda, LGC manages to collect a good portion of the cells that are not live.

### 5.1 Results

The increased effectiveness of LGC over RGC is also shown in the tables in Figure 10. The first table provides statistics regarding the analysis itself. The number of states and the analysis times are within tolerable limits. Precision of analysis refers to the percentage of dead cells that is collected by LGC, averaged over all invocations. The second table shows garbage collection statistics for RGC and LGC. LGC collects larger garbage per invocation, drags cells for lesser time and requires a smaller heap size (*MinHeap*) for program to run in comparison with RGC.

---

[9] For better clarity, visit http://www.cse.iitk.ac.in/users/karkare/fhra for coloured versions of the graphs.
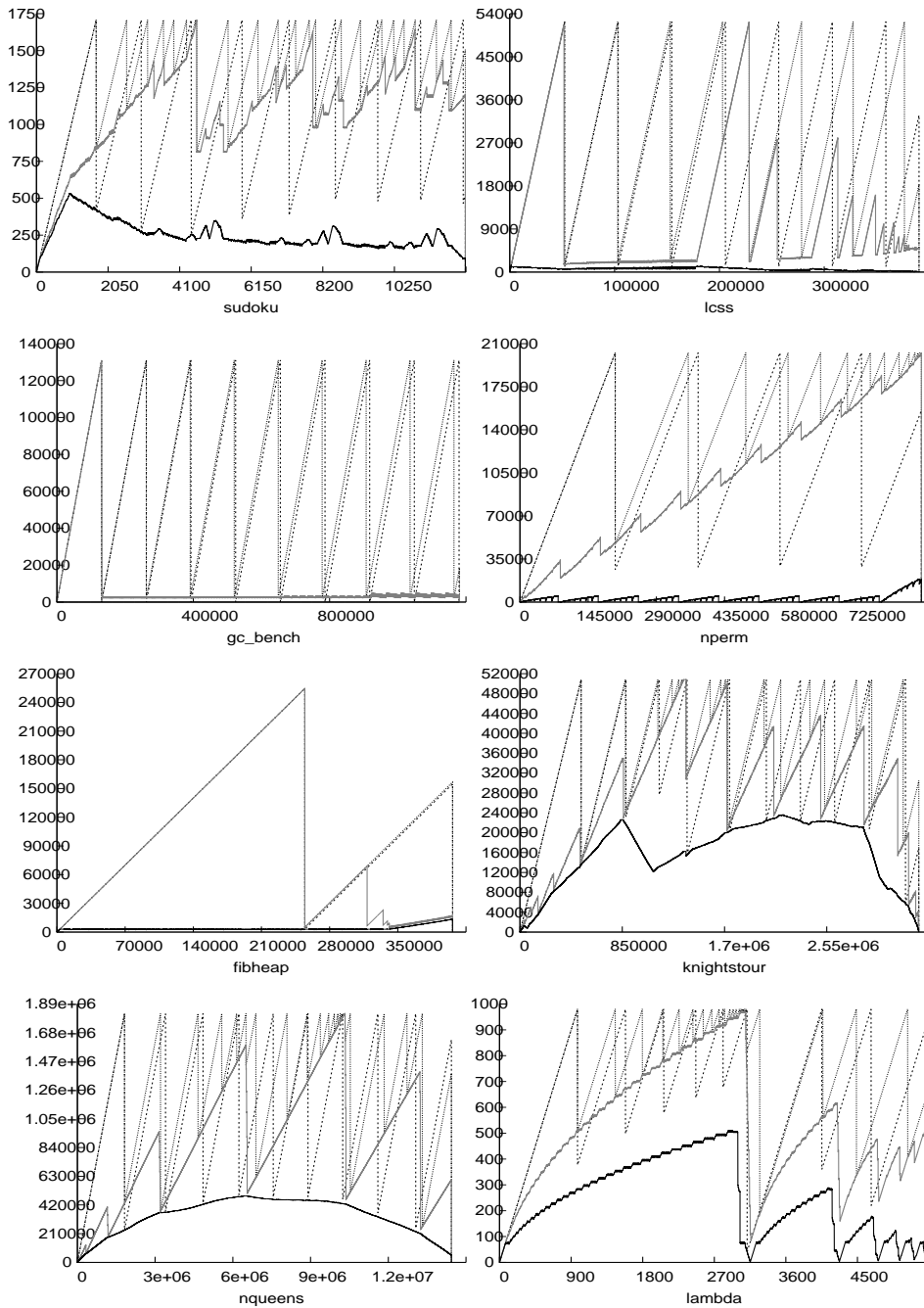
**Fig. 9.** Memory usage of programs. The dotted-grey and the dashed-black curves indicate the number of **cons** cells in the active semi-space for RGC and LGC respectively. The solid-grey curve represents the number of reachable cells and the solid-black curve represents the number of cells that are actually live (of which liveness analysis does a static approximation). x-axis is the time measured in number of **cons** cells allocated. y-axis is the number of **cons** cells.

17

| Program | sudoku | lcss | gc_bench | nperm | fibheap | knightstour | treejoin | nqueens | lambda |
|---|---|---|---|---|---|---|---|---|---|
| Time (msec) | 120.95 | 2.19 | 0.32 | 1.16 | 2.4 | 3.05 | 2.61 | 0.71 | 20.51 |
| DFA size | 4251 | 726 | 258 | 526 | 675 | 922 | 737 | 241 | 732 |
| Precision(%) | 87.5 | 98.8 | 99.9 | 87.1 | 100 | 94.3 | 99.6 | 98.8 | 83.8 |

(a)

| Program | # Collected cells per GC | | # Touched cells per GC | | #GCs | | MinHeap (#cells) | | Avg. Drag (#cells) | | GC time (sec) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | RGC | LGC | RGC | LGC | RGC | LGC | RGC | LGC | RGC | LGC | RGC | LGC |
| sudoku | 490 | 1306 | 1568 | 774 | 22 | 9 | 1704 | 589 | 858 | 146 | .028 | .122 |
| lcss | 46522 | 51101 | 6216 | 1363 | 8 | 7 | 52301 | 1701 | 5147 | 588 | .045 | .144 |
| gc_bench | 129179 | 131067 | 1894 | 4 | 9 | 9 | 131071 | 6 | 16970 | 4 | .086 | .075 |
| nperm | 47586 | 174478 | 201585 | 60882 | 14 | 4 | 202597 | 37507 | 171878 | 76618 | 1.406 | .9 |
| fibheap | 249502 | 251525 | 5555 | 2997 | 1 | 1 | 254520 | 13558 | 78720 | 0 | .006 | .014 |
| knightstour | 2593 | 314564 | 907502 | 319299 | 1161 | 10 | 508225 | 307092 | 206729 | 82112 | 464.902 | 14.124 |
| treejoin | 288666 | 519943 | 297570 | 5547 | 2 | 1 | 525488 | 7150 | 212653 | 1954 | .356 | .217 |
| nqueens | 283822 | 1423226 | 2133001 | 584143 | 46 | 9 | 1819579 | 501093 | 521826 | 39465 | 70.314 | 24.811 |
| lambda | 205 | 556 | 2072 | 90345 | 23 | 8 | 966 | 721 | 303 | 95 | .093 | 2.49 |

(b)

**Fig. 10.** Experimental results comparing RGC and LGC. Table (a) gives data related to liveness analysis, and (b) gives garbage collection data.

There are a couple of issues of concern. The garbage collection time is larger in the case of LGC for some programs. The reason is that the cost of consulting the liveness DFA may outweigh the combined benefits of fewer garbage collections and fewer markings per garbage collection. The other issue is illustrated by the program lambda. As can be seen from the table in Figure 10, the number of touched cells[10] in this example is much higher for LGC. This increase is due to excessive sharing among heap nodes in this program. Note that a node re-visited because of sharing is not explored any further during RGC. However, this curtailment cannot happen in LGC because of the possibility that the node, re-visited in a different liveness state, may mark a set of cells different from the earlier visit.

## 6 Collecting more garbage can never slow things down

Since garbage collection is effectively asynchronous to the allocator thread, one might worry as to how *robust* our measurements are. For example, while LGC would, in general, collect more garbage than RGC in the same heap state, might LGC do a larger number of collections for some programs? We prove below that this cannot happen. This result applies to classical mark-and-sweep and copying garbage collectors and, we believe, also to generational collectors.

**Lemma 1.** *For the same mutator, a liveness-based collector can never do more garbage collections than a reachability-based collector.*

*Proof.* Assume, as before, that time is measured in terms of the number of **cons** cells allocated. Now run two copies of the mutator, one with RGC and one with LGC, in

---

[10] These are the cells visited during the marking phase, often more than once due to sharing.

parallel. Memory allocations by **cons** happen simultaneously, but the times of garbage collections diverge.

To prove the lemma, it is enough to show the truth of the following statement: After every LGC invocation, the count of LGC invocations is no greater than RGC invocations. The base case holds since the first invocations of both GCs happen at the same time. Assume the statement to be true after $n$ invocations of LGC. Since LGC copies a subset of reachable cells, its heap would contain no more cells than RGC heap at the end of the $n$th invocation. Thus either RGC is invoked next before LGC, or LGC and RGC are both invoked next at the same time. In either case, the statement holds after $n+1$ invocations of LGC.

## 7  Related Work

Previous attempts to increase the space efficiency of functional programs by additional reclamation of memory fall in two broad categories. In the first, the program itself is instrumented to manage reclamation and reallocation without the aid of the garbage collector. Such attempts include: sharing analysis based reallocation [12], deforestation techniques [13–15], methods based on linear logic [16] and region analysis [17]. Closer to our approach, there are methods that enable the garbage collector to collect more garbage [18, 6] by explicitly nullifying pointers that are not live. However, the nullification, done at compile time, requires sharing (alias) analysis. Our method, in contrast, does not require sharing because of the availability of the heap itself at run time. To the best of our knowledge, this is the first attempt at liveness-based marking of the heap during garbage collection.

## 8  Conclusions

We have defined a notion of liveness on structured data; this generalises classical liveness and strong liveness. We started with a general fully context-sensitive analysis which we proved correct with respect to a minefield semantics (this models the effect of garbage collection between every evaluation step).

To avoid scalability issues (and to avoid performing part of the liveness computation at run time) we defined an 0-CFA version of this liveness analysis in which demands for function $f$ at all calling contexts are conflated into a single demand $\sigma_f$. This enabled us to treat the liveness equations symbolically obtaining context-free grammars for liveness at each GC point (calls to user functions and to **cons**). These were then converted to DFAs for run-time consultation by the garbage collector. Experiments confirm the precision of the analysis.

To obtain performance figures we compared a reachability-based garbage collector with a liveness-based collector. This showed a decrease in the number of GCs, more garbage collected per invocation. A significant benefit of LGC is that programs can run in smaller memory when compared to RGC. This is potentially useful in situations where memory is limited—as with embedded systems. For a majority of programs, the garbage collection times were reduced.

One issue we highlighted was that while fewer nodes were marked (and hence more garbage collected), sometimes **cons** cells could be visited and traversed multiple times with different sets of liveness paths to explore; this risks infinite looping if extended to languages with cyclic data structures. Avenues of further work include static analysis to avoid revisiting cells known to have been visited. One possibility is to record a representation of the liveness paths already visited from each **cons** cell; the classical mark bit indicates that all paths have been visited from the cell.

# References

1. Röjemo, N., Runciman, C.: Lag, drag, void and use—heap profiling and space-efficient compilation revisited. In: ICFP. (1996)
2. Karkare, A., Sanyal, A., Khedker, U.: Effectiveness of garbage collection in MIT/GNU Scheme. http://arxiv.org/abs/cs/0611093 (2006)
3. Shaham, R., Kolodner, E.K., Sagiv, M.: Estimating the impact of heap liveness information on space consumption in Java. In: ISMM. (2002)
4. Chatterjee, R., Ryder, B.G., Landi, W.A.: Relevant context inference. In: POPL. (1999)
5. Cherem, S., Rugina, R.: A practical escape and effect analysis for building lightweight method summaries. In: CC. (2007)
6. Lee, O., Yang, H., Yi, K.: Static insertion of safe and effective memory reuse commands into ML-like programs. Science of Computer Programming (2005)
7. Khedker, U.P., Sanyal, A., Karkare, A.: Heap reference analysis using access graphs. TOPLAS (2007)
8. Karkare, A., Khedker, U., Sanyal, A.: Liveness of heap data for functional programs. In: Heap Analysis and Verification Workshop (HAV). (2007) http://research.microsoft.com/~jjb/papers/HAV_proceedings.pdf.
9. Chakravarty, M.M.T., Keller, G., Zadarnowski, P.: A functional perspective on SSA optimisation algorithms. In: COCV. (2003)
10. Mohri, M., Nederhof, M.J.: Regular approximation of context-free grammars through transformation. In Junqua, J.C., van Noord, G., eds.: Robustness in Language and Speech Technology. Kluwer Academic Publishers (2000) 251–261
11. Hopcroft, J.E., Ullman, J.D.: Introduction To Automata Theory, Languages, And Computation. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA (1990)
12. Jones, S.B., Metayer, D.L.: Compile-time garbage collection by sharing analysis. In: FPCA. (1989)
13. Wadler, P.: Deforestation: transforming programs to eliminate trees. In: ESOP. (1988)
14. Gill, A., Launchbury, J., Jones, S.L.P.: A short cut to deforestation. In: FPCA. (1993)
15. Chitil, O.: Type inference builds a short cut to deforestation. In: ICFP. (1999)
16. Hofmann, M.: A type system for bounded space and functional in-place update. In: ESOP. (2000)
17. Tofte, M., Birkedal, L.: A region inference algorithm. TOPLAS (1998)
18. Inoue, K., Seki, H., Yagi, H.: Analysis of functional programs to detect run-time garbage cells. TOPLAS (1988)