

Memory Management for Prolog with Tabling

Bart Demoen Konstantinos Sagonas

Report CW261, April 1998



Katholieke Universiteit Leuven
Department of Computer Science
Celestijnenlaan 200A – B-3001 Heverlee (Belgium)

Memory Management for Prolog with Tabling

Bart Demoen Konstantinos Sagonas

Report CW261, April 1998

Department of Computer Science, K.U.Leuven

Abstract

Tabling can be implemented in a Prolog system by means of the SLG-WAM: consumers suspend by freezing the execution stacks. XSB is an implementation that does so. The memory model is quite complex and attempts to understand the notion of usefulness of data in XSB well enough to build a precise garbage collector have failed in the past. CAT is a recent alternative to SLG-WAM: it suspends consumers by copying parts of the execution stacks. The memory model is simpler and the design of a more precise garbage collector became feasible. CAT also provided the necessary insights in the usefulness of data in the context of the SLG-WAM. This paper describes the memory management of tabled logic programming systems, whether based on the SLG-WAM or on CAT. Since CAT can perform arbitrarily worse than SLG-WAM space-wise, also a minor garbage collection on creation of the CAT areas is described and its effectiveness is discussed.

Memory Management for Prolog with Tabling

Bart Demoen Konstantinos Sagonas

Department of Computer Science

Katholieke Universiteit Leuven

B-3001 Heverlee, Belgium

{bmd,kostis}@cs.kuleuven.ac.be

Abstract

Tabling can be implemented in a Prolog system by means of the SLG-WAM: consumers suspend by freezing the execution stacks. XSB is an implementation that does so. The memory model is quite complex and attempts to understand the notion of usefulness of data in XSB well enough to build a precise garbage collector have failed in the past. CAT is a recent alternative to SLG-WAM: it suspends consumers by copying parts of the execution stacks. The memory model is simpler and the design of a more precise garbage collector became feasible. CAT also provided the necessary insights in the usefulness of data in the context of the SLG-WAM. This paper describes the memory management of tabled logic programming systems, whether based on the SLG-WAM or on CAT. Since CAT can perform arbitrarily worse than SLG-WAM space-wise, also a minor garbage collection on creation of the CAT areas is described and its effectiveness is discussed.

1 Introduction

Incorporation of tabling (also known as memoization) in the execution model of non-deterministic programming languages such as Prolog leads to a more powerful and flexible paradigm: tabled logic programming. More specifically, through the use of tabling repeated subcomputations are avoided and more programs terminate; thus the resulting execution model allows more specifications to be executable. As a result, practical systems incorporating tabling such as XSB, have been proven useful to a wide range of application areas such as parsing, databases, program analysis, and recently verification through model checking. At the same time, tabling introduces some extra complications in the standard implementation platform for Prolog, the WAM [13]. Most of the complications are attributed to the inherent asynchrony of answer generation and consumption, or in other words to the more flexible control that tabling requires: Control, i.e. the need to *suspend* and *resume* computations, is a main issue in a tabling implementation, because some subgoals, called *generators*, generate answers that go into the tables, while other subgoals, called *consumers*, consume answers from the tables; as soon as a generator depends on a consumer, the generator and the consumer must be able to work in a coroutining fashion, something that is not readily possible in the WAM which reclaims space at backtracking. The need to suspend computations means that execution environments of these computations must be preserved. The SLG-WAM [12], the abstract machine of XSB, preserves consumer states by *freezing* them, i.e. by not reclaiming space on backtracking as is done in WAM and allocating new space above the freeze points.

Freezing is not the only possible way to implement the control of tabled execution. In [8], we introduced the ‘Copying Approach to Tabling’ abbreviated CAT. In short, CAT preserves and

reinstalls the execution environments of suspended computations through *copying*. CAT resembles other copy-based techniques, notably MUSE [2] and the implementation of backtracking in C described in [9]. Introducing backtracking into a language (like Prolog) requires an extra complication in the stacks: indeed, the WAM environment ‘stack’, is really a spaghetti stack on its own and [9] maps this to a regular (C) stack by copying. Similarly, OR-parallelism on top of Prolog, requires an even more complex stack structure than Prolog and MUSE reduces this back to the Prolog stacks by using copying to separate the areas of different workers. Also the implementation of tabling in SLG-WAM requires an extra layer of spaghetti in the stacks of the underlying Prolog engine and CAT is an attempt to reduce the complexity of the SLG-WAM stacks to the relatively well understood stacks of the WAM: indeed, the stack structure of SLG-WAM is difficult to understand and in particular the implications on memory management are very hard to grasp. The simplification that CAT offers also pays off here: reasoning about reachability and usefulness is much easier than in the freezing implementation of SLG-WAM.

In Section 2 we state the concepts and notations we use in the sequel of the paper and review the two approaches to implementing tabling: the SLG-WAM and the CAT. Section 3 explains cooperation between CAT and heap garbage collection, with a detailed discussion on early reset. We pay special attention to how the understanding of garbage collection for CAT, can be used for designing a garbage collector for the SLG-WAM implementation model. Section 4 shows further space improvements during copying for the heap and the local stack. We end with a performance evaluation of this selective copying and a conclusion.

We assume the reader to be familiar with Prolog and the WAM (see [1] for a general introduction) and to some extent with tabled execution of logic programs. Also knowledge of garbage collection techniques for Prolog will help; see [5, 7] for specific instances of heap garbage collection for Prolog and [4, 3] for a more general introduction.

2 Concepts and Terminology of Tabling and its Implementation

Tabled resolution records calls to designated *tabled* predicates and their answers in a persistent global data structure called a *table*. All other predicates are executed as in Prolog. When a call, s , to a tabled predicate is encountered, a check is made to determine whether s already exists in the table or is new to the evaluation. If s is new, s is called a *generator*, it will be entered in the table and execution will continue as in Prolog by resolving s against the program clauses of the predicate. Upon successful return from a program clause, $s\theta$ where θ denotes the accumulated bindings for variables in s constitutes an *answer* for s which is also inserted in the table if it is new to the evaluation. If, on the other hand, s already appears in the table, it will not use the program clauses as done in usual Prolog execution, but will resolve against the answers that currently exist or will be derived for the corresponding generator. In this case, s is called a *consumer* subgoal.

Tabling cannot be implemented using the purely depth-first search control of the WAM as the generation and consumption of answers are asynchronous and quite often mutually dependent processes. This calls for a mechanism to *suspend* and *resume* computations of consumers: for the purposes of this paper it is sufficient to assume that a consumer is suspended when it has consumed all answers from the table, and is possibly resumed upon determining that it has unresolved answers to consume. In implementation terms, the need for a suspend/resume mechanism means that an abstract machine for tabled execution has to alter the control — and thus the memory management — of the WAM to either preserve or be able to reconstruct execution environments of consumers until these do not need to be resumed anymore. In definite programs, this condition happens when generators have exhausted all their program clauses and all consumers have con-

sumed all answers that generators have produced and inserted in the tables: we then say that these subgoals are *complete*.¹

As Prolog execution is supported as well, the starting point for an implementation of tabling is the WAM, but we assume a WAM with four stacks: trail, choice points, local stack (the environments) and the heap. Their top of stacks are pointed by WAM registers denoted respectively as **TR**, **B**, **A** and **H** when we need them. **E** points to the most recent environment (the fact that we need both **E** and **A** in WAM is a bit unfortunate for the description). We denote by **B(E)** the environment field in the choice point **B** and likewise for the other WAM registers. In our description as well as in the figures, the stacks will grow downwards, i.e. higher in the stack means older, lower means younger.

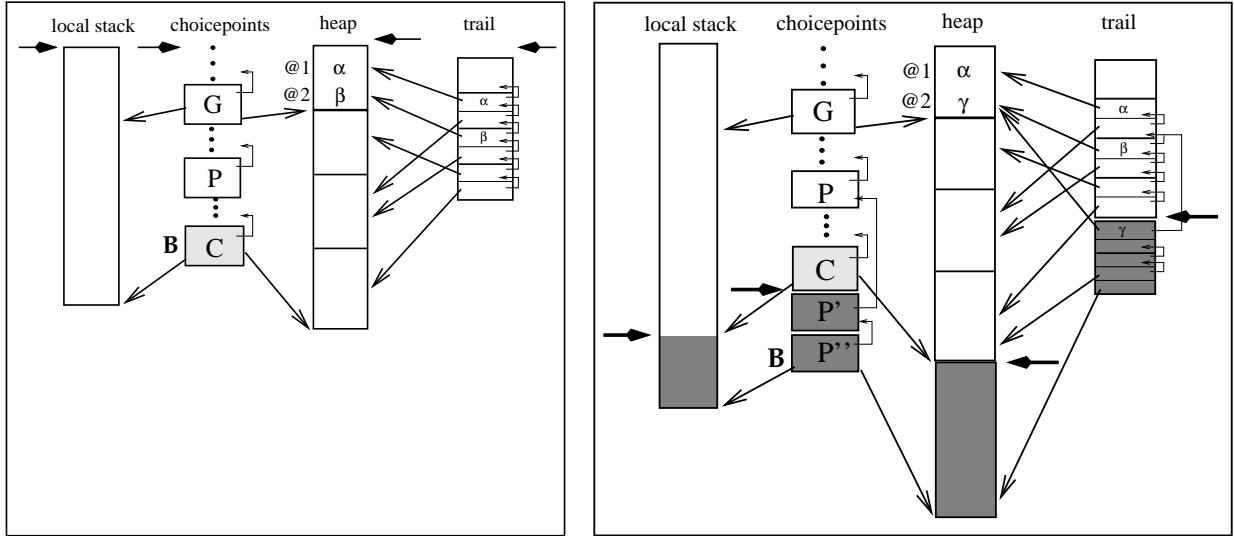
2.1 SLG-WAM: The environment-sharing implementation of tabling

Tabling can be implemented by modifying the WAM as is done by the SLG-WAM [12]. In this approach, execution environments of suspended consumers are preserved by *freezing* the WAM stacks, i.e. by not allowing backtracking to reclaim space in the stacks as is done in the WAM. In implementation terms, this means that the SLG-WAM adds an extra set of *freeze registers* to the WAM, one for each stack and allocation of new information occurs below the frozen part of the stack. Suspension of a consumer is performed in the SLG-WAM by creating a consumer choice point to backtrack through the answers in the table, setting the freeze registers to point to the current top of the stacks, and upon exhausting all answers fail back to the previous choice point without reclaiming any space. Frozen space is reclaimed *only* upon determining completion. Note that a side-effect of having frozen segments in the stacks is that the stacks actually represent trees:² for example, contrary to the WAM, choice points on the same branch of the computation may not be contiguous and the previous choice point may be arbitrarily higher in the stack.

Memory areas of the SLG-WAM and their relationships are depicted in Figure 1. Initially all freeze registers point to the beginning of the stacks; they are shown by arrows next to each stack. After executing some Prolog code and execution encounters a generator G , it creates a generator choice point for it; the execution continues and some more choice points are created and eventually a consumer C is encountered. The SLG-WAM stacks at this point are shown in Figure 1(a). The heap and the trail are shown segmented by choice points; the same segmentation is not shown for the local stack as it is a spaghetti stack. From the trail, some pointers point to cells older than the generator G : these cells have addresses @1 and @2 in the picture, and the values of the cells are α and β . One can see that a trail entry in this picture consists of two pointers and a value, while in WAM, a trail entry is just one pointer: this new trail structure will be explained later on. On encountering C the stacks are frozen by setting the freeze registers to point to the current top of the stack (cf. Figure 1(b)). After possibly returning answers to C , the execution fails out of B_C , and suppose that the youngest choice point with unexplored alternatives is B_P . As shown in Figure 1(b), allocation of new information (shown in a darker shade) takes place below the freeze registers and no memory above the freeze registers is reclaimed. Notice the conceptual tree form of e.g. the choice point stack as shown by previous pointers from choice points. Finally, note that by continuing with some other part of the computation, some cells may change value: e.g. cell @2 from β to γ .

¹Note that, contrary to what happens in the memoization of functional languages, non-determinism introduces an extra complication in the picture as it is usually not known in advance how many answers a given subgoal will produce and until when its consumers might need to be resumed and their execution environments be preserved.

²This type of WAM stacks is also called a *cactus stack* by some authors.



(a) Stacks on encountering a consumer

(b) Continuing forward execution after freezing

Figure 1: Memory areas while executing under an SLG-WAM-based implementation.

As expected, to resume a suspended computation of a consumer, the SLG-WAM needs to have a mechanism to reconstitute its execution environment. Besides resetting the WAM registers (e.g. setting B to point to the consumer choice point), the variable bindings at the time of suspension have to be restored. This can be done using what is known as a *forward trail* [14, 12]. An entry in the forward trail consists of a reference cell, a value cell, and a pointer to the previous trail entry. These entries are shown in Figure 1: entries for @1 and @2 record the values α, β , and γ and because of the previous pointers the trail is also tree-structured. Given this trail, restoring the execution environment EE from a current execution environment EE_c , is a matter of untrailing from EE_c to a common ancestor of EE_c and EE , and then using values in the forward trail to reconstitute the environment of EE .

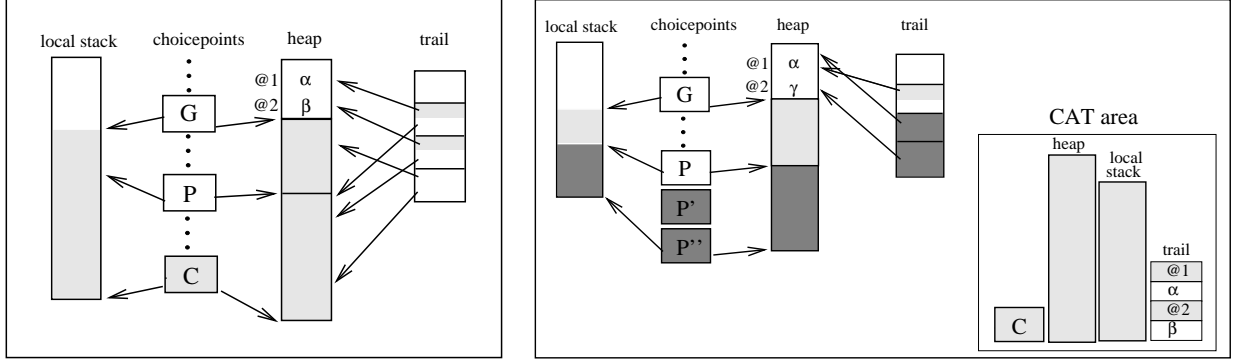
2.2 CAT: The environment-copying implementation of tabling

Instead of maintaining execution environments of suspended consumers through freezing the stacks and using an extended trail to reconstruct them, one could also preserve environments of consumers by copying all the relevant information about them in a separate memory area, let execution proceed as in the WAM, and reinstall these copies whenever the corresponding consumers need to be resumed. This is the main idea behind CAT: the ‘Copying Approach to Tabling’. An advantage of this approach is that, contrary to the SLG-WAM, Prolog execution happens as in the WAM: there is no need for a forward trail nor freeze registers and the stacks do not have a tree form. CAT selectively copies information needed to reinstall suspended environments in *CAT areas* as briefly explained below.³

The CAT area has four memory areas (containing information from each of the four WAM stacks). Figure 2 shows these memory areas in a CAT-based implementation. When execution

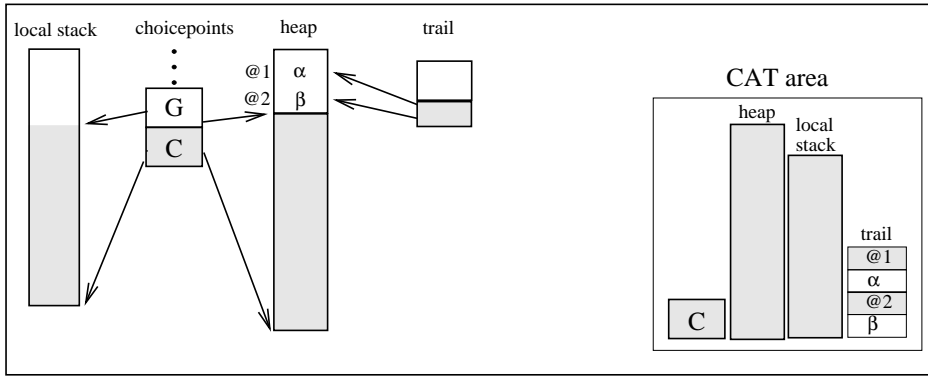
³Actually, it does so in a more incremental way, but as this is not relevant for this paper we refer those interested to [8] for more details.

encounters a consumer, a choice point C is created for it. Let the youngest generator choice point in the stack be G (the dots show possible Prolog choice points that appear in between). A CAT copy is about to be made; the situation is depicted in Figure 2(a). The shaded parts in Figure 2(a) show exactly what CAT copies. From the heap, the CAT copies the part between the current top H



(a) Stacks just before making the CAT copy

(b) Stacks and CAT area after making the CAT copy and continuing computation



(c) Memory areas upon reinstalling the CAT area for consumer C

Figure 2: Memory areas while executing under a CAT-based implementation.

and $B_G(H)$. We refer to the region of the CAT area that contains the copy of the heap as the *CAT heap* and likewise for the other areas. The part of the local stack that needs copying is between A and $B_G(E)$. One could think that from the choice point stack, CAT needs to copy from B till B_G , but [8] argues this is wrong: instead, it is correct to copy only the consumer choice point. Copying the trail is more complicated: as we do not save the part of the heap that is older than B_G and since this part can contain values that were put there during execution more recent than B_G , we need to save together with the trailed addresses also the values these trailed addresses now contain; we do not need a similar double trail for the part of the heap that is more recent than B_G , because we copy that part completely. The copied information is saved in a CAT area which is separate from the stacks (cf. Figure 2(b)) and execution continues as in the WAM by failing out of the consumer choice point. Contrary to what happens in the SLG-WAM, backtracking in CAT

now reclaims space. Figure 2(b) shows a possible situation where backtracking has taken place up to a Prolog choice point P which lied between G and C in Figure 2(a), and then an alternative path of the computation was tried (shown in a darker shade). Note that this new computation has resulted in the stacks having different contents than what is saved in the CAT areas (although as shown some parts are still intact). Also note that if P'' is a consumer choice point, another CAT area will be created at this point. Eventually, through backtracking execution will fall back to G and after G exhausts all resolution with program clauses, the evaluation reinstalls consumers with unresolved answers that have copied up to the generator G .⁴ The resulting stacks are shown in Figure 2(c): through copying, the consumer has just been reinstalled below \mathbf{B}_G and can start consuming its answers from the table. Note that after reinstalling the consumer, the choice point and trail stack are in general smaller than at the time of saving the CAT area. The CAT area itself remains in existence until it can be determined that the associated generator is complete.

At any moment, the computation that is going on and its stacks, are referred to as *current*. Space for computations that have been *frozen* (SLG-WAM terminology) or *copied* (CAT terminology) are *inactive*. Note that there is conceptually a one-to-one mapping between frozen consumers and CAT areas: this can also be seen by comparing Figure 1(b) and Figure 2(b).

3 Heap Garbage Collection

In this section we will concentrate on aspects of heap garbage collection for both CAT and SLG-WAM. The common situation is as follows: the current computation needs to be collected and there are one or more consumers suspended. We will describe the issues for a sliding collector (see [10]) but most of the issues translate easily to a compacting collector as in [5, 7]. The issues are:

- to find the reachable data
- to move it appropriately while adapting all pointers to it

In both sliding or compacting, a marking phase is performed (see [5] for why this is necessary for compacting). The root set in WAM consists of the argument registers, the saved argument registers in the choice points, the local variables in the environments reachable from the current top environment or from the choice points. Marking is done in a WAM heap garbage collector by considering this root set from newer to older (mainly because that is the way frames are linked to each other) and in the set of backtracking states⁵ from current to more in the future: this latter order corresponds to treating younger choice points (and their continuation) before older choice points. It ensures the optimal possibility for doing virtual backtracking or early reset (see [4]).

3.1 Heap garbage collection for CAT

In CAT, it is possible to collect just the current computation. A good reason for doing so, is that most of the current computation and the saved consumer are usually unrelated; so the collection of the consumer state might never be needed. A drawback is that potentially, the consumer could after its reinstallation cause very quickly a new garbage collection that might have been avoided by collecting also the CAT areas.

⁴This may depend on the scheduling strategy that is used; the strategy which governs when answers are returned to consumers — see [8].

⁵A backtracking state is the continuation of a computation which results from backtracking, possibly more than once.

Consumers share a part of the current computation: the part that is older than the generator up to which they have copied information. The CAT area of a consumer can have the following pointers to the current computation:

1. from the CAT heap to the shared heap part
2. from the CAT trail to the shared heap part
3. from the CAT local stack to the shared local stack and shared heap part

These pointers are to be followed for marking and take part in the usual relocation during the second phase of the sliding algorithm.

Moreover, the CAT area contains in general also pointers to the heap and local stack as they will be after reinstalling the consumer: conceptually, such a pointer references the CAT heap (or local stack) but in reality it points to the heap or local stack. Such a pointer we name a *CAT internal* pointer. During the marking phase, it is followed as if pointing to the CAT area itself. Since the exact place where the consumer will be reinstated will change because of the garbage collection, the CAT internal heap pointers need to be relocated by the same amount by which the heap pointer from its generator has changed. Figure 3 illustrates the situation.

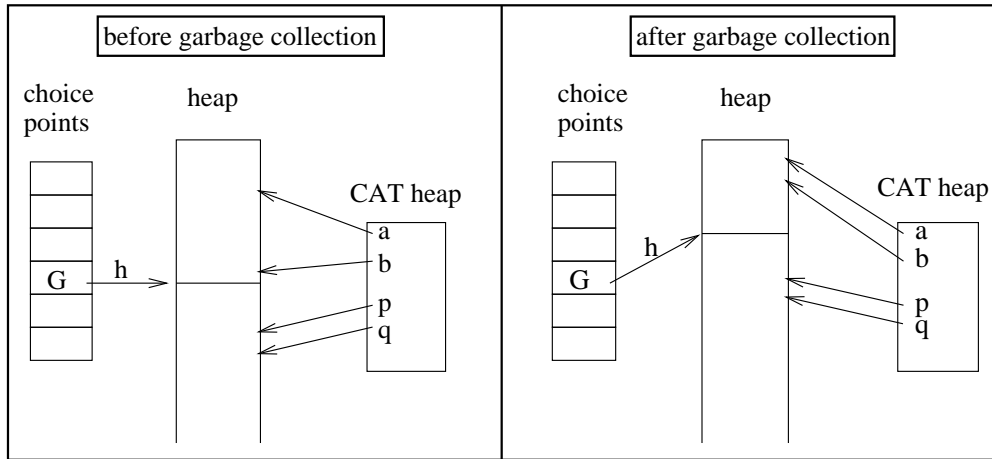


Figure 3: Some references from the CAT heap before and after garbage collection

G denotes the generator choice point up to which a consumer has been saved: from this consumer, only the CAT heap is shown. The top of heap at of the moment of creation of G is h . The CAT heap contains references a and b to the old part of the current heap, while the references p and q are CAT internal heap pointers: they do point to the heap, but the data they point is currently in the CAT heap and will be in the heap only after the consumer will be reinstated. Consequently, the references a and b are relocated according to the sliding process taking place in the old part of the heap, while the pointers p and q are shifted by the same amount as the h pointer.

Since there may exist several CAT areas at the same time (for different consumers), it is worthwhile noting that no references from one consumer to another consumer can exist: references are always to the heap and local stack in which the current computation evolves.

The above can be adapted to a copying collector: in a schema that requires marking before copying (see [5, 7]), the marking phase remains the same. Then copying [6] takes place of the reachable current heap, after which the forwarding pointers contain the information to relocate the

pointers from the CAT area to the current heap. The relocation of the CAT internal heap pointers remains as explained above.

In this schema, the CAT areas are not collected themselves; one can argue that a consumer could provoke very quickly after its reinstallation a new collection. This is obviously true. So the schema has to be made incremental, in the sense that if the part of the computation older than the generator has been collected already, it should not be collected after its consumer is reinstalled and causes overflow; only the part that corresponds to the installed consumer should be collected. For more details on incremental garbage collection for Prolog, see for instance [5]. In this way, the collection of the consumer is postponed until it is really needed and it might never be !

3.2 Heap garbage collection for SLG-WAM

In the SLG-WAM, it seems difficult to collect the current computation without collecting the frozen consumers at the same time: their data are intertwined on the stacks. So, the proposal here is to use as root set the current computation (and its backtracking states) as well as the frozen consumers. Each consumer is characterized by its choice point together with a part of the local stack and trail that is newer than the generator of its answers. In the past, the only way we could imagine marking a consumer, was by melting it first (i.e. reinstall its bindings by using the forward trail), and then mark it as if it were a current computation. However, a previous marking (during the same garbage collection, either a marking of another consumer or the current computation), will have set some mark bits, and it was not clear how they had to be taken into account.

In the setting of CAT, it was easier to come to the conclusion that a consumer can be marked without being reinstalled, that marking a consumer does not need to consider the part of the stack that is older than the generator up to which the copy was made and that even the choice points between the consumer and the generator need not be considered: indeed, in CAT, at the time of reinstalling the consumer, these intermediate choice points have already been removed by backtracking.

In the setting of SLG-WAM, two consumers can also share a part of this trail and local stack, so we expect that it pays off to avoid marking from the same roots (in this shared part) more than once. It is also common in a plain Prolog garbage collector to avoid marking from the same environment cell, which can be reachable from different choice points.

3.3 Virtual backtracking and the order of marking

Virtual backtracking, or early reset, in the context of tabled logic programming remained long a mystery to us. For a good account on this issue in Prolog, see [4], or [3] and [11]. The idea of virtual backtracking is that a trailed heap entry which is not reachable for the forward continuation of the current computation (but might be for alternative branches of the current computation, i.e. on backtracking) can be set to unbound during garbage collection and the trail entry itself can be discarded as well. The situation is recognized during marking, and it is essential that the continuation is marked before the future alternatives.

The execution of `?- main.` against the following artificial piece of Prolog code shows a potential for early reset:

```

main :- h(f(X)).
h(Z) :- Z = f([1,2,3]),
        last_use_of(Z),
        gc, % do garbage collection
        more_goals.
h(Z) :- ...

```

The binding of the variable X (in `main`) to the list `[1,2,3]` is on the trail, but since Z is not useful in forward execution after garbage collection, it can be discarded during garbage collection, even though the data structure in Z is reachable from the choice point.

For CAT, it is clear that during garbage collection, early reset can be performed for the current computation: the suspended computations, i.e. the consumers, have their own computation state saved. If this applies to CAT, it has to apply to SLG-WAM as well, as the matter of usefulness of data is related to the abstract suspension mechanism rather than to the actual realization. However, we also want to answer the question whether during marking of the consumers early reset is allowed and/or possible.

A CAT trail entry of a consumer, say C , can contain a reference to a heap entry that is not reachable from the continuation of C : there is no harm in removing this CAT trail entry. However, setting the corresponding heap entry to undefined in a non-discriminating way might be wrong, since the current computation might need the cell with its current value. In that case and if marking of the current computation was done first, the mark bit of this trailed heap entry was set and the marking phase of the consumer can decide not to make the cell undefined.

On the other hand, the mark bit does not reflect *who* (the consumer or the current computation) references the cell, so that it is possible that by marking the generator first, the consumer loses an opportunity for early reset — or at least CAT trail compaction. This means that the order in which consumers and the current computation are marked, can be crucial. For the sake of focusing the ideas, assume that at the moment of the garbage collection, there is just the current computation and one consumer. Obviously the generator choice point which can schedule the consumer (some time in the future) is in the current computation at that moment. Two reasonable orders of marking are:

1. first mark the current computation completely and then the consumer
2. interleave the marking of the current computation and the consumer in the following way: first mark the current computation up to the generator, then mark the consumer up to the generator; finally, mark the rest of the current computation. Since the consumer shares this latter part with the current computation, there is no need to go back to marking the rest of the consumer.

The second order of marking is more precise, as it is possible that by marking the current computation completely, some opportunity for early reset is lost for some consumers. For the first order, the only way to remedy this, would be to have a mark bit for each consumer, which is of course impractical since the number of consumers can be large (see Table 1 which shows the number of consumers in some benchmarks). On the other hand, method 2, has the disadvantage that when many consumers exist, they can have different generator choice points that can schedule them, meaning that the place up to which they are copied from the current computation, might be different. Since it has never been proven in the context of Prolog that early reset is really worth its while, we think that method 1 is to be preferred.

Because of mark bits set earlier during the marking of the current computation, early reset is different for the current computation and a consumer. We describe more in detail. Figure 4 shows a pointer tr which points to a CAT trail entry: the value cell of the trail entry contains a value v_1 ; the reference cell of the trail entry is a and points to a heap (or local stack) location which contains a value v_2 : in general $v_1 \neq v_2$.

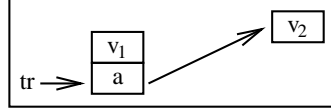


Figure 4: A CAT trail entry

Below is some pseudo-code for treating one trail entry during marking: the code is different for the current computation and for the consumer.

<i>code for current computation</i>	<i>code for consumer</i>
<pre> if (! marked(a)) { *a = UNDEF; remove_trail_entry(tr); } </pre>	<pre> if (! marked(a)) { *a = UNDEF; remove_trail_entry(tr); } else /* marking is necessary in general */ mark_object(v1); </pre>

Function `mark_object()` marks a value representing a Prolog object.

From this piece of code, one can deduce two ways in which early reset for a consumer is suboptimal; if a heap (or local stack) cell reachable from the trail is marked already, this can have two reasons:

1. the cell was reachable from the consumer itself
2. the cell was reachable from another consumer or the current computation

These two cases cannot be distinguished if there is only one mark bit. As a result, in either case, not only will early reset be prohibited, but also the data structure v_1 needs to be marked, which would be unnecessary if the cell had not been reachable from the consumer one is currently marking. It follows that only the current computation can do optimal early reset; the consumers approximate it.

One conclusion is that early reset during marking of the consumers can reset the heap entries and remove CAT trail entries.

The above analysis was made for CAT. It can be applied almost straightforwardly to SLG-WAM: after marking the current computation — performing early reset in the course of doing so — one can mark (and early reset) the state of the consumers, which is captured by the cells reachable from the consumer choice point and the part of the trail starting at the consumer up to where it meets the trail of the generator that can schedule the consumer. The part of the trail older than the generator will have been marked already, since the consumers that are scheduled by this generator share the trail with the generator. Figure 4 can also serve as a picture of an SLG-WAM forward trail entry (not including the back pointer): for the current computation, it is always true that $v_1 = v_2$. Accordingly, the code above applies to SLG-WAM as well.

4 Selective CAT: A minor garbage collection on saving consumers

One might be struck by the fact that while complete contiguous parts of the heap and the local stack are saved in the CAT area, the trail is copied selectively and from the choice point stack, only the top choice point is saved. [8] argues about the choice point stack that this is almost a necessity. For the trail, it was a conscious choice because it is clear that the CAT trail needs double entries and so one needs to inspect the trail entries while copying them: it was then natural to be more selective. We will show how heap and local stack can also be saved more selectively in CAT.

4.1 Selectively saving the heap

In general, the CAT heap contains entries that will not be reachable after the particular consumer has been reinstalled: indeed, the choice points between the consumer and the generator will no longer exist and this reduces the root set. So, the idea here is to save in the CAT heap only the part of the current heap between consumer and generator that is reachable from the current root set in *forward* execution; this takes into account that the intermediate choice points between the consumer choice point and the generator are not saved. The easiest way to understand this is: assume one performs an incremental garbage collection of the heap, where the part older than the generator is not collected and while not using as root set the environments that are *only* reachable from the intermediate choice points. The result gives exactly what needs to be saved in the CAT heap. To be more precise, the marking is based on the following pseudocode whose underlying principle can be found in [5]:

```
/* mark current computation — B points to the consumer choice point */
e = B(E); cp = B(CP); /* CP is the continuation pointer */
mark_arguments_from(B);
elimit = BG(E); /* G is the generator up to which the CAT copy is made */
while (newer(e,elimit))
  { mark_environment(e,cp);
    cp = e(CP);
    e = e(E);
  }
/* note that no other choice points are considered for finding the useful data */
/* mark all new objects pointed to from older regions */
tr = TR;
trlimit = BG(TR);
while (newer(tr,trlimit))
  { mark_object(**tr);
    tr = previous(tr);
  }
```

In the code above, the mark functions must (as usual in incremental marking) refrain from marking anything that is not in the part of the heap between the generator and the consumer: the trail acts like an exception list, which contains the references from the old generation to the new generation. For clarity, we have left out early reset from the above code.

While the above code just shows marking for the purpose of understanding which data must be saved, it can be the first phase in a sliding garbage collection as is most common in Prolog systems,

because there are good reasons to preserve the order of segments (and even of individual cells in the heap). However, as after the reinstallation of the consumer, this part of the heap is exactly one segment, a copying collector could be used as well, e.g. based on [5] or [7].

4.2 Selectively saving the local stack

Similarly to the heap, the local stack can be copied more selectively in CAT. To understand this, consider which environments are reachable after the consumer is reinstalled: only the environments in the chain starting from the environment at the moment of saving the consumer: these are exactly the environments visited by the marking code above. The environments that are *only* reachable through the choice points between the consumer and generator, will not be reachable after the consumer is reinstalled (since these choice points themselves have been removed by backtracking), so we do not need to save them in the CAT area. Once this is understood, it is reasonably straightforward to see how to implement it.

4.3 Performance evaluation of selective copying

CAT suffers from the problem that it can perform space-wise arbitrarily worse than SLG-WAM. In CAT, an arbitrarily large part of the heap might have to be copied while the SLG-WAM just freezes the stacks which is a constant time operation. This actually showed up in benchmarks in [8] where, under a particular scheduling strategy, CAT's peak memory usage was sometimes 10 times higher than for SLG-WAM. While the more selective saving of CAT heap and CAT local stack introduced above cannot avoid the worst case, it was worth investigating whether it would help in practice. After all, selective saving of the CAT areas also results in smaller reinstallation, so the cost of a garbage collection at CAT creation time might be well worth it as a CAT area is saved once but might be reinstalled several times. In preparing a full implementation, we have implemented marking, and used it as a means of getting some indication on the effectiveness of the selective copying.

benchmark	CAT	selective CAT	gain	# of consumers
cs_o	23008	10176	12832	38
cs_r	38104	20956	17148	71
disj	42812	22812	20000	77
gabriel	54116	26856	27260	80
kalah	105348	32032	73316	89
peep	3925708	739632	3186076	1721
pg	1547216	210292	1336924	497
read	5427168	2203648	3223520	3180

Table 1: Space requirements of CAT with and without selective copying.

Table 1 shows, for the same benchmark set as that used in [8], the peak memory usage of CAT versus CAT with selective copying (both heap and local stack) in the first two columns. The figures are in bytes. The next column shows the gain in space, and the last column shows the number of consumers which create a CAT area during the running of the benchmark. It therefore also indicates the number of times an incremental garbage collection has to be performed for obtaining the gain. The figures show the potential of selective saving of the CAT areas in reducing the space requirements, but on the other hand, for this set of benchmarks, they are rather

discouraging time-wise: knowing that the read benchmark runs in 0.86 seconds (on a SUN Ultra Sparc 2, 168 MHz), that most consumers never need reinstalling and the others only once, it seems hardly worthwhile to do 3180 (incremental) garbage collections. The situation might of course be completely different in programs where consumers get reinstalled multiple times.

5 Conclusion and future work

The memory model of WAM-based systems that support tabling is quite complex due to the incorporation of a suspend/resume mechanism. This, in turn makes the usefulness logic of tabling systems more complicated than that of Prolog ones. Yet, tabling systems call for even more effective and sophisticated memory management than plain Prolog systems as their space requirements are generally bigger. Our past attempts to devise an effective garbage collector for the XSB system did not enjoy much success as we found it extremely difficult to understand the reachability issues in the double spaghetti stack of the SLG-WAM model; without this understanding the best we could hope for was a conservative marking scheme, about which we had absolutely no clue how imprecise it might be. The recently introduced CAT model offered a new, simpler way to implement tabling and, more importantly for this paper, the means to reason about reachability and usefulness of data in tabled evaluations regardless of the underlying implementation.

In this paper we gave a reasonably complete account of the design decisions and most important implementation aspects of memory management in the CAT model. With this as basis, we also discussed garbage collection for the SLG-WAM model. In addition, we showed how at the moment of CAT area creation, a minor garbage collection can be performed and the space savings that can be expected by doing so. A renewed effort for writing the garbage collector for both CAT and SLG-WAM is under way.

Although our results as presented are specific to tabled logic programming implementations, we believe that the underlying concepts might prove useful, or at least give insights, to other (LP) systems that deal with asynchronous processes or implement some variant of a suspend/resume mechanism.

Acknowledgements

The second author is supported by a K.U. Leuven junior scientist fellowship.

References

- [1] H. Ait-Kaci. *Warren's Abstract Machine: A Tutorial Reconstruction*. The MIT Press, Cambridge, Massachusetts, 1991.
- [2] K. A. M. Ali and R. Karlsson. The Muse approach to OR-parallel Prolog. *International Journal of Parallel Programming*, 19(2):129–162, Apr. 1990.
- [3] K. Appleby, M. Carlsson, S. Haridi, and D. Sahlin. Garbage Collection for Prolog Based on WAM. *Communications of the ACM*, 31(6):719–741, June 1988.
- [4] Y. Bekkers, O. Ridoux, and L. Ungaro. Dynamic Memory Management for Sequential Logic Programming Languages. In Y. Bekkers and J. Cohen, editors, *Proceedings of IWMM 92: International Workshop on Memory Management*, number 637 in LNCS, pages 82–102, St. Malo, France, Sept. 1992. Springer-Verlag.

- [5] J. Beveymyr and T. Lindgren. A Simple and Efficient Copying Garbage Collector for Prolog. In M. Hermenegildo and J. Penjam, editors, *Proceedings of the Sixth International Symposium on Programming Language Implementation and Logic Programming*, number 844 in LNCS, pages 88–101, Madrid, Spain, Sept. 1994. Springer-Verlag.
- [6] C. J. Cheney. A Nonrecursive List Compacting Algorithm. *Communications of the ACM*, 13(11):677–678, Nov. 1970.
- [7] B. Demoen, G. Engels, and P. Tarau. Segment Preserving Copying Garbage Collection for WAM based Prolog. In *Proceedings of the 1996 ACM Symposium on Applied Computing*, Philadelphia, Feb. 1996. ACM Press.
- [8] B. Demoen and K. Sagonas. CAT: the Copying Approach to Tabling. Technical Report CW 260, K.U. Leuven, Apr. 1998. Submitted for publication.
- [9] P.-E. Moreau. Compiling Nondeterministic Computations. Technical Report 98-R-005, CRIN, 1998.
- [10] F. L. Morris. A Time- and Space-Efficient Garbage Compaction Algorithm. *Communications of the ACM*, 21(8):662–665, Aug. 1978.
- [11] E. Pittomvils, M. Bruynooghe, and Y. D. Willems. Towards a Real Time Garbage Collector for Prolog. In *Proceedings of the 1985 Symposium on Logic Programming*, pages 185–198, Boston, Massachusetts, July 1985. IEEE Computer Society Press.
- [12] K. Sagonas and T. Swift. An Abstract Machine for Tabled Execution of Fixed-Order Stratified Logic Programs. *ACM Transactions on Programming Languages and Systems*, 20, 1998. To appear.
- [13] D. H. D. Warren. An Abstract Prolog Instruction Set. Technical Report 309, SRI International, Menlo Park, U.S.A., Oct. 1983.
- [14] D. S. Warren. Efficient Prolog memory management for flexible control strategies. In *Proceedings of the 1984 Symposium on Logic Programming*, pages 198–202, Atlantic City, New Jersey, Feb. 1984. IEEE Computer Science Press.