

SICS/R-89/8902

GAM
An Abstract Machine for GCLA
by
Martin Aronsson

GAM

An Abstract Machine for GCLA

by

Martin Aronsson
Swedish Institute of Computer Science
Box 1263, S-164 28 Kista, Sweden

Tel: +46 8 752 15 00

Abstract

GCLA is a new programming language, which increases expressiveness compared with traditional logic programming languages and functional programming languages. The basis for the language is a generalization of the concept *inductive definitions*, called *partial inductive definitions*. The program defines a logic, which is used to make inferences to prove if a query holds or not. This report first presents a short introduction to these ideas. Then, an abstract machine, called GAM, for GCLA is presented; the instructions as well as an introduction to the compiling schema is given together with some examples. The main idea is to extend the Warren Abstract Machine (WAM), which is an abstract machine for the language Prolog.

Keywords: Logic programming, Functional programming, Inductive definitions,
Abstract machine

Contents:

1. Introduction
 2. GCLA the Language
 - 2.1 Syntax
 - 2.2 Operational Semantics
 - 2.3 Procedural Semantics
 3. An Example Deduction
 4. Some Programming Examples
 5. The Ideas Behind the Abstract Machine (GAM)
 6. An Introductory Example of Compiled Code
 7. Data Areas and Registers
 8. The Clause Orders
 9. Declarations
 10. The Procedural Instructions
 - 10.1 Calling Instructions
 - 10.2 Premise Instructions
 - 10.3 Indexing Instructions
 - 10.4 Choicepoint Instructions
 - 10.5 Guard Instructions
 - 10.6 Fail Instruction
 11. The Unification Instructions
 - 11.1 Get Instructions
 - 11.2 Put Instructions
 - 11.3 Unify Instructions
 12. Realization of the Procedural Semantics
 13. Compiling Schema
 14. Some Examples
 15. Conclusions and Future Work
- References

Appendices:

- A Formats of Different Frames:**
- A.1 Choicepoints
 - A.2 Environments
 - A.3 Continuations
 - A.4 Premises
 - A.5 Different trail-formats

1. Introduction

During the past years logic programming has proved to be a good framework for writing programs with high expressive power together with the ability to execute programs efficiently. The key to this combination is to use definite horn clauses in programs, and to use resolution as the computation rule. Generalized horn clauses (GCLA) [3,8] is a generalization of definite horn clauses. It is not a generalization within the traditional logical framework, but rather a more proof-theoretic view of horn clause programming. The program is not looked upon as a set of true facts and implications, but as a definition determining the possible inferences. One could say that in GCLA the program forms the system's mind, meaning that the inference rules are given by the program. In Prolog the inference rule is given a priori and the program is executed within this given framework. The inductive definitions [7] which form the basis of GCLA are in a sense more primitive than ordinary logic. Primitive in the sense that when writing down the definition of any kind of logic one is actually writing down an inductive definition. Several experimental interpreters for GCLA have been implemented in Prolog, and they have shown some of the expressive power of GCLA. But they have also suffered from inefficient treatment of the rules introduced by the extension of the definite hornclauses. Therefore the development of an abstract machine was interesting, to see if a suitable instruction set for GCLA could be found. This would mean increasing the expressive power without losing efficiency compared to for example Prolog. The idea is to extend the Warren Abstract Machine (WAM) [10], which is an abstract machine for Prolog. This report does not describe a garbage collection routine, although one is needed. This will not differ from already existing ones, and in fact, the garbage collecting routine for the heap (which is the only area that should be garbage collected) is the same as for the WAM (for example see [1] and [4]).

When reading this report, be aware that there is a difference between the P's. When we refer to P, we refer to a program, but when we refer to $P\text{-}$ or $\vdash\text{-}P$ we refer to a specific rule, and when we refer to $\vdash\text{-}_P$, we refer to the logic "created" or "generated" by the program P. We often omit the index P in $\vdash\text{-}_P$ although it should actually be written out to be correct.

This work is a part of a larger project, which also contains a programming calculus based on partial inductive definitions [6], theoretical foundations for Generalized Horn clauses as a programming language [8] and the definition of a new programming language [3].

In principle the same report is published as a technical report at SICS [2]. The difference is that an exact definition of the machine by transitions is given as an appendix in the technical report.

2. GCLA the Language

2.1 Syntax

The syntax for a GCLA program is extended "ordinary Prolog" syntax (i.e. DEC-10 syntax). The clauses in the program look like ordinary Prolog clauses, except that they can have a list of unifying guards as a first element in the body, and a new primitive operator, \rightarrow , can occur in the body of a clause. Another new primitive operator, \vdash , is introduced for the goals, which are sequents consisting of assumptions and a conclusion, instead of just a conclusion as in Prolog.

A difference in reading the clauses is that the operator \vdash should be read as "is defined by" rather than "if" as in Prolog in order to grasp the intuition of inductive definitions.

The syntax is the same as the one given in [3].

2.1.1 Variables

- A variable is a string beginning with an uppercase letter.

2.1.2 Functors

- A functor is a string beginning with a lower case letter. Each functor has a certain arity. A constant is a functor of arity 0.

2.1.3 Terms

- Each variable X is a term.
- If $t_1 \dots t_n$ are terms and f is an n -ary functor, then $f(t_1 \dots t_n)$ is a term.

2.1.4 Guards

If t is a term, then:

- $var(t)$ is a guard.
- $nonvar(t)$ is a guard.
- $number(t)$ is guard.
- $atom(t)$ is a guard.
- $atomic(t)$ is a guard.
- $ground(t)$ is a guard.

The only guards that are allowed are these predefined guards.

They are not of the same kind as in the parallel logic languages (GHC and others). The guards in GCLA are used to constrain variables in the head of clauses, and not to rule out other possible clauses.

2.1.5 Conditions and Clauses

- Each term is a condition.
- If C_1, \dots, C_n, C are conditions, then $(C_1, \dots, C_n) \rightarrow C$ is a condition.

- If C is a condition and n a number, then $con(n,C)$ is a condition (con stands for contraction, see 2.2.7).

- If C_1, \dots, C_n are conditions, g_1, \dots, g_m guards and t is a term, then $t :- [g_1, \dots, g_m], C_1, \dots, C_n$ is a clause.

We will refer to t as the *head* of the clause and C_1, \dots, C_n as the *body* of the clause. We will use the word *predicate* for the set of all clauses in the program that have the same principal functor in the head.

t . will be short for the clause $t :- []$.

$t :- C_1, \dots, C_n$. will be short for the clause $t :- [], C_1, \dots, C_n$.

2.1.6 Sequents

A sequent S is an expression of the form $C_1, \dots, C_n \vdash C$ where C_1, \dots, C_n, C are conditions. We call C_1, \dots, C_n *assumptions* and C a *conclusion*.

" C ." will be short for $\vdash C$.

We will also sometimes refer to the assumptions as the *premises*.

2.1.7 Programs

A program P is a finite list of clauses.

Below is a small programming example, which we will see more of later on. The program is used to determine which of the "objects" that are grey.

```
elephant (clyde) .
elephant (fido) .
elephant (P) :-
    albino_elephant (P) .

albino_elephant (karo) .

grey (P) :-
    elephant (P) ,
    (albino_elephant (P) -> false) .
```

and some possible questions

$\vdash grey(P)$	$\Rightarrow P = clyde$
	$\Rightarrow P = fido$
$grey(P) \vdash false$	$\Rightarrow P = karo$

2.2 Formal Semantics

In the following two sections we will use the term *assumption* to refer to a condition to the left of the symbol \vdash , and the term *consequent* to refer to the condition to the right of \vdash . The word *premise* will not be used because it could ambiguously refer to both the top row of a rule and to a condition to the left of the turnstile. For the same reason the word *conclusion* will not be used in these two sections. A *goal* is a sequent consisting of a consequent and a possible empty list of assumptions.

There is also a difference between GCLA and partial inductive definitions in the operational semantics. In GCLA the deductions are linear, concatenating new goals in front of the old not yet solved ones, thereby defining an order in which the goals are solved. In partial inductive definitions the inferences form an and-or-tree where the goals are not ordered in any way.

We are interested in finding substitutions σ such that certain sequents $C_1\sigma, \dots, C_n\sigma \vdash C\sigma$ hold according to a program P . The sequents $C_1\sigma, \dots, C_n\sigma \vdash C\sigma$ should intuitively be read as " $C\sigma$ follows from the assumptions $C_1\sigma, \dots, C_n\sigma$ in the program P (according to P)". To solve these sequents we have a number of rules, going from one state to another. The states are given as triples $\langle L; \sigma, P \rangle$ where L is a list of sequents, σ is an answer substitution and P is a program.

The following abbreviations will be used:

- "nil" to denote the empty list,
- C to denote a condition,
- L,M to denote (possibly empty) lists of conditions,
- C.M means that C is concatenated to M,
- L+M means appending L and M together, and
- $M_1 + t + M_2$ will indicate that the condition t occurs in a list of conditions.

$L \vdash M$ is recursively defined as follows:

$$L \vdash \text{nil} = \text{nil}$$

$$L \vdash C.M = (L \vdash C).(L \vdash M)$$

$D(t).L \vdash C$ is similarly defined by recursion on $D(t)$, i.e. $D(t).L \vdash C =_{\text{def}} \{M+L \vdash C \text{ such that } M \in D(t)\}$.

We have here slightly changed the accumulation of the answer substitution compared with the one given in [3]. This change is made because we are here working from the bottom to the top, going backwards trying to reach the Initial List (IL). Therefore, if IL is reached, the resulting answer substitution should be presented.

It should also be mentioned that there are some other rules which will not be mentioned here. Among these additional rules are the ones for asserting and retracting clauses [3]. These concern the interpreter and not the abstract machine's instruction set, and therefore we have left them out.

2.2.1 Initial List (IL)

$\langle \text{nil} ; \theta, P \rangle$

where the answer substitution could be constructed from θ .

2.2.2 Initial sequent or axiom (I)

$$\frac{\langle \text{Rest}\sigma ; \sigma\theta, P \rangle}{\langle (L_1 + t + L_2 \vdash r) . \text{Rest} ; \theta, P \rangle}$$

where σ is a mgu of t and r

2.2.3 Arrow right ($\vdash \rightarrow$)

$$\frac{\langle (M + L \vdash C) . \text{Rest} ; \theta, P \rangle}{\langle (L \vdash M \rightarrow C) . \text{Rest} ; \theta, P \rangle}$$

2.2.4 Arrow left ($\rightarrow \vdash$)

$$\frac{\langle (L_1 + L_2 \vdash M) + (C.L_1 + L_2 \vdash C') . \text{Rest} ; \theta, P \rangle}{\langle (L_1 + (M \rightarrow C) + L_2 \vdash C') . \text{Rest} ; \theta, P \rangle}$$

2.2.5 Program clauses right ($\vdash P$)

$$\frac{\langle (L\sigma \vdash M\sigma) + \text{Rest}\sigma ; \sigma\theta, P \rangle}{\langle (L \vdash t) . \text{Rest} ; \theta, P \rangle}$$

where $r := [g_1, \dots, g_n]$, M is a clause in P , σ is a mgu of t and r and $g_1\sigma \dots g_n\sigma$ are all true.

2.2.6 Program clauses left ($P \vdash$)

$$\frac{\langle (D(t\sigma).M_1\sigma + M_2\sigma) \vdash r\sigma \rangle + \text{Rest}\sigma ; \sigma\theta, P \rangle}{\langle (M_1 + t + M_2 \vdash r) . \text{Rest} ; \theta, P \rangle}$$

where σ is a t -sufficient substitution computed as:

Let $r_1 \dots r_n$ be a permutation of the heads of the program clauses in P and

$\text{mgu}(r,t)$ be the most general unifier of r and t if there is one, and nil

otherwise

Then define:

$\sigma_0 = \text{nil}$

$\sigma_{m+1} = \text{mgu}(t\sigma, r_{m+1})$

$\sigma = \bigcirc_{i \leq n} \sigma_i$

$D(t\sigma)$ is a list containing all $M\sigma$ such that $r\sigma = t\sigma$ where $r := [g_1 \dots g_n]$, M is a clause in the given program P and all $g_1\sigma \dots g_n\sigma$ are true.

In order for σ to be correct a variable check must also be performed, namely that $D(t\sigma)$

does not introduce any new variables, i.e., that the bodies of the clauses considered do not contain variables that do not occur in the heads of the clauses. These variables are existentially quantified, which gives rise to problems to the left of the symbol \vdash .

The intuitive understanding of this rule is what an assumption should mean according to the given definition, the program. An assumption holds if all objects that define it also hold. If there is no clause in the program defining this assumption, the current goal-sequent succeeds immediately.

2.2.7 Contraction

$$\frac{\langle L \vdash C \rangle + \text{Rest} ; \theta, P \rangle}{\langle L \vdash \text{con}(n, C) \rangle . \text{Rest} ; \theta, P \rangle}$$

$$\frac{\langle L_1 + L_2 \vdash C \rangle . \text{Rest} ; \theta, P \rangle}{\langle L_1 + \text{con}(0, C') + L_2 \vdash C \rangle . \text{Rest} ; \theta, P \rangle}$$

$$\frac{\langle C' + L_1 + \text{con}(n, C') + L_2 \vdash C \rangle . \text{Rest} ; \theta, P \rangle}{\langle L_1 + \text{con}(n+1, C') + L_2 \vdash C \rangle . \text{Rest} ; \theta, P \rangle}$$

2.2.8 Semantics for Guards

- $\text{var}(t)$ is true if t is currently not bound, otherwise false.
- $\text{nonvar}(t)$ is true if t is currently bound to something else than a variable, otherwise false.
- $\text{number}(t)$ is true if t is currently bound to a number, otherwise false.
- $\text{atom}(t)$ is true if t is currently bound to an atom (i.e. a functor of arity 0), otherwise false.
- $\text{atomic}(t)$ is true if t is currently bound to a number or an atom, otherwise false.
- $\text{ground}(t)$ is true if t is currently not bound to a term containing a variable, otherwise false.

2.3 Procedural Semantics

Here we are concerned about in what order different possible solutions should be tried. That is, if at some point in the execution a choicepoint is reached, a choice of which of the different possible ways that should be tried first. There are five points where a choicepoint can arise, namely:

- the choice of a rule,
- the choice of a clause when the rule is $\vdash P$,
- the choice of an assumption when the rule is I ,
- The choice of an assumption if the rule to be tried is $P\vdash$ or $\rightarrow\vdash$

- different possible permutations of the clauses in P if the rule is $P\vdash$, for a certain assumption

This list is the basis for another definition. The procedure in each inference is defined as first choosing among three different possibilities [3]. The possibilities are the Initial sequent (I), some rule to the right of the turnstile, or some rule to the left of the turnstile. These three possibilities are ordered as: first try the Initial sequent, secondly try the rules to the right, and thirdly try the rules to the left. This choice determines the *mode* for the machine, more about this later.

Depending on this choice, other choicepoints can arise:

- If rule I is tried, a suitable assumption must be chosen, which is a choicepoint.
- If some rule to the right of the turnstile is tried, the consequent uniquely determines the choice of rule. If that rule is the rule $\vdash P$ (i.e., the consequent is something else than the arrow), a suitable clause in the program P has to be chosen, which is a choicepoint.
- If some rule to the left of the turnstile should be tried, a suitable assumption has to be chosen first, which is a choicepoint. This choice of assumption uniquely determines the choice of rule in the same manner as in the case above to the right of the turnstile. If the rule is $P\vdash$, a permutation of the clauses in the program P has to be chosen, which also is a choicepoint.

So, one inference could involve zero or more choicepoints depending on the sequent considered.

To form a search strategy among all these choices, we have to order both the program's clauses and the assumptions in a sequent. The clauses are ordered from top to bottom, that is, the program is searched from top to bottom. The assumptions are ordered from left to right, that is, the assumptions are searched from left to right.

The contraction rule is handled through something called *annotation*. When the primitive $con(n,C)$ occurs as an assumption, the assumption C is annotated with the number n . This means that together with the assumption C the number n is stored. Each time this assumption is used by another rule than the rule I, its annotation number will be decreased by one. The annotation value is restored either by backtracking or when the current consequent is proved.

3. An Example Deduction

Here are two simple examples illustrating the basic behavior of GCLA's basic inference rules. Let the program P be

- a.
- b(1) :- a.
- b(Z) :- c.

The goal $a \rightarrow b(1)$ is then given by the following derivation:

$$\begin{array}{r}
 \frac{}{\langle \text{nil} ; \emptyset , P \rangle} \quad \text{(Initial sequent)} \\
 \frac{}{\langle (a \mid - a) . \text{nil} ; \emptyset , P \rangle} \quad (\neg P) \\
 \frac{}{\langle (a \mid - b(1)) . \text{nil} ; \emptyset , P \rangle} \quad (\neg \rightarrow) \\
 \langle (\mid - a \rightarrow b(1)) . \text{nil} ; \emptyset , P \rangle
 \end{array}$$

The goal $b(x) \mid - a$ may be solved in the following manner with answer substitution $\{X/1\}$:

$$\begin{array}{r}
 \frac{}{\langle \text{nil} ; \{X/1, Z/1\} , P \rangle} \quad (P\mid -) \\
 \frac{}{\langle (c \mid - a) . \text{nil} ; \{X/1, Z/1\} , P \rangle} \quad \text{(Initial sequent)} \\
 \frac{}{\langle (a \mid - a) . (c \mid - a) . \text{nil} ; \{X/1, Z/1\} , P \rangle} \quad (P\mid -) \\
 \langle (b(X) \mid - a) . \text{nil} ; \emptyset , P \rangle
 \end{array}$$

where $D(b(X/1, Z/1)) = (a, c)$ given the stated order of the program clauses and $D(c) = \text{nil}$. ($D(x)$ is defined by the rule $P\mid -$).

4. Some Programming Examples

The examples below are just small examples of code, to get some feeling of what a GCLA program could look like, and what queries could be put to the system. Remember that the program defines the system \vdash_p , so \vdash_p differs from one program to another. We start with a small example showing negation. As a false symbol we can choose any symbol that does not have a definition in the program. In the program below we have chosen the symbol `false`. This symbol corresponds to the empty set, i.e. there is no definition of what `false` should mean in terms of the other objects in the universe defined by the program.

First we define the object `clyde` as an elephant. `fido` is also defined as an elephant. Then we also define all objects that are albino elephants to also be elephants. The object `karo` is defined to be an albino elephant. Now, what constitutes a grey object? Well, in our universe it should be an elephant that is not an albino elephant.

```

elephant (clyde) .
elephant (fido) .
elephant (P) :-
    albino_elephant (P) .

albino_elephant (karo) .

```

```

grey(P) :-
    elephant(P),
    (albino_elephant(P) -> false).

```

Now, what possible questions could be put to this small program? The first question is "Which objects are elephants?", in GCLA syntax `|- elephant(Object)`, and the system will respond with the variable `Object` bound to `clyde`, and with `fido` and `karo` upon backtracking. If the query `|- grey(Object)` is put to the program, that is, asks which objects are grey, `Object` will get bound to `clyde`, and to `fido` upon backtracking, but not to `karo` because `karo` is an albino elephant.

The next example is a small toy expert system. The universe here is the set of diseases, symptoms and temperatures. Diseases are defined in terms of temperature and symptoms. The temperature is defined to be `high`, `normal` or `low`. All symptoms are true, i.e. there are no false symptoms, which is accomplished by the circular definition `symptom(X) :- symptom(X)` (these circular definitions are sometimes referred to as *opening* the predicate, also see section 9). The circular definitions should be read as "if `X` is a symptom, then it is a symptom", i.e., if `X` is assumed to be a symptom, then the program cannot inform us what it should mean that `X` is a symptom. The same is applicable for the definition of temperatures; if we assume `temp(high)`, the meaning of that is that the temperature is high, nothing more. But to assume `temp(red)` or something like that is absurd, since there is no definition of what `temp(red)` should mean (in this program).

```

disease(plague) :-
    temp(high),
    symptom(perspire).
disease(cold) :-
    temp(normal),
    symptom(cough).
disease(pneumonia) :-
    temp(high),
    symptom(multiple_chills),
    symptom(persistent_cough).

symptom(X) :- symptom(X).
temp(high) :- temp(high).
temp(normal) :- temp(normal).
temp(low) :- temp(low).

```

Then some of the possible questions are:

- i) If the person has a normal temperature and a cough, what disease does he suffer from? The query put to the system is:
 $\text{temp}(\text{normal}), \text{symptom}(\text{cough}) \mid - \text{disease}(K),$
 resulting in the answer substitution $K = \text{cold}$
- ii) We know that a person has a high temperature and multiple chills, and we are interested in additional symptoms for the possible diseases. The query is
 $\text{temp}(\text{high}), \text{symptom}(\text{multiple_chills}), \text{symptom}(\text{Symp}) \mid - \text{disease}(K)$
 resulting in the answer substitution $\text{Symp} = \text{persistent_cough}$ and $K = \text{pneumonia}$
- iii) If we assume the plague, what temperature is the patient supposed to have? In the query this time the consequent and the assumption are reordered
 $\text{disease}(\text{plague}) \mid - \text{temp}(G),$ and the system comes up with the answer substitution $G = \text{high}.$

In the real program the circular definitions of symptoms and temperatures are handled through declarations to prevent looping programs. The different possible declarations are given later.

As was mentioned before, GCLA integrates logic programming and functional programming into one single framework. This is accomplished by using the assumptions to evaluate expressions and using the Initial sequent to bind variables.

Consider the functional definition of `add`, to add two numbers. Two clauses are defined, `add` and `s`, where `s` is a substitution schema saying "if the argument to `s` can be evaluated to a term denoting the same value, continue with that new term". In this case this term is known to be canonical and cannot be reduced any further. The guiding primitive `&axiom` is used to restrict the canonical term to be used just by the axiom rule.

```
add(0,N) :- N.
```

```
add(s(M),N) :- s(add(M,N)).
```

```
s(Y) :- ((Y -> X) -> &axiom(s(X))).
```

Then some of the possible questions are:

- i) What is one plus one?

```
add(s(0),s(0)) \mid - X.
```

resulting in the answer substitution $X = s(s(0)).$

- ii) What should be added to one to get two?

```
add(X,s(0)) \mid - s(s(0)).
```

resulting in the answer substitution $X = s(0).$

5. The Ideas Behind the Abstract Machine (GAM)

One of the main ideas behind the abstract machine, called GAM, is to build upon WAM [10]. The unifying primitives have not been changed much, nor have the indexing instructions. The control instructions have been replaced by an extended set of instructions, among other things to handle the creation and deletion of assumptions. In WAM a clause is compiled into a head part and a body part. The body part consists of a number of goals to satisfy. Every clause is looked upon as a procedure, and each body goal as a procedure call. The current goal is the program counter pointing into a procedure representing this goal. In GAM there are also the assumptions which should be executable. Together with every assumption a "program counter" is stored pointing towards the code area where the execution should start if this assumption is chosen. This implies that the code must be executable for both the rules $\vdash P$ and $P \vdash$. This implies that the instructions must perform different things depending on the rule. Furthermore the code must be executable for the rule axiom.

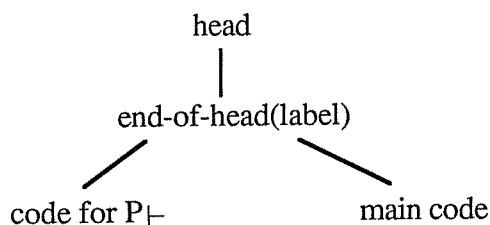
For both $P \vdash$ and $\vdash P$, the unification in the head of a clause is the same. What differs is how the body of a clause should be treated. So, each clause is compiled into one head part and two bodies, one for the $P \vdash$ rule and one for the main body. A typical clause will look like:

```

get-instructions
end_of_head(label)
create_premises-instructions      call this code "code for P⊢"
take_next_clause
label put-instructions and        call this code "main code"
call-instructions

```

which also could be illustrated by



The purpose of the left code is to create assumptions, which point towards the main code of the clause. When one of these new assumption is executed, the code in the main body will be executed.

As we mentioned earlier the machine could be in different modes. These are *rule*, *right*, *left* and *axiom*. In mode *axiom*, the rule axiom is tried on the current goal.

In mode `right` the current consequent is examined. In mode `left` an assumption is chosen and examined. In mode `rule` the other modes are tried one after another; first try the `axiom` mode, then try the `right` mode and last try the `left` mode.

The head of a clause is always executed in `right` mode or `left` mode, the code for `P` is always executed in mode `left` while the main code can be executed in `axiom` mode, `right` mode or `left` mode.

Then the indexing instructions and the choicepoints' instructions are used to link the clauses of a predicate together into one procedure.

There is also an area for handling the encoding of the procedural semantics, called the *rule code* or *guiding code*. As of the writing of this report it is not clear how the rule code will look like, but how it should act is described in the section 12 "Realization of the Procedural Semantics".

6. An Introductory Example of Compiled Code

In all the examples we have used relocatable code, that is, all labels are relative to the first instruction. As a first example of code we give the ordinary definition for `append` compiled into GAM code. The GCLA definition looks like

```
append([], Y, Y).
append([F|X], Y, [F|Z]) :-
    append(X, Y, Z).
```

and the compiled code looks like

	<code>append / 3</code>	<code>append(</code>
<i>head</i>	[[[L4], [L10]]],	
	switch-on-term(L3, L4, L10, fail),	
L3	try-me-else(L9)	
L4	get-nil(A1),	[].
	get-value(A1, A2),	Y,Y
	end-of-head(L8),)
<i>left code</i>	take-next-clause,	.
<i>main code</i> L8	end-of-clause,	.
	L9 trust-me-else-fail,	
<i>head</i>	L10 allocate(3),	[
	get-list (A0),	F
	unify-variable(A4),	X],
	unify-variable (Y0),	Y,
	get-variable(Y1, A1)	[
	get-list (A2),	F
	unify-value (A4),	

	unify-variable(Y2),	Z
	end-of-head (L21),):-
<i>left code</i>	create-premise(append / 3, 1, L21),	append(X,Y,Z)
	take-next-clause,	.
<i>main code</i> L21	put-value (Y0, A0),	append(X,
	put-value(Y1, A1),	Y,
	put-value(Y2, A2),	Z)
	execute(append / 3).	.

In GAM there is no instruction `deallocate` as there is in the WAM. The environment is deallocated, if it is possible, by the instructions `execute` and `end-of-clause` (i.e. the environment is created after the topmost choicepoint and there is no assumption referring to this environment).

7. Data Areas and Registers

A lot of the areas and registers in GAM are roughly the same as in WAM. The new areas are the area holding the assumptions and the area holding the registers for an assumption (which corresponds to the consequent's argument registers). Also the terms environment, continuation and choicepoint have the same meaning as in WAM.

The different data areas in GAM are:

- S The stack, contains environments and continuations. This area should take care of the procedural information during the execution. It corresponds roughly to the stack in PASCAL or the dump in the SECD machine [9].
- H The heap, holds global values, structures, lists, etc (The same as in the WAM). This area must be garbage collected. All objects that should be visible when leaving the current procedure should reside in this area, and therefore the answer substitutions are built from the heap.
- P The assumptions (also called *premises*) hold the defined assumptions, which form a tree. The tree arises when the $P\bar{I}$ -rule is used and there is more than one definition that is unifiable with the assumption chosen. Each assumption has a pointer to its father. The register CA always points to the leaf of the current list of defined assumptions.
- B Backtrack, a stack containing the choicepoints.
- T Trail, a stack containing things to be undone upon backtracking. The cell should be one of the following: a reference pointer to a bound variable, a pointer to an assumption whose annotation value should be increased, or a pointer to an assumption whose annotation value should be decreased.
- R The argument registers for the consequent, numbered A_1, \dots, A_n (the same as in WAM)

L The argument registers for an assumption, numbered B_1, \dots, B_n
Of course there is also the code area, split into two parts; one containing the code for predicate definitions, and one part containing the code for guiding the execution.

The *argument registers* are used to pass arguments between a calling procedure and the called procedure. They are also used as temporary registers inside a procedure when a variable does not need to be permanently stored. There are two different sets of argument registers, one for the consequents and one for the assumptions.

An *environment* is a frame holding the so called permanent variables, which are such variables that must be saved when a procedure is called. An environment also contains the value of CA (Current-Assumption, see below) to be able to "pop off" assumptions when they should not be defined any longer. It is used to restore the state as it was when this procedure was entered. The environment is created by the instruction *allocate* in the called procedure, if it needs an environment. All procedures that are not representing an atomic clause need an environment. The environment is deallocated by one of the instructions *execute* or *end-of-clause*.

A *continuation* contains information for continuing the execution in a procedure when a subgoal is successfully finished. It contains such information as the mode, the value of top of stack when the subgoal was entered, a pointer to the previous continuation, a pointer to the list of assumptions when the subgoal was entered, the environment pointer when the subgoal was entered and the continuation value of the program counter (see appendix A.3).

A continuation is created:

- every time a procedure is called by the instruction *call* in mode *right*,
- by the instruction *take-next-clause* when there is more than one clause unifiable with an assumption if the mode is *left*,
- by the instruction *push-premise* when the mode is *left*.

A *choicepoint* consists of enough information to restore a prevailing state. In GAM there are 5 different choicepoints, called *rule*, *right*, *axiom*, *left1* and *left2*. These store different kinds of information. (see appendix A.1).

An unbound variable is represented by a cell containing a reference pointer to itself.

Other types of pointers are structure pointers and list pointers. There are two data types, integers and symbols, where a symbol consists of a functor and the arity of that functor, for example *f/1* and *f/2* are different symbols. Constants are functors of arity 0.

In WAM the stack contains variables that could be popped off when a procedural call is ended, while the heap contains bindings that could not be forgotten (i.e. the bindings that are answers/instantiations of variables in previous goals). To prevent dangling references when an environment is popped off the stack there should not be any pointers from the heap to the stack. The pointers in the stack should point from younger variables to older ones. On the heap the pointers could be in either direction, but the order of the objects on the heap is still significant for efficient treatment of backtracking.

All the data areas are treated as stacks when a new object should be created. Therefore all

the areas have a pointer to the top of the area. The areas S, B and T are "true" stacks, while P is merely treated as a stack when creating new objects and when popping off objects during backtracking. The objects in P are then linked together forming a tree. H is garbage collected when the upper bound of H is (almost) reached. The names of each "top of the area" register are:

ToS	Top of stack
ToH	Top of heap
ToP	Top of premises
ToB	Top of backtrack
ToT	Top of trail

The other registers are:

Me	Mode register, holding the mode. One of: <code>rule</code> , <code>right</code> , <code>axiom</code> , <code>left</code> . For convenience two additional modes are added, <code>fail</code> for an execution that fails, and <code>finish</code> for an execution that is successfully executed.
CA	Current Assumption, points to the leftmost assumption in the list of assumptions (the last created). This assumption is one of the leaves in the assumption tree in the area P.
CAP	Current Assumption Pointer, a help register for pointing out a specific assumption, used when a premise is to be chosen by the rules <code>Initial sequent</code> and <code>P_l</code> .
TA	Treat as Absurd, a flag which is set to true every time the rule <code>P_l</code> is entered. As soon as one head in the program is unifiable with the chosen assumption, this flag is set to false, indicating that there was at least one head defining this assumption. The purpose of this register is to determine when an assumption does not have a definition in the program, in which case the current sequent trivially holds (see the rule <code>P_l</code> , page 7).
CL	Clauses Left, a register holding the clauses left to unify with under the mode <code>left</code> . Actually CL points to the first address not tried in the clause order associated with each predicate.
ST	Succeed To, points to the next continuation on the stack, which is where the execution should proceed to when this clause is successfully finished. (This register differs slightly from the CP register in the WAM in that it points towards the stack where the next continuation is stored instead of holding the next value for the program counter.)
PC	Program Counter.
Sp	Structure pointer, a help register for unifying structures. This register always points towards the heap. It is the same as the S-register in the WAM.
Um	Unification mode, read or write. The same as mode in the WAM, and

- affects only the unify instructions.
- E Environment pointer, points to the current environment in the stack.
- HB Heap Backtrack pointer, points towards the heap where the top of the heap was when the last choicepoint was created. (This register is actually not necessary because the value of HB can be found in the most recent choicepoint's ToH-cell).
- RC Right Caller, holds the functor/arity of the current consequent.

8. The Clause Orders

For the rule $P\vdash$ a suitable permutation of the clauses in the program should be chosen, where the resulting unifier of the ordered clauses should be unified with the current assumption (see 2.2.6). This choice of permutation is a choicepoint. However, a lot of these permutations result in the same unifier, so a lot of execution time could be saved if just the orders of the clauses that result in different unifiers are tried (The number of possible unifiers for a set with n terms does actually not exceed n , which is far less than $(n!)$, but we will not prove this in the paper). The permutations could be tested in advance by the compiler, thus generating code for the interesting cases.

The *clause orders* is a list, consisting of all the possible trials that could be unifiable with an assumption. These orders represents two "kinds" of failure; failure in the unification with the current chosen assumption, and failure later in some other inference step. The latter kind of failure could be due to wrong unifier, and therefore another clause order should be tried, so this is a choicepoint, while the first kind of failure does not give rise to a choicepoint.

Lets illustrate this by an example. Consider the program

```

1:   p(1,X) :- b1.
2:   p(2,X) :- b2.
3:   p(Y,3) :- b3.
4:   p(Y,4) :- b4.

```

where b_1 to b_4 could be arbitrary bodies. With this set of clauses the possible different unifiers are

[1,3], [1,4], [2,3], [2,4]

representing

$p(1,3)$, $p(1,4)$, $p(2,3)$, $p(2,4)$

respectively.

The orders are actually representing the permutations which the clauses should be tried. For example, the order [1,3] is representing all the permutations (1,3,2,4), (1,3,4,2), (3,1,2,4) and (3,1,4,2) in the algorithm for $P\vdash$, giving in section 2.2.6.

If now the assumption is unifiable with at least one of these orders (all the heads in a list must be unifiable with the assumption), a choicepoint is created for the next possible order. For example, in the goal

$p(z, 3) \text{ :- whatever}$

$p(z, 3)$ is unified with the order [1,3], where z is bound to 1, and if backtracking occurs, $p(z, 3)$ is unified with the order [2,3] giving as result z bound to 2.

But if the assumption is $p(z, 5)$ no one of the orders above is applicable, yet there is two possibilities, namely one of the heads $p(1, x)$ or $p(2, x)$. Therefore another layer is introduced, giving

$$\frac{[1,3], [1,4], [2,3], [2,4]}{[1], [2], [3], [4]}$$

representing

$$\frac{p(1, 3), p(1, 4), p(2, 3), p(2, 4)}{p(1, z), p(2, z), p(z, 3), p(z, 4)}$$

The lower layer is used if no order in the layer above is unifiable with the current assumption. Now the assumption $p(z, 5)$ will be replaced by the body b_1 , and with b_2 upon backtracking. So, in the horizontal plane, a choicepoint is created upon successful unification, and in the vertical plane the next plane is just tried if no order in the plan above were unifiable with the current assumption.

The orders are in this paper represented as lists, so the orders above will look like

[[[1,3], [2,3], [1,4], [2,4]], [[1], [2], [3], [4]]].

9. Declarations

The declarations are used to get certain effects when the rule $P\vdash$ is tried. As shown in the examples before there are situations where an assumption should not be treated as not defined, which is the same as always treating the assumption as defined, although there is no actual definition of it. For example, in the toy expert system (see page 10) symptoms are always true, that is, a person could not suffer from a non-existing symptom. To get such behaviour declarations are used. The declarations do not affect

any rule except $P\bar{-}$. They are introduced instead of writing circular clauses, to avoid having the execution going into a loop.

If subsumption were introduced on the goals these declarations would not be necessary, but subsumption is a very expensive method. The approach with subsumption was tested on an earlier version of the experimental interpreters.

A definition could be declared in three ways:

- *normal*, which does not affect anything and is the default
- *total*, which means that this definition is not applicable to the rule $P\bar{-}$. Semantically this means: suppose that a clause is named t , and it is declared *total*. This is the same as adding the clause $t \text{ :- } t$ to the program.
- *otherwise*, which theoretically means add the clause $t \text{ :- } [g_1 \dots g_n], t$. to the program. This clause's guards are such that the head is not unifiable with any other clause's head. For example assume the clauses defining $t/1$ are $t(1)$ and $t(2)$. Then the *otherwise* declaration has the same effect (theoretically) as adding $t(X) \text{ :- } t(X)$ where X is constrained not to be unified with 1 or 2. If now an assumption is not unifiable with $t(1)$ or $t(2)$, the sequent will not trivially hold (see section 2). Instead backtracking occurs. This is sometimes referred to as "open the predicate t ".

10. The Control Instructions

The control instructions can be divided into 6 groups: calling instructions, premise instructions, indexing instructions and choicepoint instructions. The guard instructions will also be treated here, as well as the instruction *fail*.

In the following, "WAM" means that the instruction is roughly the same as in the WAM (at least in mode *right*), "new" means that the instruction is new and "changed" means that the instruction has the same name as in the WAM, but does not perform the same actions.

10.1 The Calling Instructions

The calling instructions are *call*, *execute*, *end-of-premise*, *end-of-head*, *end-of-clause* and *take-next-clause*. We have also treated the instruction *allocate* here.

call(F/A) *changed*

If mode equals *right*, a continuation is saved on the stack, the PC is set to the first instruction of the rule code, F/A is put into the register RC and the execution continues in mode *rule*.

If the mode equals *axiom*, *call* is treated like *end-of-clause*, that is, the next continuation is examined.

If the mode equals *left*, the program database is searched for F/A. If F/A exists the PC

is set to the first label in the clause order list. For example, if the predicate `q` looks like

```
q/1,  
  [[4, 17], [10, 42]],  
  switch-on-term....
```

then PC is set to 4 (or actually PC + 4). A choicepoint for the clause order is created, pointing at the rest of the list (in the example above pointing at [10, 42]), and the execution continues in mode `left`. If F/A does not exist the next continuation is examined, provided that F/A is not declared `total` or otherwise. If F/A is declared `total` or otherwise, the procedure `fail` is called.

execute(F/A)

changed

This instruction is used to get tail recursion optimization (TRO), that is, it replaces the last call-instruction in the code of a clause when TRO could be performed. The only difference from `call(F/A)` is that a continuation is not pushed on the stack in mode `right`.

end-of-premise

new

This instruction is used to mark the end of the assumption list in a body goal, which is an arrow. It is executed in mode `axiom` or `right`, and for these modes it equals the instruction `end-of-clause` (this instruction could actually be replaced by `end-of-clause`).

end-of-head(L)

new

This instruction is used to mark the end of a head in a clause, and is used to split the execution depending on the mode. If the mode is `left` we just continue to the next instruction, if the mode equals `right`, PC is set to L.

end-of-clause

new

This instruction marks the end of the code for $\vdash P$ which does not end with an `execute` instruction. It examines the next continuation (e.g., the one pointed to by the register `ST`) and restores the registers according to the continuation. If the continuation is above the last choicepoint (the top of stack field in the most recent choicepoint) the continuation is popped off the stack.

allocate(N)

WAM

This instruction occurs in the head of a clause, which has at least one body goal. It can be executed in mode `right` or mode `left`. It allocates an environment containing $N+1$ cells on the stack, N permanent variables and the current value of `CA`. The variables are initialized to unbound variables (this is necessary because it is not certain that the variable is initialized by a `put-variable` instruction or a `get-variable` instruction. This is due to the fact that the assumptions could be chosen in several different orders).

take-next-clause*new*

This instruction marks the end of the code for $P\bar{-}$. It examines the current clause order if there is another clause to check, and if so the address to that clause is put into the PC register and pushes a continuation of mode *rule* on the stack. If there are no more clauses to check, the next continuation is examined (which should be a continuation of mode *rule*, pushed by the machine when the $P\bar{-}$ rule was chosen, or by a take-next-clause executed before; see the section about realization of the procedural semantics).

10.2 The Premise Instructions

Create-premise, push-premise and restore-premise belong to this group of instructions.

create-premise(L, N, F/A)*new*

This instruction is only executed in mode *left*. It creates a new assumption in the premise area. The assumption points to the address *L*, has annotation value *N* and has the name *F/A*. *N* is dereferenced, and if it is not a positive integer an error is signaled. The register *CA* is set to this new assumption, which in turn has a pointer to the former value of *CA*.

push-premise(F/A, N, L)*new*

If mode equals *right*, a new assumption is made in the premise area. The assumption has a pointer to the address *L*, an annotation *N* and the name *F/A*. *N* is dereferenced, and if it is not a positive integer an error is signaled. The register *CA* is set to this new assumption and the former *CA* value is stored in the assumption.

If mode equals *left*, a continuation of type *left* is stored on the stack and the mode is set to *right*. The reason for this is that the rule *arrow-left* is going to be executed, and the arrow's premises should be executed in mode *right* (to check if they hold). If that succeeds, the former state must be recalled when a new assumption is created for the arrows consequent.

restore-premise*new*

This instruction restores a former state where possibly not all the currently defined assumptions were defined. The value of *CA* is set to the value stored in the current environment.

10.3 The Indexing Instructions

The task of these instructions is to reduce the number of choicepoints created plus increase the efficiency by minimizing the number of applicable clauses in a predicate. These instructions are always executed in mode *right*. *Switch-on-term*, *switch-on-*

constant and switch-on-structure belong to this group.

switch-on-term(Var,Const,List,Struct) *WAM*

The value of the first argument register (A_1) is dereferenced and this value is examined. If it is a reference to an unbound variable, a jump to the label Var is made. If Var is the constant fail, backtracking occurs. An analogous action occurs if A_1 dereferences to a constant, a list pointer or a structure pointer.

switch-on-constant(Table, Default) *WAM*

The table consists of constant-label pairs of constants occurring as the first argument in the clauses of a predicate, and can be hashed for maximal efficiency. The value of the first argument register (A_1) is dereferenced and this value is examined. If the constant is in the table, the PC is set to the corresponding label, otherwise the PC is set to Default. If Default equals fail, backtracking occurs.

switch-on-structure(Table, Default) *WAM*

This instruction is analogous to the switch-on-constant instruction, but the table consist of functor/arity-label pairs, and is used to dispatch on the structures functor.

10.4 The Choicepoint Instructions

The instructions handling the choicepoints for the rule $\vdash P$ are try, retry, trust, try-me-else, retry-me-else and trust-me-else-fail, and they are therefore always executed in mode right. The choicepoints for the other rules are currently handled inside the system. (The instructions try-me-else, retry-me-else and trust-me-else-fail are in fact not necessary, the same task can be fulfilled by the instructions try, retry and trust.)

try(Label) *WAM*

This instruction creates a choicepoint with the next try set to the instruction after this one. Then the execution proceeds at Label.

retry(Label) *WAM*

The topmost choicepoint is updated with the next try set to the next instruction. Then the execution proceeds at Label.

trust(Label) *WAM*

The topmost choicepoint is popped of and the execution proceeds at Label.

try-me-else(Label) *WAM*

This instruction creates a new choicepoint on top of the backtrack area. The next try in this choicepoint is Label, and the execution proceeds to the following instruction.

ground(X)

This instruction dereferences X and examines the value. If it is currently bound to an object not containing an unbound variable the execution proceeds, otherwise backtracking occurs. If mode is `left`, X refers to one of the L-registers, otherwise X refers to one of the R-registers.

10.6 The Fail Instruction

This instruction restores a former state which was a choicepoint in the execution. Depending on the topmost choicepoint different areas and registers are updated, and the execution proceeds from the point stored in the choicepoint. The choicepoints can be of 5 different types:

- rule, where it represents the next choice of a mode (`axiom`, `right` or `left`)
- axiom, where it represents the next choice of an assumption
- right, where it represents the next choice of a clause in the program
- left1, where it represents the next choice of an order to try among the clauses in the program
- left2, where it represents the next choice of an assumption (actually the same kind of choice as axiom, but we must perform different tasks upon backtracking).

The procedure first restores all the different registers stored in the choicepoint, and the trail is unwound to the address stored in the choicepoint. Now the state is restored to what it was when the choicepoint was created. Then, depending on the type of the choicepoint, different actions are performed in order to create the next choice:

- If the type is axiom, a new assumption is searched for with the same functor as the current consequent, which is held by the register RC. If one is found, this is tried, otherwise failure to the next choicepoint occurs.
- If the type is rule, the execution simply proceeds to the instruction pointed to by the choicepoint (somewhere in the rule code).
- If the type equals right, the execution proceeds to the instruction pointed to by the choicepoint.
- If the type of the choicepoint is left1, the next possible order of the program's clauses is tried, which is stored in this choicepoint. If there is no more trials in this plane to do (see section 8), the register TA is examined. If it is false, there was at least one clause order that has succeeded in this plane, and the next choicepoint is examined. If TA is true, the next plane is tried, updating the choicepoint with the next plane number and the next order to try, which is the first one on this new plane. If there are no more planes to search, the current goal sequent holds and the next continuation is examined.
- If the type is left2, there is much to do; another assumption should be tried. The next assumption is searched for and examined. If the next assumption is the last one in the list (indicated by an empty pointer in the linking field of the assumption), this choicepoint of type left2 is popped off the stack, otherwise this choicepoint is updated with the next assumption in the list (the value of the linking field in the chosen

get-variable(V_n, X_i)

WAM

The variable V_n is bound to the dereferenced value of X_i . This instruction corresponds to a variable occurring for the first time in a clause.

get-value(V_n, X_i)

WAM

The instruction gets the value of X_i and tries to unify it with the contents of V_n . If it succeeds the dereferenced value of X_i is left in V_n if V_n is a temporary variable. If it fails to unify V_n and X_i , backtracking occurs. This instruction should occur when V_n is initialized by another instruction before.

get-constant(C, X_i)

WAM

This instruction dereferences X_i and the value is compared with C . If they are not unifiable (that is, if the value and C are not identical and the value is not an unbound variable) backtracking occurs, otherwise X_i is bound to C .

get-nil(X_i)

WAM

This instruction is the same as the instruction `get-constant`, but the constant is the constant `nil`.

get-structure($F/A, X_i$)

WAM

This instruction marks the beginning of a structure occurring in the head of a clause. The register X_i is dereferenced and the value examined. If it is an unbound variable, that variable is bound to a structure pointer pointing to the top of the heap, F/A is pushed onto the heap and the unification mode is set to write. Otherwise if the dereferenced value is a structure pointer pointing to a cell identical to F/A , the unification mode is set to read and the S -register is set to point to the beginning of the arguments of the structure, otherwise backtracking occurs.

get-list(X_i)

WAM

This instruction is analogous to the instruction `get-structure`, except that the structure is replaced by a list and the structure pointer is replaced by a list pointer.

11.2 The Put Instructions

These instructions cannot fail in mode `left` or `right`, they simply store data/references in the argument registers. They can be called in mode `left`, mode `right` and mode `axiom`. In mode `axiom` they are treated as their corresponding `get`-instructions in mode `right` (for example, in mode `axiom` the instruction `put-value` will be executed as the instruction `get-value` in mode `right`), except for the instruction `put-variable`, which is executed as `get-value`. The exception is due to the fact that the

assumptions are not guaranteed to be executed in the order the body goals are ordered. Therefore we always have to perform the more expensive unification of two variables which the instruction `get-value` does, instead of just setting one variable to the content of the other. This is also the reason why all the permanent variables in an environment have to be initialized to unbound variables when they are created by the instruction `allocate`. The description below is just for `right` mode and `left` mode.

put-variable(Y_n, X_i)

WAM

This instruction corresponds to a variable occurring for the first time in the (main) body of a clause. If the mode equals `right`, the variable Y_n is initialized to an unbound variable, and X_i is bound to Y_n . If the mode equals `left` this instruction performs the same actions as `put-value`. (This instruction is actually not necessary, because the instruction `allocate` initializes all its variables to unbound. It could be replaced by `put-value(Y_n, X_i)`.)

put-variable(X_n, X_i)

WAM

This instruction corresponds to a variable occurring in just one atomic goal (not a goal containing `->`) of a body. If the mode equals `right`, both the variable X_n and X_i are bound to a new unbound variable on the top of the heap. If the mode equals `left` this instruction performs the same actions as `put-value`.

put-value(V_n, X_i)

WAM

The instruction puts the value of V_n into X_i . It occurs when V_n has been initialized by another instruction.

put-unsafe-value(Y_n, X_i, N)

WAM

In mode `right` this instruction dereferences Y_n and examines the result. If the result is a variable in the current environment, a new global variable is created on top of the heap, the variable in the environment is set to the new variable and X_i is set to the new variable. Otherwise the instruction performs the same actions as `put-value`. In mode `left` this instruction does the same thing as `put-value`, and in mode `axiom` this instruction does the same thing as `get-value`.

put-constant(C, X_i)

WAM

This instruction puts the constant C into X_i .

put-nil(X_i)

WAM

This instruction is the same as `put-constant`, where the constant is the constant `nil`.

put-structure($F/A, X_i$)

WAM

This instruction marks the beginning of a structure occurring in a body goal. F/A is pushed onto the heap and the corresponding structure pointer is put into the register X_i . Then the unification mode is set to write.

put-list(X_i)*WAM*

This instruction is analogous to put-structure, except that the structure is replaced by a list and the structure pointer is replaced by a list pointer.

11.3 The Unify Instructions

These instructions occur in a complex argument to a head in a clause (i.e. a structure or a list). They can be executed in unification mode read or write, depending on what the instructions get-list, put-list, get-structure and put-structure have set the unification mode to.

unify-variable(V_n)*WAM*

This instruction corresponds to a variable that is not previously initialized by another instruction. If it is executed in *right* and, furthermore, if it is executed in write unification mode, it pushes a new unbound variable onto the heap and sets V_n to point to it. In read mode it stores the value of what the Sp-register points to in V_n , and increments Sp.

Otherwise if it is executed in *axiom* mode or *left* mode, it shall perform the same actions as unify-local-value.

unify-value(V_n)*WAM*

This instruction corresponds to a variable V_n which has been initialized by another instruction. If the mode equals *right* and if it is executed in write mode, it pushes the value of V_n onto the heap. If it is executed in read mode, it unifies V_n and what Sp points to, leaving the dereferenced value in V_n if V_n is a temporary variable.

Otherwise, if it is executed in mode *axiom* or mode *left* it shall perform the same actions as unify-local-value to prevent references from the heap into the stack.

unify-local-value(V_n)*WAM*

This instruction does almost the same thing as unify-value. It is used where V_n is not guaranteed to be bound to a global value (i.e. not an unbound variable on the stack). It performs the same things as unify-value, except in write unification mode where the dereferenced value of V_n is a reference to a variable on the stack. Then a new variable is created on top of the heap and the variable on the stack is bound to the new variable and V_n is set to point to the new variable if V_n is a temporary variable.

unify-constant(C)*WAM*

If this instruction is executed in write unification mode the constant is simply pushed onto the heap. If the unification mode is read, this instruction dereferences the cell pointed to by the register Sp, and if that is unifiable with C (i.e. it is identical to C or is

an unbound variable, which is bound to C) the register Sp will be incremented and the execution proceeds, otherwise backtracking occurs.

unify-nil

WAM

This instruction is the same instruction as unify-constant, except that the constant is nil.

unify-void(N)

WAM

This instruction is used for variables whose values are not needed. In read unification mode this instruction simply adds N to the register Sp. In write unification mode N unbound variables are created on top of the heap.

12. Realization of the Procedural Semantics

To get a combination of efficiency and flexibility a new code area is introduced, called the *rule code*. The idea is to specify a number of primitives, and in these primitives the operational and procedural semantics should be expressed. With this rule code there is the ability to remove one rule by simply removing some of the primitives, or add a new rule, or reorder the trials in an efficient way. If this is a good idea or not is not clear at the moment, as it has not been implemented yet. As it is now there are just the default order and the default rules described earlier in the sections about operational semantics and procedural semantics. We will here describe how choicepoints for the choice of an assumption and how choicepoints for the clause order are treated and how the machine acts before trying a rule.

The default order is to try the different modes *axiom*, *right* and last *left*. The first thing that happens when a call-instruction or an execute instruction is executed in mode *right* is that the mode is set to *rule* and the machine starts a "mode-cycle", consisting of the mode-choices. First it creates a choicepoint for the next mode to try, which is the mode *right*. Then it tries to find an assumption for the axiom to try. If there is no assumption of the same functor and arity as the consequent, or if there is no assumption at all, backtracking occurs and the most recent choicepoint is examined (which was the one for trying the mode *right*). If there is at least one assumption with the same functor and arity as the current consequent, the PC is set to the code for that assumption, and a choicepoint for the next assumption (if there is one) is made. The mode is then set to *axiom*.

Now, suppose that the axiom fails. The rule-choicepoint is then popped off from the B-area (backtrack), and a new choicepoint is created for the mode *left*. Then the code for the procedure with the name held by the register RC is searched for. If the code exists for the current consequent, the PC is set to that code and execution proceeds in mode *right*. If there is no code backtracking occurs.

Now there is no mode to try after *left*, which is to be tried now, so the rule-choicepoint is popped off the stack, and an assumption is again searched for. If there are no assumptions, backtracking occurs. If there is one assumption, this assumption is

examined by setting the mode to `left` and loading the PC with the address hold by the assumption. This address points towards the main code of a body of some clause in the program. A continuation of type rule is stored and the environment pointer is set to the value hold by the assumption, and the annotation value is decreased by one before the execution proceeds in mode `left`.

If there is more than one assumption, a choicepoint for trying the next assumption is created before the chosen assumption is examined.

Now, if the execution succeeds in mode `left`, the execution returns to mode `rule` when the saved continuation is restored and then another mode-cycle takes place (now with the new assumptions, and perhaps with some new goals on the stack to solve, which were created by `P(-)`).

13. Compiling Schema

There are a lot of things that could be done, and should be done, by the compiler. A compiler for a somewhat simplified version of this machine has been implemented. But there is certainly room for a lot of ideas.

To fulfill the operational semantics a variable check on the variables in the body of a clause should be done. When the program-left rule is chosen, a clause which introduces new variables in the body is not allowed and that try of clause order should fail. This check should be done by the compiler, which should generate code where this has been taken care of in the clause order. However, there are a lot of places where this check is not necessary and, therefore, we have so far not implemented it.

A difference with WAM is that an environment has to be generated more often. In WAM the variables in the first body goal can be temporary, and the other ones in the rest of the body goals must be permanent. In GAM an environment has to be generated for every clause that is not an atomic clause, and all variables that occur in a body goal have to be made permanent. This is done because the program-left rule saves the environment pointer `E` in the assumptions it creates, and when these assumptions are examined the environment pointer is restored again. The axiom rule (I) has the same behaviour.

There is no instruction for handling nested structures. These have to be lifted out and unified separately. For example, the code for the head argument `s(f(X))` occurring inside a head is separated to `s(Y), Y = f(X)`, which is then compiled into

```
...get-structure(s/1, Ai),
   unify-variable(Ai),
   get-structure(f/1, Aj),
   unify-variable(Ak), ...
```

One of the main tasks the compiler has to perform is to generate the clause orders, so they are built as specified in section 8.

In order to be sure that there will be no pointers from the heap to the stack, the instructions `unify-variable` and `unify-value` have the same definition as `unify-local-value` when the mode equals `left` or `axiom`. This solution is expensive and a thorough examination should be made in order to find the occasions where there is a risk for pointers in the wrong direction.

14. Some Examples

The first example of code generation is the example of negation introduced in section 3. Recall that it should decide whether an object (an elephant) was grey or not. The GCLA code for this program is:

```

elephant(clyde).
elephant(fido).
elephant(P) :-
    albino_elephant(P).

albino_elephant(karo).

grey(P) :-
    elephant(P),
    (albino_elephant(P) -> false).

```

And below is the corresponding GAM-code.

	elephant / 1,	elephant(
	[[[L19, L9], [L19, L14]], [[L19]]],	
	switch-on-term(L8, L3, L19, L19),	
L3	switch-on-constant([[clyde, L4], [fido, L6]], L19),	
L4	try(L9),	
	trust(L19),	
L6	try(L14),	
	trust(L19),	
L8	try-me-else(L13),	
L9	get-constant(clyde, X ₀),	clyde
	end-of-head(L12),)
	take-next-clause,	.
L12	end-of-clause,	.
L13	retry-me-else(L18),	
	get-constant(fido, X ₀),	fido
	end-of-head(L17),)
	take-next-clause,	.
L17	end-of-clause,	.
L18	trust-me-else-fail,	
L19	allocate(1),	
	get-variable(Y ₀ , X ₀),	X
	end-of-head(L24),):-
	create-premise(albino_elephant / 1, 1, L24),	albino_elephant(X)

L24	take-next-clause, put-value(Y ₀ , X ₀), execute(albino-elephant / 1),	albino-elephant(X)
	albino-elephant / 1, [[[L4]]], switch-on-term(L4, L3, fail, fail),	albino-elephant(
L3	switch-on-constant([[karo, L4]], fail),	
L4	get-constant(karo, X ₀), end-of-head(L7), take-next-clause,	karo)
L7	end-of-clause.	.
	grey / 1, [[[L3]]], switch-on-term(L3, L3, L3, L3),	grey(
L3	allocate(1), get-variable(Y ₀ , X ₀), end-of-head(L9), create-premise(elephant / 1, 1, L9), create-premise(-> / 2), 1, L11), take-next-clause	X):- elephant(X), (albino-elephant(X) -> false)
L9	put-value(Y ₀ , X ₀), call(elephant / 1),	elephant(X)
L11	push-premise(albino-elephant / 1, 1, L13), goto(L16, false / 0),	(albino-elephant(
L13	put-value(Y ₀ , X ₀), call(albino-elephant / 1), end-of-premise,	X) ->
L16	call(false / 0), restore-premises, end-of-clause,	false)

The next program is a program for intuitionistic propositional logic. This is an interesting example because it shows how most of the different possibilities and primitives in GCLA are compiled into GAM-code.

The program looks like:

```
t(not X) :- (t(X) -> false).
t(X => Y) :- (t(X) -> t(Y)).
t(X & Y) :- t(X), t(Y).
t(X $ Y) :- t(X).
t(X $ Y) :- t(Y).
t(contr(Num, Goal)) :- con(Num, t(Goal)).
```

and the corresponding GAM-code looks like:

(This is perhaps the most "expanding" GCLA-program that could be written...)

t / 1,	[[[L7], [L22], [L39], [L52 ,L62], [L72], otherwise]],	t(
	switch-on-term(L6, fail, fail, L3),	
L3	switch-on-structure([[not/1, L7], [=>/2, L22]	
	[&/2, L39], [\$ /2, L4]	
	[contr/2, L72], fail),	
L4	try(L52),	
	trust(L62),	
L6	try-me-else(L21),	
L7	allocate(1),	
	get-structure(not/1, X ₀),	(not
	unify-variable(Y ₀),	X)
	end-of-head(L13),	:-
	create-premise(-> /2, 1, L13),	(t(X) -> false)
	take-next-clause,	
L13	push-premise(t/1, 1, L15),	t(X)
	goto(L18, false/0),	
L15	put-value(Y ₀ , X ₀),	t(X)
	call(t/1),	
	end-of-premise,	->
L18	call(false/0),	false
	restore-premises,	
	end-of-clause,	
L21	retry-me-else(L38),	
L22	allocate(2),	
	get-structure(=>/2, X ₀),	=>(
	unify-variable(Y ₀),	X,
	unify-variable(Y ₁),	Y))
	end-of-head(L29),	:-
	create-premise(->/2, 1, L 29),	t(X) -> t(Y)
	take-next-clause,	
L29	push-premise(t/1, 1, L31),	t(X)
	goto(L34, t/1),	
L31	put-value(Y ₀ , X ₀),	t(X)
	call(t/1)),	
	end-of-premise(t/1),	->
L34	put-value(Y ₁ , X ₀),	t(Y)
	call(t/1),	
	restore-premises,	
	end-of-clause,	
L38	retry-me-else(L51),	
L39	allocate(2),	
	get-structure(&/2, X ₀),	&(
	unify-variable(Y ₀),	X,
	unify-variable(Y ₁),	Y))

	end-of-head(L47),	:-
	create-premise(t/1, 1, L47),	t(X),
	create-premise(t/1, 1, L49),	t(Y)
	take-next-clause,	.
L47	put-value(Y ₀ , X ₀),	t(Y),
	call(t/1),	.
L49	put-value(Y ₁ , X ₀),	t(Y)
	execute(t/1),	.
L51	retry-me-else(L61),	
L52	allocate(1),	
	get-structure(\$/2, X ₀),	\$(
	unify-variable(Y ₀),	X,
	unify-void(1),	Y)
	end-of-head(L59),	:-
	create-premise(t/1, 1, L59),	t(X)
	take-next-clause,	.
L59	put-value(Y ₀ , X ₀),	t(X)
	execute(t/1),	.
L61	retry-me-else(L71),	
L62	allocate(1),	
	get-structure(\$/2, X ₀),	\$(
	unify-void(1),	X,
	unify-variable(Y ₀),	Y)
	end-of-head(L69),	:-
	create-premise(t/1, 1, L69),	t(X)
	take-next-clause,	.
L69	put-value(Y ₀ , X ₀),	t(X)
	execute(t/1),	.
L71	trust-me-else-fail,	
	allocate(2),	
	get-structure(contr/2, X ₀),	contr(
	unify-variable(Y ₁),	Num,
	unify-variable(Y ₀),	Goal)
	end-of-head(L79),	:-
	create-premise(t/1, Y ₁ , L79),	con(Num,t(Goal))
	take-next-clause,	.
L79	put-value(Y ₀ , X ₀),	t(Goal)
	execute(t/1),	.

This code needs some explanation. The symbol "otherwise" is the declaration of an open predicate. When the end of the unification order is reached, the symbol "otherwise" prevents the execution from succeeding with an assumption that is not unifiable with any clause in the program.

In the last clause the treatment of contraction is shown. Y₁ corresponds to the variable "Num" in the original GCLA. This variable must be checked so that it is currently

instantiated to a number. What happens (semantically) when it is not a number is not clear. With this treatment the current goals succeed because the guard-instruction number will fail. Perhaps the best solution is that the variable Y_1 should make use of something called *freeze*, a primitive for "freezing" a variable and corresponding goals until the variable has been bound [5]. However, this has not been implemented yet.

The next example is the program for the functional definition of the arithmetic function `add`. Remember the GCLA program

```
add(0,N) :- N.
add(s(M),N) :- s(add(M,N)).

s(Y) :- ((Y -> X) -> &axiom(s(X))).
```

The corresponding GAM code is

<pre> add/2, [[[L4], [L13]]], switch-on-term(L3, L4, fail, L13), L3 try-me-else(L12), L4 allocate(1), get-constant(0, X0), get-variable(Y0, X1), end-of-head(L10), create-premise(call/0, 1, L10), take-next-clause, L10 put-value(Y0, X0), execute(call/1), L12 trust-me-else-fail, L13 allocate(2), get-structure(s/1, X0), unify-variable(Y0), get-variable(Y1, X1), end-of-head(L20), create-premise(s/1, 1, L20), take-next-clause, L20 put-structure(add/2, X0), unify-local-value(Y0), unify-local-value(Y1), execute(s/1), s / 1, [[[L2]]], L2 allocate(2), get-variable(Y0, X0), end-of-head(L7), create-premise(->/2, 1, L7),</pre>	<pre> add(0, N) :- N . N . s(M, N) :- s(add(M,N)) . s(add(M, N)) . s(Y) :- (Y -> X) -> &axiom(s(X))</pre>
--	--

```

take-next-clause,
L7  push-premise(->/2, 1, L9),          (Y -> X)
     goto(L18, s/1, 0),
L9  push-premise(call/1, 1, L11),      Y
     goto(L14, call/1),
L11 put-value(Y0, X0),              (Y
     call(call/1),
     end-of-premise,                  ->
L14 put-value(Y1, X0),              X)
     call(call/1),
     restore-premises,
     end-of-premises,                 ->
L18 put-value(Y1, X0),              s(X)
     call(s/1),
     restore-premises,
     end-of-clause.

```

This also needs some explanation. The "predicate" call is a primitive which takes a term as argument. If it is currently instantiated to a structure or constant (not a number) it performs the same actions as if that structure were the goal. If it is not instantiated the only rule that currently could be used is the axiom. It should be pointed out that this is a restriction in the current definition of GAM, and that the uninstantiated goal in the theory could, of course, be matched against the program clauses as well.

The third argument in the goto-instruction in the code for s is the annotation value for the guiding primitive &axiom. The axiom rule can always use an assumption whatever the annotation value is, and therefore if the value is set to 0 initially, the only rule that could use the assumption is the axiom rule.

15. Conclusions and Future Work

To be able to do some performance tests we have implemented both the GAM and WAM in Common lisp, and made the two implementations as similar as possible. Then the same Prolog program was executed in the two implementations, and they are indicating that pure Prolog programs are executed in two-third's of the time that it takes an "original" WAM machine executes it. If there are assumptions in the goal GAM runs the program up to two-third's of the time an "original" WAM, depending on the assumptions.

In this implementation we have not made an effort to be as efficient as possible in all the steps the machine takes. For example, if there is no assumption, the axiom will still be tried and an assumption for the rules \leftarrow and \rightarrow will still be searched for before backtracking occurs (i.e. a choicepoint for the modes `axiom` and `left` is created although there is no assumption to examine). Here is certainly room for making the machine more efficient. Some of the registers could perhaps also be removed if other test conditions were used.

The code size for GCLA program compiled into GAM code is something like 10 - 20% larger as for a Prolog program compiled into WAM code. The increase is mostly due to the new code for `Pl`.

Also in mode `left` and mode `axiom` the more expensive instruction `unify-local-value` is used even though it is not needed in all places (or actually the instructions `unify-variable` and `unify-value` have the same definition as `unify-local-value` in mode `axiom` and mode `left`). It is used to be sure that there will not be pointers from the heap to the stack. Where `unify-local-value` is needed and where it is not should be clarified, perhaps new instruction(s) should be introduced.

Another idea is to introduce three arguments in the `put` and `get` instructions. The first two are the same as in the WAM, while the third one should be used instead of the first one when the mode equals `right` or `axiom`. This should increase the possibility to use temporary registers instead of permanent registers.

Something that is also wanted is a delaying mechanism for variables, called `delay` or `freeze` [5]. With an ability to delay a goal until its variables are grounded, there is (perhaps) a possibility to implement the so-called anti-unification, or disunification, of the rule `Pl`, and also to put constraints on the variables.

Yet another idea to get more efficient code for the rule `Pl` is that the unifiers, which is generated by the compiler of all the clauses in a predicate, could be stored together with the predicate in one way or another (i.e. the unifier as it is or by some, perhaps new, instructions). This should increase the efficiency of `Pl` when there are several clauses that are unifiable with each other in a predicate. This solution should then partially replace some of the existing instructions (for example, `take-next-clause` becomes obsolete).

The next step is to systematize the code for the rules. The idea is that the user should define the inference rules himself, and that a file containing this code, the so-called operators, should be loaded together with a program into the rule code area. The user should also have the possibility to change the inference mechanism under the execution from inside the program. How this will influence the code described herein is not clear, but the idea has some very nice properties, for example, it is clear that these operators could also be treated as definitions, and thus the semantics for them is clear.

However, this experimental work with GAM has shown that GCLA could be implemented efficiently, and that there are a lot of promising ways to examine towards a more efficient and flexible implementation.

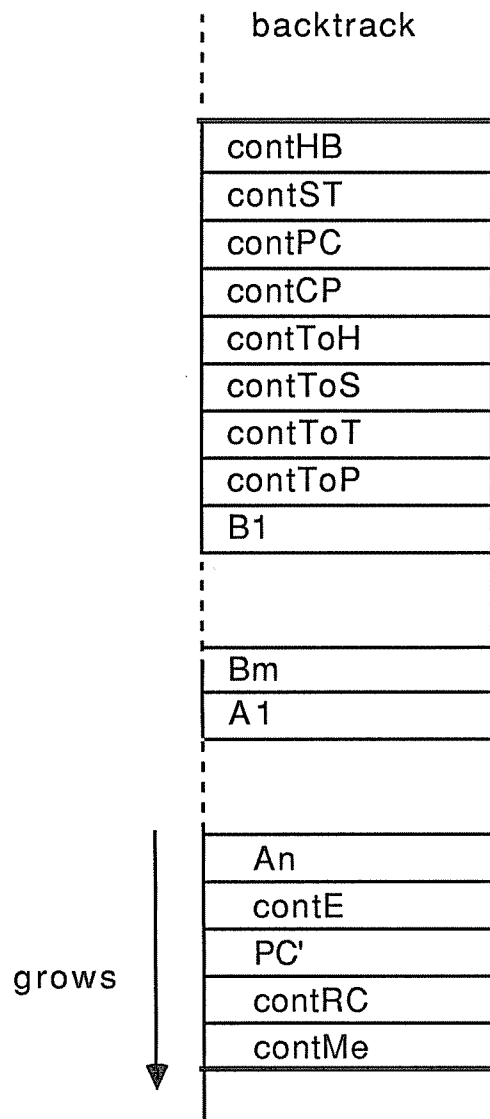
Acknowledgement

I am indebted to a lot of people at SICS for many valuable and helpful discussions. Special thanks goes to Kent Boortz and Peter Olin for many valuable and fruitful discussions, to Lars-Henrik Eriksson and Mats Carlsson for several suggested improvements to this paper, to Vicki Carleson for language support (all language errors are made by me after her careful work...), to Torbjörn Åhs at UPMAIL for comments on both an earlier version of this paper and this paper and to the project leader, Lars Hallnäs, for his never ceasing flow of new ideas and faith in the project.

References

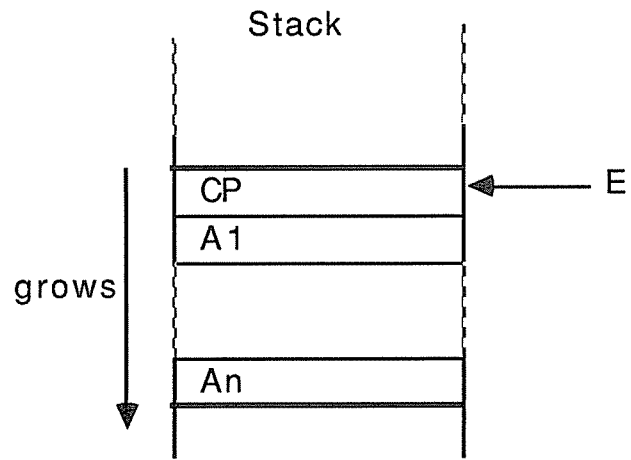
- [1] K. Appleby, M. Carlsson, S. Haridi, D. Sahlin, *Garbage Collection for Prolog Based on WAM*, SICS Research Report R86009B
- [2] M. Aronsson, *The Instruction Set for the GCLA Abstract Machine*, SICS Technical Report T89004
- [3] M. Aronsson, L. Hallnäs, *GCLA, Generalized Horn Clauses as a Programming Language*, SICS Research Report R88014
- [4] J. Barklund, *A Garbage Collection Algorithm for Tricia*, UPMAIL Technical report No. 37B
- [5] M. Carlsson, *An Implementation of dif and freeze in the WAM*, SICS Research Report R86012
- [6] L-H. Eriksson, L. Hallnäs, *A programming Calculus Based on Partial Inductive Definitions*, SICS Research Report R88013
- [7] L. Hallnäs, *Partial Inductive Definitions*, SICS Research Report R86005C
- [8]. L. Hallnäs, P. Schröder-Heister, *A Proof-Theoretic Approach to Logic Programming, I. Generalized Horn Clauses*, SICS Research Report R 88005
- [9] P. Henderson, *Functional Programming*, Prentice/Hall 1980
- [10]. D.H.D. Warren, *An Abstract Prolog Instruction Set*, SRI Technical Note 309, Menlo Park 1983

A.1 Choicepoints

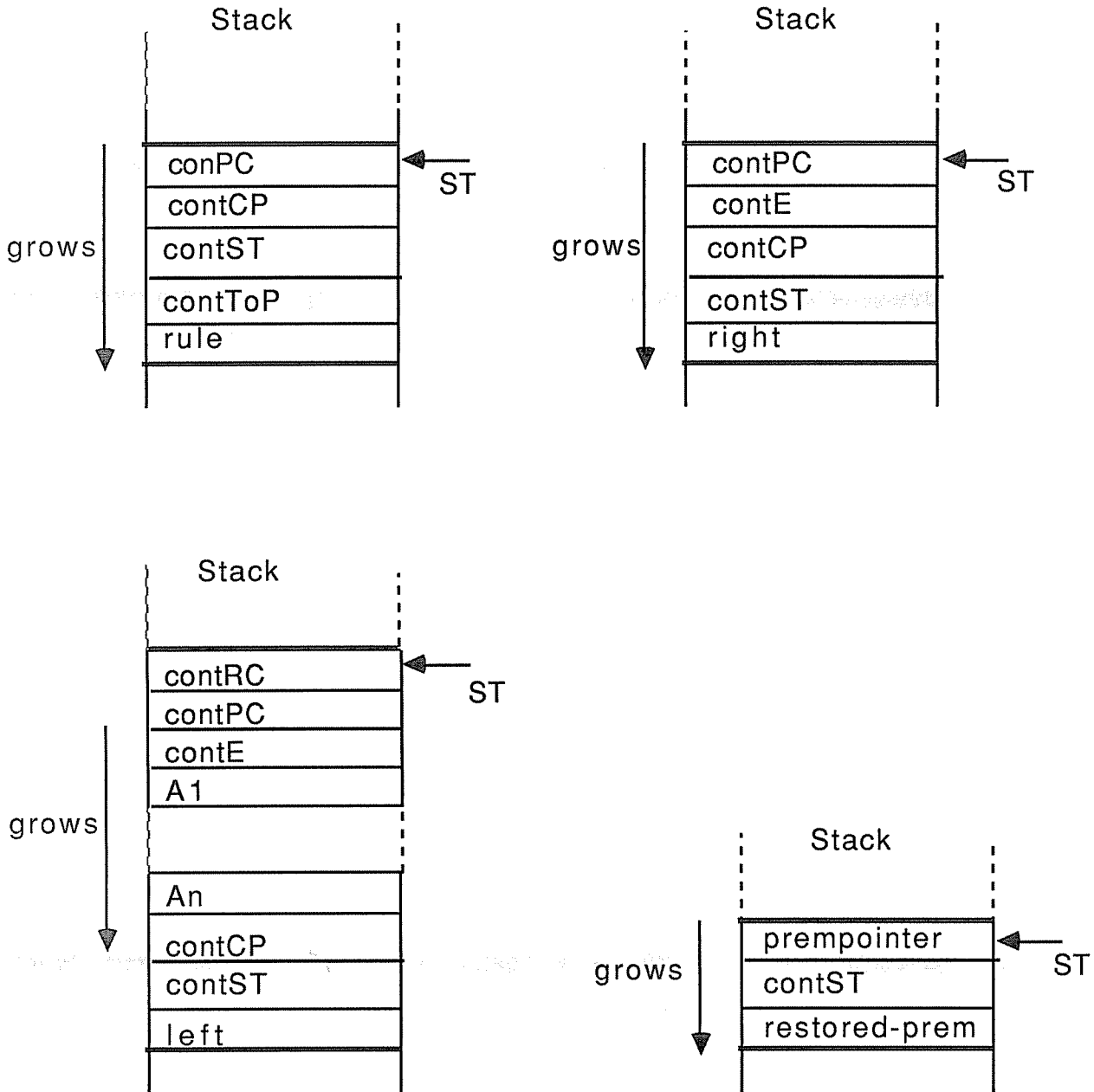


A1

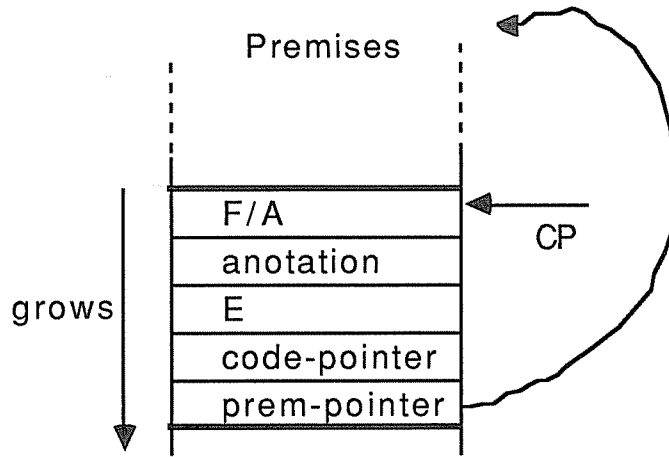
A.2 Environments



A.3 Continuations

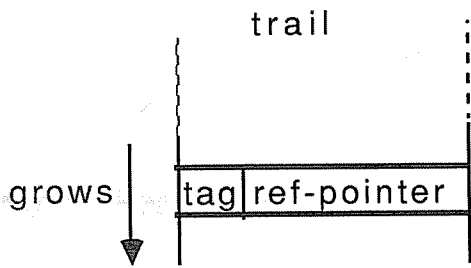


A.4 Premises

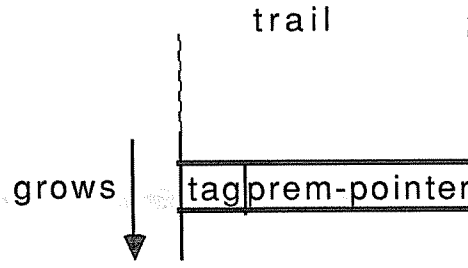


A premise

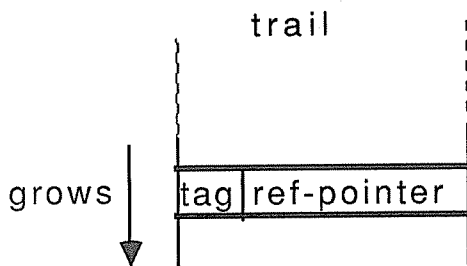
A.5 Different Trail Formats



A variable which
should be unbound



A premise whose
anotation-value
should be increased
by one



A premise whose
anotation-value
should be decreased
by one