
1983–1993: THE WONDER YEARS OF SEQUENTIAL PROLOG IMPLEMENTATION

PETER VAN ROY

- ▷ This article surveys the major developments in sequential Prolog implementation during the period 1983–1993. In this decade, implementation technology has matured to such a degree that Prolog has left the university and become useful in industry. The survey is divided into four parts. The first part gives an overview of the important technical developments starting with the Warren abstract machine. The second part presents the history and the contributions of the major software and hardware systems. The third part charts the evolution of Prolog performance since Warren’s DEC-10 compiler. The fourth part extrapolates current trends regarding the evolution of sequential logic languages, their implementation, and their role in the marketplace. ◁
-

1. INTRODUCTION

This article is a personal view of the progress made in sequential Prolog implementation from 1983 to 1993, supplemented with learning of the wise [10]. 1983 was a serendipitous year in two ways: one important and one personal. It was the year when David H. D. Warren published his seminal technical report [164] on the new Prolog engine, which was later christened the WAM (for Warren abstract machine).¹ It also marked the beginning of my research career in logic programming.

The title reflects my view that the period 1983–1993 represents the “coming of age” of sequential Prolog implementation. In 1983, most Prolog programmers (except for a lucky few at Edinburgh and elsewhere) were still using interpreters. In 1993 there are many high quality compilers, and the fastest of these are approaching or exceeding the speed of

Address correspondence to Peter Van Roy, Digital Equipment Corporation, Paris Research Laboratory, 85 Avenue Victor Hugo, 92500 Rueil-Malmaison, France. E-mail vanroy@prl.dec.com

¹The name WAM is owing to the logic programming group at Argonne National Laboratory. Received May 1993; accepted January 1994.

imperative languages. Prolog has found a stable niche in the marketplace. Commercial systems are of high quality with a full set of desirable features, and enough large industrial applications exist to prove the usefulness of the language [101, 102].

1.1. *The Influence of the WAM*

The WAM was the starting point for a veritable “gold rush” of Prolog developers, all eager for that magical moment when their very own Prolog system would be up and running.

David Warren presented the WAM in a memorable talk at UC Berkeley in October 1983. It was a talk full of mystery, and I remember being amazed at how *append/3* was compiled into WAM instructions. The sense of mystery was enhanced by the strange and melodious names of the instructions: put, get, and unify, variable and value, execute and proceed, try, retry, and trust.

The WAM is simple on the outside (a small, clean instruction set) and complex on the inside (the instructions do complex things). This simultaneously helped and hindered implementation technology. Because it is complex on the inside, for a long time many people used it “as-is” and were content with its level of performance. Because it is simple on the outside, it was a perfect environment for extensions. After a few years, people were extending the WAM left and right. Papers on yet another WAM extension for a new logic language were (and are) very common.

The quickest way to get an implementation of a new logic language is to write an interpreter in Prolog. In the past, the quickest way to get an *efficient* implementation was usually to extend the WAM. Nowadays, it is often better to compile the language into an existing implementation. For example, the QD-Janus system [39] is a sequential implementation of Janus (a flat committed-choice language) on top of SICStus Prolog (see Section 3.1.9). Performance is reasonable partly because SICStus provides efficient support for coroutines.

If the language is sufficiently different from Prolog, then it is better to design a new abstract machine. For example, the λ Prolog language [99] was implemented with MALI [19]. λ Prolog generalizes Prolog with predicate and function variables and typed λ -terms, while keeping the familiar operational and least fixpoint semantics. MALI is a general-purpose memory management library that has been optimized for logic programming systems.

1.2. *Organization of the Survey*

The survey is divided into four parts. The first part (Section 2) gives an overview from the viewpoint of implementation technology. The second part (Section 3) gives an overview from the viewpoint of the systems (both software and hardware) that were responsible for particular developments. The vantage points of the two parts are complementary, and there is some overlap in the developments that are discussed. The third part (Section 4) summarizes the evolution of Prolog performance from the perspective of the Warren benchmarks. The fourth part (Section 5) extrapolates current implementation trends into the future. Finally, Section 6 recapitulates the main developments and gives a conclusion.

A large number of Prolog systems have been developed. The subset that are mentioned in this survey were singled out because they are popular (e.g., SICStus Prolog), because they are a good example of a particular class of systems (e.g., CHIP for constraint languages), or because they are especially innovative (e.g., Parma). They all have implementations on Unix workstations. I have done my best to contact everyone who has made a significant contribution. There are Prologs that exist only on other platforms, e.g., on PCs (Arity,

LPA, Delphia) and on Lisp machines (LMI, Symbolics). There is relatively little publicly available information about these systems and, therefore, I do not cover them in this article.

2. THE TECHNOLOGICAL VIEW

This section gives an overview of Prolog implementation technology. Section 2.1 gives a brief history of the pre-WAM days (before 1983) and presents the main principle of Prolog compilation. Section 2.2 presents and justifies the WAM as Warren originally defined it. Section 2.3 explores a few of the myriad systems it has engendered. Section 2.4 highlights recent developments that break through its performance barrier. Section 2.5 presents some promising execution models different from the WAM.

Prolog systems can be divided into two categories: *structure-sharing* or *structure-copying*. The idea of structure sharing is owing to Boyer and Moore [18]. Structure copying was first described by Bruynooghe [20, 21]. The distinction is based on how compound terms are represented. In a structure-sharing representation, all compound terms are represented as a pair of pointers (called a *molecule*): one pointer to an array containing the values of the term's variables, and another pointer to a representation of the term's nonvariable part (the *skeleton*). In a structure-copying representation, all compound terms are represented as record structures with one word identifying the main functor followed by an array of words giving its arguments. It is faster to create terms in a structure-sharing representation. It is faster to unify terms in a structure-copying representation. Memory usage of both techniques is similar in practice. Early systems were mostly structure-sharing. Modern systems are mostly structure-copying. The latter includes WAM-based systems and all systems discussed in this survey, except when explicitly stated otherwise.

2.1. Before the Golden Age

The insight that deduction could be used as computation was developed in the 1960s through the work of Cordell Green and others. Attempts to make this practical failed until the conception of the Prolog language by Alain Colmerauer and Robert Kowalski in the early 1970s. It is hard to imagine the leap of faith this required back then: to consider a logical description of a problem as a program that could be executed efficiently. The early history is presented in [32], and interested readers should look there for more detail.

The work on Prolog was preceded by the Absys system. Absys (from *Aberdeen System*) was designed and implemented at the University of Aberdeen in 1967. This system was an implementation of pure Prolog [46]. For reasons that are not clear but that are probably cultural, it did not become widespread.

Several systems were developed by Colmerauer's group. The first system was an interpreter written in Algol-W by Philippe Roussel in 1972. It served to give users enough programming experience so that a much refined second system could be built. The second system was a structure-sharing interpreter written in Fortran in 1973 by Gérard Battani, Henri Meloni, and René Bazzoli, under the supervision of Roussel and Colmerauer. This system's operational semantics and its built-ins are essentially the same as in modern Prolog systems, except for the *setof/3* and *bagof/3* built-ins which were introduced by David Warren in 1980 [163]. The system had reasonable performance and was very influential in convincing people that programming in logic was a viable idea.

In particular, David Warren from the University of Edinburgh was convinced. He wrote the Warplan program during his two month stay in Marseilles in 1974 [30]. Warplan is a

general problem solver that searches for a plan (a list of actions) that transforms an initial state to a goal state. Back in Edinburgh and thinking about a dissertation topic, Warren was intrigued by the idea of building a compiler for Prolog. An added push for this idea was the fact that the parser for the interpreter was written in Prolog itself and, hence, was very slow. It took about a second to parse each clause and users were beginning to complain.

By 1977, Warren had developed DEC-10 Prolog, the first Prolog compiler [160]. This landmark system was built with the help of Fernando Pereira and Luis Pereira.¹ It is structure sharing and supports mode declarations. It was competitive in performance to Lisp systems of the day and was for many years the highest performance Prolog system. Its syntax and semantics became the de facto standard, the “Edinburgh standard.” The 1980 version of this system had a heap garbage collector and last call optimization (see Section 2.2.4) [161]. It was the first system to have either. An attempt to commercialize this system failed because of the demise of the DEC-10/20 machines and because of bureaucratic problems with the British government, which controlled the rights of all software developed with public funds.

The main principle in compiling Prolog is to simplify each occurrence of one of its basic operations (namely, unification and backtracking). This principle underlies every Prolog compiler. Compiling Prolog is possible because this simplification is possible very often. For example, unification is often used purely as a parameter passing mechanism. Most cases of this are easily detected and compiled into efficient code.

It is remarkable that this principle has continued to hold to the present day. It is valid for WAM-based systems, native code systems, and systems that use global analysis. In the WAM the simplification is done statically (at compile time) and locally [77]. The simplification can also be done dynamically (with run-time tests) and globally. An example of dynamic simplification is clause selection (see Section 2.4.3). Examples of global simplification are global analysis (see Sections 2.4.5 and 2.4.6) and the two-stream unification algorithm (see Section 2.4.2). The latter compiles the unification of a complete term as a whole, instead of compiling each functor separately like the WAM.

An important early system is the C-Prolog interpreter, which was developed at Edinburgh in 1982 by Fernando Pereira, Luis Damas, and Lawrence Byrd. It is based on EMAS Prolog, a system completed in 1980 by Luis Damas. C-Prolog was one of the best interpreters, and is still a very usable system. It did much to create a Prolog programming community and to establish the Edinburgh standard. It is cheap, robust, portable (it is written in C), and fast enough for real programs.

There were several compiled systems that bridged the gap between the DEC-10 compiler (1977–1980) and the WAM (1983) [17, 28]. They include Prolog-X and NIP (new implementation of Prolog). David Bowen, Lawrence Byrd, William Clocksin, and Fernando Pereira at Edinburgh were the main contributors in this work. These systems miss some of the WAM’s good optimizations: separate choice points and environments, argument passing in registers instead of on the stack, and clause selection (indexing). David Warren left Edinburgh for SRI in 1981. The WAM design was an outcome of his own explorations and was not influenced by this work.

2.2. *The Warren Abstract Machine (WAM)*

By 1983, Warren had developed the WAM, a structure-copying execution model for Prolog that has become the de facto standard implementation technique [164]. The WAM defines a

¹They are not related.

high-level instruction set that maps closely to Prolog source code. This section summarizes and justifies the original WAM. In particular, the many optimizations of the WAM are given a uniform justification. This section assumes a basic knowledge of how Prolog executes [84, 114, 131] and of how imperative languages are compiled [3].

For several years, Warren's report was the sole source of information on the WAM, and its terse style gave the WAM an aura of inscrutability. Many people learned it by osmosis, gradually absorbing its meaning. Nowadays, there are texts that give lucid explanations of the WAM and WAM-like systems [4, 84].

There are two main approaches to efficient Prolog implementation: emulated code and native code. Emulated code compiles to an abstract machine and is interpreted at run time. Native code compiles to the target machine and is executed directly. Native code tends to be faster and emulated code tends to be more compact. With care, both approaches are equally portable (see Section 5.1). The original WAM is designed with an emulated implementation in mind. For example, its unification instructions are more suited to emulated code (see Section 3.1.4). The two-stream unification algorithm of Section 2.4.2 is more suited to native code.

This section is divided into five parts. Section 2.2.1 situates the WAM in relationship to Prolog and imperative languages. Section 2.2.2 describes its data structures and memory organization. Section 2.2.3 presents its instruction set. Section 2.2.4 gives a detailed classification of its optimizations, most of which are based on a single principle: to minimize memory usage. Section 2.2.5 presents a simple scheme for compiling Prolog to WAM.

2.2.1. THE RELATIONSHIP OF THE WAM TO PROLOG AND IMPERATIVE LANGUAGES.

The execution of Prolog is a natural generalization of the execution of imperative languages (see Figure 2.1). It can be summarized as:

Prolog = imperative language
 + unification
 + backtracking.

As in imperative languages, control flow is left to right within a clause. The goals in a clause body are called like procedures. A goal corresponds to a predicate. When a goal is called, the clauses in the predicate's definition are chosen in textual order from top to bottom. Backtracking is chronological, that is, failure goes back to the most recently made choice and tries the next clause. Hence, Prolog is a somewhat limited realization of logic programming, but in practice its trade-offs are good enough that a logical and efficient programming style is possible [112].

The WAM mirrors Prolog closely, both in how the program executes and in how the program is compiled:

WAM = sequential control (call/return/jump instructions)
 + unification (get/put/unify instructions)
 + backtracking (try/retry/trust instructions)
 + optimizations (to use as little memory as possible).

The WAM has a stack-based structure, just as imperative language execution models. It has call and return instructions and environment (local frame) management instructions. It is extended with instructions to perform unification and instructions to perform backtracking.

Prolog	Imperative language
set of clauses	----- program
predicate; set of clauses with same name and arity	----- procedure definition; nondeterministic case statement
clause; axiom	----- one branch of a nondeterministic case statement; if statement; series of procedure calls
goal invocation	----- procedure call
unification	----- parameter passing; assignment; dynamic memory allocation
backtracking	----- conditional branching; iteration; continuation passing
logical variable	----- pointer manipulation
recursion	----- iteration

FIGURE 2.1. Correspondence between logical and imperative concepts.

These form the core of the WAM. Most of the optimizations that extend this are intended to reduce memory use.

Prolog as executed by the WAM defines a close mapping between the terminology of logic and that of an imperative language (see Figure 2.1). Predicates correspond to procedures. Procedures always have a case statement as the first part of their definition. Clauses correspond to the branches of this case statement. Variables are scoped locally to a clause.²Goals in a clause correspond to calls. Unification corresponds to parameter passing and assignment. Other features do not map directly: backtracking, the single-assignment nature, and the modification of control flow with the cut operation.

The WAM is a good intermediate language in the sense that writing a Prolog-to-WAM compiler and a WAM emulator are both straightforward tasks. A compiler and emulator can be built without a deep understanding of the internals of Prolog or the WAM. This has led to a proliferation of WAM-based systems and WAM extensions for other logic languages (a few examples are given in Section 2.3).

2.2.2. DATA STRUCTURES AND MEMORY ORGANIZATION. Prolog is a dynamically typed language, that is, variables may contain objects of any type at run time. Hence, it must be possible to determine the type of an object at run time by inspection.³In the WAM, terms are represented as tagged words: a word contains a tag field and a value field.⁴The tag field contains the type of the term (atom, number, list, or structure). The value field is used for different purposes in different types: it contains the value of integers, the address of unbound variables and compound terms (lists and structures), and it ensures that each atom has a value different from all other atoms. Unbound variables are implemented as

²Global variables and self-modifying code are possible with the *assert/1* and *retract/1* built-ins. These built-ins are potentially nonlogical and certainly inefficient and, hence, should be infrequent.

³Unless the compiler can determine that a more efficient representation is possible.

⁴See [52] for an exhaustive presentation of tagging schemes.

TABLE 2.1. Internal State of the WAM

P	Program counter
CP	Continuation pointer (top of return stack)
E	Current environment pointer (in local stack)
B	Most recent choice point (in local stack)
A	Top of local stack
TR	Top of trail
H	Top of heap
HB	Heap backtrack point (in heap)
S	Structure pointer (in heap)
Mode	mode flag (read or write)
A1, A2, ...	Argument registers
X1, X2, ...	Temporary variables

self-referential pointers, that is, they point to themselves. When two variables are unified, one of them is modified to point to the other.⁵Therefore it may be necessary to follow a chain of pointers to access a variable’s value. This is called *dereferencing* the variable.

Table 2.1 gives the internal state of the WAM (stored in registers). The purpose of most registers is straightforward. The HB register caches the value of H stored in the most recent choice point. The S register is used during unification of compound terms (with arguments): it points to an argument being unified. All arguments can be accessed one by one by successively incrementing S. Some instructions have different behaviors during read and write mode unification; the mode flag is used to distinguish between them (see Section 2.2.3). In the original WAM, the mode flag is implicit (it is encoded in the program counter).

The external state (stored in memory) is divided into six logical areas (see Figure 2.2): two stacks for the data objects, one stack (the PDL) to support unification, one stack (the trail) to support the interaction of unification and backtracking, one area as code space, and one area as a symbol table:

- **The global stack or heap.** This stack holds lists and structures, the compound data terms of Prolog.
- **The local stack.** This stack holds environments and choice points. Environments (also known as local frames or activation records) contain variables local to a clause. Choice points encapsulate the execution state for backtracking, that is, they are continuations. A variant model, the split-stack, uses separate stacks for environments and choice points. There is no significant performance difference between the split-stack and the merged-stack models. The merged-stack model uses more memory if choice points are created (see last call optimization in Section 2.2.4). Because of backtracking, control may return to a clause whose environment is deep inside the stack.
- **The trail.** This stack is used to save locations of bound variables that have to be unbound on backtracking. Saving variables is called *trailing*, and restoring them

⁵More precisely, variable-variable unification can be implemented with a Union-Find algorithm [90]. With this algorithm, unifying n variables requires $O(n\alpha(n))$ time, where $\alpha(n)$ is the inverse Ackermann function.

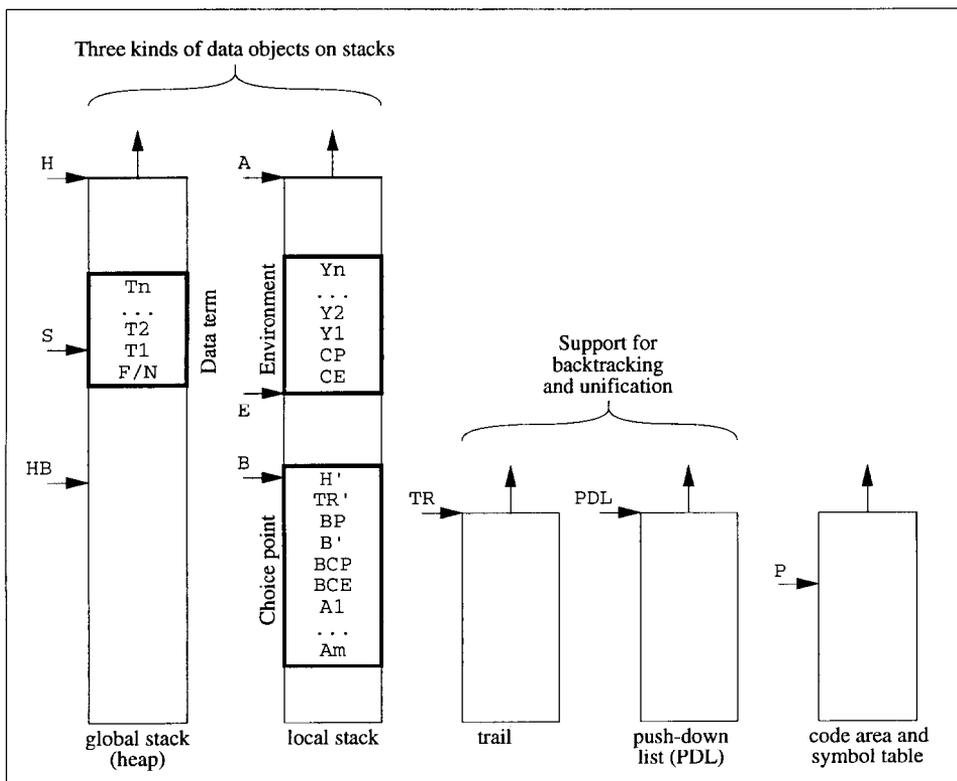


FIGURE 2.2. External state of the WAM.

to unbound is called *detrailing*. Not all variables that are bound have to be trailed. A variable must only be trailed if it continues to exist on backtracking, i.e., if its location on the global or local stack is older than the top of this stack stored in the most recent choice point. This is called the *trail condition*. Performing it is called the *trail check*.

- **The push-down list (PDL).** This stack is used as a scratch pad during the unification of nested compound terms. Often it does not exist as a separate stack, e.g., the local stack is used instead.
- **The code area.** This area holds the compiled code of a program. It is not recovered on backtracking.
- **The symbol table.** This area is not mentioned in the original article on the WAM. It holds various kinds of information about the symbols (atoms and structure names) used in the program. It is not recovered on backtracking. It contains the mapping between the internal representation of symbols and their print names, information about operator declarations, and various system-dependent information related to the state of the system and the external world. Because creating a new entry is relatively expensive, symbol table memory is most often not recovered on backtracking. It may be garbage collected. Systems that manipulate arbitrary numbers of new atoms (e.g., systems with a database interface) require it to be garbage collected to be practical.

```

append([], L, L).
append([X|L1], L2, [X|L3]) :- append(L1, L2, L3).

```

FIGURE 2.3. The Prolog code for *append/3*.

It is possible to vary the organization of the memory areas somewhat without changing anything substantial about the execution. For example, some systems have a single data area (sometimes called the *firm heap*) that combines the code area and symbol table.

2.2.3. THE INSTRUCTION SET. The WAM instruction set, along with a brief description of what each instruction does, is summarized in Table 2.2. Unification of a variable with a data term known at compile-time is decomposed into instructions to handle the functor and arguments separately (see Figures 2.3 and 2.4). There are no **unify_list** and **unify_structure** instructions; they are left out because they can be implemented using the existing instructions. The **switch_on_constant** and **switch_on_structure** instructions fall through if A_1 is not in the hash table. The original WAM report does not talk about the **cut** operation, which removes all choice points created since entering the current predicate. Implementations of cut are presented in [4, 84]. A variable stored in the current environment (pointed to by E) is denoted by Y_i . A variable stored in a register is denoted by X_i or A_i . A register used to pass arguments is denoted by A_i . A register used only internally to a clause is denoted by X_i . The notation V_i is shorthand for X_i or Y_i . The notation R_i is shorthand for X_i or A_i .

A useful optimization is the variable/value annotation. Instructions annotated with “variable” assume that their argument has not yet been initialized (i.e., it is the first occurrence of the variable in the clause). In this case, the unification operation is simplified. For example, the **get_variable** X_2, A_1 instruction unifies X_2 with A_1 . Since X_2 has not yet been initialized, the unification reduces to a move. Instructions annotated with “value” assume that their argument has been initialized (i.e., all later occurrences of the variable). In this case, full unification is done.

Figures 2.3 and 2.4 give the Prolog source code and the compiled WAM code for the predicate *append/3*. The mapping between Prolog and WAM instructions is straightforward (see Section 2.2.5). The **switch** instruction jumps to the correct clause or set of clauses depending on the type of the first argument. This implements first-argument selection (indexing). The choice point (**try**) instructions link a set of clauses together. The **get** instructions unify with the head arguments. The **unify** instructions unify with the arguments of structures.

The same instruction sequence is used to take apart an existing structure (read mode) or to build a new structure (write mode). The decision which mode to use is made in the **get** instructions, which set the mode flag. For example, if **get_list** A_i sees an unbound variable argument, it sets the flag to write mode. If it sees a list argument, it sets the flag to read mode. If it sees any other type, it fails, i.e., it backtracks by restoring the state from the most recent choice point. The **unify** instructions have different behavior in read and write mode. The **get** instructions initialize the S register and the **unify** instructions increment the S register.

Choice point handling (backtracking) is done by the **try** instructions. The **try_me_else** L instruction creates a choice point, i.e., it saves all the machine registers on the local stack. It is compiled just before the code for first clause in a predicate. It continues execution with the next instruction and backtracks to label L . The **try** L instruction is identical to

TABLE 2.2. The Complete WAM Instruction Set

Loading Argument Registers (just before a call)	
put_variable V_n, R_i	Create a new variable, put in V_n and R_i
put_value V_n, R_i	Move V_n to R_i
put_constant C, R_i	Move the constant C to R_i
put_nil R_i	Move the constant nil to R_i
put_structure $F/N, R_i$	Create the functor F/N , put in R_i
put_list R_i	Create a list pointer, put in R_i
Unifying with Registers (head unification)	
get_variable V_n, R_i	Move R_i to V_n
get_value V_n, R_i	Unify V_n with R_i
get_constant C, R_i	Unify the constant C with R_i
get_nil R_i	Unify the constant nil with R_i
get_structure $F/N, R_i$	Unify the functor F/N with R_i
get_list R_i	Unify a list pointer with R_i
Unifying with Structure Arguments (head unification)	
unify_variable V_n	Move next structure argument to V_n
unify_value V_n	Unify V_n with next structure argument
unify_constant C	Unify the constant C with next structure argument
unify_nil	Unify the constant nil with next structure argument
unify_void N	Skip next N structure arguments
Managing Unsafe Variables (an optimization; see Section 2.2.4)	
put_unsafe_value V_n, R_i	Move V_n to R_i and globalize
unify_local_value V_n	Unify V_n with next structure argument and globalize
Procedural Control	
call P, N	Call predicate P , trim environment size to N
execute P	Jump to predicate P
proceed	Return
allocate	Create an environment
deallocate	Remove an environment
Selecting a Clause (conditional branching)	
switch_on_term V, C, L, S	Four-way jump on type of A_1
switch_on_constant N, T	Hashed jump (size N table at T) on constant in A_1
switch_on_structure N, T	Hashed jump (size N table at T) on structure in A_1
Backtracking (choice point management)	
try_me_else L	Create choice point to L , then fall through
retry_me_else L	Change retry address to L , then fall through
trust_me_else <i>fail</i>	Remove top-most choice point, then fall through
try L	Create choice point, then jump to L
retry L	Change retry address, then jump to L
trust L	Remove top-most choice point, then jump to L

<i>append/3</i> :	switch_on_term $V_1, C_1, C_2, fail$	Jump if variable, constant, list, structure.
V_1 :	try_me_else V_2	Create choice point if A_1 is variable.
C_1 :	get_nil A_1	Unify A_1 with nil.
	get_value A_2, A_3	Unify A_2 and A_3 .
	proceed	Return to caller.
V_2 :	trust_me_else <i>fail</i>	Remove choice point.
C_2 :	get_list A_1	Start unification of list in A_1 .
	unify_variable X_4	Unify head: move head into X_4 .
	unify_variable A_1	Unify tail: move tail into A_1 .
	get_list A_3	Start unification of list in A_3 .
	unify_value X_4	Unify head: unify head with X_4 .
	unify_variable A_3	Unify tail: move tail into A_3 .
	execute <i>append/3</i>	Jump to beginning (last call optimization).

FIGURE 2.4. The WAM code for *append/3*.

try_me_else L , except that it continues execution at L and backtracks to the next instruction. The **retry_me_else** L instruction modifies a choice point that already exists by changing the address that it jumps to on backtracking. It is compiled with clauses after the first but not including the last. The **trust_me_else** *fail* instruction removes the topmost choice point from the stack. It is compiled with the last clause in a predicate.

2.2.4. OPTIMIZATIONS TO MINIMIZE MEMORY USAGE. The core of the WAM is straightforward. What makes it subtle is the optimizations that are done. Most of the optimizations are based on one principle: to minimize memory usage. The optimizations are explained in terms of the following classification of memory, from least to most costly to allocate, deallocate, and reuse:

- **Registers** (arguments, temporary variables). These are available at any time without overhead.
- **Short-lived memory** (environments on the local stack). This is memory that is recovered on forward execution, on backtracking, and on garbage collection.
- **Long-lived memory** (choice points on the local stack, data terms on the heap). This is memory that is recovered only on backtracking and on garbage collection. For example, choice points removed by cut are recovered in this way.
- **Permanent memory** (the code area and symbol table). This is memory that is recovered only on garbage collection.

With this ranking, the optimizations can be explained as follows:

- **Prefer registers over memory.** Three optimizations are based on this rule.

- **Argument passing.** Pass all arguments in registers. This is important because Prolog is procedure-call intensive. The most efficient way to do loops is through recursion. Backtracking can express loops as well, but it is less efficient.
 - **The return address.** Inside a procedure, the return address of the immediate caller is stored in the CP register. This optimization is closely related to the *leaf routine* calling protocol done in imperative language compilers. The return addresses of all calls higher up in the call tree are stored in environments.
 - **Temporary variables.** Temporary variables are variables that do not need to survive a call. That is, their lifetimes do not cross a call. Therefore, they may be kept in registers. This definition of temporary variables simplifies and slightly generalizes Warren's original definition.
- **Prefer short-lived memory over long-lived memory.** There are three separate optimizations in this category.
 - **Permanent variables.** Permanent variables are variables that need to survive a call. They may not be kept in registers, but must be stored in memory. They are given a slot in the environment. This is to make it easy to deallocate their memory if they are no longer needed after exiting the predicate (see Unsafe variables, below).
 - **Environment trimming (last call optimization).** Environments are stored on the local stack and recovered on forward execution just before the last call in a procedure body. This optimization is known as the *tail recursion optimization* or, more accurately, the *last call optimization*. This is based on the observation that an environment's space does not need to exist after the last call, since no further computation is done in the environment. The space can be recovered before entering the last call instead of after it returns. Because execution will never return to the procedure, this call may be converted into a jump. For recursive predicates, this converts recursion into iteration, since the jump is to the first instruction of the predicate. The WAM generalizes this optimization to be done gradually during execution of a clause: the environment size is reduced ("trimmed") after each call in the clause body, so that only the information needed for the rest of the clause is stored. Trimming increases the amount of memory that is recovered by garbage collection. If a more recent choice point exists, then trimming creates a hole in the local stack.
 - **Unsafe variables.** A variable whose lifetime crosses a call must be allocated an unbound variable cell in memory (i.e., in an environment or on the heap). If it is sure that the unbound variable will be bound before exiting the clause, then the space for the cell will not be referenced after exiting the clause. In that case the cell may be allocated in the environment and recovered with environment trimming. In the other case, if one is not sure that the unbound variable will be bound, there is a space-time trade-off. Either the variable can always be created on the heap or the variable can be created on the environment and just before trimming the environment, a test can be done to see whether the variable has been bound, and if the test is false, it can then be moved to the heap. The WAM has chosen the second alternative, and the variable being tested is referred to as an "unsafe variable." The following measurements have been done of this trade-off [81]. Lindholm measured the increase of peak heap usage for a set of programs including Chat-80 [162] and the Quintus test suite and compiler. He found that the first alternative increases peak heap usage

by 50–100% for Quintus (see Section 3.1.4). Because this leads to increased garbage collection and stack shifting, he concluded that unsafe variables are a good idea. Krall measured the increase of peak heap usage on a series of small- and medium-size programs for the VAM (see Section 2.5.1), which stores all unbound variables on the heap. He measured increases of from 4 to 26%, with an average of 15%. These measurements are for a system that is significantly faster than Quintus and that does not create unbound variables for calls. Hence, the numbers cannot be compared directly with Quintus. Because unsafe variables impose a run-time overhead (two comparisons instead of one for the trail check and run-time tests for globalizing variables), he concluded that unsafe variables are a bad idea.

- **Prefer long-lived memory over permanent memory.** Data objects (variables and compound terms) disappear on backtracking and, hence, all allocated memory for them may be put on the heap. In a typical implementation this is not quite true. The symbol table and code area are not put on the heap, because their modifications (i.e., newly interned atoms and functors and asserted clauses) are permanent.

Because of these optimizations the WAM is extremely memory efficient. For programs with sufficient backtracking, a garbage collector is not necessary.

2.2.5. HOW TO COMPILE PROLOG TO WAM. Compiling Prolog to WAM is straightforward because there is almost a bijective mapping between lexical tokens in the Prolog source code and WAM instructions. Figure 2.5 gives a scheme for compiling Prolog to WAM. For simplicity, the figure omits register allocation and peephole optimization. This compilation scheme generates suboptimal code. One can improve it by generating **switch** instructions to avoid choice point creation in some cases. For more information on WAM compilation, see [115, 149].

The clauses of predicate $p/3$ are compiled into blocks of code that are linked together with **try** instructions. Each block consists of a sequence of **get** instructions to do the unification of the head arguments, followed by a sequence of **put** instructions to set up the arguments for each goal in the body, and a **call** instruction to execute the goal. The block is surrounded by **allocate** and **deallocate** instructions to create an environment for permanent variables. The last call is converted into a jump (an **execute** instruction) because of the last call optimization (see Section 2.2.4).

2.3. WAM Extensions for Other Logic Languages

Many WAM variants have been developed for new logic languages, new computation models, and parallel systems. This section presents three significant examples: the CHIP constraint system, the `clp(FD)` constraint system, and the SLG-WAM Prolog with memoization. Information on other language families (such as concurrent languages and database languages) is given in other articles in this issue.

2.3.1. CHIP. CHIP (constraint handling in Prolog) [2, 41] is a constraint logic language developed at ECRC (see Section 3.1.7 for more information on ECRC). The system has been commercialized by Cosytec to solve industrial optimization problems. CHIP is the first compiled constraint language. In addition to Prolog terms, it adds three computation domains: finite domains, boolean terms, and linear inequalities over rationals. The CHIP

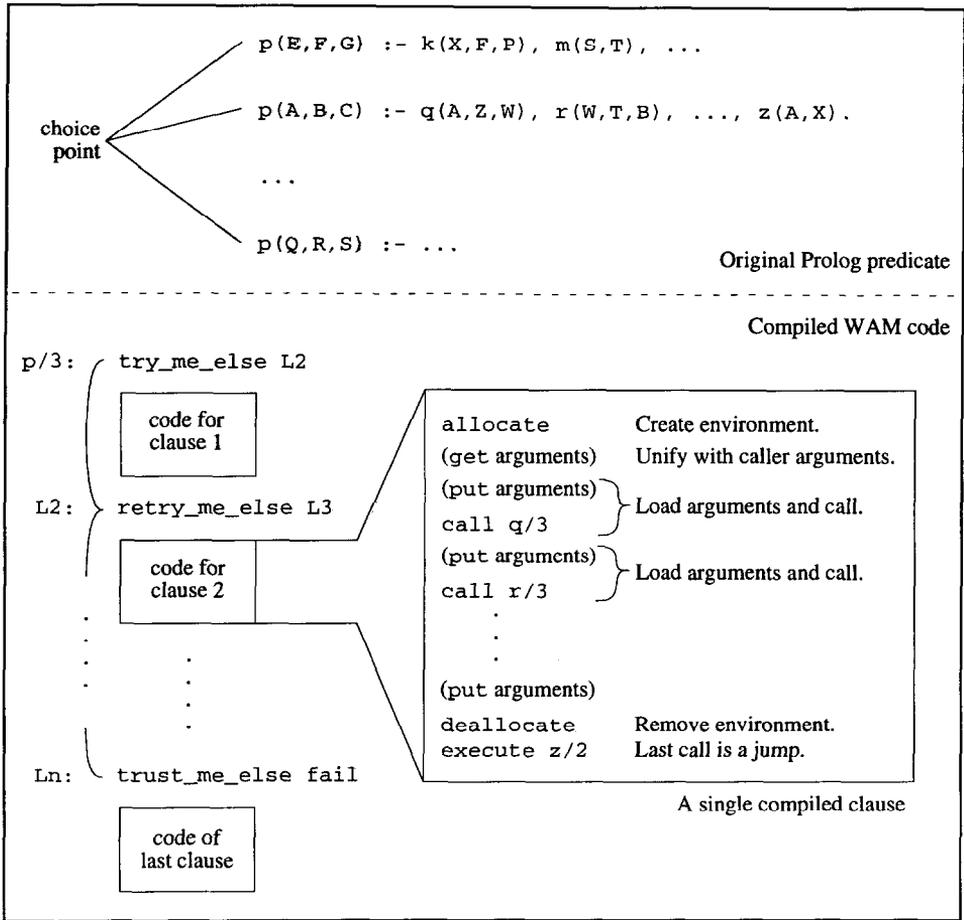


FIGURE 2.5. How to compile Prolog to WAM.

compiler is built on top of the SEPIA WAM-based Prolog compiler. The system contains a tight interface between the WAM kernel and the constraint solvers. It extends the WAM to the C-WAM (C for constraint). The C-WAM has new data structures and over 100 new instructions. Many instructions exist to solve commonly occurring constraints quickly.

Measurements of early versions of CHIP showed that an immense amount of trailing was being done, so much so that large programs quickly ran out of memory. This happened because the trailing was done with the WAM's trail condition (see Section 2.2.2). This condition is appropriate for equality constraints, which are implemented by unification in the WAM. For more complex constraints, the condition is wasteful because a variable's value is often modified several times between two choice points. The CHIP system solves this problem by introducing a different trail condition called *time-stamping* [1]. Each data term is marked with an identifier of the choice point segment it is created in (see Section 2.3.1). Trailing is only necessary if the current choice point segment is different from the segment stored in the term. Time-stamping is an essential technique for any practical constraint solver.

2.3.2. *CLP(FD)*. The *clp(FD)* system [29, 40] is a finite domain solver integrated into a WAM emulator. It was built by Daniel Diaz and Philippe Codognot at INRIA Rocquencourt, France. It uses a “glass box” approach. Instead of considering a constraint solver as a black box (in the manner of CHIP), a set of primitive operations is added that allows the constraint solver to be programmed in the user language. The resulting system outperforms the hard-wired finite domain solver of CHIP.

In *clp(FD)*, a single primitive constraint is added to the system, namely, the range constraint $X \text{ in } R$, where X is a domain variable and R is a range (e.g., $1..10$). Instead of just using constant ranges, the idea is to introduce *indexical* ranges, that is, ranges of the form $f(Y)..g(Y)$ or $h(Y)$, where $f(Y)$, $g(Y)$, and $h(Y)$ are functions of the domain variable Y . These functions are provided by the system and do local propagation. For example, the system provides the constraints $X \text{ in } \min(Y)..max(Y)$ and $X \text{ in } \text{dom}(Y)$ with the obvious meanings. Arithmetic constraints such as $X=Y+Z$ and boolean constraints such as $X=Y$ and Z can be implemented with the indexical range constraints.

Indexical range constraints are smoothly integrated into the WAM by providing support for domain variables and suspension queues for the various indexical functions [40]. The time-stamping technique of CHIP is used to reduce trailing.

2.3.3. *SLG-WAM*. Memoization is a technique that caches already computed answers to a predicate. When combined with Prolog’s resolution mechanism, this results in an execution model that can do both top-down and bottom-up execution. For certain algorithms, this means that simple logical definitions can run with a lower order of complexity. For example, the recursive definition of the Fibonacci function runs in linear time rather than exponential time. More realistic examples are parsing and dynamic programming.

One realization of this mechanism is OLDT resolution (ordered linear resolution of definite clauses with tabulation) [132]. A recent generalization, SLG resolution [27], handles negation as well. This has been implemented in an abstract machine, the SLG-WAM (previously called the OLDT-WAM), and realized in the XSB system (see Section 3.1.8). The current implementation executes untabled code with less than 10% overhead relative to the WAM as implemented in XSB, and is much faster than deductive database systems [133]. An important source of overhead is the complex trail: it is a linked list whose elements contain the address and old contents of a cell.

2.4. *Beyond the WAM: Evolutionary Developments*

The WAM was a large step forward for the execution efficiency of Prolog. From the viewpoint of theorem proving, Prolog is extremely fast. However, there is still a large gap between the efficiency of the WAM and that of imperative language implementations. After the novelty of the WAM wore off and people started using Prolog for standard programming tasks, the gap became apparent and people started to optimize their systems. This section discusses the gap and some of the clever ideas that have been developed to close it.

This section is divided into five parts. Section 2.4.1 sets the stage by listing some of the limits to Prolog performance and their causes. Section 2.4.2 presents the two-stream unification algorithm. Section 2.4.3 presents some of the problems of the WAM’s clause selection and their solutions. Section 2.4.4 discusses native code compilation without a WAM intermediate stage. Section 2.4.5 summarizes what has been done in global analysis for Prolog compilation. Section 2.4.6 shows how types are used to improve the compiled code for unification in the Aquarius system.

2.4.1. CHINKS IN THE ARMOR.

- WAM instructions are too coarse-grained to perform much optimization. For example, many WAM instructions perform an implicit dereference operation, even if the compiler can determine that such an operation is unnecessary in a particular case. In practice, dereference chains are short: dynamic measurements on real programs show that two-thirds are of length zero (no memory reference is required), one-third are of length 1, and <1% are of length 2 or greater [146]. Despite these statistics, dereferencing is expensive. For example, Aquarius on the VLSI-BAM, a high-performance system with hardware support, spends 9% of its total execution time doing dereferencing [153].
- The majority of predicates written by human programmers are intended to give at most one solution, i.e., they are deterministic. These predicates are, in effect, case statements, yet they are too often compiled in an inefficient manner using the full generality of backtracking (which implies saving the machine state and repeated failure and state restoration). The WAM's first-argument selection is inadequate to compile these predicates efficiently (see Section 2.4.3). Measurements of Prolog applications support this assertion:
 - Tick shows that shallow backtracking (backtracking from clause to clause within a single predicate) dominates even for well-written deterministic programs. Choice point references constitute about half (45–60%) of all data references [144].
 - Touati and Despain show that at least 40% of all choice point and fail operations can be removed through optimization [146].
- The single-assignment nature of Prolog (i.e., a variable can only be assigned one value in forward execution) needs to be handled well in the implementation. In a straightforward implementation it is time-consuming to modify large data structures incrementally, because the programmer may use copying of terms to represent incremental changes, and the implementation will not optimize this copying away. This problem, also known as the *copy avoidance problem*, is a special case of the general problem of efficiently representing state modification in logic. It is impossible to use large data structures with the same efficiency as in procedural languages unless the compiler is able to introduce destructive assignment (overwriting of memory locations) in the implementation. Section 5.1 gives suggestions on how to get around this problem.
- Prolog has dynamic typing (variables may contain values of any type) and dynamic memory allocation (all data objects are allocated at run time). Both of these cost time during execution. For efficiency, they should be compiled statically where possible.
- Programming style has a great effect on a program's efficiency. Prolog programming is at a high level of abstraction, so it hides many details of the implementation from the programmer, making it difficult to improve efficiency when it is important to do so. For example, adding a single cut (a language feature that increases the determinism of the program by removing choice points) can make the difference between a program that runs fast and one that thrashes. This is possible even if the cut does not change the operational semantics of the program. The thrashing behavior is caused by a pileup of choice points during deterministic (forward) execution. Because the choice points encapsulate execution states that remain accessible through potential

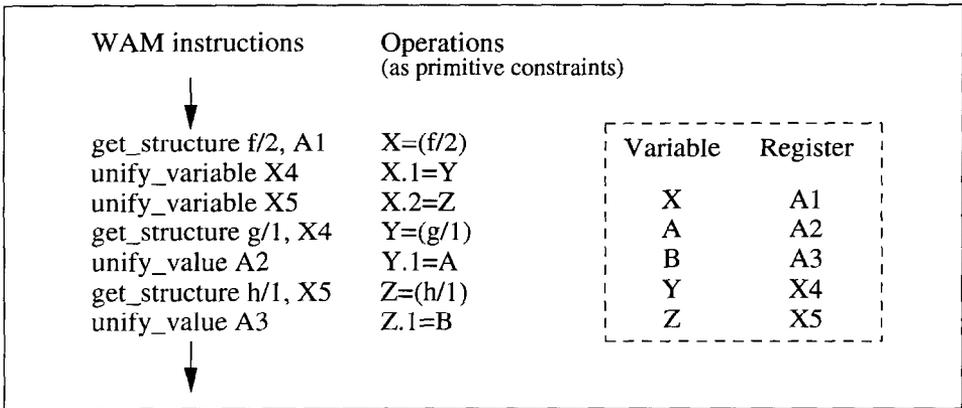


FIGURE 2.6. WAM compilation of the unification $X=f(g(A),h(B))$.

backtracking, their memory is not recovered by garbage collection.

- The apparent need for architectural support. So-called general-purpose architectures are, in fact, optimized for imperative languages and number crunching. To run Prolog equally well, either the compiler must do more work, or conceivably the architecture should be modified. Some experiments have been done with architectures optimized for Prolog (among others, the PSI-II, KCM, and VLSI-BAM; see Section 3.2), but the true architectural needs of Prolog are a moving target. They depend on the execution model and the sophistication of the compiler. As better compilers have been developed, the perceived architectural needs of Prolog have been getting smaller and smaller. One need likely to stay for a long time is a fast memory system. Prolog’s dynamic nature requires a lot of pointer chasing, and there are no compilation techniques on the horizon that are likely to reduce this (see Section 5.1).

2.4.2. HOW TO COMPILE UNIFICATION: THE TWO-STREAM ALGORITHM. This section presents the two-stream unification algorithm, an elegant scheme for compiling unification that is much more efficient than the WAM for native code implementation. Rough measurements comparing unification times of the VLSI-PLM (a microcoded WAM) and the VLSI-BAM (see Section 3.2.3) show a speedup factor of 2 to 3 [154]. This algorithm was independently reinvented at least four times by different people at about the same time: Mohamed Amraoui at the Université de Rennes I [8], André Mariën and Bart Demoen at BIM and KUL [85, 87], Kent Boortz at SICS [16], and Micha Meier at ECRC [93]. Write mode propagation was discussed earlier by Turk [147].

Figures 2.6 and 2.8 show how the unification $X=f(g(A),h(B))$ is compiled in the WAM and by the two-stream algorithm. The actions of the instructions are represented as *primitive constraints* of two kinds: functor constraints (such as $X=(f/2)$) and argument constraints (such as $X.1=Y$). Functor and argument constraints correspond to the WAM’s **get** and **unify** instructions. An important advantage of the primitive constraint representation over the WAM is that the constraints may be executed in any order. In addition to providing a pithy conceptual description of the WAM, primitive constraints are useful in compiling more advanced logic languages [6, 83, 116].

The WAM compiles unification as a single sequence of instructions (see Figure 2.6).

This has several problems:

- **Write mode is not propagated to subterms.** For example, the unification $X=f(g(a))$ is compiled as $X=f(T)$, $T=g(a)$. These two unifications are compiled independently. If X is unbound, the fact that T is created as an unbound variable in the first unification is not propagated to the second unification. This means a superfluous dereference, a superfluous trail check, and a superfluous binding.
- **Instructions have modes.** All instructions have two modes of execution, read mode and write mode. The current mode is stored in a global mode flag, which is set in `get_list` and `get_structure` instructions and tested in all `unify` instructions. Some implementations (e.g., the intended implementation of the original WAM report, and Quintus) encode the mode flag in the program counter, which avoids the testing overhead.
- **Poor translation to native code.** The straightforward method for generating native code is to macro-expand the WAM instructions. This means that the read and write mode parts are interleaved, which results in many jumps. This is less of a problem on a microcoded machine since microcode jumps are often free (the destination address is part of the microword).

The key insight is that unification should be compiled into *two* instruction streams, one for read mode and one for write mode, with conditional jumps between them. In this way, one avoids superfluous operations while keeping a linear code size. The practical problems that remain are how to configure the instructions so they work correctly despite being jumped to from different places and how to minimize bookkeeping overhead for the jumps.

Figure 2.7 illustrates a technique to execute any subsequence of a main instruction sequence with very little overhead. The idea is to give the main sequence an identifier (say, the integer 0) and the subsequence a different identifier (say, the integer 1). Then a single conditional jump is all that is required. If the subsequence is noncontiguous, then a single conditional jump is needed per contiguous segment of the subsequence to hop to the next segment. If more than one subsequence has to be selected, then a unique identifier is needed for each one. The subsequences may be overlapping.

With the idea of selective execution in mind, arrange the primitive constraints of the term according to a depth-first traversal of the term (Figure 2.8). The resulting sequence satisfies the property that each subterm corresponds to a contiguous sequence of instructions. This is all one needs to implement the algorithm. At run time, unification follows the read mode stream and selectively executes contiguous parts of the write mode stream for subterms to be created.

A reduction of bookkeeping overhead is possible based on a second property of the sequence. Nested terms correspond to nested sequences of instructions. Number each subterm with an integer representing its *nesting level*. This gives the number of levels deep the subterm is nested in the term. With this numbering, an adjacent sequence of conditional jumps back to the read mode stream can be collapsed into a single conditional jump (changing the condition from “=” to “≤”). In Figure 2.8, the two conditional jumps **if $S=1$ jump Lz'** and **if $S=0$ jump Lx'** can be rewritten as the single jump **if $S≥0$ jump Lz'** . To collapse the most jumps, reorder the arguments of all subterms to unify the most complex subterms last.

The advantages of this algorithm are:

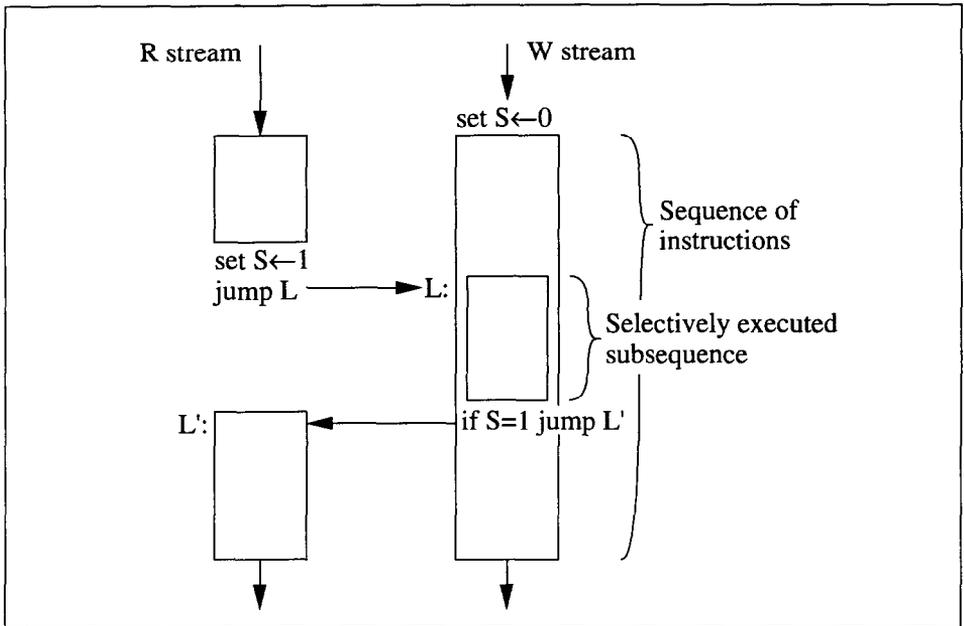


FIGURE 2.7. Executing a particular subsequence with low overhead.

- **Low overhead.** The bookkeeping overhead is a small constant factor. The only bookkeeping is the set of jumps and register moves needed to manage the selective execution of subsequences. This is small compared to the work done in the primitive constraints. There is no explicit mode flag.
- **Downward propagation of write mode.** Write mode of a term is propagated at compile time to all its subterms. There are no superfluous dereferences, trail checks, or bindings.
- **Upward propagation of read mode.** Read mode of a term is propagated at compile time to its siblings and ancestors.
- **Linear code size.** This contrasts with the algorithm of [151], which expands all cases without any sharing. That algorithm has zero bookkeeping overhead, but exponential code sizes occur in practice.
- **Efficient expansion to native code.** The number of instructions generated is about double that of the WAM, but the instructions themselves have less than half the complexity. The primitive constraints of Figure 2.8 are expanded differently in the read mode and write mode streams. Essentially, the WAM instructions' insides have been made visible and arranged in an efficient order. There are no jumps inside the primitive constraints, but only between them, and then only when it is necessary to choose between read and write mode.

2.4.3. HOW TO COMPILE BACKTRACKING: CLAUSE SELECTION ALGORITHMS. The WAM supports first-argument selection. It has instructions that can choose clauses based on the main functor of the first argument. If all of a predicate's clauses contain different main functors, then a hash table can be constructed, and calling the predicate will avoid a choice point creation if the first argument is a nonvariable. In the general case, predicates can be

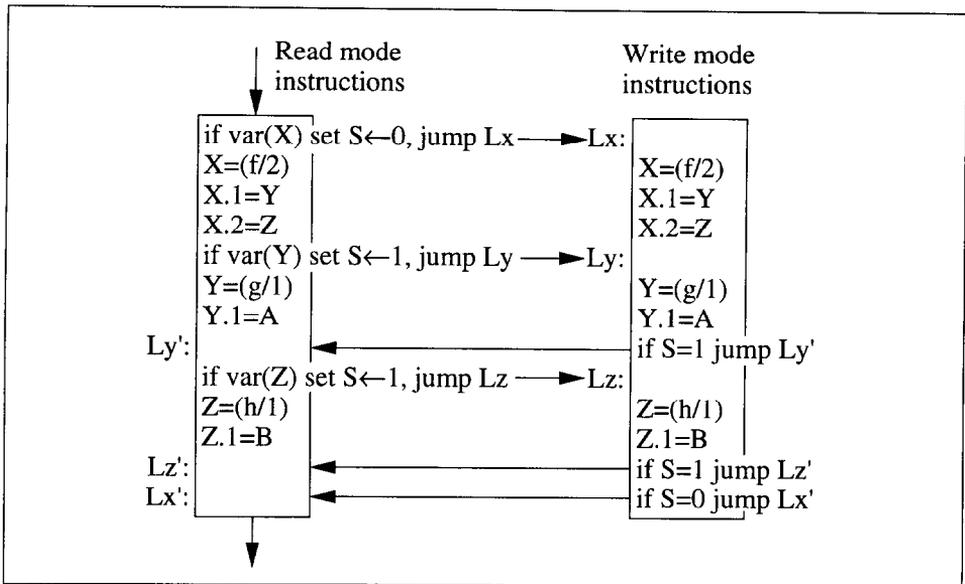


FIGURE 2.8. Two-stream compilation of the unification $X=f(g(A),h(B))$.

compiled to create at most one choice point between entry and the execution of the first clause [23, 149]. The original WAM report describes a two-level indexing scheme which creates up to two choice points [164].

Many programs cannot profit from first-argument selection. For example, selection may depend on more than one argument. The following example is extracted from an actual program. The first two arguments are integer inputs, the third is an output (all numbers are in base 2):

```
get_relop(2'001, 2'001, 2'000).
get_relop(2'001, 2'010, 2'011).
get_relop(2'001, 2'011, 2'000).
... 33 more clauses ...
```

The second example is a predicate in which selection depends on arithmetic comparisons instead of just unification:

```
max(A, B, C) :- A ≤ B, C=B.
max(A, B, C) :- A > B, C=A.
```

In general, selection is possible if the compiler can determine that only a subset of the clauses in the definition can possibly succeed, given some particular argument types at the call. An appropriate definition of *type* is given in Section 2.4.5. In such a case, the compiler should generate code to recognize these argument types and should try only the clauses that can possibly succeed. This should avoid all useless choice point creations. The resulting code size should be linear in the size of the program. Choice points, if needed, should be created incrementally. Performance degradation should be gradual if the compiler cannot

determine all information. There is no published algorithm that satisfies all these conditions.

Many techniques exist that solve part of the clause-selection problem. Some compilers accept programmer declarations of the types or modes of predicate arguments, and can use this information to improve selection. In the literature, several algorithms have been given that do better clause selection than the WAM. These are based on creating a selection tree or graph, that is, a series of tests that determine which subset of clauses to try given particular arguments (e.g., [168]). Naively generating a selection tree may result in exponential code size for predicates encountered in real-world programs. The following algorithms are noteworthy:

- Van Roy et al. [150] present a compilation algorithm that generates a naive selection tree and compiles clauses in a novel way. It creates choice points incrementally. It compiles clauses with four entry points, depending on whether or not there are alternative clauses and whether or not a previously executed clause has created a choice point. The algorithm was not implemented.
- Carlsson [24] has implemented a restricted version of the above algorithm in SICStus Prolog. Meier [91] has done a similar implementation in KCM-SEPIA. Choice point creation is split into two parts. The **try** and **try.me.else** instructions are modified to create a partial choice point that only contains P and TR. A new instruction, **neck**, is added. If a choice point exists when **neck** is executed, then it fills in the remaining registers. Two entry points are created for each clause: one when there are alternative clauses and one where there are none. A **neck** instruction is only included in the first case. In SICStus, this algorithm results in a performance improvement of 7 to 15% for four large programs, at a cost of a 5 to 10% increase in code size.
- Hickey and Mudambi [61] present compilation algorithms to generate a tree of tests and to minimize work done in backtracking. One of their selection algorithms results in a tree that has a quadratic worst-case size. They improve choice point management. The **try** instruction only stores registers needed in clauses after the first clause. The **retry** and **trust** instructions restore only those registers needed in the clause and remove the registers not needed in subsequent clauses. The latter operation lets the garbage collector recover more memory. Variants of this technique were independently invented earlier by Turk [147] and later by Van Roy [154]. The technique has not yet been quantitatively evaluated.
- Klinger [70, 71] presents a compilation algorithm that generates a directed acyclic graph of tests (a “decision graph”). The algorithm is extended by Korsloot and Tick for nondeterminate (“don’t know”) predicates [73]. The graph has two important properties. First, it never does worse than first-argument selection. Second, it has linear size in the number of clauses. This follows from the property that each clause corresponds to a unique path through the graph. Linear size is essential when compiling predicates with large numbers of clauses.
- The Aquarius system [154, 155] produces a selection graph for disjunctions containing three kinds of tests: unifications, type tests, and arithmetic comparisons. It uses heuristics to decide which tests to do first and whether to use linear search or hashing for table lookup. The nodes in the graph partition the tests occurring in the predicate. Each node corresponds to a subset of these tests. Unifications are only used as tests if it can be deduced from the predicate’s type information that they will be executed in read mode. The “type enrichment” transformation adds type information to a predicate that lacks it. The performance of the resulting code

is, therefore, always at least as good as first-argument selection. The “factoring” transformation allows the system to take advantage of tests on variables inside of terms, by performing the term unification once for all occurrences of the term. The problem with Aquarius selection is similar to that of the naive selection tree: if too much type information is given, then the selection graph may become too large.

- The Parma system [141] uses techniques similar to Aquarius. It produces efficient indexing code for the same three kinds of tests. To improve the clause selection, it uses transformations analogous to type enrichment and factoring. It uses optimal binary search for table lookup. Taylor’s dissertation discusses how to choose between linear search, binary search, jump tables, and hashing.

2.4.4. NATIVE CODE COMPILATION One way to improve the performance of a WAM-based system is to add instructions. For example, instructions can be added to do efficient arithmetic and to index on multiple arguments. Common instruction sequences can be collapsed into single instructions. This is quick to implement, but it is inherently a short-term solution. As the number of instructions increases, the system becomes increasingly unwieldy.

The main insight in speeding up Prolog execution is to represent the code in terms of simple instructions. The first published experiments using this idea were done in 1986 by Komatsu *et al.* [72, 136] at IBM Japan. Their compilation is done in three steps. The first step is to compile Prolog into a WAM-like intermediate code. In the second step this code is translated into a directed graph. The graph is optimized using rewrite rules. In the final step, the result is translated into PL.8 intermediate code and compiled with an optimizing compiler. For several small programs, the system demonstrated a fourfold performance improvement using mode hints given by the programmer.

Around 1988, Andrew Taylor and I independently set about building full systems (Parma and Aquarius) that would compile directly to a simple instruction set, using global analysis to provide information for optimizations. Both Parma and Aquarius bypass WAM instructions entirely during compilation. We were confident that the fine granularity of the instruction set would allow us to express all optimizations. Taylor presented results for his Parma system in two important papers [139, 140]. The first paper presents and practically evaluates a global analysis that reduces the need for dereferencing and trailing. The second paper presents performance results for Parma on a MIPS processor. The first results for Aquarius were presented in [63], which describes the VLSI-BAM processor and its simulated performance. A second paper measures the effectiveness of global analysis in Aquarius [153]. Both the Parma and Aquarius systems vastly outperform existing implementations. They prove the effectiveness of compiling directly to a low-level instruction set using global analysis to help optimize the code.

An important idea in both systems is *uninitialized variables*, which are essentially “unboxed” unbound variables (see Section 5.1). An uninitialized variable is defined to be an unbound variable that is *unaliased*, i.e., it has only one pointer to it. In this case it can be represented more efficiently. Beer [12] first proposed this idea after he noticed that most unbound variables in the WAM are bound soon afterwards, for example, output arguments of predicates. WAM variables are created as self-referential pointers in memory, and need to be dereferenced and trailed before being bound. This is time-consuming. Beer represents variables as pointers to memory words that have not been initialized. He introduces several new tags for these variables and keeps track of them at run time. Creation is simpler and they do not have to be dereferenced or trailed. Binding reduces to a single store operation. In Parma and Aquarius, these variables are derived by analysis at compile time. They use

the same tag as other variables.

Aquarius supports a further specialization of Beer's uninitialized variable type: the "uninitialized register" type. This idea is owing to Bruce Holmer. This type represents outputs that are passed in registers. No memory is allocated for uninitialized registers, unlike standard uninitialized variables. In addition to being much faster than uninitialized variables, this reduces the space advantage of unsafe variables. It allows Aquarius to run simple recursive integer functions faster than popular implementations of C [155].⁶In principle, all uninitialized variables can be transformed into uninitialized registers. In practice, to avoid losing last call optimization (see Section 2.2.4), only a subset is transformed [154]. The trade-off with last call optimization has not yet been quantitatively studied.

Figure 2.9 shows the Aquarius intermediate codes (kernel Prolog and BAM code) and the SPARC code generated for *append/3* in naive reverse. See Figures 2.3 and 2.4 for the Prolog source code and WAM code. Kernel Prolog is Prolog without syntactic sugar and extended with efficient conditionals, arithmetic, and cut.

The BAM (Berkeley abstract machine) is an execution model with a memory organization similar to the WAM. The BAM defines a load-store instruction set with tagged addressing modes, pragmas, and five Prolog-specific instructions (dereference, trail, general unification, choice point manipulation, and environment manipulation). Pragmas are not executable but give information that improves the translation to machine code.

In the SPARC code, tags are represented as the low two bits of a 32-bit word. This is a common representation that has low overhead for integer arithmetic and pointer dereferencing [52]. The tag of a pointer is always known at compile time (it is put in a pragma). When following a pointer, the tag is subtracted off at zero cost with the SPARC's register+displacement addressing mode. The compiler derives the following types for *append/3* (see the next section for an explanation of what they mean):

```
:- mode((append(A,B,C) :-
        ground(A), rderef(A),
        ground(B), rderef(B),
        uninit(C))).
```

This type generalizes the DEC-10 mode:

```
:- mode append(++ , ++ , -).
```

which states that the first two arguments are ground (they contain no unbound variables) and the last argument is an unbound variable.

2.4.5. GLOBAL ANALYSIS. Global analysis of logic programs is used to derive information to improve program execution. Both type and control information can be derived and used to increase speed and reduce code size. The analysis algorithms studied so far are all instances of a general method called *abstract interpretation* [34, 35, 68]. The idea is to execute the program over a simpler domain. If a small set of conditions is satisfied,

⁶I posted this result to the Internet newsgroup comp.lang.prolog in February 1991, with the comment: "Don't believe it any more that there is an inherent performance loss when using logic programming." There was a barrage of responses, ranging from the incredulous (and incorrect) comment, "Obviously, he's comparing apples and oranges, since the system must be doing memoization," to the encouraging, "That's telling 'em Peter."

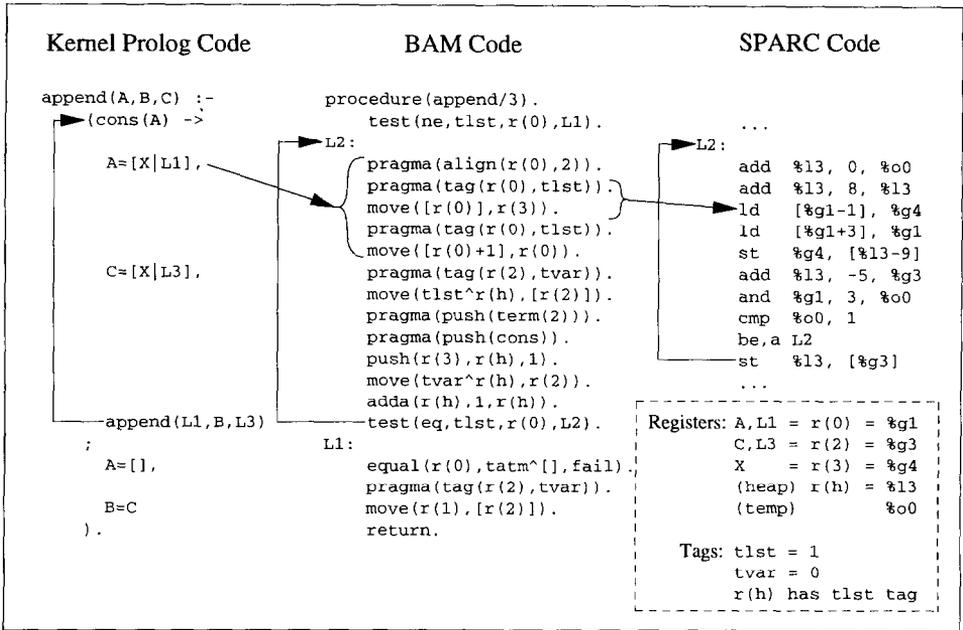


FIGURE 2.9. The Aquarius SPARC code for *append/3* in naive reverse.

this execution terminates and its results provide a correct approximation of information about the original program. Le Charlier *et al.* [78, 79] have performed an extensive study of abstract interpretation algorithms and domains and their effectiveness in deriving types. Getzinger [50] has recently presented an extensive taxonomy of analysis domains and studied their effects on execution time and code size.

Since Mellish’s early work in 1981 and 1985 [95, 97], global analysis has been considered useful for Prolog implementation. This section summarizes the work that has been done in making analysis part of a real system. From the viewpoint of the implementor, a *type* denotes any information known about a variable’s value at run time. A *mode* is a restricted type that indicates whether the variable is used as an input (nonvariable) or an output (unbound variable). The concept of mode is rather limited. Most types inferred by a good analyzer go beyond it. Useful types include argument values, compound structures, dependencies between variables, and operational information such as length of dereference chains (see also Sections 2.4.6 and 5.1).

In 1982, Naish [105] performed an experiment with automatically generated control for MU-Prolog. The MU-Prolog interpreter supports *wait* declarations. A “wait” declaration defines a set of arguments of a predicate that may not be constructed by a call (i.e., unified in write mode). If a call attempts to construct a term in any of these arguments, then it delays until the argument is sufficiently instantiated so that no construction is done (i.e., it is unified in read mode). This provides a form of coroutines. The automatic generation of declarations is based on a simple heuristic: to delay rather than guess one of an infinite number of bindings.⁷ A “wait” declaration is inserted for each recursive call that does not

⁷This heuristic is closely related to the “Andorra principle” [33, 55]. The main difference is that it is applied at analysis time rather than at run time.

progress in its traversal of a data structure. This algorithm was implemented and tested on some small examples. It significantly reduces the programmer's burden in managing control. The algorithm does not always help. If the clause head is as general as the recursive call, then no "wait" declaration is generated, even though one might be necessary.

A later system, NU-Prolog, supports *when* declarations. These are both more expressive and easier to compile into efficient code (see Section 3.1.3). A "when" declaration is a pattern containing a term with optional variables and a nested conjunction and/or disjunction of nonvariable and ground tests on these variables. Variables may not occur more than once in the term. A "when" declaration is true if unification between it and the goal succeeds and does not bind any variables in the goal. This is called *one-way* unification or *matching*. NU-Prolog contains an analyzer that derives "when" declarations.

In 1988, Warren, Hermenegildo, and Debray [60, 165] did the first measurements of the practicality of global analysis in logic programming. They measured two systems, MA³, the MCC And-parallel analyzer and annotator, and Ms, an experimental analysis scheme developed for SB-Prolog. The paper concludes that both data flow analyzers are effective in deriving types and do not increase compilation time by too much.

In 1989, Mariën et al. [86] performed an interesting experiment in which several small Prolog predicates (recursive list operations) were hand compiled with four levels of optimization based on information derivable from a global analysis. The levels progressively include unbound variable and ground modes, recursively defined types, lengths of dereference chains, and liveness information for compile-time garbage collection. Execution time measurements show that each analysis level significantly improves speed over the previous level. This experiment shows that a simple analysis can achieve good results on small programs.

Despite this experimental evidence, there was until 1993 no generally available sequential Prolog system that did global analysis, and since 1988 only a few research systems doing analysis. Why is this? I think the most important reason is that other areas of system development were considered more important. Commercial systems worked on improving their development environments: source-level debugging, a proper foreign language interface, and useful libraries. Research systems worked in other areas such as language design and parallelism. A second reason may be that the structure of the WAM (high-level compact instructions) does not lend itself well to the optimizations that analysis supports. A whole new instruction set would be needed, and the development effort involved may have seemed prohibitive given the existing investment in the WAM. A third reason is that analysis was (erroneously) considered impractical.

Currently, there are at least seven systems that do global analysis of logic programs:

- Ms, an experimental analyzer for SB-Prolog written by Debray [60, 165]. It derives ground and nonvariable types.
- MA³, the analyzer for &-Prolog written by Hermenegildo and Warren [60, 165]. The analyzer derives variable sharing (aliasing) and groundness information. This information is used to eliminate run-time checks in the And-parallel execution of Prolog. This was the first practical application of abstract interpretation to logic programs, that is, this system both derives information and uses it for optimization. PLAI, the successor to MA³, subsumes it and has been extended to analyze programs in constraint languages [49] and languages with delaying [89].
- The FCP(:,?) compiler (flat concurrent Prolog with ask and tell guards and read-only variables), written by Klinger, has a global analysis phase [71].
- The Parma system, written by Taylor, is an implementation of Prolog with global

System	Speedup Factor		Code Size Reduction	
	Small	Medium	Small	Medium
Aquarius	1.5	1.2	2.2	1.8
Parma	3.0	2.1	2.9	2.0

TABLE 2.3. Effectiveness of Analysis for Small and Medium-size Programs.

analysis targeted to the MIPS processor [141].

- The Aquarius system is an implementation of Prolog with global analysis targeted to the VLSI-BAM processor and various general-purpose processors [58, 154]. An extensive study of improved analyzers and their integration in the Aquarius system is given in [50].
- The MU-Prolog analyzer generates “wait” declarations for coroutining [105]. Its improved NU-Prolog version generates “when” declarations. See earlier in this section for more information.
- The IBM Prolog analyzer. It determines whether choice points have been created or destroyed during execution of a predicate, and whether there are pointers into the local stack. This improves the handling of unbound variables and the management of environments. There is no published information on this analyzer. The analyzer has been available since 1989 (see Section 3.1.6).

Of these systems, five were developed for sequential Prolog (the MU-Prolog and IBM Prolog analyzers, Ms, Parma, and Aquarius) and two for parallel systems (MA^3 and the $FCP(:,?)$ analyzer). Three (MA^3 , Parma, and Aquarius) have been integrated into Prolog systems and their effects on performance evaluated and published [60, 141, 154]. The analysis domains of Aquarius and Parma are shown in Figure 2.10. For both analyzers the analysis time is linear in program size and performance is adequate. Four analyzers (MA^3 , $FCP(:,?)$, Aquarius, and IBM Prolog) are robust enough for day-to-day programming.

The effect of the Aquarius and Parma analyzers on speed and code size is shown in Table 2.3. The “Small” column refers to a standard set of small benchmarks (between 10 and 100 lines). The “Medium” column refers to a standard set of medium-size benchmarks (between 100 and 1000 lines). These benchmarks are well known in the Prolog programming community [157]. They do tasks for which Prolog is well suited and are written in a good programming style. The numbers are taken from [141, 153, 154]. The numbers can be significantly improved by tuning the programs to take advantage of the analyzers’ strengths. The following two paragraphs give numbers for the medium-size benchmarks.

The Aquarius analyzer finds uninitialized, ground, nonvariable, and recursively dereferenced types in 23, 21, 10, and 17% of predicate arguments, respectively, and 56% of predicate arguments have types.⁸ One-third of the uninitialized types are uninitialized register types, so about 1/12 of all predicate arguments are outputs passed in registers. On the VLSI-BAM this means a reduction of dereferencing from 11 to 9% of execution time and a reduction of trailing from 2.3 to 1.3% of execution time.

The Parma analyzer’s domain has been split into parts and their effects on performance

⁸Arguments can have more than one type.

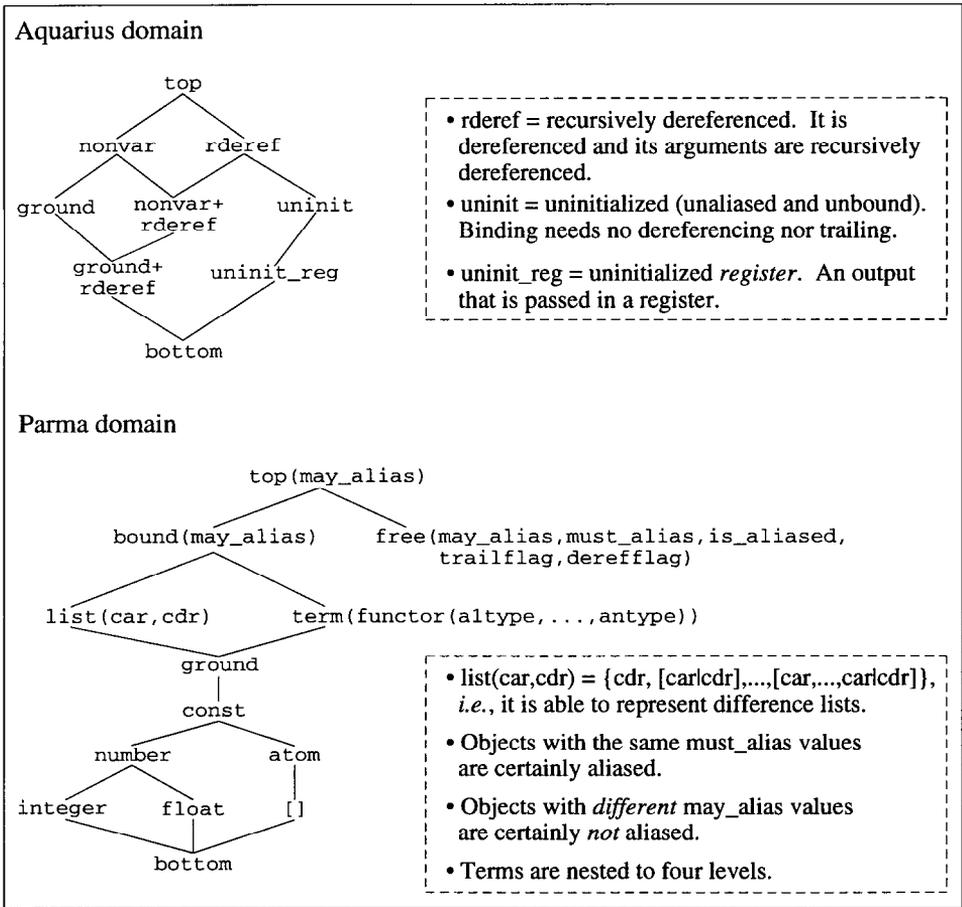


FIGURE 2.10. Analysis domains for the Aquarius and Parma systems.

measured separately. Performance is improved through dereference chain analysis by 14%, trailing analysis by 8%, structure/list analysis by 22%, and uninitialized variables by 12%. The combined benefit of two analysis features is usually not their product, since features may compete with or enhance each other. For example, uninitialized variables do not need to be trailed, and this fact will often also be determined by the trailing analysis.

Two conclusions can be drawn by studying the effects of analysis in the Parma and Aquarius systems. Analysis results in both a code size reduction and a performance improvement. The effects of analysis on code size and performance are fundamentally different. Derived types allow both *tests* and the *code* that handles other types to be removed. Tests are usually fast. The code to handle all possible outcomes of the tests can be very large. For Aquarius, the code size reduction is greater than the performance improvement. This is partly due to the lack of structure and list types in the Aquarius domain, which means that run-time type tests are still needed. For Parma, the code size reduction is about the same as the performance improvement.

A second conclusion can be drawn regarding the kinds of types that are most useful in the compiler. Deriving types that have a logical meaning is not sufficient. Performance

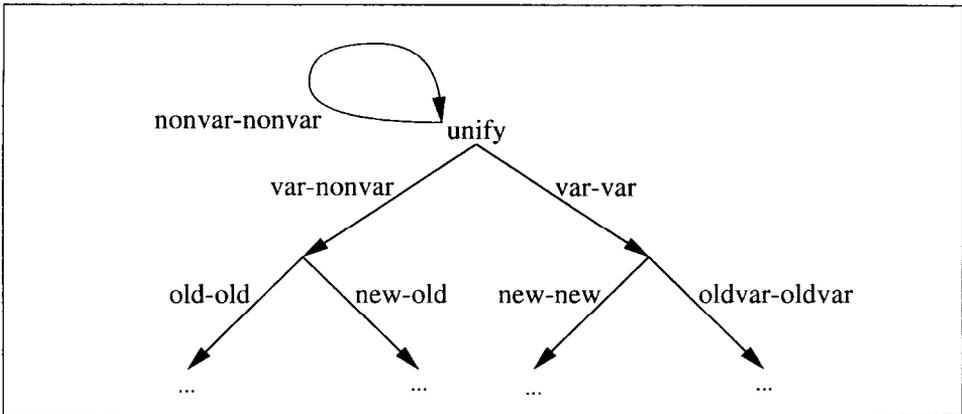


FIGURE 2.11. Case analysis in compiling unification.

increases significantly when the analysis is able to derive types that only have an operational meaning, such as dereference (reference chains), trailing, and aliasing-related types (uninitialized variables).

2.4.6. USING TYPES WHEN COMPILING UNIFICATION. It is as hard to *use* analysis in the compiler as it is to *do* analysis, yet very little has been published in this area. This section shows how unification is compiled in the Aquarius system to take maximum advantage of the types known at compile time. The code generated by the two-stream algorithm of Section 2.4.2 handles the general case when no types are known. If types are known, then compiling unification becomes a large case analysis.⁹ Even after common cases are factored out, the number of cases remains large. Figure 2.11 gives a much simplified view of the top two levels of the case analysis done in Aquarius.

Table 2.4 gives details of the case analysis done in Aquarius for the compilation of the unification $X=Y$ with type T . The compiler attempts to use all possible type information to simplify the code. A general unify instruction is only generated once (in `oldvar_oldvar`), namely, when unifying two initialized variables for which nothing is known. For simplicity, the table omits the generation of dereference and trail instructions, the handling of uninitialized memory and uninitialized register variables, the updating of type information when variables are bound, the generation of pragmas, and various less important optimizations. See Section 2.4.4 for more information.

The variable T denotes the type information known at the start of the unification. The implication $(T \Rightarrow \text{ground}(X))$ succeeds if T implies that X is bound to a ground term at run time. The conditions $\text{var}(X)$ and $(T \Rightarrow \text{var}(X))$ are very different: the first tests whether X is a variable at compile time, and the second tests whether X is a variable at run time. The condition $\text{new}(X)$ succeeds if this is the first occurrence of X in the clause or if X is uninitialized. The condition $\text{old}(X)$ is the negation of $\text{new}(X)$, that is, it is true for all initialized variables. The function $\text{atomic_value}(T, X)$ succeeds if T implies that X is an atomic term whose value is known at compile time. The function returns this atomic term. For example, if T is $(X==a)$, then the function returns the atom “a.”

⁹Compiling a goal invocation is also a large case analysis [154].

Definition of Routines		
Name	Condition	Actions
unify(X,Y)	var(X),var(Y)	var_var(X,Y)
	var(X),nonvar(Y)	var_nonvar(X,Y)
	nonvar(X),var(Y)	var_nonvar(Y,X)
	nonvar(X),nonvar(Y)	nonvar_nonvar(X,Y)
nonvar_nonvar(X,Y)		For all arguments X_i, Y_i : unify(X_i, Y_i)
var_nonvar(X,Y)	$T \Rightarrow \text{new}(X)$	new_old(X,Y)
	$T \Rightarrow \text{ground}(X)$	old_old(X,Y)
	otherwise	old_old(X,Y) (with depth limiting)
var_var(X,Y)	$T \Rightarrow (\text{old}(X), \text{old}(Y))$	oldvar_oldvar(X,Y)
	$T \Rightarrow (\text{old}(X), \text{new}(Y))$	Generate store instruction
	$T \Rightarrow (\text{new}(X), \text{old}(Y))$	Generate store instruction
	$T \Rightarrow (\text{new}(X), \text{new}(Y))$	new_new(X,Y)
new_new(X,Y)		Generate store and move instructions
new_old(X,Y)	compound(Y)	write_sequence(X,Y)
	atomic(Y)	Generate store instruction
	var(Y)	var_var(X,Y)
old_old(X,Y)	compound(Y), ($T \Rightarrow \text{nonvar}(X)$)	Test Y's type, then old_old_read(X,Y)
	atomic(Y), ($T \Rightarrow \text{nonvar}(X)$)	old_old_read(X,Y)
	nonvar(Y), ($T \Rightarrow \text{var}(X)$)	old_old_write(X,Y)
	compound(Y)	Generate switch, old_old_read(X,Y), old_old_write(X,Y)
	atomic(Y)	Generate unify_atomic instruction
oldvar_oldvar(X,Y)	var(Y)	var_var(X,Y)
	$A = \text{atomic.value}(T,X)$	unify(Y,A)
	$A = \text{atomic.value}(T,Y)$	unify(X,A)
	$T \Rightarrow (\text{atomic}(X), \text{atomic}(Y))$	Generate comparison instruction
	$T \Rightarrow (\text{var}(X), \text{nonvar}(Y))$	Generate store instruction
	$T \Rightarrow (\text{nonvar}(X), \text{var}(Y))$	Generate store instruction
otherwise	Generate general unify instruction	
old_old_write(X,Y)	compound(Y)	write_sequence(X,Y)
	atomic(Y)	Generate store instruction
old_old_read(X,Y)	compound(Y)	Test Y's functor, then for all arguments X_i, Y_i : old_old(X_i, Y_i)
	atomic(Y)	Generate comparison instruction
write_sequence(X,Y)		Generate instructions to create compound term Y in X

2.5. Beyond the WAM: Radically Different Execution Models

Some recent developments in Prolog implementation are based on novel models of execution very different from the WAM. The Vienna abstract machine (VAM) is based on partial evaluation of each call. The BinProlog system is based on explicit passing of success continuations.

2.5.1. THE VIENNA ABSTRACT MACHINE (VAM). The VAM is an execution model developed by Krall and Neumerkel at the Technische Universität Wien, Vienna, Austria [75]. The insight of the VAM is that the WAM's separation of argument setup from argument unification is wasteful. In the WAM, all of a predicate's arguments are built *before* the predicate is called. The VAM does argument setup and argument unification at the same time. During the call it combines the operations of argument setup and unification into a single operation that does the minimal work necessary. This results in considerable savings in many cases. For example, consider the call $p(X,[a,b,c],Y)$ to the definition $p(A,..,B)$. The second argument $[a,b,c]$ is not created because it is a void variable in the head of the definition. In the WAM, the second argument would be created and then ignored in the definition.

There exist two versions of the VAM: the VAM_{1p} and the VAM_{2p} . The difference is in how the argument traversal is done. In the VAM_{2p} there are two pointers. One points to the caller's arguments and one points to the definition's arguments. The operation to be performed for each argument is obtained by a two-dimensional array lookup depending on the types of the caller argument and the definition argument. This lookup operation can be made extremely fast by a technique similar to direct threaded coding, where the address of the abstract instruction is obtained by adding two offsets. In the VAM_{1p} there is a single pointer that points to compiled code representing the caller-definition pair. The code size for the VAM_{1p} is much greater than for the VAM_{2p} , since the called predicate must be compiled separately for each call. Currently, the VAM_{2p} is a practical implementation, whereas the VAM_{1p} is not because of code size explosion.

2.5.2. BINPROLOG. BinProlog is a high-performance emulator written in C developed by Tarau at the Université de Moncton, Canada [36, 137, 138]. It has two key ideas: transforming clauses to binary clauses and passing success continuations. The resulting instruction set is essentially a simplified subset of the WAM. Implementing Prolog by means of continuations is an old technique. It was used to implement Prolog on Lisp machines and in Pop-11; see, for example, [22, 96, 128]. The technique has recently received a boost by Tarau's highly efficient implementation. Functional languages have more often been implemented by means of continuations. A good example is the Standard ML of New Jersey system, which uses an intermediate representation in which all continuations are explicit ("continuation-passing style") [9].

The idea of BinProlog is to transform each Prolog clause into a binary clause, a clause containing only one body goal. Predicates that are expanded in-line (such as simple built-ins) are not considered as goals. The body goal is given an extra argument, which represents

its success continuation, that is, the sequence of goals to be executed if it completes successfully. This representation has two advantages. First, no environments are needed. Second, the continuations are represented at the source level. For example, the clauses

```
p(X, X),
p(A, B) :- q(A, C), r(C, D), s(D, B),
```

are transformed into

```
p(X, X, Cont) :- call(Cont),
p(A, B, Cont) :- q(A, C, r(C, D, s(D, B, Cont))).
```

Each predicate is given an additional argument, the continuation, and each clause is converted into a binary clause.

With a well chosen data representation, this results in a system that uses very little memory yet compiles and executes very quickly. The technique as currently implemented has two problems. First, the continuations are put on the heap (“long-lived” memory); hence, they do not disappear on forward execution as environments would in the WAM. That is, there is no last call optimization (see Section 2.2.4) if the original clause body contains more than one goal. Second, if the first goal fails, then the creation of the continuation is an overhead that is avoided in the WAM. Both of these problems are less severe than they appear at first glance. The first problem goes away with a suitable garbage collector. A copying collector has execution time proportional to the amount of active memory. A generational mark-and-sweep collector can perform even better in practice [169]. The second problem almost never occurs in real programs.

An important potential use of this technique is as a tool for source transformation in Prolog compilers. By making the continuations of the WAM explicit as data terms, a series of optimizing transformations becomes possible at the source level [109]. After doing the optimizations, a reverse transformation to standard clauses can be done.

3. THE SYSTEMS VIEW

The previous sections have summarized developments from the technical viewpoint, focusing on particular developments and only obliquely mentioning the systems that pioneered them.

This section changes the viewpoint to the systems themselves. It tells the stories of some of the more popular and influential systems, of the people and institutions behind them, and of the particular problems they encountered and how they solved them. The section is divided into two parts. Section 3.1 talks about software systems and Section 3.2 talks about hardware systems.

3.1. *Software Sagas*

Since the development of the WAM in 1983 there have been many software implementations of Prolog. As of this writing, more than 50 systems are listed in the Prolog Resource Guide [69]. The systems discussed in this section are MProlog, IF/Prolog, SNI-Prolog, MU-Prolog, NU-Prolog, Quintus, BIM, IBM Prolog, SEPIA, ECLiPSe, SB-Prolog, XSB, SICStus, and Aquarius.

All systems are substantially compatible with the Edinburgh standard. They have been released to users and used to build applications. Many have served as foundations for implementation experiments. In particular, MU-Prolog, NU-Prolog, SB-Prolog, XSB, SICStus, and Aquarius are delivered with full source code. Quintus, MProlog, IF/Prolog, and SICStus are probably the implementations that have been ported to the largest number of platforms. The most popular systems on workstations today are SICStus and Quintus. C-Prolog was also very popular in the past.

For each system we list its most important contributions to implementation technology. These lists are not exhaustive. Most of the important “firsts” have since been incorporated into many other systems. In some cases, a contribution was developed jointly or spread too fast to identify a particular system as the pioneer. Some of these contributions are mentioned in the first part of the paper. Others are mentioned here. For example, almost all commercial systems support modules. Likewise, almost all commercial systems have a full-featured foreign language interface, and many of them (including Quintus, BIM, IF/Prolog, SNI-Prolog, SICStus, and ECLiPSe) allow Prolog to call C to call Prolog, and so forth, to any level of nesting.

IF/Prolog, SNI-Prolog, IBM Prolog, SEPIA, ECLiPSe, and SICStus support rational tree unification. Rational trees account for term equations which express cycles. For example, the term equation $X=f(X)$ has a solution over rational trees, but does not over finite trees [65, 67].

All of the compiled systems except MProlog and Aquarius are based on the WAM instruction set, but modified and extended to increase performance. MProlog, BIM, IBM Prolog, SEPIA, ECLiPSe, and Aquarius support mode declarations and multiple-argument indexing. The other systems do not support mode declarations. Quintus, NU-Prolog, and XSB provide some support for multiple-argument indexing, and IF/Prolog, SNI-Prolog, and SICStus do not implement it. IBM Prolog, SEPIA, ECLiPSe, and Aquarius index on some conditions other than unification, for example, on arithmetic comparisons and type tests. Quintus, BIM, SEPIA, ECLiPSe, XSB, SB-Prolog, Aquarius, but not SICStus, compile conditionals (if-then-else) deterministically in the special case where the condition is an arithmetic comparison or a type test.

The especially interesting problems of system building are related to scalability. I call these “large program” problems because they tend to occur only when one exercises a system on large programs. They are the main obstacles on the long path between research prototype and production quality system. For each system, we list some of the more interesting of these problems that were encountered. Some of these problems by their nature occur in many systems (e.g., garbage collection bugs). In such cases, the problem is listed only once with a reference to its ubiquity.

3.1.1. MPROLOG. The first commercial Prolog system was MProlog. MProlog was developed in Hungary starting in 1978 at NIMIGUSZI (Computer Center of the Ministry of Heavy Industries) [13, 47].¹⁰ The main developer of MProlog is Péter Szeredi, aided by Zsuzsa Farkas and Péter Köves. MProlog was completed at SZKI (Computer Research and Innovation Center), a computer company set up a few years before. The implementation is based on Warren’s pre-WAM three-stack model of DEC-10 Prolog. The first public demonstration was in 1980 and the first sale was in September 1982.

MProlog is a full-featured structure-sharing system with all Edinburgh built-ins, debug-

¹⁰M for modular or Magyar.

ging, foreign language interface, and sophisticated I/O. It shows that structure sharing is as efficient as structure copying [74]. Its implementation was among the most advanced of its day. Early on, it had a native code compiler, memory recovery on forward execution (including tail recursion optimization), and support for mode declarations (including multiple-argument indexing). It had garbage collection for the symbol table and code area. It did not and does not do garbage collection for the stacks. MProlog is currently a product of IQSOFT, a company formed in 1990 from the Theoretical Lab of SZKI.

3.1.2. IF/PROLOG AND SNI-PROLOG. IF/Prolog was developed at InterFace Computer GmbH, which was founded in 1982 in Munich, Germany. Nothing has been published about the implementation of this system. The following information is owing to Christian Pichler. IF/Prolog was commercialized in 1983. The first release was an interpreter. A WAM-based compiler was released in 1985. The origin of the compiler is an early WAM compiler developed by Pichler [115]. The main developers of IF/Prolog were Preben Folkjær, Christian Reisenauer, and Christian Pichler. Many other people have contributed to this system. Siemens-Nixdorf Informationssysteme AG bought the IF/Prolog sources in 1986. They ported and extended the system, which became SNI-Prolog.

In 1990, SNI-Prolog was completely redesigned from scratch. Pichler went to Siemens-Nixdorf to help in the redesign. The main developers of SNI-Prolog are Reinhard Enders and Christian Pichler. Many other people have contributed to this system. The current system conforms to the ISO Prolog standard [121], supports constraints, is more portable, and has improved system behavior (more flexible interfaces and less memory usage). The design of the new system benefited from the fact that Siemens is one of the shareholders of ECRC. Siemens-Nixdorf bought the rights to IF/Prolog in 1993 after InterFace disappeared. They plan to integrate the best features of IF/Prolog and SNI-Prolog into a single system.

Both systems support rational tree unification. In addition, SNI-Prolog has delaying, indefinite precision rational arithmetic, and constraint solvers for boolean constraints, linear inequalities, and finite domains. It has metaterms, which allow constraint solvers to be written in the language itself (see Section 3.1.7).

Both SNI-Prolog and IF/Prolog have extensive C interoperability. In this regard they can best be compared with Quintus (see Section 3.1.4). They allow redefinition of the C and Prolog top levels and arbitrary calls between Prolog and C to any level of nesting with efficient passing of arbitrary data (including compound terms). They have configurable memory management and garbage collection of all Prolog memory areas. They are designed to interact correctly with the Unix memory system and to support signal handlers.

3.1.3. MU-PROLOG AND NU-PROLOG. MU-Prolog and NU-Prolog were developed at Melbourne University by Naish and his group [105]. Both systems do global analysis to generate delaying declarations (see Section 2.4.5). Neither system does garbage collection.

MU-Prolog is a structure-sharing interpreter. The original version (1.0) was written by John Lloyd in Pascal. Version 2.0 was written by Naish and completed in 1982. Version 2.0 supports delaying, and has a basic module system and transparent database access. Performance is slightly less than C-Prolog.

NU-Prolog is a WAM-based emulator written in C primarily by Jeff Schultz and completed in 1985. It is interesting for its pioneering implementation of logical negation, quantifiers, if-then-else, and inequality, through extensions to the WAM [104]. The delay declarations (“when” declarations) are compiled into decision trees with multiple entry points. This avoids repeating already performed tests on resumption. It results in indexing on multiple arguments in practice. NU-Prolog was the basis for many implementation

experiments, e.g., related to parallelism [106, 113], databases [117], and programming environments [107].

3.1.4. **QUINTUS PROLOG.** Quintus Prolog is probably the best-known commercial Prolog system. Its syntax and semantics have become a de facto standard, for several reasons. It is close to the Edinburgh syntax and is highly compatible with C-Prolog. It was the first widely known commercial system. Several other influential systems (e.g., SICStus Prolog) were designed to be compatible with it. The pending ISO standard for Prolog [121] will most likely be close in syntax and semantics to the current behavior of Quintus.

Quintus Computer Systems was founded in 1984 in Palo Alto, California. It is currently called Quintus Corporation, and is a wholly owned subsidiary of Intergraph Corporation. The founders of Quintus are David H. D. Warren, Lawrence Byrd, William Kornfeld, and Fernando Pereira. They were joined by David Bowen shortly thereafter, and Richard O’Keefe in 1985. Tim Lindholm was responsible for many improvements including discontinuous stacks and the semantics for self-modifying code (see below). Many other people contributed to the implementation. Quintus Prolog 1.0 first shipped in 1985.

Quintus Prolog compiles to an efficient and compact direct threaded-code representation. For portability and convenience, the emulator is written in *Progol*,¹¹ a macrolanguage which is essentially a macroassembler for Prolog using Prolog syntax. The mode flag does not exist explicitly, but is cleverly encoded in the program counter by giving the **unify** instructions two entry points.

Quintus Prolog made several notable contributions, including those listed below.

- It is the Prolog system that generates the most compact code. Common sequences of operations are encoded as single opcodes. The code size is several times smaller than native code implementations. For example, the code generated for a given input program is about one-fifth the size of that generated by the BIM compiler. It is between one-fifth and one-half the code size of Aquarius Prolog (the latter figure only when the global analysis of Aquarius performs well) [155]. For applications with large databases this property can make the difference between good and bad performance. The recent rapid increase in physical memory size makes reducing code size less of a priority, although there will always be applications with lots of knowledge (e.g., databases and natural language) that require compact code to run well.
- It was the Prolog system that first developed a foreign language interface. Since then, it is the Prolog system that has put the most effort into making the system embeddable. It is important to seamlessly integrate Prolog code with existing code. This implies a set of abilities to make the system well behaved and expressive. It is able to redefine the C and Prolog top levels. It allows arbitrary calls between Prolog and C, with efficient manipulation of Prolog terms by C and vice versa. It has an open interface to the operating system that lets one redefine the low-level interfaces to memory management and I/O. It handles signals and memory allocation correctly, e.g., it was the first system to run efficiently with discontinuous stacks. This “small footprint” version has been available since release 3.0. It carefully manages the Prolog memory area to avoid conflicts with C. It provides tools for the user including source-level debugging on compiled code and an Emacs interface.

¹¹The name is a contraction of Prolog and Algol.

- It was the first system to provide a clean and justified semantics for self-modifying code (assert and retract), namely, the *logical view* [80]. A predicate in the process of being executed sees the definition that existed at the time of the call.
- It is the system that comes with the largest set of libraries of useful utilities. More than 100 libraries are provided. This provides much extra functionality that is important to users.

3.1.5. **BIM PROLOG (PROLOG BY BIM).** BIM Prolog was developed by the BIM company in Everberg, Belgium, in close collaboration with the Catholic University of Louvain (KUL). The name has recently been changed to “ProLog by BIM” due to a copyright conflict with the prefix “BIM” in the United States.

Logic programming research at KUL started in the mid-1970s. Maurice Bruynooghe had developed one of the early Prolog systems, Pascal Prolog, which was used at BIM at that time. The BIM Prolog project started in October 1983. It was then called P-Prolog (P for professional). Its execution model was originally derived from the PLM model in Warren’s dissertation, but was quickly changed to the WAM.

The first version of BIM Prolog, release 0.1, was distributed in October 1984 and was used in an ESPRIT project. It was a simple WAM-based compiler and emulator. Meanwhile, Quintus had released their first system. The BIM team realized that they needed to go further than emulation to match the speed of Quintus, so they decided immediately to do a native code implementation through macro expansion of WAM instructions. In contrast to Quintus, which intended to cover all major platforms from the start, BIM initially concentrated on Sun and decided to do a really good implementation there. By this time (1985) the team consisted of the three main developers who are still there today: Bart Demoen, André Mariën, and Alain Callebaut. Other people have contributed to the implementation. Because BIM Prolog only ran on a few machines, it was possible for different implementation ideas to be tried over the years. For more information on the internals of BIM Prolog, see [88].

BIM Prolog made several notable contributions, including those listed below.

- It was the first system in the WAM era:
 - To do native code compilation.
 - To do heap garbage collection. The Morris constant-space pointer-reversal algorithm was available in release 1.0 in 1985.
 - To do symbol table garbage collection. This is important if the system is interfaced to an external database.
 - To support mode declarations and do multiple-argument indexing instead of indexing only on the first argument.
 - To provide modules.

These abilities were provided earlier by DEC-10 Prolog (see Section 2.1) and MProlog (see Section 3.1.1).

- It was the first system to have a source-level debugger (like dbxtool), an external database interface, and separate compilation.

3.1.6. **IBM PROLOG.** IBM Prolog was developed primarily by Marc Gillet at IBM Paris. Nothing has been published about the implementation of this system. The following information is due to Gillet and the system documentation [66]. The first version, a structure-sharing system, was written in 1983–1984 and commercialized in 1985 as VM/Prolog. A

greatly rewritten and extended version was commercialized in 1989 as IBM Prolog.¹² It runs on system 370 under the VM and MVS operating systems. The system was ported to OS/2 with a 370 emulator.

The system is WAM-based and supports delaying, rational tree unification, and indefinite precision rational arithmetic. The system does global analysis at the level of a single module (see Section 2.4.5). It supports mode declarations, but may generate incorrect code if the declarations are incorrect. The system generates native 370 code and has a foreign language interface.

3.1.7. SEPIA AND ECLIPSE. ECRC (European Computer-Industry Research Centre) was created in Munich, Germany in 1984 jointly by three companies: ICL (UK), Bull (France), and Siemens (Germany). ECRC has done research in sequential and parallel Prolog implementation, in both software and hardware. See Section 3.2.2 for a discussion of the hardware work. The constraint language CHIP was built at ECRC (see Section 2.3.1).

Several Prolog systems were built at ECRC. An early system is ECRC-Prolog (1984–1986), a Prolog-to-C compiler for an enhanced MU-Prolog. At the time, this system had the fastest implementation of delaying. The next system, SEPIA (standard ECRC Prolog integrating advanced features), first released in 1988, was a major improvement [92]. Other systems are Opium [44], an extensible debugging environment, and MegaLog [15], a WAM-based system with extensions to manage databases (e.g., persistence). The most recent system, ECLIPSe (ECRC common logic programming system) [45, 94], integrates the facilities of SEPIA, MegaLog, CHIP, and Opium. The system supports rational tree unification and indefinite precision rational arithmetic. It provides libraries that implement constraint solvers for atomic finite domains and linear inequalities.

ECLIPSe is a WAM-based emulator with extensive support for delaying [94]. This makes it easy to write constraint solvers in the language itself. It has variables with attributes (called *metaterms*). Suspensions are an opaque data type at the Prolog level. A goal can be delayed explicitly by making it into a suspension and inserting it into a list of delayed goals. The list is stored as an attribute of a variable. When the variable is unified, an event handler is invoked. It is free to manipulate the suspended goals in any way. In this manner, the wakeup order of suspended goals can be programmed by the user.

The ECLIPSe compiler is incremental and compilation time is probably the lowest of any major system. The debugger uses compiled code supplemented with debugging instructions. Because of this, the system has no need of an interpreter. ECLIPSe (and SEPIA before it) uses two-word (64-bit) data items, with a 32-bit tag and a 32-bit value field. This allows more flexibility in tag assignment, and full pointers can be stored directly in the value field. It also makes for a more straightforward C interface.

3.1.8. SB-PROLOG AND XSB. SB-Prolog is a WAM-based emulator developed by a group at SUNY (State University of New York) at Stony Brook led by David Scott Warren. The compiler was written by Saumya Debray and the system was bootstrapped with C-Prolog. After several years of development, SB-Prolog was made available by Debray from Arizona in 1986. Because it was free and portable, it became quite popular. Neither it nor XSB do garbage collection. The worst problem regarding portability was the use of the BSD Unix *syscall* system call which supports arbitrary system calls through a single interface.

¹²Curiously, both systems are written mostly in assembly code, several hundred thousand lines worth.

SB-Prolog was the basis for much serious exploration related to language and implementation (e.g., [37]): backtrackable assert, existential variables in asserted clauses, memoizing evaluation, register allocation, mode and type inferencing (see Section 2.4.5), module systems, and compilation.

The most recent system, XSB, is SB-Prolog extended with memoization (tabling) and HiLog syntax [118]. The resulting engine is the SLG-WAM (see Section 2.3.3). XSB 1.3 implements the SLG-WAM for *modularly stratified* programs, that is, for programs that do not dynamically have recursion through negation.

3.1.9. SICStus PROLOG. SICStus Prolog¹³ was developed at SICS (Swedish Institute of Computer Science) near Stockholm, Sweden [25]. SICS is a private foundation founded in late 1985 which conducts research in many areas of computer science. It is sponsored in part by the Swedish government and in part by private companies. The guiding force and main developer of SICStus is Mats Carlsson. Many other people have been part of the development team and have made significant contributions.

As of this writing, SICStus Prolog is probably the most popular high-performance Prolog system running on workstations. There are many reasons for this. It is cheap, robust, fast, and highly compatible with the “Edinburgh standard.” It has been ported to many machines. It has flexible coroutining, rational tree unification, indefinite precision integer arithmetic, and a boolean constraint solver.

The first version of SICStus Prolog, release 0.3, was distributed in 1986. SICStus became popular with the 0.5 release in 1987. Originally, SICStus was an emulated system written in C. MC680X0 and SPARC native code versions were developed in 1988 and 1991. The current version, release 2.1, has been available since late 1991.

SICStus is the first system to do path compression (“variable shunting”) of dereference chains during garbage collection [119]. The parts of a dereference chain in the same choice point segment are removed. This lets the garbage collector recover more memory. This is essential for Prologs that have *freeze* or similar coroutining programming constructs [23], since the intermediate variables in a dereference chain may contain large frozen goals that can be recovered.

Among the “large program” problems encountered during the development of SICStus are those listed below.

- Interface with malloc/free, the Unix memory allocation library. SICS wrote their own version of the malloc/free library that better handles the allocation done by their system. Increasing the size of system areas is done by calling realloc.
- Native code limitations. One problem for large programs is that the offsets in machine instructions have a limited size. For example, the SPARC’s load and store instructions use a register+displacement addressing mode with a displacement limited to 12 bits. Other native code systems (e.g., IBM Prolog) have run into the same problem.
- Garbage collection bugs. A general problem with garbage collectors is that the connection between the effect of a bug and its cause is often hard to find. Subtle bugs tend only to show up in large programs because they exercise the garbage collector more than small programs.
- The space versus time trade-off. The code size of native code implementations is

¹³The name is a pun on Quintus.

much larger than emulated implementations. This difference can be quite significant: a factor of 5 or more. For large programs, e.g., natural language parsers with large databases, this can mean the difference between a program that runs and one that thrashes. SICStus minimizes the size of its generated native code by calling little-used operations as subroutines rather than putting them in-line. For example, the dereference operation is in-lined only for a predicate's first argument.

3.1.10. **AQUARIUS PROLOG.** Aquarius Prolog was originally developed in the context of the Aquarius project at UC Berkeley as the compiler for the VLSI-BAM processor [154]. See Section 3.2.3 for the hardware side of the story. After our relationship with the hardware side of the project ended in the spring of 1991, Ralph Haygood (the main developer of the back-end, run-time system, and built-ins) and I decided to continue part-time work on the complete system so that it could be released to the general public [58]. We were joined by Tom Getzinger at USC. The system achieved 1.1 MLIPS on a SPARCstation 1+ in February 1991. It first successfully compiled itself in February 1992. It was completed and released as Aquarius Prolog 1.0 on April 1, 1993.

Aquarius Prolog made several notable contributions, including those listed below.

- It is the first system to compile to native code without a WAM-like intermediate stage. It compiles first to BAM code (see Section 2.4.4), and then macro-expands to native code.
- It is the first well-documented system to do global analysis. See Section 2.4.5 for more information on the analyzer and Sections 2.4.4 and 2.4.6 for more information on how it is used to improve code generation. Type and mode declarations are supported. They are used to supplement the information generated by analysis. The system may generate incorrect code if the declarations are incorrect.
- It is the first system in which most built-ins are written in Prolog with little or no performance penalty. A technique called *entry specialization* replaces built-ins by more specialized entry points depending on argument types known at compile time.
- It is the first system to generate code which rivals the performance of an optimizing C compiler on a nontrivial class of programs [155].

The main disadvantage of Aquarius in its current state is the compilation time. This has little to do with the sophistication of the optimizations performed, but is due primarily to the naive representation of types in the compiler. The representation was chosen for ease of development, not speed. It is user-readable and new types can be added easily.

The path from research prototype to robust system is long and tortuous, as any system developer can attest. Many long hours are spent tracking down what turns out to be trivial problems. If one is lucky, no disastrous design decision was made early on, and the project can reach a successful end. We were definitely lucky in this regard. Among the “large program” problems encountered during the development of Aquarius are those listed below.

- **Garbage collection with uninitialized variables.** Before they are bound, uninitialized variables contain unpredictable information. The garbage collector must be able to handle this correctly. In Aquarius, the garbage collector follows all pointers, including uninitialized variables. Hence, it does not recover as much memory as it could. As far as we can tell, this does not adversely affect the system, in practice. All programs we have tried, including very long running ones, have stable memory sizes.

- Interaction of memory management with malloc. The observed behavior was that the system crashed because some stdio routines were writing outside their allocated space: they had called malloc. This is incompatible with our memory manager because it expands memory size if more memory is needed. After such an expansion, the malloc-allocated memory is inside a Prolog stack. On some platforms there is a routine, `f_prealloc`, that ensures that stdio routines do all of their allocation at startup. This does not work for all platforms. Our final solution uses a public domain malloc/free package (written by Michael Schroeder) that is given its own region of memory upon startup.
- During the MIPS processor port, a bug was found in the MIPS assembler. The assembler manual states that registers `t0–t9` (\$8–\$15, \$24–\$25) are not preserved across procedure calls. The MIPS instruction scheduler apparently assumes that they need not be saved even across branches, but this is not documented. We solved the problem with the directive “`.set nobopt,`” which prevents the scheduler from moving an instruction at a branch destination into the delay slot. The problem went undiscovered until we made the system self-compiling.

3.2. *Hardware Histories*

Starting in the early 1980s there was great interest in building hardware architectures optimized for Prolog. Two events catalyzed this interest: the start of the Japanese Fifth Generation Project in 1982 and the development of the WAM in 1983. In 1984, Tick and Warren proposed a paper design of a microcoded WAM that was influential for these developments [143]. At first, the specialized architectures were mostly microcoded implementations of the WAM (e.g., the PLM and the PSI-II). Later architectures (e.g., the KCM and the VLSI-BAM) modified the WAM design.

Some of the most important efforts are the PSI and CHI machine projects primarily at ICOT, the KCM project at ECRC, the POPE project at the GMD in Berlin, the Pegasus project at Mitsubishi, the Aquarius project at UC Berkeley (with its commercial offspring, Xenologic Inc.), and the IPP project at Hitachi. All these groups built working systems.

Several designs (POPE, PUP, PLUM, and FPPM) are based on extracting fine-grain parallelism in WAM instructions. The PUP, PLUM, and FPPM are described in the context of the Aquarius project. The POPE (parallel operating Prolog engine) is a ring of up to seven tightly coupled sequential Prolog processors [11]. Parallelism is achieved at each call by interleaving argument setup with head unification. The head unification is done on the next machine in the ring. In this fashion, the machine is automatically load balancing and achieves a speed-up of up to 7. The machine was built in Berlin at the GMD (Gesellschaft für Mathematik und Datenverarbeitung).

The IPP (integrated Prolog processor) [76] is a Hitachi ECL superminicomputer of cycle time 23 ns with 3% added hardware support for Prolog. The support comprises an increased microcode memory of 2 KW and tag manipulation hardware. The IPP implements a microcoded WAM instruction set modified to reduce pipeline bubbles and memory references. Its performance is comparable to Aquarius Prolog on a SPARCstation 1+ (see Table 4.1).

In the late 1980s came the first efforts to build RISC processors for Prolog. These include Pegasus, LIBRA [100], and Carmel-2 [56] (the latter supports flat concurrent Prolog). For lack of appropriate compiler technology, these systems executed macro-expanded WAM code or hand-coded assembly code.

The Pegasus project began in 1986. They designed and fabricated three single-chip RISC microprocessors in the period 1987–1990 [124]. The first two chips were fabricated

in October 1987 and August 1988 [122, 123, 167]. The third and last chip, Pegasus-II, was fabricated in September 1990 and at 10 MHz achieves a performance comparable to the KCM (see Table 4.1). The last two chips ran the Warren benchmarks a few months after fabrication. The chips have a bank of shadow registers to improve the performance of shallow backtracking. They provide support for tagging and dereferencing with ideas similar to those of the VLSI-BAM and KCM. Pegasus-II has two remarkable features. It provides support for context-dependent execution (which the designers call “dynamic execution switching”) of read/write mode in unification (see Section 3.2.2). It provides compound instructions (pop & jump, push & jump, pop & move, push & move) to exploit data path parallelism.

By 1990, the appropriate compiler technology was developed on two RISC machines. The VLSI-BAM, a special-purpose processor, ran Aquarius Prolog [63]. The MIPS R3000, a general-purpose processor, ran Parma [140]. The VLSI-BAM has a modest amount of architectural support for Prolog (10.6% of active chip area). Parma achieved the same performance on a general-purpose processor (see Table 4.1). The major difference between the two systems is that Parma has a bigger type domain in its analysis (see Figure 2.10 and Section 2.4.5).

The experience with Aquarius and Parma proves that there is nothing inherent in the Prolog language that prevents it from being competitive in speed to imperative languages. Comparing the two systems strongly suggests that improved analysis lessens the need for architectural support.

Since 1990 the main interest in special-purpose architectures has been as experiments to guide future general-purpose designs. The interest in building special-purpose architectures for their own sake has died down. Better compilation techniques and increasingly faster general-purpose machines have taken the wind out of its sails (see also Section 5.1). This parallels the history of Lisp machines.

The rest of this section examines three projects in more detail: the PSI machine project (ICOT/Mitsubishi/Oki), the KCM project (ECRC), and the Aquarius project (UC Berkeley). I have chosen these projects because they show clearly how system performance improved as Prolog was better understood and because I have detailed information on them.

3.2.1. ICOT AND THE PSI MACHINES. The FGCS (Fifth Generation Computer System) project at ICOT (Japanese Institute for New Generation Technology) has designed and built a large number of sequential and parallel Prolog machines [135, 148]. Both in manpower and machines, this is the largest architecture project in the logic programming community. Two series of sequential machines were built: the PSI (personal sequential inference) machines (PSI-I, PSI-II, and PSI-III) and the CHI (cooperative high-performance sequential inference) machines (CHI-I and CHI-II) [54]. I will limit the discussion to the PSI machines, which were the most popular. All the PSI machines are horizontally microprogrammed and have 40-bit data words with 8-bit tag and 32-bit value fields.

The PSI-I was developed before the WAM [134]. After the development of the WAM it was followed by two WAM-based machines, the PSI-II and PSI-III. The three models were manufactured by Mitsubishi and Oki, and commercialized by Mitsubishi inside of Japan. Several multiprocessors were built at ICOT with these processors as their sequential processing elements (e.g., the PSI-II is the PE of the Multi-PSI/v2 and the PSI-III is the PE of the PIM/m).

The PSI-I was designed as a personal workstation for logic programming. It was first operational in December 1983 at a clock rate of 5 MHz. It runs ESP (extended sequential Prolog), a Prolog extended with object-oriented features. More than 100 machines were

shipped. The first ESP implementation was an interpreter written in microcode (not a WAM). A WAM emulator was later written for the PSI-I and it ran twice as fast. The main advantage of the PSI-I was not speed, but memory. It had 80 MB of physical memory, a huge amount in its day.

The PSI-II was first operational in December 1986 [108]. More than 500 PSI-II machines were shipped from 1987 until 1990 and delivered primarily to ICOT. Its clock was originally 5 MHz, but was quickly upgraded to 6.45 MHz. At the higher clock, its average performance is three to four times that of the interpreted PSI-I.

The PSI-III was first operational near the end of 1990. More than 200 PSI-III machines have been shipped. It is binary compatible with the PSI-II and has almost the same architecture with a clock rate of 15 MHz. The microcode was ported from the PSI-II by an automatic translator. Its average performance is two to three times that of the upgraded PSI-II.

3.2.2. ECRC AND THE KCM. The architecture work at ECRC culminated in the KCM (knowledge crunching machine) project, which started in 1987 [14, 111]. The KCM was probably the most sophisticated Prolog machine of the late 1980s. It had an innovative architecture and significant compiler design was done for it. It was preceded by two years of preliminary studies (the ICM, ICM3, and ICM4 architectures) [110, 166]. The KCM was built by Siemens. The first prototypes were operational in July 1988 and ran at a clock speed of 12.5 MHz. About 50 machines were delivered to ECRC and its member companies [142].

The KCM is a single user, single tasking, dedicated coprocessor for Prolog, used as a back-end to a Unix workstation. It is a tagged general-purpose design with support for Prolog, and hence is not limited to Prolog. It uses 64-bit data words, with a 32-bit tag and a 32-bit value field.

The KCM's instruction set consists of two parts: a general-purpose RISC part and a microcoded WAM-like part. Prolog compilation for the KCM is still WAM-like, but the instructions have evolved greatly from Warren's original design (see [91, 111]). The KCM supports the delayed creation of choice points. The KCM runs KCM-SEPIA, a large subset of SEPIA that was ported to it (see Section 3.1.7).

The Prolog support on the KCM improves its performance by $\approx 60\%$ [111, 142]. Most of the architectural extensions are added to the microengine. The architectural features and their effects on performance (in percent) are given in Table 3.1.

The MWAC (multi-way address calculator) is a functional unit that does a 16-way microcode branch depending on the types of two arguments. It calculates the target address during the last step of dereferencing. The MWAC is used in the execution of all unification operations. It is similar to the partial unification unit of the LIBRA [100].

Context-dependent execution uses flags in addition to the opcode during instruction decoding. Three flags are used: read/write mode for unification, choicepoint/nochoicepoint for delayed choice point creation, and deep/shallow for fast fail in shallow backtracking.

3.2.3. THE AQUARIUS PROJECT: THE PLM AND THE VLSI-BAM. In 1983, Alvin Despain and Yale Patt at UC Berkeley initiated the Aquarius project. Its main goal was to design high-performance computer systems with large symbolic and numeric components. The project continued at Berkeley until 1991.¹⁴ They decided to focus on Prolog architectures, being inspired by the FGCS project and seduced by the mathematical simplicity of

¹⁴Despain is continuing this work at USC's Advanced Computer Architecture Laboratory.

TABLE 3.1. The Benefits of Prolog-Specific Features in the KCM

Feature	Benefit (%)
Multiway tag branch (MWAC)	23.1
Context dependent execution (flags)	11.4
Dereferencing support	10.0
Trail support	7.2
Load term	5.7
Fast choice point creation/restoration	2.3
Total	59.7

Prolog. As soon as Warren presented the WAM at Berkeley, Despain turned the project to focus on hardware support for it. He proposed that I write a compiler for their architecture, the PLM. The compiler was completed and the report was delivered to the university on August 22, 1984.¹⁵ This was the first published WAM compiler [149].¹⁶

A whole series of sequential and parallel Prolog architecture designs came out of Aquarius. The sequential designs are:

- The PLM [42, 43] (1983–1987). The programmed logic machine.¹⁷ This is a microcoded WAM.
- The VLSI-PLM [129, 130] (1985–1989). This is a single-chip implementation of the PLM.
- The Xenologic X-1. This is a commercial version of the PLM, designed as a coprocessor for the Sun-3.
- The PUP [26] (1988). The parallel unification processor. This machine exploits fine-grained parallelism in unification using dynamic scheduling with the Tomasulo algorithm [145].
- The PLUM [126] (1989). The parallel unification machine. This is a second-generation PUP, more pipelined and able to use static information.
- The FPPM [127] (1990). The flow parallel prolog machine. This is a third-generation design that can exploit fine-grain parallelism for any goal, not just unification. Using a sequential processing element similar to the VLSI-BAM, and taking all overheads into account (including the cache coherence protocol), an FPPM with four processors has a simulated performance on the Warren benchmarks of about two times the VLSI-BAM.
- The VLSI-BAM [63] (1988–1991). The VLSI Berkeley abstract machine. This is a single-chip RISC processor with extensions for Prolog. This machine is described below.

¹⁵The exact day of my flight back to Belgium.

¹⁶In January 1991, I toured several German universities and research institutes to talk about Aquarius Prolog. At ECRC, a scientist from East Berlin came to me after the talk. He explained that they had *typed in* the source code of the PLM compiler from the appendix of the report.

¹⁷The name correspondence with the PLM model in Warren's dissertation [160] is a coincidence.

TABLE 3.2. The Benefits and Chip Area of Prolog-Specific Features in the VLSI-BAM

Feature	Benefit (%)	Area (%)
Fast tag logic (tagged branching)	18.9	1.6
Double-word memory port	17.1	1.9
Tag and segment mapping	10.3	4.8
Multi-cycle/conditional	9.1	0.1
Tagged-immediates	7.9	2.2
Arithmetic overflow detect	1.4	≈0.0
Total	70.1	10.6

These designs were all extensively simulated. Four were built. The PLM was wire-wrapped and ran a few small programs in 1985. The Xenologic X-1 has been running at 10 MHz since 1987. The VLSI-PLM was fabricated and ran all benchmarks at 10 MHz in June 1989. The VLSI-BAM was designed to run at 30 MHz. It was fabricated in November 1990 and ran most benchmarks of [155, 157] at 20–25 MHz on its custom cache board in November 1991.

The core of the VLSI-BAM is a RISC in the classic sense. It is a 32-bit pipelined load-store architecture with single-cycle instructions, 32 registers, and branch delay slots. The processor is extended with support for Prolog and for multiprocessing, which together form 10.6% of the active chip area and improve Prolog performance by ≈70% [63]. The VLSI-BAM executes the same Prolog program in one-third the cycles of the VLSI-PLM, a gain due to improved compilation.

The Prolog support takes the form of six architectural features and new instructions using them. The architectural features and their performance benefits and active chip area (both in percent) are given in Table 3.2. The benefit figures cannot be directly added up because the effects of the architectural features are not independent.

Except for dereference, the instructions are all single-cycle unless there is a pipeline stall or an annulled delay slot. There are two- and three-way tagged branches to support unification and a conditional push to support trailing. The instructions for data structure creation (write-mode unification) were derived automatically using constrained exhaustive search [64]. VLSI-BAM measurements [63] show that with advanced compilation techniques, multiway branches for general unification are effective only up to a three-way branch.¹⁸ Multiple-cycle (primarily dereference) and conditional instructions are implemented by logic to insert or remove opcodes in the pipeline. The opcode pipe has space for both user instructions and added “internal” instructions. The double-word memory port (with double bandwidth to cache) improves general-purpose memory operations as well as choice point creation and restoration speed.

4. EVOLUTION OF PERFORMANCE

¹⁸This does not contradict the measurements of the KCM’s MWAC since the latter is used for all unification operations, not just general unification.

TABLE 4.1. Evolution of Prolog Performance

System	Machine (Year)	Clock (MHz)	Benchmark					Mean
			N	Q	D	S	R	
†PLM compiler [149]	PLM [43] (1985)	10.	19	12	9	12	8	11
ESP	PSI-II (1986)	6.45	41	25	12	18	10	16
KCM-SEPIA [111]	KCM (1989)	12.5	83	57	37	33	15	32
†Pegasus compiler [124]	Pegasus-II (1990)	10.	91	69	39	40	19	39
†Aquarius [63]	VLSI-BAM (1991)	20.	270	260	75	57	32	72
	Machine (Architecture)							
‡DEC-10 Prolog [160]	DEC-10 (1977)		1	1	1	1	1	1
XSB 1.3	SPARCstation 1+ (SPARC)	25	7	4	2	4	3	3
Quintus 2.0 [63]	Sun 3/60 (MC68020)	20	11	4	3	4	3	4
‡MProlog 2.3	IBM PC clone (386)	33	13	6	5	5	2	5
ECLiPSe 3.3.7	SPARCstation 1+ (SPARC)	25	11	6	4	6	3	5
NU-Prolog 1.5.38	SPARCstation 1+ (SPARC)	25	22	7	5	7	2	5
SICStus 2.1	DECstation 5000/200 (R3000)	25	37	16	10	10	5	10
Quintus 2.5 [155]	SPARCstation 1+ (SPARC)	25	33	16	9	13	8	12
‡BIM 3.1 beta [155]	SPARCstation 1+ (SPARC)	25	34	21	8	16	8	13
‡SICStus 2.1	SPARCstation 1+ (SPARC)	25	39	26	15	20	8	17
‡†Aquarius [155]	SPARCstation 1+ (SPARC)	25	120	140	28	25	12	29
‡IBM Prolog	ES/9000 Model 9021 (370)		120	59	74	69	33	60
‡Aquarius 1.0	DECstation 5000/200 (R3000)	25	180	210	63	44	46	71
‡†Parma [141]	MIPS R3230 (R3000)	25	330	350	130	170	59	140

The performance of Prolog has increased about two orders of magnitude since DEC-10 Prolog. This can be equally attributed to improvements in hardware and software. Table 4.1 gives the execution time ratios relative to DEC-10 Prolog of a set of representative systems running the five Warren benchmarks [160]. For the reasons given below, the numbers in Table 4.1 do *not* generalize to large programs. *They should be seen only as indicating trends.*

Table 4.1 is split into two parts. The first five rows show the performance of specialized hardware. The following rows show general-purpose hardware. For the first five rows and for DEC-10 Prolog, the year in which the systems were first running is given. For the other systems the architecture is given. Results for the benchmarks nreverse, qsort, deriv, serialise, and query are given in columns N, Q, D, S, and R, respectively. Table 4.2 gives their absolute execution times on DEC-10 Prolog. The benchmarks were timed with a failure-driven loop. The deriv benchmark is the sum of the four benchmarks times10, log10, divide10, and ops8. The last column of Table 4.1 gives the harmonic mean of the speedup ratios. This is the correct mean when averaging speeds.

Performance is one of the few measures of a system's quality that is quantifiable. Many other measures are just as important, but are hard to quantify. For example, it is difficult to assign numbers to embeddability, robustness, debuggability, portability, and the usefulness of the available built-in operations. The overall quality of a system depends on how well it meets the needs of the task at hand. A rough indication of overall quality can be obtained from the software sagas presented earlier. This should be refined for a particular system by using it to solve a relevant problem.

The systems marked by † are research systems. The systems on general-purpose hardware that are marked by ‡ are native code systems. The others are emulated. The numbers for XSB 1.3 are within 10% of SB-Prolog 3.1. Many of the systems generate better code if the program has mode declarations. For example, IBM Prolog is about 1.5 times faster with mode declarations. MProlog 2.3 is about 1.2 times faster with mode and indexing declarations. On the same PC clone, emulated SICStus 2.1 is 1.5 times slower than MProlog

TABLE 4.2. Execution Times for the Warren Benchmarks on DEC-10 Prolog

Benchmark	N	Q	D	S	R
Time (ms)	53.7	75.0	10.1	40.2	185.0

2.3 and five times slower than native SICStus 2.1 on a SPARCstation 1+.

The Warren benchmarks were chosen because reliable performance numbers for them are available for many machines. They are *not* a good measure of the performance of real programs. A more realistic benchmark set that subsumes the Warren benchmarks is used in [141, 155] and may be obtained from [157].

The Warren benchmarks are small and many systems have been optimized to execute them fast. The speedup for *nreverse* is greater than average because more effort has been done to optimize it. The speedup for *query* is less than average because it is dominated by integer multiplication and division. Due to limitations in their analysis domains (see Section 2.4.5), *Aquarius* and *Parma* have lower performance for large programs unless the programs are tuned. Large programs are more likely to spend most of their time doing built-in operations, which are a fixed cost since they are usually implemented in a lower level language.

In older publications, a common unit in Prolog performance was the LIPS, or logical inferences per second, i.e., the number of goal invocations or procedure calls per second. Because the amount of work done by a procedure call is not constant, this number is an unreliable indicator of system performance and is not given. By convention, published LIPS numbers are measured for *nreverse*, which reverses a 30-element list in 496 logical inferences.

It is difficult to compare the performance of two systems unless they are running on identical hardware. For example, the same system can vary greatly in speed even when running the same CPU-bound program on two machines with the same processor, clock speed, and cache size. This could be the case because the write buffers are of different sizes. Among the machine-related factors that affect performance are clock speed, as well as the memory system (i.e., cache and virtual memory structure, memory size and bandwidth), the operating system (e.g., speed of I/O and context switching overhead), the implementation's data path (e.g., pipeline structure, multiple functional units, out-of-order and superscalar execution), and the implementation of various primitive operations (e.g., multiplication can vary an order of magnitude in speed even on systems with the same clock). An important difference between the SPARC-based and R3000-based systems in Table 4.1 is that the latter have a faster memory system.

5. FUTURE PATHS IN LOGIC PROGRAMMING IMPLEMENTATION

This section gives a personal view of the trends in sequential logic programming implementation. It is important to distinguish three levels of evolution. First, the low level trends. What will be the basic improvements in implementation technology for Prolog and related languages? Second, the high level trends. What will be the new tools, new languages, and programming paradigms? Finally, what will be the relation between Prolog and the mainstream computing community? See [48] for an early but still useful discussion of these issues.

5.1. Low Level Trends

There are many ways in which Prolog implementation technology can be improved. Here are some of the important ones, given in order of increasing difficulty:

- **Overlap with mainstream compiler technology.** As Prolog compilers approach imperative language performance, the standard optimizations of imperative language compilers (global register allocation, code motion, instruction reordering, and so forth) become important. Some of these are being implemented in current systems [38]. One approach is to compile to C. This shortens development time, gains portability, and (to a lesser degree) takes advantage of what the C compiler does (e.g., register allocation). This approach has traditionally had a performance loss over native code of a factor of 2 to 3. C is not a portable assembly language. This will change in the future. For example, because of its first-class labels and global register declarations, the recently released GNU C 2.X compiler has a smaller performance loss than other C compilers [36, 57]. Recent work shows that the overhead of compilation to C can be reduced to less than 30%, while keeping the system portable [98].
- **Type inference and operational types.** When writing a program, a programmer often has definite intentions about the types of predicate arguments. This includes information on the structure of compound terms (e.g., recursive types such as lists and trees) and on operational types (see Section 2.4.5). For analysis to work well with large programs as well as small benchmarks, the analysis domain has to represent this information, to track variable dependencies, and to correctly handle built-in predicates. Objects whose type is known at compile time can be represented *unboxed*, i.e., accessible without tagging or other overhead. Current systems only do this for variables (see the discussion on uninitialized variables in Section 2.4.4) and numbers within arithmetic expressions.
- **Determinism extraction.** Often, a deterministic user-defined predicate is used to select a clause. This is currently compiled by creating a choice point, executing the predicate, and backtracking if it fails. It would be more efficient to compile such a predicate as a boolean function and to do a conditional jump on its result.
- **Multiple specialization.** Different calls of the same predicate frequently have different types in the same argument. The predicate will run faster if it is compiled separately for each pattern of calling type. As a first step, this optimization can be enabled by a directive. Profiling could supply the directives. Measurements show that this is often fruitful. For example, in the `chat_parser` benchmark, the inner loop is a two-clause predicate, `terminal/5`, that is called 22 times. Making 22 copies and recompiling with analysis under Aquarius Prolog results in a 16% performance improvement. In programs with tighter inner loops, the performance improvement can be much greater. For example, the `SEND+MORE=MONEY` puzzle shows a tenfold reduction of execution time [156].
- **Compile-time garbage collection.** Prolog creates three kinds of data objects in memory: choice points, compound data terms, and environments. When a data object becomes inaccessible, a new object can often reuse part of the old one. For example, a program that uses an array can destructively update it if it is unaliased (see Section 2.4.4). Arrays with this property are called *single-threaded*. Recent developments indicate that it is more practical to enforce this condition syntactically (through source transformation) than to use a powerful analyzer-compiler combi-

nation. See, for example, the use of monads in functional programming [159] and the extended definite clause grammar notation of [152, 154], much improved in [7].

- **Dynamic to static conversion.** All data in Prolog is allocated dynamically, that is, at run time. It is accessed through tagged pointers. Often, it is necessary to follow a chain of pointers to find the data. Since CPU speed is increasing faster than memory speed [59], the overhead of memory access will become relatively more important in the future. The software and hardware approaches to speed up memory access are complementary:
 - A future compiler could statically allocate part of the dynamically allocated data to reduce access time and improve locality. This requires analysis to determine the evolution of aliasing during program execution. For example, objects that are unaliased, that exist only in one copy at any given time, and whose size is known can be allocated statically.
 - A future architecture could be designed to tolerate memory latency. If it could follow one level of tagged pointer in zero time, then the execution model of Prolog could be changed drastically and would run faster. Two techniques that help are starting to appear in existing architectures: asynchronous loads (decoupling the load request and arrival of the result) and multithreading (fast switching between register sets). These are useful for all languages.

5.2. High Level Trends

In recent years, the implementation of logic programming systems has continued in two main directions:

- **Further development of Prolog.**
 - Software engineering aspects. This development has been mostly in the area of extended usability of the system rather than performance. For example, many systems including Quintus, SICStus, BIM, and ECLiPSe, have a foreign language interface that allows arbitrary calls between Prolog and C, to any level of nesting. Quintus in particular has worked hard on allowing seamless integration of Prolog and foreign language code. Debugging has improved, and several systems now have source-level debuggers and profilers [51]. Many systems have eased Prolog's strict control flow by including corouting facilities (*freeze* and its relatives) [31]. There is an ISO standard for Prolog that is essentially complete [121].
 - “Cleaner” Prologs. These languages aim to keep the ideas and functionality of Prolog, but to replace the “dirty” operational features (such as *assert*, *var*, and *cut*) by clean declarative ones. It is not yet obvious whether this is possible without losing expressivity and performance. This group includes the MU-Prolog and NU-Prolog family [103] (see Section 3.1.3), *xpProlog* [82], and the Gödel language [62].
- **Other logic programming languages.** These can be roughly subdivided into three main families. The families overlap, but the division is still a useful rule of thumb.
 - Concurrent languages. These include the committed-choice languages [125] (e.g., *Parlog*, *FGHC*, and *FCP*) and languages based on the “Andorra principle” [33, 55] (an elegant synthesis of Prolog and committed-choice languages).

- **Constraint languages.** A language that does incremental global constraint solving in a particular domain is called a *constraint language*. These languages come in two flavors. The general-purpose languages (such as Prolog, Trilog [158], and LIFE [5]) provide domains that are useful for most programming tasks. For example, unification in Prolog handles equality constraints over finite trees. The special-purpose languages (such as Prolog III, CLP(R), and CHIP) provide specialized domains that are useful for particular kinds of problems. For example, linear arithmetic inequalities on real numbers and membership in finite domains. These languages allow practical solutions to many problems previously considered intractable such as optimization problems with large search spaces.
- **“Synthesis” languages.** There are now serious attempts to make syntheses of different styles of programming [53]. For example, λ Prolog [99] and languages based on narrowing are syntheses of logic and functional programming, LIFE is a synthesis of logic, functional, and object-oriented programming, and AKL [55] is a synthesis of concurrent and constraint languages [120]. An important principle is that a synthesis must start from a simple theoretical foundation.

Both the development of Prolog and more advanced logic languages are active areas of research. They promise many exciting new ideas and systems in the years to come.

5.3. *Prolog and the Mainstream*

As measured by the number of users, commercial systems, and practical applications, Prolog is by far the most successful logic programming language. Prolog’s closest competitors are surely the special-purpose constraint languages. However, it is true that logic programming, in particular, and declarative programming, in general, remain outside of the mainstream of computing. Two important factors that influence the widespread acceptance of Prolog are:

- **Compatibility.** Existing code works and investment in it is large. Therefore, people will not easily abandon it for new technology. Therefore, a crucial condition for acceptance is that Prolog systems be embeddable. This problem has been solved to varying degrees by commercial vendors (see Section 3.1.4).
- **Public perception.** To the commercial computing community, the terms “Prolog” and “logic programming” are, at best, perceived as useful in an academic or research setting, but not useful for industry. This image is not based on any rational deduction. Changing it requires both marketing and application development.

The ideas of logic programming will continue to be used in those areas for which it is particularly suited. This includes those areas in which program complexity is beyond what can be managed in the imperative paradigm.

6. SUMMARY AND CONCLUSIONS

This survey summarizes the technical developments in sequential Prolog implementation of the past ten years and the systems that pioneered them. Much has happened in this time, and I hope that this survey is successful in capturing most of the important developments and in pointing out some intriguing trends for the future.

The WAM was the starting shot for a proliferation of systems and ideas. It was the substrate of thought for most sequential Prolog developers for the last decade. The compilation principle that underlies the WAM and its relatives continues to hold true: to compile a logic language is to simplify each occurrence of one of its basic operations with all the information at one's disposal. The last decade has seen an increased understanding of how this can be done: by measuring actual programs to optimize frequent operations, by learning how to compile unification and backtracking, and by using simpler instruction sets and global analysis.

The Prolog language has withstood the test of time and has proved to be an elegant implementation target. The problems in the original language design have been identified and are being corrected. The language is being generalized in many ways. There have been large advances in implementation technology, but there is still plenty to do, both in implementing Prolog and its successors. The next decade promises to be as interesting as the first.

The author thanks the many developers and pioneers in the logic programming community. In particular, the following friends and colleagues helped tremendously by recollecting past events and providing critical comments: Abder Aggoun for ECRC and CHIP, Dave Bowen for Quintus Prolog, Mats Carlsson for SICS, SICStus Prolog, and its native code timing measurements, Koen De Bosschere, Saumya Debray for SB-Prolog and QD-Janus, Bart Demoen and André Mariën for Leuven and BIM Prolog, Marc Gillet for IBM Prolog and its timing measurements, Manuel Hermenegildo for MA³ and PLAI, Bruce Holmer, Tim Lindholm for Quintus Prolog, Peter Ludemann, Micha Meier for ECRC, SEPIA, the KCM, and ECLiPSe, Richard Meyer for putting up with this paper "black hole," Lee Naish and Jeff Schultz for MU-Prolog and NU-Prolog and their timing measurements, Hiroshi Nakashima, Katsuto Nakajima, Takashi Chikayama, and Kouichi Kumon for ICOT, the PSI machines, and their timing measurements, Ulrich Neumerkel for the VAM and BinProlog, Jacques Noyé for ECRC and the KCM, Fernando Pereira for Edinburgh, DEC-10 Prolog, C-Prolog, and Quintus Prolog, Christian Pichler for IF/Prolog and SNI-Prolog, Andreas Podelski, Olivier Ridoux for MALI and λProlog, Konstantinos Sagonas for SB-Prolog, XSB, and its timing measurements, Bart Sano for the VLSI-BAM, Kazuo Seo for Pegasus and the Pegasus-II timing measurements, Ashok Singhal for the FPPM, Zoltan Somogyi for NU-Prolog and databases, Vason Srinii for the VLSI-PLM, Péter Szeredi for MProlog and its timing measurements, Paul Tarau for BinProlog, Evan Tick for ICOT, David H.D. Warren for Edinburgh, DEC-10 Prolog, and Quintus Prolog, David Scott Warren for Stony Brook, SB-Prolog, the SLG-WAM, and XSB, S. Bharadwaj Yadavalli, and finally, the referees for a host of good suggestions. Special thanks to TεXnician *par excellence* Hassan Aït Kaci.

REFERENCES

1. Aggoun, A., and Beldiceanu, N., Time Stamps Techniques for the Trailed Data in Constraint Logic Programming Systems, in: *Actes du Séminaire 1990-Programmation en Logique*, CNET, Tregastel, France, May 1990.
2. Aggoun, A., and Beldiceanu, N., Overview of the CHIP Compiler System, in: *8th ICLP*, MIT Press, Cambridge, MA, June 1991, pp. 775-789.
3. Aho, A. V., Sethi, R., and Ullman, J. D., *Compilers: Principles, Techniques, and Tools*, Addison-Wesley, Reading, MA, 1986.
4. Aït-Kaci, H., *Warren's Abstract Machine, A Tutorial Reconstruction*, MIT Press, Cambridge, MA, 1991.
5. Aït-Kaci, H., and Podelski, A., Towards a Meaning of LIFE, DEC PRL Research Report 11, Digital Equipment Corporation, Paris Research Laboratory, June 1991 (revised October 1992).

6. Aït-Kaci, H., and Podelski, A., Functions as Passive Constraints in LIFE, DEC PRL Research Report 13, Digital Equipment Corporation, Paris Research Laboratory, June 1991 (revised November 1992).
7. Aït-Kaci, H., Meyer, R., Van Roy, P., Podelski, A., and Dumant, B., *The Wild Life Handbook*, Digital Equipment Corporation, Paris Research Laboratory, 1994.
8. Amraoui, M., Une Expérience de Compilation de PrologII sur MALI, Doctoral dissertation, Université de Rennes I, France, January 1988, (in French).
9. Appel, A. W., *Compiling with Continuations*, Cambridge University Press, 1992.
10. Baggins, B., and Baggins, F., *The Memoirs of Bilbo and Frodo of the Shire, Supplemented by the Accounts of Their Friends and the Learning of the Wise*, The Shire, Arnor, 3021 TA.
11. Beer, J., Concepts, Design, and Performance Analysis of a Parallel Prolog Machine, Ph.D. dissertation, Technische Universität Berlin, September 1987.
12. Beer, J., The Occur-Check Problem Revisited, *J. Logic Program.* 5(3):243–261 (1988).
13. Bendl, J., Kövcs, P., and Szeredi, P., The MPROLOG System, in: *Logic Programming Workshop*, Debrecen, Hungary, 1980, pp. 201–209.
14. Benker, H., Beacco, J. M., Bescos, S., Dorochevsky, M., Jeffre, Th., Pohimann, A., Noyé, J., Poterie, B., Sexton, A., Syre, J. C., Thibault, O., and Watzlawik, G., KCM: A Knowledge Crunching Machine, in: *16th ISCA*, IEEE Computer Society Press, New York, 1989, pp. 186–194.
15. Bocca, J., MegaLog—A Platform for Developing Knowledge Base Management Systems, in: *2nd International Symposium on Database Systems for Advanced Applications*, Tokyo, April 1991, pp. 374–380.
16. Boortz, K., SICStus Maskinkodskompilering, SICS Technical Report T91:13, August 1991 (in Swedish).
17. Bowen, D. L., Byrd, L. M., and Clocksin, W. F., A Portable Prolog Compiler, in: *Logic Programming Workshop*, Algarve, Portugal, 1983, pp. 74–83.
18. Boyer, R. S., and Moore, J. S., The Sharing of Structure in Theorem Proving Programs, in: *Machine Intelligence 7*, Edinburgh University Press, New York, 1972, pp. 101–116.
19. Brisset, P., and Ridoux, O., The Compilation of λ Prolog and its Execution with MALI, IRISA Publication Interne 687, Rennes, France, November 1992. Also published as INRIA Rapport de Recherche 1831, January 1993.
20. Bruynooghe, M., An Interpreter for Predicate Logic Programs, Report CW 10, Department of Computer Science, Katholieke Universiteit Leuven, Belgium, October 1976.
21. Bruynooghe, M., The Memory Management of Prolog Implementations, in: *Logic Programming*, K. Clark and S. Tärnlund (ed.), Academic Press, 1982, pp. 83–98.
22. Carlsson, M., On Implementing Prolog in Functional Programming, in: *1st ICLP*, IEEE Computer Society Press, New York, 1984, pp. 154–159.
23. Carlsson, M., Freeze, Indexing, and Other Implementation Issues in the WAM, in: *4th ICLP*, MIT Press, Cambridge, MA, 1987, pp. 40–58.
24. Carlsson, M., On the Efficiency of Optimising Shallow Backtracking in Compiled Prolog, in: *6th ICLP*, MIT Press, Cambridge, MA, 1989, pp. 3–16.
25. Carlsson, M., Widèn, J., Andersson, J., Andersson, S., Boortz, K., Nilsson, H., and Sjöland, T., *SICStus Prolog User's Manual*, SICS, Box 1263, 164 28 Kista, Sweden, 1991.
26. Chen, C., Singhal, A., and Patt, Y. N., PUP: An Architecture to Exploit Parallel Unification in Prolog, Report UCB/CSD 88/414, UC Berkeley, March 1988.
27. Chen, W., and Warren, D. S., Query Evaluation under the Well-Founded Semantics, in: *12th Symposium on Principles of Database Systems (PODS '93)*, ACM, New York, 1993.
28. Clocksin, W. F., Design and Simulation of a Sequential Prolog Machine, in: *J. New Generation Comput.*, 3(1):101–120 (1985)
29. Codognet, P., and Diaz, D., Boolean Constraint Solving using clp(FD), in: *10th ILPS*, MIT Press, Cambridge, MA, 1993, pp. 525–539.
30. Coelho, H., and Cotta, J. C., *Prolog by Example: How to Learn, Teach and Use It*, Springer, New York, 1988.

31. Colmerauer, A., Kanoui, H., and Caneghem, M. V., Prolog, Theoretical Principles and Current Trends, in: *Technol. Sci. Inform.* 2(4):255–292 (1983).
32. Colmerauer, A., The Birth of Prolog, in: The Second ACM-SIGPLAN History of Programming Languages Conference, *ACM SIGPLAN Not.* March:37–52 (1993).
33. Santos Costa, V., Warren, D. H. D. and Yang, R., Andorra-I: A Parallel Prolog System that Transparently Exploits both And- and Or-Parallelism, in: *Proceedings of the 3rd ACM SIGPLAN Conference on Principles and Practice of Parallel Programming*, August 1991, pp. 83–93.
34. Cousot, P., and Cousot, R., Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints, in: *4th POPL*, January 1977, pp. 238–252.
35. Cousot, P., and Cousot, R., Abstract Interpretation and Application to Logic Programs, *J. Logic Program.* 13(2-3):103–179 (1992).
36. De Bosschere, K., and Tarau, P., Continuation Passing Style Prolog-to-C Mapping at Native WAM-speed, ELIS Technical Report DG 93-15, Universiteit Gent, Vakgroep Elektronica en Informatiesystemen, November 1993. Summary in *ACM Symposium on Applied Computing (SAC '94)*, March 1994 (forthcoming).
37. Debray, S., Global Optimization of Logic Programs, Ph.D. dissertation, Computer Science Department, SUNY Stony Brook, September 1986.
38. Debray, S., A Simple Code Improvement Scheme for Prolog, in *J. Logic Program.* 13(1):57–88 (1992).
39. Debray, S., Implementing Logic Programming Systems: The Quiche-Eating Approach, in: *ICLP '93 Workshop on Practical Implementations and Systems Experience*, Budapest, Hungary, June 1993.
40. Diaz, D., and Codognet, P., A Minimal Extension of the WAM for clp(FD), in: *10th ICLP*, Budapest, Hungary, MIT Press, Cambridge, MA, 1993, pp. 774–790.
41. Dinckbas, M., Van Hentenryck, P., Simonis, H., Aggoun, A., Graf, T., and Berthier, F., The Constraint Logic Programming Language CHIP, in: *FGCS '88*, Tokyo, November 1988, pp. 693–702.
42. Dobry, T. P., Performance Studies of a Prolog Machine Architecture, in: *12th ISCA*, IEEE Computer Society Press, New York, 1985.
43. Dobry, T. P., A High Performance Architecture for Prolog. Ph.D. dissertation, Report UCB/CSD 87/352, Department of Computer Science, UC Berkeley, April 1987. Also published by Kluwer Academic Publishers, 1990.
44. Ducassé, M., Opium: An Advanced Debugging System, in: G. Comyn and N. Fuchs (eds.), *2nd Logic Programming Summer School*, Esprit Network of Excellence in Computational Logic (COMPULOG-NET), Lecture Notes in Artificial Intelligence 636, Springer, 1992.
45. ECRC, *ECLiPSe 3.2 User Manual*, August 1992.
46. Elcock, E. W., Absys: The First Logic Programming Language—A Retrospective and a Commentary, 9(1):1–17 (1990). Also published as Technical Report 210, Department of Computer Science, University of Western Ontario, July 1988.
47. Farkas, Z., Köves, P., and Szeredi, P., MProlog: An Implementation Overview, in: *ICLP '93 Workshop on Practical Implementations and Systems Experience*, Budapest, Hungary, June 1993.
48. Gallaire, H., Boosting Logic Programming, in: *4th ICLP*, Melbourne, Australia, May 1987, pp. 962–988.
49. García de la Banda, M., and Hermenegildo, M., A Practical Approach to the Global Analysis of CLP Programs, in: *10th ILPS*, MIT Press, Cambridge, MA, 1993, pp. 435–455.
50. Getzinger, T. W., Abstract Interpretation for the Compile-Time Analysis of Logic Programs, Ph.D. dissertation, Report ACAL-TR-93-09, Advanced Computer Architecture Laboratory, University of Southern California, September 1993.
51. Gorlick, M. M., and Kesselman, C. F., Gauge: A Workbench for the Performance Analysis of Logic Programs, in: *5th ICSLP*, MIT Press, Cambridge, MA, 1988, pp. 548–561.

52. Gudeman, D., Representing Type Information in Dynamically Typed Languages, Report TR93-27, University of Arizona, Department of Computer Science, September 1993.
53. Guo, Y.-K., and Lock, H. C. R., A Classification Scheme for Declarative Programming Languages, GMD-Studien Nr. 182, GMD, Germany, August 1990.
54. Habata, S., Nakazaki, R., Atarashi, A., and Umemara, M., Co-operative High Performance Sequential Inference Machine: CHI, in: *International Conference on Computer Design (ICCD '87)*, IEEE Computer Society Press, New York, 1987, pp. 601–604.
55. Haridi, S., and Janson, S., Kernel Andorra Prolog and its Computation Model, in: *7th ICLP*, MIT Press, Cambridge, MA, 1990, pp. 31–48.
56. Harsat, A., and Ginosar, R., CARMEL-2: A Second Generation VLSI Architecture for Flat Concurrent Prolog, in: *FGCS '88*, Tokyo, November 1988, pp. 962–969.
57. Hausman, B., Turbo Erlang, in: *10th ILPS*, MIT Press, Cambridge, MA, 1993, p. 662.
58. Haygood, R. C., Aquarius Prolog User Manual, in: *Aquarius Prolog 1.0 Documentation*, UC Berkeley, April 1993.
59. Hennessy, J. L., and Patterson, D. A., *Computer Architecture: A Quantitative Approach*, Morgan Kaufmann, Los Altos, CA, 1990.
60. Hermenegildo, M., Warren, R., and Debray, S., Global Flow Analysis as a Practical Compilation Tool, in: *J. Logic Program.* 13(4):349–367 (1992).
61. Hickey, T., and Mudambi, S., Global Compilation of Prolog, in: *J. Logic Program.* 7(3):193–230 (1989).
62. Hill, P. M., and Lloyd, J. W., The Gödel Programming Language, Technical Report CSTR-92-27, Department of Computer Science, University of Bristol, October 1992 (revised May 1993).
63. Holmer, B. K., Sano, B., Carlton, M., Van Roy, P., Haygood, R., Pendleton, J. M., Dobry, T. P., Bush, W. R., and Despain, A. M., Fast Prolog with an Extended General Purpose Architecture, in: *17th ISCA*, IEEE Computer Society Press, New York, 1990, pp. 282–291.
64. Holmer, B. K., Automatic Design of Computer Instruction Sets, Ph.D. dissertation, Department of Computer Science, UC Berkeley, 1993.
65. Huet, G., Résolution d'Équations dans des Langages d'Ordre 1, 2, . . . , ω , Thèse de Doctorat d'État, Université Paris VII, September 1976 (in French).
66. IBM, *IBM SAA AD/Cycle Prolog/MVS & VM Programmer's Guide and Language Reference*, Release 1, December 1992.
67. Jaffar, J., Efficient Unification over Infinite Terms, in: *J. New Generation Comput.* 2(3):207–219 (1984).
68. Janssens, G., and Bruynooghe, M., Deriving Descriptions of Possible Values of Program Variables by Means of Abstract Interpretation, in: *J. Logic Program.* 13(2-3):205–258 (1992).
69. Kantrowitz, M., *Prolog Resource Guide*, Regularly posted on Internet newsgroups comp.lang.prolog and comp.answers.
70. Klinger, S., and Shapiro, E., From Decision Trees to Decision Graphs, in: *NACLP90*, MIT Press, Cambridge, MA, 1990, pp. 97–116.
71. Klinger, S., Compiling Concurrent Logic Programming Languages, Ph.D. dissertation, Weizmann Institute, Rehovot, October 1992.
72. Komatsu, H., Tamura, N., Asakawa, Y., and Kurokawa, T., An Optimizing Prolog Compiler, in: *Logic Programming '86, Lecture Notes in Computer Science 264*, Springer, New York, 1986, pp. 104–115.
73. Korsloot, M., and Tick, E., Compilation Techniques for Nondeterminate Flat Concurrent Logic Programming Languages, in: *8th ICLP*, MIT Press, Cambridge, MA, 1991, pp. 457–471.
74. Köves, P., and Szeredi, P., Getting the Most Out of Structure-Sharing, in: *Collection of Papers on Logic Programming*, SZKI, Budapest, 1988, pp. 69–84 (revised November 1993).

75. Krall, A., and Neumerkel, U., The Vienna Abstract Machine, in: *PLILP '90, Lecture Notes in Computer Science 456*, Springer, New York, 1990, pp. 121–135.
76. Kurosawa, K., Yamaguchi, S., Abe, S., and Bando, T., Instruction Architecture for a High Performance Integrated Prolog Processor IPP, in: *5th ICSLP*, MIT Press, Cambridge, MA, 1988, pp. 1506–1530.
77. Kursawe, P., How to Invent a Prolog Machine, in: *3rd ICLP, Lecture Notes in Computer Science 225*, Springer, New York, 1986, pp. 134–148. Also in *J. New Generation Comput.* 5:97–114, (1987).
78. Le Charlier, B., Musumbu, K., and Van Hentenryck, P., A Generic Abstract Interpretation Algorithm and its Complexity Analysis (extended abstract), in: *8th ICLP*, MIT Press, Cambridge, MA, 1991, pp. 64–78.
79. Le Charlier, B., Degimbe, O., Michel, L., and Van Hentenryck, P., Optimization Techniques for General Purpose Fixpoint Algorithms: Practical Efficiency for the Abstract Interpretation of Prolog, in: *1993 Workshop on Static Analysis (WSA '93), Lecture Notes in Computer Science 724* Springer, New York, 1993, pp. 15–26.
80. Lindholm, T., and R. O'Keefe, A., Efficient Implementation of a Defensible Semantics for Dynamic Prolog Code, in: *4th ICLP*, MIT Press, Cambridge, MA, 1987, pp. 21–39.
81. Lindholm, T., Krall, A., et al., Net Talk: Term Comparisons with Variables, *ALP Newsletter* November:18–21 (1992). From Internet newsgroup comp.lang.prolog, July 1992.
82. Ludemann, P., xpProlog: High Performance Extended Pure Prolog, Master's thesis, University of British Columbia, 1988.
83. Maher, M. J., Logic Semantics for a Class of Committed-Choice Programs, in: *4th ICLP*, MIT Press, Cambridge, MA, 1987, pp. 858–876.
84. Maier, D., and Warren, D. S., *Computing with Logic—Logic Programming with Prolog*, Benjamin/Cummings, Menlo Park, CA, 1988.
85. Mariën, A., An Optimal Intermediate Code for Structure Creation in a WAM-based Prolog Implementation, Katholieke Universiteit Leuven, Belgium, May 1988.
86. Mariën, A., Janssens, G., Mulkers, A., and Bruynooghe, M., The Impact of Abstract Interpretation: An Experiment in Code Generation, in: *6th ICLP*, MIT Press, Cambridge, MA, 1989, pp. 33–47.
87. Mariën, A., and Demoen, B., A New Scheme for Unification in WAM, in: *ILPS*, MIT Press, Cambridge, MA, 1991, pp. 257–271.
88. Mariën, A., Improving the Compilation of Prolog in the Framework of the Warren Abstract Machine, Ph.D. dissertation, Katholieke Universiteit Leuven, September 1993.
89. Marriott, K., García de la Banda, M., and Hermenegildo, M., Analyzing Logic Programs with Dynamic Scheduling, in: *20th POPL*, ACM, New York, 1994.
90. Mehlhorn, K., and Tsakalidis, A., Data Structures, in: *Handbook of Theoretical Computer Science, Volume A: Algorithms and Complexity*, MIT Press/Elsevier, New York, 1990, Chap. 6, pp. 301–341,
91. Meier, M., Shallow Backtracking in Prolog Programs, Internal report, ECRC, Munich, Germany, February 1987.
92. Meier, M., Aggoun, A., Chan, D., Dufresne, P., Enders, R., Henry de Villeneuve, D., Herold, A., Kay, P., Perez, B., van Rossum, E., and Schimpf, J., SEPIA—An Extendible Prolog System, in: *Proceedings of the 11th World Computer Congress IFIP'89*, San Francisco, August 1989, pp. 1127–1132.
93. Meier, M., Compilation of Compound Terms in Prolog, in: *NACLP90*, MIT Press, Cambridge, MA, 1990, pp. 63–79.
94. Meier, M., Better Late Than Never, Internal report, ECRC, Munich, Germany, 1993, and in: *ICLP '93 Workshop on Practical Implementations and Systems Experience*, Budapest, Hungary, June 1993.
95. Mellish, C. S., *Automatic Generation of Mode Declarations for Prolog Programs*, draft, Department of Artificial Intelligence, University of Edinburgh, August 1981.

96. Mellish, C. S., and Hardy, S., Integrating Prolog in the POPLOG environment, in: J. A. Campbell, (ed.), *Implementations of PROLOG*, 1984, pp. 147–162.
97. Mellish, C. S., Some Global Optimizations for a Prolog Compiler, in *J. Logic Program.* 1:43–66 (1985).
98. Meyer, R., Private communication, Digital Equipment Corporation, Paris Research Laboratory, December 1993.
99. Miller, D., and Nadathur, G., Higher-Order Logic Programming, in *3rd ICLP, Lecture Notes in Computer Science 225* Springer, New York, 1986, pp. 448–462.
100. Mills, J. W., LIBRA: A High-Performance Balanced Computer Architecture for Prolog, Ph.D. dissertation, Arizona State University, December 1988.
101. Moss, C., and Bowen, K., (chairs), *International Conference on the Practical Application of Prolog*, ALP, London, April 1992.
102. Moss, C., and Roth, A., *The Prolog 1000 database*, Available through anonymous ftp from src.doc.ic.ac.uk in packages/prolog-progs-db/prolog1000.v1, August 1993.
103. Naish, L., *MU-Prolog 3.1db Reference Manual*, Computer Science Department, University of Melbourne, Melbourne, Australia, May 1984.
104. Naish, L., Negation and Quantifiers in NU-Prolog, in *3rd ICLP, Lectures Notes in Computer Science 225*, Springer, New York, 1986, pp. 624–634.
105. Naish, L., *Negation and Control in Prolog*, Ph.D. dissertation, University of Melbourne. also published as *Lectures Notes in Computer Science 238*, Springer, New York, 1986.
106. Naish, L., Parallelizing NU-Prolog, in: *5th ICSLP*, MIT Press, Cambridge, MA, 1988, pp. 1546–1564. Also published as Technical Report 87/17, Department of Computer Science, University of Melbourne.
107. Naish, L., Dart, P. W., and Zobel, J., The NU-Prolog Debugging Environment, in: *6th ICLP*, MIT Press, Cambridge, MA, 1989, pp. 521–536.
108. Nakashima, H., and Nakajima, K., Hardware Architecture of the Sequential Inference Machine: PSI-II, in: *SLP*, IEEE Computer Society Press, New York, 1987, pp. 104–113.
109. Neumerkel, U., Une Transformation de Programme Basée sur la Notion d'Équations entre Termes (in French). *J. Francophones Programmation Logique*, May:215–229 (1993).
110. Noyé, J., Benker, H., et al. *ICM3: The Abstract Machine*, Technical Report CA-19, ECRC, Munich, February 1987.
111. Noyé, J., An Overview of the Knowledge Crunching Machine, ECRC, Munich, Germany, 1993. Also in: *Emerging Trends in Database and Knowledge-base Machines*. IEEE Computer Society Press, New York. To appear.
112. O'Keefe, R. A., *The Craft of Prolog*, MIT Press, Cambridge, MA, 1990.
113. Palmer, D., and Naish, L., NUA-Prolog: An Extension to the WAM for Parallel Andorra, in: *8th ICLP*, MIT Press, Cambridge, MA, 1991, pp. 429–442.
114. Pereira, F. C. N., and Shieber, S. M., *Prolog and Natural-Language Analysis*, Lecture Notes Number 10, Center for the Study of Language and Information (CSLI), 1987.
115. Pichler, C., Prolog-Übersetzer, Master's thesis (Diplomarbeit), (in German) Institut für Praktische Informatik, Technische Universität Wien, November 1984.
116. Podelski, A., and Van Roy, P., The Beauty and the Beast Algorithm: Testing Entailment and Disentailment Incrementally, in: *10th ILPS*, MIT Press, Cambridge, MA, 1993 p. 653. Also Research Report, Digital Equipment Corporation, Paris Research Laboratory, 1993. To appear.
117. Ramamohanarao, K., Shepherd, J., Balbin, I., Port, G., Naish, L., Thom, J., Zobel, J., and Dart, P., The NU-Prolog Deductive Database System, *IEEE Trans. Data Eng.*, 10:4:10–19 (1987). Also in: *Prolog and Databases*, Ellis Horwood, London, 1988 and Technical Report 87/19, Department of Computer Science, University of Melbourne.
118. Sagonas, K., Swift, T., and Warren, D. S., XSB: An Overview of its Use and Implementation, Report, SUNY Stony Brook, October 1993. Available through anonymous ftp from cs.sunysb.edu in pub/TechReports/warren/xsb_overview.ps.Z.
119. Sahlin, D., and Carlsson, M., Variable Shunting for the WAM, in: *NACLP '90 Workshop on*

- Logic Programming Architectures and Implementations*, Austin, Texas, November 1990. Also available as SICS Research Report R91:07, March 1991.
120. Saraswat, V., Concurrent Constraint Programming Languages, Ph.D. dissertation, Carnegie-Mellon University, 1989. Revised version published by MIT Press, Cambridge, MA, 1992.
 121. Scowen, R., ISO Draft Prolog Standard (N72), ISO/IEC JTC1 SC22 WG17, June 1991.
 122. Seo, K., and Yokota, T., Pegasus: A RISC Processor for High-Performance Execution of Prolog Programs, in: *VLSI '87*, Vancouver, Canada, North-Holland, Amsterdam, 1987. pp. 261–274.
 123. Seo, K., and Yokota, T., Design and Fabrication of Pegasus Prolog Processor in: *VLSI '89*, Munich, Germany, North-Holland, Amsterdam, 1989, pp. 265–274.
 124. Seo, K., Study of a VLSI Architecture for the Logic Programming Language Prolog, Ph.D. dissertation, Keio University, March 1993 (in Japanese).
 125. Shapiro, E., The Family of Concurrent Logic Programming Languages, *ACM Comput. Surveys* 21(3):412–510 (1989).
 126. Singhal, A., and Patt, Y. N., A High Performance Prolog Processor with Multiple Functional Units, in: *16th ISCA*, IEEE Computer Society Press, New York, 1989.
 127. Singhal, A., *Exploiting Fine Grain Parallelism in Prolog*. Ph.D. dissertation, Report UCB/CSD 90/588, Department of Computer Science, UC Berkeley, August 1990.
 128. Sloman, A., The Evolution of Poplog and Pop-11 at Sussex University, in: J. A. D. W. Anderson (ed.), *POP-11 Comes of Age: The Advancement of an AI Programming Language*, Ellis Horwood, London, 1989, pp. 30–54.
 129. Srinivasa, V. P., Tam, J., Nguyen, T., Chen, C., Wei, A., Testa, J., Patt, Y., and Despaigne, A. M., VLSI Implementation of a Prolog Processor, in: *Stanford VLSI Conference*, March 1987.
 130. Srinivasa, V. P., Tam, J. V., Nguyen, T. M., Patt, Y. N., Despaigne, A. M., Moll, M., and Ellsworth, D., A CMOS Chip for Prolog, in *International Conference on Computer Design (ICCD '87)*, IEEE Computer Society Press, New York, 1987, pp. 605–610.
 131. Sterling, L., and Shapiro, E., *The Art of Prolog*, MIT Press, Cambridge, MA, 1986.
 132. Swift, T., and Warren, D. S., Compiling OLD T Evaluation: Background and Overview, SUNY Stony Brook Technical Report 92/04, 1992. Available through anonymous ftp from cs.sunysb.edu in pub/TechReports/warren/xwam-overview.dvi.Z, June 1993.
 133. Swift, T., and Warren, D. S., Performance of Sequential SLG Evaluation, Report, SUNY Stony Brook. Available through anonymous ftp from cs.sunysb.edu in pub/TechReports/warren/xsb-perf.ps.Z, November 1993.
 134. Taki, K., Yokota, M., Yamamoto, A., Nishikawa, H., Uchida, S., Nakashima, H., and Mitsuishi, A., Hardware Design and Implementation of the Personal Sequential Inference Machine (PSI), in: *FGCS '84*, 1984.
 135. Taki, K., Parallel Inference Machine PIM, in: *FGCS '92*, 1992.
 136. Tamura, N., Knowledge-Based Optimization in Prolog Compiler, in: *ACM/IEEE Computer Society Fall Joint Conference*, November 1986.
 137. Tarau, P., A Compiler and a Simplified Abstract Machine for the Execution of Binary Metaprograms, in: *Proceedings of the Logic Programming Conference '91*, ICOT, Tokyo, 1991, pp. 119–128.
 138. Tarau, P., *WAM-Optimizations in BinProlog: Towards a Realistic Continuation Passing Prolog Engine*, Technical Report, Université de Moncton, Canada, 1992.
 139. Taylor, A., Removal of Dereferencing and Trailing in Prolog Compilation, in: *6th ICLP*, MIT Press, Cambridge, MA, 1989, pp. 48–60.
 140. Taylor, A., LIPS on a MIPS: Results from a Prolog Compiler for a RISC, in: *7th ICLP*, MIT Press, Cambridge, MA, 1990, pp. 174–185.
 141. Taylor, A., High-Performance Prolog Implementation, Ph.D. dissertation, Basser Department of Computer Science, University of Sydney, June 1991.
 142. Thibault, O., Hardware evaluation of KCM, ECRC, Munich, Germany, May 1990 and in: *Tools for Artificial Intelligence 1990*, IEEE Computer Society Press, New York, 1990.

143. Tick, E., and Warren, D. H. D., Towards a Pipelined Prolog Processor, in: *SLP*, February 1984, pp. 29–40. Also in *J. New Generation Comput.* 2(4):323–345 (1984).
144. Tick, E., Memory- and Buffer-referencing Characteristics of a WAM-based Prolog, *J. Logic Program.* 12:133–162 (1991).
145. Tomasulo, R. M., An Efficient Algorithm for Exploiting Multiple Arithmetic Units, *IBM J.*, 11:25–33 (1967). Also in: Siewiorek, Bell, and Newell, (eds.), *Computer Structures: Principles and Examples*, McGraw-Hill, New York, 1982, Chap. 19, pp. 293–302.
146. Touati, H., and Despain, A. M., An Empirical Study of the Warren Abstract Machine, in: *SLP*, IEEE Computer Society Press, New York, 1987, pp. 114–124.
147. Turk, A. K., Compiler Optimizations for the WAM, in: *3rd ICLP, Lecture Notes in Computer Science 225* Springer, New York, 1986, pp. 657–662.
148. Uchida, S., Summary of the Parallel Inference Machine and its Basic Software, in: *FGCS '92*, 1992.
149. Van Roy, P., A Prolog Compiler for the PLM, M.S. Thesis, Report UCB/CSD No. 84/203, UC Berkeley, November 1984.
150. Van Roy, P., Demoen, B., and Willems, Y. D., Improving the Execution Speed of Compiled Prolog with Modes, Clause Selection, and Determinism, in: *TAPSOFT '87, Lectures Notes in Computer Science 250* Springer, New York, 1987, pp. 111–125. Also CW Report 51, K. U., Leuven.
151. Van Roy, P. An Intermediate Language to Support Prolog's Unification, in: *NACLP 89*, MIT Press, Cambridge, MA, October 1989, pp. 1148–1164.
152. Van Roy, P., A Useful Extension to Prolog's Definite Clause Grammar notation, *ACM SIGPLAN Not.* November:132–134 (1989). See also Extended DCG Notation: A Tool for Applicative Programming in Prolog, Report UCB/CSD 90/583, UC Berkeley, July 1990.
153. Van Roy, P., and Despain, A. M., The Benefits of Global Dataflow Analysis for an Optimizing Prolog Compiler, in: *NACLP 90*, MIT Press, Cambridge, MA, 1990, pp. 491–515.
154. Van Roy, P., Can Logic Programming Execute as Fast as Imperative Programming?, Ph.D. dissertation, Report UCB/CSD 90/600, Department of Computer Science, UC Berkeley, December 1990. Revised version published as *Fast Logic Program Execution* by Intellect Books.
155. Van Roy, P., and Despain, A. M., High-Performance Logic Programming with the Aquarius Prolog Compiler, *IEEE Computer* 25(1):54–68 (1992).
156. Van Roy, P., How to Get the Most Out of Aquarius Prolog, in: *Aquarius Prolog 1.0 documentation*, Digital Equipment Corporation, Paris Research Laboratory, April 1993.
157. Van Roy, P., et al. *Aquarius Benchmarks*, Available through anonymous ftp from gatekeeper.dec.com in pub/plan/prolog/AquariusBenchmarks.tar.Z.
158. Voda, P. J., Types of Trilogy, In *5th ICSLP*, MIT Press, Cambridge, MA, 1988, pp. 580–589.
159. Wadler, P., The Essence of Functional Programming, in: *19th POPL*, ACM Press, New York, 1992, pp. 1–14.
160. Warren, D. H. D., Applied Logic—Its Use and Implementation as a Programming Tool, Ph.D. dissertation, DAI Research Reports 39 & 40, University of Edinburgh, 1977; also SRI Technical Report 290, 1983.
161. Warren, D. H. D., Prolog on the DECsystem-10, in: *Expert Systems in the Micro-Electronic Age*, Edinburgh University Press, 1979.
162. Warren, D. H. D., and Pereira, F. C. N., An Efficient Easily Adaptable System for Interpreting Natural Language Queries, *Am. J. Computational Linguistics*, 8(3-4):110–122 (1982).
163. Warren, D. H. D., Higher-Order Extensions to Prolog—Are They Needed?, in: *Machine Intelligence 10*, Ellis Horwood, London, 1982.
164. Warren, D. H. D., *An Abstract Prolog Instruction Set*, Technical Note 309, SRI International Artificial Intelligence Center, October 1983.

165. Warren, R., Hermenegildo, M., and Debray, S. K., On the Practicality of Global Flow Analysis of Logic Programs, in *5th ICSLP*, MIT Press, Cambridge, MA, 1988, pp. 684–699.
166. Watzlawik, G., Benker, H., Noyé, J., et al. *ICM4*, Technical Report CA-25, ECRC, Munich, Germany, February 1987.
167. Yokota, T., and Seo, K., Pegasus—An ASIC Implementation of High-Performance Prolog Processor, in: *EURO ASIC '90*, Paris, France, IEEE Computer Society Press, New York, 1990, pp. 156–159.
168. Zhou, N.-F., *Backtracking Optimizations in Compiled Prolog*, Ph.D. dissertation, Kyushu University, Fukuoka, Japan, November 1990.
169. Zorn, B. G., Comparative Performance Evaluation of Garbage Collection Algorithms, Ph.D. dissertation, Report UCB/CSD 89/544, Department of Computer Science, UC Berkeley, December 1989.