

Parallelizing irregular and pointer-based computations automatically: Perspectives from logic and constraint programming

Manuel Hermenegildo

*School of Computer Science, Technical University of Madrid (UPM), 28660 Boadilla del Monte,
Madrid, Spain*

Abstract

Irregular computations pose some of the most interesting and challenging problems in automatic parallelization. Irregularity appears in certain kinds of numerical problems and is pervasive in symbolic applications. Such computations often use dynamic data structures, which make heavy use of pointers. This complicates all the steps of a parallelizing compiler, from independence detection to task partitioning and placement. Starting in the mid 80s there has been significant progress in the development of parallelizing compilers for logic programming (and more recently, constraint programming) resulting in quite capable parallelizers. The typical applications of these paradigms frequently involve irregular computations, and make heavy use of dynamic data structures with pointers, since logical variables represent in practice a well-behaved form of pointers. This arguably makes the techniques used in these compilers potentially interesting. In this paper, we introduce in a tutorial way, some of the problems faced by parallelizing compilers for logic and constraint programs and provide pointers to some of the significant progress made in the area. In particular, this work has resulted in a series of achievements in the areas of inter-procedural pointer aliasing analysis for independence detection, cost models and cost analysis, cactus-stack memory management, techniques for managing speculative and irregular computations through task granularity control and dynamic task allocation (such as work-stealing schedulers), etc.

1. Introduction

Multiprocessing hardware is already available, which offers significant advantages in either performance or cost/performance over uniprocessors. For example, departmental servers using fast, inexpensive off-the-shelf processors are currently offered at a fraction of the cost of the mainframes they replace, and even multiprocessor workstations are now not uncommon. Faster and more ubiquitous high-speed networks increase the potential of exploiting distributed execution.

One of the recurring facts that hamper the progress of widespread use of parallelism is that in practice, beyond some manually parallelized high volume applications and scientific codes, still comparatively few programs are written or transformed to exploit parallelism. The traditional argument that parallelization is a difficult and error-prone task (see, e.g., [52]) seems to remain valid [3], and still points to the necessity of improving the tools used in the process. This includes developing languages that offer better support for parallel programming, improved libraries for supporting parallel programming on conventional languages, and significant progress in support tools, from parallelizing compilers to performance analyzers.

Herein, we concentrate on the issue of automatic parallelization. While manual parallelization may of course always have a place, parallelizing compilers are interesting in that they have the potential to dramatically lessen the parallelization burden and there is hope that one day they may eliminate it altogether. However, despite much progress, it appears that significant challenges still remain in the area of automatic parallelization, including dealing well with both regular and irregular computations, performing efficient partitioning for both types of computations, dealing with data structures with pointers, handling speculative computations, automatically changing data structures for more efficient exploitation of parallelism, and developing parallelization techniques for new, higher level programming paradigms.

The goal of developing effective parallelizing compilers is being sought after concurrently and, unfortunately, somewhat independently in the context of different programming paradigms or even individual languages. As a result of the characteristics of the typical applications of such paradigms or languages, the amount of progress made on the different topics involved made differs.

For example, some very significant progress has been made in parallelizing compilers for regular, numerical computations, generally based on the FORTRAN language (see, e.g., [7,79]). This research has resulted in well-known concepts and techniques including a well-understood notion of independence (based on the Bernstein conditions or, for example, more recent notions of “semantic independence” [9]), sophisticated syntactic loop transformations, transformations based on polytope models, extensive work on partitioning and placement, etc. On the other

hand, the applicability of these techniques has remained comparatively limited for irregular or symbolic computations, and still few practical systems deal well with parallelization across procedure calls or with irregular computations. Also, the techniques used often rely on the relative cleanliness of FORTRAN as a programming language and additional work is needed in order to extend them to other mainstream languages like C or C++. These languages include features such as dynamic, recursive data structures and pointer manipulation which complicate the detection of independence among statements or procedure calls and much current work is aimed at developing the related independence analyzers. An important example is pointer aliasing analysis (see, e.g., [4,68], and their references).

We argue that, despite the apparent differences among imperative, functional, logic, constraint, and object-oriented languages, the fundamental issues being tackled are quite similar. Thus, we believe that progress towards more effective parallelizing compilers for all programming paradigms can be sped-up by cross fertilization of the results obtained in different paradigms. It is with this thought in mind (and without aspiring to being exhaustive, which is impossible given the space available and unnecessary to make the point) that we present in the following a brief overview of some of the problems, which appear in the area of automatic parallelization of logic and constraint programs and provide pointers to the some of the solutions and significant achievements of the area.

2. Logic and constraint programming

Due to space limitations, we will present only a brief overview of logic and constraint programming, specifically tailored to the objective of our presentation (the reader is referred for example to [6,50,56,72] for details). We warn the reader that this cannot in any way be considered a fair introduction to the topic, since we completely overlook aspects of logic and constraint programming, which are widely perceived as important. These include the declarative nature and the logical semantics: programs in these languages are often not only the coding of an algorithm, but also a logical statement of a problem, which is very close to a specification. In the following, we take a fully operational view – the same one that the parallelizing compiler takes.

The basic “statements” of a constraint logic program are *constraints*. Constraints relate (logical) *variables* (variable identifiers start with upper case while constants and data structure descriptors – functors, see later – start with lower case). Such variables can be *free*, or they can be *constrained* to a certain value or set of values. For example, the statement $X = Y + Z$ establishes that the given constraint must hold among those variables (we assume for example that the variables range over floating point numbers). Such constraints are kept in the *store*. Assume Y and Z have a “known” value at the time of executing this constraint (for example, the store contains $Y = 2$ and $Z = 3$). Then, the operational semantics of such a constraint is very similar to that in any other language: the statement implies an addition ($2 + 3$) and an “assignment” of the result (5) to X . This can also be seen as *telling* (posting)

the constraint $X = 5$. Assume instead that such values are not known. Then, executing the statement involves placing the constraint in the store for later solution if/when another constraint is executed. Sequences of constraints are separated by commas. Assume again an empty initial store and the sequence of constraints “ $Y = 2, X = Y + Z$ ”. After executing this sequence the store would contain “ $Y = 2, X = 2 + T1, Z = T1$ ”. Here, we are making the assumption that sequences of constraints execute sequentially in the order in which they appear and that the store is always kept as “fully solved” as possible and in a normalized form (see [50] for details).

Constraint logic programming also provides a method for *procedure abstraction*. For example, code segment (a) below:

$\text{foo}(Z, X) \text{ :- } Y = 2, \quad (a)$ $\quad \quad \quad X = Y + Z.$	$\text{main} \text{ :- } \text{foo}(K, W), \quad (b)$ $\quad \quad \quad K = 3,$ $\quad \quad \quad \text{write}(W).$
-----------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------

defines a two-argument procedure foo . A procedure defines a local dynamic invocation context in the usual way, i.e., upon entering the procedure Y is a new local variable, while X and Z are formal parameters. The calling regime is not unlike “call by reference” (see the discussion later about logical variables being essentially pointers). For example, the effect of calling $\text{foo}(3, W)$ is that upon return $W = 5$ is added to the calling context. Note that the procedure is syntactically not very different from what one would write in a functional or imperative language, and its behavior is essentially the same for calls such as $\text{foo}(3, W)$. However, the complete operational behavior of the constraint programming procedure is richer because it allows other “calling modes”. For example, a call to $\text{foo}(K, 5)$ succeeds and upon return $K = 3$ is added to the calling context. Furthermore, a call to $\text{foo}(K, W)$ also succeeds and upon return the constraint $W = 2 + K$ is added to the calling context. In some ways, the statements and procedures in constraint programs can be seen as “reversible” versions of their syntactic counterparts in conventional languages. Note that also the declarative meaning of such programs is richer because it defines a complete logical *relation* (rather than a function) among its arguments. Procedure calls can appear in the bodies of procedures interspersed with constraints. For example, code segment (b) above would produce “5” on the standard output.

Procedures can have multiple definitions, which represent different *alternatives*. Establishing a somewhat inaccurate parallel with conventional languages, a set of procedure definitions can be seen as an “undoable” form of case statement or conditional. When such a procedure is entered it is said to create a *choice*. Such alternatives are tried in the textual order in which they appear in the program, i.e., the first definition of a procedure is tried first and, if that results in a *failure*, then the next one is tried (again, we follow the default execution strategy used in most practical constraint programming languages). A failure occurs when a constraint is executed, which makes the store unsolvable (i.e., it is incompatible with the current state of the store). This is not unlike the case of a test evaluating to *false* in a

conditional. When a failure occurs, the system *backtracks* to the last choice left behind and tries the next alternative in that choice. Since procedure calls can be nested, a stack of choices is kept by the system. A choice is pushed on the stack every time a procedure with several alternatives is invoked. When a failure occurs, execution continues at the next alternative of the choice on top of the choice stack. When the last alternative of a choice is entered, the choice itself is popped from the stack.

For example, the following program:

<pre>main :- bar(K,W), K > 2, write(W).</pre>	<pre>bar(X,Y) :- X < 0, Y = -10. bar(X,Y) :- X >= 0, Y = 10.</pre>
--------------------------------------------------------------	--------------------------------------------------------------------------

prints “10”. The first alternative of `bar` is tried first, resulting in $W = -10$ and $K < 0$, but executing $K > 2$ produces a failure since the store now has no solution. After trying the second alternative of `bar`, $K > 2$ succeeds (the store is then $K > 2$, $W = 10$) and the program terminates after printing the value of W .¹

The following, slightly more interesting example defining the Fibonacci relation illustrates the use of recursion:

<pre>fib(0, 0). fib(1, 1).</pre>	<pre>fib(N, F1+F2) :- N>1, F1>=0, F2>=0, fib(N-1, F1), fib(N-2, F2).</pre>
----------------------------------	---------------------------------------------------------------------------------------------------------------------

(in this example we have used a more convenient syntax where input parameter normalization is done automatically by the system – i.e., “`fib(0, 0).`” is a shorthand for “`fib(X, Y) :- X=0, Y=0.`” and “`fib(N, F1+F2) :- ...`” a shorthand for “`fib(N, X) :- X=F1+F2, ...`”). Calling `fib(8, Y)` establishes $Y = 21$, and calling `fib(X, 21)` establishes $X = 8$. Calling `fib(X, Y)` produces as *alternatives* the constraints $(X=0, Y=0)$, $(X=1, Y=1)$, $(X=2, Y=1)$, etc.

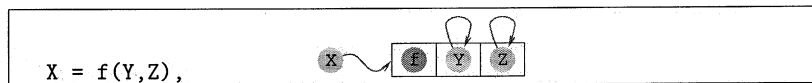
In the previous examples, we have been using a certain *constraint system*: essentially, equalities and inequalities involving linear arithmetic expressions over the (pseudo-) real-numbers. In many cases, the operations of constraint programs can be compiled directly into standard machine operations. However, in others (when actual constraint solving is involved) a *constraint solving algorithm* needs to be applied. Thus, the definition of each constraint system must include a decidable and (hopefully) efficient “solver”. Practical languages typically include several constraint systems.

¹ Of course, an optimizing compiler would compile away much of the behavior described in this very simple case.

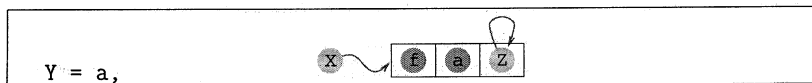
A particularly interesting constraint system present in almost all constraint languages is that of “equality relations over data structures” (i.e., finite trees). This is generally referred to as the *Herbrand domain* (and is the “working domain” of the Prolog language). This domain is crucial because it allows building and processing data structures with (single assignment) pointers in a very simple and declarative way. For example, the following program:

```
main :-
    X = f(Y, Z),
    Y = a,
    W = Z,
    W = g(K),
    X = f(a, g(b)).
```

first builds (dynamically) a new two-argument structure whose constructor symbol is f (in other words, a tree whose root node is f and which has two open branches)



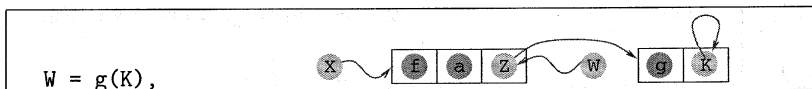
The variables Y and Z are *pointers* to the two arguments of the structure. The statement



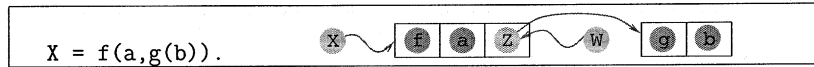
“binds” the first argument of the structure to the constant a (i.e., at this time X points to $f(a, Z)$). The following statement:



aliases the pointers W and Z (e.g., W points to Z). Therefore, the result of the statement



is to “bind” the second argument of the structure to $g(K)$ (and as a result X now points to $f(a, g(K))$). The last statement



finally binds X to the constant b . This last statement illustrates how open arguments inside a structure can also be accessed by traversing the structure using a process not unlike the “pattern matching” available in modern functional programming languages (except that it is again a “reversible” version of it). The algorithm capable of solving all such equality constraints over data structures is *unification* [57,61,66]. One of the nice characteristics of this constraint system is that there exist very efficient algorithms for performing unification.² As mentioned before, Prolog, one of the most popular logic programming languages, is essentially a constraint logic programming language, which uses exclusively the Herbrand domain. It is no surprise that Prolog is considered very well-suited for the easy manipulation of data structures with pointers.³

3. Parallelization of constraint logic programs

One of the main thesis of this paper is that logic programming and constraint programming languages offer a particularly interesting case study for the area of automatic parallelization. On one hand, these programming paradigms pose significant challenges to the parallelization task, which relate closely to the more difficult challenges faced in imperative language parallelization. Such challenges include highly irregular computations and dynamic control flow (due to the symbolic nature of many of their applications), non-trivial notions of (semantic) independence, the presence of dynamically allocated, complex data structures containing pointers, and having to deal with speculation.

On the other hand, due to their high-level nature these languages also facilitate the study of parallelization issues. As we have seen, logical variables are actually a quite “well-behaved” version of pointers, in the sense that no castings or pointer arithmetic (other than array indexing through the `arg/3` builtin) is allowed. Thus, pointers in these languages are not unlike those allowed, for example, in “clean” versions of C (or, to a lesser extent, in Java). In addition, similarly to functional

² Furthermore, there are also very successful compilation techniques, which (specially if global analysis of the program is performed) can translate sequences of operations such as those in the program above into a number of machine instructions that is essentially the same as if a lower-level language had been used to express the same data structure and pointer creation and binding operations. The reader is referred to [74] for an overview of progress in such compilation techniques.

³ Modern logic and constraint programming languages have many other features, such as support for higher order and meta programming, module and object systems, aggregation procedures, different sets of libraries, etc. with interesting implications on the automatic parallelization process. However, space limitations prevent us from considering these additional issues.

languages, logic and constraint languages allow coding in a way, which expresses the desired algorithm in a way that reflects more directly the structure of the problem (i.e., staying closer to the specifications). This makes the parallelism available in the problem more accessible to the compiler. The relatively clean semantics of these languages also makes it comparatively easy to use formal methods and prove the transformations performed by the parallelizing compiler both correct (in terms of computed outputs) and efficient (in terms of computational cost).⁴ Quite significant progress has been made in the past decade in the area of automatic program parallelization for logic programs and some of the challenges have been tackled quite effectively. In the following we touch upon a few of them (see, for example, [19] for an overview of the area).

3.1. Where the parallelism can be found

There are several types of parallelism, which are traditionally exploited in logic and constraint programs. For example, in applications involving extensive *search* (which is a frequent case in general search problems or in the enumeration part of constraint problems) the choices represented by alternative procedure definitions are often “deep” i.e., a number of steps are typically executed before a failure implies exploring an alternative definition. In this case different processors can execute simultaneously the different procedure definitions (i.e., the different branches of this *search space*). The resulting parallelism is called *or-parallelism*. This type of parallelism is present for example in the following program:

<pre> money(S,E,N,D,M,O,R,Y) :- digit(S), digit(E), ..., carry(I), ..., N is E+O-10*I, ..., </pre>	<pre> digit(0). digit(1). ... digit(9). carry(0). carry(1). </pre>
--------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------

The calls to `digit` and `carry` in the body of `money` are choices. Each alternative of these choices creates a branch that includes all the continuation (the rest of the body of `money` as well as the rest of the environment in which `money` was called). These branches can be executed in parallel.

⁴ Functional programming is another paradigm, which also facilitates exploitation of parallelism. However, it can be argued that the lack of certain features, such as pointers and backtracking, while making the parallelization problem easier, also precludes studying some interesting problems.

An alternative strategy is to parallelize the statements and/or procedure calls in procedure bodies, in the same way as in more traditional languages.⁵ This kind of parallelism is referred to as *and-parallelism*. A typical example of and-parallelism is the parallel execution of the two recursive calls in the definition of the Fibonacci relation given before. Another example is the following definition of the *quick-sort* program (where the functor “:” is used as list constructor) for example the two recursive calls to `qsort` can be executed in parallel:

```
qsort(nil, nil).
qsort(X:L, R) :-          partition(L, X, L1, L2),
                           qsort(L2, R2),
                           qsort(L1, R1),
                           append(R1, X:R2, R).
```

As and-parallelism corresponds to the traditional parallelism exploited in loop parallelization, divide and conquer algorithms, etc., we will concentrate our discussion on it. Also, and-parallelism is the only kind of parallelism that can be exploited in applications where choices are “shallow” (i.e., they correspond more closely to standard conditionals). It turns out that there are strong relationships between these forms of parallelism and the traditional notion of “data-parallelism” (see [10,11,41]).

3.2. Correctness and efficiency of the parallelization

As in any other programming paradigm, the objective of the parallelizing compiler is to uncover as much as possible of the available parallelism, while guaranteeing that the correct results are computed (*correctness*) and that other observable characteristics of the program, such as execution time, are improved (*speedup*) or, at the minimum, preserved (*no-slowdown*) – *efficiency*. A central issue is, then, under which conditions statements in a constraint logic program can be correctly and efficiently parallelized.

For comparison, consider the following segments of programs in (a) a traditional imperative language, (b) a (strict) functional language, and (c) a constraint logic programming language (we assume that the values of W and Z are initialized to some value before execution of these statements):

⁵ In fact, at a finer level of granularity, also *parts* of body statements can be executed in parallel. However, for simplicity, and without loss of generality, we assume parallelization at the *goal level*, meaning that the units scheduled will be body statements and procedure calls. Note also that the concurrency expressed by *concurrent logic programming languages* is between “and-parallel tasks”. See [42] for an extended discussion on this topic. Interesting models for exploiting and-parallelism at a finer level of granularity are, for example, [16,40,51,69,77].

s_1	$Y := W+2;$	$(+ (+ W 2)$	$Y = W+2,$
s_2	$X := Y+Z;$	$Z)$	$X = Y+Z,$
	(a)	(b)	(c)

For simplicity, we will reason about the correctness and efficiency of parallelism using the instrumental technique of considering reorderings (interleavings). Statements s_1 and s_2 in (a) are generally considered to be *dependent* because reversing their order would yield an *incorrect* result, i.e., it violates the *correctness* condition above (this is an example of a *flow-dependency*).⁶ A slightly different, but closely related situation occurs in (b): reversing the order of function application would result in a run-time error (one of the arguments to a function would be missing). Interestingly, reversing the order of statements s_1 and s_2 in (c) does yield the correct result. In fact, this is an instance of a more general rule: if no side effects are involved, reordering statements does not affect correctness in a constraint logic program. As another example, consider the following program (which uses only the Herbrand domain, i.e., it is a Prolog program, and which we will call program (d)):

main:-		p(X) :- X=a.
s_1	p(X),	q(X) :- X=b, <i>large computation</i> .
s_2	q(X),	q(X) :- X=a.
	write(X).	

Note that, again, reversing statements s_1 and s_2 produces the same result ($X = a$).

The fact that (at least in pure segments of programs) the order of statements in constraint logic programming does not affect the result⁷ led in early models to the proposal of execution strategies where parallelism was exploited “fully” (i.e., all statements were eligible for parallelization). However, the problem is that such parallelization often violates the principle of efficiency: for a finite number of processors, the parallelized program can be arbitrarily slower than the sequential program, even under ideal assumptions regarding run-time overheads. For instance, in the last example, reversing the order of the calls to p and q in the body of `main` implies that the call $q(X)$ (X at this point is free, i.e., a pointer to an empty cell) will first enter its first alternative, performing the large computation. Upon return of q

⁶ To complete the discussion above, note that output-dependencies do not appear in functional or logic and constraint programs because single assignment is generally used – we consider this a minor point of difference since one of the standard techniques for parallelizing imperative programs is to perform a transformation to a single assignment program before performing the parallelization.

⁷ Note that in practical languages, however, termination characteristics may change, but termination can actually also be seen as an extreme effect of the other problem to be discussed: efficiency.

(with X pointing to the constant b) the call to p will *fail* and the system will backtrack to the second alternative of q , after which p will succeed with $X = a$. On the other hand, the sequential execution would terminate in two or three steps, without performing the large computation. The fundamental observation is that, in the sequential execution, p *affects* q , in the sense that it *prunes* (limits) its choices. Executing q before executing p results in performing *speculative choices* with respect to the sequential execution. Note that this is in fact very related to executing conditionals in parallel (or ahead of time) in traditional languages (note that q above could also be (loosely) written as “ $q(X) : - \text{if } X = b \text{ then } \textit{largecomputation} \text{ else if } X = a \text{ then true else fail.}$ ”).

Something very similar occurs in case (c) above: while execution of the two constraints in the original order involves two additions and two assignments (the same of operations as those of the imperative or functional programs), executing them in reversed order involves first adding an equation to the system, corresponding to statement s_2 , and then solving it against s_1 , which is more expensive. The obvious conclusion is that, in general, arbitrary parallelization does not guarantee that the *two* conditions above are met.⁸

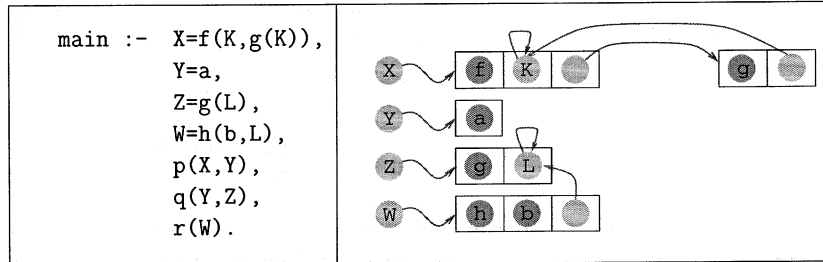
3.3. Notions of independence

Contrary to early the beliefs held in the field, most work in the last decade has considered that violating the efficiency condition is as much a “sign of dependence” among statements as violating the correctness condition. As a result, novel notions of independence have been developed, which capture these two issues of correctness and efficiency at the same time: independent statements as those whose run-time behavior, if parallelized, produces the same results as their sequential execution and an increase (or, at least, no decrease) in performance. As seen before, dealing with correctness is a matter of correctly sequencing side-effects (plus low-level issues, of course, such as locking). The techniques developed to this end are interesting, but, due to space limitations, we will concentrate on the arguably more interesting issue of guaranteeing efficiency. To separate issues better, we will discuss the issue under the assumption of ideal run-time conditions, i.e., no task creation and scheduling overheads (we will deal with overheads later). Note that, even under these ideal conditions, the statements in (c) and (d) are clearly *dependent*.

A fundamental question then is how to guarantee independence (without having to actually run the statements, as suggested by the definition given above). A fundamental result in this context is the fact that, if only the Herbrand constraint system is used (as in the Prolog language), a statement or procedure call, q , *cannot be*

⁸ In fact, a similar phenomenon occurs in or-parallelism where arbitrarily parallelizing branches of the search does not produce incorrect results, but, if looking for only one solution to a problem (or, more generally, in the presence of *pruning operators* – operators which control the search, which are pervasive in practical programs) results in speculative computations which can have a negative effect of efficiency. However, due to space limitations we concentrate our discussion on and-parallelism, because of its more direct relation to the parallelism that is usually exploited in conventional programs.

affected by another, p , unless there are free pointers (pointers to empty structure fields) from the run-time data structures passed to q from the data structures passed to p . This condition is called *strict independence* [30,45,47].⁹ For example, in the following program:



p and q are *strictly independent*, because, at the point in execution just before calling p (the situation depicted in the right part of the figure), X and Z point to data structures, which do not point to each other, and, even though Y is a pointer, which is shared between p and q , Y points to a fixed value, which p cannot change (note again that we are dealing with single assignment languages). As a result, the execution of p cannot affect q in any way and q can be safely run ahead of time in parallel with p (and, again assuming no run-time overheads, no-slowdown is guaranteed). Furthermore, no locking or copying of the intervening data structures is required (which helps bring the implementation closer to the ideal situation). Similarly, q and r are not strictly independent, because there is a pointer in common (L) among the data structures they have access to and thus the execution of q could affect that of r .

Unfortunately, the compiler cannot always determine independence by simply looking at one procedure, as above. For example, in the program (a) below

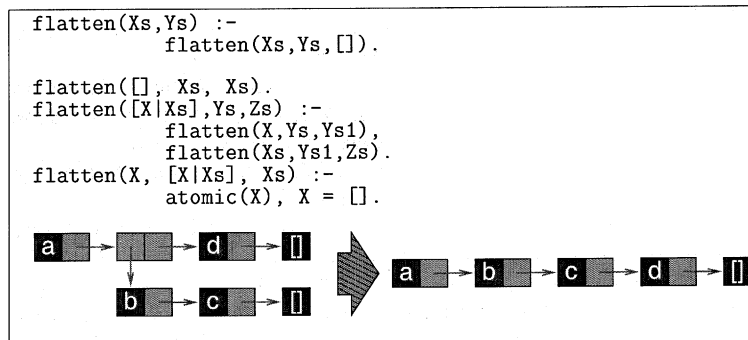
<pre> main :- t(X,Y), p(X), (a) q(Y). </pre>	<pre> main :- t(X,Y), (indep(X,Y) (b) -> p(X) & q(Y) ; p(X), q(Y)). </pre>
--------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------

it can determine that p and q are not (strictly) independent of t , since, upon entering the body of the procedure, X , Y , and Z are free pointers which are shared with t . On the other hand, after execution of t the situation is unknown since perhaps the structures created by t (and pointed to by X and Y) have no free pointers to each

⁹ To be completely precise, in order to avoid creating speculative parallelism, some non-failure conditions are also required of the goals executed in parallel, but we knowingly overlook this issue at this point to simplify the discussion.

other. Unfortunately, in order to determine this for sure a global (inter-procedural) analysis of the program must be performed. An alternative is to compile in a *run-time test* just after the execution of t . This has the undesirable side-effect that then the no-slowdown property does not automatically hold, because of the overhead involved in the test, but it is still potentially useful. The compilation of such a test can be seen as a source to source transformation of the program as shown in program (b) above (where, following the &-Prolog [43] notation, “&” represents parallel execution, and $(a \rightarrow b; c)$ is Prolog’s syntax for “(if a then b else c)”).

Furthermore, it is also sometimes possible to determine directly that in fact the operations that t performs on X and Y do not affect the execution of p and q . This kind of independence is called *non-strict independence* [46]. It cannot be determined in general *a priori* (i.e., by inspecting the state of the computation prior to executing t , p , and q) and thus necessarily requires a global analysis of the program. However, it is very interesting because it appears often in programs which manipulate “open” data structures (difference lists, dictionaries, etc.). An example of this is the following `flatten` example, which eliminates nestings in lists ($[X|Xs]$ represents the list whose head is X and whose tail is Xs and $[]$ represents the empty list):



This program unnests a list without copying by creating open-ended lists and passing a pointer to the end of the list ($Ys1$) to the recursive call. Since this pointer is not bound by the first call to `flatten/3` in the body of the recursive clause, the calls to `flatten(X, Ys, Ys1)` and `flatten(Xs, Ys1, Zs)` are (non-strictly) independent and *all the recursions can be run in parallel*.

An even more interesting case occurs if other constraint systems are used in addition to or in place of the Herbrand domain. Consider for example the parallelization of two procedure calls $p(X)$, $q(Z)$ in the following two situations:

- (a) `main :- X > Y, Z > Y, p(X) & q(Z), ...`
- (b) `main :- X > Y, Y > Z, p(X) & q(Z), ...`

In case (a) the store contains $(X > Y, Z > Y)$ before calling q and q , whereas in case (b) the store contains $(X > Y, Y > Z)$. The simple pointer aliasing reasoning implied by the definition of strict independence does not apply directly. However, p

cannot in any way affect q in case (a), while this could be possible in case (b), i.e., the two calls are clearly independent in case (a) while they are (potentially) dependent in case (b).

Notions of independence, which apply to general constraint programming (and can thus deal with the situation above) have been proposed recently [22,35]. For example, two goals p and q are independent if all constraints posed during the execution of q are consistent with the output constraints of p .¹⁰ The following is a sufficient condition for the previous definition but which only needs to look at the state of the store prior to the execution of the calls to be parallelized (for example, using run-time tests which explore the store c), in the same spirit as the strict-independence condition for the Herbrand case. Assuming the calls are $p(\bar{x})$ and $q(\bar{y})$ then the condition is

$$(\bar{x} \cap \bar{y} \subseteq \text{def}(c)) \text{ and } (\exists_{\bar{x}} c \wedge \exists_{\bar{y}} c \rightarrow \exists_{\bar{y} \cup \bar{x}} c),$$

where \bar{x} is the set of arguments of p , $\text{def}(c)$ the set of variables constrained to a unique value in c , and $\exists_{\bar{x}}$ represents the projection of the store on the variables \bar{x} (the notion of projection is predefined for each constraint system). The first condition states that the variables, which are shared between the goals in the program text must be bound at run-time to unique values. The second condition is perhaps best illustrated through an example. In the two cases above, for (a) $c = \{X > Y, Z > Y\}$ we have $\exists_{\{X\}} c = \exists_{\{Z\}} c = \exists_{\{X,Z\}} c = \text{true}$ and therefore p and q are independent. For (b) $c = \{X > Y, Y > Z\}$ we have $\exists_{\{X\}} c = \exists_{\{Z\}} c = \text{true}$ while $\exists_{\{X,Z\}} c = X > Z$ and therefore p and q are not independent. While checking these conditions accurately and directly can be inefficient in practice, the process can be approximated at compile-time via analysis or at run-time via simplified checks on the store.

Other interesting notions of independence which have been proposed are based on “determinacy” (i.e., lack of choices) [67]: two computations that have no choices (i.e., “do not backtrack”) are independent (provided, as before, that they can be guaranteed not to fail). Note that this is in general also captured by the notion of constraint independence given above.

3.4. The parallelization process

Experiments have shown that parallelization using only local analysis and generating run-time tests results in an excessive amount of overhead that severely limits speedups (see [15] for a recent comparison of actual speedups obtained by several parallelization methods). On the other hand, it has also been observed that there exist programs that obtain better speedups if a limited amount of run-time checking of independence is used than if only static decisions are made. Thus, a parallelization

¹⁰ This actually implies a better result even for Prolog programs since its projection on the Herbrand domain is a strict generalization of previous notions of non-strict independence, e.g., the sequence $p(X), q(X)$ can be parallelized if p is defined for example as $p(a)$ and q is defined as $q(a)$.

methodology is generally used which can accommodate both static analysis and run-time checking.

One of the more widely used approaches is illustrated in Fig. 1, representing the parallelization of “ $g_1(\dots), g_2(\dots), g_3(\dots)$ ”. The bodies of procedures are explored looking for statements and procedure calls, which are candidates for parallelization. As in many other parallelizers, a dependency graph is first built, which in principle reflects the total ordering of statements and calls given by the sequential semantics. To control the complexity of the process these graphs are limited to one body of one procedure (if the body is too long, the body can be partitioned in segments, but this does not happen often in constraint logic programs). Each edge in the graph is then labeled with the independence condition (the run-time check) that would guarantee independence of the statements or calls joined by the edge. A global analysis of the program then tries to prove these conditions statically true or false. If a condition is proved to be true, then the corresponding edge in the dependency graph is eliminated. If proved false, then an unconditional edge (i.e., a static dependency) is left. Still, in other edges conditions may remain (possibly simplified). The annotation process then encodes the resulting graph in the target parallel language (a variant of the source language). The techniques proposed for performing this process depend on many factors including whether the target language allows arbitrary parallelism or just fork-join structures and whether run-time independence tests are allowed or not. As an example, Fig. 1 presents two possible encodings in &-Prolog of the (schematic) dependency graph obtained after analysis. The parallel expressions generated in this case use only fork-join structures, one with run-time checks and the other one without them. Interesting techniques have been developed for compilation of *conditional* non-planar dependency graphs into fork-join structures, in addition to other, non-graph-based techniques [14,31,59].

The global analysis required to simplify the conditional graphs has to perform, among other tasks, inter-procedural pointer analyses, not unlike those recently

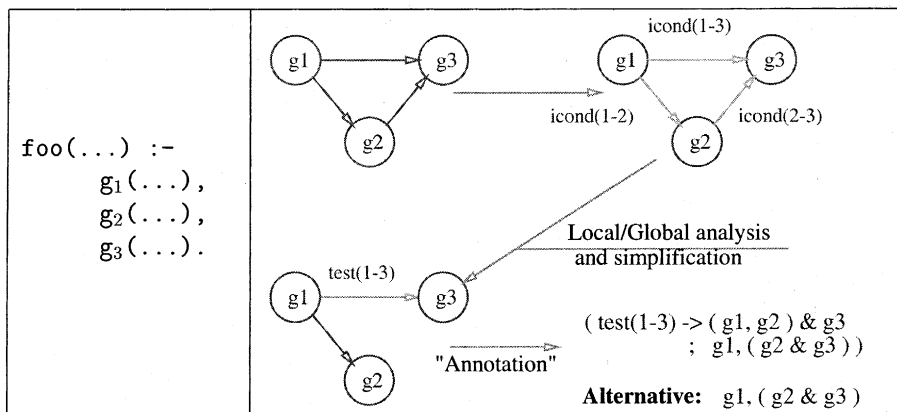


Fig. 1. Parallelizing “ $g_1(\dots), g_2(\dots), g_3(\dots)$ ”

proposed for clean versions of C or C++. Early proposals based on traditional data flow analysis techniques pointed in the right direction but proved imprecise [18]. The presence of recursion and dynamic data structures has fueled the development of quite sophisticated, incremental inter-procedural analyzers based on abstract interpretation [21]. This has required the development of efficient analysis algorithms as well as abstract domains for accurately and efficiently keeping track of sharing patterns and pointer aliasing in recursive data structures [15,49,58,60]. These analyses have been applied to the detection of both strict and non-strict independence [15,17] (for example, the `flatten` program of Section 3.3 is parallelized automatically by the system described in [17]). Analyses have been developed also to derive other important properties beyond variable instantiation states such as determinism [29], non-failure [26], and number of answers [13]. These parallelization techniques have also recently been extended to support “dependent” and-parallelism [63] (which, as mentioned before, really refers to exploiting independence at a finer level of granularity than goals [42]).

3.5. Dealing with overheads and irregularity – scheduling and memory management

The preceding discussion has on purpose avoided the issue of run-time overheads. The obvious practical implication of the existence of overheads (task creation, scheduling, data movement, etc.) is that even if a task is known to be independent, its parallel execution may still render a slow-down. This can happen if the task does not represent a sufficient amount of computation with respect to the overheads incurred in its parallelization. In the case of constraint logic programming the problem is compounded by the fact that, because of the symbolic nature of the applications typically coded, the number of tasks generated at run-time (as well as the computational cost and dynamic memory demands of each such task) depends on run-time parameters, i.e., the computations are typically highly irregular.

Two main approaches have been explored in order to overcome these problems. The first one is to combine dynamic task allocation policies with compilation techniques (abstract machines), which reduce as much as possible the overhead involved in the parallel execution of tasks. The best results have been obtained by performing low level “micro-task” scheduling, independently of the operating system threads [38,43,55], and generally based on non-centralized, “task stealing” approaches. Micro tasks are often represented simply by two pointers, one pointing to the procedure call or statement and another to the relevant invocation record. The tasks are executed by a number of instances of (a parallel version of) the conceptual abstract machines, which have been shown to provide the best performance for sequential implementation [1,37,55,75]. Interesting techniques have also been proposed for parallel dynamic memory management (using “cactus stacks” [2,12,37,44,55]). These techniques support, for example, efficient memory recovery during parallel back-tracking search. Some interesting examples of these dynamic scheduling and memory management techniques are presented in [37,43,62,64,71] for and-parallelism and in [2,20,32,55,76] for or-parallelism.

3.6. Dealing with overheads and irregularity – granularity control

The techniques mentioned above have proven sufficient for keeping the overheads of communication, scheduling, and memory management low and obtaining significant speedups in a wide variety of applications on shared memory multiprocessors (starting from the early paradigmatic examples: the sequent balance and symmetry series). However, current trends point towards larger multiprocessors but with less uniform shared memory access times. Controlling in some way the granularity (execution time and space) of the tasks to be executed in parallel can be a useful optimization in such machines, and is in any case a necessity when parallelizing for machines with slower interconnections. The latter include, for example, networks of workstations or distribution of work over the Internet.

This area of *granularity control* (task partitioning) has also received a significant amount of attention in the context of logic program parallelization. The idea of granularity control is to replace parallel execution with sequential execution or vice-versa based on knowledge (actual data, bounds, or estimations) of task size and overheads. The problem is challenging because, while the basic communication overhead parameters of a system can be determined experimentally, the computational cost of the tasks (e.g., procedure calls) being parallelized, as well as the amount of data that needs to be transferred before and after a parallel call, usually depend on dynamic characteristics of the input data. In the following example, we consider for parallel execution q (which assuming it is called with X bound to a list of numbers, adds one to each element of the list):

```
..., r(X) & q(X, Y), ...  
q([], []).  
q([I|Is], [I+1|Os]):- q(Is, Os).
```

The computational cost of a call to q (and also the communication overheads) are obviously proportional to the number of elements in the list. The characterization of input data required has made the problem difficult to solve (well) completely at compile-time.

One of the solutions which has been explored is to derive at compile time complexity cost functions, which give *upper and lower bounds* on task execution time as a function of certain measures of input data [24,25,27,28,53,54] (alternative solutions are given in, e.g., [70,73]; see also [48] in the context of functional languages). Interestingly, some of the analyses used in the derivation of such functions (e.g., [28]) make use of some techniques developed in the context of imperative program parallelization, such as the Omega test [65]. Programs are then transformed at compile-time into semantically equivalent counterparts but which automatically control granularity at run-time based on such functions. In the example above, these tools derive cost functions such as, for example, $2 * \text{length}(X) + 1$ for q (i.e., the unit of

cost is in this case a procedure call, where the addition is counted for simplicity as one procedure call). If we assume that we should parallelize when the total computation cost is larger than “100”, then we can transform the parallel call to p and q above into

$\dots, \text{Cost} = 2 * \text{length}(X) + 1, (\text{Cost} > 100$	$\rightarrow r(X) \ \& \ q(X, Y)$
	$; r(X),$
	$q(X, Y)), \dots$

(again, using an if-then-else). Clearly, many issues arise. For example, the cost of performing granularity control can be factored into the decisions. The cost functions can be simplified and related back to data structure sizes – list length in the case above, i.e., the call will only be parallelized if the length of the list is larger than a statically pre-computed value

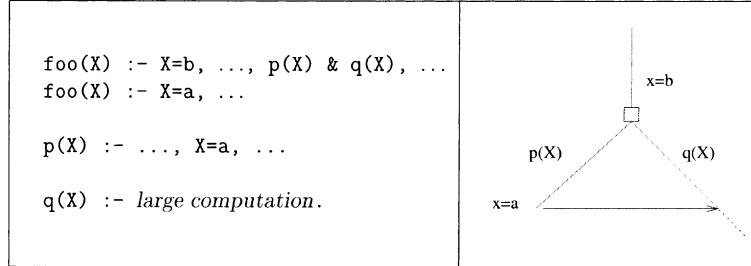
$\dots, (\text{length_greater_than} (X, 50)$	$\rightarrow r(X) \ \& \ q(X, Y)$
	$; r(X), q(X, Y)), \dots$

This in turn has inspired the development of algorithms for keeping track of data sizes at run-time. Also, the same techniques used for cost bounding allow deriving upper- and lower-bounds on the sizes of the structures being passed as arguments. This information can be factored into parallelization decisions (it affects the threshold). For example, in the example above, the argument size analysis (assuming that C is the cost of sending one element of a list, and a distributed setting where data is sent and returned eagerly) will infer that the communication cost is $2 * \text{length}(X) * C$. Interestingly, the *Computation > Overhead* condition ($2 * \text{length}(X) + 1 > 2 * \text{length}(X) * C$) can be determined statically to be always true (and parallelize unconditionally) or false (and never parallelize) depending only on the value of C , which in turn can perhaps be determined experimentally in a simple way. Performance improvements have been shown to result from the incorporation of this type of grain size control, specially for systems with medium to large parallel execution overheads [54]. Clearly, there are many interesting issues involved: techniques for derivation of data measures, data size functions, and task cost functions, program transformations, program optimizations, etc. Typically, the techniques are proved correct, again typically using the notions of approximation and bounding, formalized as abstract interpretations.

3.7. Dealing with speculation

Finally, also quite interesting techniques have been developed for controlling speculation, for both and- and or-parallelism. Explaining these issues in detail is

beyond the scope of this paper, but we will illustrate briefly with an example how speculation appears in and-parallelism



In the situation above, the first clause of `foo`, after binding `X` to `b`, executes `p` and `q` in parallel. However, the execution of `p` eventually fails when it poses the constraint `X = a` and execution must continue with the second clause of `foo`. Since `p` and `q` are in conjunction, the execution of `q` must now be discarded (i.e., starting `q` ahead of time was *speculative*). A combination of “left-biased scheduling” (ensuring that a processor has taken `p` before another can take `q`) and “instantaneous killing of siblings” (e.g., of `q` above) at least ensures no-slowdown [37,45,47]. No-slowdown (and even theoretical speedup) can also be guaranteed by determining statically that the tasks involved in a parallel conjunction (except the leftmost one) will not fail (techniques for this have been proposed in [26]). Many other interesting techniques for dealing with speculation have been developed (specially in the context of or-parallelism), including sophisticated schedulers, dynamic throttling of speculative tasks, etc. [8,26,36,38].

4. Conclusions: Towards cross-fertilization

As a result of the work outlined in previous sections, quite robust, publicly available compilers and run-time systems have been available for some time now, generally for Prolog, which automatically exploit parallelism in complex applications. Such systems have been shown to provide speedups over the state of the art sequential implementations available at the time of their development. The speed and robustness of these compilers has also been instrumental in demonstrating that abstract interpretation provides a very adequate framework for developing provably correct, powerful, and efficient global analyzers and, consequently, parallelizers [15,63,78]. More recently, techniques and practical tools have also been developed for the analysis of general constraint logic programs [34] as well as for their parallelization [33]. Prototypes incorporating the granularity control techniques mentioned above are also starting to be available. However, much work still remains to be done in these areas, and we believe there may be good opportunity at this time for increased transference of techniques across programming paradigms.

It can be argued that particularly strong progress has been made in the context of (constraint) logic programming in inter-procedural analysis of programs with

dynamic data structures and pointers, in parallelization using conditional dependency graphs (and possibly generating run-time independence tests), in the definition of the advanced notions of independence that are needed in the presence of speculative computations or languages, which include constraints, in the development of efficient task representation techniques and dynamic scheduling algorithms to deal with irregularity and speculation, and in the static inference of task cost functions for controlling granularity.

On the other hand, the techniques developed in the area of constraint logic program parallelization are certainly weaker than those developed in the context of numerical computing for regular problems. For example, logic programming parallelizers can discover the parallelism in complex recursive traversals of data structures, but do not handle well traversals that are based on integer (i.e., array subscript) arithmetic, for which much work exists in the area of imperative languages. Also, while current parallel constraint logic programming systems are reasonably good at dealing with tasks with dynamic costs, the techniques currently used are again comparatively weaker for the static case than the partitioning and placement algorithms used in imperative program parallelization [10,11,23,41]. Ideally, a parallelizing compiler should perform good partitioning and placement for any kind of architecture, using static techniques when possible and dynamic techniques when unavoidable. It thus appears that it would be quite interesting to merge the complementary work done in these areas by the different communities. Some progress has been made in one direction in the context of “data parallelism” [10,23,41], but it still seems like a very promising avenue for future research.

Constraint logic programming extends the high-level programming paradigm that logic programming offers in symbolic applications to numerical domains. We believe it offers a natural platform in which to study the combination of the parallelization techniques used in the numerical and symbolic programming fields. Independently of the convenience of using constraint programming languages directly (as is being done with significant commercial success in difficult problem areas such as scheduling or resource allocation), we also believe that many features of these languages, such as the use of constraints (“reversible statements”) or the embedded search capabilities, will slowly make their way into the designs of mainstream languages. In the same way, other features of symbolic languages (such as dynamic data structure creation and garbage collection, or bytecode compilation) have already made it into widely used languages such as Java. Current proposals in this direction include ILOG (a commercially successful library which extends C++ and Java with constraint handling capabilities) and [5], an imperative language with search capabilities.¹¹

¹¹ Of course, there are no scientific reasons not to use constraint logic languages directly, and this is indeed currently being done routinely with great commercial success by several companies working in difficult problem areas such as scheduling or resource allocation. However, it is entirely possible that the pure constraint logic programming languages, as so many other products of computer science, may remain powerful tools used by literate users, certainly making their impact on the mainstream, but in an indirect way.

References

- [1] H. Ait-Kaci, Warren's Abstract Machine, A Tutorial Reconstruction, MIT Press, Cambridge, MA, 1991.
- [2] K.A.M. Ali, R. Karlsson, The muse or-parallel prolog model and its performance, in: The 1990 North American Conference on Logic Programming, MIT Press, October 1990, pp. 757–776.
- [3] G. Almasi, A. Gottlieb (Eds.), Highly Parallel Computing, Benjamin Cummins, Menlo Park, CA, 1994.
- [4] L.O. Andersen, Binding-time analysis and the taming of C pointers, in: Proceedings of the Symposium on Partial Evaluation and Semantics-Based Program Manipulation, ACM Press, Copenhagen, Denmark, 1993, pp. 47–58.
- [5] K. Apt, A. Shaerf, Search and Imperative Programming, in: POPL'97: 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, ACM Press, Paris, France, January 1997, pp. 67–79.
- [6] K.R. Apt, Introduction to logic programming, in: J. van Leeuwen (Ed.), Handbook of Theoretical Computer Science, volume B: Formal Model and Semantics, Elsevier, Amsterdam and The MIT Press, Cambridge, MA, 1990, pp. 495–574.
- [7] D. Bacon, S. Graham, O. Sharp, Compiler transformations for high-performance computing, Computing Surveys 26 (4) (1994) 345–420.
- [8] T. Beaumont, D.H.D. Warren, Scheduling speculative work in or-parallel prolog systems, in: Proceedings of the 10th International Conference on Logic Programming, MIT Press, CA, June 1993, pp. 135–149.
- [9] E. Best, C. Lengauer, Semantic independence, Science of Computer Programming 13 (1990) 23–50.
- [10] J. Bevenmyr, T. Lindgren, H. Millroth, Exploiting recursion-parallelism in prolog, in: Proceedings of the PARLE'93, Springer, Berlin, 1993.
- [11] J. Bevenmyr, T. Lindgren, H. Millroth, Reform prolog: the language and its implementation, in: Proceedings of the 10th International Conference on Logic Programming, MIT Press, Cambridge, MA, 1993.
- [12] P. Borgwardt, D. Rea, Distributed semi-intelligent backtracking for a stack-based and-parallel prolog, in: International Symposium on Logic Programming, IEEE Computer Society, Silver Spring, MD, 1986, pp. 211–222.
- [13] C. Braem, B. Le Charlier, S. Modart, P. Van Hentenryck, Cardinality analysis of prolog, in: Proceedings of the International Symposium on Logic Programming, MIT Press, Ithaca, NY, November 1994, pp. 457–471.
- [14] F. Bueno, M. García de la Banda, M. Hermenegildo, A comparative study of methods for automatic compile-time parallelization of logic programs, in: The First International Symposium on Parallel Symbolic Computation, World Scientific Publishing Company, Singapore, September 1994, pp. 63–73.
- [15] F. Bueno, M. García de la Banda, M. Hermenegildo, Effectiveness of abstract interpretation in automatic parallelization: a case study in logic programming, ACM Trans. Program. Languages Syst. 21 (2) (1999) 189–238.
- [16] F. Bueno, M. Hermenegildo, U. Montanari, F. Rossi, Partial order and contextual net semantics for atomic and locally atomic CC programs, Sci. Comput. Program. 30 (1998) 51–82 Special CCP95 Workshop issue.

- [17] D. Cabeza, M. Hermenegildo, Extracting non-strict independent and-parallelism using sharing and freeness information, in: The 1994 International Static Analysis Symposium, number 864 in LNCS, Namur, Springer, Belgium, September 1994, pp. 297–313.
- [18] J.-H. Chang, A.M. Despain, D. Degroot, And-Parallelism of logic programs based on static data dependency analysis, in: Compcon Spring '85, February 1985, pp. 218–225.
- [19] J. Chassin, P. Codognet, Parallel logic programming systems, *Comput. Surveys* 26 (3) (1994) 295–336.
- [20] J. Chassin, J. Syre, H. Westphal, Implementation of a parallel prolog system on a commercial multiprocessor, in: Proceedings of the Ecai, August 1988, pp. 278–283.
- [21] P. Cousot, R. Cousot, Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints, in: The Fourth ACM Symposium on Principles of Programming Languages (1977) 238–252.
- [22] M. García de la Banda, F. Bueno, M. Hermenegildo, Towards independent and-parallelism in CLP, in: Programming Languages: Implementation, Logics, and Programs, number 1140 in LNCS, Springer, Aachen, Germany, September 1996, pp. 77–91.
- [23] S. Debray, M. Jain, A simple program transformation for parallelism, in: The 1994 International Symposium on Logic Programming, MIT Press, Cambridge, MA, November 1994, pp. 305–319.
- [24] S.K. Debray, N.-W. Lin, M. Hermenegildo, Task granularity analysis in logic programs, in: Proceedings of the 1990 ACM Conference on Programming Language Design and Implementation, ACM Press, New York, June 1990, pp. 174–188.
- [25] S.K. Debray, N.W. Lin, Cost analysis of logic programs, *ACM Trans. Program. Lang. Syst.* 15 (5) (1993) 826–875.
- [26] S.K. Debray, P. López-García, M. Hermenegildo, Non-failure analysis for logic programs, in: The 1997 International Conference on Logic Programming, MIT Press, Cambridge, MA, June 1997, pp. 48–62.
- [27] S.K. Debray, P. López-García, M. Hermenegildo, N.-W. Lin, Estimating the computational cost of logic programs, in: Static Analysis Symposium, SAS'94, number 864 in LNCS, Springer, Namur, Belgium, September 1994, pp. 255–265.
- [28] S.K. Debray, P. López-García, M. Hermenegildo, N.-W. Lin, Lower bound cost estimation for logic programs, in: The 1997 International Logic Programming Symposium, MIT Press, Cambridge, MA, October 1997, pp. 291–305.
- [29] S.K. Debray, D.S. Warren, Functional computations in logic programs, *ACM Trans. Program. Languages Syst.* 11 (3) (1989) 451–481.
- [30] D. DeGroot, Restricted AND-parallelism, in: The International Conference on Fifth Generation Computer Systems, Tokyo, November 1984, pp. 471–478.
- [31] D. DeGroot, A technique for compiling execution graph expressions for restricted AND-parallelism in logic programs, in: The International Supercomputing Conference, Springer, Athens, 1987, pp. 80–89.
- [32] European Computer Research Center, Eclipse User's Guide, 1993.
- [33] M. García de la Banda, F. Bueno, M. Hermenegildo, Towards independent And-parallelism in CLP, in: Programming Languages: Implementation, Logics, and Programs, number 1140 in LNCS, Springer, Aachen, Germany, September 1996, pp. 77–91.
- [34] M. García de la Banda, M. Hermenegildo, M. Bruynooghe, V. Dumortier, G. Janssens, W. Simoens, Global analysis of constraint logic programs, *ACM Trans. Program. Languages Syst.* 18 (5) (1996) 564–615.
- [35] M. García de la Banda, M. Hermenegildo, K. Marriott, Independence in CLP Languages, *ACM Trans. Program. Languages Syst.* 22 (2) (2000) 269–339.
- [36] B. Hausman, Handling speculative work in or-parallel prolog: evaluation results, in: North American Conference on Logic Programming, Austin, TX, October 1990, pp. 721–736.
- [37] M. Hermenegildo, An Abstract Machine for Restricted AND-parallel execution of Logic Programs, in: The Third International Conference on Logic Programming, number 225 in Lecture Notes in Computer Science, Imperial College, Springer, Berlin, July 1986, pp. 25–40.
- [38] M. Hermenegildo, Relating goal scheduling, precedence, and memory management in AND-parallel execution of logic programs, in: The Fourth International Conference on Logic Programming, University of Melbourne, MIT Press, Cambridge, MA, May 1987, pp. 556–575.

- [39] M. Hermenegildo, Automatic parallelization of irregular and pointer based computations: perspectives from logic and constraint Programming, in: *Proceedings of EUROPAR'97*, vol. 1300 of LNCS, Springer, Berlin, August 1997, pp. 31–46 (invited).
- [40] M. Hermenegildo, D. Cabeza, M. Carro, Using attributed variables in the implementation of concurrent and parallel logic programming systems, in: *Proceedings of the 12th International Conference on Logic Programming*, MIT Press, Cambridge, MA, June 1995, pp. 631–645.
- [41] M. Hermenegildo, M. Carro, Relating data-parallelism and (And)- parallelism in logic programs, *Comput. Languages J.* 22 (2/3) (1996) 143–163.
- [42] M. Hermenegildo and The CLIP Group, Some methodological issues in the design of CIAO – A generic, parallel, concurrent constraint system, in: *The Principles and Practice of Constraint Programming*, number 874 in LNCS, Springer, Berlin, May 1994, pp. 123–133.
- [43] M. Hermenegildo, K. Greene, The &-prolog system: exploiting independent And-parallelism, *New Generation Computing* 9 (3,4) (1991) 233–257.
- [44] M. Hermenegildo, R.I. Nasr, Efficient management of backtracking in AND-parallelism, in: *The Third International Conference on Logic Programming*, number 225 in LNCS, Imperial College, Springer, Berlin, July 1986, pp. 40–55.
- [45] M. Hermenegildo, F. Rossi, On the correctness and efficiency of independent And-parallelism in logic programs, in: 1989 North American Conference on Logic Programming, MIT Press, Cambridge, MA, October 1989, pp. 369–390.
- [46] M. Hermenegildo, F. Rossi, Non-Strict Independent And-parallelism, in: 1990 International Conference on Logic Programming, MIT Press, Cambridge, MA, June 1990, pp. 237–252.
- [47] M. Hermenegildo, F. Rossi, Strict and non-strict independent And-parallelism in logic programs: correctness, efficiency, and compile-time conditions, *J. Logic program.* 22 (1) (1995) 1–45.
- [48] L. Huelsbergen, J.R. Larus, A. Aiken, Using run-time list sizes to guide parallel thread creation, in: *Proceedings of the ACM Conference on Lisp and Functional Programming*, June 1994.
- [49] D. Jacobs, A. Langen, Accurate and efficient approximation of variable aliasing in logic programs, in: 1989 North American Conference on Logic Programming, MIT Press, Cambridge, MA, October 1989.
- [50] J. Jaffar, M.J. Maher, Constraint logic programming: a survey, *J. Logic Program.* 19/20 (1994) 503–581.
- [51] S. Janson S. Haridi, Programming Paradigms of the Andorra Kernel Language, in: 1991 International Logic Programming Symposium, MIT Press, Cambridge, MA, 1991, pp. 167–183.
- [52] A.H. Karp, R.C. Babb, A Comparison of 12 Parallel Fortran Dialects, *IEEE Software*, September 1988.
- [53] A. King, K. Shen, F. Benoy, Lower-bound time-complexity analysis of logic programs, in: 1997 International Logic Programming Symposium, MIT Press, Cambridge, MA, October 1997, pp. 261–275.
- [54] P. López-García, M. Hermenegildo, S.K. Debray, A methodology for granularity based control of parallelism in logic programs, *Journal of Symbolic Computation*, Special Issue on Parallel Symbolic Computation 22 (1996) 715–734.
- [55] E. Lusk, et al., The aurora Or-parallel prolog system, *New Generation Computing* 7 (2,3) (1990).
- [56] K. Marriot, P. Stuckey, *Programming with Constraints: An Introduction*, MIT Press, Cambridge, MA, 1998.
- [57] A. Martelli, U. Montanari, An efficient unification algorithm, *ACM Transactions on Programming Languages and Systems* 4 (3) (1982) 258–282.
- [58] K. Muthukumar, M. Hermenegildo, Determination of variable dependence information at compile-time through abstract interpretation, in: 1989 North American Conference on Logic Programming, MIT Press, Cambridge, MA, October 1989, pp. 166–189.
- [59] K. Muthukumar, M. Hermenegildo, The CDG, UDG, and MEL methods for automatic compile-time parallelization of logic programs for independent And-parallelism, in: *The International Conference on Logic Programming*, MIT Press, June 1990, pp. 221–237.
- [60] K. Muthukumar, M. Hermenegildo, Combined determination of sharing and freeness of program variables through abstract interpretation, in: 1991 International Conference on Logic Programming, MIT Press, Cambridge, MA, June 1991, pp. 49–63.

- [61] M.S. Paterson, M. Wegman, Linear unification, *Journal of Computer and System Sciences* 16 (2) (1978) 158–167.
- [62] E. Pontelli, G. Gupta, M. Hermenegildo, &ACE: A high-performance parallel prolog system, in: *International Parallel Processing Symposium*, IEEE Computer Society Technical Committee on Parallel Processing, IEEE Computer Society, Silver Spring, MD, April 1995, pp. 564–572.
- [63] E. Pontelli, G. Gupta, F. Pulvirenti, A. Ferro, Automatic compile-time parallelization of prolog programs for dependent And-parallelism, in: *Proceedings of the 14th International Conference on Logic Programming*, MIT Press, Cambridge, MA, July 1997, pp. 108–122.
- [64] E. Pontelli, G. Gupta, D. Tang, M. Carro, M. Hermenegildo, Improving the efficiency of nondeterministic And-parallel systems, *The Computer Languages Journal* 22 (2/3) (1996) 115–142.
- [65] W. Pugh, A practical algorithm for exact array dependence analysis, *Communications of the ACM* 35 (8) (1992) 102–114.
- [66] J.A. Robinson, A machine oriented logic based on the resolution principle, *Journal of the ACM* 12 (23) (1965) 23–41.
- [67] V. Santos-Costa, D.H.D. Warren, R. Yang, Andorra-I: a parallel prolog system that transparently exploits both And- and Or-parallelism, in: *Proceedings of the Third ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ACM, April 1991, *SIGPLAN Notices* vol 26(7), July 1991, pp. 83–93.
- [68] M. Shapiro, S. Horwitz, Fast and accurate flow-insensitive points-to analysis, in: *POPL'97: 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, ACM, Paris, France, January 1997, pages 1–14.
- [69] K. Shen, Overview of DASWAM: exploitation of dependent And-parallelism, *Journal of Logic Programming* 29 (1–3) (1996) 245–293.
- [70] K. Shen, V.S. Costa, A. King, Distance: a new metric for controlling granularity for parallel execution, in: Joxan Jaffar (Ed.), *Joint International Conference and Symposium on Logic Programming*, MIT Press, Cambridge, MA, 1998, pp. 85–99.
- [71] K. Shen, M. Hermenegildo, Flexible scheduling for non-deterministic, And-parallel execution of logic programs, in: *Proceedings of EuroPar'96*, number 1124 in LNCS, Springer, Berlin, August 1996, pp. 635–640.
- [72] L. Sterling, E. Shapiro, *The Art of Prolog*, MIT Press, Cambridge, MA, 1986.
- [73] E. Tick, Compile-Time Granularity Analysis of Parallel Logic Programming Languages, in: *International Conference on Fifth Generation Computer Systems*, Tokyo, November 1988.
- [74] P. Van Roy, 1983–1993: the wonder years of sequential prolog implementation, *Journal of Logic Programming* 19/20 (1994) 385–441.
- [75] D.H.D. Warren, An Abstract Prolog Instruction Set, Technical Report 309, Artificial Intelligence Center, SRI International, 333 Ravenswood Ave, Menlo Park CA 94025, 1983.
- [76] D.H.D. Warren, OR-Parallel Execution Models of Prolog, in: *Proceedings of TAPSOFT '87*, Lecture Notes in Computer Science, Springer, Berlin, March 1987.
- [77] D.H.D. Warren, The Extended Andorra Model with Implicit Control, Presented at ICLP'90 Workshop on Parallel Logic Programming, Eilat, Israel, June 1990 (unpublished).
- [78] R. Warren, M. Hermenegildo, S.K. Debray, On the practicality of global flow analysis of logic programs, in: *The Fifth International Conference and Symposium on Logic Programming*, MIT Press, Cambridge, MA, August 1988, pp. 684–699.
- [79] M. Wolfe, *High Performance Compilers for Parallel Computing*, Addison, Reading, MA, 1996.