

Parallel Logic Programming Systems on Scalable Architectures*

Vitor Santos Costa[†], Ricardo Bianchini, and Inês de Castro Dutra

Department of Systems Engineering and Computer Science
Federal University of Rio de Janeiro, Rio de Janeiro, Brazil
e-mail: {vitor,ricardo,ines}@cos.ufrj.br

Abstract

Parallel logic programming (PLP) systems are sophisticated examples of symbolic computing systems. They address problems such as dynamic memory allocation, scheduling irregular execution patterns, and managing different types of implicit parallelism. Most PLP systems have been developed for bus-based shared-memory architectures. The complexity of PLP systems and the large amount of data they process raises the question of whether logic programming systems can still obtain good performance on scalable architectures, such as distributed shared-memory systems.

In this work we use execution-driven simulation of a DASH-like architecture to investigate the access patterns and caching behaviour exhibited by a parallel logic programming system, Andorra-I. We first show that, without modifications, the system obtains reasonable performance, but that it does not scale well. By studying the behaviour of the major data structures in Andorra-I in detail, we conclude that this result is largely a consequence of the scheduling and work manipulation implementation used in the system. Our detailed analysis exposes several opportunities for improvements to Andorra-I, such as changing the layout of certain data structures. Based on the analysis of the caching behaviour of all Andorra-I data structures we optimised the Andorra-I code using 5 different techniques. We present the isolated and combined performance improvements provided by these optimisations. Our results show that the techniques provide significant performance improvements, leading to the conclusion that the system can and should perform well on modern scalable multiprocessors. Moreover, since Andorra-I shares its main data-structures with other PLP systems, we further conclude that the methodology and techniques used in our work can greatly benefit most PLP systems.

1 Introduction

Parallel logic programming (PLP) systems are sophisticated examples of symbolic computing systems. They address problems such as dynamic memory allocation, scheduling irregular execution patterns, and managing different types of implicit parallelism. In fact, one of the most important advantages of logic programming is the availability of several forms of implicit parallelism that can be naturally exploited on shared-memory multiprocessors. These forms include: or-parallelism, as exploited in the systems Aurora [20] and Muse [3]; independent and-parallelism, as in &-Prolog [14] and &-ACE [13]; dependent and-parallelism, as in Parlog's JAM [9], KLIC [30], and DASWAM [27, 28]; data-parallelism, as in Reform Prolog [6]; and combined and-or parallelism, as in Andorra-I [26] and Penny [21].

*This work was sponsored by CNPq, Brazilian Research Council.

[†]On leave from the Universidade do Porto, Portugal.

All these systems have been able to obtain good performance on bus-based systems, such as the Sequent Symmetry multiprocessors [19].

The complexity of PLP systems and the large amount of data they process raises the question of whether logic programming systems can still obtain good performance on scalable architectures, such as distributed shared-memory systems (DASH [18], Alewife [1], Exemplar [8] or Origin [17]). In distributed shared-memory systems, performance depends heavily on the read miss rates and may be limited by the communication overhead caused by coherence messages needed to maintain local copies of shared data consistent.

Sharing in parallel logic programming systems occurs under several circumstances. The use of logical variables for communication in dependent and-parallel applications, for instance, is an example of producer-consumer sharing of data, where the processor that instantiates a logical variable is the writer and the processors that consult the binding the readers.

A second major form of sharing, migratory sharing, arises from synchronisation between processors. Synchronisation occurs in tasks such as fetching work from other processors, and on being the leftmost goal or branch to execute cuts or side-effects. As an example, because of the high cost of suspending and restarting processes, it is very common that idle processors will be cycling through shared data structures searching for work. A processor that produces a piece of work writes to one of these data structures, which later will be read and then modified by one of the idle processors.

In this work we experiment with a sophisticated PLP system, Andorra-I [26], that exploits both dependent and-parallelism and or-parallelism. Andorra-I is one of the most sophisticated PLP systems that have been built to date. Original work with Andorra-I obtained good performance on the Sequent Symmetry, but our initial experience with Andorra-I running on modern multiprocessors demonstrates indeed that scalability suffers greatly on these architectures [24].

This paper addresses the question of whether the poor scalability of Andorra-I is inherent to the complex structure of PLP systems or can be improved through careful analysis and tuning. In order to answer this question, we analyse the caching behaviour of all Andorra-I data areas when applied to several different logic programs. The analysis pinpoints the areas that are responsible for most misses and the main sources for these misses. Based on this analysis we remove the main performance limiting factors in Andorra-I through a set of optimisations that did not require a redesign of the system. More specifically, we optimise Andorra-I using 5 different techniques: trimming of shared variables, data layout modification, privatisation of shared data structures, lock distribution, and elimination of locking in scheduling.

We present the isolated and combined performance improvements provided by the optimisations on a simulated DASH-like multiprocessor with up to 24 processors. In isolation, shared variable trimming and the modification of the data layout produced the greatest improvements. The improvements achieved when all optimisation techniques are combined are substantial. A few of our programs approach linear speedups as a consequence of our modifications. In fact, for one program the speedup of the modified Andorra-I is a factor of 3 higher than that of the original version of the system on 24 processors. Our main conclusion is then that, even though Andorra-I is indeed complex and irregular, it can and should scale well on modern scalable multiprocessors.

This conclusion can be extrapolated to several other PLP systems, since they share most of Andorra-I's data structures. More specifically, single processor execution in Andorra-

I is based on Warren’s WAM [33, 2], the abstract machine of most other sequential and PLP systems. Thus, Andorra-I’s major abstract machine data structures can be found in these systems. Parallel execution in Andorra-I applies the implementation technology developed for or-parallel systems such as Aurora [20], and for and-parallel systems such as PARLOG [9]. Thus, Andorra-I’s major parallelism-related data structures can be found in these PLP systems. Results obtained from Andorra-I are therefore of immediate interest to most PLP systems, making Andorra-I the most interesting system for studies such as ours.

Our approach contrasts with previous studies of the performance of parallel logic programming systems. Tick and Hermenegildo [31] studied caching behaviour of independent and-parallelism in bus-based multiprocessors. Other researchers have studied the performance of parallel logic programming systems on scalable architectures, such as the DDM [23]. Our previous work investigated the impact of different cache coherence protocols [24, 25], and the impact of different architectural parameters such as cache sizes and cache block sizes on the Andorra-I system [29].

The remainder of the paper is organised as follows. Section 2 presents the methodology used to obtain our results. We describe the Andorra-I PLP system, the simulator we used to perform the experiments, and how Andorra-I was ported to the simulator. In Section 3 we analyse in detail the initial speedup performance of each application we consider. Section 4 introduces the major optimisations and discusses their individual contribution. We present the performance of the optimised system in Section 5. Finally, Section 6 draws our conclusions and suggests future work.

2 Methodology

In this section we detail the methodology used in our experiments. The experiments consisted of the simulation of the parallel execution of Andorra-I, compiled for the MIPS architecture [16].

2.1 Multiprocessor Simulation

We use a detailed on-line, execution-driven simulator that simulates a 24-node, DASH-like [18], directly-connected multiprocessor. Each node of the simulated machine contains a single processor, a write buffer, cache memory, local memory, a full-map directory, and a network interface. The simulator was developed at the University of Rochester and uses the MINT front-end [32] (developed by Veenstra and Fowler) to simulate the MIPS architecture, and a back-end [7] (developed by Bianchini, Kontothanassis, and Veenstra) to simulate the memory and interconnection systems.

In our simulated machine, each processor has a 128-KB direct-mapped data cache with 64-byte cache blocks. All instructions and read hits are assumed to take 1 cycle. Read misses stall the processor until the read request is satisfied. Writes go into a 16-entry write buffer and take 1 cycle, unless the write buffer is full, in which case the processor stalls until an entry becomes free. Reads are allowed to bypass writes that are queued in the write buffers. Shared data are interleaved across the memories at the block level.

A memory bus clocked at half of the speed of the processor connects the main components of each machine node. A new bus operation can start every 34 processor cycles. A memory module can provide the first word of a cache line 20 processor cycles after the request is issued. The other words are delivered at 2 cycles/word bandwidth.

The interconnection network is a bi-directional wormhole-routed mesh, with dimension-ordered routing. The network clock speed is the same as the processor clock speed. Switch nodes introduce a 4-cycle delay to the header of each message. Network paths are 16-bit wide, which matches the memory bandwidth. In these networks contention for links and buffers is captured at the source and destination of messages.

All hardware characteristics mentioned above are common in actual modern parallel architectures.

In order to keep caches coherent we used a write-invalidate (WI) protocol [12]. In the WI protocol, whenever a processor writes a data item, copies of the cache block containing the item in other processors' caches are invalidated. If one of the invalidated processors later requires the same item, it will have to fetch it from the writer's cache. Our WI protocol keeps caches coherent using the DASH protocol with release consistency [18].

2.2 Andorra-I

The Andorra-I parallel logic programming system is based on the Basic Andorra Model [35]. The system was developed at the University of Bristol by Beaumont, Dutra, Santos Costa, Yang, and Warren [26, 36]. To the best of the authors' knowledge, Andorra-I was the first parallel logic programming system that exploited both and- and or-parallelism, and yet could run real-world applications with significant parallel performance.

Andorra-I employs a very interesting method for exploiting and-parallelism in logic programs, namely to execute *determinate* goals first and concurrently, where determinate goals are the ones that match at most one clause in a program. Thus, Andorra-I exploits determinate dependent and-parallelism. Eager execution of determinate goals can result in a reduced search space, because unnecessary choicepoints are eliminated. The Andorra-I system also exploits or-parallelism that arises from the non-determinate goals. Its implementation is influenced by JAM [10] when exploiting and-parallelism, and by Aurora [20] when exploiting or-parallelism.

The Andorra-I system consists of several components. The preprocessor is responsible for compiling the program and for the sequencing information necessary to maintain the correct execution of Prolog programs. The engine is responsible for the execution of the Andorra-I programs. The two schedulers manage and- and or-work. The reconfigurer allows workers to migrate between teams to find better sources of work.

A processing element that performs computation in Andorra-I is called a *worker*. In practice, each worker corresponds to a separate processor. Andorra-I is designed in such a way that workers are classified into *masters* and *slaves*. One master and zero or more slaves form a *team*. Each master in a team is responsible for creating a new choicepoint, while slaves are managed and synchronised by their master. Workers in a team cooperate with each other in order to share available and-work. Different teams of workers cooperate to share or-work. Note that workers arranged in teams share the same set of data structures used in a given branch of the search tree.

Most of the execution time of workers should be spent executing *engine* code [37], i.e., performing reductions. Andorra-I is designed in such a way that data corresponding to each worker is as local as possible, so that each worker tries to find its own work without interfering with others. Scheduling in Andorra-I is demand-driven, that is, whenever a worker runs out of work, it enters a *scheduler* to find another piece of available work.

The or-scheduler is responsible for finding or-work, i.e., an unexplored alternative in the

or-tree. Our experiments used the Bristol or-scheduler [5], originally developed for Aurora.

The and-scheduler is responsible for finding eligible and-work, which corresponds to a goal in the run queue (list of goals not yet executed) of a worker in the same team. Each worker in a team keeps a run queue of goals. This run queue of goals has two pointers. The pointer to the head of the queue is only used by the owner. The pointer to the tail of the queue is used by other workers to “steal” goals when their own run queues are empty. If all the run queues are empty, the slaves wait either until some other worker (in our implementation, the master) creates more work in its run queue or until the master detects that there are no more determinate goals to be reduced and it is time to create a choicepoint.

Finally, the reconfigurer [11] is responsible for arranging the workers into teams in a way that allows both and- and or-parallelism to be freely exploited when they are available.

2.3 The MIPS Port

In order to use Andorra-I with the simulator we needed to port the system to the MIPS architecture. We used the FSF’s `gcc 2.7.2` C compiler and `binutils-2.6` assembler and linker under a Solaris 2.5 environments as cross development tools for this purpose. Andorra-I was compiled with `-O2`.

Most of the port was straightforward. The only difficulties arose with shared memory allocation and locking. For shared memory allocation we use the `shmalloc` library supported by MINT. For locking in modern MIPS machines, we would use the `ll`, `sc`, and `sync` machine instructions [16] to implement locks and atomic operations. Unfortunately, these instructions are not yet supported by the back-end. The alternative is to use the lock library routines, as implemented by the simulator, which allow us to control the synchronisation overhead. In our simulations, these routines were implemented as atomic instructions.

To ensure correct execution under the release consistency model, we guarantee that all accesses to shared data are surrounded by lock and unlock operations. The only exception to this rule is the detection of the end of the determinate phase. Here we maintained the original protocol because any action after detection requires the slaves to grab a lock previously released by the master.

3 Analysing the Caching Behaviour of Andorra-I

The major goal of our work was to determine which factors affect caching in Andorra-I, as this is one of the most important requirements for obtaining good performance in parallel systems. A parallel logic programming system includes several shared data structures, whose importance depends on the type and amount of parallelism available in the application. Our first goal was therefore to classify them and study their individual contributions to the overall caching behaviour.

3.1 Data Areas in Andorra-I

Andorra-I combines the techniques used to implement parallel committed-choice systems [9] and the techniques used to implement or-parallel systems [20], themselves based on sequential Prolog [33]. The shared memory areas that implement the functionality of Andorra-I can be classified as:

- **Or-Scheduler Data Structures.** The Bristol Or-scheduler uses three different data structures: a set of fields on each choice-point, and two data structures with global variables. One of the two latter structures is shared by every worker, and the other is replicated by every worker. We grouped these data structures into a single area.
- **Worker Data Structures.** These include a copy of the abstract machine registers, the variables used to synchronise within a team, and the run-queues for each worker. Most accesses to these data structures occur when performing and-scheduling, that is when fetching work from other workers in a team, or when synchronising workers in a team.
- **Lock Array.** The lock array is used to establish a mapping between a shared memory position (such as a variable in the heap) and a lock. It is used whenever we need to lock a variable. This area is required because the simulator does not implement the MIPS instructions to synchronise access to shared memory. We therefore use hashing to a lock array as the mechanism to guarantee synchronisation.
- **Code Space.** The code space includes the compiled code for every procedure. During execution of the benchmarks it is read-only, but it could be updated by programs that perform `assert` or `record`.
- **Heap Space.** As in other Prolog systems, this space stores structured terms and variables. It grows during forward execution, and contracts during backtracking.
- **Goal Frame Space.** This data area stores goal frames, which consist of the goal's arguments, multi-assignment variables used to link the goal frames, and several control fields. The engine tries to reuse goal frames during forward execution.
- **Choicepoint Stack.** Choicepoints include pointers to the top of stacks, and flags that are updated as processors move around searching for work. Differently from Prolog, no arguments need to be stored. Note that as we backtrack and move forward, the same choice point stack space may be used several times for different choicepoints.
- **Trail Stack.** The trail stack records any conditional binding to logical or multi-assignment variables (used in the implementation of Andorra-I [37]). It can be considered as an extension of a choicepoint that records non-determinate bindings and is only important in applications with or-parallelism.
- **Binding Arrays.** They are used to implement the SRI model [34] for or-parallelism, by storing conditional bindings. They should be private in or-parallel applications, and almost never written in and-parallel-only applications. Even during an execution without choice points (a determinate execution), Andorra-I will access this data structure to verify whether variables have conditional bindings.
- **Miscellaneous Shared Variables.** This area includes the shared I/O data structures, such as open stream descriptors. It also includes information on the clauses compiled into the system, and pointers to the determinacy and non-determinacy code. A set of flags and counters is globally manipulated by the system. This area also includes the data structures used by the reconfigurer.

We instrumented the simulator to report data separately on each region. We then ran the applications to collect statistics for each data area in turn.

3.2 Workload and Original Performance

The benchmarks we used in this work are applications representing predominantly and-parallelism, predominantly or-parallelism, and both and- and or-parallelism. We next discuss application performance for the original Andorra-I. Note that our results correspond to the *first* run of an application; results would be somewhat better for other runs.

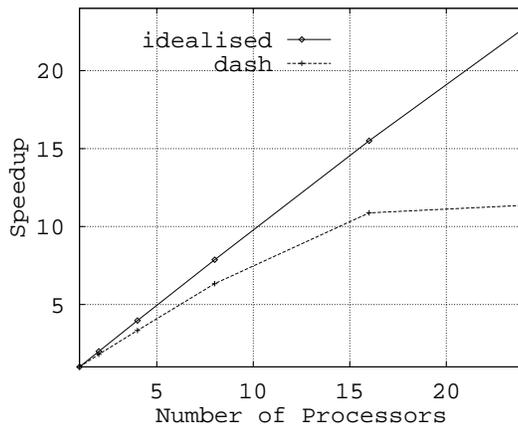


Figure 1: Speedups for `bt-cluster`

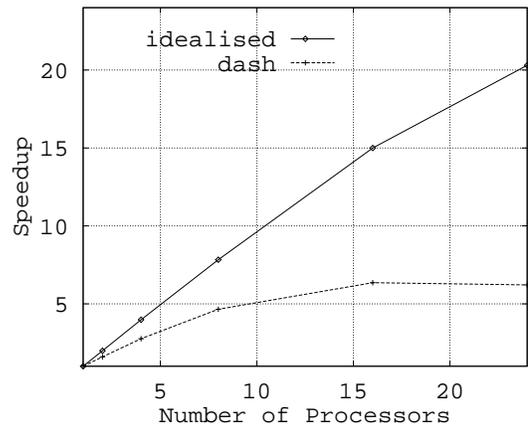


Figure 2: Speedups for `tsp`

And-Parallel Applications. We use two example and-parallel applications, the clustering algorithm for network management from British Telecom, `bt-cluster`, and a program to calculate approximate solutions to the traveling salesperson problem, `tsp`. The clustering program receives a set of points in a three dimensional space and groups these points into *clusters*. Basically, three points belong to the same cluster if the distance between them is smaller than a certain limit. And-parallelism in this case naturally stems from running the calculations for each point in parallel. The test program uses a cluster of 400 points as input data. This program has very good and-parallelism, and, being completely determinate, no or-parallelism. The traveling salesperson program is based on a Reform Prolog [6] benchmark that finds an approximate solution for the TSP problem in a graph with 24 nodes. To obtain best performance, the Andorra-I team rewrote the original applications to make them determinate-only computations.

Figure 1 shows the `bt-cluster` speedups for the simulated architecture as compared to an idealised shared-memory machine, where data items can always be found in cache. The `idealised` curve shows that the application has excellent and-parallelism and can achieve almost linear speedups up to twenty four processors. Unfortunately, performance for the DASH-like machine is barely acceptable. The DASH curve starts with an efficiency of 85% for

2 processors. Efficiency smoothly decreases as the number of processors increases, reaching 70% for 16 processors. Efficiency is less than 50% for 24 processors. Figure 2 shows that the `tsp` application achieves worse speedups than `bt-cluster` on a modern multiprocessor. The maximum speedup actually decreases for 24 processors, whereas the ideal machine would achieve a speedup of 20 for 24 processors.

Figure 3 illustrates the number and sources of cache misses per data area in the `bt-cluster` application running on 16 processors as a representative example. The figure shows that the overall miss rate of `bt-cluster` is dominated by true and false sharing misses from the `Worker` and `Misc` areas. This suggests that the system could be much improved by reducing false sharing and studying activity in the `Worker` and `Misc` areas.

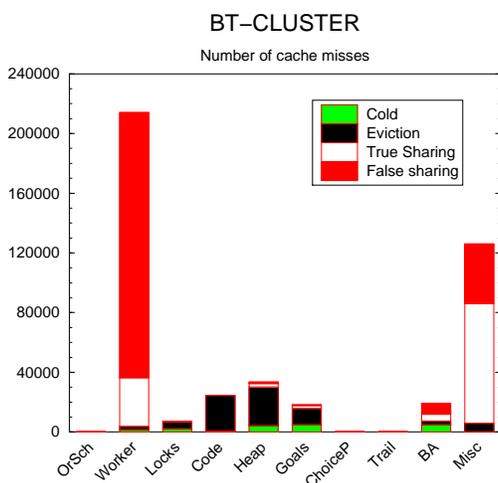


Figure 3: Misses by data area for `bt-cluster`

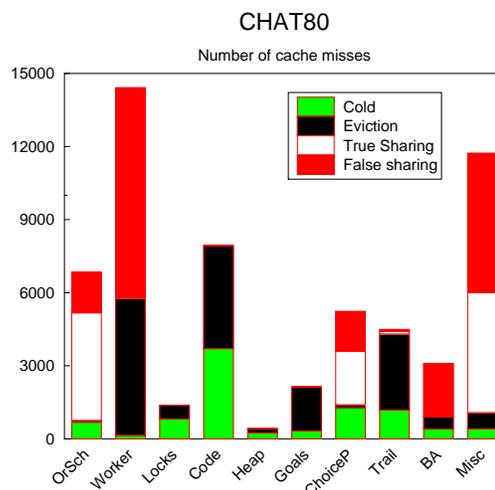


Figure 4: Misses by data area for `chat80`

Or-Parallel Applications. We use two or-parallel applications. Our first application, `chat80`, is an example from the well-known natural language question-answering system `chat-80`, written at the University of Edinburgh by Pereira and Warren. This version of `chat-80` operates on the domain of world geography. The program `chat80` makes queries to the `chat-80` database. This is a small scale benchmark with good or-parallelism, and it has been traditionally used as one of the or-parallel benchmarks for both the Aurora and Muse systems. The second application, `fp`, is an example query for a knowledge-based system for the automatic generation of floor plans. This application should at least in principle have significant or-parallelism. Figure 5 shows the speedups for the `chat80` application from 1 to 24 processors. These speedups are very similar to those obtained by Andorra-I on the Sequent Symmetry architecture. In contrast, the DASH curve reaches a maximum speedup of 4.2 for 16 processors. Figure 6 shows the speedups for the `fp` application. The theoretical

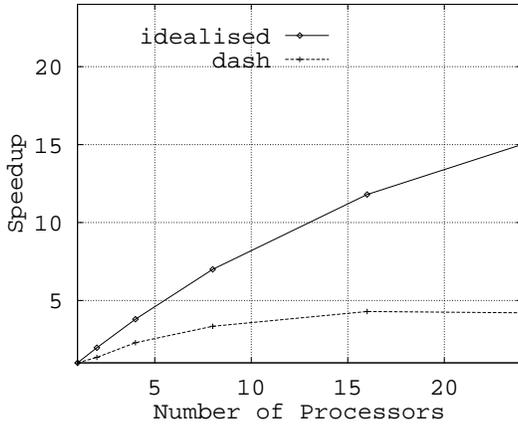


Figure 5: Speedups for `chat80`

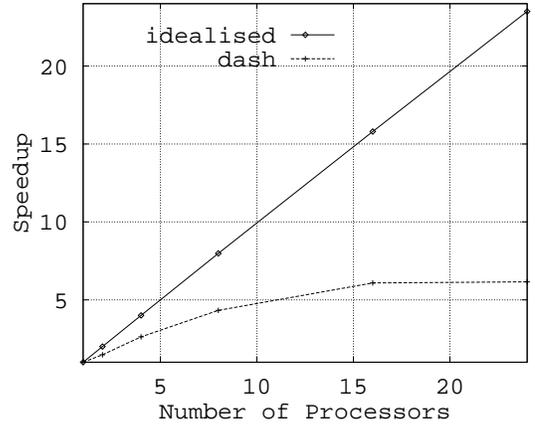


Figure 6: Speedups for `fp`

speedup is very good, in fact quite close to linear, in sharp contrast to the actual speedup for the DASH-like machine. Figure 4 shows the number and source of misses for `chat80` running on 16 processors, again as an example of this type of application. Note that `chat80` does not have enough parallelism to feed 16 processors, suggesting that most sharing misses should result from or-parallel scheduling areas, `OrSch` and `ChoiceP`. Indeed, the figure shows that these areas are responsible for a large number of sharing misses, but the areas causing the most misses are `Worker` and `Misc` as in the and-parallel applications, indicating again that these two areas should be optimised.

Note also the large number of eviction misses in the `Trail`. The trail is used to store conditional bindings, and is used in backtracking and when searching for work. The high number of misses in this area indicates that workers often go here, and suggests work is fine-grained.

And/Or-Parallel Application. As an example of and/or application we used a program to generate naval flight allocations, based on a system developed by Software Sciences and the University of Leeds for the Royal Navy. It is an example of a real-life resource allocation problem. The program allocates airborne resources (such as aircraft) whilst taking into account a number of constraints. The problem is solved by using the technique of active constraints as first implemented for Pandora [4]. In this technique, the co-routining inherent to the Andorra model is used to activate constraints as soon as possible. The program has both or-parallelism, arising from the different possible choices, and and-parallelism, arising from the parallel evaluation of different constraints. To test the program, we ask it to schedule 11 aircrafts, 36 crew members and 10 flights. The degree of and- and or-parallelism in this program varies according to the queries, but this query gives rise to more and-parallelism than or-parallelism. The system uses the reconfigurer to dynamically adapt to

the available sources of parallelism.

Figure 7 shows the speedups for `pan2`. The `idealised` curve shows that the application has less parallelism than all other applications; the ideal speedup does not even reach 12 on 24 processors. When run on the DASH simulator, `pan2` exhibits unacceptable speedups for all numbers of processors; speedup starts out at about 1.8 for 2 processors and slowly improves to a maximum of only 4.8 for 24 processors. Figure 8 shows the distribution of cache misses by the different Andorra-I data areas for 16 processors. In this case, the `Worker` area clearly dominates, since the contribution from the `Misc` area is not as significant as in the and-parallel benchmarks. Note that there is more true than false sharing activity in `Worker`. The true sharing probably results from idle processors looking for work.

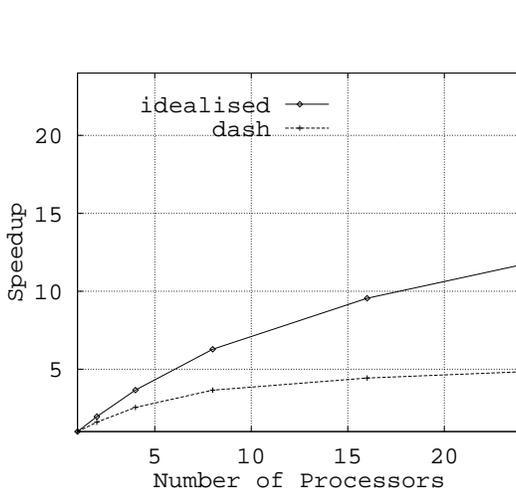


Figure 7: Speedups for `pan2`

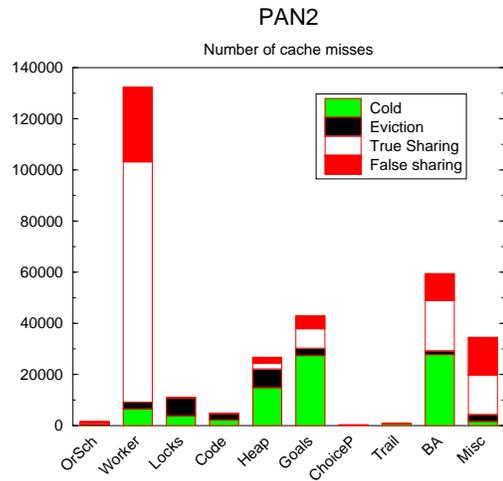


Figure 8: Misses by data area for `pan2`

4 Optimisation Techniques and Performance

The previous analysis suggests that relatively high miss rates may be causing the poor scalability of Andorra-I. It is interesting to note that most misses come from fixed layout areas, such as `Worker` and `Misc`, and not from the execution stacks, as one would assume.

We next discuss how several optimisations can be applied to the system, particularly in order to improve the utilisation of the `Worker` and `Misc` areas. The first two optimisations were prompted by our simulation-based analysis of caching behaviour, and they are the ones that give the best improvement. The other three were based on our original intuitions regarding the system, and were the ones we would have performed first without simulation data. We studied performance for three applications, `bt-cluster`, `chat80` and `pan2`.

In the remainder of this section, we study the impact of each of the techniques studied when applied in isolation.

Variable Trimming. In this technique we investigate the areas that have unexpected sharing, and try to eliminate this sharing if possible. For two of our applications, the `Misc` area gave a surprisingly significant contribution to the number of misses. The area is mostly written at the beginning of the execution to set up execution parameters. During execution it is used by the reconfigurer and to keep reduction and failure counters. By investigating each component in the area, we detected that the counters were the major source of misses. As they are only used for research and debugging purposes, we were able to eliminate them from the Andorra-I code.

Figure 9 shows the result of this optimisation for three benchmarks, `bt-cluster`, `chat80` and `pan2`. The figure plots the variation in speedup as a function of the number of processors. The results show that `chat80` benefits the most from this optimisation. This is because the failure counter is never updated by and-parallel applications, but often updated by this or-parallel application. The optimisation does not impact the and-parallel benchmarks as much, leading to less than 10% speedup improvements.

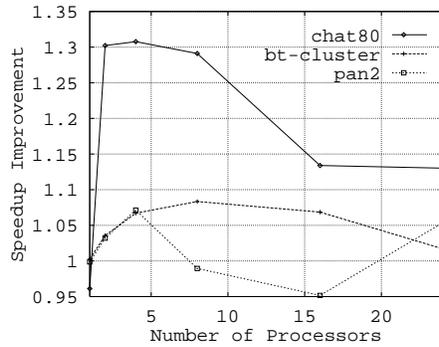


Figure 9: Impact of Counter Elimination Optimisations

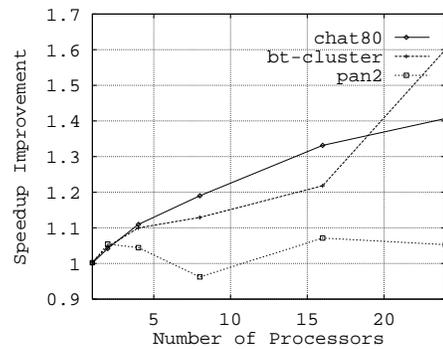


Figure 10: Impact of Data Layout Optimisations

Data Layout Modification. Figure 10 shows the impact of the data layout modifications on the speedup of our applications. The `bt-cluster` and `chat80` applications benefit the most from this optimisation. Note that, on machine configurations of up to 16 processors, `chat80` benefits from this optimisation more than `bt-cluster` does, suggesting that the false sharing was having a big negative impact on this application until its lack of enough parallelism started to take over. On the largest machine configuration however, the optimisation actually becomes more beneficial to `bt-cluster`, an and-parallel application with a much greater amount of parallelism.

The `pan2` application benefits only slightly from this optimisation. Note that we again have a slowdown, now for 4 processors, probably because of slightly different scheduling.

All benchmarks but `pan2` exhibit a high false sharing rate, showing a need for this technique. On 16 processors, 15% of the misses in `pan2` are false sharing misses, whereas in the other applications false sharing causes between 40% (`chat80`) and 51% (`bt-cluster`) of all misses. These results suggest that improving false sharing is of paramount importance. According to our detailed analysis of caching behaviour, false sharing misses are concentrated in the `Worker`, `OrSch`, `ChoiceP` and `BA` areas, besides the `Misc` area optimised by the previous technique.

The `Worker` and `OrSch` data areas are allocated statically. This indicates that we can effectively reduce false sharing. We applied two common techniques to tackle false sharing, *padding* between fields that belonged to different workers or that were logically independent, and *field reordering* to separate fields that were physically close but logically distinct. Although these are well-known techniques, padding required careful analysis, as it increases eviction misses significantly. The field reordering technique was not easily applied either, as the relationships between fields are quite complex.

Padding may lead to serious performance degradation for the dynamic data areas, such as `ChoiceP` and `BA`. This restricted our options for layout modification to just field reordering for these areas. The `BA` area was the target of one final data layout modification, since the analysis of cache behaviour surprised us with a high number of false sharing misses in this area for `chat80`. Further investigation showed that this was a memory allocation problem. The engines' top of stacks were being *shmalloc'ed* separately and appeared as part of the `BA` area in the analysis. This increased sharing misses in the area and was especially bad for the or-parallel applications, as different processor's top of stacks would end up in the same cache line. We addressed the problem by moving these pointers into the `Worker` area, where they logically belong. Our results show that the `bt-cluster` and `chat80` applications benefit the most from this optimisation; speedup improvements can be as significant as 60%. In contrast, the `pan2` application achieves improvements of less than 10% from this optimisation.

Privatisation of Shared Variables. This technique reduces the number of shared memory accesses by making local copies in each node of the machine. In the best case, the shared variables are read-only and hence local copies can actually be allocated in private memory. The high number of references to `Worker` suggested that privatisation could be applied there. In fact, Andorra-I did already use private copies of the variables in `Worker` and there was little room for improvement. The `Locks` and `Code` data areas are the major candidates to privatisation in Andorra-I. The `Locks` area only includes pointers to the actual locks, is thus read-only during execution, and can be easily privatised. Another area that is also read-only during parallel execution of our benchmarks is `Code`. Unfortunately, logic programs in general can change the database and, therefore, update `Code`, making privatisation complex. In the port of the related logic programming system Aurora to the Butterfly, an early scalable machine, it was necessary to privatise the `Code` area to obtain good performance [22].

Figure 11 shows the impact of privatisation on the speedup of Andorra-I. The figure shows that privatisation improves speedups by up to 10% at most and that the impact of this optimisation decreases as the number of processors increases.

Lock Distribution. This technique was considered to reduce contention on accesses to logical variables, and-scheduling, or-scheduling, and stack management. The original implementation used a single array of locks to implement these operations. In the worst case, several workers would contend for the same lock causing contention. To improve scalability, we implemented different lock data structures for different purposes. We expected best

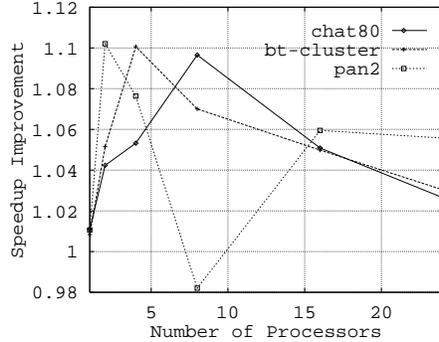


Figure 11: Impact of Privatisation Optimisations

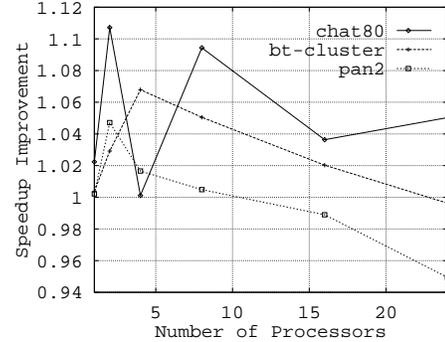


Figure 12: Impact of Lock Distribution Optimisations

results for or-parallel applications, as the optimisation prevents different teams from contending on accesses to logical variables. The cost of this optimisation is that, if the arrays of locks are shared, there will be more expensive remote cache misses. Our results show that the `bt-cluster` and `chat80` applications benefit somewhat from this optimisation, but that the `pan2` application already exhibited a significant number of misses in the `Locks` area and suffers a slowdown.

Figure 12 shows the impact of lock distribution in the absence of data privatisation. The `bt-cluster` and `chat80` applications benefit somewhat from this optimisation. On the other hand, the `pan2` application already had significant misses from the `Locks` area. Increasing the number of shared locks also increases the number of remote misses and leads to worse speedups on 16 and 24 processors. Note that the number of misses would not have increased, if this optimisation had been applied together with privatisation. Section 5 will present results for the combination of these two optimisations.

Elimination of Locking in Scheduling. This technique improves performance in benchmarks with significant and-parallelism by testing whether there is available work, before actually locking the work queue. This modification is equivalent to replacing a `test_and_set` lock with a `test_and_test_and_set` lock. This optimisation provides a small speedup improvement for `pan2`, as it avoids locking when there is no and-work.

For `bt-cluster` the technique does not improve speedups as this application exhibits enough and-work to keep processors busy. `chat80` is not affected by this optimisation as it has or-parallelism only.

Figure 13 shows the variations in speedup entailed by this optimisation. The figure shows that the locking elimination has no impact on the or-parallel benchmark, `chat80`, as this application never uses the and-scheduler. The optimisation does not improve the speedup of `bt-cluster` either, since this application has plenty of and-work to keep processors busy. Finally, locking elimination provides a small increase in speedup for `pan2`, as it avoids locking when there is no and-work.

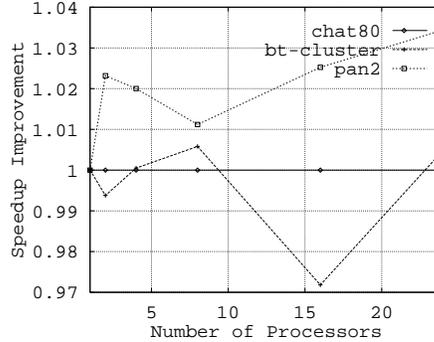


Figure 13: Impact of Scheduling Optimisations

5 Combined Performance of Optimisation Techniques

We next discuss the overall system performance with all optimisations combined. We compare speedups against the **idealised** and **original** results. The **idealised** speedups were recalculated for the new version of Andorra-I, but, as it is shown in the figures, the optimisations did not have any significant impact for the **idealised** machine. Figures 14 and 15 show the speedups for the two and-parallel applications running on top of the modified Andorra-I system. The maximum speedup for **bt-cluster** jumped from 12 to 20, whereas the maximum speedup for **tsp** jumped from 6.3 to 19. This indicates that the realistic machine is now able to exploit the available parallelism more fully. The explanation for the better speedups is a significant decrease in miss rates. For **bt-cluster**, the new version of Andorra-I exhibits a miss rate of only 0.6% for 16 processors, versus the 1.6% of the previous version. In the case of **tsp**, the optimisations decreased the miss rate from 3% to 1.2% again on 16 processors.

Figure 16 shows the number and source of misses for **bt-cluster** on 16 processors. Note that the figure keeps the same Y-axis as in Figure 3 to simplify comparisons against the cache behaviour of the original version of Andorra-I. The figure shows that the number of misses in the **Worker** area was reduced by a factor of 4, while the number of misses in the **Misc** area was reduced by an order of magnitude. The figure also shows that there is still significant true sharing in **Worker**, but false sharing is much less significant. The number of misses from **Misc** is now almost irrelevant.

The or-parallel benchmarks also show remarkable improvements due to the combination of the optimisation techniques we applied. Figure 18 shows the speedups for **chat80** and Figure 19 shows the speedups for **fp**. The maximum speedup for **chat80** almost doubles from one version of the system to the other. Note that speedups for the optimised system still flatten out on 16 processors, but at a much better efficiency. The other benchmark, **fp**, displays our most impressive result. The speedup for 24 processors jumps from 6.2 with the original Andorra-I system to 20 when all our optimisations are applied. This result

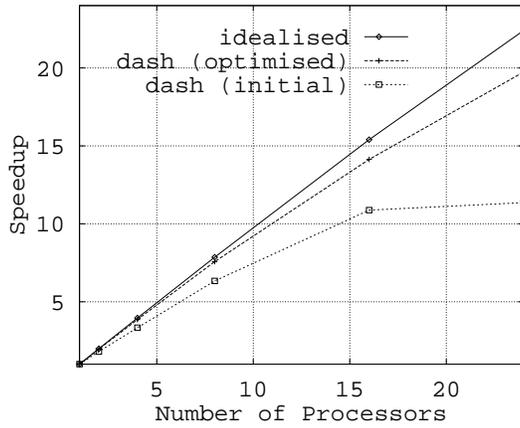


Figure 14: Speedups for `bt-cluster`

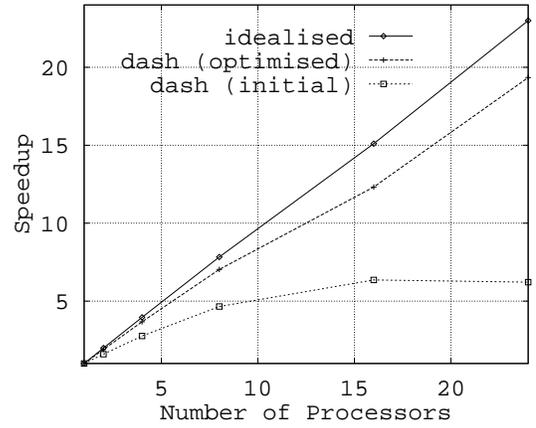


Figure 15: Speedups for `tsp`

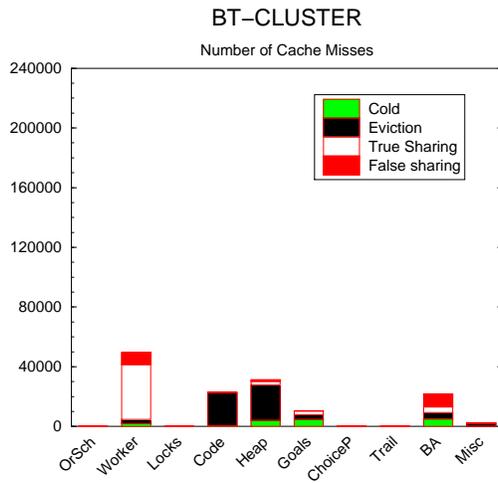


Figure 16: Misses by data area for `bt-cluster`

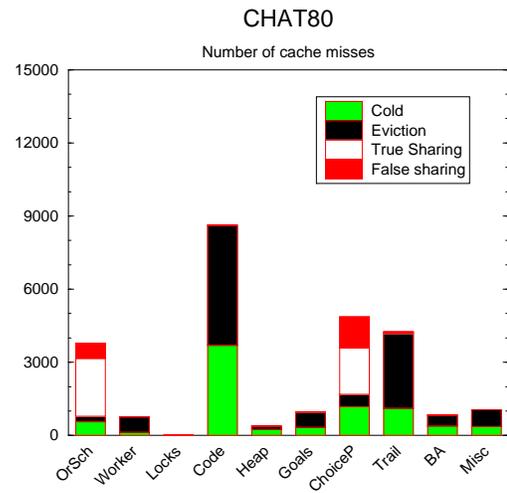


Figure 17: Misses by data area for `chat80`

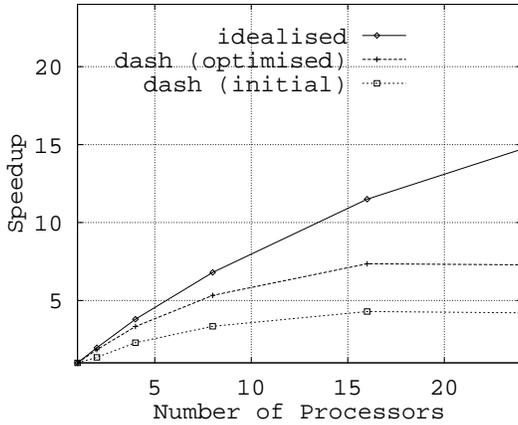


Figure 18: Speedups for `chat80`

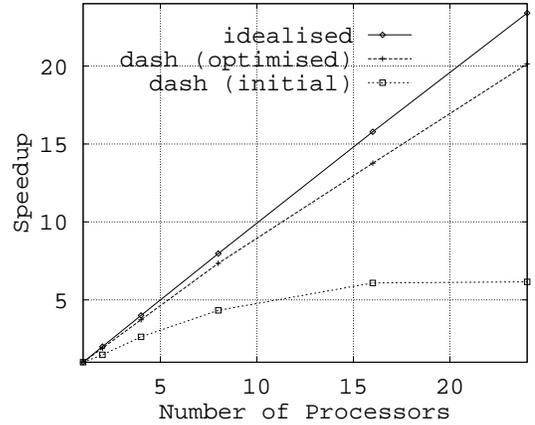


Figure 19: Speedups for `fp`

represents more than a three-fold improvement. Figure 17 shows the distribution of misses for `chat80` with 16 processors. The figure demonstrates that the number of misses in the **Worker** and **Misc** areas was reduced by an order of magnitude. The large number of eviction and cold start misses in the **Code** area remains however. Sharing misses are now concentrated in the **OrSch** and **ChoiceP** areas, as they should.

Figure 20 shows the speedups of the new version of Andorra-I for the `pan2` benchmark. In this case, the improvement resulting from our optimisations was quite small. Figure 21 shows the cache miss distribution for the optimised Andorra-I. The main source of misses was true sharing originating from the **Worker** region. A more detailed analysis proved that these misses originate from lack of work. Workers are searching each other's queues and generating misses. We are investigating more sophisticated scheduling strategies to address this problem.

6 Conclusions and Future Work

Andorra-I is an example of an and/or-parallel system originally designed for traditional bus-based shared-memory architectures. We have demonstrated that the system can also achieve good performance on scalable shared-memory systems. The key to these results was the extensive data available from detailed simulations of Andorra-I. The data we obtained from these simulations showed that significant performance improvements can be accrued without restructuring the system. Instead, performance can be dramatically improved by focusing on accesses to shared data. Our study concentrated on Andorra-I, but it can be directly extrapolated to several other PLP systems, since they share most of Andorra-I's data structures.

We believe there is potential for improving the performance of PLP systems even further.

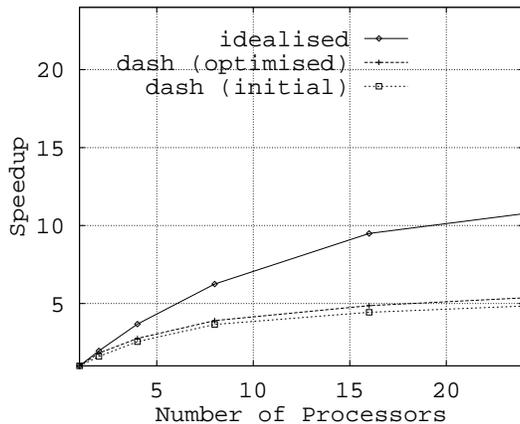


Figure 20: Speedups for pan2

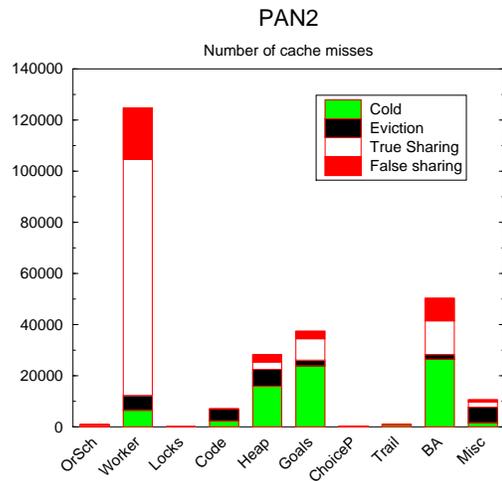


Figure 21: Misses by data area for pan2

To prove so will require more radical changes to data structures within Andorra-I itself, as the system was simply not designed for such large numbers of processors. In addition to making these changes, we are also interested in evaluating the performance of other parallel logic programming systems, such as the ones that exploit independent and-parallelism, on scalable architectures.

Acknowledgements

The authors would like to thank Leonidas Kontothanassis and Jack Veenstra for their help with the simulation infrastructure used in this paper. The authors would also like to thank Rong Yang, Tony Beaumont, D. H. D. Warren for their work in Andorra-I that made this work possible. Vitor Santos Costa would like to thank the University of Porto for granting his period of leave to perform this work, and also would like to thank support from the PROLOPPE and MELODIA projects. Inês Dutra would like to thank the APPELO project.

References

- [1] A. Agarwal, R. Bianchini, D. Chaiken, K. L. Johnson, D. Kranz, J. Kubiawicz, B. Lim, K. Mackenzie, and D. Yeung. The MIT Alewife Machine: Architecture and Performance. In *22nd Annual International Symposium on Computer Architecture (ISCA'95)*, June 1995.
- [2] Hassan Ait-Kaci. *Warren's Abstract Machine — A Tutorial Reconstruction*. MIT Press, 1991.
- [3] Khayri A. M. Ali and Roland Karlsson. The Muse Or-parallel Prolog Model and its Performance. In *Proceedings of the 1990 North American Conference on Logic Programming*, pages 757–776. MIT Press, October 1990.

- [4] Reem Bahgat. Solving Resource Allocation Problems in Pandora. Technical report, Imperial College, Department of Computing, 1990.
- [5] Anthony Beaumont, S. Muthu Raman, and Péter Szeredi. Flexible Scheduling of Or-Parallelism in Aurora: The Bristol Scheduler. In Aarts, E. H. L. and van Leeuwen, J. and Rem, M., editor, *PARLE91: Conference on Parallel Architectures and Languages Europe*, volume 2, pages 403–420. Springer Verlag, June 1991. Lecture Notes in Computer Science 506.
- [6] Johan Bevenmyr, Thomas Lindgren, and Hakan Millroth. Reform Prolog: The Language and its Implementation. In *Proceedings of the Tenth International Conference on Logic Programming*, pages 283–298. MIT Press, June 1993.
- [7] R. Bianchini and L. I. Kontothanassis. Algorithms for Categorizing Multiprocessor Communication Under Invalidate and Update-Based Coherence Protocols. In *Proceedings of the 28th Annual Simulation Symposium*, April 1995.
- [8] Convex Computer Corp. *Convex Exemplar Architecture*, November 1993.
- [9] J. A. Crammond. *Implementation of Committed Choice Logic Languages on Shared Memory Multiprocessors*. PhD thesis, Heriot-Watt University, Edinburgh, May 1988. Research Report PAR 88/4, Dept. of Computing, Imperial College, London.
- [10] J. A. Crammond. The Abstract Machine and Implementation of Parallel Parlog. Technical report, Dept. of Computing, Imperial College, London, June 1990.
- [11] I. C. Dutra. Strategies for Scheduling And- and Or-Work in Parallel Logic Programming Systems. In *Proceedings of the 1994 International Logic Programming Symposium*, pages 289–304. MIT Press, 1994. Also available as technical report CSTR-94-09, from the Department of Computer Science, University of Bristol, England.
- [12] James R. Goodman. Using Cache Memory to Reduce Processor-Memory Traffic. In *Proceedings of the 10th International Symposium on Computer Architecture*, pages 124–131, 1983.
- [13] Gopal Gupta, Enrico Pontelli, and Manuel Hermenegildo. &ACE: A High Performance Parallel Prolog System. In *Proceedings of the First International Symposium on Parallel Symbolic Computation, PASCOS'94*, 1994.
- [14] M. V. Hermenegildo and K. Greene. &-Prolog and its Performance: Exploiting Independent And-Parallelism. In *Proceedings of the Seventh International Conference on Logic Programming*, pages 253–268. MIT Press, June 1990.
- [15] Markus Hitz and Erich Kaltofen, editors. *Proceedings of the Second International Symposium on Parallel Symbolic Computation, PASCOS'97*, July 1997.
- [16] Gerry Kane and Joe Heinrich. *MIPS RISC Architecture*. Prentice Hall, 1992.
- [17] J. Laudon and D. Lenoski. The SGI Origin: A CC-NUMA Highly Scalable Server. In *Proceedings of the 24th Annual International Symposium on Computer Architecture*, pages 241–251, June 1997.
- [18] D. Lenoski, J. Laudon, K. Gharachorloo, A. Gupta, and J. Hennessy. The Directory-Based Cache Coherence Protocol for the DASH Multiprocessor. *Proceedings of the 17th ISCA*, pages 148–159, May 1990.
- [19] T. Lovett and S. Thakkar. The Symmetry Multiprocessor System. In *International Conference of Parallel Processing*, pages 303–310, 1988. University Park, Pennsylvania.
- [20] Ewing Lusk, David H. D. Warren, Seif Haridi, et al. The Aurora Or-parallel Prolog System. *New Generation Computing*, 7(2,3):243–271, 1990.
- [21] Johan Montelius. Penny, A Parallel Implementation of AKL. In *ILPS'94 Post-Conference Workshop in Design and Implementation of Parallel Logic Programming Systems, Ithaca, NY, USA*, November 1994.
- [22] Shyam Mudambi. Performances of aurora on NUMA machines. In *Proceedings of the Eighth International Conference on Logic Programming*, pages 793–806. MIT Press, June 1991.
- [23] S. Raina, D. H. D. Warren, and J. Cownie. Parallel Prolog on a Scalable Multiprocessor. In Peter Kacsuk and Michael J. Wise, editors, *Implementations of Distributed Prolog*, pages 27–44. Wiley, 1992.
- [24] V. Santos Costa, R. Bianchini, and I. C. Dutra. Evaluating the Impact of Coherence Protocols on Parallel Logic Programming Systems. In *Proceedings of the 5th EUROMICRO Workshop on Parallel and Distributed Processing*, pages 376–381, 1997. Also available as technical report ES-389/96, COPPE/Systems Engineering, May, 1996.

- [25] V. Santos Costa, R. Bianchini, and I. C. Dutra. Parallel Logic Programming Systems on Scalable Multiprocessors. In *Proceedings of the 2nd International Symposium on Parallel Symbolic Computation, PASC0'97 [15]*, pages 58–67, July 1997.
- [26] V. Santos Costa, D. H. D. Warren, and R. Yang. Andorra-I: A Parallel Prolog System that Transparently Exploits both And- and Or-Parallelism. In *Third ACM SIGPLAN Symposium on Principles & Practice of Parallel Programming*, pages 83–93. ACM press, April 1991. SIGPLAN Notices vol 26(7), July 1991.
- [27] Kish Shen. Initial Results from the Parallel Implementation of DASWAM. In Michael Maher, editor, *Proceedings of the 1996 Joint International Conference and Symposium on Logic Programming*. The MIT Press, 1996.
- [28] Kish Shen. Overview of DASWAM: Exploitation of Dependent And-parallelism. *J. of Logic Prog.*, 29(1–3), 1996.
- [29] Márcio G. Silva, I. C. Dutra, Ricardo Bianchini, and Vítor Santos Costa. The Influence of Computer Architectural Parameters on Parallel Logic Programming Systems. In *Workshop on Practical Aspects of Declarative Languages (PADL99)*, January 1999. Also available as Technical Report ES/477-98, COPPE Systems Engineering, Sep/98.
- [30] T. Chikayama, T. Fujise, and H. Yashiro. A Portable and Reasonably Efficient Implementation of KLI. In *Proceedings of the Eleventh International Conference on Logic Programming*, June 1993.
- [31] Evan Tick. *Memory Performance of Prolog Architectures*. Kluwer Academic Publishers, Norwell, MA 02061, 1987.
- [32] J. E. Veenstra and R. J. Fowler. MINT: A Front End for Efficient Simulation of Shared-Memory Multiprocessors. In *Proceedings of the 2nd International Workshop on Modeling, Analysis and Simulation of Computer and Telecommunication Systems (MASCOTS '94)*, 1994.
- [33] David H. D. Warren. An Abstract Prolog Instruction Set. Technical Note 309, SRI International, 1983.
- [34] David H. D. Warren. The SRI Model for Or-Parallel Execution of Prolog—Abstract Design and Implementation Issues. In *Proceedings of the 1987 International Logic Programming Symposium*, pages 92–102, 1987.
- [35] David H. D. Warren. The Andorra model. Presented at Gigalips Project workshop, University of Manchester, March 1988.
- [36] Rong Yang, Tony Beaumont, Inês Dutra, Vítor Santos Costa, and David H. D. Warren. Performance of the Compiler-Based Andorra-I System. In *Proceedings of the Tenth International Conference on Logic Programming*, pages 150–166. MIT Press, June 1993.
- [37] Rong Yang, Vítor Santos Costa, and David H. D. Warren. The Andorra-I Engine: A parallel implementation of the Basic Andorra model. In *Proceedings of the Eighth International Conference on Logic Programming*, pages 825–839. MIT Press, 1991.