

New OPENMP Directives for Irregular Data Access Loops

J. Labarta, E. Ayguadé, José Oliver (CEPBA, Barcelona)
and D. Henty (EPCC, Edinburgh)

1 Abstract

Many scientific applications involve array operations that are sparse in nature, ie array elements depend on the values of relatively few elements of the same or another array. When parallelised in the shared-memory model, there are often inter-thread dependencies which require that the individual array updates are protected in some way. Possible strategies include protecting all the updates, or having each thread compute local temporary results which are then combined globally across threads. However, for the extremely common situation of sparse array access, neither of these approaches is particularly efficient. The key point is that data access patterns usually remain constant for a long time, so it is possible to use an inspector/executor approach. When the sparse operation is first encountered, the access pattern is inspected to identify those updates which have potential inter-thread dependencies. Whenever the code is actually executed, only these selected updates are protected. We propose a new OPENMP clause, `indirect`, for parallel loops that have irregular data access patterns. This is trivial to implement in a conforming way by protecting every array update, but also allows for an inspector/executor compiler implementation which will be more efficient in sparse cases. We describe efficient compiler implementation strategies for the new directive. We also present timings from a Discrete Element Modelling application code where the inspector/executor approach has been implemented by hand using standard OPENMP directives. The results demonstrate that the method can be extremely efficient in practice.

2 Introduction

Many codes in science and engineering involve array operations that are sparse in nature, i.e. array elements are updated based on the values of relatively few elements of the same or another array. Examples include Molecular-Dynamics (MD) simulations with short-ranged interactions (the force depends on the positions of a few nearby particles) and Finite-Element (FE) calculations (eg the values of a node depending on those nodes directly connected to it). As these opera-

tions often comprise the most computationally intensive part of the code it is essential to parallelise them efficiently. A typical example is computing the total force `force(i)` on each particle `i` in an MD code where the neighbouring particle pairs are stored in a list:

```
do ipair = 1, npair
  i = pairlist(1, ipair)
  j = pairlist(2, ipair)
  fij = pairforce(x(i), x(j))
  force(i) = force(i) + fij
  force(j) = force(j) - fij
end do
```

Similar access patterns occur in FE codes where the main loop is over the edges, faces or elements of a mesh, but the actual computation involves updating the nodes attached to them. The unstructured nature of the mesh again means that there are potential inter-thread dependencies, and that it is difficult for the programmer to predict in advance where they will occur for an arbitrary number of threads.

The inter-thread dependencies require that the individual array updates are protected in some way to ensure program correctness. If the above loop is split across threads then, although threads will always have distinct values of `ipair`, the values of `i` and `j` may simultaneously have the same values on different threads. As a result, there is a potential problem with updating the `force` array. Simple solutions to this problem include making all the updates atomic, or having each thread compute temporary results which are then combined across threads (in cases like the above this amounts to an array reduction which, although not currently supported, is proposed for inclusion in OPENMP 2.0). However, for the extremely common situation of sparse array access neither of these approaches is particularly efficient.

Due to the sparse access pattern, the majority of updates can actually take place with no protection; making every update atomic therefore incurs an unnecessary overhead. If temporary arrays are created for each thread, there is a bottleneck when they are combined across threads; even if this is implemented in parallel (i.e. not using the naive approach of a critical section) the heavy load on the memory system may produce poor scaling. When using many threads there will

also be a very high memory requirement for the temporary arrays, most of which will in practice be filled with zeroes.

The key point is that the irregular access patterns usually remain constant for a long time (sometimes the whole simulation) so it is possible to use a two-pass or inspector/executor approach. In the first pass, the data access pattern is inspected to identify those updates which have potential inter-thread dependencies and their locations are stored in a lookup table. When the calculation is actually executed, only these selected updates are protected (e.g. by an atomic update).

Unfortunately, the only way to do this in OPENMP is to write the code by hand. Even if the array reduction features of OPENMP 2.0 could be used, a compiler would be forced to use an inefficient method as there is no way of knowing at compile time that the data access patterns are constant (there is little benefit in using the inspector/executor approach if the lookup tables cannot be reused). We therefore propose a new OPENMP directive, `indirect`, for loops that have irregular data access patterns. This is trivial to implement in a conforming way (e.g. by making every update in the MD loop atomic) but also allows for a two-pass compiler implementation which will be more efficient in sparse cases. The directive will provide an efficient alternative to array reduction (as in the above MD example), as well as being of use in other situations such as the use of ordered sections.

3 Related Work

Similar problems have previously been addressed in the context of High Performance Fortran. The simplest constructs, eg specific cases of irregular array reductions, are dealt with by special routines in the HPF library such as `SUM_SCATTER` and `MAXVAL_SCATTER`. Directives for the specification of more general irregular data access loops have been used previously in some versions of HPF [8]. These directives have influenced the one proposed in this paper. The main difference is, however, that in the context of HPF this directive is used to compute long-lived communication patterns, while our proposal focuses on iteration scheduling in the context of shared-memory multiprocessors.

4 OPENMP extensions

This section presents the two proposed extensions to OPENMP: the `indirect` clause and the `schedule` directive.

4.1 The indirect clause

The problem of irregular data access patterns situations in loops can be described in simple words: they are access patterns in which, given a loop, two different iterations of the loop modify the same data. This situation prevents the loop from being parallel, and thus serializes its execution. OPENMP provides two mechanisms that help in the parallelization of such loops: `atomic` and `critical` directives. These mechanisms, however, are excessively expensive in codes where not all the accesses need to be done in mutual exclusion. This situation arises also in two other situations, which are irregular reductions with low level of shared updates, and loops in which only some iterations need to be executed in sequential ordered (and thus, the `ordered` clause is too restrictive). Such codes can benefit from the use of the new proposed clause (`indirect`). The `indirect` clause can be applied to the `parallel do` directive in one of the three following situations:

1. In the presence of the `reduction` clause in the directive: this tells the compiler that the reduction that is being computed in the body of the loop has an irregular data access pattern, and that some different iterations can modify the same data.
2. In the presence of the `ordered` clause in the directive: this tells the compiler that the loop is partially ordered; that is: some iterations access the same data, and need to be executed sequentially in order to guarantee sequential consistency, but some other iterations access completely different data, and thus can be executed in parallel.
3. In the presence of some `critical` or `atomic` directive in the body of the loop: this is quite similar to the case in the `ordered` clause: not all the elements need to be computed in mutual exclusion; only those being accessed by different iterations (in this case, order is not important).

The `indirect` clause accepts a list of expressions as parameters. These expressions are those causing the irregular data access pattern:

```
$!omp parallel do [reduction|ordered]  
$!omp& indirect([expr1,...exprN])
```

4.2 Irregular Reductions

Reduction operations are frequently found in the core of scientific applications. Simplest reductions are those in which the final result can be computed as a combination of partial results (ie an associative/commutative operation). As result, these computations can be computed in parallel, since each thread can compute its own partial result which will be combined later (this can be

also done in parallel) with the partial results from the other threads.

Some scientific applications, however, need to perform reduction operations which are not directly parallelizable. These kind of reductions are what we call irregular reductions. The point that makes these reductions non-parallelizable is that the update index for the element is not the induction variable of the loop, but a function (f) of it, having the property that for two given values of the induction variable ($i, j, i \neq j$), it can happen that $f(i) = f(j)$. The following code shows an example of such a case:

```
$!omp parallel do reduction (+:force)
!$omp&          indirect(i,j)
do ipair = 1, npair
  i = pairlist(1,ipair)
  j = pairlist(2,ipair)
  fij = pairforce(x(i),x(j))
  force(i) = force(i) + fij
  force(j) = force(j) - fij
end do
$!omp end parallel do
```

As can be observed, the updated elements of `force` depend on the lookup table `pairlist`, which can give the same values for `i` or `j` for different values of `ipair`. The only way to parallelize the previous code using OPENMP is to protect the update of the `force` vector either with `atomic` or with `critical` directives. These solutions, however, are excessively expensive in codes where not all the accesses need to be done in mutual exclusion. The inclusion of the `indirect` clause in the `parallel do` directive in the presence of a `reduction` clause tells the compiler that the reduction being performed in the parallel loop has an irregular data access pattern, but some parts of it can be executed in parallel, thus enabling the compiler to generate code to deal with this situation. The naive implementation of this `indirect` clause is to surround the updates of `force` either with `critical` or `atomic` directives, but more complex implementations can generate code for a two-pass implementation of the reduction.

4.3 Non-complete Mutual Exclusion

As said in the previous subsection, a naive implementation of the `indirect` clause for irregular reductions can be done with the use of the critical sections (achieved using the `critical` directive). This also implies that the use of the `indirect` clause can also help in the code generation for cases in which the `critical` directive has been used to parallelize a loop without a reduction, but with some shared updates between iterations executed by different processors. An example of this case follows:

```
$!omp parallel do indirect(j,k)
do i = 1,n
  j = jindex(i)
  k = kindex(i)
$!omp critical
  a(j) = a(k) op expression_1
$!omp end critical
end do
$!omp end parallel do
```

4.4 Partially Ordered Loops

A special case of the previous example occurs when the shared updates need not only to be performed in mutual exclusion, but also in an ordered way. The use of the `indirect` clause in this case tells the compiler that the only iterations that actually need to be executed in an ordered way are those which are updating the same data. An example of code using the `indirect` clause in this manner is the following:

```
$!omp parallel do ordered indirect(j,k)
do i = 1,n
  j = jindex(i)
  k = kindex(i)
$!omp ordered
  a(j) = a(k) op expression_1
$!omp end ordered
end do
$!omp end parallel do
```

4.5 Scheduling re-use: the `schedule directive`

The typical structure of most scientific applications follows an iterative outer loop which implies a number of parallel computations, most of them over the same data. In that scenario, it could be useful to communicate scheduling information from one parallel region to others. This idea has been applied previously to data communication in HPF [8], with the use of the `SCHEDULE` clause, and can be also applied to OPENMP.

Two new directives can be introduced, in order to define/undefine scheduling patterns:

- Schedule definition:

```
$!omp schedule schedule_name
```

This directive tells the compiler that the next parallel loop (dynamically executed) will define (if not yet already done) a scheduling pattern that will be associated to the symbolic name `schedule_name`

- Schedule undefinition:

```
$!omp reset schedule_name
```

This directive tells the compiler that the schedule associated to the symbolic name `schedule_name` is not valid any more.

Finally, the indirect clause syntaxes can be expanded to accept a symbolic name for an schedule:

```

$!omp parallel do [reduction|ordered]
$!omp& indirect([expr1,...exprN]
!$omp&          [:schedule_name])

```

This tells the compiler that the given parallel loop follows the same scheduling policy as the one that defined the schedule `schedule_name`, and thus this scheduling information can be used, avoiding the overhead of computing the schedule for this loop. An example of the use of the `schedule` directive combined with the `indirect` could be something like the following code:

```

C outer sequential loop
do step = 1,nsteps
.
.
$!omp schedule S
$!omp parallel do reduction (+:force)
!$omp          indirect(i,j:S)
  do ipair = 1, npair
    i = pairlist(1,ipair)
    j = pairlist(2,ipair)
    fij = pairforce(x(i),x(j))
    force(i) = force(i) + fij
    force(j) = force(j) - fij
  end do
$!omp end parallel do
.
.
$!omp parallel do reduction (+:force)
!$omp          indirect(i,j:S)
  do ipair = 1, npair
    i = pairlist(1,ipair)
    j = pairlist(2,ipair)
    fij = new_force(x(i),x(j))
    force(i) = force(i) + fij
    force(j) = force(j) - fij
  end do
$!omp end parallel do
.
.
  if (some_condition(x)) then
    recalculate(pairlist)
  end if
$!omp reset schedule S
end if
.
.
end do

```

5 Compiler Implementations

This section outlines some implementations for the `indirect` clause in the context of an irregular reduction or non-complete mutual exclusion (using intervals) and reviews some proposed implementations that can be applied in the context of a partially ordered loop.

5.1 Intervals

This implementation is based on the inspector/executor model. The first time that a parallel loop marked with an `indirect` clause is reached, the program starts the execution of the inspector code, which has been generated by the compiler. This code takes note of the behaviour of each processor, and builds some data structure that keeps information of it. Following executions of the parallel loop make use of that data structure in order to execute in parallel as many of the iterations as possible.

The following code presents a simple example of parallel loop with non-complete mutual exclusion that has been parallelized using the `indirect` clause.

```

!
! v(1:20)=(1,2,3,10,20,4,2,5,6,8,9,
!          3,10,3,11,12,13,14,15,16)
!$omp parallel do schedule (static)
!$omp&          indirect (v(i),i)
do i = 1, 20
  a(v(i)) = a(i) * 2
end do

```

As can be seen, there are some iterations of the loop that modify the same data (positions 2,3 and 10 of the vector `v`). But most of the iterations do not modify the same data, and thus they can be executed in parallel. Figure 1 shows an example of a data structure generated by the compiler during the inspector phase. For each processor, there is a bitmap that shows which elements it is accessing in the `v` vector. After that, the inspector computes the shared elements vector, which is a bitmap that tells, for each element in the `v` vector whether it is shared or not (a position of `v` is shared if the bitmap count for the column is more than one; that is: it is being modified by more than one processor). Once this bitmap is computed, the inspector can build the final data structure that will be used by the executor phase: the intervals for each processor. The iteration space of each processor is split up into intervals. An interval is a range of consecutive iterations that are all either completely shared or not shared. The data structures for the previous example result in the following:

```

nintervals(1:4) = (3,3,3,1)
shared(1) = (false,true,false)
shared(2) = (false,true,false)
shared(3) = (false,true,false)

```

```

shared(4) = (false)
lower(1) = (1,2,4)
lower(2) = (6,7,8)
lower(3) = (11,12,15)
lower(4) = (16)
upper(1) = (1,3,5)
upper(2) = (6,7,10)
upper(3) = (11,14,15)
upper(4) = (20)

```

Subsequent executions of the parallel loop will use the compiler’s generated executor code, which makes use of the previous data structures to execute the parallel loop. This code executes in mutual exclusion only those iterations (actually intervals) that conflict with some other processor. Those intervals that are accessed only by one processor are executed in parallel. The executor’s code for the previous example could look like this:

```

!$omp parallel do schedule (static)
do p = 1, numprocs
  my_proc = omp_get_threadnum()+1
  do interval = 1, nintervals(my_proc)
    if (shared(my_proc, interval)) then
      do i = lower(my_proc, interval), &
        upper(my_proc, interval)
!$omp critical
      a(v(i)) = a(i) * 2
!$omp end critical
    end do
    else
      do i = lower(my_proc, interval), &
        upper(my_proc, interval)
      a(v(i)) = a(i) * 2
    end do
  end if
end do
end do

```

It is important to point out that this method can be further optimised by improving the executor phase in order to classify the shared segments depending on the processors that are accessing them (shared sets), and utilizing different locks to access each one of the sets, in order to reduce contention when *hot spots* are detected.

5.2 Partially ordered loops

Partially ordered loops have been a case of study in many previous works. One possible approach to the parallelization of non-completely parallel loops is to detect data dependencies at runtime, again using an inspector/executor model. The basic idea is to move data dependence detection from compile-time (where data available for analysis is restricted) to runtime. In most of the reviewed implementations, at runtime, the inspector builds data dependence graphs based on the addresses accessed by the loop and, based on those graphs,

schedules iterations in *wavefronts* (sets of iterations which are dependence-free among them) [1, 2, 3, 4, 5, 6, 7]. Many of these implementations can be applied by the compiler to implement the semantics of the indirect clause applied to an ordered directive.

6 Experiments

The results presented here come from a code whose major computational loop is of the same form as that discussed in Section 2. The DEMONS application — Discrete Element Modelling on SMP clusters — is a Fortran 90 test code used to evaluate the efficiency of hybrid message-passing and shared-memory parallelism on clusters of shared-memory machines [9]. Here we do not utilise the MPI capabilities of the code, and run on a single process to investigate the performance of the pure OpenMP implementation with varying numbers of threads T .

Discrete Element Models (DEMs) simulate the behaviour of systems of particles which interact via pairwise forces. The particle positions are evolved in time using an iterative loop comprising many small, discrete time-steps. A typical DEM might study the formation of “sand piles” as grains of sand are dropped onto a solid surface. The general properties of a DEM are: there are many particles; particles remain relatively static; the inter-particle force is very short-ranged; the time taken to compute the forces dominates the simulation. It is therefore very important to parallelise the force calculation efficiently.

To avoid having to loop through all pairs of particles, the DEMONS code maintains a list of pairs of all those particles which are close together and therefore likely to interact. As the force is short-ranged, this simple and commonly used technique reduces the complexity of the calculation from $O(N^2)$ to $O(N)$. As particles remain relatively stationary, this list only requires to be recalculated fairly infrequently; in a real simulation, the same list can typically be used for many hundreds of iterations. Calculating the forces involves looping over this list of pairs, computing the inter-particle force for each pair, and updating the force on the two particles accordingly. The loop is parallelised using a simple PARALLEL DO directive. Since the loop is over particle pairs (not the particles themselves), and each particle belongs to many pairs, care must be taken when updating the forces on each particle due to potential inter-thread dependencies.

6.1 Implementation of force calculation

Three basic approaches were investigated for resolving the inter-thread dependencies in the force calculation

- making every force update atomic

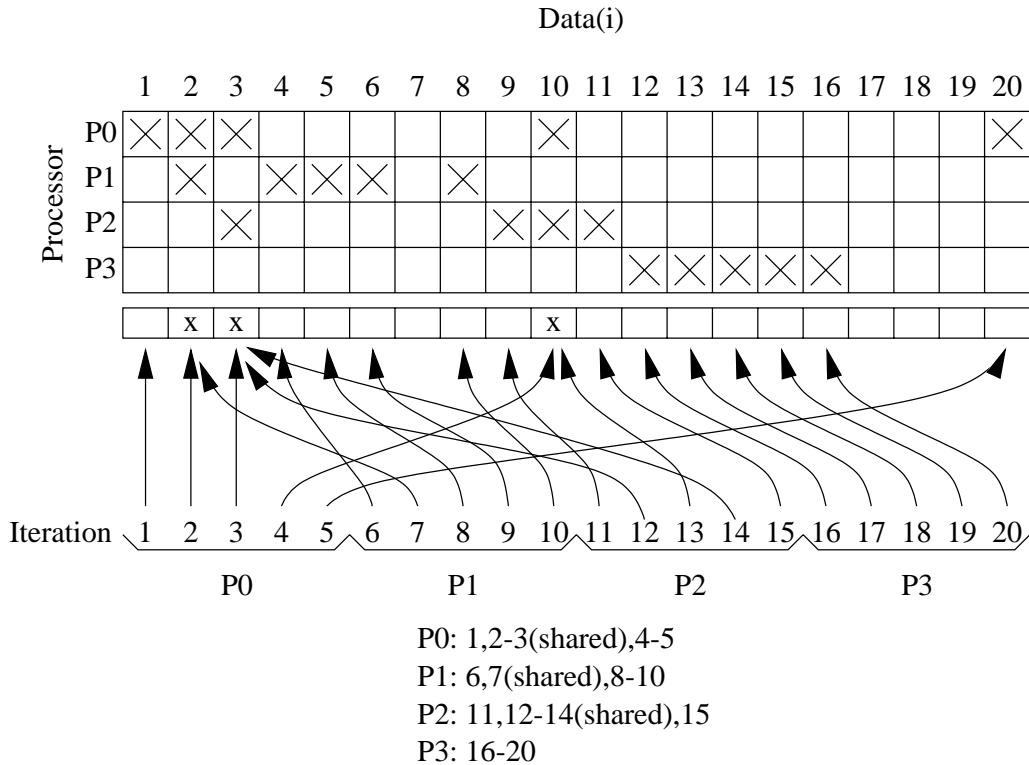


Figure 1: Example of intervals construction using a bitmap.

- using array reduction
- using a two-pass inspector / executor approach

Since array reduction is not part of the OPENMP 2.0 standard we implemented it by hand in several ways

- first accumulating into private temporary arrays and performing the final global array sum across threads in a critical region
- accumulating into private arrays and performing the global sum in T sections, striped across threads and separated by barriers, so that there are no inter-thread dependencies
- accumulating into different rows of a shared matrix of dimension $T \times N$, then performing the global sum in parallel over the transpose column direction

The use of critical sections gave very poor performance; the other two methods gave very similar timings. Here we quote results from the third “transpose” approach. We implement the inspector / executor approach in a very straightforward manner. Whenever a new list of pairs is created, we identify particles receiving force updates from more than one thread. When calculating the force, only these updates are protected by an ATOMIC directive.

6.2 Timings

The code for computing the force is almost identical to that in Section 2 except that the calculation is done in three dimensions so there are three position coordinates x rather than just one. For the tests we use a million particles, with the force law requiring one floating-point inverse and one square root per call. All calculations were performed in double precision.

The benchmark platforms were a Sun HPC 3500 (eight \times 400 MHz UltraSPARC-II CPUs) at EPCC and a Compaq ES40 (four \times 500 MHz Alpha EV6 CPUs). For compilation on the Sun we used the Kuck and Associates (KAI) Guide system version 3.7, which uses source-to-source translation with calls to a parallel runtime library; on the Compaq, OpenMP is part of the standard f90 compiler. Sun’s own native OpenMP compiler, part of the recently-released Workshop Forte 6.0 environment, was not available in time for this paper.

We plot the parallel efficiencies of the three implementations of the force calculation on the Sun and Compaq in Figures 2 and 3. The timings for the serial code, $t(T = 1)$, were 4.16 and 1.89 seconds respectively.

Using the KAI compiler on the Sun, the atomic updates are done using software locks which are rather expensive. Contention for these locks makes using atomic updates for all force calculations completely infeasible. However, using the two-pass approach is more effective than using array reduction despite the high cost associated with atomic locks.

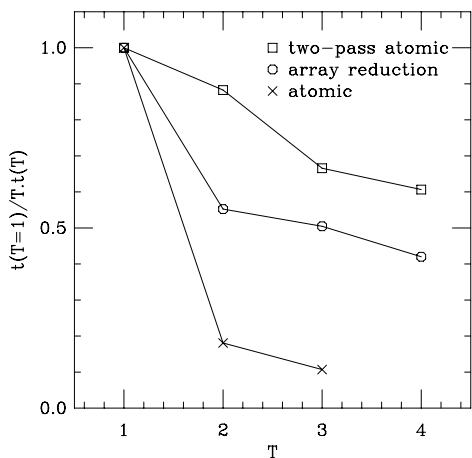


Figure 2: Parallel efficiency against threads T on Sun

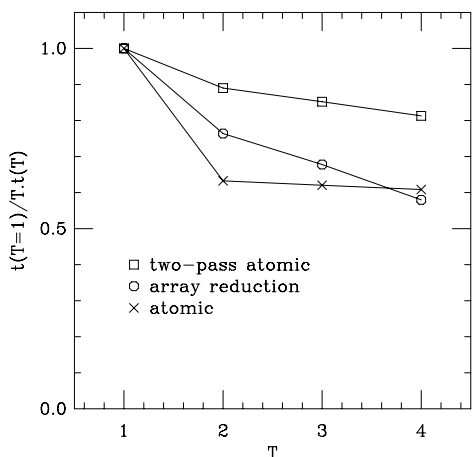


Figure 3: Parallel efficiency against T on Compaq

On the Compaq, atomic updates are much more efficient but they still incur a significant overhead. For $T < 4$ the cost is such that array reduction is more efficient, although the atomic approach is slightly better for $T = 4$. However, using the two-pass approach is always the most efficient method.

7 Conclusions

The results shown in Section 6 demonstrate that the two-pass technique for irregular reductions outperforms the alternative methods in a real code. For the DEMONS code, the calculation of the force is relatively expensive; it was possible to use a very naive approach to implement the two-pass method where the decision as to whether or not to use `atomic` was taken for each individual update. In other situations the amount of work per update may be much smaller, for example a single addition, and reducing the overhead of the implementation will be important. Here, the technique of dividing the loop into intervals, as proposed in Section 5, will be extremely beneficial to performance. We con-

clude that our proposed new OPENMP directives would enable straightforward and efficient automatic parallelisation of a wide range of scientific applications.

8 Acknowledgements

This research was partially supported by the Improving the Human Potential Programme, Access to Research Infrastructures, under contract HPRI-1999-CT-00071 “Access to CESSA and CEPBA Large Scale Facilities” established between The European Community and CESSA-CEPBA. We acknowledge the use of the PPARC-funded Compaq MHD Cluster in St. Andrews.

References

- [1] L. Rauchwerger, N. M. Amato and D. A. Padua, “Run-time methods for parallelizing partially parallel loops”, In proceedings of the 1995 International Conference on Supercomputing, July 3–7, 1995 Barcelona (Spain)
- [2] C. Zhu and P. C. Yew, “A scheme to enforce data dependence on large multiprocessor systems”, IEEE Transactions on Software Engineering, 13(6):726–739, 1987.
- [3] S. Midkiff and D. Padua, “IEEE Transactions on Computers”, C-36(12):1485–1495, 1987.
- [4] D. K. Chen, P. C. Yew and J. Torrellas, “An efficient algorithm for the run-time parallelization of doacross loops”, In proceedings of Supercomputing ’94, Nov. 1994.
- [5] V. Krothapalli and P. Sadayappan, “An approach to synchronization of parallel computing”, In proceedings of the 1988 International Conference on on Supercomputing, pp. 573–581, June 1988.
- [6] J. Wu, J. Saltz, S. Hiranandani and H. Berryman, “Runtime compilation methods for multicomputers”, In Dr. H.D. Schewtman, editor, Proceedings of the 1991 International Conference on Parallel Processing, p. 26–30. CRC Press, Inc., 1991, Vol. II—Software
- [7] C. Polychronopoulos, “Compiler optimizations for enhancing parallelism and their impact on architecture design”, IEEE Transactions on Computers, C-37(8):991–1004, Aug. 1998.
- [8] S. Benkner, K. Sanjari, V. Sipkova and B. Velkov, “Parallelizing irregular applications with the Vienna HPF+ Compiler VFC”
- [9] D.S. Henty, “Performance of Hybrid Message-Passing and Shared-Memory Parallelism for Discrete Element Modeling”, In proceedings of Supercomputing 2000, Nov. 2000 (to appear).