

# SUPPLE: an Efficient Run–Time Support for Non–Uniform Parallel Loops

Salvatore Orlando<sup>a</sup>, Raffaele Perego<sup>b</sup>

<sup>a</sup> *Dip. di Matematica Applicata ed Informatica, Università Ca' Foscari di Venezia, via Torino 155,  
30173 Venezia Mestre, Italy*

<sup>b</sup> *CNUCE - C.N.R., via S. Maria 36, 56126 Pisa, Italy*

---

## Abstract

This paper presents SUPPLE (SUPort for Parallel Loop Execution), an innovative run–time support for the execution of parallel loops with regular stencil data references and non–uniform iteration costs. SUPPLE relies upon a static block data distribution to exploit locality, and combines static and dynamic policies for scheduling non–uniform iterations. It adopts, as far as possible, a static scheduling policy derived from the owner computes rule, and moves data and iterations among processors only if a load imbalance actually occurs. SUPPLE always tries to overlap communications with useful computations by reordering loop iterations and prefetching remote ones in the case of workload imbalance. The SUPPLE approach has been validated by many experimental results obtained by running a multi-dimensional flame simulation kernel on a 64–node Cray T3D. We have fed the benchmark code with several synthetic input data sets built on the basis of a load imbalance model. We have compared our results with those obtained with a CRAFT Fortran implementation of the benchmark.

*Keywords:* Data parallelism, parallel loop scheduling, load balancing, run–time supports, compiler optimizations.

---

## 1 Introduction

Data parallelism is the most common form of parallelism exploited to speed-up scientific applications. In the last few years, high level data parallel languages such as High Performance Fortran (HPF) [4], Vienna Fortran [25] and Fortran D [6] have received considerable interest because they facilitate the expression of data parallel computations by means of a simple programming abstraction. In particular HPF, whose definition is the fruit of several research proposals, has been recommended as a standard by a large Forum of universities and industries. Using HPF, programmers only have to provide a few directives to specify processor and data layouts, and the compiler translates the source program into an SPMD code with explicit interprocessor communications and synchronizations for the target multicomputer. Data parallelism can be expressed above all

by using collective operations on arrays [22], e.g. collective Fortran 90 operators, or by means of parallel loop constructs, i.e. loops in which iterations are declared as independent and can be executed in parallel. In this paper we are interested in run-time supports for parallel loops with regular stencil data references. In particular, we consider *non-uniform parallel loops*, i.e. parallel loops in which the execution time of each iteration varies considerable and cannot be predicted statically.

The typical HPF run-time support for parallel loops exploits a static data layout of arrays onto the network of processing nodes, and a static scheduling of iterations which depends on the specific data layout. Arrays are distributed according to the directives supplied by programmers, and computations, i.e. the various loop iterations, are assigned to processors following a given rule which depends on the data layout (e.g. the *owner computes rule*). A BLOCK distribution is usually adopted to exploit data locality: computations mapped on a given processor by the compiler mainly use the data block allocated to the corresponding local memory. Conversely, CYCLIC distribution is usually adopted when load balancing issues are more important than locality exploitation. A combination of both the distributions, where smaller array blocks are scattered on processing nodes, can be adopted to find a tradeoff between locality exploitation and load balancing. It is worth noting, however, that the choice of the best distribution is still up to programmers, and that the right choice depends on the features of the particular application. The general problem of finding an optimal data layout is in fact NP-complete [9].

The adoption of a static policy to map data and computations reduces run-time overheads, because all mapping and scheduling decisions are taken at compile-time. While it produces very efficient implementations for regular concurrent problems, the code produced for *irregular problems* (i.e. problems where some features cannot be predicted until run-time) may be characterized by poor performance. Many researches have been conducted in the field of run-time supports and compilation methods to efficiently implement irregular concurrent problems, and these researches are at the basis of the new proposal of the HPF Forum for HPF2 [5]. The techniques proposed are mainly based on run-time codes which collect information during the first phases of computation, and then use this information to optimize the execution of the following phases of the same computation. An example of these techniques is the one adopted by the CHAOS support [20] to implement non-uniform parallel loops. The idea behind this feature of the CHAOS library is the run-time *redistribution* of arrays, and the consequent re-mapping of iterations. Redistribution is carried out synchronously between the execution of subsequent executions of parallel loops, and is decided on the basis of information (mainly, timing information) collected at run-time during a previous execution of the loop. If the original load distribution is not uniform, the data layout is thus modified in order to balance the processor loads. In [18] a dialect of HPF has been extended with new language constructs which are interfaced with the CHAOS library to support irregular computations.

In this paper we address non-uniform parallel loop implementations, introducing a truly innovative support called SUPPLE (SUPort for Parallel Loop Execution). SUPPLE is not a general support on which we can compile every HPF parallel loop, but only non-uniform (as well as uniform) parallel loops with regular stencil data references. Since stencil references are regular and known at compile time, optimizations such as *message vectorization*, *coalescing* and *aggregation*, as well as *iteration reordering* can be carried out to reduce overheads and hide communication latencies [7].

```

C Loop on the time steps
DO k= 1,K
C   Convection phase over a structured mesh
C   Nearest neighbor communications - Constant computation time
FORALL (i= 2:N1-1, j= 2:N2-1)
  A(i,j) = A(i,j) + F(B(i,j), B(i-1,j), B(i,j-1), B(i+1,j), B(i,j+1), C(i,j))
END FORALL
B = A
C   Chemical reaction phase - No communications - High load imbalance
C   Workload distribution evolves during simulation
FORALL (i= 1:N1, j= 1:N2) C(i,j) = Reaction(A(i,j))
END DO

```

Fig. 1. HPF-like code of a simplified two-dimensional flame modeling code.

SUPPLE is able to initiate loop computation using a statically chosen BLOCK distribution, thus starting iteration scheduling according to the *owner computes rule*. During the execution, however, if a load imbalance actually occurs, SUPPLE adopts a dynamic scheduling strategy that migrates iterations and data from overloaded to underloaded processors in order to increase processor utilization and improve performances. All decisions about migration are mainly local, in order not to introduce global communication/synchronizations. SUPPLE overlaps most overheads deriving from dynamic scheduling, i.e. messages to monitor loads and move data, with useful computations. The main features that allow us to overlap the otherwise large overheads are (1) *prefetching* of remote loop iterations on underloaded processors, and (2) a complete asynchronous support that does not introduce barriers between consecutive iterations of loops. Below we describe our support in detail, and we show some results obtained by running the kernel of a flame simulation code on a Cray T3D. In this application, the characteristics of the input data set strongly influence the time needed to compute loop iterations and thus may cause workload imbalances. For this reason, we ran our tests on synthetic data sets, built on the basis of a simple load imbalance model. To compare SUPPLE with an HPF-style static approach, we used the Cray CRAFT-Fortran, whose language directives and compilers were partially adapted from Fortran-D [6] and Vienna Fortran [25]. Even if SUPPLE is presented here as a run-time support to implement non-uniform parallel loops, it can be profitably used [15] for a less narrow class of applications, i.e. to compile “uniform” parallel loops on heterogeneous and/or time-shared systems, where we can have high load imbalance even if the specific application is regular.

The paper is organized as follows. Section 2 describes the flame simulation code we used as a benchmark to validate the SUPPLE approach. Section 3 presents the main features of SUPPLE and gives some details on its implementation. Assessments of the experiment results are reported and discussed in depth in Section 4, which also describes the load imbalance model used to build the synthetic data sets. Finally, Section 5 provides a summary of related work, and Section 6 draws the conclusions.

## 2 The benchmark

The benchmark we used to validate the SUPPLE approach is directly derived from a class of actual scientific codes used to carry out detailed flame simulations [21]. Flame simulations have attracted

growing interest in the last few years: highly detailed simulations can in fact give scientists important information on fuel flammability limits and on the interactions between chemical and physical processes. One example of code that performs a time-dependent multi-dimensional simulation of hydrocarbon flames is described in detail in [18], and is included in the HPF-2 draft [5] as one of its motivating applications. The code is split into two distinct computational phases executed at each time step. The first phase computes fluid *convection* over a structured mesh. The computation at each point has similar costs, and involves the point itself and the nearest neighboring elements of the simulation mesh. Stencil communications are thus required at each time step to carry out this phase on a distributed memory machine. The second phase simulates chemical *reaction* and subsequent energy releases represented by ordinary differential equations. The solution at each grid point only depends on the value of the point itself, but its computational costs may vary significantly since the chemical combustion proceeds at different rates across the simulation space. The reaction phase globally requires more than half of the total execution time, and most of this time is spent on a small fraction of the mesh points. Furthermore, the workload distribution evolves as the simulation progresses. The HPF-like code of a simplified two-dimensional flame modeling code is reported in Fig. 1. The interesting characteristic of flame simulation codes is that, in order to efficiently exploit a highly parallel machine, two conflicting requirements have to be dealt with: the need to adopt a regular block partitioning to take advantage of data locality during the first convection phase, and, conversely, the need to balance the highly variable workloads during the reaction phase. This second point would entail a different distribution of the simulation grid on the processor local memories. Unfortunately, since the workload distribution evolves and depends on the specific features of the input data set, we cannot devise at compile time a distribution that is suitable for both phases of the simulation. The following are therefore possible solutions.

- (i) CYCLIC distribution of the simulation grid thus losing data locality in favor of a better balance of the processor workloads. This solution must be followed to obtain acceptable performances with HPF-like data-parallel languages, which currently do not support irregular applications [5]. A combination of BLOCK and CYCLIC distribution may also be used to find a tradeoff between locality exploitation and load balancing. However, static optimizations like these may be unsuccessful when the workload distribution varies unpredictably at run time.
- (ii) Using BLOCK distribution for the first loop, and CYCLIC for the second. This solution would entail using an executable redistribution directive. Although several HPF implementations do not support yet this directive, we can simulate it by assigning one array distributed CYCLIC to another distributed BLOCK, and vice versa.
- (iii) BLOCK partitioning the grid in the first convection phase, and then suitably redistributing data and iterations to balance the workloads before starting the reaction computation. This approach is different from the previous one, because data and computation reallocation is minimized, and is carried out at run time on the basis of information collected during previous phases of the computation. This implementation scheme is followed by Ponnusamy *et al.* in [21,18]. They exploit the CHAOS run-time library [20], whose functionalities are used to build at run-time an irregular data redistribution plan resulting in a nearly balanced reaction computation. Experiments were conducted by the authors on up to 64 Intel iPSC/860 nodes [18]. Redistribution for achieving load balance takes about 15-20% of the total execution

time, but results in execution times 4–10 times better than without load balance actions. The authors do not compare the approach with the first, much simpler solutions.

- (iv) BLOCK partitioning the grid to exploit data locality, and dynamically redistributing the expensive calculations of the reaction loop to balance the load. This is the approach discussed in this paper. Non-uniform loops are executed by adopting, as far as possible, a static scheduling policy, and dynamically moving data and iterations among processors only if a load imbalance is actually detected. SUPPLE builds the balanced schedule from scratch at every loop execution: this is useful to efficiently implement highly adaptive application codes in which the workload varies at each iteration.

### 3 Overview of SUPPLE

SUPPLE is a run-time support for the efficient implementation of data parallel loops with regular, *stencil* references [19], characterized by per-iteration computation times which may be non-uniform, thus generating a workload imbalance. SUPPLE only supports BLOCK array distributions to exploit locality deriving from stencil computations. Since stencil references are regular and known at compile time, several optimizations such as *message vectorization*, *coalescing* and *aggregation* [7] can be carried out to reduce communication overheads. These optimizations combine element messages which are then sent to the same processor in one single message without redundant data. Another optimization technique implemented in SUPPLE is *iteration reordering*. The local loops assigned to every processor according to the BLOCK distribution and the owner computes rule are split into multiple localized loops: those accessing only local data, and those accessing remote data as well. The loops accessing only local data are scheduled between the asynchronous *sends* and the corresponding *receives* which are used to transmit remote data: in this way the execution of the iterations accessing only local data allows communication latencies to be hidden [7]. Besides uniform loops, SUPPLE also supports non-uniform loops through the adoption of an innovative dynamic load balancing technique. Dynamic scheduling is combined with the optimized static scheduling illustrated above. This hybrid (static + dynamic) strategy only applies to loop iterations that access local data. Loop iterations that need to fetch remote data are simply scheduled as in the static case. This limitation in the dynamic scheduling is not very strong, and does not affect the ability of our technique to balance the overall workload. Consider, in fact, that the number of iterations that access remote data is usually a small fraction of the total when data sets of interesting sizes are BLOCK distributed and stencil computations apply. One of the reasons for this limitation is to avoid the introduction of irregularities in the stencil data references, thus preventing the possibility of optimizing communications. The hybrid scheduling strategy works as follows: at the beginning iterations are scheduled statically according to the optimized ordering mentioned above, and start being dynamically migrated only if a load imbalance is actually detected at run-time. The load balancing heuristic adopts local policies of imbalance detection to avoid global synchronizations, and prefetching of remote iterations to hide communication latencies.

This section presents SUPPLE in detail. In particular, we describe the optimized implementation

of stencil communications, and both the static and hybrid (static + dynamic) scheduling strategies that can be adopted. SUPPLE supports, however, different processor layouts, multi-dimensional arrays distributed in any dimensions, and parallel loops with different stencil references. To this end, SUPPLE uses several descriptors to store information about the arrays and loops involved, and provides routines to fill the descriptor fields used at run-time by the implementation code.

**Management of stencil references.** SUPPLE provides the support to implement uniform loops by adopting the scheduling and the communication optimizations that can be employed when data are accessed with regular stencils [7,19]. It allocates on each processing node enough memory to host the array partitions assigned according to the BLOCK distribution, and a surrounding *ghost* (or *overlap*) area. The *ghost* area width depends on the shape of the stencil, and is used to buffer the *perimeter* areas of the adjacent partitions owned by neighboring processors, and accessed through non-local references.

In the final SPMD program which implements the loop, an equal fraction of loop iterations is assigned to each processor according to the owner computes rule. Each processor executes a *localized* loop, whose boundaries and array references refer to the local array block. Before executing the loop, however, non local data must be fetched from the nearest neighboring nodes. SUPPLE allow communications required to implement stencil computations to be optimized, by avoiding sending several messages or replicated data to the same processor [7]. However, if all communication are scheduled before the localized loop execution, we do not hide communication latency. To overlap communications with useful computations, SUPPLE allows iteration to be reordered [7]. To this end, if one considers the iterations executed by a localized loop, she/he can distinguish those assigning the perimeter area and those assigning the remaining *inner* area of the array blocks. Loop iterations assigning the inner area only refer to local data, while the others, which also access the ghost area, need to wait for non local data. Hence, we can reorder the loop iterations by executing the iterations that update the inner area between the asynchronous sending of messages and the corresponding receiving. Finally, in the SPMD code that implements the parallel loop we can distinguish four sequential steps:

1. < pack and send perimeter areas to neighboring processors >
2. < execute localized loops that update inner areas >
3. < receive from neighboring processors and unpack ghost areas >
4. < execute localized loops that update perimeter areas >

Steps 2 and 4 include the execution of the new parallel loops obtained by iteration reordering, while steps 1 and 3 correspond to communications that implement stencil data references. Nevertheless, if a parallel loop accesses several arrays, steps 1 and 3 only exchange data for those arrays for which stencil references apply. While step 4 is straightforward, because it simply requires the execution of the iterations which modify data elements belonging to the perimeter area, step 2 may be implemented by exploiting either a static or a hybrid (static + dynamic) scheduling technique, where the right choice depends on whether the loop to be implemented is uniform or not. In order to support hybrid scheduling, the iterations that assign the elements lying on the inner area are statically grouped into *chunks* by *tiling* the iteration space. Tiling is necessary because SUPPLE migrates chunks of iterations on the basis of the workload distribution. Tiling is not strictly needed

```

Static_Scheduler (queue  $Q$ )
 $T$ : chunk;
begin
  foreach  $T \in Q$ 
    while (* the full/empty flag of  $T$  is set *) do
      Wait_for_coherency( $T$ );
    end while
    Execute ( $T$ );
  end foreach
end

```

Fig. 2. Pseudo-code of step 2 when static scheduling is adopted.

if static scheduling is exploited, but we adopted it in this case too because it has been proved that it may improve locality exploitation [11]. Another problem considered in SUPPLE is data coherence when hybrid scheduling is adopted: in fact some tiles of a BLOCK partition may be updated by processors other than the owner of the partition itself and must be returned to the owner. To this end SUPPLE associates a *full/empty-like* flag [1] with each data tile. If a chunk that modifies a given tile is migrated, then the associated flag is set. The flag will be unset when the updated data tile is received back from the remote processor. Flags are checked only when the corresponding tiles need to be read. This mechanism does not introduce large overheads, because flags are maintained for every tile and not for every array element. Moreover, the mechanism is useful to avoid the introduction of a global synchronization between executions of distinct loops, where, for example, the first loop updates an array that is read by a following one. Note that if the first loop adopts hybrid scheduling, the following loop has to check flags even if it is uniform and thus chunks can be scheduled statically. Fig. 2 shows the SPMD pseudo code that implements step 2 of the general technique described above. All the code is included in the routine *Static\_Scheduler()*, which statically schedules the chunks of iterations by accessing local data only. The queue  $Q$  contains this chunk sequence. Note the subroutine *Wait\_for\_coherency()*, which guarantees the correct program semantics by waiting until the data tile associated with  $T$  becomes coherent, and the subroutine *Execute (T)*, which is a simple localized loop that updates the associated data tile. The results discussed in Section 4 show that the routine *Wait\_for\_coherency(T)* does not actually block computation because, when it is invoked, the corresponding coherence message has already been received. Note that, however, the **while** loop which waits for coherency can be completely removed if the data tiles of the arrays involved are not scheduled dynamically by other parallel loops. As far the benchmark illustrated in Section 2 is concerned, we have used such a static scheduling scheme to implement the uniform *Convection* loop (see Fig. 1). The width of the perimeter area is in this case equal to one, while a ghost area of the same width surrounds the blocks of array B. The resulting inner area is tiled, and the chunks are scheduled statically. Each chunk updates the array A and reads both arrays B and C. Before the execution of each chunk, according to the code in Fig. 2, we wait for the coherence of the corresponding tile of array C. Coherence is checked because the array C is written by the previous non-uniform *Reaction* loop, which is implemented by adopting the hybrid scheduling scheme described below.

**Non-uniform loop implementation.** The innovative feature of SUPPLE is its support for non-uniform loops, which is attained by adopting a hybrid (static + dynamic) scheduling strategy to implement step 2 of the general technique described in the previous section. At the beginning, to reduce overheads, iterations are scheduled statically, by sequentially executing the chunks stored

in queue  $Q$ , hereinafter called *local* queue. Once a processor understands that its local queue is *becoming empty*, the dynamic part of our scheduling policy is initiated. Computations are migrated at run-time to balance the workload. Work migration requires both chunk identifiers and corresponding data tiles to be transmitted over the interconnection network. If loops access an array with a given stencil data reference, tiles with a suitable surrounding area must be migrated. In order to overlap computation with communication, an aggressive prefetching policy is adopted to avoid processors becoming idle while waiting for further work. Migrated chunks are stored by each receiving processor in a queue  $RQ$ , called *remote*. The size  $g$  of chunks is fixed and is decided statically. This size gives only a lower bound on the amount of work migrated at run-time, because the support adopts a heuristics to decide at run-time how many chunks must be moved from an overloaded processor to a underloaded one. On the other hand, since SUPPLE uses a *polling* technique to probe message arrivals,  $g$  directly influences the time elapsed between two consecutive *polls* (see Section 4). The dynamic scheduling algorithm is fully distributed and asynchronous: chunk migration decisions are taken on the basis of local information only, and, since it is the *receiver* of remote chunks that asks overloaded processors for work migration, it can be classified as a *receiver initiated* load balancing technique. According to the framework proposed by Willebeek-LeMair and Reeves [24], the load balancing technique can be characterized by considering the strategies exploited for processor load evaluation, load balancing profitability determination, task migration, and task selection.

**Processor Load Evaluation.** Each load balancer policy requires a reliable workload estimation to detect load imbalances and take, hopefully, correct decisions about workload migration. Our technique is based on the assumption that each processor can derive the expected, average execution cost of its own chunks stored in the local queue  $Q$ . Since the first part of our technique is static, at the beginning each processor only executes chunks belonging to  $Q$ . During this phase, the average chunk cost is derived through a non-intrusive code instrumentation (the parallel system used, a Cray T3D, allows to read the clock cycle count in a single, inexpensive, machine instruction). Moreover, when chunks are migrated toward underloaded processors, an estimate of their cost is communicated as well. Each processor is thus aware of the average costs of remote chunks stored in  $RQ$ . On the basis of the cost estimates of chunks still stored in  $Q$  and  $RQ$ , each processor can therefore measure its own *current load*.

**Load Balancing Profitability Determination.** This strategy is used to evaluate the potential speedup obtainable through load migration weighted against the consequent overheads. In SUPPLE, each processor detects a possible load imbalance, and starts the dynamic part of the scheduling technique on the basis of local information only. It compares the value of its *current load* with a machine-dependent *Threshold*. When the estimated local load becomes lower than *Threshold*, the processor begins to ask other processors for remote chunks. The same comparison of the *current load* with *Threshold* is used to decide whether a chunk migration request should be granted or not. A processor  $p_j$ , which receives the migration request from a processor  $p_i$ , will grant the request only if its *current load* is higher than *Threshold*. Note that the *Threshold* parameter is chosen high enough to prefetch remote chunks, thus avoiding underloaded processors becoming idle while waiting for chunk migrations.



**Task Migration Strategy.** Sources and destinations for task migration are determined. In our case, since the load balancing technique is receiver initiated, the source (*sender*) of a load migration is determined by the destination (*receiver*) of the chunks. The *receiver* selects the processor to be asked for further chunks by using a *round-robin* policy. This criterion was chosen for its simplicity and also because tests showed its effectiveness. SUPPLE uses the same *Threshold* parameter mentioned above to reduce overheads of our task migration strategy (overheads should derive from requests for remote chunks which cannot be served because the *current load* of the partner that has been asked for is too low). Thus each processor, when its *current load* becomes lower than *Threshold*, broadcasts a so-called *termination message*. Our round-robin strategy can thus select a processor from those that have not yet communicated their termination. Termination messages are also needed to end the execution of a parallel loop. When a processor has already received a termination message from all the other processors, and both its queues  $Q$  and  $RQ$  are empty, then it can locally terminate the parallel loop since no further chunk can be fetched from both its local queues and the remote processors.

**Task Selection Strategy.** Source processors select the most suitable tasks for effectively balancing the load, and send them to the destinations that asked for load migration. In terms of our support, this means choosing the most appropriate number of chunks that must be moved to grant a given migration request. We use a modified *Factoring* scheme to determine this number [8]. *Factoring* is a *Self Scheduling* heuristics formerly proposed to address the efficient implementation of parallel loops on shared-memory multiprocessors. It provides a way to determine the appropriate number of iterations that each processor must fetch at each access from a central queue storing the indexes of unscheduled loop iterations. Clearly, the larger this number is, the lower the contention overheads for accessing the shared queue, and the higher the probability of introducing load imbalances when the fetched iterations are actually executed. *Factoring* requires that, if  $u$  is the number of remaining unscheduled iterations and  $P$  is the number of processors involved, the next  $P$  requests for new work are granted with a bunch of  $\frac{u}{2 \cdot P}$  iterations. Consequently, at the beginning large bunches of iterations are scheduled, and this number is reduced when the shared queue is going to be emptied. In our case, instead of having a single queue, we have multiple shared queues, one for each processor. When a loop is unbalanced only some of these queues, i.e. those owned by overloaded processors, will be involved in granting remote requests for further work. The other difference regards our scheduling unit, which is a chunk of iterations instead of a single iteration. The modified heuristic we used is thus the following: an overloaded processor replies to a request for further work by sending  $\frac{k}{2 \cdot P}$  chunks, where  $k$  is the number of chunks currently stored in  $Q$ , and  $P$  is the number of processors. Note that, in order to exactly apply *Factoring*, each processor should have the global knowledge of all the chunks stored into the local queues of all the overloaded processors, and also of all the possible underloaded processors that can ask for further work. If this knowledge was available and the original *Factoring* technique was thus applied, the number of chunks returned for each request would be larger than  $\frac{k}{2 \cdot P}$ . Unfortunately, this knowledge cannot be achieved without introducing global synchronizations due to the asynchronous behavior of processors involved. However, due to the multiple and concurrent requests from underloaded processors, if we moved more than  $\frac{k}{2 \cdot P}$  chunks, we would risk prefetching too much chunks toward underloaded processors.

```

Hybrid_Scheduler (queue  $Q$ )
 $RQ$ : queue;  $T$ : chunk;
begin
  Initialize_Data_Structures ();
  while (not Terminated&Empty_queues ( $Q$ ,  $RQ$ )) do
    if ( $My\_Load$  ( $Q$ ,  $RQ$ ) <  $Threshold$ ) then
      Prefetch_Chunk ();
      if (* not already done *) then
        Comm_termination ();
      end if
    end if
     $T = Extract(Q)$ ;
    if ( $Empty$  ( $T$ )) then
       $T = Extract$  ( $RQ$ );
      Execute ( $T$ );
      Send_Results_Coherence ( $T$ );
    else Execute ( $T$ );
    end if
    Request_Handler ();
  end while
end

```

Fig. 3. Pseudo-code of step 2 when hybrid scheduling is adopted.

### 3.1 SPMD code of the load balancer.

Fig. 3 shows the SPMD pseudo code to implement step 2 of the general technique above. The code exploits the hybrid (scheduling + dynamic) scheduling used by SUPPLE to implement non-uniform parallel loops. All the code is included in the routine *Hybrid\_Scheduler()*, which schedules the iteration chunks of the local  $Q$  either locally, according to a static sequential ordering, or remotely, according to dynamic scheduling decisions made only if load imbalance occurs. The core part of the code is a **while** construct, which ends looping when the function *Terminated&Empty\_queues()* returns a *True* value. This function checks the termination condition: a processor that is executing this SPMD code terminates as soon as its queues  $Q$  and  $RQ$  are empty, and it has received a termination message from all the other processors. The function *My\_load()* estimates the current *workload* of a processor. It returns the product of the number of chunks still stored in  $Q$  times the average local chunk execution cost, plus the the cost estimated for any chunks stored in  $RQ$ .  $Q$  chunk average cost is determined by monitoring the execution of the local chunks, while estimated costs of remote chunks are communicated by sender processors along with data tiles. At the beginning, when no chunk executions have been monitored, these costs are initialized with high values to avoid wrong chunk requests and migrations.

The prefetching policy adopted to reduce overheads of chunk migration is driven by a machine-dependent *Threshold* parameter: a processor sends a request for remote chunks by calling subroutine *Prefetch\_Chunk()* only if its *current load* is lower than *Threshold*. Likewise, the processor that has been asked for, grants the request only if its own load is greater than *Threshold*. Since latencies of chunk migrations may be very large, we may risk prefetching too much because we might send many requests before the first remote chunks arrive. To avoid this, we fix a limit on the number of requests sent by a processor and not yet granted. *Prefetch\_Chunk()* sends a request message only if this limit has still not been reached. *Prefetch\_Chunk()* employs a simple local *round-robin* policy to select the next processor to be asked for a remote chunk: processor  $p_i$ , whose round-robin

counter is initialized to  $p_{i+1}$  at the beginning, chooses the next partner from those which have not yet communicated their termination by means of a call to subroutine *Comm\_termination()*. The function *Extract()* returns a chunk from a queue. The subroutine *Execute()* executes a chunk  $T$ , and, if  $T$  is local, measures the time taken to complete it and consequently updates the estimate average cost of local chunks. We adopt an asynchronous *coherence protocol* to migrate updated data from the processor that actually executed the chunk to the processor that is the owner of the corresponding data partition. The updated data tiles are transmitted by the subroutine *Send\_Results\_Coherence()*, which also implements message vectorization by buffering more data tiles to be sent to the same processor. The subroutine inserts data tiles in a buffer, and transmits the buffer when either it becomes full or termination is detected. A *full/empty-like* technique [1] is used by our coherence protocol to avoid processing invalid data. When processor  $p_i$  sends a chunk  $b$  to  $p_j$ , it sets a flag marking the data tile to be modified by  $p_j$  as invalid. The next time  $p_i$  needs to access the same data, for example during the execution of another loop accessing the same array,  $p_i$  checks the flag and, if the flag is still set, waits for the updated data tile from node  $p_j$ . We have already seen how this check is performed in the static scheduling code of Fig. 2. The same job is performed in the code of Fig. 3 by the function *Extract()* when applied to the local queue  $Q$ . Such a control of data consistency does not entail stopping waiting for coherence messages at the end of parallel loop execution: a new loop can be started and useful computation can be overlapped with these communications. Finally, the subroutine *Request\_Handler()* performs many of the above tasks. This subroutine whose code, for the sake of simplicity, is not shown in detail, polls and handles all the messages that can arrive at a processor. For example:

- requests for work migration, which entail deciding about migration profitability and possibly extracting a bunch of chunks from  $Q$  and sending it to the requesting processor;
- coherence messages transporting updated data tiles, which entail copying the received data in the local array partitions and unsetting the corresponding full/empty flags;
- remote chunk bunches sent by overloaded processors, which have to be inserted into  $RQ$ ;
- termination messages which modify the behavior of our task migration strategy and are needed to manage termination detection.

## 4 Experimental results

We have tested SUPPLE on a 64 node Cray-T3D. The message passing layer used by SUPPLE is MPI, release CRI/EPCC 1.3, developed by the Edinburgh Parallel Computing Centre in collaboration with Cray Research Inc. The experiments are concerned with a SUPPLE implementation of the benchmark illustrated in Section 2 on a set of synthetic data sets presented in Section 4.1. We have also compared the SUPPLE implementation with an HPF-style one. For the HPF implementation we have used the Cray CRAFT-Fortran, whose language directives and compiler have been adapted in part from Rice University’s Fortran-D project [6] and Vienna Fortran [25].

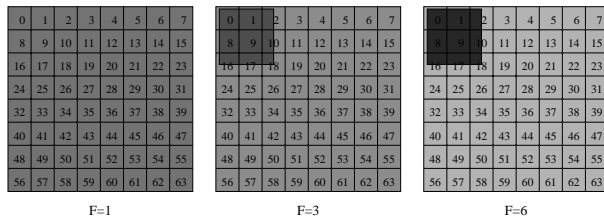


Fig. 4. Data layout of some data sets characterized by different  $F$  on a  $8 \times 8$  processor grid.

#### 4.1 The synthetic input data sets

We have adopted a set of synthetic data sets because of the need to validate SUPPLE as a general support. Conversely, if we tested SUPPLE on a single problem instance, we might have come to the wrong conclusions due to the specific features of the data set chosen for the test. According to the code skeleton shown in Fig. 1, to build the different data sets we make the following assumptions:

- the simulation grid is  $1024 \times 1024$  points. Due to the blocking data distribution,  $128 \times 128$  contiguous points are assigned to each of the 64 processors of the target machine;
- convection and reaction phases require  $\frac{1}{4}$  and  $\frac{3}{4}$  of the total execution time, respectively;
- $\mu$  is the average time needed to process a grid point during the reaction phase. Due to the previous assumption, the average per-point convection time is thus  $\frac{\mu}{3}$ ;
- each input data set is built according to a simple model of load imbalance [14] discussed in the following. Note that the imbalance features of the data set are important for the behavior of the reaction phase, where the per-point execution time depends on the values of the point itself.

**Model of load imbalance.** Given  $\mu$ , the average execution time, our model of load imbalance assumes that the workload is not distributed uniformly, and that there exists a region of the simulation grid, hereafter called *loaded region*, in which the average per-point execution time is greater than  $\mu$ . The load imbalance model is thus based on two parameters,  $d$  and  $t$ . The parameter  $d$ ,  $0 < d \leq 1$ , is a fraction of the whole simulation grid, and determines the dimension of the loaded region. The parameter  $t$ ,  $0 < t \leq 1$ , is a fraction of the total workload  $T$ , and determines the whole workload concentrated on the loaded region. Therefore, it follows that  $t \geq d$ . Note that the imbalance is directly proportional to  $t$  for a given value of  $d$ , since larger values of  $t$  correspond to larger fractions of  $T$  concentrated on the loaded region. Similarly, the imbalance is inversely proportional to  $d$  for a given value of  $t$ . From these remarks, we can derive  $F$ , called *factor of imbalance*, defined as  $F = \frac{t}{d}$

**Sample data sets.** Fig. 4 shows a representation of some different data sets we have used to force an unbalanced workload during the reaction phase of our benchmark. We employed data sets characterized by values of  $F$  ranging from 1 to 9, where the size of the loaded region is kept fixed ( $d = 0.1$ ). The grey levels used in Fig. 4 represent the workloads associated with grid points, where a darker grey stands for a heavier workload. Note that the case  $F = 1$  corresponds to a data set that does not introduce any workload imbalance, since the computational cost of each reaction point simulation is equal to  $\mu$ . The figure shows also the blocking distribution of the data

Table 1

Reaction phase: per-point execution time for some values of  $\mu$  and  $F$ 

$\mu$	$F$	Loaded points	Unloaded points
0.038 msec	2	0.075 msec	0.033 msec
0.3 msec	2	0.6 msec	0.267 msec
0.038 msec	5	0.187 msec	0.021 msec
0.3 msec	5	1.498 msec	0.167 msec
0.038 msec	9	0.337 msec	0.004 msec
0.3 msec	9	2.697 msec	0.033 msec

sets on a grid of  $8 \times 8$  processors. Due to the position of the grey loaded region, processors 0, 1, 8, and 9 appear to be overloaded, while processors 2, 10, 16, 17, and 18 are only partially loaded. We will refer to this figure to explain the overheads introduced by our support during the dynamic load balancing phase. All the results reported in this paper were obtained by keeping fixed the position of the loaded region during all the simulation iterations. We also carried out many tests in which the position of the loaded region is different or is moved at each iteration to simulate the adaptivity of the code. The difference between the results of these experiments and those reported in the paper are not appreciable, because SUPPLE does not exploit any previous knowledge about which the overloaded or the underloaded processors are.

#### 4.2 The experimental parameters

The experiments were conducted on 64 processors for different values of  $F$  and  $\mu$ , and specifically for  $F \in \{1, \dots, 9\}$ , and  $\mu \in \{0.038 \text{ msec}, 0.075 \text{ msec}, 0.15 \text{ msec}, 0.3 \text{ msec}\}$ . The external sequential loop of the flame simulation code was iterated for 10 time steps. The theoretical Optimal Completion Time (OCT) is thus 8.2 sec, 16.4 sec, 32.8 sec, and 65.5 sec, for  $\mu$  equal to 0.038 msec, 0.075 msec, 0.15 msec, and 0.3 msec, respectively. This time, which does not take into account any overheads, is computed by multiplying the average time required to process a single point of the simulation grid ( $\mu + \mu/3$ ), times the number of points assigned to each processor ( $128 \times 128$ ), times the number of simulation time steps (10). Table 1 reports some examples regarding the imbalance introduced by our synthetic data sets. In particular, it shows the times needed to compute each point during the reaction phase, and distinguishes between points belonging to the unloaded and loaded region.

**Tuning the dynamic load balancing parameters.** The only parameters that can be modified and tuned in our support are the size  $g$  of a chunk of iteration, i.e. our scheduling unit, and the prefetching *Threshold*. The parameter  $g$  affects the scheduling overhead, and the polling mechanism used to check for message arrivals. More specifically, a larger  $g$  reduces the scheduling overhead since it reduces the number of scheduling units stored in each local queue. On the other hand, a larger  $g$  may reduce the ability of our support to effectively balance the workload, mainly because, due to the polling mechanism adopted, it makes the response time of overloaded processors slower when they are asked for by underloaded ones. A tradeoff must be found, and

this depends on both the architecture and the features of the problem. The tuning of the other parameter, *Threshold*, mainly depends on the specific architecture. It should be larger on architectures where communications are not very efficient. We found, however, that it is also related to the granularity of each chunk, which in turn depends on  $g$  but also on the features of the specific problem, i.e. on the parameters  $\mu$  and  $F$ . From the remarks above it follows that the two parameters  $g$  and *Threshold* are highly interdependent. One strategy is to fix one and tune the other parameter on the basis of the first setting. Hence, all the experiments were carried out by fixing  $g = 21$ , so that each chunk assigns the iterations which update a tile of  $7 \times 3$  points of the simulation grid. The number of chunks scheduled by each parallel loop is thus 756. We changed the granularity of each chunk by using data sets characterized by different  $\mu$ . We observed that for larger  $\mu$ , i.e. for chunks with a larger granularity, it is better to increase the value of *Threshold*. This is not surprising, since for large chunk granularities, overloaded processors become lazier in answering incoming requests due to our polling implementation of input message checking. Thus, for larger  $\mu$ , prefetching further work should be begun earlier to avoid emptying both the local and the remote queues. To do this we have to increase the value of *Threshold*. Hence, we used values of *Threshold* ranging from 2 to 40 *msec*, where larger *Threshold* values are used for larger  $\mu$ . We found, however, that a heuristic resulting in a small decrease in performance with respect to the case in which the correct value of *Threshold* is established before running the program, is the following:

- fix a chunk dimension  $g$ , which depends on the size of the problem. The number of chunks must be large enough to allow the load balancing algorithm to work properly: a small number should only allow a few very large bunch of iterations to be moved, thus introducing probable wrong workload balances. Conversely, a very large number of chunks can increase the overhead introduced by our algorithm, though for some applications, we may have cache performance benefits from scheduling loop iterations using small tiles [11];
- start the load balancing algorithm with a small value of *Threshold*, and increase its value at runtime if necessary. When an underloaded processor begins to receive remote chunks along with their expected execution times, on the basis of this value it can increase, if necessary, the current *Threshold* to take into account the possible lazy behavior of remote overloaded processors. Both the initial value of *Threshold* and the amount of increase depend on the specific architecture.

### 4.3 Result assessment

Fig. 5.(a) shows several curves. Each curve is related to a distinct  $\mu$ , and plots the difference between the completion time of the slowest processor and the theoretical OCT for various factors of load imbalance  $F$ . The difference shown in this and in the following plots can be considered as the sum of all overheads due to communications, residual workload imbalance, and loop housekeeping. It therefore also includes the time to send/receive messages by implementing the stencil communications of the first convection loop. There is an overhead of almost 0.5sec even in the balanced case, and the overhead becomes larger when the input data set introduces a more extensive load imbalance. For  $F = 9$  and  $\mu = 0.15$  *msec*, SUPPLE balances the load with a completion time

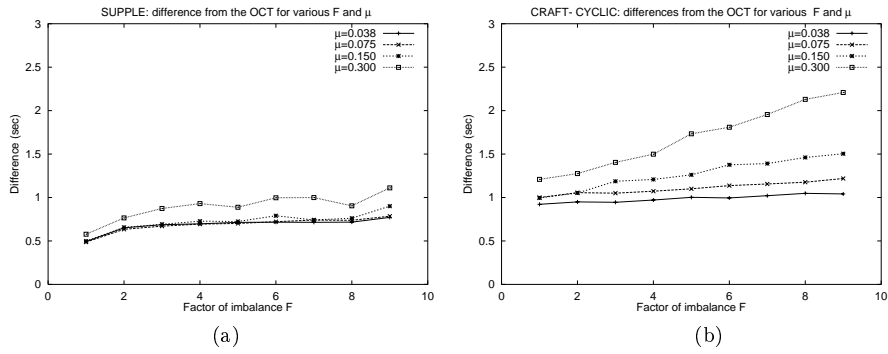


Fig. 5. Differences from the OCT for various  $\mu$  as a function of the factor of imbalance  $F$ : (a) SUPPLE results, (b) CRAFT Fortran results.

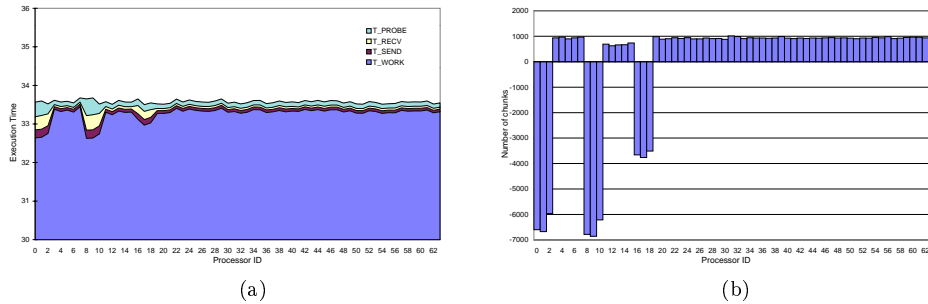


Fig. 6. Case  $F = 8$  and  $\mu = 0.15$ : (a) overheads due to *send*, *receive*, and *probe* MPI primitives, (b) chunks sent and received by each processor.

which is very close to the OCT: the difference is nearly 1 *sec*. Note that, in the absence of load balancing actions, the completion time is around 229 *sec*, with a difference from the OCT of more than 196 *sec*. If we look at the same results by considering their optimality in terms of percentage with respect to the optimum time, then for  $\mu = 0.038$  this percentage ranges from 6% to 9.4%, while for  $\mu = 0.3$  it ranges between 0.88% and 1.69%.

**Execution profile.** We have profiled the execution of some tests carried out with different values of  $F$  and  $\mu$ . Fig. 6.(a) shows, for all 64 processors, the time spent by each of them doing: (1) useful work on local or remote chunks, but also some small service computations requested by our load balancing algorithm; (2) sending, receiving, and checking for arrivals of messages. Fig. 4 helps to understand this graph, which is relative to the case  $F = 8$  and  $\mu = 0.15$  *msec*, since it identifies the overloaded processors with respect to the input data set and the data layout. Recall that processors identified by 0, 1, 8, and 9 are overloaded, while processors 2, 10, 18, 16, and 17 are partially loaded. Going back to the graph of Fig. 6.(a), one can note a larger send/receive/probe overhead for the processors above. The reason for this overhead is the large amount of messages that they must dispatch to give away chunks, and to receive migration requests and coherence messages. Moreover, note that, due to these overheads, the processors that execute less useful work are the same processors that are overloaded at the beginning of computation. The graph in Fig. 6.(b) plots, for each processor, the expression  $(\text{Chunks\_executed} - \text{Chunks\_assigned})$ .

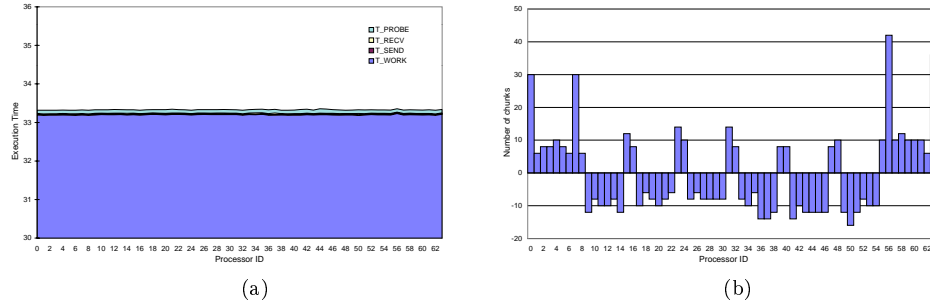


Fig. 7. Case  $F = 1$  and  $\mu = 0.15$ : (a) overheads due to *send*, *receive*, and *probe* MPI primitives, (b) chunks sent and received by each processor.

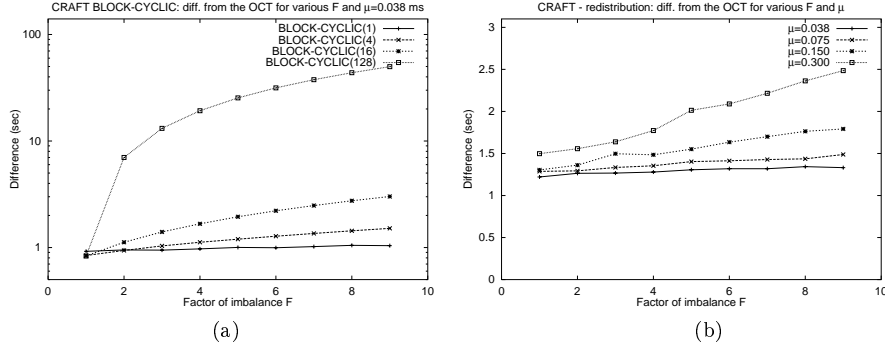


Fig. 8. CRAFT Fortran differences from the OCT: (a) for different data layouts with  $\mu = 0.038$  msec, (b) for different  $\mu$  with array redistribution.

*Chunks\_assigned* is a constant: it is the number of chunks assigned to the local queue  $LQ$  of each processor during all the iterations of the reaction parallel loops. This number, considering that the simulation is iterated 10 times, is equal to 7560. *Chunks\_executed* is the interesting datum, and corresponds to the number of chunks (local + remote) that a processor actually executes. The expression is positive if a processor executes more chunks than those statically assigned: this is the case, as Fig. 6.(b) shows, of the underloaded processors, which receive further remote chunks. The same expression is negative if a processor executes less chunks than the initial ones: this is the case of the overloaded processors, which give away a lot of chunks to underloaded processors. Fig. 7.(a) illustrates the moderate and uniform overheads introduced by SUPPLE for a balanced case characterized by  $F = 1$  and  $\mu = 15$  msec. Fig. 7.(b) shows, for the same case but with a different scale, the expression ( $Chunks\_executed - Chunks\_assigned$ ). We can note a sort of symmetry between processors that export and import chunks. If we look at Fig. 4 (in this case, since  $F = 1$ , consider that the workload is uniformly distributed), we can discover that those processors that receive a few remote tasks are the same as those that are assigned blocks on the borders of the simulation grid, and processors on the corners are those that import the most chunks. This behavior is due to the first convection loop of our benchmark: this loop does not compute, in fact, the elements on the boundaries of the simulation grid, and thus introduces a low load imbalance which is detected and managed by SUPPLE during the reaction phase of the computation.



**Comparison with CRAFT.** We implemented the same benchmark with CRAFT Fortran, by feeding it with the same synthetic data sets. We tested several different data layouts. The best results for unbalanced problem instances were obtained by using a pure CYCLIC distribution. One of the reasons of these performance results is the computational weight of the reaction loop, which is three times that of the first convection loop. Thus, due to this feature of our benchmark, it is more important to try to balance the workload through a pure CYCLIC distribution than to increase locality exploitation with a BLOCK distribution. Fig. 5.(b) shows the results attained by CRAFT Fortran by adopting a CYCLIC distribution. The differences from the OCT are larger than those obtained by SUPPLE, even for the completely balanced case (nearly a half second more than SUPPLE). These bad results for  $F = 1$  are also due to the CYCLIC distribution, which should be turned to BLOCK to take advantage of data locality in the first convection loop. The results are worse for larger  $\mu$  and  $F$ . For example, for  $F = 9$ , the CRAFT differences from the optimal times are 1.04, 1.22, 1.50, 2.21 *sec* for  $\mu$  equal to 0.038, 0.075, 0.15, 0.3 *msec* respectively, against the 0.77, 0.78, 0.90, 1.11 *sec* obtained on the same tests by SUPPLE. The CYCLIC distribution, besides losing data locality, is not able to perfectly balance the workload, so that, for larger  $\mu$ , the workload differences between overloaded and underloaded processors increase. We also conducted other tests with CRAFT Fortran to evaluate combinations of BLOCK and CYCLIC distributions. By adopting a pure BLOCK distribution for the completely balanced case ( $F = 1$ ) we obtained, as expected, the best CRAFT results. Also in this case, however, the SUPPLE results are better than the CRAFT ones. We believe that this is due to our very efficient implementation of stencil communications, the absence of barrier synchronizations, the small overheads introduced by our load balancing technique, and some migrations of chunks toward the processors holding the borders of the simulation grid (see Fig. 7.(b)). Some of the results obtained by adopting a combination of BLOCK and CYCLIC distributions are shown in Fig. 8.(a). All the curves, which refer to the case  $\mu = 0.038$  *msec*, plots the differences from the optimum time for different values of  $F$ . The scale of the ordinate axis is logarithmic due to the very large differences from the optimal times when adopting a BLOCK data distribution for unbalanced cases. Finally, we tested the redistribution of the simulation grid between consecutive parallel loop executions. CRAFT Fortran does not support run-time redistribution directives, so we used two different arrays, the former distributed BLOCK to use in the first convection loop, and the latter distributed CYCLIC to use in the second reaction phase. We re-assign one array to the other before the execution of each parallel loop. The results, which are shown in Fig. 8.(b) are worse than those obtained by using a pure CYCLIC distribution for both loops. The reason for this worse performance depends on the huge volume of data moved simultaneously between non adjacent processors which may cause network congestion problems. With regard to the communication volumes, we have measured the amounts of data transferred during 10 time steps of the whole simulation for  $\mu = 0.15$  *msec*. SUPPLE transfers from 2.5 to 19.4 MB for  $F = 1$  and  $F = 9$ , respectively. Independently from the factor of imbalance, the communication volumes are 2.3, 320 and 159.8 MB for the BLOCK, CYCLIC, and REDISTRIBUTION CRAFT versions of the benchmark, respectively. Note that for CRAFT CYCLIC, most of the communications derive from the implementation of the first convection loop, while for CRAFT REDISTRIBUTION it derives from the assignment statements used to redistribute the arrays. Data volumes moved in the case of a CYCLIC data layout are larger than the ones for REDISTRIBUTION. However, while in the

CYCLIC implementation messages are exchanged only among nearest neighbor processors, array redistribution requires that each processor communicates with all the others.

## 5 Related work

The parallel loop scheduling problem has been investigated in depth by researchers working on shared-memory multiprocessors. Most proposals address the efficient implementation of loops by defining dynamic *Self Scheduling* policies which reduce synchronizations among processors by enlarging parallel task granularity. The main goal of these works is to determine the optimal size for the chunks fetched by each processor at each access to a shared queue which stores iteration indexes. Clearly, the larger the chunk size is, the lower the contention overheads for accessing the shared queue, and the higher the probability of introducing load imbalances. Polychronopoulos and Kuck have proposed *Guided Self Scheduling*, according to which  $\frac{u}{P}$  iterations, where  $u$  is the number of remaining unscheduled iterations and  $P$  is the number of processors involved, are fetched at each time by an idle processor [17]. *Trapezoid Self Scheduling* [23] has been proposed by Tzen and Ni to reduce the number of synchronizations by linearly decreasing the chunk size. Hummel, Schonberg and Flynn have presented *Factoring* [8], the policy adopted also in SUPPLE to implement the *task selection* strategy of our load balancer (see Section 3). The introduction of large caches on shared-memory multiprocessors makes the schema above unsuitable because they do not guarantee the exploitation of locality. Markatos and LeBlanc [13] have investigated locality and data reuse to obtain scalable and efficient implementations of non-uniform parallel loops. They have explored a scheduling strategy, based on a static partitioning of iterations, which initially assigns iterations for *affinity* with previously assigned ones. Affinity regards the presence of accessed data in the processor caches. The dynamic part of the technique, which performs runtime load balancing, is postponed until a load imbalance occurs. Liu and Saletore have worked on Self Scheduling techniques for distributed-memory machines [12]. They have designed hierarchical and distributed implementations of a shared queue manager, and have investigated partial data replication as a way to solve the data locality problem. Plata and Rivera [16] presents another centralized solution on distributed-memory systems. Differently from their proposal, SUPPLE scheduling policy is fully distributed and does not introduce bottlenecks which may jeopardize the efficiency when many processors are used.

Willebeek-LeMair and Reeves [24] have presented several general load balancing strategies for multicomputers. They have introduced the framework that is used in Section 3 to characterize our load balancer. Another interesting work that presents general load balancing techniques is the one by Kumar et al. [10]. They introduce a technique that adopts a *global round-robin* policy to select the processing node to which a further task must be requested. The technique does not assume any knowledge of the load, and thus it may not be very accurate in scheduling decisions, but it does not waste any time evaluating the best load balancing choices. Kumar et al. have shown that the technique is actually very scalable, provided that an efficient contention-free implementation of the *global round-robin* is adopted. On the other hand, SUPPLE adopts a local round-robin technique to select the source of a task migration (*Task Migration* strategy), where the round-robin

counter of each processor  $p_i$  is initialized to  $p_{i+1}$  at the beginning. This solution avoids contention problems, but it would result in less accurateness if underloaded processors that select the partners to be requested for further work did not have any knowledge about their load. To this end, to avoid underloaded processors to ask other underloaded ones for, asynchronous broadcast termination signals are sent by those processors whose local chunk queues are becoming empty: these signals states that those processors made a transition from the overloaded state to the underloaded one. A lot of research has been carried out on run-time supports and compilation methods for distributed-memory multiprocessors. The Fortran D project at Rice and Syracuse universities [6], the Vienna Fortran project at the University of Vienna [25], the Fx Fortran compiler developed at Carnegie Mellon [3], and the CHAOS project at the University of Maryland [20], are only some of the most important projects addressing these topics. A number of compile-time optimizations have been designed to increase the performances obtained on distributed-memory multiprocessors [25,6,7]. The solution to problems whose irregularities prevent compile-time optimizations, is generally addressed by exploiting *inspector-executor* codes which collect information at run-time, and then use this information to optimize the successive computations [2]. The CHAOS/PARTI library [20] developed by the group headed by Joel Saltz at the University of Maryland, adopted in the Vienna Fortran Compilation System and in other HPF-like compilers, is the most notable example of this approach. It provides support for irregular array distribution, loop scheduling based upon the *almost owner computes* rule, run-time data redistribution, data accessed through indirection arrays, thus offering general support for irregular/adaptive codes whose behavior cannot be anticipated until run-time. Non-uniform parallel loops can be implemented with the CHAOS library by collecting information about iteration execution time at run-time, and by consequently building an irregular data layout resulting in balanced processor workloads [21,18]. Data partitioning and redistribution have to be repeated several times if the computational costs associated with the data items involved vary as execution progresses.

## 6 Conclusions

We have described SUPPLE, an innovative run-time support, which can be used to compile either uniform or non-uniform parallel loops with regular stencil data references. SUPPLE exploits general BLOCK distributions of the involved arrays to favor data locality exploitation, and uses the static knowledge of the data layout and of the regular stencil references to perform well-known optimizations such as *message vectorization*, *coalescing* and *aggregation*. *Iteration reordering* is exploited to hide the latencies of the communications needed to fetch remote data. SUPPLE supports two different loop scheduling strategies. Parallel loops characterized by iterations whose costs are uniform, can be scheduled according to a pure static strategy based on the *owner computes rule*. A hybrid scheduling strategy is used instead if the load assigned to the processors according to the owner computes rule is unbalanced. In this case loop iterations are initially scheduled statically, and, if a load imbalance is actually detected, an efficient dynamic scheduling, which requires data tiles and iteration indexes to be migrated, is started. We have shown that SUPPLE hides most dynamic scheduling overheads produced by the need for monitoring the load and moving data, by

overlapping them with useful computations. To this end it exploits an aggressive prefetching technique that tries to avoid underloaded processors from becoming idle while waiting for workload migration requests to be satisfied. Nevertheless, since data may be migrated and updated remotely due to dynamic scheduling decisions, and the same data must be kept coherent on the original owner, SUPPLE adopts a fully asynchronous coherence protocol that allows useful computations to be executed while waiting for data to become coherent. At our best known, SUPPLE is the first distributed memory run-time support that exploits this kind of dynamic techniques to compile data parallel languages. All the previous approaches to the compilation of loops adopt more static and synchronous approaches. A typical example is CHAOS, which collects performance information during previous execution of a loop, and prepare an irregular data redistribution plan to obtain a more balanced execution. This solution require synchronizations between distinct loop execution, and, in addition, does not work when a non-uniform loop must be executed only one time. Moreover, SUPPLE can be profitably used for a less narrow class of applications [15], i.e. to compile “uniform” parallel loop on heterogeneous and/or time-shared systems, where load imbalance is due to the target architecture rather than to the specific application. The results of our experiments conducted on a Cray T3D machine have shown that the performances obtained are very close to the optimal ones. To validate SUPPLE as a general support for non-uniform loops we employed a simplified flame simulation application and different synthetic data sets built on the basis of a simple load imbalance model. These data sets, fed in input to the flame simulation code, introduce different workload imbalances. In most of the experiments we obtained completion times which differ by less than 1 *sec* from the optimal completion times. For data sets characterized by an average iteration execution times of 0.3 *msec*, the overheads due to communications and the residual load imbalance range between 0.88% and 1.69% of the optimal times. In addition, we have compared the SUPPLE results with those obtained running HPF-like implementations of the same benchmark. In all the cases tested, SUPPLE obtained better performances.

## References

- [1] R. Alverson et al. The Tera computer system. In *Proc. of ACM ICS*, pages 1–6, 1990.
- [2] R. Das, M. Uysal, J. Saltz, and Y.-S. Hwang. Communication optimizations for irregular scientific computations on distributed memory architectures. *JPDC*, 22(3):462–479, 1994.
- [3] T. Gross, D. O’Hallaron, and J. Subhlok. Task parallelism in a high performance fortran framework. *IEEE Parallel and Distributed Technology*, 2(2):16–26, 1994.
- [4] High Performance Fortran Forum. *HPF Language Specification*, Ver. 1.1.
- [5] High Performance Fortran Forum. *HPF-2 Scope of Activities and Motivating Appl.*, Ver. 0.8.
- [6] S. Hiranandani, K. Kennedy, and C. Tseng. Compiling Fortran D for MIMD Distributed-Memory Machines. *CACM*, 35(8):67–80, Aug. 1992.
- [7] S. Hiranandani, K. Kennedy, and C. Tseng. Evaluating Compiler Optimizations for Fortran D. *JPDC*, 21(1):27–45, April 1994.

- [8] S.F. Hummel, E. Schonberg, and L.E. Flynn. Factoring: A Method for Scheduling Parallel Loops. *CACM*, 35(8):90–101, Aug. 1992.
- [9] C. W. Kessler. Pattern-driven automatic program transformation and parallelization. In *Proc. of the 3rd EUROMICRO Work. on Parallel and Distr. Proc.*, pages 76–83, Jan. 1995.
- [10] V. Kumar, A.Y. Grama, and N. Rao Vempaty. Scalable Load Balancing Techniques for Parallel Computers. *JPDC*, 22:60–79, 1994.
- [11] M. S. Lam, E. E. Rothberg, and M. E. Wolf. The cache performance and optimizations of block algorithms. In *Proc. of ASPLOS IV*, pages 63–74, Santa Clara, CA, April 1991.
- [12] J. Liu and V. A. Saletore. Self-Scheduling on Distributed-Memory Machines. In *Proc. of Supercomputing '93*, pages 814–823, 1993.
- [13] E.P. Markatos and T.J. LeBlanc. Using Processor Affinity in Loop Scheduling on Shared-Memory Multiprocessors. *IEEE TPDS*, 5(4):379–400, April 1994.
- [14] S. Orlando and R. Perego. A template for non-uniform parallel loops based on dynamic scheduling and prefetching techniques. In *Proc. of ACM ICS*, pages 117–124, 1996.
- [15] S. Orlando and R. Perego. Scheduling Data-Parallel Computations on Heterogeneous and Time-Shared Environments. In *Proc. of EuroPar '98*. In printing.
- [16] O. Plata and F. F. Rivera. Combining static and dynamic scheduling on distributed-memory multiprocessors. In *Proc. of ACM ICS*, pages 186–195, 1994.
- [17] C. Polychronopoulos and D.J. Kuck. Guided Self-Scheduling: A Practical Scheduling Scheme for Parallel Supercomputers. *IEEE TOC*, 36(12), Dec. 1987.
- [18] R. Ponnusamy, J. Saltz, A. Choudary, Y-S Hwang, and G. Fox. Runtime Support and Compilation Methods for User-Specified Irregular Data Distributions. *IEEE TPDS*, 6(8):815–831, Aug. 1995.
- [19] D. A. Reed, L. M. Adams, and M. L. Patrick. Stencils and problem partitionings: Their influence on the performance of multiple processor systems. *IEEE TOC*, 36(7):845–858, July 1987.
- [20] J. Saltz et al. Runtime and Language Support for Compiling Adaptive Irregular Programs on Distributed Memory Machines. *Software Practice and Experience*, 25(6):597–621, June 1995.
- [21] J. Saltz et al. Runtime Support and Dynamic Load Balancing Strategies for Structured Adaptive Applications. In *Proc. of the 1995 SIAM Conf on Par. Proc. for Scientific Computing*, Feb. 1995.
- [22] J.M. Sipelstein and E. Blelloch. Collection-Oriented Languages. *Proc. of the IEEE*, 79(4):504–523, April 1991.
- [23] T.H. Tzen and L.M. Ni. Dynamic Loop Scheduling on Shared-Memory Multiprocessors. In *Proc. of ICPP - Vol II*, pages 247–250, 1991.
- [24] M.H. Willebeek-LeMair and A.P. Reeves. Strategies for Dynamic Load Balancing on Highly Parallel Computers. *IEEE TPDS*, 4(9):979–993, Sept. 1993.
- [25] H.S. Zima and B.M. Chapman. Compiling for Distributed-Memory Systems. *Proc. of the IEEE*, pages 264–287, Feb. 1993.