

Fortran 90D/HPF Compiler for Distributed Memory MIMD Computers: Design, Implementation, and Performance Results*

Zeki Bozkus[†], Alok Choudhary, Geoffrey Fox, Tomasz Haupt, and Sanjay Ranka

Northeast Parallel Architectures Center

3-201, Center for Science and Technology

Syracuse University

Syracuse, NY 13244-4100

{zbozkus, choudhar, gcf, haupt, ranka}@npac.syr.edu

April 1, 1993

Abstract

Fortran 90D/HPF is a data parallel language with special directives to enable users to specify data alignment and distributions. This paper describes the design and implementation of a Fortran90D/HPF compiler. Techniques for data and computation partitioning, communication detection and generation, and the run-time support for the compiler are discussed. Finally, initial performance results for the compiler are presented which show that the code produced by the compiler is portable, yet efficient. We believe that the methodology to process data distribution, computation partitioning, communication system design and the overall compiler design can be used by the implementors of HPF compilers.

*This work was supported in part by NSF under CCR-9110812 (Center for Research on Parallel Computation) and DARPA under contract # DABT63-91-C-0028. The content of the information does not necessarily reflect the position or the policy of the Government and no official endorsement should be inferred.

[†]Corresponding Author: Zeki Bozkus, NPAC, 111 College Place, Rm. 3-201, Syracuse University, Syracuse, NY 13244-4100

1 Introduction

Distributed memory multiprocessors are increasingly being used for providing high performance for scientific applications. Distributed memory machines offer significant advantages over their shared memory counterparts in terms of cost and scalability, though it is widely accepted that they are difficult to program given the current status of software technology. Currently, distributed memory machines are programmed using a node language and a message passing library. This process is tedious and error prone because the user must perform the task of data distribution and communication for non-local data access.

There has been significant research in developing parallelizing compilers. In this approach, the compiler takes a sequential program, e.g. a Fortran 77 program as input, applies a set of transformation rules, and produces a parallelized code for the target machine. However, a sequential language, such as Fortran 77, obscures the parallelism of a problem in sequential loops and other sequential constructs. This makes the potential parallelism of a program more difficult to detect by a parallelizing compiler. Therefore, compiling a sequential program into a parallel program is not a natural approach. An alternative approach is to use a programming language that can naturally represent an application without losing the application's original parallelism. Fortran 90 [1] (with some extensions) is such a language. The extensions may include the *forall* statement and compiler directives for data partitioning, such as decomposition, alignment, and distribution. Fortran 90 with these extensions is what we call "Fortran 90D", a Fortran 90 version of the Fortran D language [2]. We developed the Fortran D language with our colleagues at Rice University. There is an analogous version of Fortran 77 with compiler directives and other constructs, called Fortran 77D. Fortran D allows the user to advise the compiler on the allocation of data to processor memories. Recently, the High Performance Fortran Forum, an informal group of people from academia, industry and national labs, led by Ken Kennedy, developed a language called HPF (High Performance Fortran) [3] based on Fortran D. HPF essentially adds extensions to Fortran 90 similar to Fortran D directives. Hence, Fortran 90D and HPF are very similar except a few differences. For this reason, we call our compiler the Fortran 90D/HPF compiler.

From our point of view, Fortran90 is not only a language for SIMD computers [4, 5], but it is a natural language for specifying parallelism in a class of problems called *loosely synchronous* problems [6]. In Fortran 90D/HPF, parallelism is represented with parallel constructs, such as array

operations, *where* statements, *forall* statements, and intrinsic functions. This gives the programmer a powerful tool to express the data parallelism natural to a problem.

This paper presents the design of a prototype compiler for Fortran 90D/HPF. The compiler takes as input a program written in Fortran 90D/HPF. Its output is SPMD (Single Program Multiple Data) program with appropriate data and computation partitioning and communication calls for MIMD machines. Therefore, the user can still program using a data parallel language but is relieved of the responsibility to perform data distribution and communication.

The goals of this paper are to present the underlying design philosophy, various design choices and the reasons for making these choices, and to describe our experience with the implementation. That is, in contrast to many other compiler papers which present specific techniques to perform one or more functions, our goal is to describe the overall architecture of our compiler. We believe that the presented design will provide directions to the implementors of HPF compilers.

The rest of this paper is organized as follows. The compiler architecture is described in Section 2. Data partitioning, and computation partitioning are discussed in Sections 3, and 4. Section 5 presents the communication primitives and communication generation for Fortran 90D/HPF programs. In Section 6, we present the runtime support system including the intrinsic functions. Some optimization techniques are given in Section 7. Section 8 summarizes our initial experience using the current version of the compiler. It also presents a comparison of the performance with hand written parallel code. Section 9 presents a summary of related work. Finally, summary and conclusions are presented in Section 10.

2 Compiler System Diagram

Our Fortran90D/HPF parallel compiler exploits only the parallelism expressed in the data parallel constructs. We do not attempt to *parallelize* other constructs, such as *do* loops and *while* loops, since they are used only as naturally sequential control constructs in this language. The foundation of our design lies in recognizing commonly occurring computation and communication patterns. These patterns are then replaced by calls to the optimized run-time support system routines. The run-time support system includes parallel intrinsic functions, data distribution functions, communication primitives and several other miscellaneous routines. This approach represents a significant departure from traditional approaches where a compiler needs to perform in-depth dependency

analyses to recognize parallelism, and embed all the synchronization and low-level communication functions inside the generated code.

Figure 1 shows the components of the basic Fortran 90D/HPF compiler. Given a syntactically correct Fortran90D/HPF program, the first step of the compilation is to generate a parse tree. The front-end to parse Fortran 90 for the compiler was obtained from ParaSoft Corporation. In this module, our compiler also transforms each array assignment statement and *where* statement into equivalent *forall* statement with no loss of information [7]. In this way, the subsequent steps need only deal with *forall* statements.

The partitioning module processes the data distribution directives; namely, decomposition, distribute and align. Using these directives, it partitions data and computation among processors.

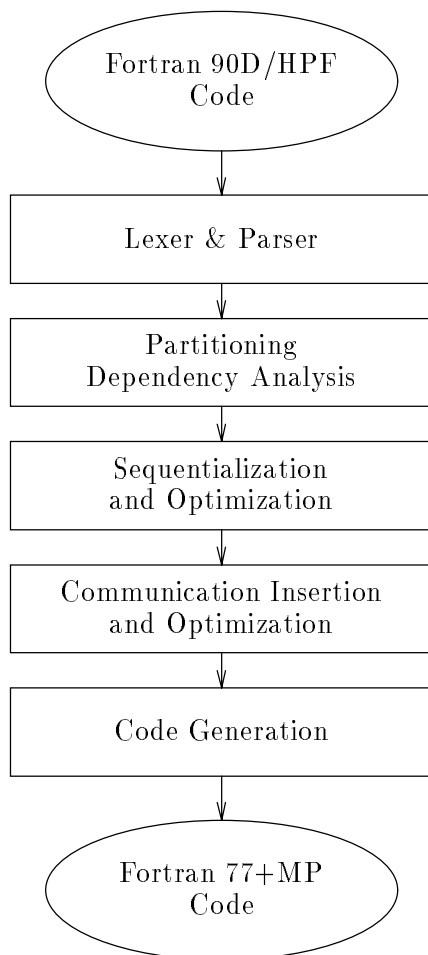


Figure 1: Diagram of the compiler.

After partitioning, the parallel constructs in the node program are sequentialized since it will be

executed on a single processor. This is performed by the sequentialization module. Array operations and *forall* statements in the original program are transferred into loops or nested loops. The communication module detects communication requirements and inserts appropriate communication primitives.

Finally, the code generator produces *loosely synchronous* [6] SPMD code. The generated code is structured as alternating phases of local computation and global communication. Local computations consist of operations by each processor on the data in its own memory. Global communication includes any transfer of data among processors, possibly with arithmetic or logical computation on the data as it is transferred (e.g. reduction functions). In such a model, processes do not need to synchronize during local computation. But, if two or more nodes interact, they are implicitly synchronized by global communication.

3 Data Partitioning

The distributed memory system solves the memory bottleneck of vector supercomputers by having separate memory for each processor. However, distributed memory systems demand high locality for good performance. Therefore, the distribution of data across processors is of critical importance to the efficiency of a parallel program in a distributed memory system.

Fortran D provides users with explicit control over data partitioning with both data *alignment* and *distribution* specifications. We briefly overview directives of Fortran D relevant to this paper. The complete language is described elsewhere [2]. The DECOMPOSITION directive is used to declare the name, dimensionality, and the size of each problem domain. We call it “template” (the name “template” has been chosen to describe “DECOMPOSITION” in HPF [3]). The ALIGN directive specifies fine-grain parallelism, mapping each array element onto one or more elements of the template. This provides the minimal requirement for reducing data movement. The DISTRIBUTE directive specifies coarse-grain parallelism, grouping template elements and mapping them to the finite resources of the machine. Each dimension of the template is distributed in either a block or cyclic fashion. The selected distribution can affect the ability of the compiler to minimize communication and load imbalance in the resulting program.

The Fortran 90D/HPF compiler maps arrays to physical processors by using a three stage mapping as shown in Figure 2 which is guided by the user-specified Fortran D directives.

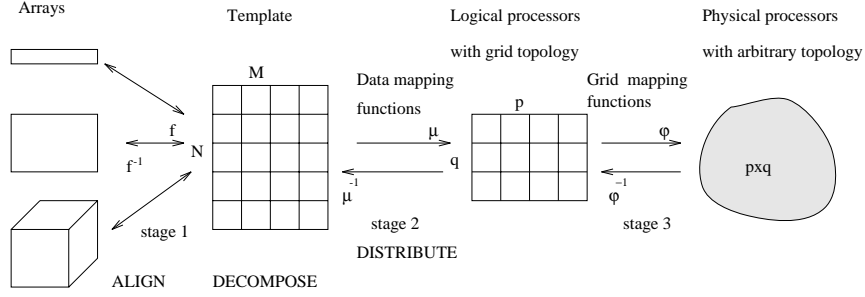


Figure 2: **Three stage array mapping**

Stage 1 : The alignment of arrays to template is determined by their subscript expressions in the ALIGN directive. The compiler computes f and f^{-1} function from the directive and applies f functions for the corresponding array indices to bring them onto common template index domain. The original indices can be calculated by f^{-1} if they are required. The algorithm to compile align directive can be found in [8].

Stage 2 : Each dimension of the template is mapped onto the logical processor grid, based on the DISTRIBUTE directive attributes. *Block* divides the template into contiguous chunks. *Cyclic* specifies a round-robin division of the template. The mapping functions μ and μ^{-1} to generate relationship between global and local indices are computed.

Stage 3 : The logical processor grid is mapped onto the physical system. The mapping functions ϕ and ϕ^{-1} can change from one system to another but the data mapping onto the logical processor grid does not need to change. This enhances portability across a large number of architectures.

By performing the above three stage mapping, the compiler is decoupled from the specifics of a given machine or configuration. Compilation of distribution directives is discussed in detail in [8].

4 Computation Partitioning

Once the data is distributed, there are several alternatives to assign computations to processing elements (PEs) for each instance of a *forall* statement. One of the most common methods is to use the *owner computes rule*. In the owner computes rule, the computation is assigned to the PE owning the *lhs* data element. This rule is simple to implement and performs well in a large number of cases. Most of the current implementations of parallelizing compilers uses the owner computes rule [9, 10].

However, it may not be possible to apply the owner computes rule for every case without extensive overhead. The following examples describe how our compiler performs computation partitioning.

Example 1 (canonical form) Consider the following statement, taken from the Jacobi relaxation program

```
forall (i=1:N, j=1:N)
&   B(i,j) = 0.25*(A(i-1,j)+A(i+1,j)+A(i,j-1)+A(i,j+1))
```

In the above example, as in a large number of scientific computations, the *forall* statement can be written in the canonical form. In this form, the subscript value in the *lhs* is identical to the *forall* iteration variable. In such cases, the iterations can be easily distributed using the owner computes rule. Furthermore, it is also simpler to detect structured communication by using this form (This will be elaborated in Section 5.2.).

Figure 3 shows the possible data and iteration distributions for the $lhs_I = rhs_I$ assignment caused by iteration instance I . Cases 1 and 2 illustrate the order of communication and computation arising from the owner computes rule. Essentially, all the communications to fetch the off-processor data required to execute an iteration instance are performed before the computation is performed. The generated code will have the following communication and computation order.

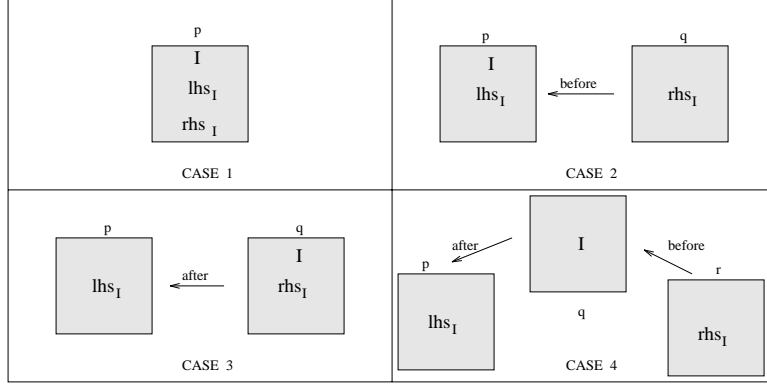
```
Communications ! some global communication primitives
Computation   ! local computation
```

Example 2 (non-canonical form) Consider the following statement, taken from an FFT program

```
forall (i=1:incrm, j=1:nx/2)
&   x(i+j*incrm*2+incrm) = x(i+j*incrm*2) - term2(i+j*incrm*2+incrm)
```

The *lhs* array index is not in the canonical form. In this case, the compiler equally distributes the iteration space on the number of processors on which the *lhs* array is distributed. Hence, the total number of iterations will still be the same as the number of *lhs* array elements being assigned. However, this type of *forall* statement will result in either Case 3 or Case 4 in Figure 2. The generated code will be in the following order.

```
Communications ! some global communication primitives to read off-processor values
Computation   ! local computation
Communication ! a communication primitive to write the calculated values to off-processors
```



CASE 1: No communications
CASE 2: Communication before computation to fetch non-local rhs
CASE 3: Communication after computation to store non-local data lhs
CASE 4: Communication before and after computation to fetch and store non-locals

Figure 3: **I** shows the processor on which the computation is performed. lhs_I and rhs_I show the processors on which the lhs and rhs of instance **I** reside.

For reasonably simple expressions, the compiler can transform such index expressions into the canonical form by performing some symbolic expression operations [11]. However, it may not always be possible to perform such transformations for complex expressions.

Example 3 (vector-valued index) Consider the statement

```
forall (i=1:N) A(U(i)) = B(V(i)) +C(i)
```

The iteration i causes an assignment to element $A(U(i))$, where $U(i)$ may only be known at run-time. Therefore, if iterations are statically assigned at compile time to various PEs, iteration i is likely to be assigned to a PE other than the one owning $A(U(i))$. This is also illustrated in cases 3 and 4 of Figure 3. In this case, our compiler distributes the computation i with respect to the owner of $A(i)$.

Having presented the computation partitioning alternatives for various reference patterns of arrays on the lhs , we now present a primitive to perform global to local transformations for loop bounds.

```
set_BOUND(l1b,lub,lst,glb,gub,gst,DIST,dim) ! computes local lb, ub, st from global ones
```

The `set_BOUND` primitive takes a global computation range with global lower bound, upper bound and stride. It distributes this global range statically among the group of processors specified

by the *dim* parameter on the logical processor dimension. The *DIST* parameter gives the distribution attribute such as *block* or *cyclic*. The *set_BOUND* primitive computes and returns the local computation range in local lower bound, local upper bound and local stride for each processor. The algorithm to implement this primitive can be found in [7].

The other functionality of the *set_BOUND* primitive is to mask inactive processors by returning appropriate local bounds. For example, such a case may arise when the global bounds do not specify the entire range of the *lhs* array. If the inputs for this primitive are compile-time constants, the compiler can calculate the local bounds at compile-time.

In summary, our computation and data distributions have two implications.

- The processor that is assigned an iteration is responsible for computing the *rhs* expression of the assignment statement.
- The processor that owns an array element (*lhs* or *rhs*) must communicate the value of that element to the processors performing the computation.

5 Communication

Our Fortran 90D/HPF compiler produces calls to collective communication routines [12] instead of generating individual processor send and receive calls inside the compiled code. There are three main reasons for using collective communication to support interprocessor communication in the Fortran 90D/HPF compiler.

1. *Improved performance estimation of communication costs.* Our compiler takes the data distribution for the source arrays from the user as compiler directives. However, any future compiler will require a capability to perform automatic data distribution and alignments [13, 14, 15]. Such techniques usually require computing trade-offs between exploitable parallelism and the communication costs. The costs of collective communication routines can be determined more precisely, thereby enabling the compiler to generate better distributions automatically.
2. *Improved performance of Fortran 90D/HPF programs.* To achieve good performance, interprocessor communication must be minimized. By developing a separate library of interpro-

cessor communication routines, each routine can be optimized. This is particularly important given that the routines will be used by many programs compiled through the compiler.

3. *Increased portability of the Fortran 90D/HPF compiler.* By separating the communication library from the basic compiler design, portability is enhanced because to port the compiler, only the machine specific low-level communication calls in the library need to be changed.

5.1 Communication Primitives

In order to perform a collective communication on array elements, the communication primitive needs the following information 1-) send processors list, 2-) receive processors list, 3-) local index list of the source array and, 4-) local index list of the destination array.

There are two ways of determining the above information. 1) Using a preprocessing loop to compute the above values or, 2) based on the type of communication, the above information may be implicitly available, and therefore, not require preprocessing. We classify our communication primitives into *unstructured* and *structured* communication.

Our structured communication primitives are based on a logical grid configuration of the processors. Hence, they use grid-based communications such as shift along dimensions, broadcast along dimensions etc. The following summarizes some of the structured communication primitives implemented in our compiler.

- **transfer:** Single source to single destination message.
- **multicast:** broadcast along a dimension of the logical grid.
- **overlap_shift:** shifting data into overlap areas in one or more grid dimensions. This is particularly useful when the shift amount is known at compile time. This primitive uses that fact to avoid intra processor copying of data and directly stores data in the overlap areas [16].
- **temporary_shift:** This is similar to overlap shift except that the data is shifted into a temporary array. This is useful when the shift amount is not a compile time constant. This shift may require intra-processor copying of data.
- **concatenation:** This primitive concatenates a distributed array and the resultant array ends up in all the processors participating in this primitive.

We have implemented two sets of unstructured communication primitives: 1) where the communicating processors can determine the send and receive lists based only on local information, and hence, only require preprocessing that involves local computations [17], and 2) where to determine the send and receive lists preprocessing itself requires communication among the processors [18]. The primitives are as follows.

- **precomp_read:** This primitive is used to bring all non-local data to the place it is needed before the computation is performed.
- **postcomp_write:** This primitive is used to store remote data by sending it to the processors that own the data after the computation is performed. Note that these two primitives requires only local computation in the preprocessing loop.
- **gather:** This is similar to *precomp_read* except that preprocessing loop itself may require communication.
- **scatter:** This is similar to *postcomp_write* except that preprocessing loop itself may require communication.

5.2 Communication Detection

The compiler must recognize the presence of collective communication patterns in the computations in order to generate the appropriate communication calls. Specifically, this involves a number of tests on the relationship among subscripts of various arrays in a forall statement. These tests should also include information about array alignments and distributions. We use pattern matching techniques similar to those proposed by Chen [19] and also used by Gupta [20]. Further, we extend the above tests to include unstructured communication.

Consider the following forall statement to illustrate the steps involved in communication detection.

FORALL (i1=11:u1:s1, i2= ..., ...) LHS(f_1, f_2, \dots, f_n) = RHS1(g_1, g_2, \dots, g_m) + ...

where g_i and f_j , $1 \leq i \leq m$, $1 \leq j \leq n$, are functions of index variables or are indirection arrays.

The steps involved in determining a communication pattern are summarized in Algorithm 1.

The algorithm first attempts to detect structured communication if the arrays are aligned to the same template. For each array on the RHS, the following processing is performed. Each

Algorithm 1 (Detecting the communication for the forall statement.)*Input:* Forall statement with untagged array and array subscripts*Output:* Forall statement with arrays and array subscripts tagged with communication primitives.*Method:*

1. **for** each RHS array **do**
2. **if** (is_aligned_same_template(LHS,RHS)) **then**
3. **for** each subscript g_i of RHS **do**
4. find f_j such that g_i and f_j are aligned with the same dimension of a template
5. **if** the pair (f_j, g_i) is in Table 1
- tag the subscript g_i with the corresponding structured communication primitive.
6. **end do**
7. **end if**
8. • **if** an untagged distributed dimension of array reference pattern is in Table 2,
- tag the RHS array with the unstructured primitives to read RHS before computation.
9. **end do**
10. • **If** a distributed dimension of LHS reference pattern is in Table 2
- tag the LHS array with the unstructured primitives to write LHS after computation
11. • **if** LHS array is not distributed
- tag the distributed RHS array with concatenation primitive.

subscript of the array is coupled with the corresponding subscript on the LHS array such that both subscripts are aligned with the same dimension of the template. For each such pair, the algorithm attempts to find a structured communication pattern in that dimension according to Table 1. If a structured communication pattern is found then the subscript on the RHS from this pair is tagged with indicating the appropriate communication primitive.

If any distributed dimension of an array on the RHS is left untagged then the array is marked with one of the unstructured communication primitives (the third column of Table 2) depending on the reference pattern given in the second column of Table 2.

The algorithm tags the LHS array as *postcomp_write* or *scatter* according to the reference patterns given in Table 2 if one or more of the distributed dimension's subscript is in non-canonical form, is vector-valued or is *unknown* at compiler time. Note that any pattern that can not be classified according to Tables 1 or 2, is marked as *unknown* (such subscripts involving more than one forall index, e.g $I+J$) so that scatter and gather can be used to parallelize any forall statement.

Table 1: Structured communication primitives based on the relationship between LHS and RHS array subscript reference patterns for block distribution. (c : compile time constant, s, d : scalar). Similar structured primitives for cyclic distributions are defined but are not presented here.

Steps	(lhs,rhs)	Comm. primitives
1	(i, s)	multicast
2	$(i, i + c)$	overlap_shift
3	$(i, i - c)$	overlap_shift
4	$(i, i + s)$	temporary_shift
5	$(i, i - s)$	temporary_shift
6	(d, s)	transfer
7	(i, i)	no_communication

Table 2: Unstructured communication primitives to read RHS data before the computation is performed and to write non-local LHS data after the computation is performed (f : invertible function, V : indirection array).

Steps	Reference pattern	Comm. primitives to read RHS	Comm. primitive to write LHS
1	$f(i)$	precomp_read	postcomp_write
2	$V(i)$	gather	scatter
3	<i>unknown</i>	gather	scatter

5.3 Communication Generation

Having recognized the type of communication in each dimension of an array for structured communication or each array for unstructured communication in a forall statement, the compiler needs to perform the appropriate program transformations. We now illustrate these transformations with the aid of some examples.

5.3.1 Structured Communication

All the examples discussed below have the following mapping directives.

```
C$ PROCESSORS(P,Q)
C$ DISTRIBUTE TEMPL(BLOCK,BLOCK)
C$ ALIGN A(I,J) WITH TEMPL(I,J)
C$ ALIGN B(I,J) WITH TEMPL(I,J)
```

Example 1 (transfer) Consider the statement

```
FORALL(I=1:M) A(I,8)=B(I,3)
```

The first subscript of B is marked as *no_communication* because A and B are aligned in the first dimension and have identical indices. The second dimension is marked as *transfer*.

```
1.  call set_BOUND(lb,ub,st,1,M,1) ! compute local lb, ub, and st
2.  call set_DAD(B_DAD,.....)      ! put information for B into B_DAD
3.  call transfer(B, B_DAD, TMP, source=global_to_proc(8), dest=global_to_proc(3))
4.  DO I=lb,ub,st
5.    A(I,global_to_local(8)) = TMP(I)
6.  END DO
```

In the above code, the *set_BOUND* primitive (line 1) computes the local bounds for computation assignment based on the iteration distribution (Section 4). In line 2, the primitive *set_DAD* is used to fill the Distributed Array Descriptor (DAD) associated with array B so that it can be passed to the *transfer* communication primitive at run-time. The DAD has sufficient information for the communication primitives to compute all the necessary information including local bounds, distributions, global shape etc. Note that *transfer* performs one-to-one send-receive communication based on the logical grid. In this example, one column of grid processors communicate with another column of the grid processors as shown in Figure 4 (a).

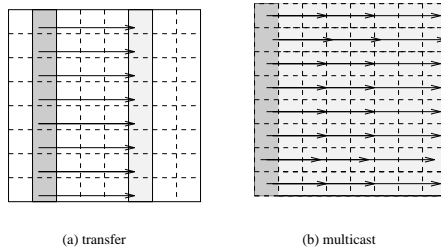


Figure 4: Structured communication on logical grid processors.

Example 2 (multicast) Consider the statement

```
FORALL(I=1:N,J=1:M) A(I,J)=B(I,3)
```

The second subscript of B marked as *multicast* and the first as *no_communication*.

```
1. call set_BOUND(lb,ub,st,1,N,1) ! compute local lb, ub, and st
2. call set_BOUND(lb1,ub1,st1,1,M,1) ! compute local lb, ub, and st
3. call set_DAD(B_DAD,...) ! put information for B into B_DAD
4. call multicast(B, B_DAD, TMP,source_proc=global_to_proc(3), dim=2)
5. DO I=lb,ub,st
6. DO J=lb1,ub1,st1
7. A(I,J) = TMP(I)
8. END DO
```

Line 4 shows a broadcast along dimension 2 of the logical processor grid by the processors owning elements $B(I,3)$ where $1 \leq I \leq N$ (Figure 4 (b).)

Example 3 (multicast_shift) Consider the statement

```
FORALL(I=1:N,J=1:M) A(I,J)=B(3,J+s)
```

The first subscript of array B is marked as *multicast* and the second subscript is marked as *temporary_shift*. The above communication can be implemented as two separate communication steps: multicast along the first dimension of logical grid *TEMPL* and *temporary_shift* along the second dimension of the logical grid. Alternatively, the two communication patterns can be composed together to obtain a better communication primitive such as the *multicast_shift* primitive.

```
call set_BOUND(lb,ub,st,1,N,1) ! compute local lb, ub, and st
call set_BOUND(lb1,ub1,st1,1,M,1) ! compute local lb, ub, and st
multicast_shift(B, B_DAD,TMP, source=global_to_proc(3),
& shift=s, multicast_dim=1, shift_dim=2)
```

```

DO I=lb,ub,st
DO J=lb1,ub1,st1
  A(I,J)=TMP(J)
END DO
END DO

```

Combining two primitives eliminates the need for creating temporary storage and eliminates some of intra processor copying, message-packing, and unpacking.

5.3.2 Unstructured Communication

In distributed memory MIMD architectures, there is typically a non-trivial communication latency or startup cost. Hence, it is attractive to vectorize messages to reduce the number of startups. For unstructured communication, this optimization can be achieved by performing the entire preprocessing loop before communication so that the schedule routine can combine the messages to the maximum extent. The preprocessing loop is also called the “inspector” loop [21, 22].

Example 1 (precomp_read) Consider the statement

```
FORALL(I=1:N) A(I)=B(2*I+1)
```

The array B is marked as *precomp_read* since the distributed dimension subscript is written as $f(i) = 2 * i + 1$ which is invertible as $g(i) = (i - 1)/2$.

```

1   count=1
2   call set_BOUND(lb,ub,st,1,N,1) ! compute local lb, ub, and st
3   DO I=lb,ub,st
4     receive_list(count)=global_to_proc(f(i))
5     send_list(count)= global_to_proc(g(i))
6     local_list(count) = global_to_local(g(i))
7     count=count+1
8   END DO
9   isch = schedule1(receive_list, send_list, local_list, count)
10  call precomp_read(isch, tmp,B)
11  count=1
12  DO I=lb,ub,st
13    A(I) = tmp(count)
14    count= count+1
15  END DO

```


The preprocessing loop is given in lines 1-9. Note that this preprocessing loop executes concurrently in each processor. It fills out the *receive_list* as well as the *send_list* of processors. Each processor also fills the local indices of the array elements which are needed by that processor.

The schedule *isch* can also be used to carry out identical patterns of data exchanges on several different but identically distributed arrays or array sections. The same schedule can be reused repeatedly to carry out a particular pattern of data exchange on a single distributed array. In these cases, the cost of generating the schedules can be amortized by only executing it once. This analysis can be performed at compile time. Hence, if the compiler recognizes that the same schedule can be reused, it does not generate code for scheduling but it passes a pointer to the already existing schedule. Furthermore, the preprocessing computation can be moved up as much as possible by analyzing *definition-use chains* [23]. Reduction in communication overhead can be significant if the scheduling code can be moved out of one or more nested loops by this analysis.

In the above example, *local_list* (line 6) is used to store the index of one-dimensional array. However, in general, *local_list* will store indices from a multi-dimensional Fortran array by using the usual column-major subscript calculations to map the indices to a one-dimensional index.

The *precomp_read* primitive performs the actual communication using the schedule. Once the communication is performed, the data is ordered in a one dimensional array, and the computation (lines 12-15) uses this one dimensional array.

Example 2 (gather) Consider the statement

```
FORALL(I=1:M) A(I)=B(V(I))
```

The array *B* is marked as requiring *gather* communication since the subscript is only known at runtime. The receiving processors can know what non-local data they need from other processors, but a processor may not know what local data it needs to send to other processors. For simplicity, in this example, we assume that the indirection array *V* is replicated. If it is not replicated, the indirection array must also be communicated to compute the receive list on each processor.

```

1   count=1
2   call set_BOUND(lb,ub,st,1,M,1) ! compute local lb, ub, and st
3   DO I=lb,ub,st
4       receive_list(count)=global_to_proc(V(i))
5       local_list(count) = global_to_local(V(i))
6       count=count+1
7   END DO
```

```

9   isch = schedule2(receive_list, local_list, count)
10  call gather(isch, tmp,B)
11  count=1
12  DO I=lb,ub,st
13     A(I) = tmp(count)
14     count= count+1
15  END DO

```

Once the scheduling is completed, every processors knows exactly which non-local data elements it needs to send to and receive from other processors. Recall that the task of *scheduler2* is to determine exactly which send and receive communications must be carried out by each processor. The scheduler first figures out how many messages each processor will have to send and receive during the data exchange. Each processor computes the number of elements (*receive_list*) and the local index of each element it needs from all other processors. In *schedule2* routine, processors communicate to combine these lists (a fan-in type of communication). At the end of this processing, each processor contains the send and receive list. After this point, each processor transmits a list of required array elements (*local_list*) to the appropriate processors. Each processor now has the information required to set up the send and receive messages that are needed to carry out the scheduled communication. This is done by the gather primitives.

Example 3 (scatter) Consider the statement

```
FORALL(I=1:N) A(U(I))=B(I)
```

The array *A* is marked as requiring *scatter* primitive since the subscript is only known at runtime. Note that owner computes rule is not applied here. The processor performing the computation knows the processor and the corresponding local-offset at which the resultant element must be written.

```

1   count=1
2   call set_BOUNDS(lb,ub,st,1,N,1) ! compute local lb, ub, and st
3   DO I=lb,ub,st
4     send_list(count)=global_to_proc(U(i))
6     local_list(count) = global_to_local(U(i))
7     count=count+1
8   END DO
9   isch = schedule3(proc_to, local_to, count)
10  call scatter(isch, A, B)

```

Unlike the *gather* primitive, in this case each processor computes a *send_list* containing processor ids and *local_list* containing the local index where the communicated data must be stored. The *schedule3* is similar to *schedule2* of gather primitives except that *schedule3* does not need to send local index in a separate communication step.

The gather and scatter operations are powerful enough to provide the ability to read and write distributed arrays with vectorized communication facility. These two primitives are available in PARTI (Parallel Automatic Runtime Toolkit at ICASE) [21] designed to efficiently support irregular patterns of distributed array accesses. The PARTI and other communication primitives and intrinsic functions form the run-time support system of our Fortran 90D compiler.

6 Run-time Support System

The Fortran 90D compiler relies on a very powerful run-time support system. The run-time support system consists of functions which can be called from the node programs of a distributed memory machine.

Intrinsic functions support many of the basic data parallel operations in Fortran 90. They do not only provide a concise means of expressing operations on arrays, but also identify parallel computation patterns that may be difficult to detect automatically. Fortran 90 provides intrinsic functions for operations such as shift, reduction, transpose, and transpose, and matrix multiplication. The intrinsic functions that may induce communication can be divided into five categories as shown in Table 3.

Table 3: Fortran90D Intrinsic Functions

1. Structured communication	2. Reduction	3. Multicasting	4. Unstructured communication	5. Special routines
CSHIFT EOSHIFT	DOTPRODUCT ALL, ANY COUNT MAXVAL, MINVAL PRODUCT SUM MAXLOC, MINLOC	SPREAD	PACK UNPACK RESHAPE TRANSPOSE	MATMUL

The first category requires data to be transferred using with less overhead structured shift communications operations. The second category of intrinsic functions require computations based on local data followed by use of a reduction tree on the processors involved in the execution of the intrinsic function. The third category uses multiple broadcast trees to spread data. The fourth category is implemented using unstructured communication patterns. The fifth category is implemented using existing research on parallel matrix algorithms [12]. Some of the intrinsic functions can be further optimized for the underlying hardware architecture. Our Fortran 90D/HPF compiler has more than 500 parallel run-time support routines and the implementation details can be found in [24].

Arrays may be redistributed across subroutine boundaries. A dummy argument which is distributed differently than its actual argument in the calling routine is automatically redistributed upon entry to the subroutine by the compiler, and is automatically redistributed back to its original distribution at subroutine exit. These operations are performed by the redistribution primitives which transform from *block* to *cyclic* or vice versa.

When a distributed array is passed as an argument to some of the run-time support primitives, it is also necessary to provide information such as its size, distribution among the nodes of the distributed memory machine etc. All this information is stored into a structure which is called *distributed array descriptor* (DAD) [24].

In summary, parallel intrinsic functions, communication routines, dynamic data redistribution primitives and others are part of the run-time support system.

7 Optimizations

Several types of *communication* and *computation* optimizations can be performed to generate a more efficient code. In terms of *computation* optimization, it is expected that the scalar node compiler performs a number of classic scalar optimizations within basic blocks. These optimizations include common subexpression elimination, copy propagation (of constants, variables, and expressions), constant folding, useless assignment elimination, and a number of algebraic identities and strength reduction transformations. However, to use parallelism within the single node (e.g. using attached vector units), our compiler propagates information to the node compiler using node directives. Since there is no data dependency between different loop iteration in the original data parallel

constructs such as *forall* statement, vectorization can be performed easily by the node compiler.

Our compiler performs several optimizations to reduce the total cost of communication. Some of *communication* optimizations [19, 25, 26] are as follows.

1. *Vectorized communication.* Vectorization combines messages for the same source and destination into a single message to reduce communication overhead. Since we are only parallelizing array assignments and forall statements in Fortran 90D/HPF, there is no data dependency between different loop iterations. Thus, all the required communication can be performed before or after the execution of the loop on each of the processors involved.
2. *Eliminate unnecessary communications.* In many cases, communication required for two different operands can be replaced by their union. For example, the following code may require two *overlapping_shifts*. However, with a simple analysis, the compiler can eliminate the shift of size 2.

```
FORALL(I=1:N) A(I)=B(I+2)+B(I+3)
```

3. *Reuse of scheduling information.* *Unstructured* communication primitives are required by computations which require the use of a preprocessor. As discussed in Section 5.3.2, the schedules can be reused with appropriate analysis.
4. *Code movement.* The compiler can utilize the information that the run-time support routines do not have procedural side effects. For example, the preprocessing loop or communication routines can be moved up as much as possible by analyzing *definition-use chains* [23]. This may lead to moving of the scheduling code out of one or more nested loops which may reduce the amount of communication required significantly. We are incrementally incorporating many more optimizations in the compiler.

8 Experimental Results

A prototype compiler is complete (it was demonstrated at Supercomputing'92). In this section, we describe our experience in using the compiler.

8.1 Portability of the Fortran 90D/HPF Compiler

One of the principal requirements of the users of distributed memory MIMD systems is some “guarantee” of the portability for their code. Express parallel programming environment [27] guarantees this the portability on various platforms including, Intel iPSC/860, nCUBE/2, networks of workstations etc. We should emphasize that we have implemented a collective communication library which is currently built on the top of Express message passing primitives. Hence, in order to change to any other message passing system such as PVM [28] (which also runs on several platforms), we only need to replace the calls to the communication primitives in our communication library (not the compiler). However, it should be noted that a penalty must be paid to achieve portability because portable routines are normally built on top of the system routines. Therefore, the performance also depends on how efficient are the communication primitives on the top of which the communication library is built.

As a test application we use Gaussian Elimination, which is a part of the FortranD/HPF benchmark test suite [29]. Figure 5 shows the execution times obtained to run the same compiler generated code on a 16-node Intel/860 and nCUBE/2 for various problem sizes. Due to space limitations, we do not present performance of many other programs, and some of them can be found in [30].

8.2 Performance Evaluation

Table 4 shows a comparison between the performance of the hand-written Fortran 77+MP code with that of the compiler generated code. We can observe that the performance of the compiler generated code is within 10% of the hand-written code. This is due to the fact that the compiler generated code produces an extra communication call that can be eliminated using optimizations. However as Figure 6 shows, the gap between the performance of the two codes increases as the number of processors increases. This is because the extra communication step is a broadcast which is almost $O(\log(P))$ for a P processor hypercube system.

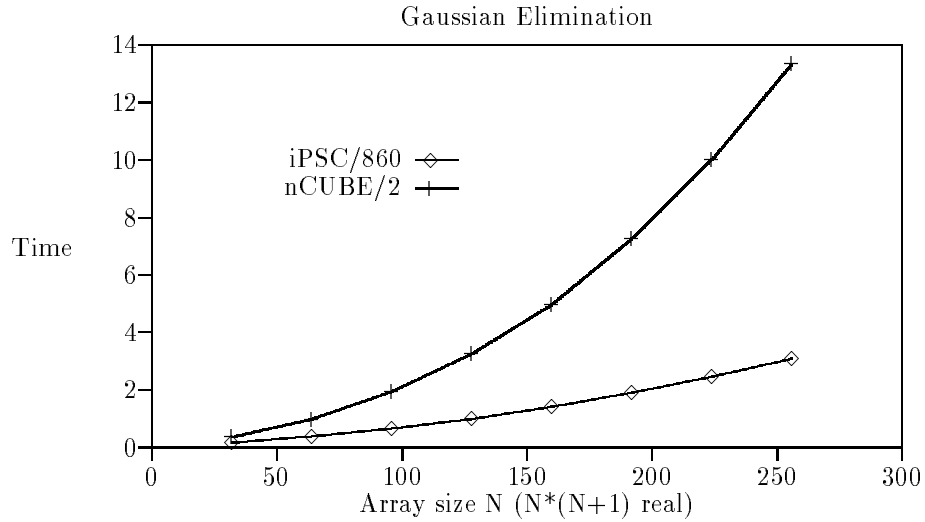


Figure 5: Execution time of Fortran 90D compiler generated code for Gaussian Elimination on a 16-node Intel iPSC/860 and nCUBE/2 (time in seconds).

Table 4: Comparison of the execution times of the hand-written code and Fortran 90D compiler generated code for Gaussian Elimination. Matrix size is 1023x1024 and it is column distributed.(Intel iPSC/860, time in seconds).

	Number of PEs				
	1	2	4	8	16
Hand Written	623.16	446.60	235.37	134.89	79.48
Fortran 90D	618.79	451.93	261.87	147.25	87.44

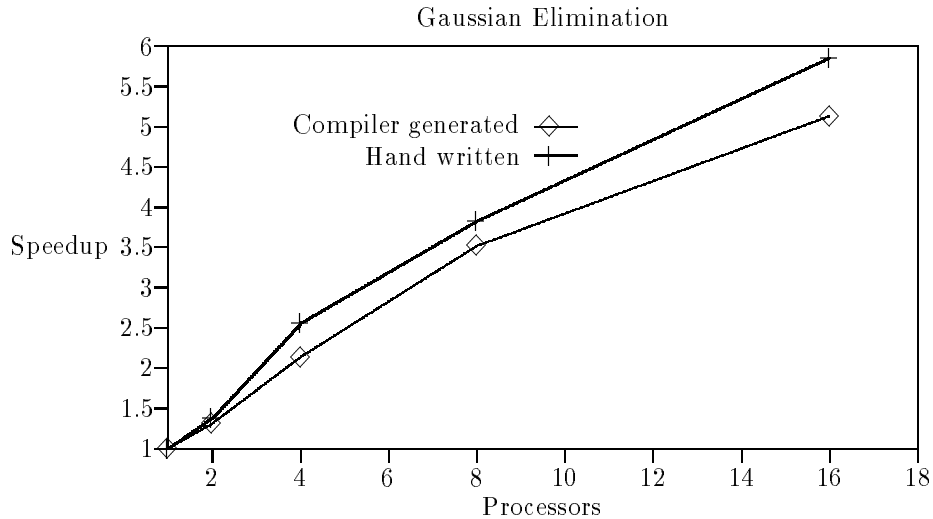


Figure 6: **Speed-up against the sequential code (corresponds to Table 4 of the hand-written code and Fortran 90D compiler generated code for Gaussian Elimination).**

9 Summary of Related Work

The compilation technique of Fortran 77 for distributed memory systems has been addressed by Callahan and Kennedy [10]. Currently, a Fortran 77D compiler is being developed at Rice [25, 31]. Superb [9] compiles a Fortran 77 program into a semantically equivalent parallel SUPRENUM multiprocessor. Koelbel and Mehrotra [22, 17] present a compilation method where a great deal of effort is put on run-time analysis for optimizing message passing in implementation of Kali. Quinn *et al.* [32, 33] use a data parallel approach for compiling C* for hypercube machines. The ADAPT system [34] compiles Fortran 90 for execution on MIMD distributed memory architectures. The ADAPTOR [35] is a tool that transform data parallel programs written in Fortran with array extension and layout directives to explicit message passing. Chen [19, 36] describes general compiler optimization techniques that reduce communication overhead for Fortran-90 implementation on massively parallel machines. Many techniques especially for unstructured communication of Fortran 90D compiler are adapted from Saltz *et al.* [37, 26, 18]. Gupta *et al.* [20, 38] use collective communication on automatic data partitioning on distributed memory machines. Due to space limitations, we do not elaborate on various other related projects.

10 Conclusions

In this paper, we presented design, implementation and performance results of our Fortran 90D/HPF compiler for distributed memory machines. Specifically, techniques for processing distribution directives, computation partitioning, communication detection and generation were presented. We also showed that our design is portable, yet efficient.

We believe that the methodology presented in this paper to compile Fortran 90D/HPF can be used by the designers and implementors for HPF language.

Acknowledgments

We are grateful to Parasoft for providing the Fortran 90 parser and Express without which the prototype compiler could have been delayed. We would like to thank the other members of our compiler research group I. Ahmad, R. Bordawekar, R. Ponnusamy, R. Thakur, and J. C. Wang for their contribution in the project including the development of the run-time library functions, testing, and help with programming. We would also like to thank K. Kennedy, C. Koelbel, C. Tseng and S. Hiranandani of Rice University for many inspiring discussions and inputs that have greatly influenced this work.

References

- [1] American National Standards Institute. Fortran 90: X3j3 internal document s8.118. *Submitted as Text for ISO/IEC 1539:1991*, May 1991.
- [2] G. C. Fox, S. Hiranadani, K. Kennedy, C. Koelbel, U. Kremer, C. Tseng, and M. Wu. Fortran D Language Specification. Technical report, Rice and Syracuse University, 1992.
- [3] High Performance Fortran Forum. High performance fortran language specification version 1.0. *Draft, Also available as technical report CRPC-TR92225 from the Center for Research on Parallel Computation, Rice University.*, Jan. 1993.
- [4] The Thinking Machine Corporation. *CM Fortran User's Guide version 0.7-f*, July 1990.
- [5] Maspar Computer Corporation. *MasPar Fortran User Guide version 1.1*, Aug. 1991.
- [6] G. Fox. The architecture of problems and portable parallel software systems. Technical Report SCCS-78b, Syracuse University, 1991.
- [7] Z. Bozkus et al. Compiling the FORALL statement on MIMD parallel computers. Technical Report SCCS-389, Northeast Parallel Architectures Center, July 1992.

- [8] Z. Bozkus et al. Compiling Distribution Directives in a Fortran 90D Compiler. Technical Report SCCS-388, Northeast Parallel Architectures Center, July 1992.
- [9] H. Zima, H. Bast, and M. Gerndt. Superb: A tool for semi Automatic SIMD/MIMD Parallelization. *Parallel Computing*, January 1988.
- [10] D. Callahan and K. Kennedy. Compiling programs for Distributed Memory Multiprocessors. *The Journal of Supercomputing*, pages 171–207, 1988.
- [11] M. Wu and G. Fox et al. Compiling Fortran 90 programs for distributed memory MIMD parallel computers. Technical Report SCCS-88, Northeast Parallel Architectures Center, May 1991.
- [12] G. C. Fox, M.A. Johnson, G.A. Lyzenga, S. W. Otto, J.K. Salmon, and D. W. Walker. In *Solving Problems on Concurrent Processors*, volume 1-2. Prentice Hall, May 1988.
- [13] K. Knobe, J. D. Lukas, and G. L. Steele. Data optimization: Allocation of arrays to reduce communication on SIMD machines. *Journal of Parallel and Distributed Computing*, pages 102–118, Feb 1990.
- [14] J. Li and M. Chen. The data alignment phase in compiling programs for distributed-memory machines. *Journal of Parallel and Distributed Computing*, pages 213–221, Oct 1991.
- [15] S. Chatterjee, J.R. Gilbert, R. Schreiber, and S.H Tseng. Automatic Array Alignment in Data-Parallel Programs. *Twentieth Annual ACM SIGACT/SIGPLAN Symposium on Principles of Programming Languages*, January 1993.
- [16] M. Gerndt. Updating distributed variables in local computations. *Concurrency: Practice and Experience*, September 1990.
- [17] C. Koelbel and P. Mehrotra. Supporting Compiling Global Name-Space Parallel Loops for Distributed Execution. *IEEE Transactions on Parallel and Distributed Systems*, October 1991.
- [18] H. Berryman J. Saltz, J. Wu and S. Hiranandani. Distributed Memory Compiler Design for Sparse Problems. *Interim Report ICASE, NASA Langley Research Center*, 1991.
- [19] J. Li and M. Chen. Compiling Communication -Efficient Programs for Massively Parallel Machines. *IEEE Transactions on Parallel and Distributed Systems*, pages 361–376, July 1991.
- [20] M. Gupta and P. Banerjee. Demonstration of Automatic Data Partitioning Techniques for Parallelizing Compilers on Multicomputers. *IEEE: Transaction on Parallel and Distributed Systems*, pages 179–193, March 1992.
- [21] R. Das, J. Saltz, and H. Berryman. A Manual For PARTI Runtime Primitives. *NASA, ICASE Interim Report 17*, May 1991.
- [22] C. Koelbel, P. Mehrotra, and J. V. Rosendale. Supporting Shared Data Structures on Distributed Memory Architectures. *PPoPP*, March 1990.
- [23] A.V. Aho, R. Sethi, and J.D Ullman. *Compilers Principles, Techniques and Tools*. March 1988.
- [24] I. Ahmad, R. Bordawekar, Z. Bozkus, A. Choudhary, G. Fox, K. Parasuram, R. Ponnusamy, S. Ranka, and R. Thakur. Fortran 90D Intrinsic Functions on Distributed Memory Machines: Implementation and Scalability. Technical Report SCCS-256, Northeast Parallel Architectures Center, March 1992.

- [25] S. Hiranandani, K. Kennedy, and C. Tseng. Compiler optimization for Fortran D on MIMD distributed-memory machines. *Proc. Supercomputing'91*, Nov 1991.
- [26] R. Mirchandaney J. Saltz, K. Crowley and H. Berryman. Run-time scheduling and execution of loops on message passing machines. *Journal of Parallel and Distributed Computing*, December 1991.
- [27] ParaSoft Corp. *Express Fortran reference guide Version 3.0*, 1990.
- [28] A. Beguelin, J. Dongarra, A. Geist, R. Mancheck, and V. Sunderam. A Users Guide to PVM Parallel Virtual Machine. Technical Report ORNL/TM-11826, Oak Ridge National Laboratory, July 1991.
- [29] A. G. Mohamed, G. Fox, G. V. Laszewski, M. Parashar, T. Haupt, K. Mills, Y. Lu, N. Lin, and N. Yeh. Application Benchmark Set for Fortran-D and High Performance Fortran. Technical Report SCCS-327, Northeast Parallel Architectures Center, May 1992.
- [30] Z. Bozkus et al. Compiling Fortran 90D/HPF for Distributed Memory MIMD Computers. Technical Report SCCS-444, Northeast Parallel Architectures Center, 1993.
- [31] S. Hiranandani, K. Kennedy, and C.W. Tseng. Compiler support for machine-indepentet Parallel Programming in Fortran D. *Compiler and Runtime Software for Scalable Multiprocessors*, 1991.
- [32] Michael Quinn, Philip Hatcher, and Karen Jourdenais. Compiling C* Programs for a Hypercube Multicomputer. *Parallel Computing Laboratory, University of New Hampshire*, PCL-87-12, December 1987.
- [33] Philip Hatcher, Anthony Lapadula, Robert Jones, Michael Quinn, and Ray Anderson. A Production-Quality C* Compiler for Hypercube Multicomputers. *Third ACM SIGPLAN symposium on PPOPP*, 26:73–82, July 1991.
- [34] J.H Merlin. Techniques for the Automatic Parallelisation of 'Distributed Fortran 90'. Technical Report SNARC 92-02, Southampton Novel Architecture Research Centre, 1992.
- [35] T. Brandes. ADAPTOR Language Reference Manual. Technical Report ADAPTOR-3, German National Research Center for Computer Science, 1992.
- [36] M. Chen and J.J Wu. Optimizing FORTRAN-90 Programs for Data Motion on Massively Parallel Systems. Technical Report YALEU/DCS/TR-882, Yale University, Dep. of Comp. Sci., 1992.
- [37] H. Berryman J. Saltz and J. Wu. Multiprocessors and run-time compilation. *Concurrency: Practice and Experience*, December 1991.
- [38] M. Gupta. Automatic Data Partitioning on Distributed Memory Multicomputers. Technical Report PhD thesis, University of illinois at Urbana-Champaign, 1992.