

# Unified Compilation of Fortran 77D and 90D

ALOK CHOUDHARY,\* GEOFFREY FOX,\* SEEMA HIRANANDANI,†  
KEN KENNEDY,† CHARLES KOELBEL,† SANJAY RANKA,\* and  
CHAU-WEN TSENG†

---

We present a unified approach to compiling Fortran 77D and Fortran 90D programs for efficient execution on MIMD distributed-memory machines. The integrated Fortran D compiler relies on two key observations. First, array constructs may be *scalarized* into FORALL loops without loss of information. Second, *loop fusion*, *partitioning*, and *sectioning* optimizations are essential for both Fortran D dialects.

Categories and Subject Descriptors: C.1.2 [**Processor Architectures**]: Multiple Data Stream Architectures (Multiprocessors)—*multiple-instruction-stream, multiple-data-stream processors (MIMD), parallel processors*; D.3.2 [**Programming Languages**]: Language Constructs and Features—*concurrent programming structures*; D.3.4 [**Programming Languages**]: Processors—*code generation; compilers; optimization; preprocessors*

General Terms: Languages, Performance

Additional Key Words and Phrases: Fortran D, parallel languages, parallel programming

---

## 1. INTRODUCTION

Parallel computing on distributed-memory machines is scalable and cost-effective; however, it is hindered by both the difficulty of parallel programming and lack of portability of the resulting programs. We propose to solve this problem using Fortran D, a version of Fortran extended with data decomposition specifications that can be applied to Fortran 77 and Fortran 90 [ANSI 1990] to produce Fortran 77D and Fortran 90D, respectively. Fortran D has contributed to the development of High Performance Fortran (HPF), an informal Fortran standard for programming massively parallel machines [High Performance Fortran Forum 1993].

---

This research was supported by the Center for Research on Parallel Computation (CRPC), a National Science Foundation Science and Technology Center. CRPC is funded by NSF through Cooperative Agreement Number CCR-9120008. Additional support was provided by DARPA under contract DABT63-91-C-0028. Alok Choudhary is also supported by an NSF Young Investigator Award CCR-9357840.

Authors' addresses: \*Northeast Parallel Architectures Center, 111 College Place, Syracuse University, Syracuse, NY 13244-4100.

†Center for Research on Parallel Computation, Rice University, P. O. Box 1892, Houston, TX 77251-1892.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1993 ACM 1057-4514/93/0300-0095\$03.50

ACM Letters on Programming Languages and Systems,

Vol. 2, Nos. 1-4, March-December 1993, Pages 95-114.

We are developing the compiler technology needed to automate translation of Fortran D to different parallel architectures. Our goal is to establish a machine-independent programming model for data-parallel programs that is easy to use, yet performs with acceptable efficiency on different parallel architectures. In this paper, we describe a unified strategy for compiling both Fortran 77D and Fortran 90D into efficient *single-program, multiple-data* (SPMD) message-passing programs. In particular, we concentrate on the design of a prototype Fortran 90D compiler for the Intel iPSC/860 and Thinking Machines CM-5, two *multiple-instruction, multiple-data* (MIMD) distributed-memory machines.

The principal issues involved in compiling Fortran 90D are *partitioning* the program across multiple nodes and *scalarizing* it for execution on each individual node. Previous work has described both the scalarization [Allen and Kennedy 1992] and partitioning process [Hall et al. 1992; Hiranandani et al. 1991; 1992a; 1992b]. The contributions of this paper are to:

- utilize the FORALL loop to preserve the semantics of Fortran 90 array constructs in a common intermediate form;
- provide a compilation framework for integrating program partitioning with Fortran 90 scalarization;
- show that an efficient, portable run-time library can ease the task of compiling Fortran D;
- validate the effectiveness of a unified compilation approach through empirical case studies.

Because our main interest is in the interaction between scalarization and partitioning optimizations, in this paper we only focus on a subset of Fortran 90, namely: array operations, array constructs, and the FORALL loop.

The remainder of this paper presents a brief overview of the Fortran D language and compilation strategy, then describes the Fortran 90D and 77D front ends and the common Fortran D back end. The design of the run-time library is discussed, and an example is used to illustrate the compilation process. We conclude with a discussion of related work.

## 2. FORTRAN D LANGUAGE

We briefly overview aspects of Fortran D relevant to this paper. These extensions can be added to either Fortran 77 or Fortran 90. The complete language is described elsewhere [Fox et al. 1990].

### 2.1 Data Alignment and Distribution

In Fortran D, the DECOMPOSITION statement declares an abstract problem or index domain. The ALIGN statement maps each array element onto one or more elements of the decomposition. This provides the minimal requirement for reducing data movement for the program given an unlimited number of processors. The DISTRIBUTE statement groups decomposition elements, mapping them and any array elements aligned with them to the finite resources of the physical machine. Each dimension of the decomposition is distributed

in a block, cyclic, or block-cyclic manner; the symbol “:” marks dimensions that are not distributed. Because the alignment and distribution statements are executable, dynamic data decomposition is possible.

## 2.2 FORALL

Fortran D provides FORALL loops to permit the user to specify difficult parallel loops in a deterministic manner. Its semantics are borrowed from the Myrias PARDO and previous FORALL constructs [Albert et al. 1991; Beltrametti et al. 1988; Lundstrom and Barnes 1980]. In a FORALL loop, each iteration *uses only values defined before the loop or within the current iteration*. When a statement in an iteration of the FORALL loop accesses a memory location, it will not get any value written by a different iteration of the loop. Instead, it will get the *old* value at that memory location (*i.e.*, the value at that location before the execution of the FORALL loop), or it will get some new value written on the current iteration. Similarly, a merging semantics ensures that a deterministic value is obtained after the FORALL if several iterations assign to the same memory location.

Another way of viewing the FORALL loop is that it has copy-in/copy-out semantics. In other words, each iteration gets its own copy of the entire data space that exists before the execution of the loop, and writes its results to a new data space at the end of the loop. Since no values depend on other iterations, the FORALL loop may be executed in parallel without synchronization. However, communication may still be required before the loop to acquire nonlocal values, and after the loop to update or merge nonlocal values. Single-statement Fortran D FORALL loops are identical to those supported in CM FORTRAN [Albert et al. 1988] and HPF [High Performance Fortran Forum 1993], but multistatement Fortran D FORALL loops are different from those found in HPF.

## 3. FORTRAN D COMPILATION STRATEGY

### 3.1 Overall Strategy

Our strategy for parallelizing Fortran D programs for distributed-memory MIMD computers is illustrated in Figure 1. In brief, we transform both Fortran 77D and Fortran 90D to a common intermediate form, which is then compiled to code for the individual nodes of the machine. We have several pragmatic and philosophical reasons for this strategy:

- Sharing a common back end for both the Fortran 77D and Fortran 90D avoids duplication of effort.
- Decoupling the Fortran 77D and Fortran 90D front ends allows them to become machine independent.
- Providing a common intermediate form helps us experiment with defining an efficient compiler/programmer interface for programming the nodes of a massively parallel machine.

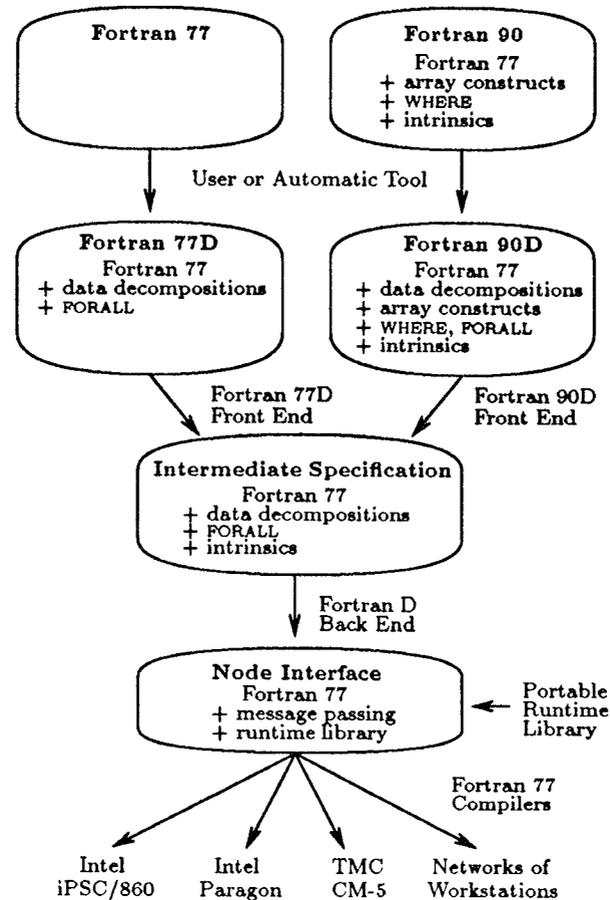


Fig. 1. Fortran D compilation strategy.

### 3.2 Intermediate Form

To compile both dialects of Fortran D using a single back end, we must select an appropriate intermediate form. In addition to standard computation and control flow information, the intermediate form must capture three important aspects of the program:

- Data decomposition information, telling how data is aligned and distributed among processors.
- Parallelization information, telling when operations in the code are independent.
- Communication information, telling what data must be transferred between processors.

In addition, we believe that the primitive operations of the intermediate form should be relatively low-level operations that can be translated simply for single-processor execution.

We have chosen Fortran 77 with data decompositions, FORALL, and intrinsic functions to be the intermediate form for the Fortran D compiler. We show later that this form preserves all of the information available in a Fortran 90 program, but maintains the flexibility of Fortran 77. Parallelism and communication can be determined by the compiler for simple computations, and specified by the user using FORALL and intrinsic functions for complex computations.

### 3.3 Node Interface

Another topic of interest in the overall strategy is the node interface: the node program produced by the Fortran D compiler. It must be both portable and efficient. In addition, the level of the node interface should be neither so high that efficient translation to object code is impossible, nor so low that its workings are completely opaque to the user. We have selected Fortran 77 with calls to communication and run-time libraries based on Express, a collection of portable message-passing primitives [Parasoft Corp. 1989]. Evaluating our experiences with this node interface is the first step toward defining an “optimal” level of support for programming individual nodes of a parallel machine.

## 4. FORTRAN D COMPILER

The Fortran D compiler thus consists of three parts. The Fortran 90D and 77D front ends process input programs into the common intermediate form. The Fortran D back end then compiles this to the SPMD message-passing node program. The Fortran D compiler is implemented in the context of the ParaScope programming environment [Cooper 1993].

### 4.1 Fortran 90D Front End

The function of the Fortran 90D front end is to *scalarize* the Fortran 90D program, translating it to an equivalent Fortran 77D program. This is necessary because the underlying machine executes computations sequentially, rather than on entire arrays at once as specified in Fortran 90. For the Fortran D compiler we find it useful to view scalarization as three separate tasks:

- Scalarizing Fortran 90 constructs.* Many Fortran 90 features are not present in our intermediate form. They must be translated into equivalent Fortran 77D statements.
- Fusing loops.* Simple scalarization results in many small loop nests. Fusing these loop nests can improve the locality of data accesses, simplify partitioning, and enable other program transformations.
- Sectioning.* Fortran 90 array operations allow the programmer to access and modify entire arrays atomically, even if the underlying machine lacks this capability. The Fortran D compiler must divide array operations into *sections* that fit the hardware of the target machine [Allen and Kennedy 1992].

We defer both loop fusion and sectioning to the Fortran D back end. Loop fusion is deferred because even hand-written Fortran 77 programs can benefit significantly [Carr et al. 1992; Kennedy and McKinley 1992; McKinley 1992]. Sectioning is needed in the back end because `FORALL` loops may also be present in Fortran 77D.

We assign to the Fortran 90D front end the remaining task, scalarizing Fortran 90 constructs that have no equivalent in the Fortran 77D intermediate form. There are three principal Fortran 90 language features that must be scalarized: array constructs, `WHERE` statements, and intrinsic functions [ANSI 1990].

*Array constructs.* Fortran 90 *array constructs* allow entire arrays to be manipulated atomically, enhancing the clarity and conciseness of the program. Previous research has shown that efficiently implementing array constructs for scalar machines may be difficult [Allen and Kennedy 1992]. One problem arises when Fortran 90 array constructs are used in assignment statements, the entire right-hand side (*rhs*) must be evaluated before storing the results in the left-hand side (*lhs*). Without adequate analysis, *rhs* array elements would need to be stored in temporary buffers to ensure that they are not overwritten before their values are used. The Fortran 90 front end can defer this problem by relying on a key observation: the `FORALL` loop possesses copy-in/copy-out semantics identical to Fortran 90 assignment statements utilizing array constructs. Such statements may thus be translated into equivalent `FORALL` loops with no loss of information.

*Where statement.* Another Fortran 90 feature is the `WHERE` statement. It takes a boolean argument that is used to *mask* array operations, inhibiting assignments to array elements whose matching boolean flag has the value *false*. Fortunately, the `WHERE` statement may be easily translated into equivalent `IF` and `FORALL` statements. Consider the following example where `A` is assumed to be a 1D  $N$ -element array. Because of `FORALL` copy-in/copy-out semantics, it is unnecessary at this point to explicitly store the value of the boolean argument to prevent it from being overwritten.

```

WHERE (A .EQ. 0)  FORALL i = 1, N
  A = 1.0          IF (A(i) .EQ. 0) THEN
ELSEWHERE        =>  A(i) = 1.0
  A = 0.0          ELSE
ENDWHERE         A(i) = 0.0
                  ENDIF
                  ENDFOR

```

*Intrinsic functions.* Intrinsic functions are fundamental to Fortran 90. They 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 matrix multiplication. Additional intrinsics are described in Table I. To avoid excessive complexity and machine-dependence in the Fortran D compiler, we convert most Fortran 90 intrinsics into calls to customized run-time library functions.

Table I. Representative Intrinsic Functions of Fortran 90D

Data Movement	Reductions	Irregular Operations	Special Routines
CSHIFT EOSHIFT SPREAD RESHAPE TRANSPOSE	DOTPRODUCT ALL, ANY, COUNT MAXVAL, MINVAL SUM, PRODUCT MAXLOC, MINLOC	PACK / UNPACK	MATMUL

The strategy used by the Fortran 90D front end is thus to preserve all intrinsic functions, passing them to the Fortran D compiler back end. However, some processing is necessary. Like the WHERE statement, some intrinsic functions accept a mask expression that restricts execution of the computation. The Fortran 90D front end may need to evaluate the expression and store it in a temporary boolean array before performing the computation, so the mask can be passed as an argument to the run-time library.

For example, consider the following reduction operation, where X is a scalar and A, B are arrays:

```
X = MAXVAL(A, A .EQ. B)
```

It should return the value of the element of A that is the maximum of all elements for which element of A is equal to the corresponding element of B. The Fortran 90D front end translates this to:

```
FORALL i = 1, N
  TMP(i) = A(i) .EQ. B(i)
ENDFOR
X = MAXVAL(A, TMP)
```

TMP can then be passed as an argument to the run-time routine MAXVAL. Temporary arrays may also be introduced when intrinsic functions return a value that is part of a Fortran 90 expression.

*Temporary arrays.* When the Fortran 90D front end needs to create temporary arrays, it must also generate appropriate Fortran D data decomposition statements. A temporary array is usually aligned and distributed in the same manner as its master array. For example, in the previous example the temporary logical array TMP is aligned and distributed in the same manner as A and B. If A and B are distributed differently, automatic data decomposition techniques must be applied to select efficient data distributions for compiler-generated temporary arrays [Chatterjee et al. 1993; Kennedy and Kremer 1991; Knobe et al. 1990].

#### 4.2 Fortran 77D Front End

The Fortran 77D front end does not need to perform much work since Fortran 77D is very close to the intermediate form. Its only task is to detect complex high-level parallel computations, replacing or annotating them by their equivalent Fortran 90 intrinsics. These intrinsic functions help the compiler recognize complex computations such as reductions and scans that are supported by the run-time library. With advanced program analysis, some opera-

tions such as DOTPRODUCT, SUM, TRANSPOSE, or MATMUL can be detected automatically with ease. Others computations such as COUNT or PACK may require user assistance.

### 4.3 Fortran D Back End

The Fortran D back end performs two main functions—it partitions the program onto the nodes of the parallel machine and completes the scalarization of Fortran D into Fortran 77. We find that the desired order for compilation phases is to apply loop fusion first, followed by partitioning and sectioning.

Loop fusion is performed first because it simplifies partitioning by reducing the need to consider interloop interactions. It also enables optimizations such as *strip-mining* and *loop interchange* [Allen and Kennedy 1987; Wolfe 1989]. In addition, loop fusion does not increase the difficulty of later compiler phases. On the other hand, sectioning is performed last because it can significantly disrupt the existing program structure, increasing the difficulty of partitioning analysis and optimization.

**4.3.1 Loop Fusion.** Loop fusion is particularly important for the Fortran D back end because scalarized Fortran 90 programs present many single-statement loop nests. Fusing such loops simplifies the partitioning process and enables additional optimizations.

Data dependence is a concept developed for vectorizing and parallelizing compilers to characterize memory access patterns at compile time [Allen and Kennedy 1987; Kuck et al. 1981; Wolfe 1989]. A true dependence indicates definition followed by use, while an antidependence shows use before definition. Data dependences may be either loop-carried or loop-independent. Loop fusion is legal if it does not reverse the direction of any data dependence between two loop nests [Allen and Kennedy 1992; Wolfe 1989].

The current Fortran D back end fuses all adjacent loop nests where legal, if no loop-carried true dependences are introduced. This heuristic does not adversely affect the parallelism or communication overhead of the resulting program, and should perform well for the simple cases found in practice. More sophisticated algorithms are discussed elsewhere [McKinley 1992].

Loop fusion also has the added advantage of being able to improve memory reuse in the resulting program. Modern high-performance processors are so fast that memory latency and bandwidth limitations become the performance bottlenecks for most scientific programs. Transformations such as loop fusion promote memory reuse and can significantly improve program efficiency for both scalar and vector machines [Allen and Kennedy 1992; Kuck et al. 1981; McKinley 1992; Sarkar and Gao 1991]. For instance, consider the following example.

```

FORALL i = 1, N      FORALL i = 1, N
  A(i) = i           A(i) = i
ENDFOR              ⇒ B(i) = A(i)*A(i)
FORALL i = 1, N      ENDFOR
  B(i) = A(i)*A(i)
ENDFOR

```

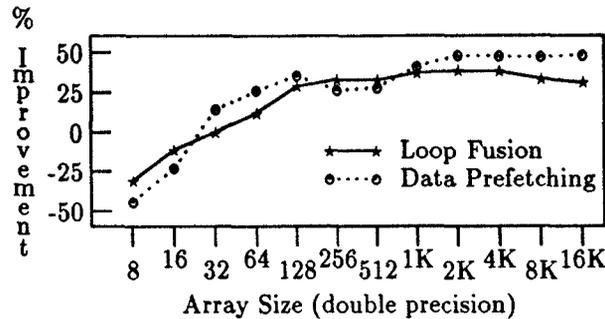


Fig. 2. Effect of scalarization optimizations.

The occurrences of  $A(i)$  in separate loops mean that the memory location referenced by  $A(i)$  in the first loop is likely to have been flushed from the cache by the reference in the second loop. If the two loops are fused, all accesses to  $A(i)$  occur in the same loop iteration, allowing the value to be reused in a register or cache. For this example, we measured improvements of up to 30% for some problem sizes on an Intel i860, as shown in Figure 2. Additional transformations to enhance memory reuse and increase unit-stride memory accesses are also quite important; they are described elsewhere [Kennedy and McKinley 1992; McKinley 1992].

**4.3.2 Program Partitioning.** The major step in compiling Fortran D for MIMD distributed-memory machines is to partition the data and computation across processors, introducing communication where needed. We present a brief overview of the Fortran D compilation process below; details are discussed elsewhere [Hall et al. 1992; Hiranandani et al. 1992a; 1992b; Tseng 1993].

- Analyze program.* Symbolic and data dependence analysis is performed.
- Partition data.* Fortran D data decomposition specifications are analyzed to determine the decomposition of each array in a program.
- Partition computation.* The compiler partitions computation across processors using the “owner computes” rule: where each processor only computes values of data it owns [Callahan and Kennedy 1989; Rogers and Pingali 1989; Zima et al. 1988].
- Analyze communication.* Based on the work partition, references that result in nonlocal accesses are marked.
- Optimize communication.* Nonlocal references are examined to determine optimization opportunities. The key optimization, message vectorization, uses the level of loop-carried true dependences to combine element messages into vectors [Balasundaram et al. 1990; Zima et al. 1988].
- Manage storage.* “Overlaps” [Zima et al. 1988] or buffers are allocated to store nonlocal data.

—*Generate code.* Information gathered previously is used to generate the SPMD program with explicit message passing that executes directly on the nodes of the distributed-memory machine.

Two extensions are needed in the Fortran D back end to handle FORALL loops and intrinsics. During communication optimization, the Fortran D compiler treats all true dependences carried by FORALL loops as antidependences. This reflects the semantics of the FORALL loop and ensures that the message vectorization algorithm will place all communication outside the loop. In addition, during code generation intrinsic functions are translated into calls to the run-time library. Parameters are added where necessary to provide necessary data-partitioning information.

**4.3.3 Sectioning.** The final phase of the Fortran D back end completes the scalarization process. After partitioning is performed, the compiler applies *sectioning* to convert FORALL loops into DO loops [Allen and Kennedy 1992] in the node program. The Fortran D back end detects cases where temporary storage may be needed using data dependence analysis. True dependences carried on the FORALL loop represent instances where values are defined in the loop and used on later iterations; they point out where the copy-in/copy-out semantics of the FORALL loop is being violated.

During simple translation of Fortran 90 array constructs or FORALL loops, arrays involved in loop-carried true dependences must be saved in temporary buffers to preserve their old values. For instance, consider the translation of the following concise Fortran 90 formulation of the Jacobi algorithm:

```

A(2:N - 1) = 0.5 * (A(1:N - 2) + A(3:N))
  ↓
FORALL i = 2, N - 1
  A(i) = 0.5 * (A(i - 1) + A(i + 1))
ENDFOR
  ↓
DO i = 1, N - 2
  TMP(i) = A(i - 1)
ENDDO
DO i = 2, N - 1
  A(i) = 0.5 * (TMP(i) + A(i + 1))
ENDFOR

```

A loop-carried true dependence exists between the definition to  $A(i)$  and the use of  $A(i - 1)$ . A temporary array TMP is needed so that the old values of  $A(i - 1)$  are not overwritten before they are used. The values of  $A(i + 1)$  do not need to be buffered since they are used before being redefined.

The previous example is problematic because temporary storage is required for the values of  $A(i - 1)$ . In some cases, the Fortran D compiler can eliminate buffering through program transformations such as *loop reversal*. In other cases, the compiler can reduce the amount of temporary storage required through *data prefetching* [Allen and Kennedy 1992; Bromley et al. 1991]. For instance, in the Jacobi example a more efficient translation would

Table II. Performance of Some Fortran 90 Intrinsic Functions

Processor	Time (milliseconds)						
	ALL	ANY	MAXVAL	PRODUCT	TRANSPOSE		
	1K × 1K	1K × 1K	1K × 1K	256K	256 × 256	512 × 512	1K × 1K
1	580.6	606.2	658.8	90.1	58	299	-
2	291.0	303.7	330.4	50.0	118	575	-
4	146.2	152.6	166.1	25.1	87	395	-
8	73.84	77.1	84.1	13.1	61	224	1039
16	37.9	39.4	43.4	7.2	41	140	539
32	19.9	20.7	23.2	4.2	36	85	316

result in:

```

X = A(1)
DO i = 2, N - 1
  Y = 0.5 * (X + A(i + 1))
  X = A(i)
  A(i) = Y
ENDFOR

```

This reduces the temporary memory required significantly, from an entire array to two scalars. For this version of Jacobi, we measured improvements of up to 50% for certain problem sizes on an Intel i860, as shown in Figure 2.

## 5. RUN-TIME LIBRARY

Fortran 90 intrinsic functions represent computations (such as TRANSPOSE and MATMUL) that may have complex communication patterns. It is possible to support these functions at compile time, but we have chosen to implement these functions in the run-time library instead to reduce the complexity and machine-dependence of the compiler. The Fortran D compiler translates intrinsics into calls to run-time library routines using a standard interface. Additional information is passed describing bounds, overlaps, and partitioning for each array dimension. The run-time library developed at Syracuse is built on top of the Express communication package to ensure portability across different architectures [Parasoft Corp. 1989].

Table II presents some sample performance numbers for a subset of the intrinsic functions on an iPSC/860; details are presented elsewhere [Ahmad et al. 1992]. The times in the table include both the computation and communication times for each function. These measurements are also displayed in Figure 3. Timings in seconds are plotted logarithmically along the Y-axis. The number of processors is plotted along the X-axis. Different lines in the graph correspond to timings of individual functions in the run-time library for different problem sizes. For large problem sizes, we were able to obtain almost linear speedups. In the case of the TRANSPOSE function, going from one processor to two or four degrades execution time due to increased communication. However, speedup improves as the number of processors increases.

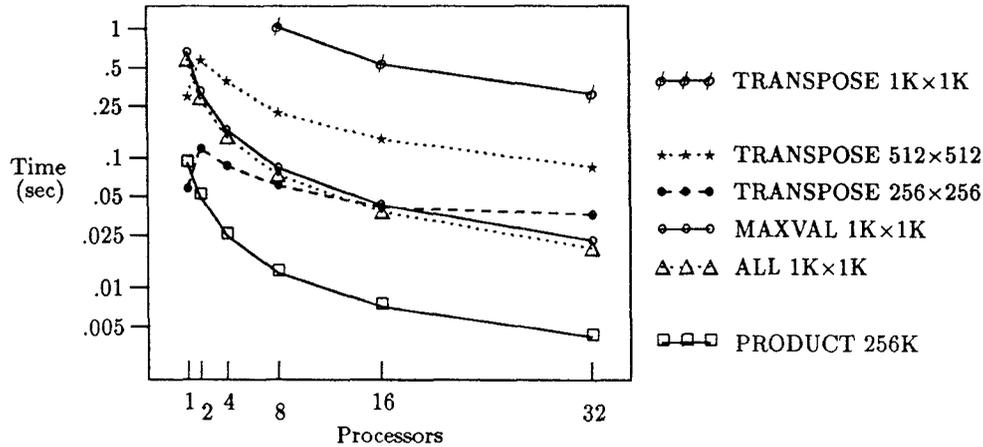


Fig. 3. Performance of run-time library.

## 6. FORTRAN 90D COMPILATION EXAMPLE

### 6.1 Compilation

Figure 4 shows a code fragment implementing one sweep of ADI integration on a 2D mesh, a typical (if short) numerical algorithm. Conceptually, the code is solving a tridiagonal system (represented by the arrays  $A$  and  $B$ ) along each row of the matrix  $X$ . The tridiagonal systems are solved by a sequential method, but separate columns are independent and may be solved in parallel. The full version of ADI integration sweeps each dimension of the mesh, preventing completely parallel execution for any static data decomposition.

In the example, Fortran D data decomposition statements are used to partition the 2D array into blocks of columns. For clarity, we declare the number of processors (NSPROC) to be 32 at compile time. The Fortran 90D example is concise and convenient for the user, since it can be written for a single address space without requiring explicit communication. However, additional compilation techniques are required to generate efficient code. First, the Fortran 90D front end translates the program into intermediate form as shown in Figure 5, converting all array constructs into FORALL loops. Since no true dependences are carried on the FORALL loops, they may be directly replaced with DO loops.

The compilation process for the Fortran D back end merits closer examination. First, array bounds are reduced to the local sections plus overlaps. The local processor number is determined using *myproc*( ), a library function; it is used to compute expressions for reducing loop bounds. Analysis determines that both  $I$  and  $J$  are *cross-processor* loops—loops carrying true dependences that sequentialize the computation across processors. To exploit pipeline parallelism, the Fortran D compiler interchanges such loops inward. We call this technique *fine-grain pipelining* [Hiranandani et al. 1992a; 1992b].

```

PARAMETER (N = 512, N$PROC = 32)
REAL X(N,N), A(N,N), B(N,N)
DECOMPOSITION DEC(N,N)
ALIGN X, A, B WITH DEC
DISTRIBUTE DEC(:,BLOCK)
DO I = 2,N
  X(1:N,I) = X(1:N,I) - X(1:N,I-1)*A(1:N,I)/B(1:N,I-1)
  B(1:N,I) = B(1:N,I) - A(1:N,I)*A(1:N,I)/B(1:N,I-1)
ENDDO
X(1:N,N) = X(1:N,N) / B(1:N,N)
DO J = N-1,1,-1
  X(1:N,J) = (X(1:N,J)-A(1:N,J+1)*X(1:N,J+1))/B(1:N,J)
ENDDO

```

Fig. 4. ADI integration in Fortran 90D.

```

PARAMETER (N = 512)
REAL X(N,N), A(N,N), B(N,N)
DO I = 2,N
  FORALL K = 1,N
    X(K,I) = X(K,I) - X(K,I-1)*A(K,I)/B(K,I-1)
  ENDFOR
  FORALL K = 1,N
    B(K,I) = B(K,I) - A(K,I)*A(K,I)/B(K,I-1)
  ENDFOR
ENDDO
FORALL K = 1,N
  X(K,N) = X(K,N)/B(K,N)
ENDFOR
DO J = N-1,1,-1
  FORALL K = 1,N
    X(K,J) = (X(K,J)-A(K,J+1)*X(K,J+1))/B(K,J)
  ENDFOR
ENDDO

```

Fig. 5. ADI in intermediate form.

For this version of ADI integration, data dependences permit the Fortran D compiler to interchange the  $J$  loop inward. However, if loop fusion is not performed, the imperfectly nested  $K$  loops inhibit loop interchange for loop  $I$ , forcing it to remain in place. During code generation, true dependences for nonlocal references carried on the  $I$  and  $J$  loop cause calls to *send* and *recv* to be inserted to provide communication and synchronization. Figure 6 shows the resulting program. Unfortunately, the computation in the  $I$  loop has been sequentialized, since each processor has to wait for its predecessor to complete. Note that this is not due to communication placement; the values needed by the succeeding processor are simply computed last.

If loop fusion is enabled, the Fortran D back end will fuse the two inner  $K$  loops. This is legal because the dependence between the definition and use of  $B$  is carried on the  $I$  loop and is thus unaffected. Fusion is also conservative because it does not introduce any true dependences carried by the  $K$  loop. Fusing the  $K$  loops promotes reuse of  $A$  and  $B$ , but its main benefit is to enable the Fortran D back end to interchange the  $I$  and  $K$  loops, exposing pipeline parallelism. The resulting program is displayed in Figure 7. For

```

REAL X(512,0:17), A(512,17), B(512,0:16)
MY$P = myproc()      { * 0...31 * }
LB1 = MAX((MY$P*16)+1,2) - MY$P*16
UB1 = MIN((MY$P+1)*16,511) - MY$P*16
IF (MY$P .GT. 0) recv(X(1:N,0),B(1:N,0),MY$P-1)
DO I = LB1, 16
  DO K = 1,N
    X(K,I) = X(K,I) - X(K,I-1)*A(K,I)/B(K,I-1)
  ENDDO
  DO K = 1,N
    B(K,I) = B(K,I) - A(K,I)*A(K,I)/B(K,I-1)
  ENDDO
ENDDO
IF (MY$P .LT. 31) send(X(1:N,16),B(1:N,16),MY$P+1)
IF (MY$P .EQ. 31) THEN
  DO K = 1,N
    X(K,16) = X(K,16)/B(K,16)
  ENDDO
ENDIF
IF (MY$P .GT. 0) send(A(1:N,1),MY$P-1)
IF (MY$P .LT. 31) recv(A(1:N,17),MY$P+1)
DO K = 1,N
  IF (MY$P .LT. 31) recv(X(K,17),MY$P+1)
  DO J = UB1,1,-1
    X(K,J) = (X(K,J)-A(K,J+1)*X(K,J+1))/B(K,J)
  ENDDO
  IF (MY$P .GT. 0) send(X(K,1),MY$P-1)
ENDDO

```

Fig. 6. ADI without loop fusion.

```

REAL X(512,0:17), A(512,17), B(512,0:16)
MY$P = myproc()      { * 0...31 * }
LB1 = MAX((MY$P*16)+1,2) - MY$P*16
DO K = 1,N
  IF (MY$P .GT. 0) recv(X(K,0),B(K,0),MY$P-1)
  DO I = LB1,16
    X(K,I) = X(K,I) - X(K,I-1)*A(K,I)/B(K,I-1)
    B(K,I) = B(K,I) - A(K,I)*A(K,I)/B(K,I-1)
  ENDDO
  IF (MY$P .LT. 31) send(X(K,16),B(K,16),MY$P+1)
ENDDO

```

Fig. 7. ADI with fine-grain pipelining (first loop only).

simplicity, only the first loop is shown. The remaining loops are compiled in a similar manner as before. Note that *loop distribution* applied to the  $I$  loop can enable loop interchange, but is prevented by the recurrence for  $B$ .

To reduce communication overhead, we can also apply strip-mining in conjunction with loop interchange to adjust the granularity of pipelining. We call this technique *coarse-grain pipelining* [Hiranandani et al. 1992a; 1992b]. In the ADI example, we strip-mine the  $K$  loop by four, then interchange the resulting loop outside the  $I$  loop. For this study, the desirable strip size was derived by hand; we describe elsewhere how it may be automatically calculated for a given architecture based on communication and computation costs [Hiranandani et al. 1992b]. After strip-mining, messages inserted outside the

```

REAL X(512,0:17), A(512,17), B(512,0:16)
MY$P = myproc()      { * 0...31 * }
LB1 = MAX((MY$P*16)+1,2) - MY$P*16
DO KK = 1,N,4
  IF (MY$P .GT. 0) THEN
    recv(X(KK:KK+3,0),B(KK:KK+3,0),MY$P-1)
  ENDIF
  DO I = LB1,16
    DO K = KK,KK+3
      X(K,I) = X(K,I) - X(K,I-1)*A(K,I)/B(K,I-1)
      B(K,I) = B(K,I) - A(K,I)*A(K,I)/B(K,I-1)
    ENDDO
  ENDDO
  IF (MY$P .LT. 31) THEN
    send(X(KK:KK+3,16),B(KK:KK+3,16),MY$P+1)
  ENDIF
ENDDO
    
```

Fig. 8. ADI with coarse-grain pipelining (first loop only).

Table III. Performance of ADI Integration (in seconds)

Problem Size	$P$	ADI w/o Loop Fusion	Fine-grain Pipelining	Coarse-grain Pipelining	Data Redistribution	Perfect Speedup	
256	1	1.22	1.45	1.32	-	1.22	
	2	1.13	0.78	0.69	1.32	0.61	
	×	4	1.02	0.45	0.38	0.83	0.30
	256	8	0.96	0.28	0.21	0.52	0.15
	16	0.93	0.20	0.12	0.32	0.08	
	32	0.96	0.17	0.08	0.25	0.04	
512	1	5.44	6.26	5.93	-	5.44	
	2	4.79	3.18	2.98	6.17	2.72	
	×	4	4.29	1.72	1.53	3.72	1.36
	512	8	4.04	0.97	0.81	2.02	0.68
	16	3.94	0.59	0.44	1.18	0.34	
	32	3.94	0.41	0.26	0.68	0.17	
1K	1	21.74	-	-	-	21.74	
	4	17.13	6.44	5.98	-	5.44	
	×	8	16.09	3.42	3.07	8.95	2.72
	1K	16	15.61	1.91	1.62	4.59	1.36
		32	15.47	1.17	0.89	2.58	0.68

$K$  loop allow each processor to reduce communication costs at the expense of some parallelism, resulting in Figure 8. Except for coarse-grain pipelining, all these versions of ADI integration were generated automatically by the Fortran D compiler.

## 6.2 Performance Results

To validate these methods, we executed these codes on an iPSC/860. The programs were compiled under O4 using Release 3.0 of *if77*, the iPSC/860 compiler. Timings were taken for three double-precision problem sizes using *clock*( ) on a 32-node Intel iPSC/860 with 8 Meg of memory per node. Results for three problem sizes are tabulated in Table III. Timings are not provided where problem size exceeds available memory. We also graphically display

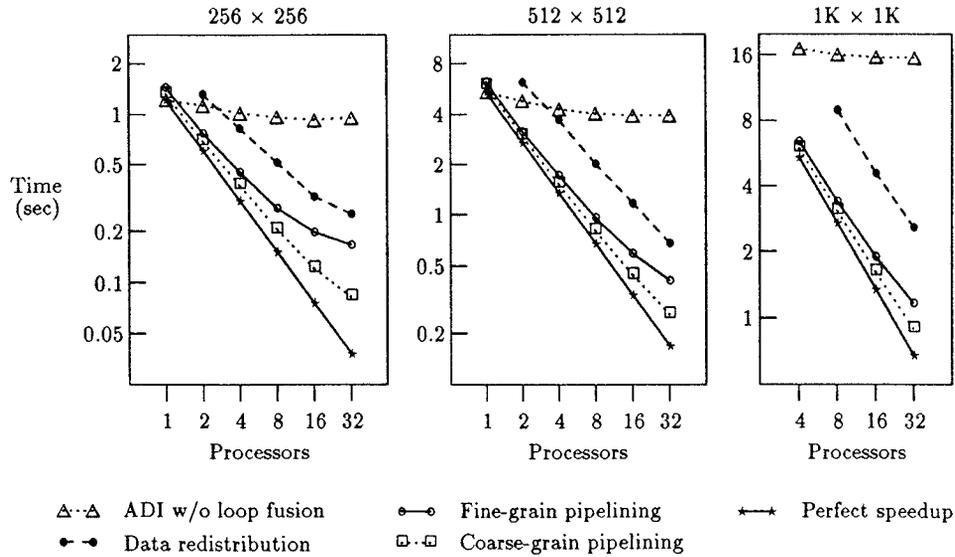


Fig. 9. Execution times for ADI integration (double-precision).

the timings in Figure 9. Execution times in seconds are plotted logarithmically along the Y-axis. The number of processors used is plotted logarithmically along the X-axis. Numbers for perfect or ideal speedup (sequential execution time divided by number of processors) are provided for comparison.

The original version of ADI (Figure 6) exploits pipeline parallelism in the  $J$  loop, but shows limited speedup, since the  $I$  loop is sequentialized. Fusing the  $K$  loops to improve memory reuse provides very little improvement in this case, yielding nearly identical results. Applying loop interchange after fusion to enable fine-grain pipelining (Figure 7) parallelizes the  $I$  loop as well, yielding significant speedup. Strip-mining to apply coarse-grain pipelining can improve efficiency an additional 10-50% (Figure 8). Pipelining comes closest to perfect speedup for large problems on a small number of processors.

We also compared the efficiency of pipelining versus dynamic data decomposition. By changing the distribution of data at run-time from columns to rows, all dependences in each sweep of ADI may be internalized, enabling completely parallel execution. Data must be redistributed twice, once to achieve the desired distribution, then a second time to return it to its original configuration. The cost of redistributing is approximated by the performance of the TRANSPOSE routine shown in Table II.

Our results show that on the iPSC/860, dynamic data decomposition for this formulation of ADI integration achieves speedup. However, the resulting program is significantly slower than pipelining, even for small problems distributed over large numbers of processors, the expected best case for dynamic data decomposition. Our experiences show that some common algorithms, such as ADI integration, require significant amounts of optimization to compete with hand-crafted code.

## 7. RELATED WORK

The Fortran D compiler is a second-generation distributed-memory compiler that integrates and extends many previous analysis and optimization techniques. Many distributed-memory compilers reduce communication overhead by aggregating messages outside of parallel loops [Ikudome et al. 1990; Koelbel and Mehrotra 1991] or parallel procedures [Hatcher et al. 1991; Rosing et al. 1991], while others rely on functional language [Li and Chen 1991] or single assignment semantics [Rogers and Pingali 1989]. In comparison, the Fortran D compiler uses dependence analysis to automatically exploit parallelism and extract communication even from sequential loops such as those found in ADI integration.

Several other projects are also developing Fortran 90 compilers for MIMD distributed-memory machines. ADAPT [Merlin 1991] and ADAPTOR [Brandes 1993] propose to scalarize and partition Fortran 90 programs using a run-time library for Fortran 90 intrinsics. The CM FORTRAN compiler compiles Fortran 90 with alignment and layout specifications directly to the physical machine, and can optimize floating-point register usage [Albert et al. 1988; Bromley et al. 1991]. The FORTRAN-90-Y compiler uses formal specification techniques to generate efficient code for the CM-2 and CM-5 [Chen and Cowie 1992]. FORGE90, formerly MIMDIZER, is an interactive parallelization system for MIMD shared- and distributed-memory machines from Applied Parallel Research [Applied Parallel Research 1992]. It performs data flow and dependence analyses, and also supports loop-level transformations. Our compiler resembles the VIENNA FORTRAN 90 compiler [Benkner 1992] derived from SUPERB [Zima et al. 1988]. It has also been influenced by a proposal by Wu and Fox that discussed program generation and optimization using a test-suite approach [Wu and Fox 1992].

A number of researchers have studied techniques to reduce storage and promote memory reuse [Allen and Kennedy 1992; Kennedy and McKinley 1992; Kuck et al. 1981; McKinley 1992; Sarkar and Gao 1991]. These optimizations have proved useful for both scalar and parallel machines. The goal of the Fortran 90D compiler is to integrate these scalarization techniques with advanced communication and parallelism optimizations.

## 8. CONCLUSIONS

This paper presents an integrated approach to compiling both Fortran 77D and 90D based on a few key observations. First, using FORALL preserves information in Fortran 90 array constructs. Dividing the scalarization process into translation, loop fusion, and sectioning allows it to be easily integrated with the partitioning performed by the Fortran D compiler. A portable run-time library can also reduce the complexity and machine-dependence of the compiler. All optimizations except data prefetching have been implemented in the current Fortran D compiler prototype.

Compiling for MIMD distributed-memory machines is only a part of the Fortran D project. We also are working on Fortran 77D and Fortran 90D compilers for SIMD machines, translations between the two Fortran dialects,

support for irregular computations, and environmental support for static performance estimation and automatic data decomposition [Balasundaram et al. 1990; 1991; Hiranandani et al. 1991; Kennedy and Kremer 1991].

#### ACKNOWLEDGMENTS

We are grateful to the ParaScope and Fortran D research groups for their assistance, and to Parasoft for providing the Fortran 90 parser and Express. Use of the Intel iPSC/860 was provided by the CRPC under NSF Cooperative Agreement Nos. CCR-8809615 and CDA-8619893 with support from the Keck Foundation. The content of this information does not necessarily reflect the position or policy of the government, and no official endorsement should be inferred.

#### REFERENCES

- AHMAD, I., CHOUDHARY, A., FOX, G., PARASURAM, K., PONNUSAMY, R., RANKA, S., AND THAKUR, R. 1992. Implementation and scalability of Fortran 90D intrinsic functions on distributed memory machines. Technical Report SCCS-256, NPAC, Syracuse Univ., Mar.
- ALBERT, E., KNOBE, K., LUKAS, J., AND STEELE, JR., G. 1988. Compiling Fortran 8x array features for the Connection Machine computer system. In *Proceedings of the ACM SIGPLAN Symposium on Parallel programming: Experience with Applications, Languages, and Systems (PPEALS)* (New Haven, CT, July). ACM, New York.
- ALBERT, E., LUKAS, J., AND STEELE, JR., G. 1991. Data parallel computers and the FORALL statement. *J. Parallel Distrib. Comput.* 13, 2 (Oct.), 185–192.
- ALLEN, J. R., AND KENNEDY, K. 1992. Vector register allocation. *IEEE Trans. Comput.* 41, 10 (Oct.), 1290–1317.
- ALLEN, J. R., AND KENNEDY, K. 1987. Automatic translation of Fortran programs to vector form. *ACM Trans. Prog. Lang. Syst.* 9, 4 (Oct.), 491–542.
- ANSI X3J3 / S8.115. 1990. Fortran 90, June. ANSI, New York.
- APPLIED PARALLEL RESEARCH. 1992. *Forge 90 Distributed Memory Parallelizer: User's Guide*, version 8.0 ed. Placerville, CA.
- BALASUNDARAM, V., FOX, G., KENNEDY, K., AND KREMER, U. 1991. A static performance estimator to guide data partitioning decisions. In *Proceedings of the 3rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (Williamsburg, VA, Apr.). ACM, New York.
- BALASUNDARAM, V., FOX, G., KENNEDY, K., AND KREMER, U. 1990. An interactive environment for data partitioning and distribution. In *Proceedings of the 5th Distributed Memory Computing Conference* (Charleston, SC, Apr.).
- BELTRAMETTI, M., BOBEY, K., AND ZORBAS, J. 1988. The control mechanism for the Myrias parallel computer system. *Comput. Arch. News* 16, 4 (Sept.), 21–30.
- BENKNER, S., CHAPMAN, B., AND ZIMA, H. 1992. Vienna Fortran 90. In *Proceedings of the 1992 Scalable High Performance Computing Conference* (Williamsburg, VA, Apr.).
- BRANDES, T. 1993. Automatic translation of data parallel programs to message passing programs. In *Proceedings of AP'93 International Workshop on Automatic Distributed Memory Parallelization, Automatic Data Distribution and Automatic Parallel Performance Prediction* (Saarbrücken, Germany, Mar.).
- BROMLEY, M., HELLER, S., MCNERNEY, T., AND STEELE, JR., G. 1991. Fortran at ten gigaflops: The Connection Machine convolution compiler. In *Proceedings of the SIGPLAN '91 Conference on Program Language Design and Implementation* (Toronto, Canada, June). ACM, New York.
- CALLAHAN, D., AND KENNEDY, K. 1988. Compiling programs for distributed-memory multiprocessors. *J. Supercomput.* 2 (Oct.), 151–169.
- CARR, S., KENNEDY, K., MCKINLEY, K. S., AND TSENG, C. 1992. Compiler optimizations for improving data locality. Tech. Rep. TR92-195, Dept. of Computer Science, Rice Univ., Nov.
- ACM Letters on Programming Languages and Systems, Vol. 2, Nos. 1–4, March–December 1993

- CHATTERJEE, S., GILBERT, J., SCHREIBER, R., AND TENG, S. 1993. Automatic array alignment in data-parallel programs. In *Proceedings of the 20th Annual ACM Symposium on the Principles of Programming Languages* (Charleston, SC, Jan.). ACM, New York.
- CHEN, M., AND COWIE, J. 1992. Prototyping Fortran-90 compilers for massively parallel machines. In *Proceedings of the SIGPLAN '92 Conference on Program Language Design and Implementation* (San Francisco, CA, June). ACM, New York.
- COOPER, K., HALL, M. W., HOOD, R. T., KENNEDY, K., MCKINLEY, K. S., MELLOR-CRUMMEY, J. M., TORCZON, L., AND WARREN, S. K. 1993. The ParaScope parallel programming environment. *Proc. IEEE* 81, 2 (Feb.), 244–263.
- FOX, G., HIRANANDANI, S., KENNEDY, K., KOELBEL, C., KREMER, U., TSENG, C., AND WU, M. 1990. Fortran D language specification. Tech. Rep. TR90-141, Dept. of Computer Science, Rice Univ., Dec.
- HALL, M. W., HIRANANDANI, S., KENNEDY, K., AND TSENG, C. 1992. Interprocedural compilation of Fortran D for MIMD distributed-memory machines. In *Proceedings of Supercomputing '92* (Minneapolis, MN, Nov.).
- HATCHER, P., QUINN, M., LAPADULA, A., SEEVERS, B., ANDERSON, R., AND JONES, R. 1991. Data-parallel programming on MIMD computers. *IEEE Trans. Parallel Distrib. Syst.* 2, 3 (July), 377–383.
- HIGH PERFORMANCE FORTRAN FORUM. 1993. High Performance Fortran language specification, version 1.0. Tech. Rep. CRPC-TR92225, Center for Research on Parallel Computation, Rice Univ., Houston, TX, Jan.
- HIRANANDANI, S., KENNEDY, K., KOELBEL, C., KREMER, U., AND TSENG, C. 1991. An overview of the Fortran D programming system. In *Languages and Compilers for Parallel Computing, 4th International Workshop* (Santa Clara, CA, Aug.), U. Banerjee, D. Gelernter, A. Nicolau, and D. Padua, Eds. Springer-Verlag, New York.
- HIRANANDANI, S., KENNEDY, K., AND TSENG, C. 1992a. Compiling Fortran D for MIMD distributed-memory machines. *Commun. ACM* 35, 8 (Aug.), 66–80.
- HIRANANDANI, S., KENNEDY, K., AND TSENG, C. 1992b. Evaluation of compiler optimizations for Fortran D on MIMD distributed-memory machines. In *Proceedings of the 1992 ACM International Conference on Supercomputing* (Washington, DC, July). ACM, New York.
- KUDOME, K., FOX, G., KOLAWA, A., AND FLOWER, J. 1990. An automatic and symbolic parallelization system for distributed memory parallel computers. In *Proceedings of the 5th Distributed Memory Computing Conference* (Charleston, SC, Apr.).
- KENNEDY, K., AND KREMER, U. 1991. Automatic data alignment and distribution for loosely synchronous problems in an interactive programming environment. Tech. Rep. TR91-155, Dept. of Computer Science, Rice Univ., Apr.
- KENNEDY, K., AND MCKINLEY, K. S. 1992. Optimizing for parallelism and data locality. In *Proceedings of the 1992 ACM International Conference on Supercomputing* (Washington, DC, July). ACM, New York.
- KNOBE, K., LUKAS, J., AND STEELE, JR., G. 1990. Data optimization: Allocation of arrays to reduce communication on SIMD machines. *J. Parallel Distrib. Comput.* 8, 2 (Feb.), 102–118.
- KOELBEL, C., AND MEHROTRA, P. 1991. Compiling global name-space parallel loops for distributed execution. *IEEE Trans. Paralle. Distrib. Syst.* 2, 4 (Oct.), 440–451.
- KUCK, D., KUHN, R., PADUA, D., LEASURE, B., AND WOLFE, M. J. 1981. Dependence graphs and compiler optimizations. In *Conference Record of the 8th Annual ACM Symposium on the Principles of Programming Languages* (Williamsburg, VA, Jan.). ACM, New York.
- LI, J., AND CHEN, M. 1991. Compiling communication-efficient programs for massively parallel machines. *IEEE Trans. Paralle. Distrib. Syst.* 2, 3 (July), 361–376.
- LUNDSTROM, S., AND BARNES, G. 1980. Controllable MIMD architectures. In *Proceedings of the 1980 International Conference on Parallel Processing* (St. Charles, IL, Aug.).
- MCKINLEY, K. S. 1992. Automatic and interactive parallelization. Ph.D. thesis, Dept. of Computer Science, Rice Univ., Apr.
- MERLIN, J. 1991. ADAPTING Fortran-90 array programs for distributed memory architectures. In *First International Conference of the Austrian Center for Parallel Computation* (Salzburg, Austria, Sept.).
- PARASOFT CORP. 1989. *Express User's Manual*.

- ROGERS, A., AND PINGALI, K. 1989. Process decomposition through locality of reference. In *Proceedings of the SIGPLAN '89 Conference on Program Language Design and Implementation* (Portland, OR, June). ACM, New York.
- ROSIING, M., SCHNABEL, R., AND WEAVER, R. 1991. The DINO parallel programming language. *J. Parall. Distrib. Comput.* 13, 1 (Sept.), 30-42.
- SARKAR, V., AND GAO, G. 1991. Optimization of array accesses by collective loop transformations. In *Proceedings of the 1991 ACM International Conference on Supercomputing* (Cologne, Germany, June). ACM, New York.
- TSENG, C. 1993. An optimizing Fortran D compiler for MIMD distributed-memory machines. Ph.D. thesis, Dept. of Computer Science, Rice Univ., Jan.
- WOLFE, M. J. 1989. *Optimizing Supercompilers for Supercomputers*. The MIT Press, Cambridge, MA.
- WU, M., AND FOX, G. 1992. A test suite approach for Fortran 90D compilers on MIMD distributed memory parallel computers. In *Proceedings of the 1992 Scalable High Performance Computing Conference* (Williamsburg, VA, Apr.).
- ZIMA, H., BAST, H.-J., AND GERNDT, M. 1988. SUPERB: A tool for semi-automatic MIMD/SIMD parallelization. *Parallel Comput.* 6, 1-18.