

# **SMS USER'S GUIDE**

**Tom Henderson**

**Dan Schaffer**

**Mark Govett**

**Leslie Hart**

Advanced Computing Branch  
Aviation Division  
NOAA/Forecast Systems Laboratory  
325 Broadway  
Boulder, Colorado 80303

September 2001  
SMS Software Version: 2.5

<http://www-ad.fsl.noaa.gov/ac/sms.html>.

<b>1</b>	<b>INTRODUCTION.....</b>	<b>5</b>
1.1	ORGANIZATION OF THIS DOCUMENT .....	5
1.2	TERMS AND CONVENTIONS .....	6
<b>2</b>	<b>GETTING STARTED .....</b>	<b>7</b>
2.1	BASIC PARALLELIZATION STEPS .....	7
2.2	A VERY SIMPLE PROGRAM .....	7
2.3	SIMPLE COMPUTATION ON A REGULAR GRID.....	9
2.3.1	<i>Parallelization by Domain Decomposition.....</i>	<i>10</i>
2.3.2	<i>Parallel Printing .....</i>	<i>16</i>
2.3.3	<i>Reduction.....</i>	<i>17</i>
2.4	BOUNDARY INITIALIZATION.....	18
2.5	A SIMPLE FDA PROGRAM .....	21
2.6	WRITING OUTPUT TO DISK.....	27
2.7	USING MULTIPLE DECOMPOSITIONS .....	27
<b>3</b>	<b>DECOMPOSING ARRAYS AND PARALLELIZING LOOPS.....</b>	<b>30</b>
3.1	CHOOSING DECOMPOSITIONS .....	30
3.2	TWO-DIMENSIONAL DECOMPOSITIONS .....	31
3.3	CONTROLLING PROCESS LAYOUT .....	32
3.4	USING STATICALLY ALLOCATED MEMORY.....	35
3.5	LOAD BALANCING.....	38
3.6	MORE ABOUT DISTRIBUTE .....	40
3.6.1	<i>Further Detail on DISTRIBUTE Syntax.....</i>	<i>40</i>
3.6.2	<i>Using DISTRIBUTE to Define Decomposed Boundary Arrays.....</i>	<i>41</i>
3.7	MORE ABOUT PARALLEL .....	42
3.8	ARRAYS WITH NON-UNIT LOWER BOUNDS.....	44
<b>4</b>	<b>TRANSLATING ARRAY INDICES .....</b>	<b>45</b>
4.1	TRANSLATING LOCAL INDICES TO GLOBAL INDICES.....	45
4.2	TRANSLATING GLOBAL INDICES TO LOCAL INDICES INSIDE LOOPS .....	46
4.3	USING TO_LOCAL TO GENERATE PROCESS-LOCAL SIZES AND LOOP BOUNDS.....	47
4.4	USING GLOBAL INDEX TO HANDLE BOUNDARY CONDITIONS .....	49
<b>5</b>	<b>HANDLING ADJACENT DEPENDENCIES.....</b>	<b>52</b>
5.1	FURTHER DETAILS ON EXCHANGE .....	52
5.1.1	<i>Using EXCHANGE in the Case of Two-Dimensional Decompositions....</i>	<i>52</i>
5.1.2	<i>Larger Stencils .....</i>	<i>56</i>
5.1.3	<i>Exchanges in Static Memory Models .....</i>	<i>58</i>
5.1.4	<i>Miscellaneous.....</i>	<i>59</i>
5.2	PERFORMANCE OPTIMIZATIONS .....	59
5.2.1	<i>Limited-Thickness Exchanges .....</i>	<i>59</i>
5.2.2	<i>Aggregating Exchanges.....</i>	<i>60</i>
5.2.3	<i>Exchanging Array Sections .....</i>	<i>62</i>

5.2.4	<i>Trading Communications for Computations Using HALO_COMP</i> .....	63
5.2.5	<i>Using HALO_COMP and TO_LOCAL To Make Subroutines Do Redundant Computations</i> .....	67
<b>6</b>	<b>HANDLING COMPLEX DEPENDENCIES USING TRANSFER</b> .....	<b>69</b>
6.1	FURTHER DETAILS ABOUT TRANSFER .....	69
6.2	APPLYING TRANSFER TO SPECTRAL MODELS .....	70
<b>7</b>	<b>HANDLING GLOBAL DEPENDENCIES USING REDUCE</b> .....	<b>72</b>
7.1	MORE ON STANDARD REDUCTIONS .....	72
7.2	BIT-WISE EXACT REDUCTIONS .....	74
<b>8</b>	<b>INCREMENTAL PARALLELIZATION USING SERIAL</b> .....	<b>77</b>
8.1	IMPROVING THE PERFORMANCE OF SERIAL .....	77
8.2	LIMITATIONS OF SERIAL .....	79
<b>9</b>	<b>PERIODIC BOUNDARY CONDITIONS</b> .....	<b>81</b>
<b>10</b>	<b>NESTING AND COUPLING: TRANSFER-INTERPOLATION</b> .....	<b>84</b>
10.1	USING SET_TRANSFER_INTERPOLATE TO DEFINE INTERPOLATIONS .....	86
10.1.1	<i>Transforming Serial Code to “Stencil-Oriented” Form</i> .....	86
10.1.2	<i>Defining an Interpolation</i> .....	87
10.2	USING TRANSFER TO INTERPOLATE BETWEEN GRIDS .....	89
10.3	USING SET_NEST_LEVEL TO SWITCH BETWEEN GRIDS .....	90
<b>11</b>	<b>I/O</b> .....	<b>92</b>
11.1	UNFORMATTED I/O .....	92
11.2	FORMATTED I/O .....	97
11.2.1	<i>Formatted Input</i> .....	97
11.2.2	<i>Formatted Output</i> .....	97
11.3	I/O PERFORMANCE TUNING .....	100
11.3.1	<i>General Guidelines</i> .....	101
11.3.2	<i>The SMS Server Process</i> .....	102
11.3.3	<i>Serverless I/O</i> .....	103
11.3.4	<i>The FLUSH_OUTPUT Directive</i> .....	103
11.3.5	<i>Improving Output Performance</i> .....	104
11.3.6	<i>Improving Input Performance</i> .....	106
<b>12</b>	<b>PROGRAM TERMINATION</b> .....	<b>108</b>
12.1	AUTOMATIC CODE GENERATION FOR TERMINATION .....	108
12.2	EXIT DIRECTIVE .....	108
12.3	MESSAGE DIRECTIVE.....	109
<b>13</b>	<b>DEBUGGING</b> .....	<b>110</b>
13.1	USING COMPARE_VAR To FIND PARALLELIZATION ERRORS.....	110
13.2	DEBUGGING ADJACENT DEPENDENCIES: CHECK_HALO.....	113
<b>14</b>	<b>BUILDING A PARALLEL PROGRAM</b> .....	<b>114</b>

14.1	OVERVIEW .....	114
14.2	PPP-GENERATED OUTPUT FILES.....	114
14.3	BUILDING SMS PARALLEL SOURCE CODE.....	114
14.3.1	<i>PPP Command Line Options</i> .....	114
14.3.2	<i>Examples</i> .....	115
14.4	BUILDING SMS PROGRAMS .....	118
14.4.1	<i>Makefile Compiler and Linker Options</i> .....	119
14.4.2	<i>Include File Handling</i> .....	119
14.4.3	<i>Building the Executable</i> .....	120
14.5	PPP ERROR REPORTING .....	120
14.5.1	<i>Parsing Errors</i> .....	120
14.5.2	<i>PPP Diagnostic Messages</i> .....	121
14.6	COMPILATION ERRORS .....	122
<b>15</b>	<b>RUNNING AN SMS PROGRAM.....</b>	<b>123</b>
15.1	INTRODUCTION.....	123
15.2	OPTIONAL COMMAND LINE ARGUMENTS.....	123
15.3	RUN-TIME ENVIRONMENT VARIABLES .....	124
15.4	RUN-TIME ERROR MESSAGES.....	125

# 1 Introduction

This document describes the Scalable Modeling System (SMS) and shows how SMS directives can be used to parallelize a serial Fortran program for distributed or shared memory machines. SMS is intended for use with programs that perform computations on regular gridded data sets. The primary application area thus far has been weather, ocean, and climate models. SMS has been used to parallelize models that use finite difference approximation (FDA) or the spectral transform method. SMS is general enough to be useful for parallelizing similar programs in other application areas.

Before reading this document, the reader should first read the companion overview document "SMS: A Directive-Based Parallelization Tool for Shared and Distributed Memory High Performance Computers". It is assumed that the reader of this User's Guide is familiar with the concepts and terms introduced in the overview document. The reader should also be familiar with basic parallel processing concepts such as distributed and shared memory, message latency and bandwidth, the Single Program Multiple Data (SPMD) programming model, and dependence analysis. The overview document describes these concepts briefly and contains references for further reading. After reading this User's Guide, the reader should have a good understanding of the steps that need to be taken to parallelize a serial program using SMS directives. If more detailed information about any directive is needed, the reader should refer to the companion reference document, "SMS Reference Manual". Answers to common questions and detailed discussions of problems not covered here may be found on the SMS FAQ web site at:

[http://www-ad.fsl.noaa.gov/ac/SMS\\_FAQ.html](http://www-ad.fsl.noaa.gov/ac/SMS_FAQ.html)

## 1.1 Organization of this Document

The SMS User's Guide begins by introducing the SMS directives in their simplest form. Section 2 introduces the most fundamental SMS directives with simple example programs that use the method of finite difference approximation. This section also introduces other SMS directives that are useful in transform-based programs such as spectral numerical weather prediction (NWP) models. The remaining sections describe in detail how the SMS directives are used in more complex situations. Section 3 explains how to divide work among multiple processes by the method of data decomposition and how to parallelize loops. Additional loop index translations needed during parallelization are described in Section 4. Sections 5, 6 and 7 cover further details about the inter-process communication directives introduced in Section 2. Section 8 describes a method by which parallelization can be done incrementally. Section 9 addresses periodic boundary conditions. Section 10 describes SMS support for mesh refinement (nesting) and coupling between different grids. Section 11 discusses parallel I/O. Directives that control program termination are dealt with in Section 12. Section 13 discusses debugging tools. Sections 14 and 15 explain how to build and run parallel SMS programs.

## 1.2 Terms and Conventions

Throughout most of this document, the term "process" is used instead of "processor" or "CPU". "Process" is slightly more general because it is possible to run more than one process on a single "processor" (and this may actually make sense on some types of CPU's that provide direct hardware support for multi-threaded applications). However, on most machines there will be a one-to-one mapping of processes to processors.

Fortran source code will appear in *courier* font. When program source code appears inside the main body of text, it will also be *italicized*. Large blocks of code will include line numbers to simplify discussions. Commands will also appear in *courier* font and will be preceded by a generic command line prompt, ">>". The results of commands will appear in *courier* font as well. Warning messages output by SMS will be ***courier bold***. File names will appear in *italics* when not in code examples or command lines. SMS directives will appear in **bold** in code examples. When directive arguments appear in the text they will be ***courier font, bold and italicized***. Sometimes example code will be a slightly modified version of a previous example. In that case, the changed lines will be **highlighted**.

## 2 Getting Started

### 2.1 Basic Parallelization Steps

The first step in any parallelization effort is to understand the performance characteristics of the serial program. Program components that take little time to run may not need to be parallelized at all. The second step is to perform dependence analysis to identify the places in the code where inter-process communication may be required. Dependencies will be discussed as relevant SMS directives are introduced. A strategy for dividing the work among the processes must then be chosen. SMS uses the method of domain decomposition in which portions of large arrays, and their associated computations, are assigned to each process. The dependence analysis is used to help pick optimal decompositions that will minimize inter-process communication. The final step is to add SMS directives to parallelize the code.

To build the parallel program, the Parallelizing Pre-Processor (PPP, a component of SMS) is first run to translate the source code with directives into new parallel source code. The translated source code is then compiled and linked with the SMS library to produce an executable program that can be run on multiple processes. The *smsRun* command is used to run the parallel program. The debugging features of SMS can then be used to test the parallel program.

SMS supports ANSI standard Fortran77 and Fortran90 language features such as full array assignment, allocatable arrays, namelist, pointer, include, do-endo, automatic arrays, and while statements. Partial support of modules is also offered (modules may contain variable declarations but not subroutines). Currently, only fixed-format source code is accepted. A more detailed description of supported language features can be found at the following SMS web site:

[http://www-ad.fsl.noaa.gov/ac/SMS\\_Supported\\_Fortran\\_Features.html](http://www-ad.fsl.noaa.gov/ac/SMS_Supported_Fortran_Features.html)

### 2.2 A Very Simple Program

Below is a simple Fortran program that prints a message on the screen:

```
program basic_ex1
  print *, 'HELLO'
end
```

If this program were stored in a file named *basic\_ex1.f*, it could be built using the following command:

```
>> f90 -o basic_ex1 basic_ex1.f
```

The above command assumes that the Fortran compiler is named “f90”. When run, the program produces the expected output on the screen:

```
>> basic_ex1
```

```
HELLO
```

This program is simple enough that a parallel version can be built directly without adding any SMS directives. To build with SMS, first run the Parallel Pre-Processor (PPP) to convert the print statements into parallel print statements:

```
>> ppp basic_ex1.f
```

The above command assumes that the SMS environment variable has been correctly set and that `$SMS/bin` is in the current path. For example, if SMS is located in the directory `/usr/local/sms/` then (assuming a c-shell environment) the SMS environment variable should be set as follows:

```
>> setenv SMS /usr/local/sms
```

The path could be modified using a command like this:

```
>> set path= ( $SMS/bin $path )
```

See Section 15.3 for details about setting other environment variables used by SMS. SMS translates the serial code in `basic_ex1.f` into parallel code and places the result in file `basic_ex1_sms.f`. Depending on the configuration of SMS, other temporary files may also be created. The next step is to compile `basic_ex1_sms.f` and link it to the SMS library.

```
>> f90 -c -I$SMS/include basic_ex1_sms.f
>> f90 -o basic_ex1_sms -I $SMS/include basic_ex1_sms.o -L$SMS/lib \
-lsms -lmpi
```

The above example assumes common behavior for f90 options "-I" (specify path for include files) and "-L" (specify path for libraries). Some Fortran compilers handle these options in slightly different ways. Note that link argument "-lmpi" links to the Message Passing Interface (MPI) library. SMS uses MPI to perform underlying low-level inter-process communication on most supported machines. Some machines may require different linkers or linker arguments to link to their MPI libraries.

The next step is to run the parallel program:

```
>> smsRun -np 1 basic_ex1_sms
```

The `smsRun` command shown above runs program `basic_ex1_sms` on 1 process. The output written to the screen will look something like this:

```
SMS:: Program started: 1999:12:02::15:55:22
SMS: BITWISE EXACT reductions will NOT be used.
HELLO
SMS:: Program complete, exiting: 1999:12:02::15:55:22 Elapsed Time = 0
sec.
```

All output lines beginning with "SMS::" are diagnostic messages from the SMS run-time system. The first and last output lines are time-stamps printed by SMS when a program begins and when it ends. These time-stamps are a useful guide for measuring wall-clock run times. The second text line is another message from SMS that indicates default behavior of reduction operations to be discussed in Section 7.2. Henceforth, diagnostic messages from SMS will usually be omitted for brevity. The remaining line contains the text that was output when this program was run as a serial Fortran code.

The program can be run on 3 processes using the *smsRun* command like this:

```
>> smsRun -np 3 basic_ex1_sms
```

The following text appears on the screen:

```
HELLO
```

This looks just like the run made on one process. Why? By default, SMS prints only one message per Fortran print (or write) statement to mimic the behavior of the original serial code as closely as possible. SMS also provides other "parallel print" modes, as described in Section 2.3.2 and in detail in Section 11.2.

## 2.3 Simple Computation on a Regular Grid

Example 2-1 illustrates a very simple code that initializes an array, performs a simple computation, and prints results on the screen. It consists of two parts: include file *basic.inc* and source file *basic\_ex2.f*.

```
[Include file: basic.inc]
```

```
integer im, jm  
common /sizes_com/ im, jm
```

```
[Source file: basic_ex2.f]
```

```
program basic_ex2  
include 'basic.inc'  
im = 10  
jm = 10  
call compute  
end  
  
subroutine compute  
include 'basic.inc'  
integer i, j, xsum  
integer x(im,jm)  
do 100 j=1,jm  
do 100 i=1,im  
x(i,j) = 1  
100 continue  
xsum = 0  
do 200 j=1,jm  
do 200 i=1,im
```

```

        xsum = xsum + x(i,j)
200 continue
    print *, 'xsum = ', xsum
    return
end

```

**Example 2-1: A simple serial code to initialize an array and print a global sum.**

This program initializes array  $x$ , sums the elements of  $x$ , and prints the result on the screen as shown below:

```

>> basic_ex2
xsum = 100

```

Notice that this program uses automatic (dynamically allocated) arrays instead of traditional Fortran77 static array declarations. The SMS directives support both dynamic and static memory allocation schemes. Examples with dynamic memory allocation are shown first because they are slightly simpler. Static allocation examples appear in Section 3.4.

### 2.3.1 Parallelization by Domain Decomposition

Programs such as this one that involve computations on regular grids are often best parallelized using the method of domain decomposition. Arrays and the computations performed on them are "decomposed" (divided up) among the processes as evenly as possible. For example, Figure 2-1, Figure 2-2, and Figure 2-3 show how array  $x$  might be decomposed in the  $i$  dimension over one, two and three processes.

Note that the sub-domains of array  $x$  become smaller as the number of processes increases. These sub-domains are referred to as "local" arrays because they cannot be accessed by other processes on a distributed memory machine. In SMS terms, the original array  $x$  in the serial code is sometimes referred to as a "global array". Indices used to access a global array are called "global indices" while indices used to access a local array are called "local indices". Similarly, sizes of the dimensions of a global array are called "global sizes" and sizes of the dimensions of a local array are called "local sizes". For dynamic memory code, the local and global indices are identical. We will see in Section 3.4 that the global and local indices differ from each other for static memory codes.

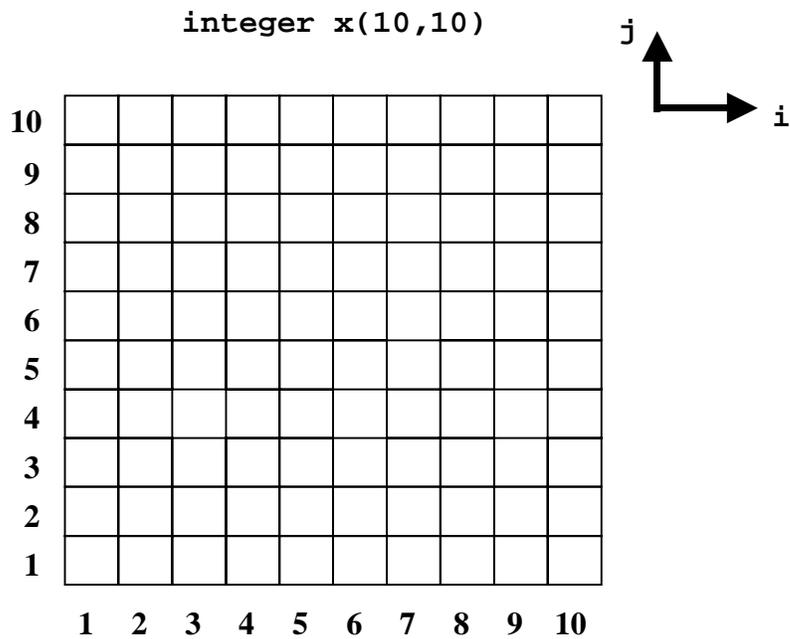


Figure 2-1: A graphical representation of a non-decomposed 10 by 10 integer array.

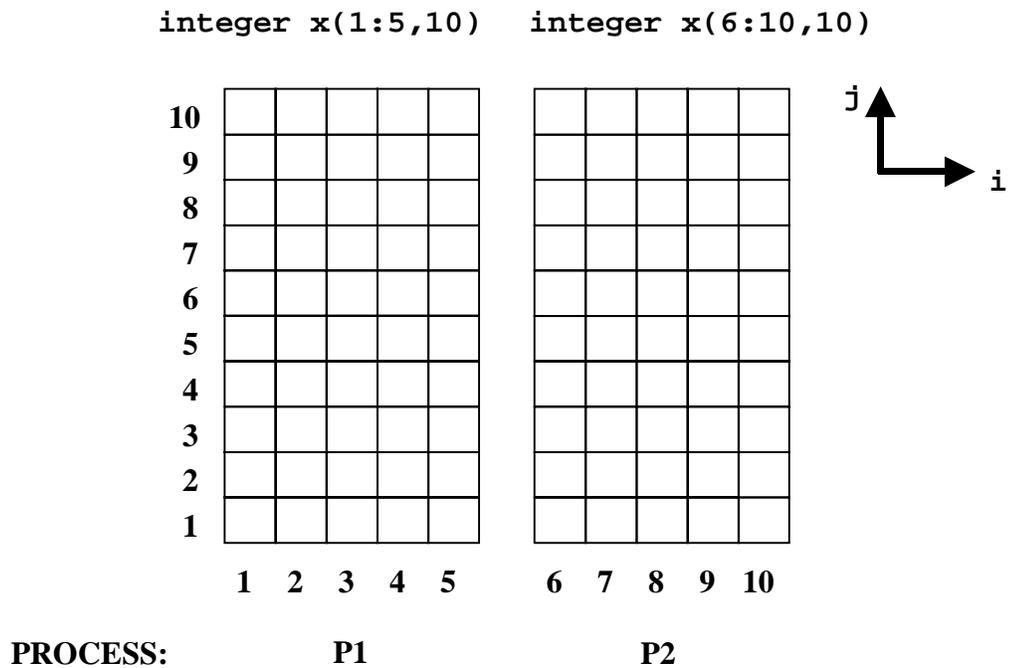
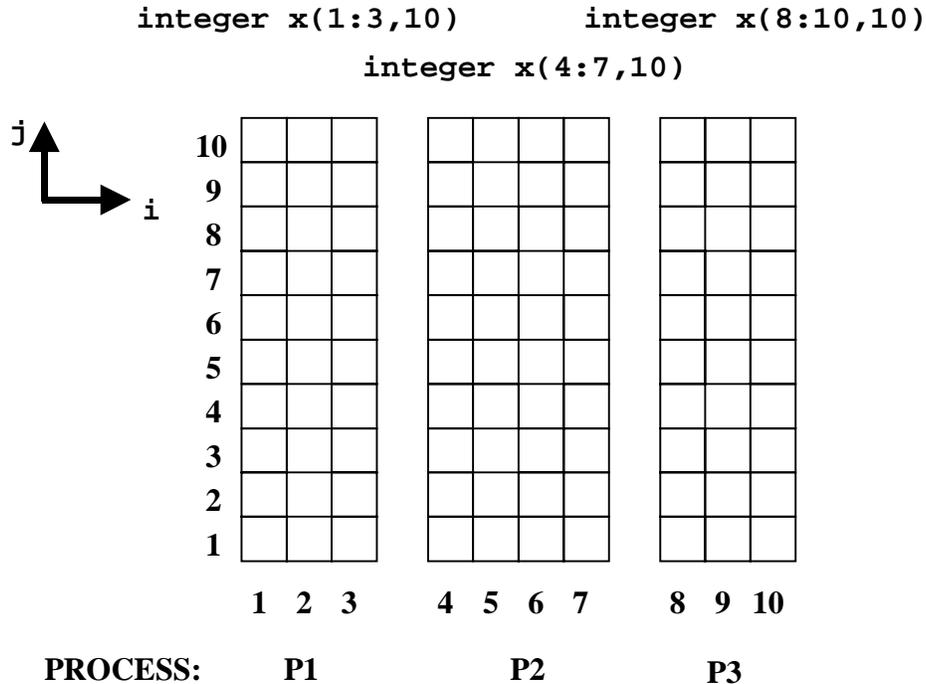


Figure 2-2: An illustration of a 10 by 10 array decomposed over two processes. These integer arrays are now local arrays declared by each process. When dynamic memory is used, global addressing is used to access local array elements. Thus, on process P2, the first dimension ranges from 6 to 10.



**Figure 2-3: A 10 by 10 array decomposed over three processes. In this example, the locally declared size of process P2 is larger than the sizes of P1 or P3.**

In this program, domain decomposition of array  $x$  requires three basic steps. First, the way in which  $x$  will be decomposed must be described. For this simple example, we choose to decompose only in the  $i$  dimension (decompositions of two dimensions are discussed in Section 3.2). Second, the declarations of array  $x$  should be modified to reflect smaller local sizes. Finally, the start and stop indices of each relevant loop must be changed to reflect the smaller range of local indices. These three steps are accomplished using four SMS directives. The `DECLARE_DECOMP` and `CREATE_DECOMP` directives are used to describe a decomposition. Array declarations are modified using the `DISTRIBUTE` directive while loop start and stop indices are changed using the `PARALLEL` directive. These directives have been inserted into the serial program as shown in Example 2-2:

```
[Include file: basic.inc]
1      integer im, jm
2      common /sizes_com/ im, jm
3      CSMS$DECLARE_DECOMP(DECOMP_I, 1)

[Source file: basic_ex2.f]

1      program basic_ex2
2      include 'basic.inc'
3      im = 10
4      jm = 10
5      CSMS$CREATE_DECOMP(DECOMP_I, <im>, <0>)
6      call compute
7      end
8
```

```

 9      subroutine compute
10      include 'basic.inc'
11      integer i, j, xsum
12  CSMS$DISTRIBUTE(DECOMP_I, 1) BEGIN
13      integer x(im,jm)
14  CSMS$DISTRIBUTE END
15  CSMS$PARALLEL(DECOMP_I,<i>) BEGIN
16      do 100 j=1,jm
17      do 100 i=1,im
18          x(i,j) = 1
19      100 continue
20      xsum = 0
21      do 200 j=1,jm
22      do 200 i=1,im
23          xsum = xsum + x(i,j)
24      200 continue
25  CSMS$PARALLEL END
26      print *, 'xsum = ', xsum
27      return
28      end

```

**Example 2-2: A simple serial code with comment-based SMS directives added.**

Notice that each of the SMS directives begins with five characters "**CSMS\$**" which makes it a Fortran comment. This is true for all SMS directives. Also, note that both the **DISTRIBUTE** and **PARALLEL** directives come as **BEGIN-END** pairs. When an SMS directive appears in this form, its scope consists of all lines of code between the "**BEGIN**" and "**END**" directives. Some SMS directives, such as **TRANSFER** (Section 6) and **REDUCE** (Section 7) may be used either alone or as a **BEGIN-END** pair. The text translation effects of a **BEGIN-END** directive pair do not extend into called subroutines.

The first directive, **DECLARE\_DECOMP**, is used to give a name to the SMS decomposition that will be used to divide among the processes the work done in loops 100 and 200. Its first argument, **DECOMP\_I**, is the user-chosen name for the decomposition. Any valid Fortran variable name (up to 20 characters long) may be used to name a decomposition provided it does not conflict with any variable in the serial code. The second argument, **1**, is an integer that indicates how many dimensions are decomposed. This argument is omitted if static memory allocation is used (see Section 3.4) or if the decomposed arrays have non-unit lower bounds (see Section 3.8).

Next, the **CREATE\_DECOMP** directive is used to describe what kind of decomposition **DECOMP\_I** will be. The first argument is the decomposition name **DECOMP\_I** specified in the **DECLARE\_DECOMP** directive. The second argument, **<im>**, describes the decomposition as a 1-dimensional decomposition where the number of data points is the global size of the original serial dimension. The last argument, **<0>**, indicates that this decomposition will have no halo regions (halo thickness = 0). Halo regions are introduced later in this section and are described in detail in Section 5.1.

The third directive, **DISTRIBUTE**, associates array **x** with the decomposition **DECOMP\_I**. The second argument is used to indicate how array dimension(s) correspond to the dimensions of the decomposition named **DECOMP\_I**. In this simple one-dimensional decomposition, **1** indicates that the first dimension of the array **x** will be decomposed as described by the single dimension of the SMS decomposition named

**DECOMP\_I.** The distinction between "dimension of an array" and "dimension of an SMS decomposition" will become more clear in the two-dimensional decomposition examples shown later in Section 3.2.

The DISTRIBUTE directive does two things. First, it identifies array declarations that will be translated to use local sizes. In the above example program, the DISTRIBUTE directive will cause the declaration of *x* to be translated to a local declarations such as one of those shown in Figure 2-1, Figure 2-2, or Figure 2-3 (depending on the number of processes). The second task of DISTRIBUTE is to provide information about how each array is decomposed to other SMS directives and to support automatic parallelization of some operations (such as unformatted I/O). These features are described in detail in later sections.

Finally, the PARALLEL directive identifies loops that must be modified to span the smaller local arrays during translation. The second argument, *<i>*, indicates that loops with loop index *i* should be translated to span the decomposed dimension of array *x*. For example, if the program in Example 2-2 is run on two processes then *i* loops 100 and 200 will span local indices 1 through 5 on each process.

Building the SMS parallel code is a bit more complicated than the previous example due to the presence of the include file that contains a directive. Two commands are now needed. The first translates the include file:

```
>> ppp --header basic.inc
```

The "--header" option to the PPP command indicates that the file is an include file and must be translated differently than a standard Fortran source file. In the command above, include file *basic.inc* will be translated into new SMS include file *basic.inc.SMS*. The second command requires PPP option "--Finclude" to translate the Fortran source file:

```
>> ppp --Finclude=basic.inc basic_ex2.f
```

The "--Finclude" option to the PPP command indicates that file *basic.inc* is an include file that has been translated by PPP. During translation of source file *basic\_ex2.f*, inclusions of this file will be translated from

```
    include 'basic.inc'  
to  
    include 'basic.inc.SMS'
```

to ensure that the translated include file is used.

Running this program on one process produces the expected result.

```
>> smsRun -np 1 basic_ex2_sms  
xsum = 100
```

However, when this program is run on two and three processes, the values of  $xsum$  differ from the serial run.

```
>> smsRun -np 2 basic_ex2_sms
xsum = 50

>> smsRun -np 3 basic_ex2_sms
xsum = 30
```

Why did the parallel program produce incorrect results? The answer lies in the computations made in loop 200. In this loop, all of the elements of array  $x$  are summed and the result is placed in variable  $xsum$ . However, when the program is run on two or three processes, each process sums only its own local sub-domain of  $x$  as illustrated in Figure 2-4, and Figure 2-5. To reproduce the result of the original serial code, we will need the REDUCE (see Section 2.3.3).

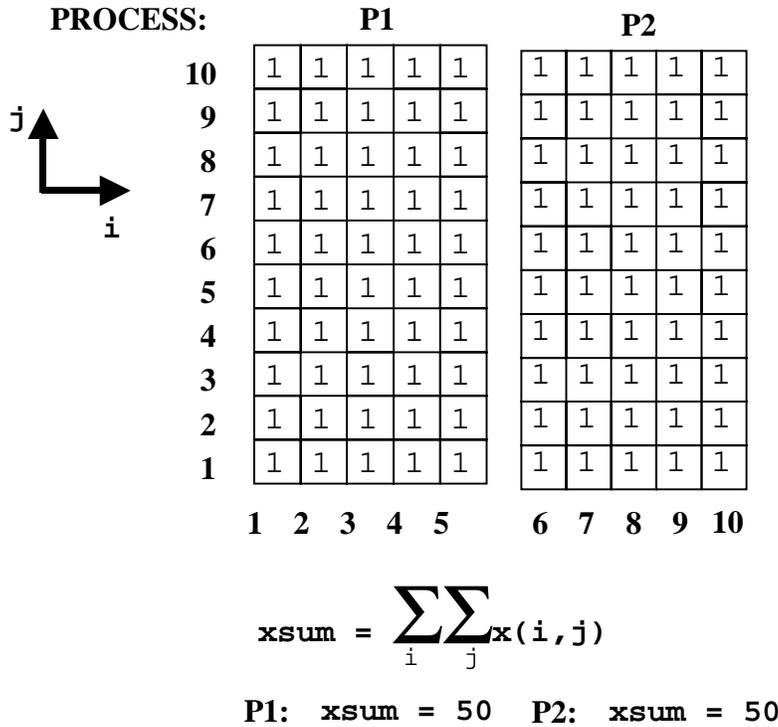


Figure 2-4: Each process sums its local portion of the array  $x$ .

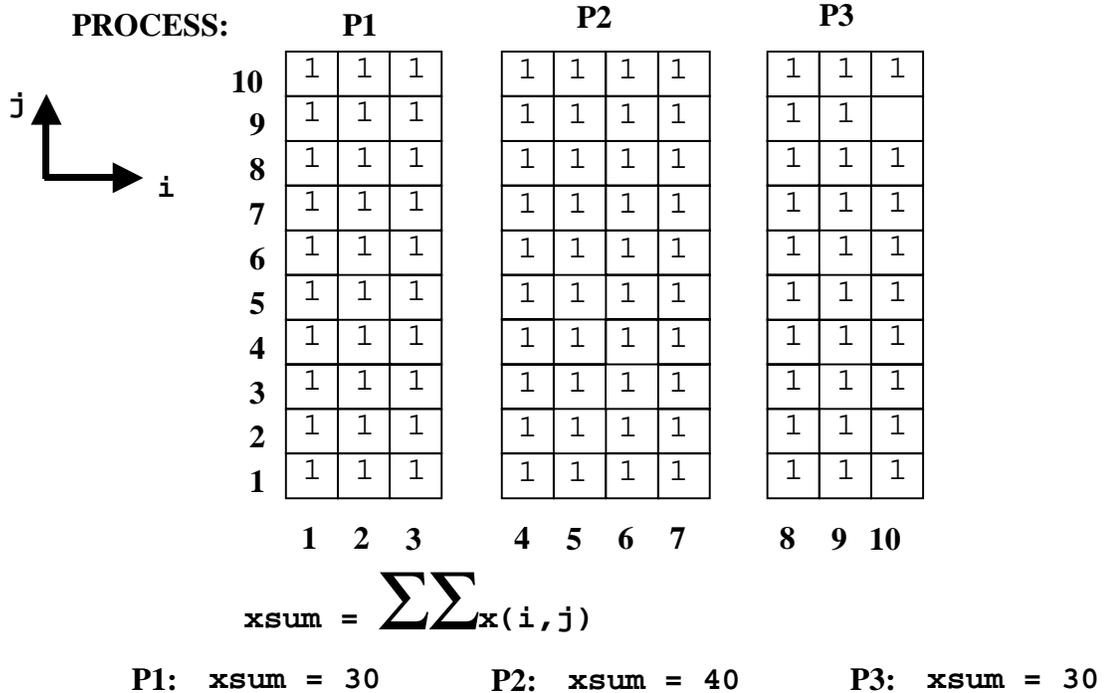


Figure 2-5: In this example, local sums are produced on each of the three processes.

### 2.3.2 Parallel Printing

In SMS, by default, only one process will print a message when a print statement is encountered. Therefore, the value of  $xsum$  printed is the value of  $xsum$  computed locally only on the printing process. We can see the value of  $xsum$  on every process by changing the default print behavior with the `PRINT_MODE` directive. The print statement on line 26 of the program in Example 2-2 would be modified as shown below:

```
CSMS$PRINT_MODE(ASYNC) BEGIN
    print *, 'xsum = ', xsum
CSMS$PRINT_MODE END
```

This `PRINT_MODE` directive changes the print mode from the default mode to "asynchronous" mode. When a print statement is encountered in asynchronous print mode, each process will print a message to the screen. When run on two processes, the following results are printed:

```
>> smsRun -np 2 basic_ex2_sms
xsum = 50
xsum = 50
```

Clearly, each process has computed the correct sum for its local half of array  $x$ . When run on three processes we may see any of the following results:

```

>> smsRun -np 3 basic_ex2_sms
xsum = 40
xsum = 30
xsum = 30

>> smsRun -np 3 basic_ex2_sms
xsum = 30
xsum = 40
xsum = 30

>> smsRun -np 3 basic_ex2_sms
xsum = 30
xsum = 30
xsum = 40

```

In the asynchronous print mode, the messages printed by each process may come out in any order. Another parallel print mode supported by SMS is the "ORDERED" print mode which preserves process order. Section 11.2 describes the SMS print modes in more detail.

### 2.3.3 Reduction

We have seen that each process has computed the correct sum for its local sub-domain of array  $x$ . To reproduce the same result as the original serial code, the local sums must be added together as shown in Figure 2-6. In more general terms, the computed value of  $xsum$  depends on all of the values of array  $x$ . This is known as a "global dependence" because the result of the computation depends on every element of global array  $x$ .

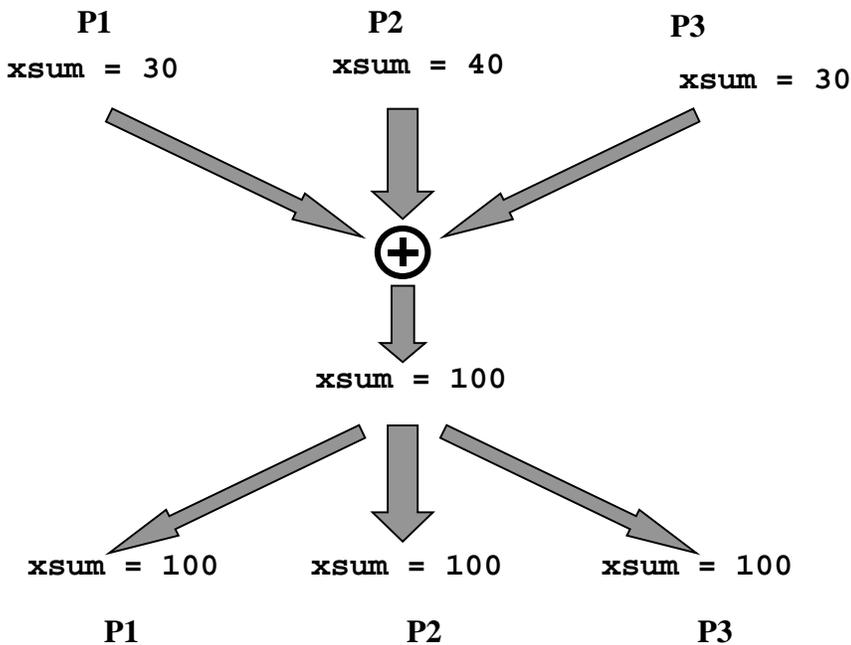


Figure 2-6: In this example, the reduction gathers the local sums, computes a global sum and then broadcasts the result out to the processes.

The REDUCE directive is used to resolve this dependence. To compute a global sum, insert the following line immediately before the print statement on line 26 of Example 2-2:

```
CSMS$REDUCE(xsum, SUM)
```

The REDUCE directive performs communications necessary to reduce the local values of a variable on each process to a single value that is identical on all processes. A specified operator is used to combine the values from each process. The first argument indicates that *xsum* is the name of the variable to be reduced. The second argument, **SUM**, specifies that the local values of *xsum* will be summed during reduction. Reductions are described in more detail in Section 7. The parallel program now produces the expected results when run on various numbers of processes (assuming the PRINT\_MODE directives used in Section 2.3.2 are removed):

```
>> smsRun -np 2 basic_ex2_sms
  xsum = 100
>> smsRun -np 3 basic_ex2_sms
  xsum = 100
```

## 2.4 Boundary Initialization

In Example 2-2, all elements of array *x* were initialized to the same value. Usually, it is desirable to initialize array elements differently depending on their location. This occurs often in models where elements near the model boundaries may be treated differently than other array elements. Example 2-3 below shows a variant of subroutine *compute* from Example 2-2 (changes are highlighted) that sets elements on the array boundaries where  $i=1$  or  $i=im$  to 2 and all other elements to 1. This is illustrated in Figure 2-7.

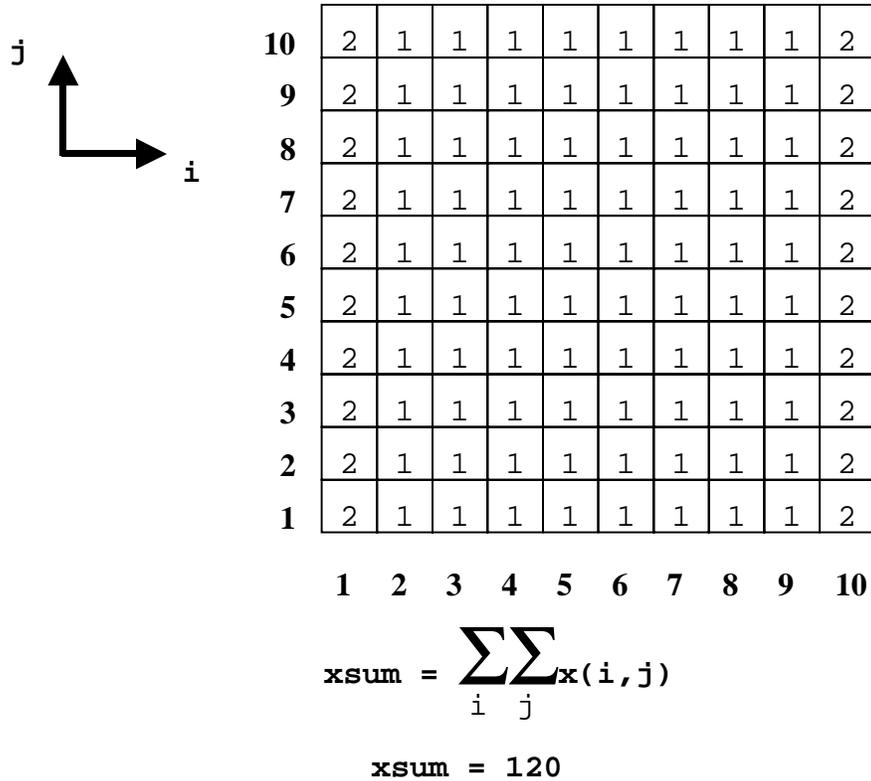


Figure 2-7: An illustration of a boundary initialization where edge point values are different than interior points.

```

1      subroutine compute
2      include 'basic.inc'
3      integer i, j, xsum
4      CSMS$DISTRIBUTE(DECOMP_I, 1) BEGIN
5      integer x(im,jm)
6      CSMS$DISTRIBUTE END
7      CSMS$PARALLEL(DECOMP_I,<i>) BEGIN
8      do 100 j=1,jm
9      do 100 i=1,im
10     x(i,j) = 1
11     100 continue
12     do 110 j=1,jm
13     x( 1,j) = 2
14     x(im,j) = 2
15     110 continue
16     xsum = 0
17     do 200 j=1,jm
18     do 200 i=1,im
19     xsum = xsum + x(i,j)
20     200 continue
21     CSMS$PARALLEL END
22     CSMS$REDUCE(xsum,SUM)
23     print *, 'xsum = ', xsum
24     return
25     end

```

**Example 2-3: Boundary initialization requires special handling.**

When the serial version of Example 2-3 is run, the following results are printed on the screen:

```
>> basic_ex3
xsum = 120
```

However, when the parallel code is run on more than one process, results are unpredictable:

```
>> smsRun -np 2 basic_ex3_sms
<core dump>
```

The reason for these erroneous results can be seen by examining new loop 110 in detail. Line 14 in loop 110 contains the following statement:

```
x(im,j) = 2
```

This statement will perform the following assignments:

```
x(10, 1) = 2
x(10, 2) = 2
...
x(10,10) = 2
```

On process 1 of a 2 process run, array  $x$  is dimensioned  $x(1:5, 1:10)$  (see Figure 2-2) so  $x(10,10)$  is out of bounds. The behavior of any program that performs such assignments is unpredictable. Similarly, line 13 causes an out-of-bounds assignment on process 2.

To address this problem, the assignment statements must be modified so they are only executed on the processes that contain the specified global indices in their local sub-domains. The GLOBAL\_INDEX directive solves these problems as shown below:

```
do 110 j=1,jm
CSMS$GLOBAL_INDEX(1) BEGIN
    x( 1,j) = 2
    x(im,j) = 2
CSMS$GLOBAL_INDEX END
110 continue
```

The GLOBAL\_INDEX directive ensures the enclosed statements are only executed on the appropriate processes. Now only process 1 will execute line 13 and only process 2 will execute line 14. The argument in the GLOBAL\_INDEX directive, **1**, indicates that these translations will be applied to array indices that correspond to the first (and in this case only) decomposed dimension. In this case, the decomposed dimension corresponds to the  $i$  dimension of array  $x$ . (The concept of "decomposed dimension" is explained in detail in Section 3.) The effects of the GLOBAL\_INDEX directives on the assignments of  $x(1, j)$  and  $x(im, j)$  are shown for the two process case in Figure 2-8.

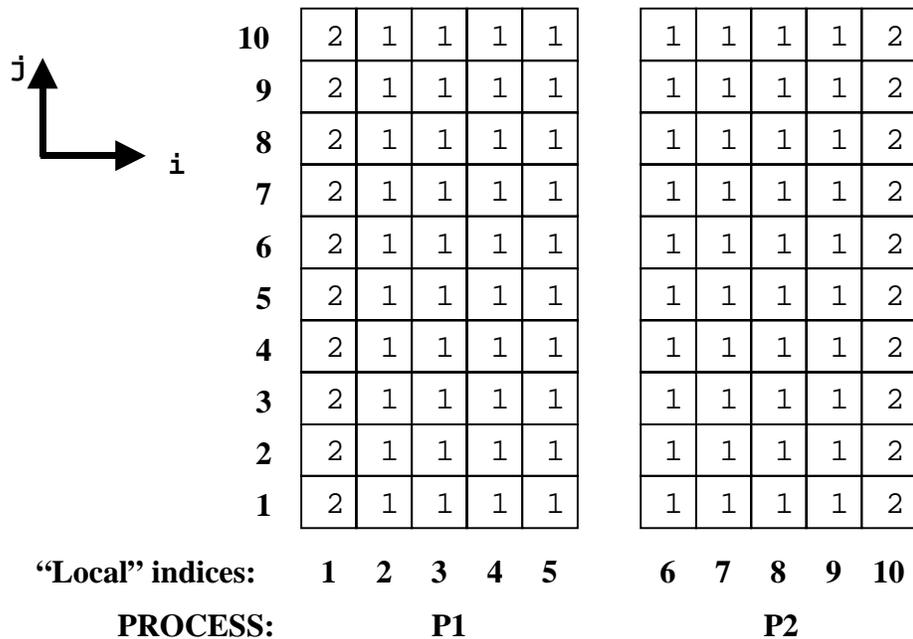


Figure 2-8: GLOBAL\_INDEX is used to initialize the boundaries of the array  $\mathbf{x}$ .

Now when the parallel code is run, results match the serial code:

```

>> smsRun -np 2 basic_ex3_sms
xsum = 120
>> smsRun -np 3 basic_ex3_sms
xsum = 120
  
```

## 2.5 A Simple FDA Program

The following example is a FDA program that solves Laplace's equation on a two-dimensional surface with fixed boundaries using Jacobi relaxation. On a two-dimensional surface, Laplace's equation takes the form:

$$\frac{\partial^2 f}{\partial^2 x} + \frac{\partial^2 f}{\partial^2 y} = 0$$

A simple approach is to discretize the two-dimensional space and use a finite difference approximation to the derivatives to seek a numerical solution. The discrete equation is:

$$4*f(i,j) - f(i-1,j) - f(i+1,j) - f(i,j-1) - f(i,j+1) = 0$$

The initial state is  $f$  on the boundaries. The boundaries are constant and non-periodic. The above equation is solved for  $f(i,j)$  iteratively until it converges. The solution is said to converge when the difference between successive solutions is less than a specified

threshold. The difference between values of  $f(i,j)$  in two successive iterations is the following:

$$df(i,j) = (1/4) * (f(i-1,j) + f(i+1,j) + f(i,j-1) + f(i,j+1)) - f(i,j)$$

Using the method of Jacobi relaxation, the value of  $f(i,j)$  during an iteration is calculated from the value of  $f(i,j)$  computed in the previous iteration as follows:

$$fnew(i,j) = fold(i,j) + df(i,j)$$

In Example 2-4 below, boundary elements of array  $f$  are initially set to 2.0 (lines 25-31). Laplace's equation is then solved and diagnostic messages are printed on the screen. Previously described SMS directives have already been inserted.

```
[Source file:  laplace.f]
1      program laplace
2      include 'basic.inc'
3      im = 10
4      jm = 10
5      CSMS$CREATE_DECOMP(DECOMP_I, <im>, <0>)
6      call laplace
7      end
8
9      subroutine laplace
10     include 'basic.inc'
11     integer i, j, iter
12     real max_error
13     real tolerance
14     parameter (tolerance = 0.001)
15     CSMS$DISTRIBUTE(DECOMP_I, 1) BEGIN
16     real f(im,jm), df(im,jm)
17     CSMS$DISTRIBUTE END
18     CSMS$PARALLEL(DECOMP_I,<i>) BEGIN
19     do 100 j=1,jm
20     do 100 i=1,im
21     f(i,j) = 0.0
22     100 continue
23     do 110 j=1,jm
24     CSMS$GLOBAL_INDEX(1) BEGIN
25     f( 1,j) = 2.0
26     f(im,j) = 2.0
27     CSMS$GLOBAL_INDEX END
28     110 continue
29     do 120 i=1,im
30     f(i, 1) = 2.0
31     f(i,jm) = 2.0
32     120 continue
33     iter = 0
34     max_error = 2.0 * tolerance
35     C main iteration loop...
36     do while ((max_error .gt. tolerance) .and. (iter .lt. 1000))
37     iter = iter + 1
38     max_error = 0.0
39     do 200 j=2,jm-1
40     do 200 i=2,im-1
41     df(i,j) = 0.25*(f(i-1,j) + f(i+1,j) + f(i,j-1) + f(i,j+1))
42     &          - f(i,j)
43     200 continue
44     do 300 j=2,jm-1
```

```

45         do 300 i=2,im-1
46             if (max_error .lt. abs(df(i,j))) then
47                 max_error = abs(df(i,j))
48             endif
49         300 continue
50     CSMS$REDUCE(max_error, MAX)
51         do 400 j=2,jm-1
52             do 400 i=2,im-1
53                 f(i,j) = f(i,j) + df(i,j)
54         400 continue
55     enddo
56     CSMS$PARALLEL END
57     print *, 'Solution required ',iter,' iterations'
58     print *, 'Final error = ', max_error
59
60     return
61     end

```

**Example 2-4: Serial code plus directives illustrate a parallel solution to Laplace's equation. This solution, using a one-dimensional decomposition, produces incorrect results.**

Notice that the REDUCE directive generates the global maximum error from the local maxima on each process.

When the serial program is run, the following messages are printed on the screen:

```

>> laplace
    Solution required  85  iterations
    Final error =  9.9968910E-4

```

When the parallel program is run on more than one process, results are incorrect:

```

>> smsRun -np 2 laplace_sms
    Solution required  45  iterations
    Final error =  9.9253654E-4

>> smsRun -np 3 laplace_sms
    Solution required  131 iterations
    Final error =  9.9420547E-4

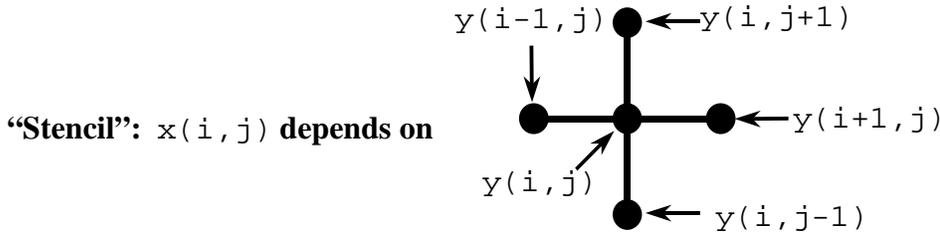
```

Why do results change for different numbers of processes? The answer lies in the computations made on lines 41 and 42:

$$df(i,j) = 0.25*(f(i-1,j) + f(i+1,j) + f(i,j-1) + f(i,j+1)) - f(i,j)$$

Here, each  $df(i,j)$  is computed from  $f(i-1,j)$ ,  $f(i+1,j)$ ,  $f(i,j-1)$ ,  $f(i,j+1)$ , and  $f(i,j)$ . This type of dependence is called an "adjacent dependence" because the computation at point  $(i,j)$  depends on data at adjacent (or "nearby") points. Adjacent dependencies are often represented graphically using a "stencil" as shown in Figure 2-9.

$$x(i, j) = y(i, j) + y(i+1, j) + y(i-1, j) + y(i, j-1) + y(i, j+1)$$



**Figure 2-9:** This five-point stencil illustrates the dependencies of the array  $y$  on the computation of  $x$ .

In Figure 2-10 stencils have been overlaid on graphical representations of the sub-domains assigned to each process during a run made on three processes. The stencil centered at point  $(2, 2)$  on process P1 illustrates that computations at this grid point require values from points  $(2, 2)$ ,  $(2, 1)$ ,  $(1, 2)$ ,  $(2, 3)$ , and  $(3, 2)$ . These array elements are all inside the local sub-domain of process P1. Similarly, computations at point  $(5, 8)$  depend only on array elements inside the local sub-domain of process P2. However, computations on sub-domain boundaries cannot be performed so easily. For example, the stencil centered at point  $(7, 5)$  on process P2 depends on the element at point  $(8, 5)$  which is located in the local sub-domain of process P3. Similarly, the stencil centered at point  $(8, 2)$  on process P3 requires an element from process P2. The results of the parallel program above are incorrect because no data is sent between processes to resolve the adjacent dependence in loop 200.

It is possible to solve this problem by sending single data points between processes. However, on high-latency machines, sending messages that contain only one array element is very inefficient compared to sending messages that contain many array elements. The most common approach to handle adjacent dependencies is to create "halo" or "ghost" regions to store these data as shown in Figure 2-11. Each halo point corresponds to an interior point of a neighboring process. For example, in Figure 2-11, halo point  $(8, 5)$  in process P2 corresponds to interior point  $(8, 5)$  in process P3. When data in these regions are needed, the halo regions are updated by swapping columns (or larger blocks) of data between processes as shown in Figure 2-12. This form of inter-process communication is called "exchange" and is supported by the EXCHANGE directive.

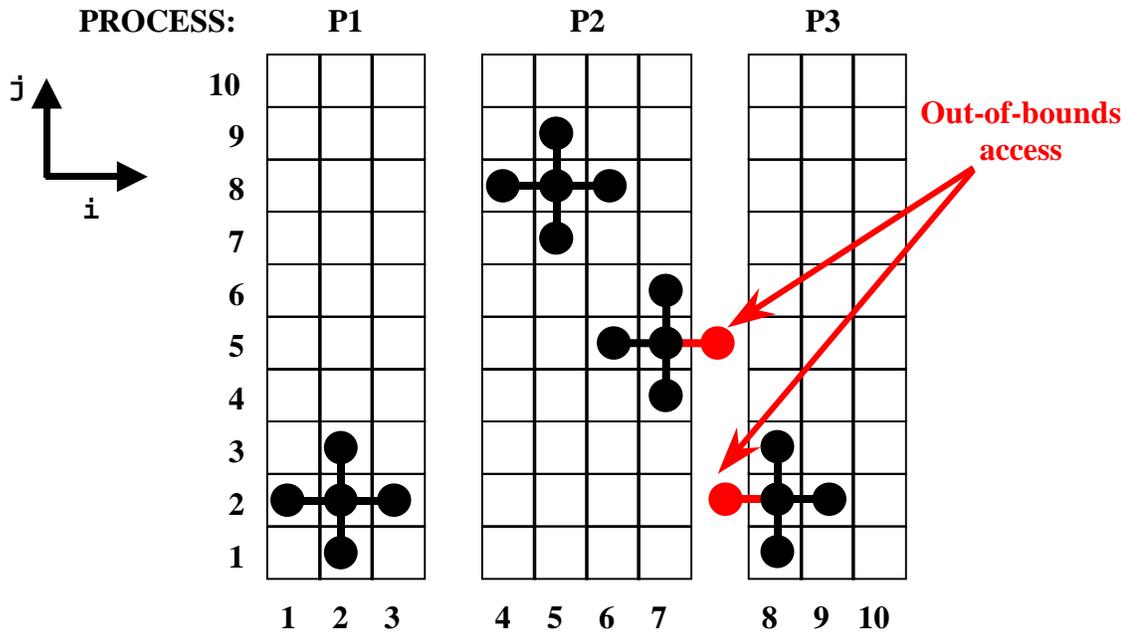


Figure 2-10: Illustration of how an adjacent dependence causes out-of-bounds data references on processes P2 and P3.

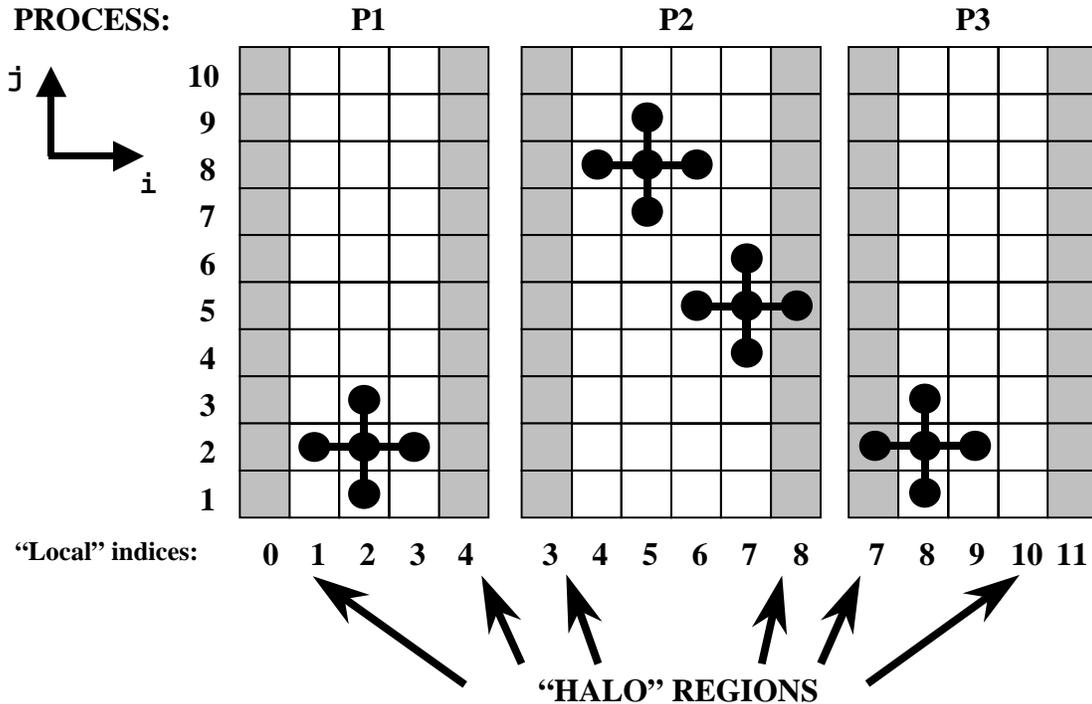


Figure 2-11: Halo regions eliminate the out-of-bounds array references. Notice the distinction between interior points (in white) and halo points (in gray). The local indices of the halo points on the domain edges actually lie outside the serial domain range (1 to 10). These edge halo points are only used for problems that have periodic boundary conditions as described in Section 9.

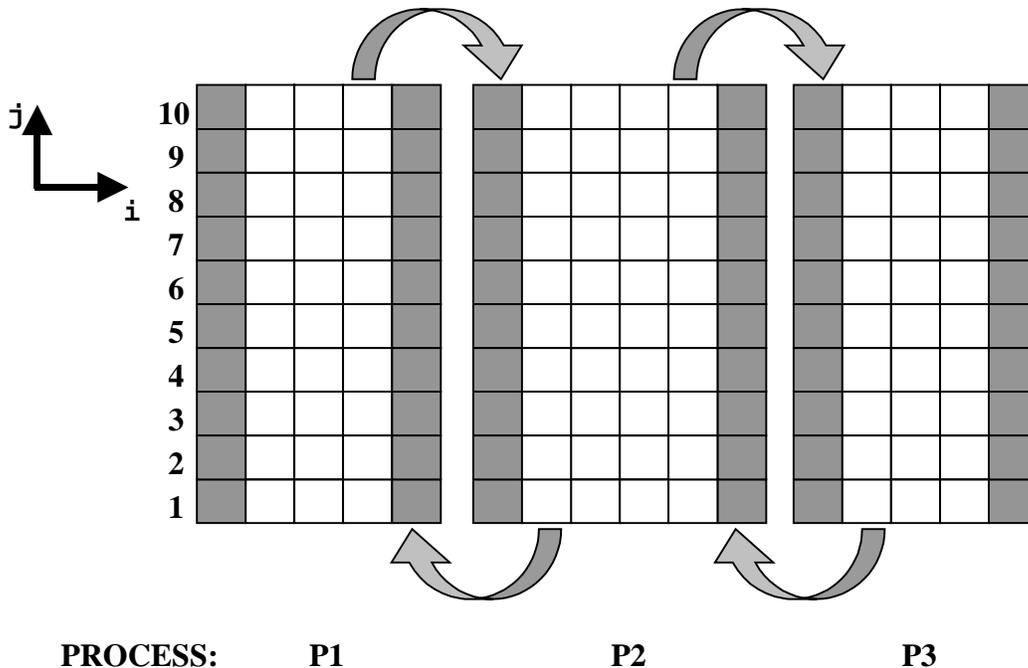


Figure 2-12: Halo regions are updated by exchanging data between adjacent processes.

The laplace program in Example 2-4 can be corrected by modifying line 5 to specify one halo point

```
CSMS$CREATE_DECOMP(DECOMP_I, <im>, <1>)
```

and by adding the following directive before line 39:

```
CSMS$EXCHANGE(f)
```

The third argument of the CREATE\_DECOMP directive has been changed to <1>. This indicates that all arrays decomposed using *DECOMP\_I* will have a halo region one point thick added in the first decomposed dimension (the *i* dimension in this case). The only argument of the EXCHANGE directive is the name of the variable (*f*) to be exchanged. The directive is placed immediately before loop 200 to ensure that halo regions of *f* are updated prior to the computations that need them. The EXCHANGE directive is described in more detail in section 5.1.

Now the parallel program produces the correct results on more than one process:

```
>> smsRun -np 2 laplace_sms
Solution required 85 iterations
Final error = 9.9968910E-4

>> smsRun -np 3 laplace_sms
Solution required 85 iterations
Final error = 9.9968910E-4
```

## 2.6 Writing Output to Disk

The Laplace solver in Example 2-4 would be more useful if the final state of array  $f$  could be written to disk. This is easily done by adding the following code fragment immediately before the `return` statement (line 60) in subroutine `laplace`:

```
open(10, file='f.out', form='unformatted')
write(10) f
close(10)
```

When the serial program is run, file `f.out` is written. For the SMS parallel program, no additional directives are required to handle this output. By default, SMS automatically generates `f.out` in exactly the same format as the serial program, for any number of processes. However, SMS can also produce other file formats as discussed in Section 11.1.

## 2.7 Using Multiple Decompositions

So far, we have seen how to parallelize a program that only requires a single domain decomposition. However, many programs require the use of different decompositions at different times to run efficiently in parallel. The TRANSFER directive provides the means to transform arrays between decompositions. Spectral models are prime candidates for application of TRANSFER (see Section 6.2).

In this section, we present a simple case where two different decompositions are needed. In Example 2-5, the statement at line 42 contains a dependence called a "recurrence relation". In this statement, an update to  $x(i, j)$  depends on  $x(i, j-1)$  which was updated in the previous loop iteration. SMS does not provide directives that directly support parallelization of a recurrence relation if the array dimension is decomposed. Since the second ( $j$ ) dimension of  $x$  is decomposed, SMS cannot handle this statement. Similarly, the statement at line 63 prevents decomposition in  $i$ . One solution, shown in Example 2-5, is to decompose  $x$  in  $i$  and  $y$  in  $j$ .

```
[transfer.inc]
1      integer im, jm
2      common /sizes_com/ im, jm
3
4      CSMS$DECLARE_DECOMP(DECOMP_I, 1)
5      CSMS$DECLARE_DECOMP(DECOMP_J, 1)
6

[transfer.f]
1      program TRANSFER1
2      implicit none
3
4      include 'transfer.inc'
5
6      integer i
7      integer j
8
9      im = 60
10     jm = 90
```

```

11
12 CSMS$CREATE_DECOMP(DECOMP_I, <im>, <0>)
13 CSMS$CREATE_DECOMP(DECOMP_J, <jm>, <0>)
14
15     call DO_IT
16
17     end
18
19
20
21     subroutine DO_IT
22     include 'transfer.inc'
23
24 CSMS$DISTRIBUTE(DECOMP_I, 1) BEGIN
25     real x(im,jm)
26 CSMS$DISTRIBUTE END
27
28 CSMS$DISTRIBUTE(DECOMP_J, 2) BEGIN
29     real y(im,jm)
30 CSMS$DISTRIBUTE END
31
32 C BEGIN
33
34     x = 1.0
35
36 CSMS$PARALLEL(DECOMP_I, <i>) BEGIN
37
38 C dependence in the j dimension that
39 C SMS does not provide directives to parallelize
40     do j = 2, jm
41         do i = 1, im
42             x(i,j) = x(i,j) + x(i,j-1)
43         end do
44     end do
45 CSMS$PARALLEL END
46
47 CSMS$TRANSFER(<X, Y>) BEGIN
48     do j = 1, jm
49         do i = 1,im
50             y(i,j) = x(i,j)
51         end do
52     end do
53 CSMS$TRANSFER END
54
55
56
57 CSMS$PARALLEL(DECOMP_J, <j>) BEGIN
58
59 C dependence in the i dimension that
60 C SMS does not provide directives to parallelize
61     do j = 1, jm
62         do i = 2, im
63             y(i,j) = y(i,j) + y(i-1,j)
64         end do
65     end do
66 CSMS$PARALLEL END
67
68     open(10,file='f1',form='unformatted')

```

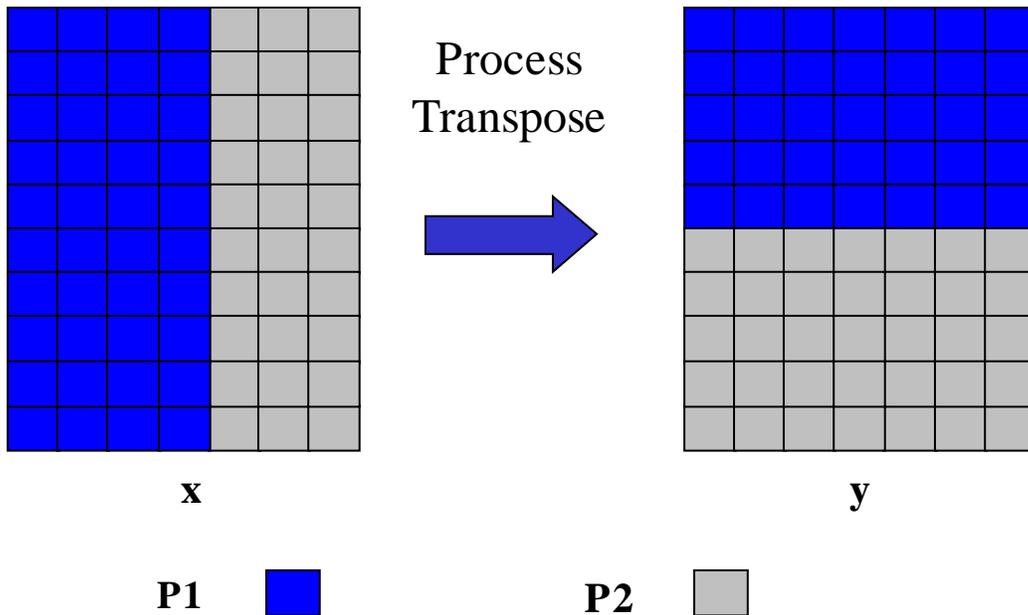
```

69      write(10) y
70      close(10)
71
72      return
73      end

```

**Example 2-5: A simple SMS parallel program that requires two data decompositions due to recurrence relations in  $i$  for array  $x$  and  $j$  for array  $y$ .**

Example 2-5 contains two `DECLARE_DECOMP` and `CREATE_DECOMP` directives. The `DISTRIBUTE` directive at line 24 uses `DECOMP_I` to decompose  $x$  in  $i$ . The `DISTRIBUTE` directive at line 28 uses `DECOMP_J` to decompose  $y$  in  $j$ . The `TRANSFER` directive at line 47 causes SMS to replace the serial code between the `BEGIN` and `END TRANSFER` directives (a simple copy) with communication that re-distributes (transposes) the data among the processes as illustrated in Figure 2-13.  $x$  is referred to as the source array of the `TRANSFER` directive and  $y$  is referred to as the destination array.



**Figure 2-13: An illustration of the data movement required between processes P1 and P2 for a transposition operation.**

### 3 Decomposing Arrays and Parallelizing Loops

#### 3.1 Choosing Decompositions

In order to choose domain decompositions that will yield optimal performance, the dependencies of arrays on one another must be analyzed. Usually, several decomposition options are possible. Decompositions of 3D arrays supported by SMS are shown in Figure 3-1. Dependence analysis is used to help pick optimal decompositions that will minimize inter-process communication. Typical FDA models will be optimally decomposed in one or both of the horizontal dimensions as illustrated by "a", "b", or "d" of Figure 3-1. All of these decompositions may be used by spectral models which are described in Section 6.2.

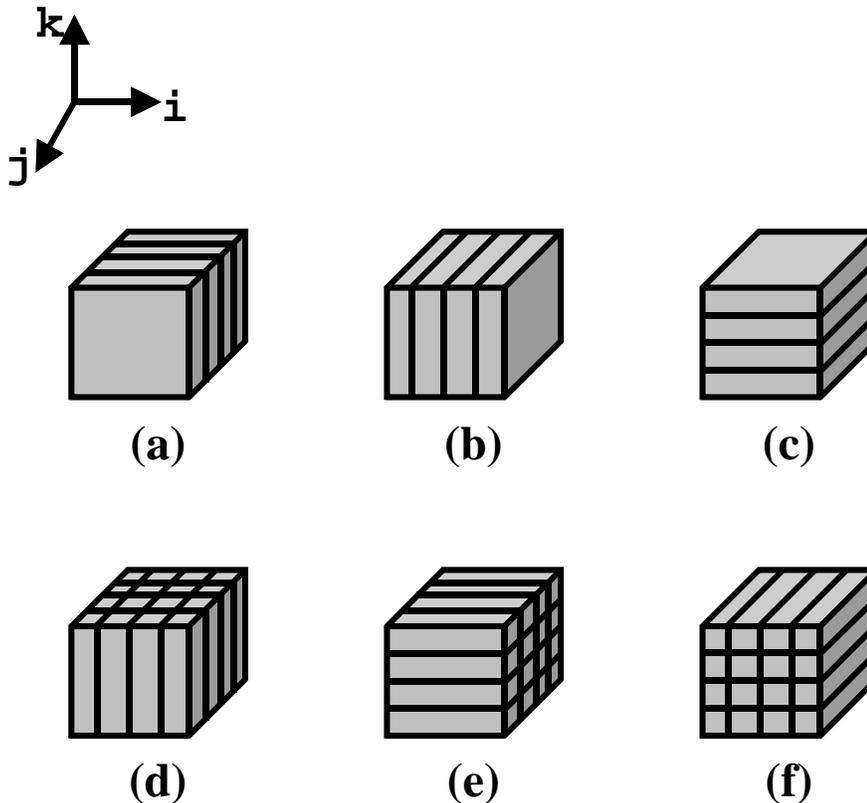


Figure 3-1: Three-dimensional decompositions supported by SMS.

Other issues to consider when selecting decompositions are the architecture of the machine on which the program will most likely be run and how many processes will be available. For vector machines, it is best to leave the inner dimension non-decomposed when possible to maximize vector lengths. On cache-based machines, it may be best to decompose the inner dimension instead. For example, in Figure 3-1, decomposition "a" would preserve long vector lengths while decomposition "b" would not. In addition,

when the number of processes available is larger than the number of grid points in the single decomposed dimension, two dimensions should be decomposed.

## 3.2 Two-Dimensional Decompositions

The full power of the DECLARE\_DECOMP, CREATE\_DECOMP, DISTRIBUTE, and PARALLEL directives becomes more apparent when two dimensions are decomposed. Consider the following example:

```
[Include file:  decomp_ex1.inc]

1      integer im, jm, km
2      common /sizes_com/ im, jm, km
3      CSMS$DECLARE_DECOMP(DECOMP_IJ, 2)

[Source file:  decomp_ex1.f]

1      program decomp_ex1
2      include 'decomp_ex1.inc'
3      im = 15
4      jm = 10
5      km = 2
6      CSMS$CREATE_DECOMP(DECOMP_IJ, <im, jm>, <0,0>)
7      call compute
8      end
9
10     subroutine compute
11     include 'decomp_ex1.inc'
12     integer i, j, k
13     CSMS$DISTRIBUTE(DECOMP_IJ, 1, 2) BEGIN
14     integer z(im,jm,km)
15     CSMS$DISTRIBUTE END
16     integer zsum
17     CSMS$PARALLEL(DECOMP_IJ,<i>,<j>) BEGIN
18     do 100 k=1,km
19     do 100 j=1,jm
20     do 100 i=1,im
21     z(i,j,k) = 1
22     100 continue
23     zsum = 0
24     do 200 k=1,km
25     do 200 j=1,jm
26     do 200 i=1,im
27     zsum = zsum + z(i,j,k)
28     200 continue
29     CSMS$PARALLEL END
30     CSMS$REDUCE(zsum, SUM)
31     print *, 'zsum = ', zsum
32     return
33     end
```

**Example 3-1: An SMS program that uses a two dimensional decomposition.**

When run, the serial version of this program prints the following:

```
>> decomp_ex1
zsum = 300
```

Directives `CREATE_DECOMP`, `DISTRIBUTE`, and `PARALLEL` now have more complex arguments than in the simple examples from Section 2.3. The second argument to `CREATE_DECOMP`, `<im, jm>`, indicates that the decomposition named `DECOMP_IJ` has two decomposed dimensions and that the global size of the first decomposed dimension is `im` and the global size of the second decomposed dimension is `jm`. The third argument, `<0,0>`, indicates that `DECOMP_IJ` has no halo regions in either decomposed dimension.

The second argument to `DISTRIBUTE`, `1`, indicates that the first dimension of array `z` is decomposed as described by the first decomposed dimension of `DECOMP_IJ`. The third argument, `2`, indicates the second dimension of array `z` is decomposed as described by the second decomposed dimension of `DECOMP_IJ`. The third dimension of array `z` will not be decomposed. This is decomposition "d" in Figure 3-1. More details about `DISTRIBUTE` can be found in Section 3.6.1.

The second argument to `PARALLEL`, `<i>`, is used to identify loop indices for loops spanning the first decomposed dimension of `DECOMP_IJ`. Similarly, the third argument, `<j>`, is used to identify loop indices for loops spanning the second decomposed dimension of `DECOMP_IJ`. The `PARALLEL` directive will translate both `i` and `j` dimensions of loops 100 and 200 to local loop bounds.

When this code is run on 2 or 3 processes, we see the expected results:

```
>> smsRun -np 2 decomp_ex1_sms
SMS: Using default process layout (2 x 1) for decomposition decomp_ij
     zsum = 300
>> smsRun -np 3 decomp_ex1_sms
SMS: Using default process layout (3 x 1) for decomposition decomp_ij
     zsum = 300
```

Note that SMS prints an additional diagnostic message for two-dimensional decompositions. This message describes how many processes are assigned to each decomposed dimension, which can be useful for debugging or performance analysis. For brevity, this message will not be shown again.

### 3.3 Controlling Process Layout

By default, SMS uses a pre-determined set of rules to decide how data points are assigned to each process and how many processes are allocated to each decomposed dimension. Roughly speaking, the data points are distributed evenly among the processes along a given decomposed dimension as seen in Figure 2-2 and Figure 2-3. Also, processes are divided among the decomposed dimensions so as to minimize the amount of data moved during an `EXCHANGE` operation.

SMS also provides a mechanism for the user to specify the assignment of data points using a process configuration file in the form of a Fortran namelist. For a given decomposition, the user defines how many processes are assigned to each decomposed dimension. SMS will then choose how many data points are assigned to each process as

before. The user can also optionally define how many data points are assigned to each process. Example 3-2 shows how process configuration works.

```
[file config.f]
1      program CONFIGURE
2      integer, parameter :: im = 12
3      integer, parameter :: jm = 12
4      CSMS$DECLARE_DECOMP(dh1, 2)
5      CSMS$DECLARE_DECOMP(dh2, 2)
6      CSMS$CREATE_DECOMP(dh1, <im,jm>, <0,0>)
7      CSMS$CREATE_DECOMP(dh2, <im,jm>, <0,0>)
8      CSMS$DISTRIBUTE(dh1, 1, 2) begin
9      real, allocatable :: a1(:, :)
8      CSMS$DISTRIBUTE end
10
11     CSMS$DISTRIBUTE(dh2, 1, 2) begin
12     real, allocatable :: a2(:, :)
13     CSMS$DISTRIBUTE end
    .
    .
    .

[file my_config]
&decomp
decomp1_name='dh1',
decomp1_nps=2 2,
decomp2_name='dh2',
decomp2_nps=4 1,
decomp2_ddim1_sizes=2 4 4 2/
```

**Example 3-2: Sample program and associated configuration file illustrating how the user can tell SMS how to distribute the data among the processes.**

The program in the example defines two decompositions. The configuration file, *my\_config*, specifies how the data are distributed among the processes. In this case, it indicates that two processes will be assigned to each decomposed dimension of decomposition *dh1*. However, the user leaves it up to SMS to determine how many data points are assigned to each process. Figure 3-2 shows the memory layout of *a1*.

For *dh2*, the user specifies that the 4 processes will be assigned to the first decomposed dimension. The user further specifies how many points are assigned to each process. Figure 3-3 shows the memory layout of *a2*.

The user tells SMS to use configuration file *my\_config* as follows:

```
>> smsRun -cf my_config my_program
SMS: For processor layout of decomposition dh1, using config file :
my_config
SMS: For processor layout of decomposition dh2, using config file :
my_config
```

Notice that the *smsRun* command does not specify how many processes are requested. SMS figures how many are needed from the configuration file. Also notice the diagnostic message from SMS that describes which decompositions are defined in the configuration file. For brevity, this message will not be shown again.

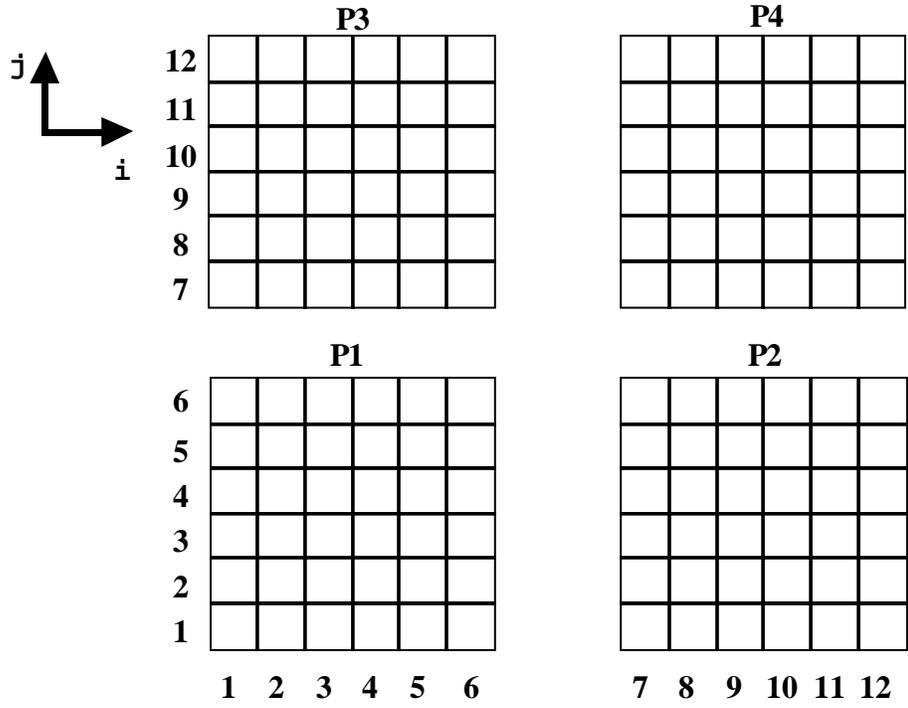


Figure 3-2: Memory layout of array *a1* in program *CONFIGURE*.

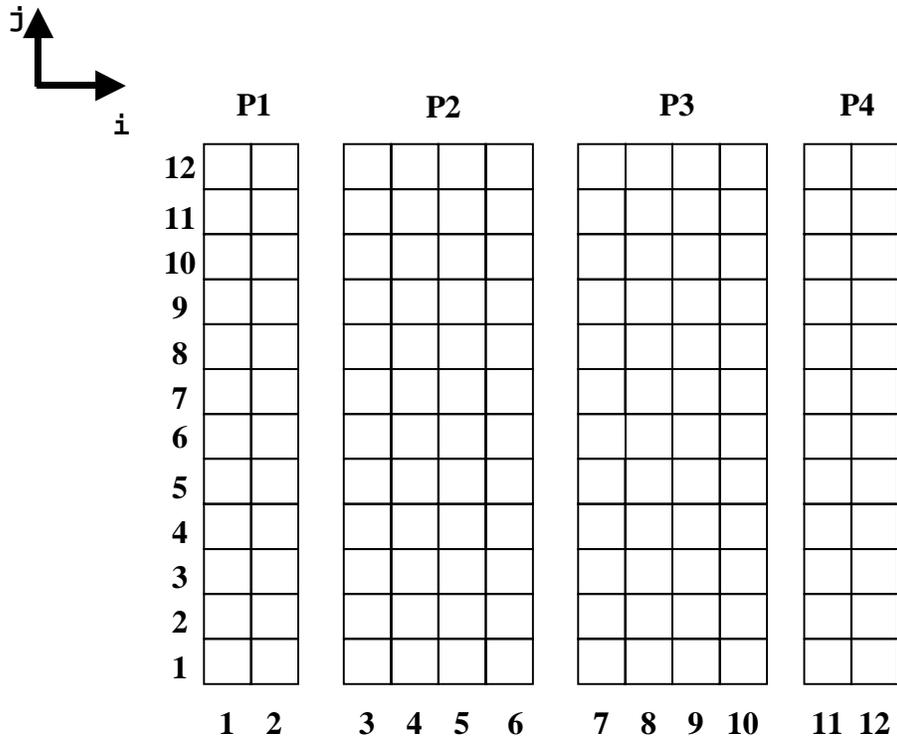


Figure 3-3: Memory layout of array *a2* in program *CONFIGURE*.

### 3.4 Using Statically Allocated Memory

When dynamic memory allocation is used, SMS automatically sets local array sizes at run-time. However, static memory codes require the local array sizes to be declared by the programmer. In addition, the local and global indices differ (Figure 3-4 below), often necessitating conversions between the two (see Section 4).

Example 3-3 illustrates a program using static memory allocation. In this example, the `DECLARE_DECOMP` directive requires a new second argument, `<(im/2)+1, jm/2>`. This informs SMS that the decomposition named `DECOMP_IJ` has two decomposed dimensions and specifies declared local sizes for each. The declared sizes will be used to translate declarations of static arrays enclosed by `DISTRIBUTE` directives. For instance, `z` will have a size of `(im/2 + 1, jm/2, km)` in the translated version of the code in Example 3-3.

```
[Include file: decomp_ex4.inc]
1      integer im, jm, km
2      parameter (im = 15, jm = 10, km = 2)
3      CSMS$DECLARE_DECOMP(DECOMP_IJ, <(im/2)+1, jm/2>)

[Source file: decomp_ex4.f]
4      program decomp_ex4
5      include 'decomp_ex4.inc'
6      CSMS$DISTRIBUTE(DECOMP_IJ, 1, 2) BEGIN
7      integer z(im,jm,km)
8      CSMS$DISTRIBUTE END
9      integer zsum, i, j, k
10     CSMS$CREATE_DECOMP(DECOMP_IJ, <im, jm>, <0,0>)
11     CSMS$PARALLEL(DECOMP_IJ,<i>,<j>) BEGIN
12         do 100 k=1,km
13             do 100 j=1,jm
14                 do 100 i=1,im
15                     z(i,j,k) = 1
16             100 continue
17             zsum = 0
18             do 200 k=1,km
19                 do 200 j=1,jm
20                     do 200 i=1,im
21                         zsum = zsum + z(i,j,k)
22             200 continue
23     CSMS$PARALLEL END
24     CSMS$REDUCE(zsum, SUM)
25     print *, 'zsum = ', zsum
26     end
```

**Example 3-3:** An SMS program that uses static memory allocation requires the local sizes be declared in the `DECLARE_DECOMP` directive. In this example, these local sizes are: `(im/2)+1` and `jm/2`.

In static memory cases, where the number of processes assigned to a decomposed dimension does not evenly divide the global size of that dimension, the declared local sizes specified in the `DECLARE_DECOMP` directive must be set for the process(es) that use(s) the most memory. For a 4-process run, the term `(im/2)+1` (Example 3-3, line 3) ensures there will be sufficient local memory for all processes even though two require



In addition to wasting memory, performance of the 8-process run might not be optimal on a cache-based machine because the data used in each array are scattered over a block of memory twice the needed size. This will likely result in more cache misses and may significantly degrade performance. Further, this effect becomes more severe as the number of processes increases. For example, if the program were run on 32 processes, the largest local array sizes required on any process would be only  $(2, 3, 2)$ . Therefore, it is especially important to declare arrays using the smallest possible sizes for large numbers of processes.

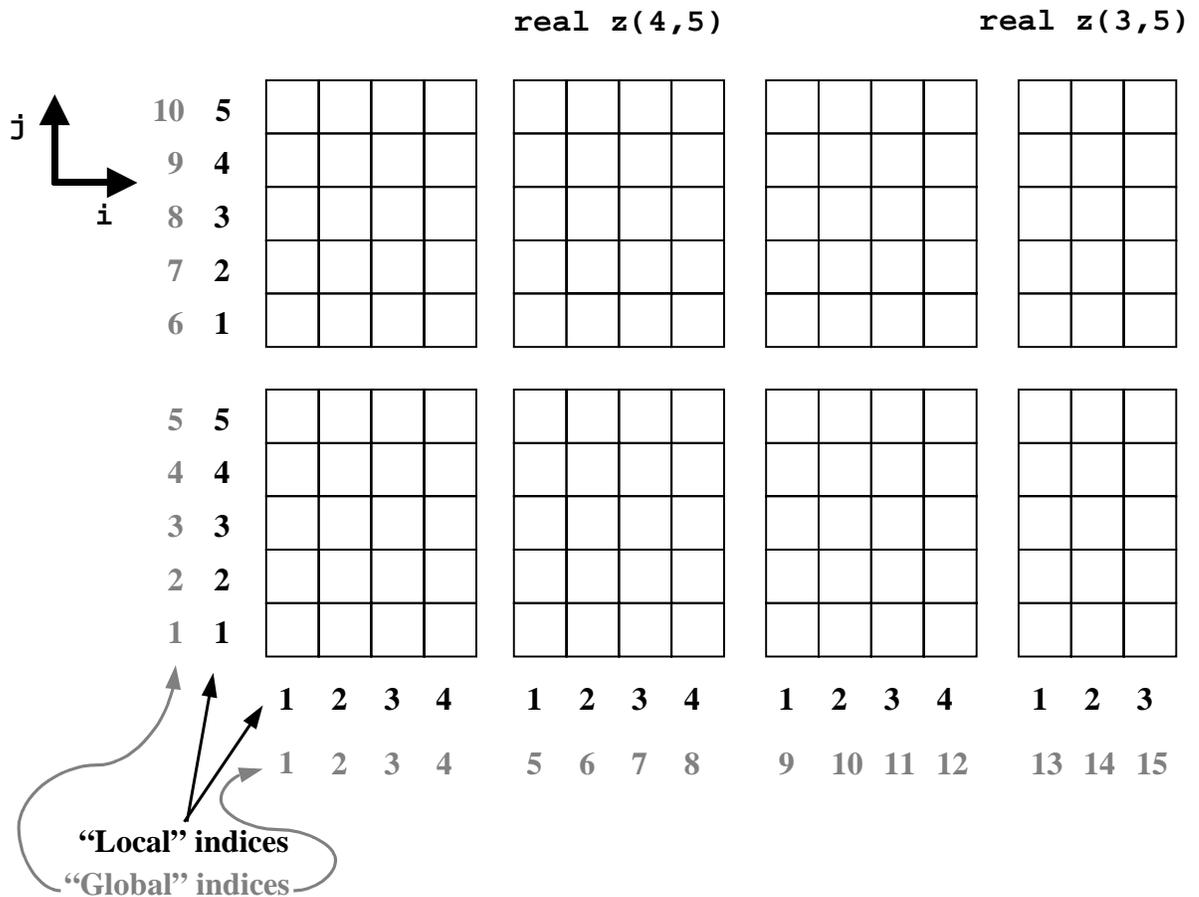


Figure 3-5: Memory layout for 8-process run.

Running Example 3-3 on 2 processes produces the following:

```
>> smsRun -np 2 decomp_ex4_sms
Process: 1 Error at: ./decomp_ex4_sms.f:10.1
Process: 1 Error status= -2202 : USER DECLARED STATIC ARRAY IS TOO
SMALL.
Process: 1 Aborting...
```

What happened? By default, the two processes are distributed along the  $i$  dimension so the largest local array sizes required on any process for the two-process run is  $(8, 10, 2)$ . However, the `DECLARE_DECOMP` directive set local array sizes to  $((im/2)+1, jm/2, km) = (8, 5, 2)$  which is too small for the two process run (see Figure 3-4). SMS detects this error at run time, prints the error messages, and aborts the program.

To provide sufficient memory for the local arrays in a two-process run, we can modify the sizes in the `DECLARE_DECOMP` directive as follows:

```
CSMS$DECLARE_DECOMP(DECOMP_IJ, <(im/2)+1, jm>)
```

If the following `DECLARE_DECOMP` directive were used

```
CSMS$DECLARE_DECOMP(DECOMP_IJ, <im, jm>)
```

all translated arrays would be declared full-size. This code could then be run on any number of processes (provided each process has enough memory). This is very useful during debugging because it allows comparison of results for runs made on different numbers of processes. Once debugging is complete, the `DECLARE_DECOMP` directives should be changed to minimize memory use.

Determining the proper local sizes for static memory models that need `EXCHANGES` will be discussed in Section 5.1.3.

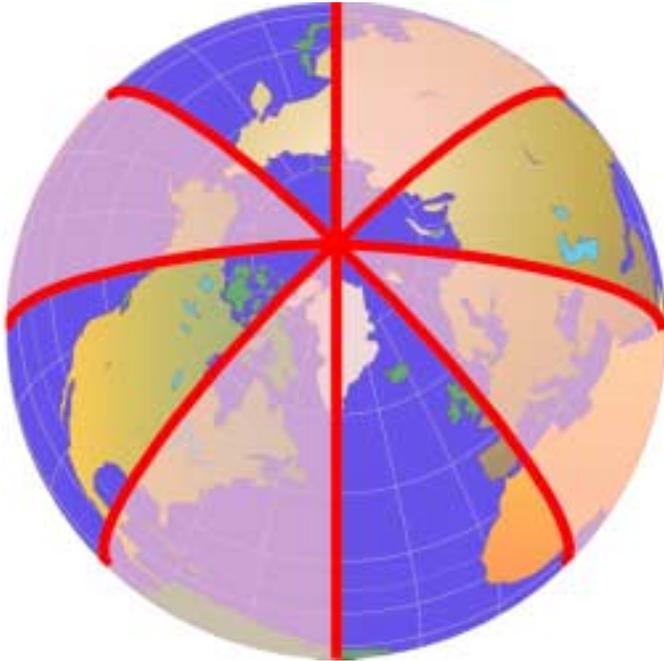
### 3.5 Load Balancing

Ideally, each process will have exactly the same amount of work to do. In practice, most models have computations that vary spatially so some processes may have more work than others. This is commonly known as load imbalance. Load imbalances slow down a parallel program because some processes with less work are forced to wait for processes with more work to catch up.

Load imbalances can be fixed (static) or variable (dynamic) in time. Several types of static load imbalances are found in weather and ocean models. One type is due to the fact that edge processes have boundary condition calculations while interior processes do not. We saw in Section 3.3 how the user can specify the number of grid points allocated to each process. By giving more points to interior processes, the imbalance is mitigated.

Global atmospheric models can have longitudinal imbalances due to differences in the computations required for day/night points and latitudinal imbalances stemming from summer/winter computational differences. SMS provides means to mitigate both of these effects using a technique called index scrambling. Index scrambling, accessible via an option to the `CREATE_DECOMP` directive, moves adjacent row/column pairs to other processes. Figure 3-6 shows how the data might be distributed following longitudinal scrambling. Due to the complexity of these redistributions, they are not permitted for

models that have adjacent dependencies (because the EXCHANGEs would perform poorly).



**Figure 3-6: Longitude scrambling is used to reduce load imbalances due to computational differences stemming from day/night cycles in a global weather model. In this case, the model is run using 2 processes. One process has the brightly colored segments; the other has the darker colored segments. The effect is to give each process half the day-time points and half the night-time points.**

To use index scrambling, a fourth argument is added to the CREATE\_DECOMP as shown in the code fragments below:

```
CSMS$CREATE_DECOMP(DECOMP_J, <jm>, <0> : <SCRAMBLE_LAT_STRATEGY>)
```

```
CSMS$CREATE_DECOMP(DECOMP_I, <im>, <0> : <SCRAMBLE_LON_STRATEGY>)
```

In the first case, argument **<SCRAMBLE\_LAT\_STRATEGY>** indicates that the first decomposed dimension of *DECOMP\_J* will be scrambled using a method appropriate for balancing load among latitudes in a global model. In the second case, argument **<SCRAMBLE\_LON\_STRATEGY>** indicates that the first decomposed dimension of *DECOMP\_I* will be scrambled using a method appropriate for balancing load among longitudes in a global model. No other code changes are required to use the scrambling feature. For this reason, it is convenient to add this feature as a performance optimization once debugging of the non-scrambled parallel code is complete.

A future SMS release will provide means to mitigate dynamic load imbalances.

## 3.6 More about DISTRIBUTE

### 3.6.1 Further Detail on DISTRIBUTE Syntax

This section explains the distinction between “dimension of an array” and “dimension of an SMS decomposition”. The DISTRIBUTE directive can decompose several types of arrays as shown in the following code fragments:

```
CSMS$DISTRIBUTE(DECOMP_IJ, 1, 2) BEGIN
    integer x(im,jm,km)
CSMS$DISTRIBUTE END
```

Here, the first dimension of array  $x$  is decomposed as described by the first decomposed dimension of *DECOMP\_IJ* and the second dimension of array  $x$  is decomposed as described by the second decomposed dimension of *DECOMP\_IJ*. The third dimension of array  $x$  is not decomposed.

```
CSMS$DISTRIBUTE(DECOMP_IJ, 1, 3) BEGIN
    real a(im,km,jm)
CSMS$DISTRIBUTE END
```

The numbers *1* and *3* refer to array dimensions. The order in which they appear determines the dimensions of the decomposition to which they refer. Here, the first dimension of array  $a$  is decomposed as described by the first decomposed dimension of *DECOMP\_IJ* and the third dimension of array  $a$  is decomposed as described by the second decomposed dimension of *DECOMP\_IJ*. The second dimension of array  $a$  is not decomposed.

```
CSMS$DISTRIBUTE(DECOMP_IJ, 3, 2) BEGIN
    real b(km,jm,im), avg
CSMS$DISTRIBUTE END
```

Here, the third dimension of array  $b$  is decomposed as described by the first decomposed dimension of *DECOMP\_IJ* and the second dimension of array  $b$  is decomposed as described by the second decomposed dimension of *DECOMP\_IJ*. The first dimension of array  $b$  is not decomposed. *avg* is not decomposed at all because it is a scalar variable.

The user can also specify how variables are distributed by using variable name tags instead of dimension numbers. For example,

```
CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
    real x(im, jm, km)
CSMS$DISTRIBUTE END
```

again indicates the first dimension of array  $x$  is decomposed as described by the first decomposed dimension of *DECOMP\_IJ* and the second dimension of array  $x$  is decomposed as described by the second decomposed dimension of *DECOMP\_IJ*.

```

CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
    real b(km, jm, im)
CSMS$DISTRIBUTE END

```

indicates that the third dimension of  $b$  is distributed based on the first decomposed dimension and the second dimension of  $b$  is distributed based on the second decomposed dimension.

Using this syntax, it is possible to enclose the last two arrays inside the same distribute directive:

```

CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
    real x(im, jm, km)
    real b(km, jm, im)
CSMS$DISTRIBUTE END

```

### 3.6.2 Using DISTRIBUTE to Define Decomposed Boundary Arrays

Regional weather forecast models and ocean models often require boundary condition data. A code segment handling western boundary conditions might look as shown in Example 3-4.

```

1      subroutine UPDATE_BOUNDARIES(u)
2      integer, parameter :: im = 10
3      integer, parameter :: jm = 20
4      integer, parameter :: km = 30
5      csms$declare_decomp(dh, 2)
6
7      csms$distribute(dh, 1, 2) begin
8      real u(im, jm, km)
9      csms$distribute end
10
11     csms$distribute(dh, , 1) begin
12     real ubw(jm, km)
13     csms$distribute end
14
15     open(10, file='west_bdy', form='unformatted')
16     read(10) ubw
17     close(10)
18
19     csms$parallel(dh, , <j>) begin
20     csms$global_index(1) begin
21     do k = 1, km
22     do j = 1, jm
23     u(1, j, k) = (u(1, j, k) + ubw(j,k))/2.0
24     end do
25     end do
26     csms$global_index end
27     csms$parallel end
28
29     return
30     end

```

**Example 3-4:** Subroutine showing how boundary condition arrays can be handled in SMS.

The second DISTRIBUTE statement on line 11 defines an unusual kind of decomposed array. Its first dimension is decomposed according to the second dimension of

decomposition  $dh$  but none of the array dimensions are decomposed based on the first decomposed dimension. Since the exact manner in which this array is distributed is somewhat ill-defined, it can only be read in from disk. REDUCE, SERIAL, TRANSFER, and COMPARE\_VAR (discussed later in Section 13.1) cannot be applied to this type of array. Output is also not allowed. The PARALLEL directive syntax on line 19 will be discussed in Section 3.7. Note that for a dynamic memory code like Example 3-4, it is possible to avoid these problems by simply leaving the ill-defined array non-decomposed. In this case, removing the DISTRIBUTE statement on line 11 would be sufficient. The only drawback to this solution is a possible performance degradation due to wasted memory (each process will only use a small portion of globally-sized array  $ubw$ ). This solution will not work for a static memory code unless TO\_GLOBAL directives are used to handle index translation (see Section 4.1).

### 3.7 More About PARALLEL

The PARALLEL directive will translate serial loops correctly provided the upper and lower loop bounds are valid global indices. For example, the  $i$  and  $j$  loops below would all be correctly translated:

```

CSMS$PARALLEL(DECOMP_IJ,<i>,<j>) BEGIN
  do 100 k=1,km
    do 200 j=3,jm-2
      do 200 i=3,im-2
        z(i,j,k) = x(i,j,k) + y(i,j,k)
200  continue

      do 210 j=1,2
        do 210 i=1,im
          z(i,j,k) = 0
210  continue

      do 220 j=jm-1,jm
        do 220 i=1,im
          z(i,j,k) = 0
220  continue

      do 230 j=1,jm
        do 230 i=1,2
          z(i,j,k) = 0
230  continue

      do 240 j=1,jm
        do 240 i=im-1,im
          z(i,j,k) = 0
240  continue

100  continue
CSMS$PARALLEL END

```

In this code fragment, notice that the translated version of loop 210 would only be executed on processes that contain global indices  $j=1$  or  $j=2$ . The PARALLEL directive ensures that other processes will skip loop 210. Similar translations will occur for the other loops.

Recall the syntax seen on lines 19-22 of Example 3-4.

```
CSMS$PARALLEL(dh, , <j>) BEGIN
  do k = 1, km
    do j = 1, jm
```

It indicates that no enclosed loops correspond to the first decomposed dimension but any loops that use index  $j$  correspond to the second decomposed dimension and should be translated.

There is no run-time performance penalty for using a PARALLEL directive because processes are not synchronized. Also, PARALLEL directives may enclose any valid Fortran executable statements. Therefore, a program that uses only one decomposition will usually require no more than one BEGIN-END pair of PARALLEL directives for each program unit (subroutine, function, or main program).

In tagging loop indices to be translated, some care is required. First, indices can sometimes be used for non-decomposed loops as well as for loops that span decomposed dimensions. This is the case in the following fragment:

```
CSMS$PARALLEL(DECOMP_IJ,<i>,<j>) BEGIN
  do 200 k=1,km
    do 200 j=1,jm
      do 200 i=1,im
        z(i,j,k) = x(i,j,k) + y(i,j,k)
200  continue
      do 500 i=1,3
        call smooth(z)
500  continue
CSMS$PARALLEL END
```

In this case, loop 500 is used to repeatedly call subroutine *smooth* which performs some computations on decomposed array  $z$ . This loop should NOT be translated because  $i$  is being used as an iteration count, not as an index into a decomposed dimension. This is easily fixed either by using a different loop index in loop 500, by moving the PARALLEL END directive to exclude loop 500, or by placing the loop inside an IGNORE directive as shown below in order to tell SMS not to translate the loop

```
CSMS$IGNORE BEGIN
  do 500 i = 1, 3
    call smooth(z)
500  continue
CSMS$IGNORE END
```

Second, make sure that all loops manipulating decomposed arrays are enclosed inside PARALLEL directives. During translation, PPP will generate a warning message whenever it finds a loop that is not enclosed by PARALLEL directives if that loop contains a decomposed array:

**This variable, decomposed by CSMS\$DISTRIBUTE, is being used outside of a parallel region.**

### 3.8 Arrays with Non-Unit Lower Bounds

When arrays in the serial code are declared with non-unit lower bounds, the SMS decomposition must reflect this fact. Consider the following variant of Example 3-1:

```
[Include file:  decomp_ex6.inc]

    integer im, jm, km
    common /sizes_com/ im, jm, km
CSMS$DECLARE_DECOMP(DECOMP_IJ : <0,0>)

[Source file:  decomp_ex6.f]

    program decomp_ex6
    include 'decomp_ex6.inc'
    im = 15
    jm = 10
    km = 2
CSMS$CREATE_DECOMP(DECOMP_IJ, <im, jm>, <0,0>)
    call compute
    end

    subroutine compute
    include 'decomp_ex6.inc'
    integer i, j, k
CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
    integer z(0:im-1,0:jm-1,0:km-1), zsum
CSMS$DISTRIBUTE END
CSMS$PARALLEL(DECOMP_IJ,<i>,<j>) BEGIN
    do 100 k=0,km-1
    do 100 j=0,jm-1
    do 100 i=0,im-1
        z(i,j,k) = 1
100 continue
    zsum = 0

    do 200 k=0,km-1
    do 200 j=0,jm-1
    do 200 i=0,im-1
        zsum = zsum + z(i,j,k)
200 continue

CSMS$PARALLEL END
CSMS$REDUCE(zsum, SUM)
    print *, 'zsum = ', zsum
    return
    end
```

In this program, array  $z$  is declared so the first index (lower bound) is zero in each dimension instead of the Fortran default of one. The bounds of loops 100 and 200 now start at zero. The only difference between the directives in this example and those in Example 3-1 is `DECLARE_DECOMP`. The new final argument, `<0,0>` indicates that array declarations have a lower bound of zero in both decomposed dimensions.

## 4 Translating Array Indices

We have seen that, in static memory models, local and global indices are different whenever more than one process is used so conversions between them are required. TO\_GLOBAL (Section 4.1) and TO\_LOCAL (Section 4.2) provide support for these conversions. In addition, it is sometimes desirable to generate process-local loop start and end indices and array sizes to simplify parallelization of subroutines in both dynamic and static memory codes. The TO\_LOCAL (Section 4.3) directive does this as well. Finally, boundary condition calculations must be restricted to processes containing boundary points. GLOBAL\_INDEX (Section 4.4) handles these cases.

### 4.1 Translating Local Indices to Global Indices

For a static memory code, when a loop has been translated using the PARALLEL directive, the value of the index is now process-local as illustrated in Figure 3-4. If the intent of the program is to access the global value, this index will need to be translated back to a global value. This point is illustrated in Example 4-1.

```
1      program tran_index1
2      implicit none
3      integer i, j
4      integer, parameter :: im = 5
5      integer, parameter :: jm = 3
6  CSMS$DECLARE_DECOMP(DECOMP_IJ, <im, jm>)
7
8  CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
9      integer x(im,jm)
10 CSMS$DISTRIBUTE END
11
12 CSMS$CREATE_DECOMP(DECOMP_IJ, <im, jm>, <0,0>)
13
14 CSMS$PARALLEL(DECOMP_IJ,<i>,<j>) BEGIN
15
16      do 100 j=1,jm
17      do 100 i=1,im
18          x(i,j) = (100 * i) + j
19  100 continue
20
21 CSMS$SERIAL BEGIN
22      do j = 1, jm
23          write(*,'(16i5)') (x(i,j),i=1,im)
24      end do
25 CSMS$SERIAL END
26 CSMS$PARALLEL END
27
28      end
```

**Example 4-1:** An SMS parallel program that incorrectly initializes the array  $x$  inside subroutine *compute*.

This program initializes array  $x$  in loop 100 (lines 16-19). Each element of array  $x$  is then printed on the screen. When the serial code is run, the following is printed on the screen:

```
>> tran_index1
 101  201  301  401  501
 102  202  302  402  502
 103  203  303  403  503
```

The same result is seen when the SMS version is run on one process. However, the results are incorrect when two processes are used:

```
>> smsRun -np 2 tran_index1_sms
 101  201  301  101  201
 102  202  302  102  202
 103  203  303  103  203
```

Why are the results incorrect? The PARALLEL directive has translated the  $i$  and  $j$  indices used to compute  $x$  in loop 100 using local indices. However, correct operation requires that  $x$  be initialized using global indices as in the original serial code. The solution is to use the TO\_GLOBAL directive to translate the local indices to global indices. In this case, line 18 would be replaced with the following code:

```
CSMS$TO_GLOBAL(<1,i>, <2,j>) BEGIN
      x(i,j) = (100 * i) + j
CSMS$TO_GLOBAL END
```

The first argument in the TO\_GLOBAL directive,  $\langle 1, i \rangle$ , indicates that array index  $i$  is an index in the first decomposed dimension. The second argument,  $\langle 2, j \rangle$ , indicates that array index  $j$  is an index in the second decomposed dimension. All occurrences of indices  $i$  and  $j$  inside the TO\_GLOBAL directives that are not array references will be converted to their global equivalents in the first and second decomposed dimensions, respectively. Running the new parallel code on various numbers of processes will now yield the same result as the serial run.

Note that the TO\_GLOBAL directive must appear within a PARALLEL directive. Directives TO\_LOCAL and GLOBAL\_INDEX, introduced later in this section, also have this restriction. Also note that since  $x$  is decomposed, the SERIAL directive is required to handle the write statement on line 24 as will be explained in Section 8.

## 4.2 Translating Global Indices to Local Indices Inside Loops

Sometimes, indices that have been translated to global values need to be translated back to local values to be used as indices into decomposed arrays in a static memory code. The TO\_LOCAL directive is used for this translation. Consider the following code fragment that uses computed indices to avoid out-of-bounds references:

```
CSMS$PARALLEL(DECOMP_IJ, <i>, <j>) BEGIN
  do j=1,jm
    do i=1,im
CSMS$TO_GLOBAL(<1,i>) BEGIN
CSMS$TO_LOCAL(<1,im1,ip1>) BEGIN
      im1 = max( 1,i-1)
      ip1 = min(im,i+1)
CSMS$TO_LOCAL END
CSMS$TO_GLOBAL END
```

```

        x(i,j) = y(i,j) - y(im1,j) - y(ip1,j)
    end do
end do
CSMS$PARALLEL END

```

The max and min functions compare index  $i$  with global index values  $1$  and  $im$ . Therefore, the TO\_GLOBAL directive must be used. The TO\_GLOBAL directive will convert  $i-1$  and  $i+1$  to global values so  $ip1$  and  $im1$  will be computed as global indices. However,  $ip1$  and  $im1$  are then used as indices into decomposed array  $x$ , so they must be converted back from global to local values to avoid out-of-bounds array references for multi-process runs. The TO\_LOCAL directive shown accomplishes this. The first argument in the TO\_LOCAL directive,  $\langle 1, im1, ip1 \rangle$ , indicates that array indices  $im1$  and  $ip1$  are both used in loops that span the first decomposed dimension. In this example, occurrences of either index in code enclosed by the TO\_LOCAL directives that are not array references will be converted to their local equivalents in the first decomposed dimension.

### 4.3 Using TO\_LOCAL to Generate Process-Local Sizes and Loop Bounds

In many models, large sections of code contain no dependencies that require communications (typically weather model physics routines). If the array bounds and loop limits are passed into these routines, SMS provides a means to parallelize them without inserting directives into the code. Example 4-2 shows such a case.

```

1      program AVOID_DIRECTIVES
2      implicit none
3      integer i
4      integer, parameter :: im = 8
5
6      CSMS$DECLARE_DECOMP(dh, 1)
7
8      integer istart, iend
9
10     CSMS$DISTRIBUTE(dh, 1) BEGIN
11         integer, allocatable :: x(:), y(:)
12     CSMS$DISTRIBUTE END
13
14     CSMS$CREATE_DECOMP(dh, <im>, <2>)
15
16         allocate(x(im))
17         allocate(y(im))
18         x = 0.0
19
20     CSMS$PARALLEL(dh, <i>) BEGIN
21
22         do i=1,im
23     CSMS$TO_GLOBAL(<1,i>) BEGIN
24         x(i) = i
25     CSMS$TO_GLOBAL END
26         end do
27
28         y = 0.0
29
30
31     csms$to_local(<1, istart : lbound>, <1, iend : ubound>) begin

```

```

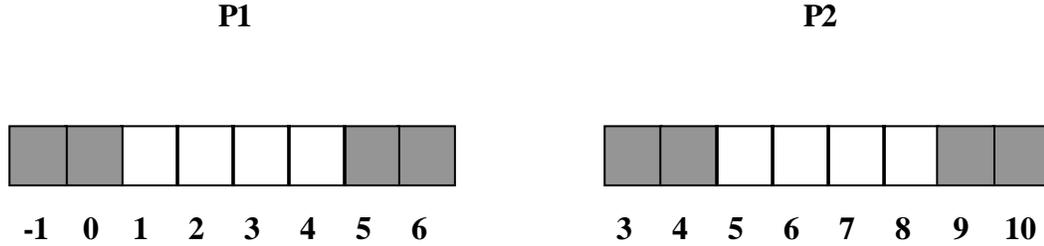
32     istart = 1
33     iend   = im - 1
34 csms$to_local end
35
36     call physics(x, lbound(x,1), ubound(x,1), istart, iend, y)
37
38 CSMS$SERIAL BEGIN
39     write(*,'(8i5)') (y(i),i=1,im)
40 CSMS$SERIAL END
41
42 CSMS$PARALLEL END
43     end
44
45     subroutine physics(arr_in, i_mem_start, i_mem_end,
46 &                      dim1_start, dim1_end,
47 &                      arr_out)
48     implicit none
49     integer i_mem_start, i_mem_end
50     integer arr_in(i_mem_start:i_mem_end)
51     integer dim1_start, dim1_end
52     integer arr_out(i_mem_start:i_mem_end)
53
54     integer i
55     do i = dim1_start, dim1_end
56         arr_out(i) = 2.0*arr_in(i)
57     end do
58     return
59     end

```

**Example 4-2: Sample code that shows how TO\_LOCAL can be used to pass local array bounds and start/end loop limits to subroutines so that no directives need to be added to them.**

Program *AVOID\_DIRECTIVES* calls subroutine *physics* (line 36), passing the arrays *x* and *y*, the starting and ending addresses of those arrays, and the starting and ending loop limits (*istart*, *iend*) over which the loops in *physics* will span. The TO\_LOCAL directive at lines 31-34 converts the dimensions and loop limits to their process local values. During source code translation, the syntax **<1, istart : lbound>** causes replacement of instances of *istart* with the local index of the first interior point for the first decomposed dimension for the given process. Figure 4-1 shows all the sizes and bounds for this case, assuming the program is run on 2 processes.

The result is that, inside subroutine *physics*, *i\_mem\_start*, *i\_mem\_end*, *dim1\_start*, *dim1\_end*, *dim2\_start*, and *dim2\_end* all have the correct process-local values. Consequently, subroutine *physics* produces the right answer for any process decomposition, even though it contains no SMS directives.



Process	<code>i_mem_start</code>	<code>i_mem_end</code>	<code>dim1_start</code>	<code>dim1_end</code>
P1	-1	6	1	4
P2	3	10	5	8

**Figure 4-1: Process layout and memory bounds of arrays `arr_in` and `arr_out` and loop bounds for a 2-process run of Example 4-2.**

## 4.4 Using Global Index to Handle Boundary Conditions

Consider the following code fragment that is enclosed in a `PARALLEL` directive but is not inside a loop:

```

id = 5
jd = 4
x(id,jd) = 10

```

The following use of `TO_LOCAL` would be incorrect:

```

CSMS$TO_LOCAL(<1,id>,<2,jd>) BEGIN
  id = 5
  jd = 4
CSMS$TO_LOCAL END
  x(id,jd) = 10

```

The translation of `id` and `jd` from global values to process-local values will work fine on the process that "owns" global point  $(5, 4)$ . However, the translation will be erroneous on processes that do not own global point  $(5, 4)$  because there is no valid local equivalent of these global coordinates on those processes. In order to restrict the execution of these statements to the process that owns the data, the `GLOBAL_INDEX` directive must be used as shown below:

```

  id = 5
  jd = 4
CSMS$GLOBAL_INDEX(1,2) BEGIN
  x(id,jd) = 10
CSMS$GLOBAL_INDEX END

```

The GLOBAL\_INDEX directive ensures that the execution of the enclosed assignment statement will only be permitted on the process that owns the global point ( $id, jd$ ). In addition, if index translation is needed,  $id$  and  $jd$  will be translated to process-local equivalents. The first argument in the directive, **1**, indicates that all array indices corresponding to the first decomposed dimension are affected. The second argument, **2**, indicates that all array indices corresponding to the second decomposed dimension are affected.

Consider the following example that initializes the boundaries of an array that is decomposed in two dimensions:

```

1      subroutine compute(im,jm)
2      integer im,jm
3  CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
4      integer x(im,jm)
5  CSMS$DISTRIBUTE END
6      integer i, j
7  CSMS$PARALLEL(DECOMP_IJ,<i>,<j>) BEGIN
8      do 100 j=1,jm
9      do 100 i=1,im
10     x(i,j) = (100 * i) + j
11  100 continue
12     do 110 j=2,jm-1
13  CSMS$GLOBAL_INDEX(1) BEGIN
14     x( 1,j) = 0
15     x(im,j) = 0
16  CSMS$GLOBAL_INDEX END
17  110 continue
18     do 120 i=2,im-1
19  CSMS$GLOBAL_INDEX(2) BEGIN
20     x(i, 1) = 0
21     x(i,jm) = 0
22  CSMS$GLOBAL_INDEX END
23  120 continue
24  CSMS$GLOBAL_INDEX(1,2) BEGIN
25     x( 1, 1) = 0
26     x(im, 1) = 0
27     x( 1,jm) = 0
28     x(im,jm) = 0
29  CSMS$GLOBAL_INDEX END
30
31  CSMS$PARALLEL END
32
33  CSMS$SERIAL BEGIN
34     print *, 'ARRAY x:'
35     print *, x
36  CSMS$SERIAL END
37     return
38     end

```

**Example 4-3: An SMS subroutine that illustrates the use of the GLOBAL\_INDEX directive to initialize array boundaries.**

This subroutine initializes array  $x$  as in previous examples (lines 8-11). It is assumed this is a dynamic memory code so TO\_GLOBAL is not required. It then proceeds to set the boundary values of  $x$  to zero in lines 12 through 28. Three pairs of GLOBAL\_INDEX directives handle the necessary translations. The first pair deals with global indices  $1$  and  $im$  in loop 110 while the second pair deals with global indices  $1$  and  $jm$  in loop 120.

The third pair handles global indices in the four assignment statements on lines 25 through 28. In each case, the enclosed statements are only executed on the appropriate processes. The SERIAL directive on line 33 will be discussed in Section 8.

When the serial and parallel codes are run, the following is printed on the screen (assuming values of *im* and *jm* as in previous examples):

```
ARRAY x:
  0   0   0   0   0
  0 202 302 402  0
  0   0   0   0   0
```

## 5 Handling Adjacent Dependencies

### 5.1 Further Details on EXCHANGE

In Section 2.5, we saw how the EXCHANGE directive was used to implement communications needed to resolve adjacent dependencies for a simple one-dimensional decomposition. In this sub-section, we expand on that discussion by examining the treatment of two-dimensional decompositions and larger stencils, and by discussing other miscellaneous details about EXCHANGE.

#### 5.1.1 Using EXCHANGE in the Case of Two-Dimensional Decompositions

We begin by modifying the Laplace program (Example 2-4) introduced in Section 2.5 so that a two dimensional decomposition is used. Two-dimensional data decompositions allow parallel programs to scale to a large number of processes.

```
1      program basic_ex_2d_decomp
2      include 'basic.inc'
3      im = 10
4      jm = 10
5      CSMS$CREATE_DECOMP(DECOMP_I, <im,jm>, <1,1>)
6      call laplace
7      end
8
9      subroutine laplace
10     include 'basic.inc'
11     integer i, j, iter
12     real max_error
13     real tolerance
14     parameter (tolerance = 0.001)
15     CSMS$DISTRIBUTE(DECOMP_I, 1, 2) BEGIN
16     real f(im,jm), df(im,jm)
17     CSMS$DISTRIBUTE END
18     CSMS$PARALLEL(DECOMP_I,<i>, <j>) BEGIN
19     do 100 j=1,jm
20     do 100 i=1,im
21     f(i,j) = 0.0
22     100 continue
23     do 110 j=1,jm
24     CSMS$GLOBAL_INDEX(1) BEGIN
25     f( 1,j) = 2.0
26     f(im,j) = 2.0
27     CSMS$GLOBAL_INDEX END
28     110 continue
29     do 120 i=1,im
30     CSMS$GLOBAL_INDEX(2) BEGIN
31     f(i, 1) = 2.0
32     f(i,jm) = 2.0
33     CSMS$GLOBAL_INDEX END
34     120 continue
35     iter = 0
36     max_error = 2.0 * tolerance
37     C main iteration loop...
38     do while ((max_error .gt. tolerance) .and. (iter .lt. 1000))
39     iter = iter + 1
```

```

40         max_error = 0.0
41 CSMS$EXCHANGE(f)
42         do 200 j=2,jm-1
43           do 200 i=2,im-1
44             df(i,j) = 0.25*(f(i-1,j) + f(i+1,j) + f(i,j-1) +
45             &                f(i,j+1)) - f(i,j)
46         200 continue
47           do 300 j=2,jm-1
48             do 300 i=2,im-1
49               if (max_error .lt. abs(df(i,j))) then
50                 max_error = abs(df(i,j))
51               endif
52             300 continue
53 CSMS$REDUCE(max_error, MAX)
54           do 400 j=2,jm-1
55             do 400 i=2,im-1
56               f(i,j) = f(i,j) + df(i,j)
57           400 continue
58         enddo
59
60 CSMS$PARALLEL END
61         print *, 'Solution required ',iter,' iterations'
62         print *, 'Final error = ', max_error
63
64         return
65         end

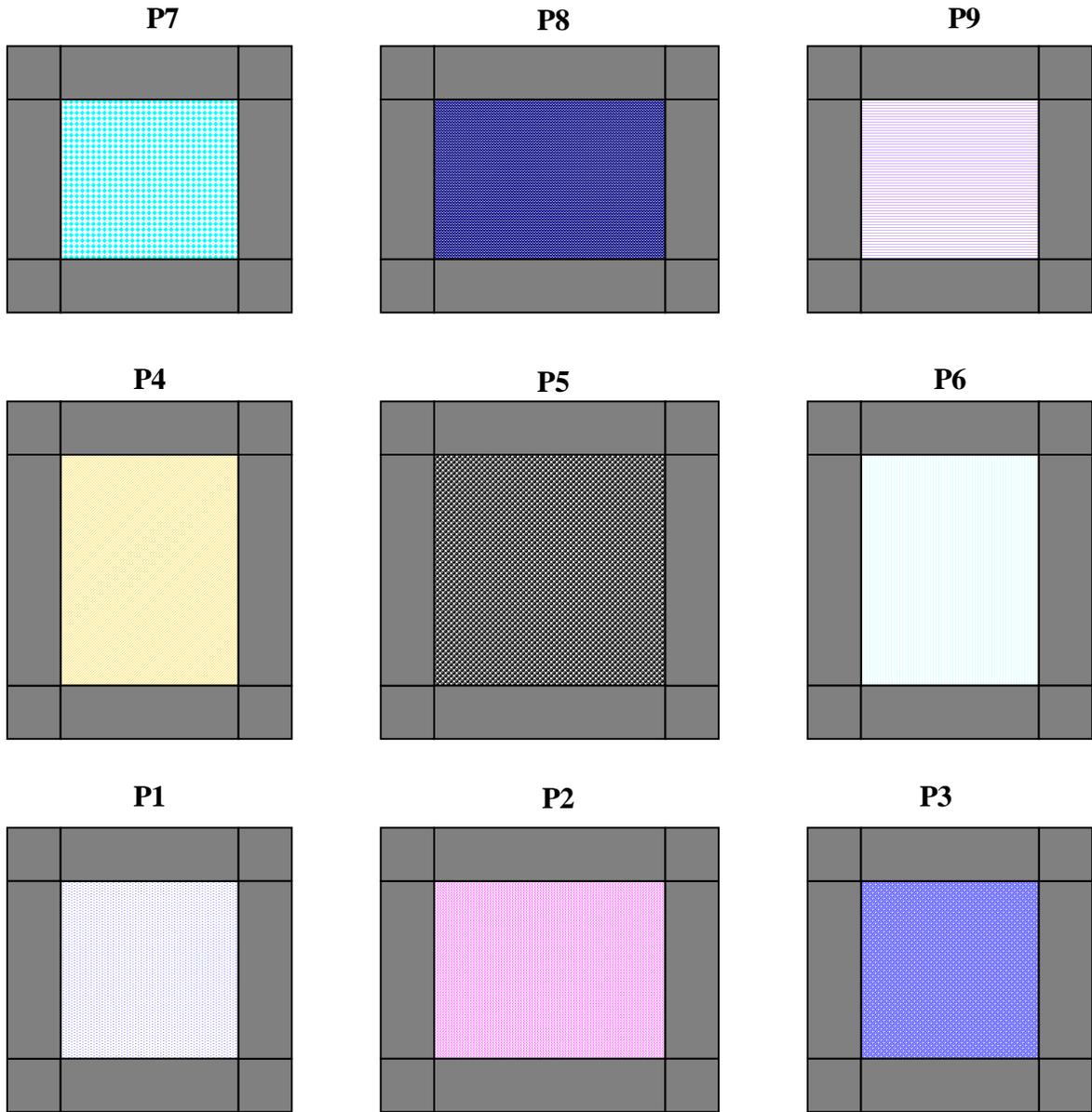
```

**Example 5-1: Two-dimensional decomposition version of Example 2-4**

The CREATE\_DECOMP directive now lists two decomposed dimensions (with global sizes  $im$  and  $jm$ ). The halo width for each dimension is 1 in this case. As discussed in Section 3.2, the DISTRIBUTE, PARALLEL, and GLOBAL\_INDEX directives are modified to handle the 2-D decompositions. Although the communication patterns required to support 2-dimensional decompositions are more complex than the 1-dimensional case, SMS hides all of these details. Thus, the EXCHANGE directive is unchanged. Figure 5-1 shows some sample stencils overlaid on a 3x3 process decomposition of the problem. The stencil centered at global coordinate (3,2) only requires P1 communicate with P2. However, the stencil centered at global coordinate (4,4) requires P5 communicate with both P2 and P4. Figure 5-2 and Figure 5-3 show the full communications pattern for a 2-D exchange. Notice that the corner halo points of the center process are filled with data from the corresponding corner processes in the drawings.

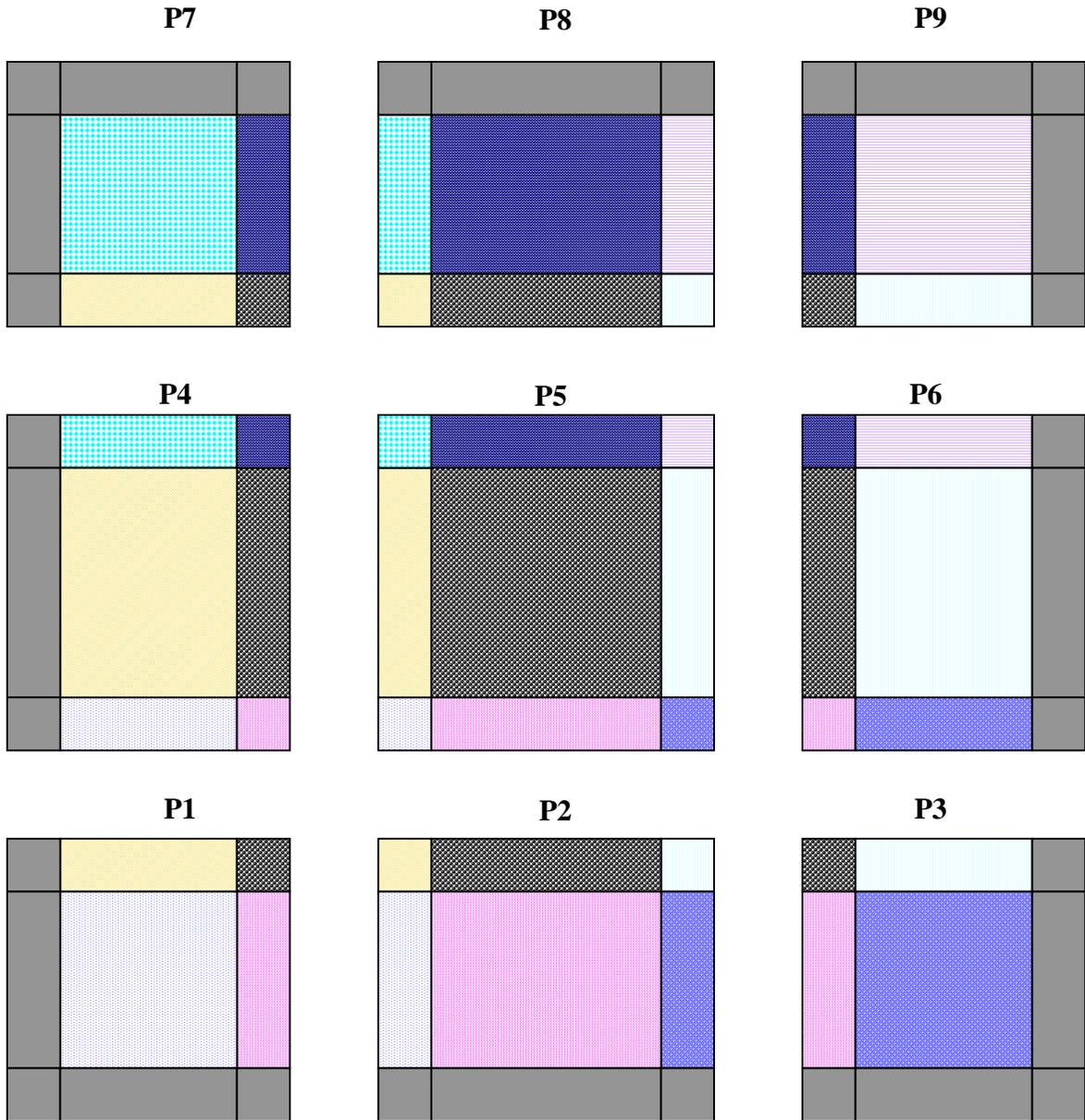


**BEFORE EXCHANGE**



**Figure 5-2: Schematic of how data are distributed among 9 processes just prior to an exchange operation. The big boxes contain the interior data. The boxes on the edges are the halo regions.**

**AFTER EXCHANGE**



**Figure 5-3: Illustration of the data distribution just after a 2-dimensional exchange for a problem with non-periodic boundaries.**

**5.1.2 Larger Stencils**

In Figure 2-10, the widths of the stencil for the calculation of  $d\mathcal{F}$  in the laplace program are one point in each direction. Since this is the only computation in Laplace requiring "exchange", it is clear that the halo widths specified by CREATE\_DECOMP must be 1 in the  $i$  dimension (line 5). However, suppose we modify Example 2-4 by adding

additional calculations of x that step 2 points into the halo region (lines 60-64 in Example 5-2 below).

```

1      program basic_ex_halo2
2      include 'basic.inc'
3      im = 10
4      jm = 10
5      CSMS$CREATE_DECOMP(DECOMP_I, <im>, <2>)
6      call laplace
7      end
8
9      subroutine laplace
10     include 'basic.inc'
11     integer i, j, iter
12     real max_error
13     real tolerance
14     parameter (tolerance = 0.001)
15     CSMS$DISTRIBUTE(DECOMP_I, <im>) BEGIN
16     real f(im,jm), df(im,jm)
17     CSMS$DISTRIBUTE END
18     CSMS$PARALLEL(DECOMP_I,<i>) BEGIN
19     do 100 j=1,jm
20     do 100 i=1,im
21     f (i,j) = 0.0
22     df(i,j) = 0.0
23     100 continue
24     do 110 j=1,jm
25     CSMS$GLOBAL_INDEX(1) BEGIN
26     f( 1,j) = 2.0
27     f(im,j) = 2.0
28     CSMS$GLOBAL_INDEX END
29     110 continue
30     do 120 i=1,im
31     f(i, 1) = 2.0
32     f(i,jm) = 2.0
33     120 continue
34     iter = 0
35     max_error = 2.0 * tolerance
36     C main iteration loop...
37     do while ((max_error .gt. tolerance) .and. (iter .lt. 1000))
38     iter = iter + 1
39     max_error = 0.0
40     CSMS$EXCHANGE(f)
41     do 200 j=2,jm-1
42     do 200 i=2,im-1
43     df(i,j) = 0.25*(f(i-1,j) + f(i+1,j) + f(i,j-1) + f(i,j+1))
44     & - f(i,j)
45     200 continue
46     do 300 j=2,jm-1
47     do 300 i=2,im-1
48     if (max_error .lt. abs(df(i,j))) then
49     max_error = abs(df(i,j))
50     endif
51     300 continue
52     CSMS$REDUCE(max_error, MAX)
53     do 400 j=2,jm-1
54     do 400 i=2,im-1
55     f(i,j) = f(i,j) + df(i,j)
56     400 continue
57     enddo
58
59     CSMS$EXCHANGE(df)
60     do j = 1, jm

```

```

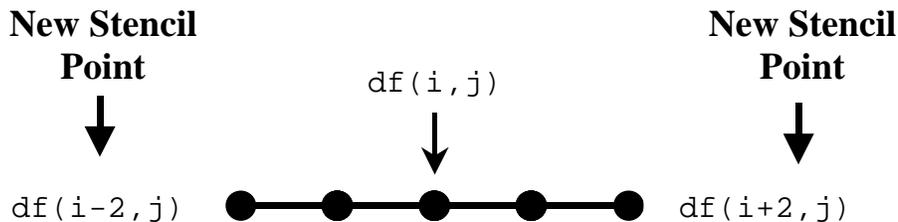
61     do i = 3, im-2
62         f(i,j) = f(i,j) + 2.0*df(i,j) - df(i-2,j) - df(i+2,j)
63     end do
64 end do
65
66 CSMS$PARALLEL END
67     print *, 'Solution required ',iter,' iterations'
68     print *, 'Final error = ', max_error
69
70     end

```

**Example 5-2: Modified version of Example 2-4 with additional code that has a stencil width of 2 in the *i* direction.**

For the calculations starting at line 60, the width of the stencil is 2 in the *i* direction as shown in Figure 5-4.

$$df2(i,j) = 2.0*df(i,j) - df(i-2,j) - df(i+2,j)$$



**Figure 5-4: Modified stencil for additional calculations in Example 5-2. This time the stencil width is 2 in the *i* direction.**

This program now has two calculations involving the same dimension of the same decomposition with different stencil widths. SMS handles this by requiring the programmer to make the halo width of the decomposition equal to the larger of the two widths. It is up to the programmer to determine the width of the largest stencil required in each dimension for every decomposition. The CREATE\_DECOMP directive (line 5) shows the correct halo width specification (<2>).

### 5.1.3 Exchanges in Static Memory Models

For static memory models that require exchanges, the process-local array sizes specified in the DECLARE\_DECOMP directive must be large enough to include the halo regions. In the program fragment below, the halo size is one. Since halo regions are on each side, the declared local array size is the global size (*im*) divided by the number of processes (4) plus 2 to accommodate the halo regions and plus 1 since 4 does not divide 30 evenly.

```

program STATIC_MEMORY_EXCHANGE
implicit none
integer im
parameter(im = 30)

```

```

integer jm
parameter(jm = 5)
CSMS$DECLARE_DECOMP(my_dh, <im/4 + 2 + 1>)

```

### 5.1.4 Miscellaneous

Another point about EXCHANGE is that, for both static and dynamic memory models, the number of processes used must be small enough to ensure the size of the smallest interior region is greater than the halo width in each decomposed dimension. SMS will issue the following run-time error message if this condition is violated:

```

Process: 0 Error status= -2100 MSG IS: NNT_DECOMP_ERR

```

Also, we point out that EXCHANGE automatically implements the process synchronization required for the parallel code to produce correct results. A process scheduled to receive data from a neighbor will wait until the data have fully arrived before proceeding with the next set of calculations. A side effect of this synchronization is that the EXCHANGE directive cannot be used inside a decomposed loop because the number of iterations may not be the same on every process, causing deadlock.

## 5.2 Performance Optimizations

In this section, some optimizations are described that can be employed to reduce the number of exchanges and the amount of data exchanged in a parallel SMS program.

### 5.2.1 Limited-Thickness Exchanges

Choosing a single halo width could mean some data are communicated unnecessarily. The exchange at line 40 in Example 5-2 is an example of such inefficiency. The stencil of the computations in loop 200 is still one point wide in the  $i$  direction. However, since the halo width of  $f$  is now 2 in this dimension, one extra halo point on each side for each  $j$  index will be communicated unnecessarily. This extra communication can be eliminated by using a variant of the EXCHANGE directive that only exchanges part of the halo region:

```

CSMS$EXCHANGE(f <1,1>)

```

This option to EXCHANGE tells SMS to exchange only the first halo point in the lower and upper halo regions.

If we were to modify Example 5-2 to use a two dimensional decomposition, the CREATE\_DECOMP directive would look as follows:

```

CSMS$CREATE_DECOMP(DECOMP_IJ, <im,jm>, <2,1>)

```

Now, the maximum stencil width is 2 in the first decomposed dimension (for the exchange at line 59) and 1 in the second decomposed dimension (for the exchange at line

40). If the exchange at line 40 only requires 1 point in each direction, it could be optimally written as:

```
CSMS$EXCHANGE(f <1,1> <1,1>)
```

Since the exchange at line 59 is only needed to update points in the  $i$  dimension, it would optimally be written as:

```
CSMS$EXCHANGE(df <2,2> <0,0>)
```

## 5.2.2 Aggregating Exchanges

The program SLOW in Example 5-3, uses a statically declared one-dimensional decomposition (line 10) to distribute the arrays  $a$ ,  $b$  and  $c$ . In this example, a halo thickness of one is defined by CREATE\_DECOMP (line 24). After a series of iterations (line 39) a global sum is produced with the REDUCE directive (line 63).

```
1      program SLOW
2      implicit none
3      integer im
4      parameter(im = 30)
5      integer jm
6      parameter(jm = 5)
7      integer iterations
8      parameter(iterations = 3)
9
10     CSMS$DECLARE_DECOMP(my_dh, <im/3 + 2>)
11
12     CSMS$DISTRIBUTE(my_dh, <im>) BEGIN
13         real a(im)
14         real b(im,jm)
15         real c(im,jm)
16     CSMS$DISTRIBUTE END
17
18     real ysum
19
20     integer i
21     integer j
22     integer iter
23
24     CSMS$CREATE_DECOMP(my_dh, <im>, <1>)
25
26     ysum = 0.0
27     b = 0.0
28     c = 0.0
29
30     do j = 1, jm
31
32     CSMS$PARALLEL(my_dh, <i>) BEGIN
33         do i = 1, im
34     CSMS$TO_GLOBAL(<1, i>) BEGIN
35         a(i) = real(3*i + 2 + j)
36     CSMS$TO_GLOBAL END
37         end do
38
39         do iter = 1, iterations
40
41     CSMS$EXCHANGE(a)
42
```

```

43         do i = 2, im-1
44             b(i,j) = a(i+1) + a(i-1)
45             c(i,j) = b(i,j) + c(i,j)
46         end do
47
48     CSMS$EXCHANGE(b)
49     CSMS$EXCHANGE(c)
50
51         do i = 2, im-1
52             a(i) = b(i+1,j) + b(i-1,j) + c(i+1,j) - c(i-1,j)
53         end do
54
55     end do
56
57     do i = 2, im - 1
58         ysum = ysum + a(i)
59     end do
60
61 end do
62
63 CSMS$REDUCE(ysum, SUM)
64
65     print *, 'ysum is ', ysum
66 CSMS$PARALLEL END
67     end

```

**Example 5-3: A sub-optimal version of a program parallelized using SMS.**

SMS provides the capability to aggregate the exchanges of multiple variables. If lines 48-49 are replaced with

```
CSMS$EXCHANGE(b,c)
```

then SMS will combine the communications of the corresponding halo regions of *b* and *c* as shown in Figure 5-5. Since the number of messages sent is halved, performance on high-latency machines will improve.

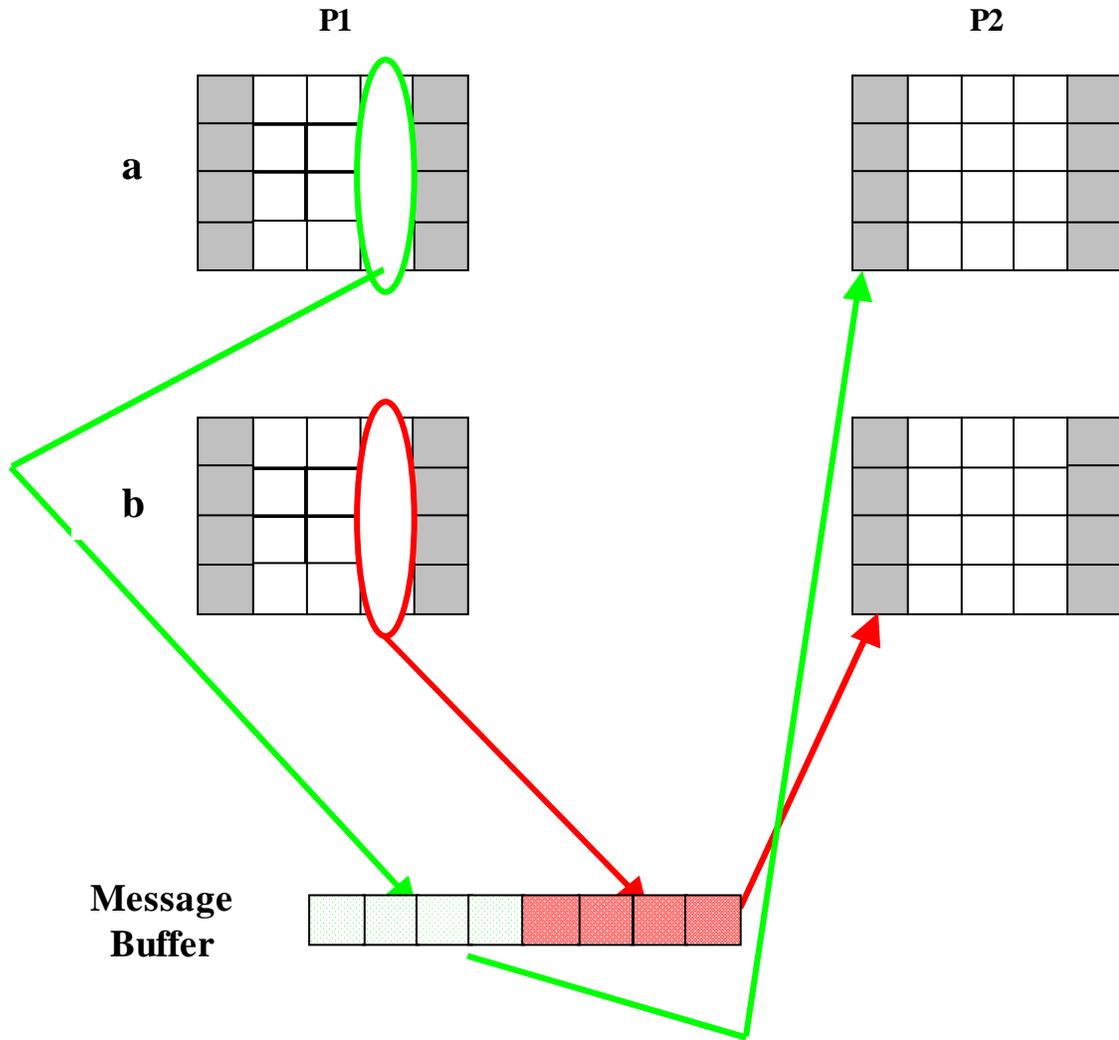


Figure 5-5: An illustration of how communications are aggregated to reduce latency for a portion of the exchange of *a* and *b*. The last column of process P1's variables are communicated as a single message to P2 where they are unpacked into the corresponding halo regions.

### 5.2.3 Exchanging Array Sections

Sometimes, it is not necessary to exchange an entire array. For example, in the following code fragment an adjacent dependence may only apply to some of the vertical levels of a 3D array:

```

CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
  real x(im,jm,km)
  real y(im,jm,km)
CSMS$DISTRIBUTE END

```

```

CSMS$EXCHANGE(x <0,1> <0,0>)
  do 100 k = 1,2
  do 100 j = 1,im
  do 100 i = 1,im
    y(i,j,k) = x(i+1,j,k) - x(i,j,k)
  100 continue

```

However, the exchange directive will exchange array  $x$  at all  $k$  levels even though the dependence exists only for  $k=1$  and  $k=2$ . This exchange directive can be optimized using standard Fortran array syntax as shown below:

```

CSMS$EXCHANGE(x(:, :, 1:2) <0,1> <0,0>)

```

Now, only the  $k=1$  and  $k=2$  levels of array  $x$  will be exchanged. Note that use of array section syntax will only improve performance for subsections in non-decomposed dimensions.

## 5.2.4 Trading Communications for Computations Using HALO\_COMP

Example 5-3 can be further optimized by trading communication for redundant computations in the halo region as is briefly discussed in the SMS overview paper. This is done using the HALO\_COMP directive to modify the ranges of parallel loops to include computations in the halo regions. These extra computations can eliminate the need for some exchanges.

Figure 5-6, Figure 5-7, and Figure 5-8 illustrate how redundant computations work. Without the HALO\_COMP directive,  $b$  and  $c$  are only computed in interior points using stencils like that shown in Figure 5-6. Halo regions of  $b$  and  $c$  must then be updated via an exchange for  $a$  to be properly computed as shown in Figure 5-7. A computation one step into the halo region (Figure 5-8) requires that  $a$  have a halo size of two instead of one. Since process P1 now computes points such as  $b(4,2)$  and  $c(4,2)$ , the computation of  $a(3,2)$  shown in Figure 5-7 can proceed without having exchanged  $b$  and  $c$ . However, extra computations are done since process P2 must also perform exactly the same computation for its corresponding interior points  $b(4,2)$  and  $c(4,2)$ ,

Use of the HALO\_COMP directive in this example reduces latency because the exchanges of  $b$  and  $c$  are no longer required. In addition, communication bandwidth is reduced. Although the amount of data communicated by the exchange of  $a$  has doubled, this is more than offset by the elimination of the exchanges of  $b$  and  $c$ .

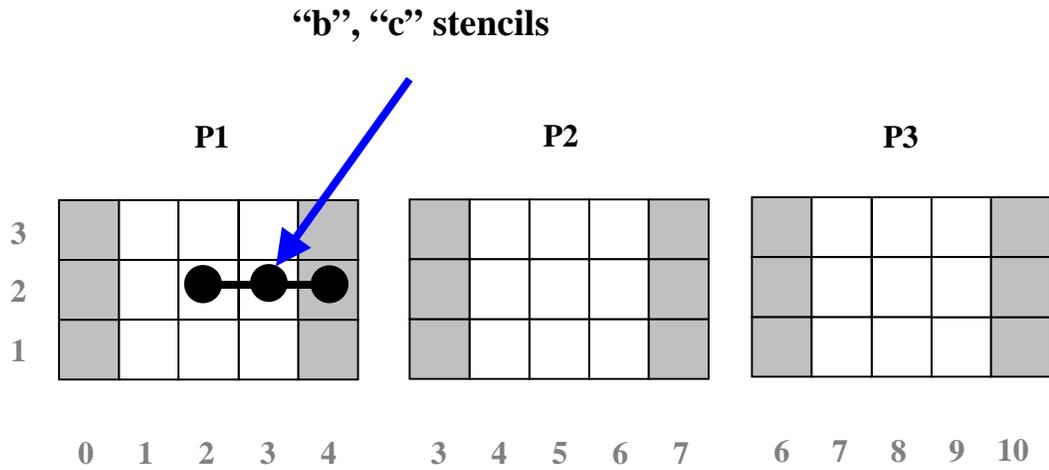


Figure 5-6: Memory layout of *a* (assuming  $im=9, jm=3$ ) with sample stencil for calculations of *b* and *c* overlaid.

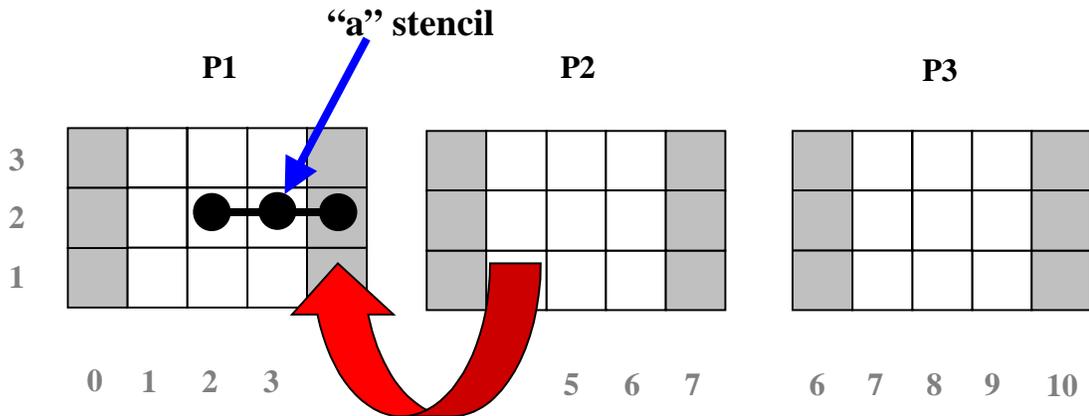
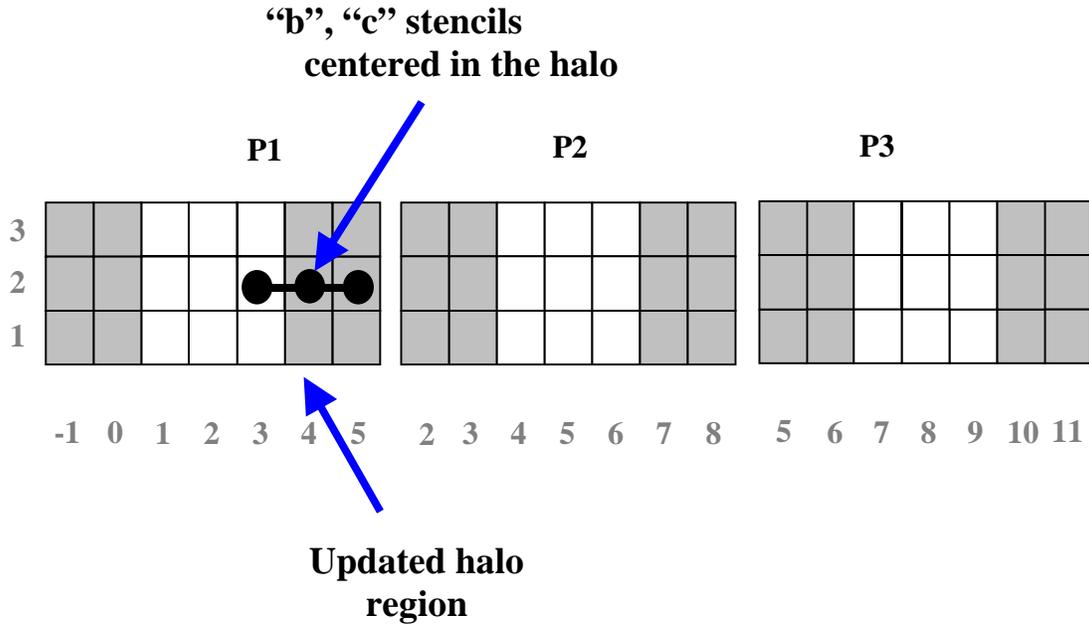


Figure 5-7: Memory layout of *b* and *c* with sample stencil for calculation of *a* overlaid. The halo regions of *b* and *c* must be updated via exchange for the calculation of *a* to be executed correctly.



**Figure 5-8: Modified memory layout of *a* with new sample stencil centered in the halo region. The computation of point  $b(4,2)$  and  $c(4,2)$  effectively updates the halo regions of *b* and *c* so that the computation of *a* in Figure 5-7 can be made without an exchange.**

A net improvement in performance by this technique will only be realized if the cost of the additional computation by each process is less than the cost of exchanging *b* and *c*. Whether or not the code runs faster will, in general, depend on the communication patterns in the program, the number of processes used, and the target hardware platform.

A version of Example 5-3 that implements redundant calculations is shown in Example 5-4. The HALO\_COMP directive on line 43 tells SMS that the enclosed loop should be executed 1 step into the halo region in each direction. This updates *b* and *c* sufficiently to satisfy the dependencies in the loop at lines 52-54. DECLARE\_DECOMP and CREATE\_DECOMP have been modified to accommodate the new halo size of 2. The exchanges of *b* and *c* have been eliminated.

```

1      program FASTER
2      implicit none
3      integer im
4      parameter(im = 30)
5      integer jm
6      parameter(jm = 5)
7      integer iterations
8      parameter(iterations = 3)
9
10     CSMS$DECLARE_DECOMP(my_dh, <im/3 + 4>)
11
12     CSMS$DISTRIBUTE(my_dh, <im>) BEGIN
13         real a(im)
14         real b(im,jm)

```

```

15      real c(im,jm)
16  CSMS$DISTRIBUTE END
17
18      real ysum
19
20      integer i
21      integer j
22      integer iter
23
24  CSMS$CREATE_DECOMP(my_dh, <im>, <2>)
25
26      ysum = 0.0
27      b = 0.0
28      c = 0.0
29
30      do j = 1, jm
31
32  CSMS$PARALLEL(my_dh, <i>) BEGIN
33      do i = 1, im
34  CSMS$TO_GLOBAL(<1, i>) BEGIN
35      a(i) = real(3*i + 2 + j)
36  CSMS$TO_GLOBAL END
37      end do
38
39      do iter = 1, iterations
40
41  CSMS$EXCHANGE(a)
42
43  CSMS$HALO_COMP(<1,1>) BEGIN
44      do i = 2, im-1
45          b(i,j) = a(i+1) + a(i-1)
46          c(i,j) = b(i,j) + c(i,j)
47      end do
48  CSMS$HALO_COMP END
49
50
51
52      do i = 2, im-1
53          a(i) = b(i+1,j) + b(i-1,j) + c(i+1,j) - c(i-1,j)
54      end do
55
56      end do
57
58      do i = 2, im - 1
59          ysum = ysum + a(i)
60      end do
61
62      end do
63
64  CSMS$REDUCE(ysum, SUM)
65
66      print *, 'ysum is ', ysum
67
68  CSMS$PARALLEL END
69
70      end

```

**Example 5-4:** A version of Example 5-3 that has been optimized by trading communications for redundant calculations in the halo region.

## 5.2.5 Using HALO\_COMP and TO\_LOCAL To Make Subroutines Do Redundant Computations

We saw in Section 4.3 how the TO\_LOCAL directive can be used to parallelize subroutines without requiring directives inside the subroutine code. The approach works by making the subroutines operate on the interior of the process local arrays. Now, suppose we want those called routines to do redundant computations in the halo region to avoid communication. Example 5-5 shows a modified version of Example 4-2, illustrating how this is done.

```
1      program tran_index7
2      implicit none
3      integer i
4      integer, parameter :: im = 8
5
6      CSMS$DECLARE_DECOMP(dh, 1)
7
8      integer istart, iend
9
10     CSMS$DISTRIBUTE(dh, 1) BEGIN
11     integer, allocatable :: x(:), y(:), z(:)
12     CSMS$DISTRIBUTE END
13
14     CSMS$CREATE_DECOMP(dh, <im>, <2>)
15
16     allocate(x(im))
17     allocate(y(im))
18     allocate(z(im))
19     x = 0.0
20     CSMS$PARALLEL(dh,<i>) BEGIN
21
22     CSMS$HALO_COMP(<1,1>) BEGIN
23     do i=1,im
24     CSMS$TO_GLOBAL(<1,i>) BEGIN
25     x(i) = i
26     CSMS$TO_GLOBAL END
27     end do
28
29     y = 0.0
30
31
32     CSMS$TO_LOCAL(<1, istart : lbound>, <1, iend : ubound>) BEGIN
33     istart = 1
34     iend = im - 1
35     CSMS$TO_LOCAL END
36     CSMS$HALO_COMP END
37
38     call physics(x, lbound(x,1), ubound(x,1), istart, iend, y)
39
40     do i = 1, im - 1
41     z(i) = y(i) + y(i+1)
42     end do
43
44     CSMS$SERIAL BEGIN
45     write(*,'(8i5)') (z(i),i=1,im)
46     CSMS$SERIAL END
47
48     CSMS$PARALLEL END
49     end
```

**Example 5-5: Modified version of Example 4-2 that passes lower and upper bounds into subroutine *physics* so that it does redundant computations for one point in the halo region for each dimension and for each direction.**

Since the calculations of *istart* and *iend* are now contained within both a *TO\_LOCAL* and *HALO\_COMP* directive, the effect is to change the lower and upper bounds passed to the *physics* so that it will do redundant computations for one point in the halo region for each direction. Figure 5-9 shows the new table of lower and upper bounds (compare to the table in Figure 4-1). Now, following the call to *physics*, variable *y* contains valid data one point into the halo region. Consequently, an *EXCHANGE* directive is not need prior to the loop at lines 39-41.

Process	i_mem_start	i_mem_end	dim1_start	dim1_end
P1	-1	6	1	5
P2	3	10	4	8

**Figure 5-9: Table of memory and computational bounds for Example 5-5. Compare the *dim1\_start* and *dim1\_end* values to those the table in Figure 4-1. The memory start and end values are unchanged.**

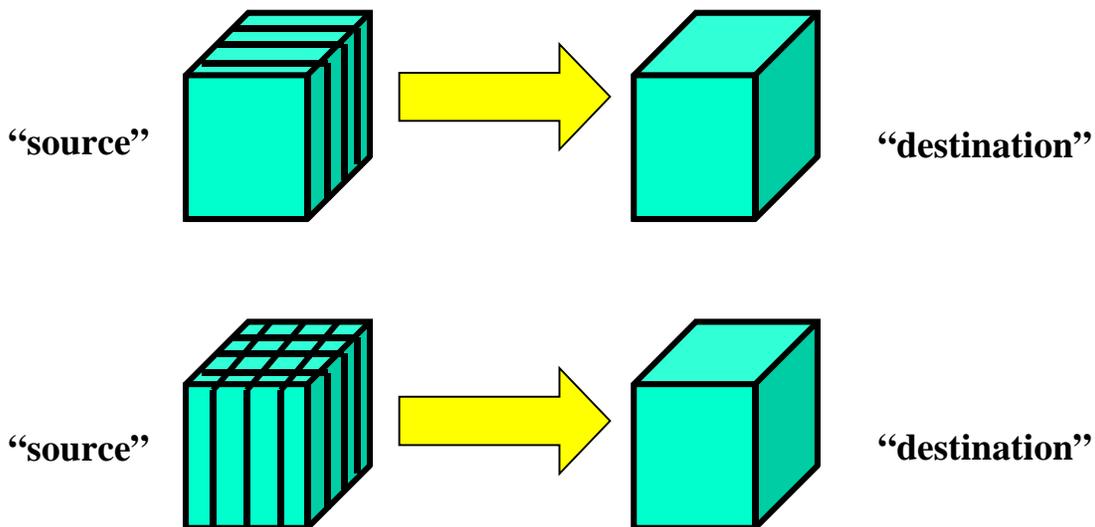
## 6 Handling Complex Dependencies Using TRANSFER

Section 2.7 introduced the TRANSFER directive and explained how it could be used to handle complex dependencies in more than one dimension (see Example 2-5). In Section 6.1, we show how TRANSFER can be used when either the source or destination array are non-decomposed. In Section 6.2, we examine how TRANSFER can be applied to the parallelization of spectral models. TRANSFER is also used for inter-grid interpolation as described in Section 10.2.

Like EXCHANGE, TRANSFER automatically implements the process synchronization required for the parallel code to produce correct results. A side effect of this synchronization is that the TRANSFER directive cannot be used inside a decomposed loop because the number of iterations may not be the same on every process, causing deadlock.

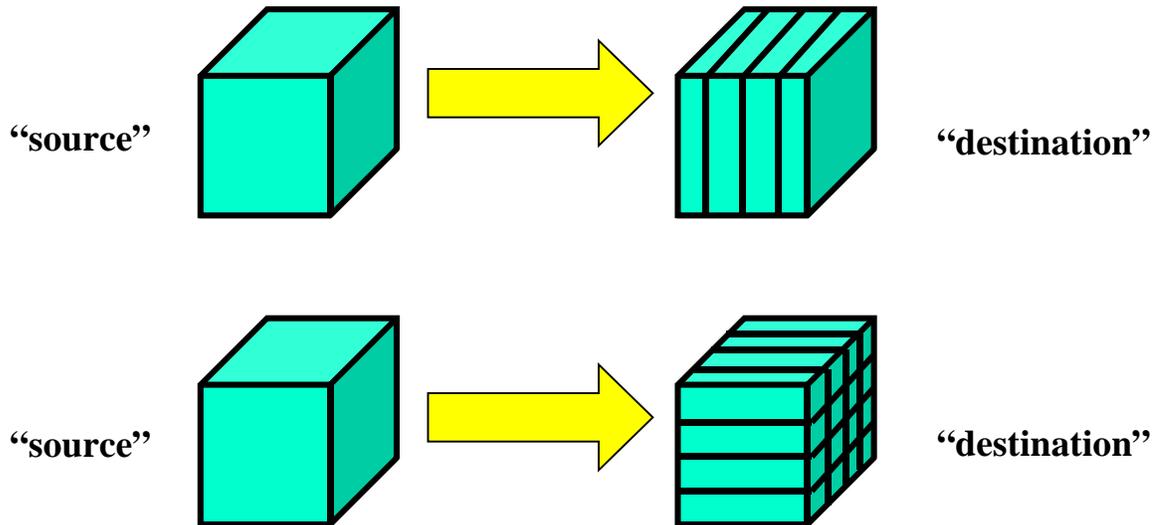
### 6.1 Further Details about TRANSFER

While TRANSFER can be used to generate communications to transpose arrays decomposed in one or more dimensions, it can also be used when either the source or destination arrays are not decomposed. If the destination array is not decomposed but the source is, then the TRANSFER directive effectively implements a “gather” of the source into the destination as illustrated in Figure 6-1. After the transfer, the entire array is replicated on each process. Since the local data for each process must be communicated to all other processes, this operation can be quite expensive.



**Figure 6-1:** Schematic of the behavior of TRANSFER when the source array is decomposed and the destination array is NOT decomposed. The effect is to “gather” the process-local data from the source array into the globally-sized destination array.

If the source array is not decomposed but the destination array is decomposed then TRANSFER performs an "extract" of data from the source into the destination as shown in Figure 6-2. Note that no communication is needed in this case since each process has access to all needed data to begin with. Note that the SERIAL directive also performs "gather" and "extract" operations and is often easier to use than TRANSFER (see Section 8).



**Figure 6-2:** Schematic of the behavior of TRANSFER when the source array is NOT decomposed and the destination array is decomposed. The effect is to “extract” the appropriate data from the globally sized source array into the process-local destination array.

As in the case of EXCHANGE, TRANSFERS can be aggregated as follows to reduce latency:

```
CSMS$TRANSFER(<source1, destination1>, <source2, destination2>) BEGIN
    Serial code to be replaced...
CSMS$TRANSFER END
```

Note that, for all TRANSFER directives, the type and rank of the source and destination arrays must be the same. However, the array sizes may differ.

## 6.2 Applying TRANSFER to Spectral Models

Many spectral models have multiple phases of computation that repeat in a fixed pattern. Phases often have different optimal decompositions. Therefore, performance may be maximized by using multiple decompositions and transferring between them. Consider the case of one-dimensional decompositions for these models. The physical parameterizations contain complex dependencies in the vertical. This makes it efficient to decompose in one of the horizontal dimensions. At the same time, many computer system vendors provide highly optimized assembly FFT libraries that far out-perform anything that can be done with hand-tuned Fortran code.

Taking advantage of this serial code requires decomposing in a dimension other than  $i$ . Typically, the data are decomposed in the  $j$  dimension during physics and FFT computations (decomposition "a" in Figure 3-1). The Legendre transformations contain complex dependencies in the  $j$  dimension. Therefore, a second decomposition in either  $i$  (decomposition "b" in Figure 3-1) or  $k$  (decomposition "c" in Figure 3-1) is needed for optimal performance during these calculations. The TRANSFER directive provides the means to transpose the data from decomposition "a" to ("b" or "c") and back again. For large numbers of processes, 2D decompositions are needed to ensure that all processes have work to do. In this case, physics computations may be done using decomposition "d", FFT computations may be done using decomposition "e", and Legendre transform computations may be done using decomposition "f".

## 7 Handling Global Dependencies Using REDUCE

In Section 2.3, we saw how the REDUCE directive was used to implement communication needed to do global summations and maxima. In this section we examine the REDUCE directive in more detail. In addition to giving more examples of the directive, Section 7.1 also shows that the form of REDUCE introduced in Section 2.3.3 (which will be referred to as "Standard Reductions") does not necessarily produce the bit-wise exact same answer as the serial code for global summations of floating point numbers. Section 7.2 introduces a second form of REDUCE called "Bit-wise Exact" that does produce the bit-wise exact same answer, regardless of the number of processes. Although quite useful for debugging, the second form has the drawback that it runs slowly.

Like EXCHANGE, both forms of REDUCE automatically implement the process synchronization required for the parallel code to produce correct results. A side effect of this synchronization is that the REDUCE directive cannot be used inside a decomposed loop because the number of iterations may not be the same on every process, causing deadlock.

### 7.1 More on Standard Reductions

In addition to global summations and maxima, the REDUCE directive can be used to generate global minima and to reduce arrays as seen in Example 7-1. Global minima are generated by specifying the keyword **MIN** (line 52). Also notice that reductions can be aggregated in the same way as exchanges (line 50). One of the variables reduced is the non-decomposed array *xsum* (line 50). The summation of *xsum* looks like the following:

```
Xsum_global(1) = Xsum_local1 (1) + Xsum_local2 (1) + ...
Xsum_global(2) = Xsum_local1 (2) + Xsum_local2 (2) + ...
.
.
.
```

where  $Xsum\_local(j)$  is the value of process-local  $xsum(j)$  on process P and  $Xsum\_global$  is the value of  $xsum$  after the global summation is complete.

```
1      program REDUCTIONS
2      implicit none
3      include 'basic.inc'
4
5      im = 50
6      jm = 2
7
8      CSMS$CREATE_DECOMP(DECOMP_I, <im>, <0>)
9
10     call DO_THEM
11
12     end
13
```

```

14      subroutine DO_THEM
15      implicit none
16      include 'basic.inc'
17
18      CSMS$DISTRIBUTE(DECOMP_I, <im>) BEGIN
19          real x(im,jm)
20          real y(im,jm)
21      CSMS$DISTRIBUTE END
22
23          real xsum(jm)
24          real ysum
25          real xmin
26          real xmax
27
28          integer i
29          integer j
30
31          open (10, file='reduce_data', form='unformatted')
32          read (10) x, y
33          close(10)
34
35      CSMS$PARALLEL(DECOMP_I, <i>) BEGIN
36          xsum = 0.0
37          ysum = 0.0
38          xmax = -999.0
39          xmin = 999.0
40
41          do j = 1, jm
42              do i = 1, im
43                  xsum(j) = xsum(j) + x(i,j)
44                  ysum = ysum + y(i,j)
45                  xmax = max(xmax, x(i,j))
46                  xmin = min(xmin, x(i,j))
47              end do
48          end do
49
50      CSMS$REDUCE(xsum, ysum, SUM)
51      CSMS$REDUCE(xmax,      MAX)
52      CSMS$REDUCE(xmin,      MIN)
53
54          print *
55          print *, 'Global values'
56          do j = 1, jm
57              write(*,100) j, xsum(j)
58          end do
59          write(*,150) ysum
60          write(*,200) xmax
61          write(*,300) xmin
62
63      100 format('j ', i2, ' xsum = ', F13.5)
64      150 format('ysum = ', F13.5)
65      200 format('xmax = ', F13.5)
66      300 format('xmin = ', F13.5)
67
68      CSMS$PARALLEL END
69
70          return
71      end

```

**Example 7-1: Program showing additional examples of how the REDUCE directive can be used.**

If we were to modify Example 7-1 so that the *j* dimension is also decomposed and were to make *xsum* a decomposed variable, then the reduction of *xsum* would FAIL. This

would happen because SMS does not currently support reductions that produce decomposed variables. (This would require doing the reduction over a subset of the processes.)

When run with 2 processes, program REDUCTIONS yields the following results:

```
Global values
j 1 xsum = 1258.28589
j 2 xsum = 1310.71448
ysum = -2464.28540
xmax = 100.00000
xmin = -100.00000
```

However, when run with 4 processes, the results are:

```
Global values
j 1 xsum = 1258.28577
j 2 xsum = 1310.71436
ysum = -2464.28613
xmax = 100.00000
xmin = -100.00000
```

Notice that the values for *xsum* and *ysum* are slightly different between the 2 and 4 process runs. We will now see why this is the case.

## 7.2 Bit-wise Exact Reductions

The differences in results in Example 7-1 are due to round-off error in the floating-point addition. The numbers are added in a different order in the 4-process case as compared to the 2-process case because the sums are first computed locally before being combined (see Section 2.3.3). In weather and climate models (which are non-linear systems), if the global sums feed back into the main model equations, these slight errors can grow and propagate; potentially yielding completely different model predictions for runs with differing numbers of processes.

For debugging purposes, it is useful to avoid these round-off errors. In fact, this is necessary for correct operation of the COMPARE\_VAR debugging feature (see Section 13.1). To do this, SMS offers a form of REDUCE that produces the bit-wise exact same answer for any number of processes. Example 7-2 below shows how this works.

```
1      program EXACT_REDUCTIONS
2      implicit none
3      include 'basic.inc'
4
5      im = 50
6      jm = 2
7
8      CSMS$CREATE_DECOMP(DECOMP_I, <im>, <0>)
9
10     call DO_THEM
11
12     end
13
```

```

14      subroutine DO_THEM
15      implicit none
16      include 'basic.inc'
17
18      CSMS$DISTRIBUTE(DECOMP_I, <im>) BEGIN
19      real x(im,jm), y(im,jm)
20      CSMS$DISTRIBUTE END
21
22      real ysum
23
24      integer i
25      integer j
26
27      open (10, file='reduce_data', form='unformatted')
28      read (10) x, y
29      close(10)
30
31      CSMS$PARALLEL(DECOMP_I, <i>) BEGIN
32
33      CSMS$REDUCE(ysum, SUM) BEGIN
34      ysum = 0.0
35      do j = 1, jm
36      do i = 1, im
37      ysum = ysum + y(i,j)
38      end do
39      end do
40      CSMS$REDUCE END
41
42      print *
43      print *, 'Global values'
44      write(*,150) ysum
45
46      150 format('ysum = ', F13.5)
47
48      CSMS$PARALLEL END
49
50      return
51      end

```

**Example 7-2: Program illustrating the bit-wise exact form of the REDUCE directive.**

The modified REDUCE syntax can be seen on lines 33 and 40. The syntax requires a BEGIN and END directive. SMS replaces the calculations between the REDUCE BEGIN and END with code that gathers each process's piece of  $y$  into a globally-sized (replicated) variable and then sums the result in the correct order. Conceptually, the generated parallel code would be:

```

call GATHER(y, y_global)
ysum = 0.0
do j = 1, jm
  do i = 1, im
    ysum = ysum + y_global(i,j)
  end do
end do

```

The "gather" operation is done in the same way as TRANSFER was used to gather variables as discussed in Section 6.1. Since the gather operation requires communicating the entire contents of  $y$  to all processes, this form of global sum is significantly less efficient than the "standard" form of reduction. In that case, only the process-local scalar sums are communicated to all the processes.

The bit-wise exact form of the REDUCE directive will only produce exact sums if an environment variable called SMS\_BITWISE is set to the value EXACT. Running EXACT\_REDUCE in a c-shell environment might look as follows:

```
>> setenv SMS_BITWISE EXACT
>> smsRun -np 2 exact_reduce
SMS: BITWISE EXACT reductions will be used when requested.
Global values
ysum = -2464.28418
```

Notice that the message printed by SMS regarding reductions now indicates that bit-wise exact reductions will be used.

If SMS\_BITWISE is NOT set to EXACT then the effect of the REDUCE directive is the same as in the "standard" reduction; each process computes a local sum of  $y$  and the resulting scalars are summed across the processes.

```
>> setenv SMS_BITWISE INEXACT
>> smsRun -np 2 exact_reduce
SMS: Standard reductions will be used.
Global values
ysum = -2464.28540
```

In summary, the "bit-wise exact" form of global summation is valuable for testing purposes, particularly for non-linear systems. Its use is required for correct operation of the COMPARE\_VAR debugging feature (see Section 13.1). However, for long model runs, when optimal performance is important, the "standard" form of REDUCE will likely be more appropriate because it is much faster. The programmer can use the bit-wise exact form of REDUCE in the code and then decide at run-time, with the SMS\_BITWISE environment variable, which reduction to use.

## 8 Incremental Parallelization Using SERIAL

The SERIAL directive is useful when other SMS directives cannot be easily applied to a piece of serial code or when efficient performance is not critical. One example is initialization. For long model runs, the effects of inefficient code during initialization become negligible. Diagnostic print messages are another case. If the user can turn off diagnostic messages when high performance is needed then the presence of inefficient parallel code that generates these messages does not pose a problem. Another use of SERIAL is for incremental parallelization of a large code. Using this technique, SERIAL directives are inserted around large sections of code. The directives are then removed one-by-one as each section is parallelized. This simplifies testing and debugging because each parallel code section can be tested and debugged separately.

### 8.1 Improving the Performance of SERIAL

Any code segment enclosed by SERIAL BEGIN and SERIAL END directives is called a “serial region”. When SMS encounters a serial region, it automatically gathers all decomposed arrays on a single process, executes the enclosed code segment on the process, and scatters all decomposed arrays back to all processes (see Figure 8-1). In addition, all non-decomposed variables are broadcast to all the processes. By default, SMS gathers/scatters all decomposed variables and broadcasts all non-decomposed variables referenced in the serial region. These communications cause the code to run even more slowly than the original serial code. To improve performance, the user can specify variables to be scattered, gathered, or broadcast using keywords “*IN*”, “*OUT*”, or “*INOUT*” (when both operations are required). Also, when no communication is required, the “*IGNORE*” keyword can be used to suppress communication.

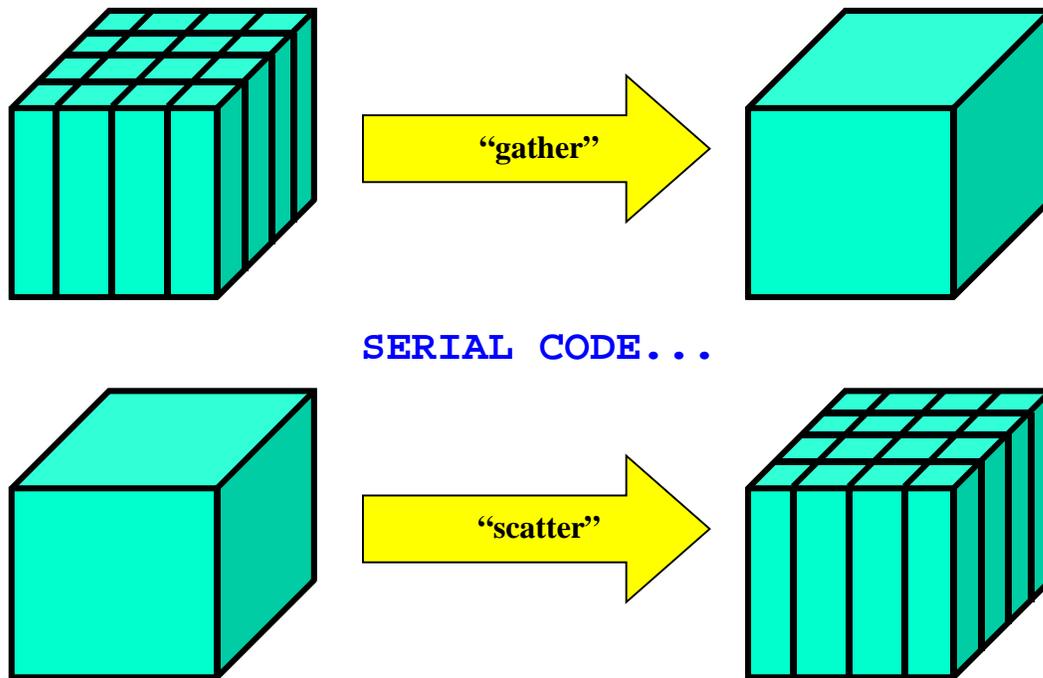


Figure 8-1: Gather and scatter of decomposed data at the beginning and end of a serial region.

In Example 8-1,  $x$  and  $y$  are decomposed while  $z$  is not. The subroutine calls at lines 39-40 read in  $x$  and  $z$  using C language routines (the C code is not shown). Since C routines cannot be handled by SMS, the SERIAL directive is used to generate code that gathers  $x$  and  $y$  into global variables. A single process then executes the code at lines 39-41. Finally, the generated code scatters  $x$  and  $y$  and broadcasts the value of  $z$ . Scalar variables  $im$  and  $jm$  are not broadcast since the default is set to "IGNORE". When high performance is desired, the user can avoid this poorly performing code segment by setting `ENABLE_DIAGS` to `.FALSE..`

```
[Include file: serial.inc]
1     integer im,jm
2     common /sizes_com/ im,jm
3     CSMS$DECLARE_DECOMP(DECOMP_IJ, 2)

[Source file: serial1.f]
1     program SERIAL
2
3     include 'serial.inc'
4
5     integer i
6     integer j
7
8     im = 5
9     jm = 4
10
11    CSMS$CREATE_DECOMP(DECOMP_IJ, <im,jm>, <0,0>)
12
13    call DO_IT
14
15    end
```

```

16
17     subroutine DO_IT
18     include 'serial.inc'
19
20 CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
21     real x(im,jm)
22     real y(im,jm)
23 CSMS$DISTRIBUTE END
24     real z
25     logical ENABLE_DIAGS
26     ENABLE_DIAGS = .true.
27
28     open(10, file='yin', form='unformatted')
29     read(10) y
30     close(10)
31
32 C Some parallel computations
33 C     .
34 C     .
35 C     .
36
37     if (ENABLE_DIAGS) then
38 CSMS$SERIAL(<x, y, INOUT>, <z, OUT> : DEFAULT=IGNORE) BEGIN
39     call READ_2D_ARRAY_USING_C(x, im, jm)
40     call READ_SCALAR_USING_C(z)
41     print *, 'y(2,2), z ', y(2,2), z
42 CSMS$SERIAL END
43     end if
44 C More parallel calculations
45     .
46     .
47     return
48     end

```

**Example 8-1:** A sample program showing how the SERIAL directive can be used to generate correct parallel code in a simple fashion when other SMS directives will not suffice.

## 8.2 Limitations of SERIAL

Some care must be taken when using the SERIAL directive. First, SMS does not perform inter-procedural analysis so calling a routine that uses common block variables from within a serial region can produce erroneous results. Suppose we insert the following code after line 38 in Example 8-1:

```
call sub1
```

Further suppose *sub1* looks like this:

```

subroutine sub1
real xc
common /com1/ xc
xc = 2.0
return
end

```

SMS has no way of knowing that *xc* has to be broadcast before the end of the serial region because it does no inter-procedural analysis. A solution here would be to include */com1/* in subroutine *DO\_IT* and specify *xc* as an “**OUT**” variable in the directive:

```
CSMS$SERIAL(<x, y, INOUT>, <xc, z, OUT> : DEFAULT=IGNORE) BEGIN
```

Second, care must be taken when constants are passed into subroutines that use SERIAL. In Example 8-2, the constant 2 is passed to subroutine *DO\_IT*. Since *DO\_IT* calls a C routine that uses dummy argument *n*, a SERIAL directive would normally be required to handle this. However the SERIAL directive generates a broadcast of dummy argument *n*, which will attempt to modify its value. Since *n* is the constant 2, the result will be unpredictable (with luck, a core dump). This problem can be corrected by using the “**IGNORE**” keyword to ensure that the SERIAL directive will not attempt to broadcast *n*.

```
1      program SERIAL
2
3      include 'serial.inc'
4
5      integer i
6      integer j
7
8      im = 5
9      jm = 4
10
11     CSMS$CREATE_DECOMP(DECOMP_IJ, <im,jm>, <0,0>)
12
13     call DO_IT(2)
14
15     end
16
17     subroutine DO_IT(n)
18
19     integer n
20
21     CSMS$SERIAL BEGIN
22     call c_routine(n)
23     CSMS$SERIAL END
24
25     return
26     end
```

**Example 8-2:** Code where use of the SERIAL directive generates parallel code that fails to run properly.

Care must also be taken when using SERIAL in combination with statements that alter execution sequence (such as IF, GOTO, etc.). To avoid problems, program execution should only enter a serial region from the statement immediately before the SERIAL BEGIN directive. Similarly, program execution should only exit a serial region from the statement immediately before the SERIAL END directive.

Finally, like EXCHANGE, SERIAL automatically implements the process synchronization required for the parallel code to produce correct results. A side effect of this synchronization is that the SERIAL directive cannot be used inside a decomposed loop because the number of iterations may not be the same on every process, causing deadlock.

## 9 Periodic Boundary Conditions

By default, SMS assumes that models have non-periodic boundary conditions. However, SMS allows the user to specify periodic boundary conditions using the *PERIODIC* option to the CREATE\_DECOMP directive as shown in Example 9-1. This option causes exchanges to fill in the halo regions on the outer model boundaries as shown in Figure 9-1.

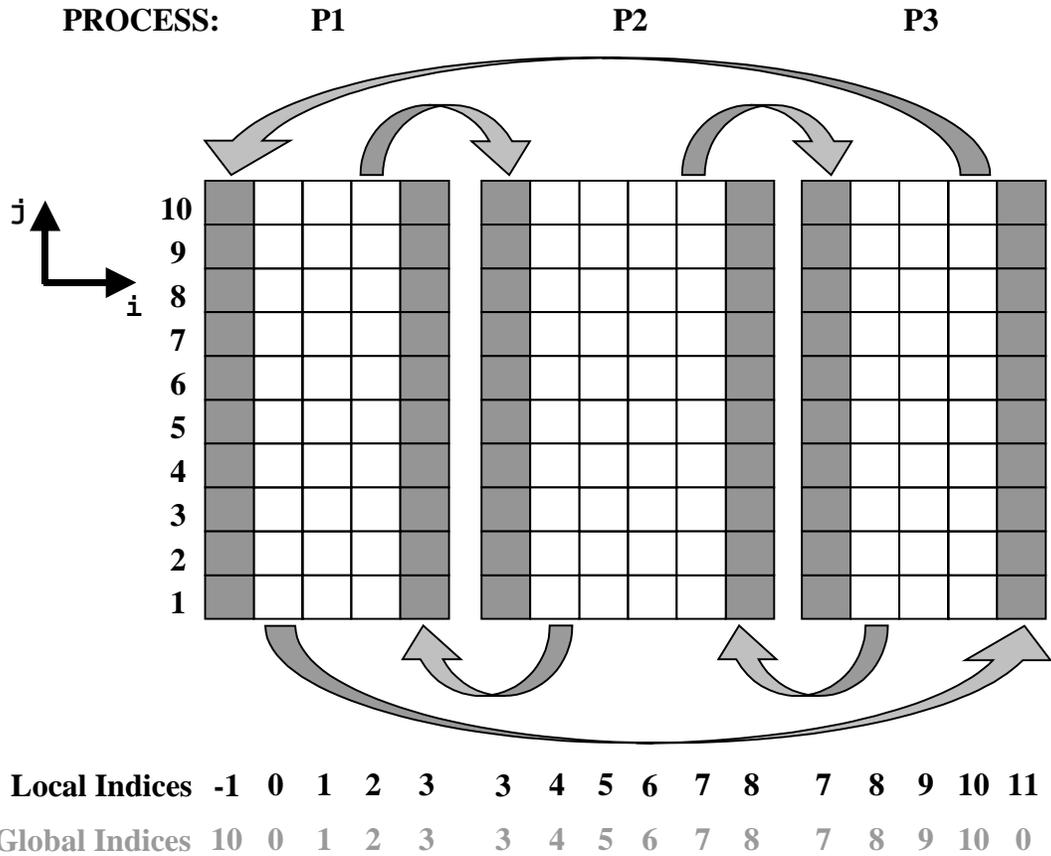
```
1      program laplace_periodic
2      implicit none
3
4      integer im
5      parameter(im=100)
6      integer jm
7      parameter(jm=99)
8
9      CSMS$DECLARE_DECOMP(dh : <0,0>)
10
11     real global_error
12     real tolerance
13     parameter (tolerance = 0.002)
14     integer iter
15     integer i
16     integer j
17
18     CSMS$DISTRIBUTE(dh, 1, 2) BEGIN
19     real, allocatable :: f(:, :)
20     real, allocatable :: df(:, :)
21     CSMS$DISTRIBUTE END
22
23
24     CSMS$CREATE_DECOMP(dh, <im+2, jm+2>, <2,2> : <PERIODIC, PERIODIC>)
25
26     allocate(f (0:im+1, 0:jm+1))
27     allocate(df(0:im+1, 0:jm+1))
28
29     CSMS$PARALLEL(dh, <i>, <j>) BEGIN
30
31     CSMS$SERIAL BEGIN
32     do j = 0, jm+1
33     do i = 0, im + 1
34     f(i,j) = sqrt(real(i+j)) - sqrt(real(i)) - sqrt(real(j))
35     end do
36     end do
37     CSMS$SERIAL END
38
39     iter = 0
40     global_error = 1.0
41
42     do while ((global_error .gt. tolerance) .and. (iter .lt. 1000))
43     iter = iter + 1
44
45     global_error = 0.0
46
47     CSMS$EXCHANGE(f)
48
49     CSMS$GLOBAL_INDEX(1) BEGIN
50     do j = 1, jm
51     f(0 ,j) = f(im, j)
52     f(im+1,j) = f(1 , j)
```

```

52         end do
53 CSMS$GLOBAL_INDEX END
54
55 CSMS$GLOBAL_INDEX(2) BEGIN
56     do i = 1, im
57         f(i,0 ) = f(i, jm)
58         f(i,jm+1) = f(i, 1 )
59     end do
60 CSMS$GLOBAL_INDEX END
61
62     do j = 1, jm
63         do i = 1, im
64             df(i,j) = 0.25*(f(i-1,j) + f(i+1,j) + f(i,j-1) + f(i,j+1))
65 &             - f(i,j)
66         end do
67     end do
68
69     do j = 1, jm
70         do i = 1, im
71             if (global_error .lt. abs(df(i,j))) then
72                 global_error = abs(df(i,j))
73
74             end if
75         end do
76     end do
77
78
79     do j = 1, jm
80         do i = 1, im
81             f(i,j) = f(i,j) + df(i,j)
82         end do
83     end do
84
85 CSMS$REDUCE(global_error, MAX)
86
87     end do
88
89     print *, 'Ended with iter : ', iter
90     print *, 'Global_error: ', global_error
91
92 CSMS$PARALLEL END
93
94     end

```

**Example 9-1: Version of the Laplace program with periodic boundary conditions.**

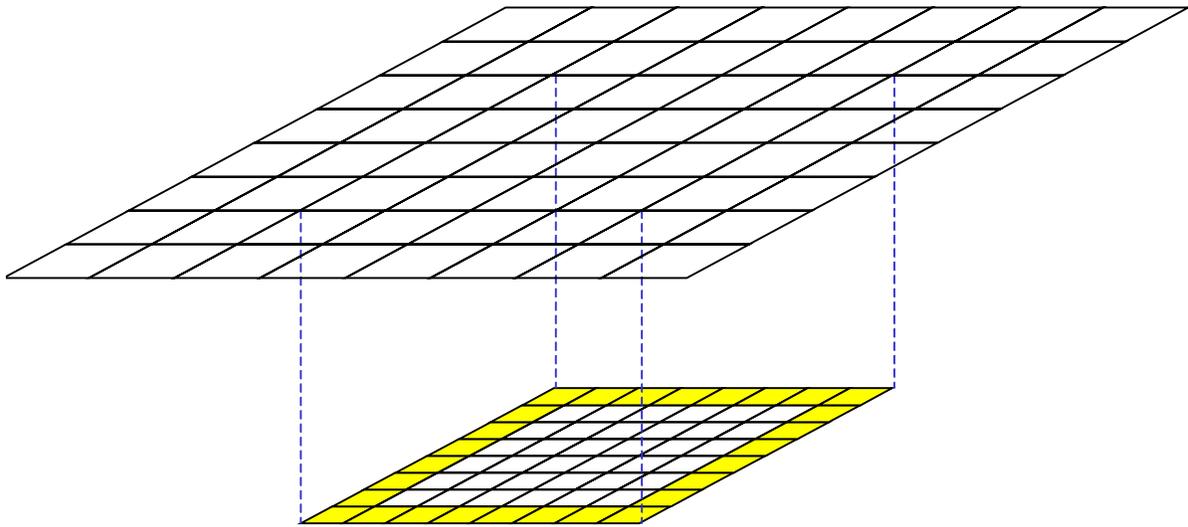


**Figure 9-1: Exchanges in a periodic model.**

With the model boundary halo regions properly updated, the periodic boundary conditions shown in lines 48-60 will produce the correct answer. Notice that even in the dynamic memory case, the local and global indices are not identical as shown in Figure 9-1. With *im* set to 9, local index -1 on process P1 corresponds to global index *im+1* (10). Local index 11 on process P3 corresponds to global index 0. However, explicit index translation using TO\_LOCAL and TO\_GLOBAL is still not needed in a periodic code that uses dynamic memory because index differences only occur at the periodic boundaries where GLOBAL\_INDEX will automatically handle any needed translations.

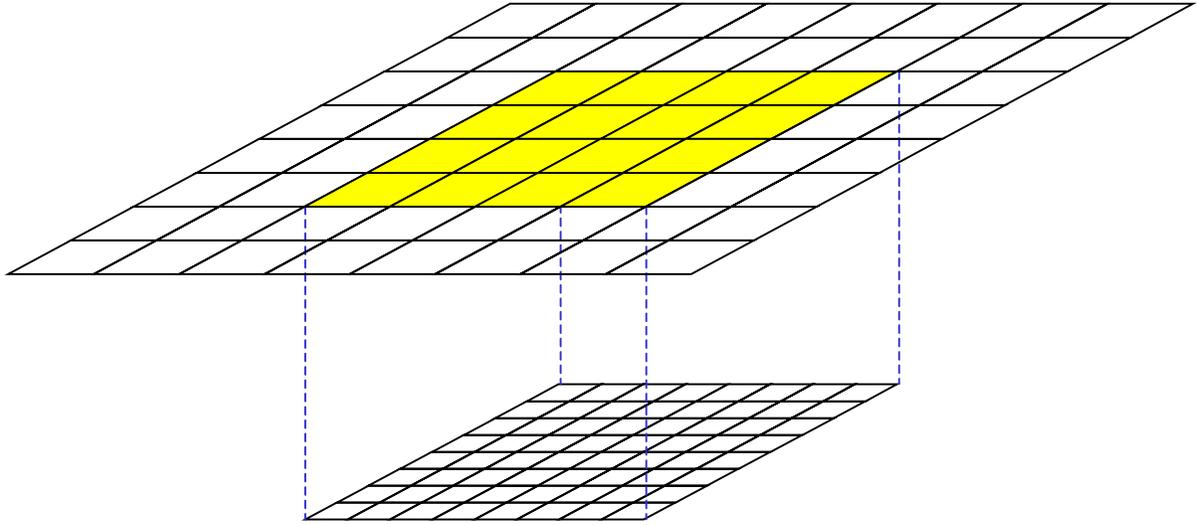
## 10 Nesting and Coupling: Transfer-Interpolation

In some programs it is necessary to interpolate data from one grid onto another grid. This kind of “inter-grid interpolation” occurs, for example, in nested atmospheric codes that use mesh refinement techniques to improve resolution in critical areas. Data from a coarse “parent” grid may be used to compute boundaries of a fine “child” grid (see Figure 10-1). Conversely, data from a fine “child” grid may be used to compute overlapping values in its coarse “parent” (see Figure 10-2). Also, it may be necessary to reconcile overlapping portions of two “sibling” grids (see Figure 10-3). In addition, inter-grid interpolation is used to couple models that are based on different grids (see Figure 10-4). A common example is coupling of atmospheric and oceanic models.

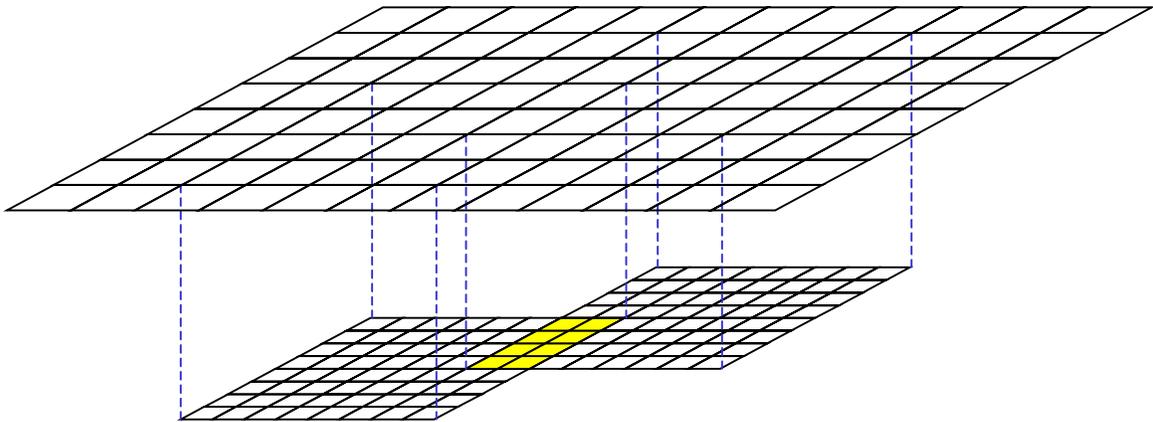


**Figure 10-1: Computing boundary points in a high-resolution “child” grid from points in a low-resolution “parent” grid.**

In a parallel program, inter-grid interpolations often require inter-process communication that can be quite complex. SMS encapsulates this complexity by combining both interpolation and communication into a single “*transfer-interpolation*” operation. SMS transfer-interpolation can be used to parallelize any inter-grid interpolation that uses weighted sums of values in a source grid to compute values in a destination grid.



**Figure 10-2: Computing overlapping points in a low-resolution “parent” grid from points in a high-resolution “child” grid.**



**Figure 10-3: Reconciling overlapping points in “sibling” grids.**

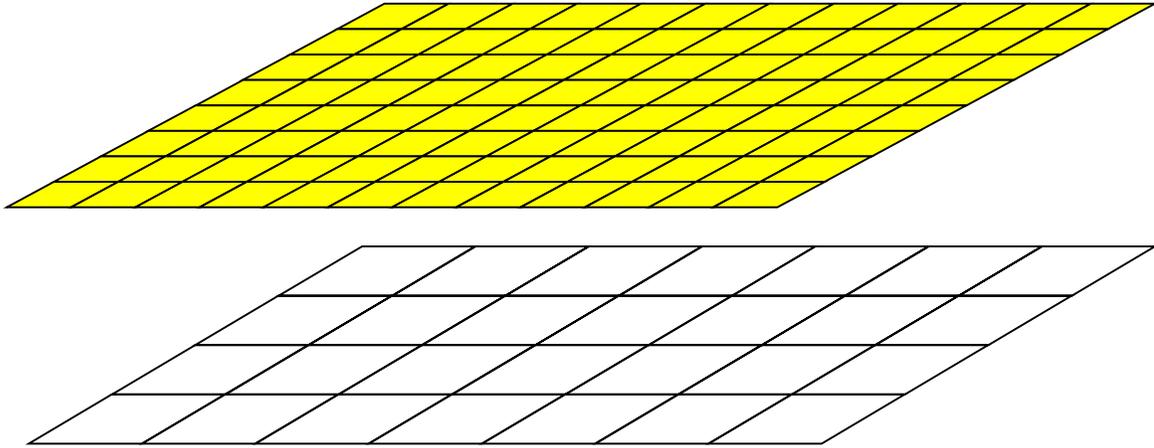


Figure 10-4: Coupling different grids.

## 10.1 Using SET\_TRANSFER\_INTERPOLATE to Define Interpolations

The SET\_TRANSFER\_INTERPOLATE directive is used to specify any inter-grid interpolation scheme that can be expressed as weighted sums. Before SET\_TRANSFER\_INTERPOLATE can be used, the serial code must be transformed so that source grid coordinates, destination grid coordinates, and weights are each stored in arrays. Code in this form will be referred to as “*stencil-oriented*”.

### 10.1.1 Transforming Serial Code to “Stencil-Oriented” Form

Example 10-1 shows a simple stencil-oriented serial code fragment that is used to compute elements of fine-grid array *X\_FINE* from weighted sums of elements of coarse-grid array *X\_COARSE*. In this code, *NUM\_FINE\_POINTS* is the number of elements to be computed in array *X\_FINE*. *MAX\_STENCIL\_POINTS* is the maximum number of elements from *X\_COARSE* that will be used to compute any element in *X\_FINE* by weighted sum. The coordinates of each element of *X\_FINE* to be computed are stored in array *fine\_indices*. The coordinates of each element of *X\_COARSE* that are used to compute each element of *X\_FINE* are stored in array *stencil\_indices*. The corresponding weights are stored in array *stencil\_weights*. (Initializations of *fine\_indices*, *stencil\_indices*, and *stencil\_weights* are not shown). The actual weighted-sum computations are performed in the loop beginning at line 13. Figure 10-5 illustrates details of the weighted sum used to compute one of the elements of *X\_FINE*.

```

1      real X_COARSE(imc,jmc)
2      real X_FINE(imf,jmf)
3
4      integer   fine_indices(2,NUM_FINE_POINTS)
5      integer stencil_indices(2,MAX_STENCIL_POINTS,
6      &          NUM_FINE_POINTS)
7      real      stencil_weights(MAX_STENCIL_POINTS,
8      &          NUM_FINE_POINTS)
9

```

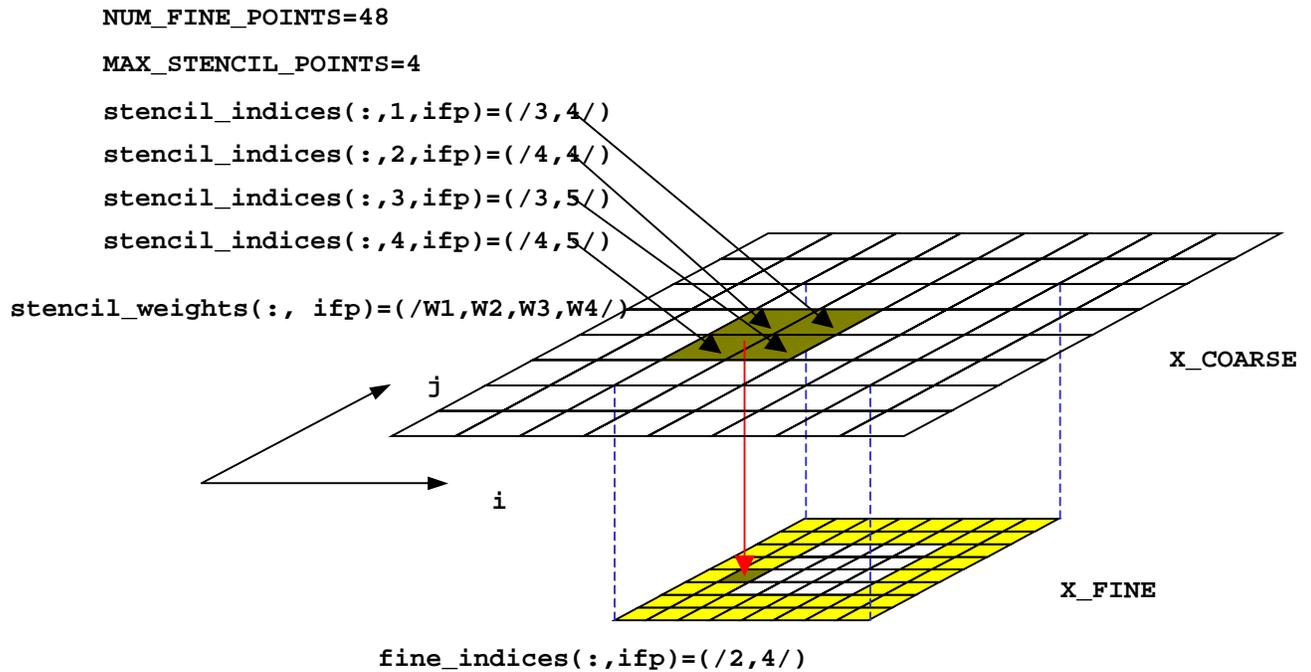
```

10 ...Initialize fine_indices, stencil_indices, and stencil_weights...
11
12 C NOTE "stencil-oriented" form of serial interpolation code
13   do ifp = 1,NUM_FINE_POINTS
14     i = fine_indices(1,ifp)
15     j = fine_indices(2,ifp)
16     X_FINE(i,j) = 0.0
17     do icp = 1,MAX_STENCIL_POINTS
18       X_FINE(i,j) = X_FINE(i,j) +
19         &         (stencil_weights(icp,ifp) *
20         &         X_COARSE(stencil_indices(1,icp,ifp),
21         &         stencil_indices(2,icp,ifp)))
22     enddo
23   enddo

```

**Example 10-1: Serial inter-grid interpolation code in “stencil-oriented” form.**

$$X\_FINE(2,4) = W1 * X\_COARSE(3,4) + W2 * X\_COARSE(4,4) + W3 * X\_COARSE(3,5) + W4 * X\_COARSE(4,5)$$



**Figure 10-5: Detail of coarse-to-fine interpolation computations. Integer *ifp* is the loop index used on line 13 of Example 10-1.**

### 10.1.2 Defining an Interpolation

Once the serial code is in “stencil-oriented” form, the SET\_TRANSFER\_INTERPOLATION directive can be used to define inter-grid interpolations. Example 10-2 shows how this is done for the inter-grid interpolation introduced in Example 10-1. Arrays *X\_COARSE* and *X\_FINE* are distributed using decompositions *dhCoarse* and *dhFine*, respectively. The SET\_TRANSFER\_INTERPOLATION directive begins on line 17.

```

1 CSMS$DISTRIBUTE(dhCoarse, <imc>, <jmc>) BEGIN
2     real X_COARSE(imc, jmc)
3 CSMS$DISTRIBUTE END
4 CSMS$DISTRIBUTE(dhFine, <imf>, <jmf>) BEGIN
5     real X_FINE(imf, jmf)
6 CSMS$DISTRIBUTE END
7
8     integer    fine_indices(2, NUM_FINE_POINTS)
9     integer stencil_indices(2, MAX_STENCIL_POINTS,
10    &                NUM_FINE_POINTS)
11     real      stencil_weights(MAX_STENCIL_POINTS,
12    &                NUM_FINE_POINTS)
13
14 CSMS$PARALLEL(dhFine, <i>, <j>) BEGIN
15 ...Initialize fine_indices, stencil_indices, and stencil_weights...
16
17 CSMS$SET_TRANSFER_INTERPOLATION( dhCoarse, dhFine, 2,
18 CSMS$> NUM_FINE_POINTS, fine_indices, MAX_STENCIL_POINTS,
19 CSMS$> stencil_indices, stencil_weights, INTERP_C F)
20
21 CSMS$TRANSFER( <X_COARSE, X_FINE : INTERP_C F> ) BEGIN
22 C NOTE "stencil-oriented" form of serial interpolation code
23     do ifp = 1, NUM_FINE_POINTS
24         i = fine_indices(1, ifp)
25         j = fine_indices(2, ifp)
26         X_FINE(i, j) = 0.0
27         do icp = 1, MAX_STENCIL_POINTS
28             X_FINE(i, j) = X_FINE(i, j) +
29             &                (stencil_weights(icp, ifp) *
30             &                X_COARSE(stencil_indices(1, icp, ifp),
31             &                stencil_indices(2, icp, ifp)))
32         enddo
33     enddo
34 CSMS$TRANSFER END

```

**Example 10-2:** The SET\_TRANSFER\_INTERPOLATION directive is used to define the inter-grid interpolation introduced in Example 10-1. The TRANSFER directive performs the computations and inter-process communications required to interpolate between grids.

The first two arguments in the SET\_TRANSFER\_INTERPOLATION directive are the names of the decompositions used by the source and destination arrays. The third argument is the rank of the interpolation, which cannot be larger than the rank of the source or destination arrays. The fourth argument is the number of elements in the destination array that will be computed during interpolation and the fifth argument is a list of coordinates for each element. The maximum number of elements from the source array that will be used to compute any element of the destination array is indicated by the sixth argument. The seventh argument is an array of lists of coordinates of elements from the source array that will be used to compute each element of the destination array. The weights used to interpolate each element of the destination array are indicated by the eighth argument. Finally, the last argument is a user-defined name for the interpolation. This name is used again in the TRANSFER directive to select the interpolation scheme. Finally, the code that initializes *fine\_indices*, *stencil\_indices*, and *stencil\_weights* must appear inside a PARALLEL directive (see line 14 of Example 10-2).

Note that any array dimension may be interpolated whether or not it is decomposed. The independence of interpolated dimensions and decomposed dimensions allows full

flexibility to support any weighted-sum interpolation scheme with any decomposition. However, to minimize inter-process communication, interpolation of non-decomposed dimensions should be decoupled and handled separately whenever possible (see Example 10-4).

## 10.2 Using TRANSFER to Interpolate Between Grids

The TRANSFER directive has special syntax for inter-grid interpolations as shown on line 21 of Example 10-2. The source and destination arrays (*X\_COARSE* and *X\_FINE*) are enclosed in angle brackets (see Section 6). In addition, a third argument (*INTERP\_C\_F*) references the name of the interpolation, defined by the SET\_TRANSFER\_INTERPOLATION directive (line 19 of Example 10-2). Note that the serial code enclosed between the TRANSFER “BEGIN” and “END” must include all interpolation computations to ensure that the serial code works correctly. During source code translation, code between the “BEGIN” and “END” is replaced with its parallel equivalent.

Transfer-interpolation operations are aggregated to reduce latency in the same way aggregation is done for standard transfers. Note that aggregation of multiple interpolations is done even if the source or destination arrays are not decomposed in the same way or if different interpolations are used. Aggregation of transfer-interpolation operations is illustrated in Example 10-3 for a simple case where arrays *Y\_COARSE* and *Y\_FINE* are decomposed like *X\_COARSE* and *X\_FINE*. Aggregation is achieved by combining the transfer-interpolations into a single TRANSFER directive (lines 1 and 2).

```

1 CSMS$TRANSFER( <X_COARSE, X_FINE : INTERP_C_F>,
2 CSMS$>          <Y_COARSE, Y_FINE : INTERP_C_F>) BEGIN
3     do ifp = 1,NUM_FINE_POINTS
4         i = fine_indices(1,ifp)
5         j = fine_indices(2,ifp)
6         X_FINE(i,j) = 0.0
7         Y_FINE(i,j) = 0.0
8         do icp = 1,MAX_STENCIL_POINTS
9             X_FINE(i,j) = X_FINE(i,j) +
10            &          (stencil_weights(icp,ifp) *
11            &          X_COARSE(stencil_indices(1,icp,ifp),
12            &          stencil_indices(2,icp,ifp)))
13            Y_FINE(i,j) = Y_FINE(i,j) +
14            &          (stencil_weights(icp,ifp) *
15            &          Y_COARSE(stencil_indices(1,icp,ifp),
16            &          stencil_indices(2,icp,ifp)))
17         enddo
18     enddo
19 CSMS$TRANSFER END

```

**Example 10-3:** This form of the TRANSFER performs inter-grid interpolations to compute elements of both *X\_FINE* and *Y\_FINE*. When a single TRANSFER directive is used for multiple interpolations, messages are aggregated. This will improve performance on some machines by reducing message latency. Note that the serial code enclosed by the TRANSFER directive must perform interpolation computations for all destination arrays.

In some circumstances, it is convenient to use an interpolation for computations on arrays with higher rank. For example, suppose that arrays *X\_COARSE* and *X\_FINE* from

Example 10-2 each had a third, non-decomposed dimension with size *km*. Furthermore, suppose that interpolation between *X\_COARSE* and *X\_FINE* was independent of the new third dimension. The same interpolation, *INTERP\_C\_F* could still be used as shown in Example 10-4. The syntax ***INTERP\_C\_F(1,2)*** in the TRANSFER directive on line 1 indicates that the first two dimensions of the arrays correspond to the first two dimensions of the interpolation. SMS will automatically replicate interpolation computations in the third array dimension.

```

1 CSMS$TRANSFER( <X_COARSE, X_FINE : INTERP_C_F(1,2)> ) BEGIN
2     do k = 1,km
3         do ifp = 1,NUM_FINE_POINTS
4             i = fine_indices(1,ifp)
5             j = fine_indices(2,ifp)
6             X_FINE(i,j,k) = 0.0
7             do icp = 1,MAX_STENCIL_POINTS
8                 X_FINE(i,j,k) = X_FINE(i,j,k) +
9                 &             (stencil_weights(icp,ifp) *
10                &             X_COARSE(stencil_indices(1,icp,ifp),
11                &             stencil_indices(2,icp,ifp),k))
12             enddo
13         enddo
14     enddo
15 CSMS$TRANSFER END

```

**Example 10-4:** Interpolations may be used for arrays with higher rank. Here, interpolation of 3D array *X\_FINE* from 3D array *X\_COARSE* is done using 2D decomposition *INTERP\_C\_F* (defined on lines 14-16 of Example 10-2).

### 10.3 Using SET\_NEST\_LEVEL to Switch Between Grids

In the examples shown so far, the source and destination decompositions have had different names (*dhCoarse* and *dhFine* in Example 10-2). However, it is often desirable to define a single decomposition name and use it to describe both the source and destination decompositions. In Example 10-5, subroutine *SOLVE* is called twice, once with *X\_COARSE* (line 29) and once with *X\_FINE* (line 36). However, *X\_COARSE* and *X\_FINE* are decomposed differently. To avoid passing SMS-specific decomposition information through the interface of subroutine *SOLVE* (which would unnecessarily complicate the serial code), SMS allows a single decomposition name to refer to more than one decomposition. Individually, these component decompositions are called “*nests*”.

In Example 10-5, decomposition *dh* is declared to have two nests by using the syntax “***dh(2)***” in the DECLARE\_DECOMP directive (line 4). The nests are then referred to individually using “***dh(1)***” (coarse nest) and “***dh(2)***” (fine nest) on lines 10, 11, 21, and 24. SMS directives inside subroutine *SOLVE* do not explicitly reference a specific nest (***dh(1)*** or ***dh(2)***). Instead, the nest level is selected before each call to *SOLVE* by SET\_NEST\_LEVEL directives on lines 28 and 35. Unlike most other SMS directives, the effects of the SET\_NEST\_LEVEL directive carry over into called subroutines. The selected nest will remain selected until another SET\_NEST\_LEVEL directive is encountered.

```
1 [include file params.incl]
```

```

2         integer imc,jmc,imf,jmf
3         common /sizes_com/ imc,jmc,imf,jmf
4 CSMS$DECLARE_DECOMP(dh(2), 2)
5
6
7 [main program source file]
8 C ...
9         include 'params.inc'
10 CSMS$CREATE_DECOMP(dh(1), <imc, jmc>, <3,3>)
11 CSMS$CREATE_DECOMP(dh(2), <imf, jmf>, <3,3>)
12 C ...
13 CSMS$SET_TRANSFER_INTERPOLATION( dh(1), dh(2), 2,
14 CSMS$> NUM_FINE_POINTS, fine_indices, MAX_STENCIL_POINTS,
15 CSMS$> stencil_indices, stencil_weights, INTERP_C_F)
16
17
18 [subroutine source file]
19 C ...
20         include 'params.inc'
21 CSMS$DISTRIBUTE(dh(1), <imc>, <jmc>) BEGIN
22         real X_COARSE(imc,jmc)
23 CSMS$DISTRIBUTE END
24 CSMS$DISTRIBUTE(dh(2), <imf>, <jmf>) BEGIN
25         real X_FINE(imf,jmf)
26 CSMS$DISTRIBUTE END
27
28 CSMS$SET_NEST_LEVEL( dh, 1 )
29         CALL SOLVE(X_COARSE,imc,jmc)
30
31 CSMS$TRANSFER( <X_COARSE, X_FINE : INTERP_C_F> ) BEGIN
32 C ...
33 CSMS$TRANSFER END
34
35 CSMS$SET_NEST_LEVEL( dh, 2 )
36         CALL SOLVE(X_FINE,imf,jmf)
37
38
39 [another subroutine source file]
40         SUBROUTINE SOLVE(X,im,jm)
41         include 'params.inc'
42 CSMS$DISTRIBUTE(dh, <im>, <jm>) BEGIN
43         real X(im,jm)
44 CSMS$DISTRIBUTE END
45         integer i,j
46 CSMS$PARALLEL(dh, <i>, <j>) BEGIN
47 C ...
48 CSMS$PARALLEL END
49         END

```

**Example 10-5:** Subroutine *SOLVE* must be called with either *X\_COARSE* or *X\_FINE* as the first argument. Since these two arrays are decomposed differently, the best way to do this is to create a single decomposition with two “nests”, *dh(2)*. The *SET\_NEST\_LEVEL* directive selects which nest is used.

## 11 I/O

One of the most powerful features of SMS is its ability to support most types of I/O without directives. The fact that communication patterns for I/O of decomposed and non-decomposed arrays differ is hidden from the programmer. SMS automatically generates the communication needed to read or write data to or from disk in the same sequence as the serial code would have done it regardless of the number processes used. Unformatted I/O is discussed in Section 11.1 while formatted I/O is covered in Section 11.2. Printing in a parallel program often requires additional decisions from the programmer. SMS allows the programmer to make these decisions by providing several print modes that are introduced in Section 11.2.2. Finally, Section 11.3 discusses methods for improving I/O performance.

### 11.1 Unformatted I/O

Figure 11-1 illustrates dependencies for read and write of a simple one-dimensional decomposed array. During a read, data from a single file must be parceled out to each process. This type of communication pattern is called "scatter". During a write, data from each process must be combined in the "proper order" and written to disk. This type of communication pattern is a different form of "gather" than the form used for TRANSFER and bit-wise exact REDUCE. Instead of gathering data into a global variable that is replicated in memory on all processes, data is gathered into a single file on disk. "Proper order" means the data must be read from or written to disk in the same sequence that the serial code uses. Though it appears quite simple in Figure 11-1, the data reorganization required to match serial ordering in files can be quite complex, especially for two-dimensional decompositions or when the decomposed arrays have halo regions (Figure 11-2). Additionally, when variables being input have halo regions associated with them, these regions are automatically updated by SMS.

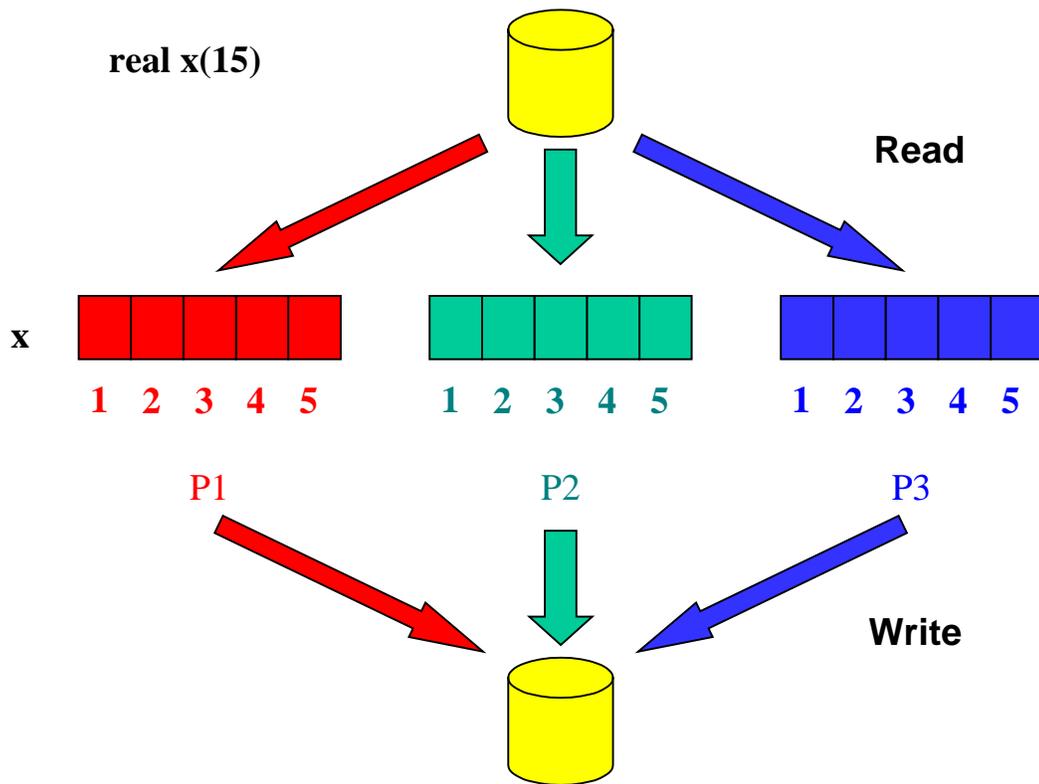


Figure 11-1: Conceptual schematic of the input and output of a decomposed array. On input, one process reads the global data from disk. The appropriate sections of the global array are then “scattered” to each process. On output, the decomposed data are gathered into a global array and then written to disk. The underlying implementation may use a more efficient scheme on some machines.

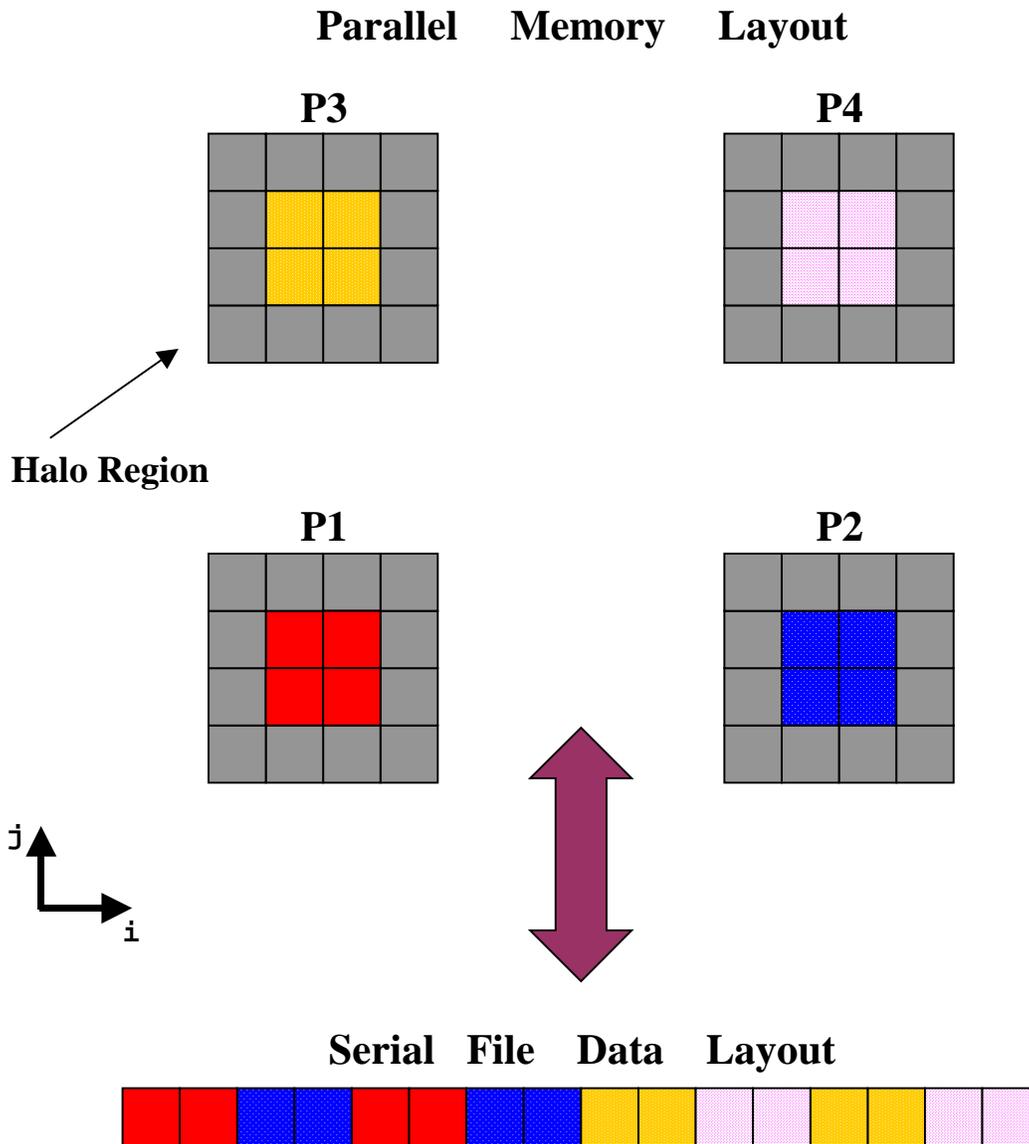
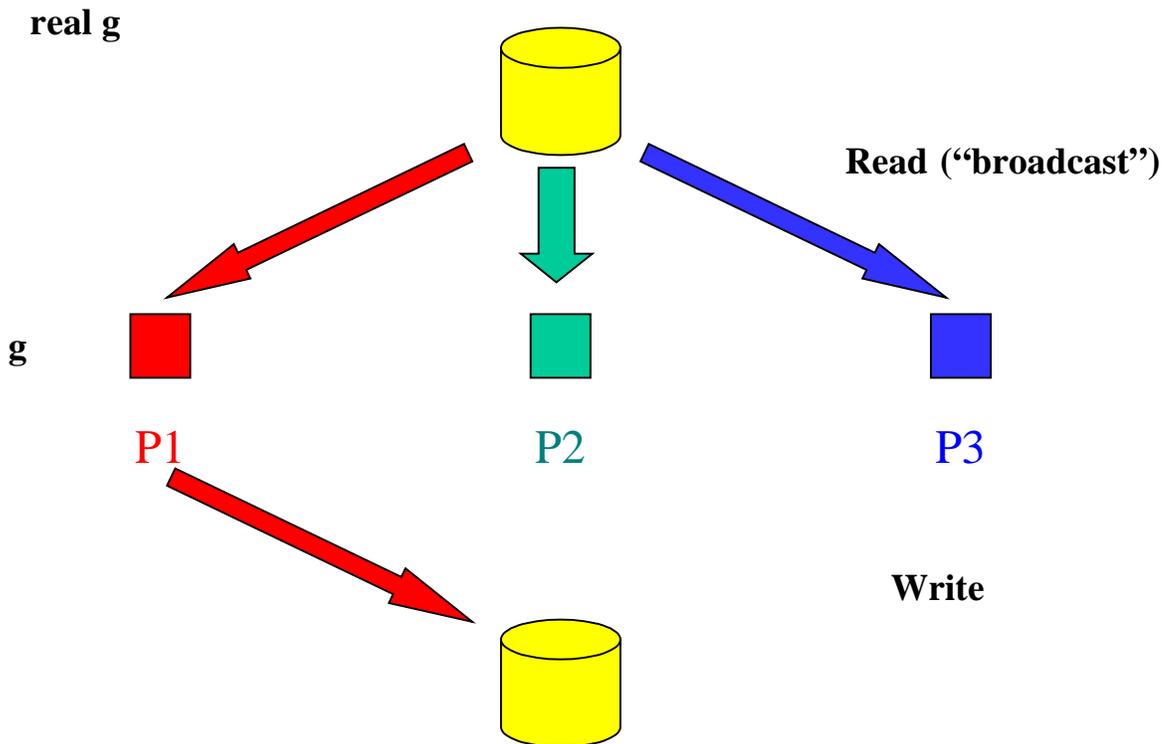


Figure 11-2: Conceptual schematic of the re-ordering required to write and read two-dimensionally decomposed data to disk in the same order as the serial code would write it. Special care has to be taken to write the only the interior of each process-local domain and not the halo data. The halo regions are filled during the read operations.

Figure 11-3 illustrates dependencies for read and write of a non-decomposed variable. During a read, a copy of data from a single file must be sent to each process. This type of communication pattern is called "broadcast". During write, it is only necessary to write data from a single process because each process should have an identical copy of the variable.



**Figure 11-3: Schematic of the input and output of a non-decomposed array. On input, one process reads the data from disk. The data are then replicated on all other processes. On output, a single process writes the data to disk.**

Example 11-1 demonstrates unformatted I/O of both decomposed and non-decomposed variables. SMS automatically translates the read (line 32) and write (line 28) statements for both decomposed arrays  $x$  and  $y$  and non-decomposed scalar  $scale$  to the appropriate parallel I/O operations. When automatically generating parallel I/O operations, SMS uses information in the DISTRIBUTE directives to determine how to scatter, gather, or broadcast data. Notice that any types of decomposed or non-decomposed variables can be mixed in a single write or read statement. It is not necessary to reorganize existing serial read or write statements to take advantage of automatic parallelization by SMS.

```
[Include file: io.inc]
1      integer im, jm
2      common /sizes_com/ im, jm
3      CSMS$DECLARE_DECOMP(DECOMP_IJ, 2)
```

```
[Source file: binary.f]
1      program binary_io
2      include 'io.inc'
3      im = 10
4      jm = 5
```

```

5  CSMS$CREATE_DECOMP(DECOMP_IJ, <im,jm>, <1,0>)
6      call write_data
7      end
8
9      subroutine write_data
10     include 'io.inc'
11     integer i, j
12     real scale
13  CSMS$DISTRIBUTE(DECOMP_IJ, <im>, <jm>) BEGIN
14     integer x(im,jm), y(im,jm)
15  CSMS$DISTRIBUTE END
16  CSMS$PARALLEL(DECOMP_IJ,<i>,<j>) BEGIN
17     do j=1,jm
18     do i=1,im
19  CSMS$TO_GLOBAL(<1,i>, <2,j>) BEGIN
20     x(i,j) = (100 * i) + j
21     y(i,j) = mod(i,2)
22  CSMS$TO_GLOBAL END
23     end do
24     end do
25  CSMS$PARALLEL END
26     scale = -1.0
27     open (17,file='iol_out.dat',form='unformatted')
28     write (17) x, y, scale
29     close (17)
30
31     open (18,file='iol_out.dat',form='unformatted')
32     read (18) x, y, scale
33     close (18)
34     return
35     end

```

**Example 11-1: This program does output and input of both decomposed and non-decomposed data. No additional directives are required for the correct I/O to be performed, regardless of the number of processes.**

By default, SMS assumes unformatted files are stored in native FORTRAN binary format. The default behavior can be modified using the following environment variables:

```

SMS_READ_FORMAT
SMS_WRITE_FORMAT

```

The currently available (case insensitive) formats are:

```

IBM
SUN
SGI
FUJITSU
HP
DEC
COMPAQ
IA32
MPI
MPI_EXTERNAL
EXTERNAL
SMS

```

In some cases, file formats with different names are actually the same format. For example, *SGI* and *SUN* are really the same. Note that *MPI*, *MPI\_EXTERNAL*, *EXTERNAL*, and *SMS* all refer to the portable MPI I/O external format. The advantage to

using this format is that any file written by an SMS program may be read by any other SMS program on any other machine. This is true regardless of the number of processes used on either machine because SMS preserves serial data ordering. To convert data files from one format to another, simply write a serial program that reads and writes the data, compile and link with SMS and then set the above environment variables appropriately.

## 11.2 Formatted I/O

SMS supports the formatted input and output of non-decomposed data without directives. Also supported is unformatted I/O via namelist of both decomposed and non-decomposed data. Formatted I/O of decomposed variables by means other than namelist is not currently supported, so code segments that include this kind of I/O must be enclosed by a SERIAL directive (see Section 8). In a future release of SMS, this use of SERIAL will become optional.

### 11.2.1 Formatted Input

Formatted input including namelists is handled automatically by SMS. The user does not need to add any directives. The only caveat is that input variables cannot be decomposed arrays unless a namelist is used. In this case, a work-around is to enclose the formatted read statements within a SERIAL directive. Since formatted reads typically occur infrequently during the course of a model run, this approach usually does not incur a significant performance penalty.

### 11.2.2 Formatted Output

The simple task of printing a message on the screen becomes complicated in a parallel program. Consider the following simple print statement:

```
print *, 'HELLO'
```

There are no clear standard definitions of what will appear on the screen when a "parallel" print statement is executed. Will each process print a separate message? Will the separate messages appear on different lines on the screen? Will all processes be forced to wait until the print is complete before useful work can continue? If the statement were executed on three processes, we might see any of the following output:

```
HELLO
```

```
HELLO  
HELLO  
HELLO
```

```
HHHEEELLLLLLLOOO
```

```
HELLHEHLEOLOLLO
```

During the brief history of parallel computing, each of these possibilities has been implemented on at least one parallel machine.

SMS simplifies this situation by providing three "print modes" that allow the user to control the behavior of parallel print. The modes are default, ASYNC, and ORDERED. These modes are selected using the PRINT\_MODE directive. The PRINT\_MODE directive may only be used for formatted output to stdout. This may be accomplished either by using *PRINT* statements or *WRITE* statements when the unit number is either 6 or \*. (Note that units 0, 5, and 6 should not be opened in an SMS program because many Fortran compilers behave strangely when these units are connected to named files.) The print modes are illustrated in Example 11-2.

```

1      program print_modes
2      implicit none
3      integer, parameter :: im = 12
4      integer xmax, i
5
6      CSMS$DECLARE_DECOMP(dh, 1)
7
8      CSMS$DISTRIBUTE(dh, 1) BEGIN
9          integer, allocatable :: x(:)
10     CSMS$DISTRIBUTE END
11
12     CSMS$CREATE_DECOMP(dh, <im>, <0>)
13         allocate(x(im))
14
15     CSMS$SERIAL BEGIN
16         do i = 1, im
17             x(i) = i
18         end do
19     CSMS$SERIAL END
20
21     CSMS$PARALLEL(dh,<i>) BEGIN
22         xmax = 0
23         do i = 1,im
24             xmax = max(xmax,x(i))
25
26             if (x(i) .ge. 12) then
27     CSMS$PRINT_MODE(ASYNC) BEGIN
28         print *,'WARNING:  x .ge. 12 !!  '
29     CSMS$PRINT_MODE END
30             endif
31
32         end do
33     CSMS$PARALLEL END
34
35     CSMS$PRINT_MODE(ORDERED) BEGIN
36     CSMS$INSERT      print *,'DEBUG:  local maximum value = ',xmax
37     CSMS$PRINT_MODE END
38
39     CSMS$REDUCE(xmax,MAX)
40
41     write( *,900) 'maximum value = ',xmax
42     900 format(a, i4)
43     end

```

**Example 11-2: Program that illustrates the various print modes supported by SMS.**

When the serial code versions *print\_modes* is run, the following is printed on the screen:

```
>> print_modes
```

```
WARNING: x .ge. 12 !!
maximum value = 12
```

When the parallel code is run on 1 process, the following is printed on the screen:

```
>> smsRun -np 1 print_modes_parallel
WARNING: x .ge. 12 !!
DEBUG: local maximum value = 12
maximum value = 12
```

For 4 processes:

```
>> smsRun -np 4 print_modes_parallel
WARNING: x .ge. 12 !!
DEBUG: local maximum value = 3
DEBUG: local maximum value = 6
DEBUG: local maximum value = 9
DEBUG: local maximum value = 12
maximum value = 12
```

The write statement on line 41 in Example 11-2, is printed using the default print mode. The default print mode is used for any print statement that is not enclosed by a PRINT\_MODE directive. It will cause the parallel code to print the same messages as the serial code in most cases. Only one system-dependent designated process (the “root” process) will execute the print statement; the others will skip it and can immediately continue with useful computations.

The print statement on line 36 is executed using the ORDERED print mode. The ORDERED print mode may only be selected using the PRINT\_MODE directive. This mode causes one message to be printed on the screen for each process and guarantees that the messages always appear in the same order. It is most useful for debugging purposes. However, in order to guarantee message ordering, no process can continue until all processes have executed the print statement. This means care must be taken that all processes will ALWAYS execute an ordered print or the program will hang. For example, suppose we use the ORDERED print mode at line 27:

```
        if (x(i) .ge. 12) then
CSMS$PRINT_MODE(ORDERED) BEGIN
        print *, 'WARNING: x .ge. 12 !! '
CSMS$PRINT_MODE END
        endif
```

In this case, we see the same results for the one-process run. However, the four-process run produces the following results:

```
>> smsRun -np 4 print_modes_parallel
DEBUG: local maximum value = 3
DEBUG: local maximum value = 6
DEBUG: local maximum value = 9
WARNING: x .ge. 12 !!
```

In this case, the program hangs (deadlocks) before the final message can be printed because the warning print statement is now an ordered-mode print that has been executed

by only one process. The program will wait forever for the other processes to enter this print statement. The default print mode is also inappropriate here because the warning message would not be printed if point 12 were not on the root process. Deadlock would not occur, but the warning message would not be printed.

The ASYNC (asynchronous) mode is the proper mode to use in cases like the printed warning statement on line 28 because there is no guarantee that all processes will execute the print statement. In this mode, one message will appear on the screen for each process that executes the print statement. Like the default mode, there is no process synchronization during asynchronous prints. As a result, ordering of print statements may vary from one run to the next when ASYNC mode is used. Like the ORDERED mode, the ASYNC print mode may only be selected using the PRINT\_MODE directive. For example, suppose we use the ASYNC mode for line 35 instead of ORDERED.

```
CSMS$PRINT_MODE(ASYNC) BEGIN
CSMS$INSERT      print *,'DEBUG: local maximum value = ',xmax
CSMS$PRINT_MODE END
```

Running with four processes two times might produce the following results:

```
>> smsRun -np 4 print_modes_parallel
DEBUG: local maximum value = 3
DEBUG: local maximum value = 6
DEBUG: local maximum value = 9
WARNING: x .ge. 12 !!
DEBUG: local maximum value = 12
maximum value = 12

>> smsRun -np 4 print_modes_parallel
DEBUG: local maximum value = 6
DEBUG: local maximum value = 3
WARNING: x .ge. 12 !!
DEBUG: local maximum value = 9
DEBUG: local maximum value = 12
maximum value = 12
```

Note that the ASYNC-mode prints can appear in any order and can even appear out-of-order with other non-ASYNC-mode prints. This can be confusing in some cases. Also, the ASYNC mode does not work properly when the SMS program is being run in “serverless” mode (see Section 11.3.3).

### 11.3 I/O Performance Tuning

This section discusses ways the user of SMS can optimize the I/O performance of their codes. These optimizations require a good understanding of how input and output operations are handled in SMS. Several environment variables that can be used to tune SMS I/O at run time are described below. The following table offers suggested values for these environment variables when an SMS program is run with an I/O server (see Section 11.3.2). When a server is not used (see Section 11.3.3), default values are recommended.)

CASE I: Input files will fit in server memory

	Server without Cachfers	Server with Cachfers
	-----	-----
SMS_RBS	Size of largest input file in bytes divided by (SMS_RBC-1)	Size of largest input file in bytes divided by (SMS_RBC-1)
SMS_RBC	Default (16)	Default (16)
SMS_WBS	Size of largest output file in bytes	Default
SMS_CLOSE_MODE	Require-flush	Require-flush
SMS_IOC_SIZE	N/A	Size of largest output file in bytes divided by the number of cache processes and multiplied by 2
SMS_RAN_RSTYLE	File	File

CASE II: Input files will NOT fit in server memory (only affects input)

	Server without Cachfers	Server with Cachfers
	-----	-----
SMS_RBS	Size of largest input variable in bytes divided by (SMS_RBC-1)	Size of largest input variable in bytes divided by (SMS_RBC-1)
SMS_RBC	Default (16)	Default (16)
SMS_RAN_RSTYLE	One-var	One-var

**Figure 11-4: Suggested values for SMS environment variables that affect I/O performance. “Cachfers” refers to SMS write-cache processes (see Section 11.3.5).**

### 11.3.1 General Guidelines

Two general guidelines should always be considered to improve both serial and parallel I/O performance. First, the user should do as little I/O as possible. Since I/O operations do not scale well, their effect on parallel performance will increase as the number of processes increase. What is an insignificant amount of run-time for 2 processes may be quite significant for 200 processes.

One optimization that can be very useful is to optionally turn off all print statements. Many serial codes already allow users to turn off some or even all print statements by setting a flag at run-time. This may speed up the serial code in some cases. The optimization is very useful in a parallel code where, on some machines, disabling prints can result in significant performance improvements. The following code fragment illustrates this common optimization:

```

if (print_enabled) then
  print *, 'whatever...'
endif

```

In this case, “print\_enabled” could be input through a namelist at the beginning of program execution.

A second general guideline to improve I/O performance is to combine I/O operations whenever possible. For example,

```
read(10) u
read(10) v
```

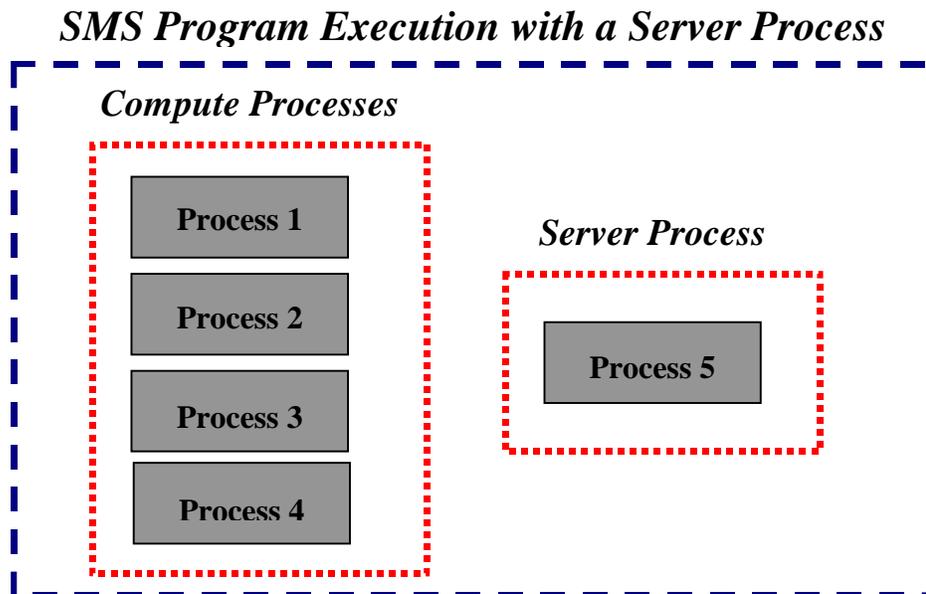
could be combined into a single read statement:

```
read(10) u,v
```

This will maximize the size of data blocks read from or written to disk and minimize I/O latency. Both unformatted and formatted statements should be considered for these optimizations.

### 11.3.2 The SMS Server Process

By default, SMS designates an additional process, called the server process, to manage the other processes and to handle all formatted and unformatted I/O operations. This allows computations to be done concurrent with I/O operations and can improve the overall performance of SMS program execution. Figure 11-5 illustrates a program run using four compute processes and a SMS server process.



**Figure 11-5:** In this example, four processes are requested to run the program. By default, an additional process, called the server process, will be used by SMS for process management and I/O operations.

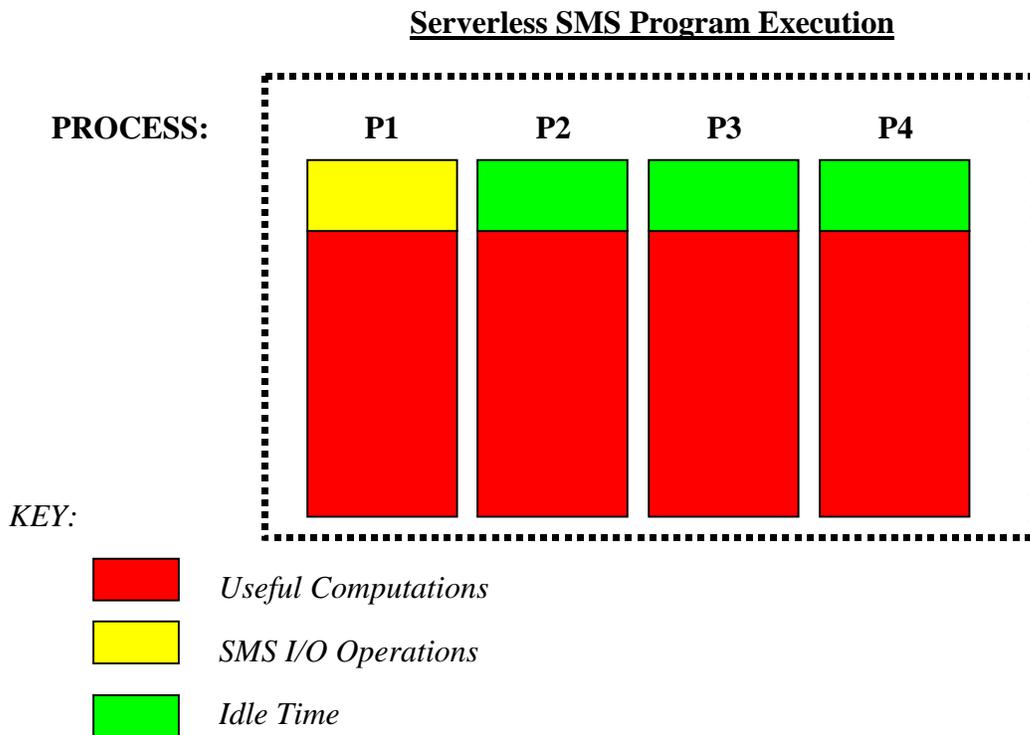
### 11.3.3 Serverless I/O

For small numbers of processors (less than 8), it may be beneficial to combine the server process functions with one of the computational processes. This type of operation is called serverless I/O and is illustrated in Figure 11-6.

If serverless I/O is used, the I/O functions that would normally be run on a separate process will be combined with one of the compute processes. Serverless SMS can be requested through an environment variable given by the command:

```
>> setenv SMS_SERVER_MODE serverless
```

On most machines, where there is a one-to-one correspondence between processes and processors, serverless I/O will improve performance by making one more processor available to do computations. However, when large numbers of processors are available, program execution will usually be faster when a server process is used.



**Figure 11-6:** An illustration of four SMS processes used to run a program without a server process. In this example, process P1 must handle both program computations and SMS server functions that include I/O operations. While these operations occur, the other processes will be idle.

### 11.3.4 The FLUSH\_OUTPUT Directive

During write operations, the I/O server process buffers the data to be output in memory, re-orders the decomposed data into serial order, and then writes it out in large blocks to

disk. By default, any write to disk will be delayed until the buffer is full or the file is closed. When this happens, buffers are "flushed" and their contents written to disk in large blocks. While buffers are being flushed, any processes requesting I/O services will have to wait until the flush operation is complete.

The environment variable `SMS_CLOSE_MODE` can be set to "require-flush" for improved user control of when SMS output buffers are flushed. Then, performance improvement can be gained by controlling when these buffers are flushed using the `FLUSH_OUTPUT` directive. This directive instructs the SMS I/O server process to flush the buffers immediately. If `FLUSH_OUTPUT` is placed so no other I/O requests are made during the flush operation, then no process will have to wait for the flush. If any I/O request is encountered, it must wait until the flush operation is complete thus minimizing the effectiveness of `FLUSH_OUTPUT`.

The following code fragment shows how this directive can be used:

```
        open (17,file='main_fields.dat',form='unformatted')
        write (17) u,v,w,p,t
        close (17)
c useful computations ...

        open (17,file='moisture.dat',form='unformatted')
        write (17) qs,qi,qj,qg,qw
        close (17)
CSMS$FLUSH_OUTPUT
c more useful computation ...
```

**Example 11-3: Proper placement of a `FLUSH_OUTPUT` directive.**

In this example, two files are written. As long as no other I/O (unformatted, formatted, or print) operations occur while the flush instruction is being processed, useful computations will proceed at full speed while data is simultaneously re-ordered and written to disk. This ability to overlap I/O with useful computation is key to achieving scalable I/O performance on many machines. However, any I/O statement that occurs soon after the flush operation will be sufficient to make the directive ineffective. For example, when a print statement appears just after the `FLUSH_OUTPUT`, it will force one of the processes to wait until the flush operation completes. Most likely, all other processes will eventually end up waiting for this process and useful computation will quickly come to a halt until the flush completes:

```
        open (17,file='diagnostics.dat',form='unformatted')
        write (17) x1,x2
        close (17)
CSMS$FLUSH_OUTPUT
        print *,'bad idea to print something here...'
... more useful computation
```

**Example 11-4: Improper placement of a `FLUSH_OUTPUT` directive.**

### 11.3.5 Improving Output Performance

To further increase the performance of output operations, two options are available. First, SMS allows the user to designate at run-time any number of processes to serve as write-cache processes. For example, Figure 11-7 illustrates a program that is run using twelve computational processes, two write-cache processes, and a server process.

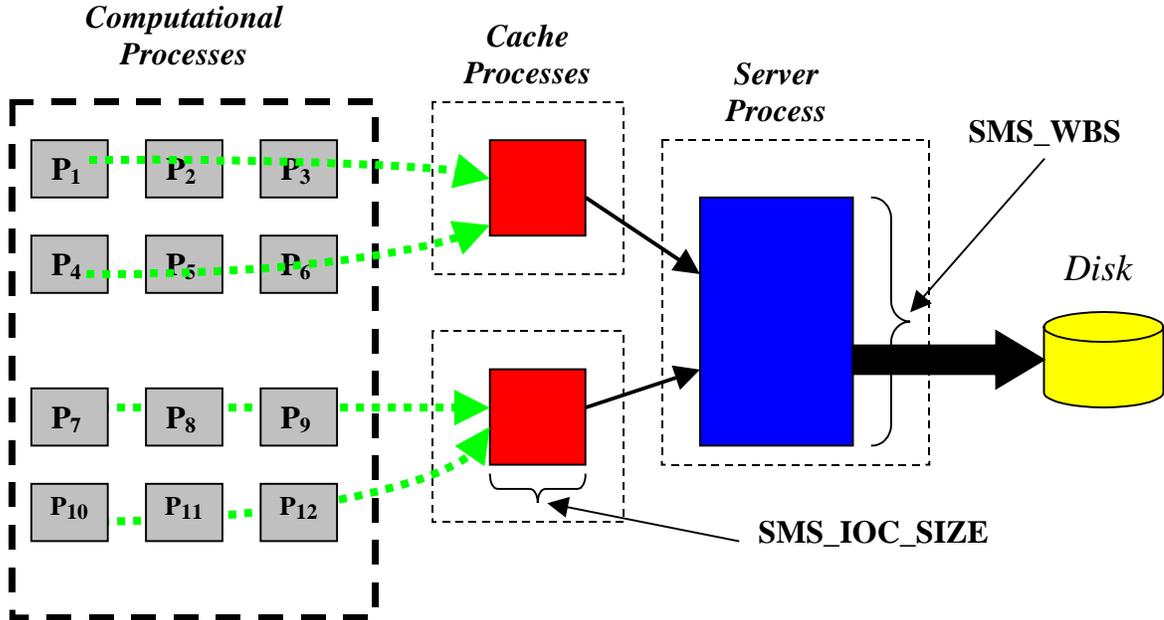


Figure 11-7: An illustration of SMS output when write-cache processes and a server process are used. SMS output operations pass data from the computational domain to the write-cache processes (if specified). Data is re-ordered on the write-cache processes before being passed through the server process to disk. The amount of memory allocated to the write-cache processes and the server process can be controlled using `SMS_IOC_SIZE` and `SMS_WBS` respectively.

The function of write-cache processes is to temporarily store data being output so it can be reordered and written to disk. The computational processes can write their data to multiple write-cache processes at high speeds and although the write-cache processes will proceed at relatively slow speeds, total execution time is not affected because disk writes can be done at the same time as computations. Further, write-cache processes provide more memory capacity to temporarily store the data before it is written to disk. The number of write-cache processes can be requested at run-time using the `-wc` option to the `smsRun` command. For example:

```
>> smsRun -np nprocs execname -wc <ncachers>
```

executes a program where `nprocs` is the number of computational processes, `execname` is the name of the executable, and `ncachers` is the number of write-cache processes to be used in the run. Refer to Section 12 more details about running an SMS program.

For optimal performance, there should be enough write-cache processes to store all data to be output at one time. By default, SMS allocates 8 Mbytes of memory for each write-

cache process. However, the environment variable `SMS_IOC_SIZE` is provided to allow the user to set the amount of memory (in bytes) they wish to allocate on each write-cache process. The command:

```
>> setenv SMS_IOC_SIZE 1000000
```

will allocate one million bytes of buffer space on each write-cache process. Since up to 50 percent of the buffer space can be lost to the overhead required to store the data segments, a recommended size for this field is double the size of the expected output. For example, assume we wish to output the following array:

```
real*4 big_array(100,200,300)
```

It will require 24 Mbytes of memory to output this array ( $4 \times 100 \times 200 \times 300$ ). This figure should be doubled to 48Mbytes to account for SMS overhead costs. If each write-cache process had 10 Mbytes of memory available for SMS caching, we would need to allocate five write-cache processes to output this array efficiently.

A second way to improve output performance is to change the size of the I/O buffer on the server process. By default the size of this buffer is 256 Mbytes, however this value can be changed through the SMS environment variable `SMS_WBS`. If write-cache processes are not used, then this variable should be set to the size of the largest output file when possible, otherwise output performance could degrade.

### 11.3.6 Improving Input Performance

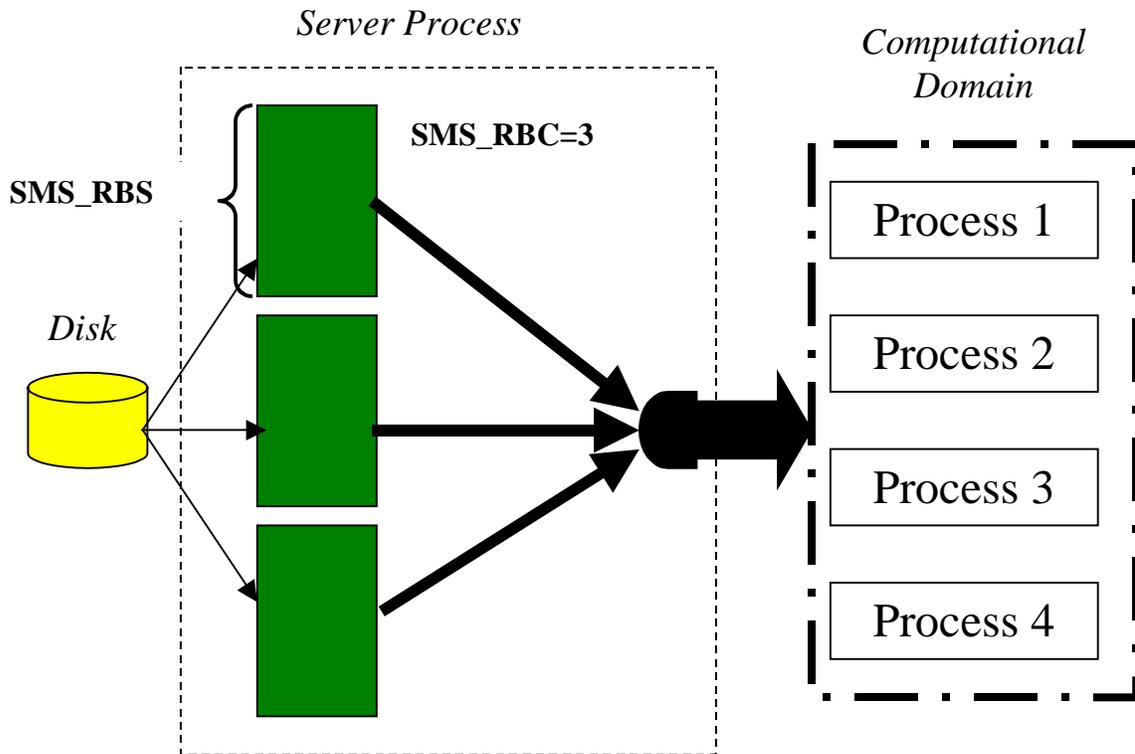
Only the server process is used to read formatted and unformatted input data; write-cache processes are not used. If data are decomposed, they are scattered to the other processes; otherwise, they are broadcast to the other processes.

By default, three environment variables can be used to control SMS input performance: `SMS_RBS`, `SMS_RBC`, and `SMS_RAN_RSTYLE`. `SMS_RBS` determines the size of each block that will be allocated to store input variables read from disk. `SMS_RBC` defines the number of blocks of size `SMS_RBS` that will be used for input. Finally, `SMS_RAN_STYLE` determines if files or individual variables will be input at one time. Figure 11-8 illustrates how these variables are used for input operations.

If a single file is input, the environment variable `SMS_RBS` should be set to the size of that file and `SMS_RBC` should be set to one. If multiple files (e.g. initial conditions and boundary conditions) are input with differing sizes, `SMS_RBS` should be set to a common factor of the size of each input file. For example suppose two files are required; an initial conditions file of size 53 Mbytes and a boundary conditions of size 16 Mbytes. An approximate common factor for these two files is 8 Mbytes ( $8 \times 2 = 16$ ,  $8 \times 6 = 54$ ). Therefore, good starting values would be: `SMS_RBS=8Mbytes`, `SMS_RBC=6`.

Using these variables, the total size of each input file should be considered when optimizing for performance. For example the execution of a program may be handled

with two files: input of initial conditions, followed by the input of boundary conditions. There should be sufficient memory on a single process to store the entire contents of each input file.



**Figure 11-8: All input will pass from disk, through the server process, to individual processes within the computational domain. Two SMS environment variables can be set to control the size of two data structures within the server process: the number of input buffers (SMS\_RBC) and the size of each buffer (SMS\_RBS).**

If not enough memory is available to store all input on a single process, SMS\_RAN\_RSTYLE should be set to “one-var”. This will force SMS to read each variable into a buffer that resides on the server process, transfer that data to the server process for distribution among the compute processes, and then read the next variable. In this case, the quantity:  $SMS\_RBS * (SMS\_RBC - 1)$  should be set to the size of the largest input variable.

The techniques described above are useful for reducing execution time when performance analysis indicates that run-time is limited by I/O time. Exact values of environment variables and number of write-cache processes are best determined by experimentation.

## 12 Program Termination

Parallel programs using the SMS run-time system require special handling to ensure all processes exit normally. An SMS control process is often used to manage the child processes that are spawned through the *smsRun* command to execute a program. Two types of program termination are supported by SMS: a normal exit and an abort. When a program exits normally, the SMS control process will wait until every processes' computations, communications and I/O are complete before exiting. A program abort will not guarantee the completion of outstanding operations or an orderly termination of processes.

### 12.1 Automatic Code Generation for Termination

By default, SMS automatically generates code to abort whenever a Fortran “stop” statement is encountered. SMS also generates a normal exit whenever a program “end” statement is encountered. Consider the following program:

```
program main

do ii=0, num_iter
  call time_steps(ii,status)
  if (status .eq. ABORT) then
    print *, ' Model Run failed at iteration: ',ii
    stop
  endif
enddo

print *, ' Model Run Successfully Completed'
stop
end
```

**Example 12-1: Automatic Code Generation by SMS will cause this program to always abort.**

Since the Fortran “stop” appears before the line before the end program statement, SMS will generate code to abort the parallel run. During code translation the following warning message will appear when source contains a Fortran stop statement:

**WARNING: Program abort detected.**

Since the intent of the original code in this case is to exit normally from the program, two actions can be taken to ensure this happens in the SMS-generated source. Either the second “stop” statement (above the “end”) should be removed, or the EXIT directive should be used as illustrated in the next section.

### 12.2 EXIT Directive

EXIT is used to control the run-time behavior of an SMS program. This directive, when inserted just before a “stop” statement, will instruct SMS to generate code to exit rather than abort. The proper placement of this directive is illustrated in Example 12-2. In this example, SMS will generate an ABORT at line 7 and a normal exit at line 12.

```

1      program main
2
3      do ii=0, num_iter
4          call time_steps(ii,status)
5          if (status .eq. ABORT) then
6              print *, ' Model Run failed at iteration: ',ii
7              stop
8          endif
9      enddo
10
11     print *, ' Model Run Successfully Completed'
12 CSMS$EXIT
13     stop
14     end

```

**Example 12-2: Using CSMS\$EXIT to override automatic translations**

### 12.3 MESSAGE Directive

MESSAGE, is used to send a message to the user at run-time and optionally terminate execution of the program when it is encountered. This directive is useful when the user wishes to avoid unnecessary parallelization of code they believe is never executed. Three run-time actions are available to the user of MESSAGE: ABORT, terminates execution after writing the given message to stderr, WARN writes the given text to stderr, and INFORM writes the text to stdout.

```

      if (condition_ever_met) then
CSMS$MESSAGE(ABORT, 'COMPS: THIS CODE HAS NOT BEEN PARALLELIZED BY SMS')
          call comps(a,b,c,d,NX,NY)
      endif

```

**Example 12-3: Using MESSAGE to output run-time messages.**

In this example, the programmer believes the subroutine *comps* is never executed so rather than parallelizing it, MESSAGE is used. Since ABORT is specified, SMS will terminate the execution of this program after the message is output to stderr.

## 13 Debugging

SMS provides two directives that aid the debugging process. COMPARE\_VAR enables the user to compare the values of variables between simultaneous runs of an SMS program on different numbers of processes. This helps the user quickly pinpoint the source of an error that causes solutions to diverge. Sometimes, errors are due to missing or incorrectly placed exchanges. The CHECK\_HALO directive helps identify these cases by flagging halo regions that are not updated as expected.

### 13.1 Using COMPARE\_VAR To Find Parallelization Errors

Example 13-1 shows an application of the COMPARE\_VAR directive. In this code, a one process run of the code yields:

```
Running program check_var
i, y(i)
      9  34.00000
     10  38.00000
     11  42.00000
     12  46.00000
```

However, a 2 process run yields a different answer:

```
Running program check_var
i, y(i)
      9  34.00000
     10  17.00000
     11  23.00000
     12  46.00000
```

```
1      program check_var
2      parameter (IM = 20)
3  CSMS$DECLARE_DECOMP(dh, 1)
4  CSMS$DISTRIBUTE(dh, 1) BEGIN
5      real, allocatable :: x(:)
6      real, allocatable :: y(:)
7  CSMS$DISTRIBUTE END
8
9  CSMS$CREATE_DECOMP(dh, <im>, <1>)
10
11     print *, 'Running program check_var'
12     allocate(x(im))
13     allocate(y(im))
14     x = 0.0
15
16  CSMS$PARALLEL(dh, <i>) BEGIN
17     do i = 1, im
18         x(i) = i*2 - 1
19     end do
20
21     y = 0.0
22
23     do i = 2, im-1
24         y(i) = x(i-1) + x(i+1)
25     end do
```

```

26
27
28 CSMS$SERIAL BEGIN
29     print *, 'i, y(i)'
30     do i = 9, 12
31         print *, i, y(i)
32     end do
33 CSMS$SERIAL END
34
35 CSMS$PARALLEL END
36     end

```

**Example 13-1: Application of the COMPARE\_VAR directive.**

To track down the problem, the user can insert a COMPARE\_VAR directive at line 26 as follows:

```
CSMS$COMPARE_VAR(y(2:im-1), 'after y assignment')
```

Then the user can run the code as follows:

```
>> setenv SMS_SERVER_MODE serverfull
>> smsRun -np 1 check_var -np 2 check_var -cv
```

This tells SMS to simultaneously launch a 1-process and 2-process run of the program *check\_var* and compare results (“-cv”). When the code generated by the COMPARE\_VAR directive is reached, each run gathers *y* into a global equivalent, exchanges its values with the other run, and then compares them in the specified range (2:im-1). (If no range is specified, all elements of *y* are compared.) If the variable is a scalar or a non-decomposed array then it is immediately compared; no gather operation is required. In this example, since the variables differ, the program terminates with the following error message:

```

NP=1: Running program check_var
NP=2: Running program check_var
COMPARE_VAR failed :y after y assignment
Variable values for first, second run: 38.00000 17.00000
Incorrect at indices = 10

```

The difference occurs because an exchange is needed prior to the loop starting at line 23. The character string, “after y assignment”, helps pinpoint the location of the difference. The error message also identifies the name of the variable that is in error in case more than one variable is specified in the directive. It also prints out the global array location and variable values of the first point in the array that differs. Both runs immediately exit when an error is found. Notice that the print statement at line 11 appears twice. One instance is labeled with “NP=1” to indicate it came from the one-process run; the other is labeled with “NP=2” for the two-process run.

Several additional points should be made about COMPARE\_VAR. First, if the user runs the program in the standard fashion:

```
smsRun -np 2 check_var
```

then the result is the same as if the COMPARE\_VAR directive had not been included in the code at all. Second, when running with COMPARE\_VAR enabled, any output to files is turned off to avoid generation of incorrect data that could occur when the two runs simultaneously write to the same file. Third, SMS must be run with a server process for the comparison to be made properly. Fourth, if the user wishes to use configuration files while doing a comparison run, the SMS launch line looks as follows:

```
smsRun -cf config1 check_var -cf config2 check_var -cv
```

Fifth, the user should avoid putting COMPARE\_VAR directives inside a decomposed loop. For example, consider the following code fragment:

```
CSMS$PARALLEL(dh, <i>) BEGIN
  do i = 1, 2
    x(i) = x(i) + y(i)
  CSMS$COMPARE_VAR(x, "Compare 1")
  end do
```

Suppose simultaneous runs on 1 and two processes are compared. Each process in the 2 process run will call COMPARE\_VAR once. However, the single process in the 1 process run will execute the COMPARE\_VAR code twice. This will cause a hang or possibly an error message if a subsequent COMPARE\_VAR directive appears in the code. Sixth, as a reminder, bitwise-exact reductions should be enabled when using COMPARE\_VAR. Otherwise, round-off differences in summations will cause COMPARE\_VAR to indicate the presence of a spurious error. Seventh, more than one variable can be compared in one directive. For example:

```
CSMS$COMPARE_VAR(x, y, "Compare 2")
```

It is also possible to use COMPARE\_VAR to verify that the parallel and serial code solutions are identical. To do this, translate and compile the code as before:

```
>> ppp prog.f
>> f90 prog_sms.f -o par_code
```

Next, translate the code with the --CompareOnly option and compile:

```
>> ppp --CompareOnly prog.f
>> f90 prog_sms.f -o serial_code
```

The --CompareOnly flag tells SMS to ignore all directives except for COMPARE\_VAR. I/O statements are also left un-translated. Finally, the two executables can be simultaneously launched as follows:

```
>> smsRun -np 1 serial_code -np 2 par_code -cv
```

As before, when a difference is found, it will be flagged and the programs will terminate. WARNING: this use of COMPARE\_VAR may not work on all systems!

## 13.2 Debugging Adjacent Dependencies: CHECK\_HALO

The analysis of adjacent dependencies in a serial code and the process of accurately placing EXCHANGE and HALO\_COMP directives are highly prone to error. To help the user track down such errors, the CHECK\_HALO directive and associated SMS\_CHECK\_HALO environment variable can be used to check if all or part of a halo variable is up-to-date. Suppose, in Example 5-4, the user forgot to include the HALO\_COMP directives on lines 43 and 48. When the program is run, it does not produce the correct answer for *y<sub>sum</sub>*. The user can observe that the loop on lines 52-54 requires one point of the lower and upper halo regions of *b* and *c* be up-to-date. To check this assumption, the following directive can be added at line 51:

```
CSMS$CHECK_HALO(b<1,1>, c<1,1>, 'LOOP 52')
```

If the SMS\_CHECK\_HALO environment variable is set to "ON", the generated code checks if the afore-mentioned halo points are up-to-date. In this case, since the halo regions are not up-to-date, the SMS program will generate the following error message and terminate:

```
LOOP 52 Halo check failed for var : b
```

Suppose the HALO\_COMP directives are included as shown on lines 43 and 48. This time the check passes so no error messages are generated and the program continues. Suppose the user includes the HALO\_COMP directives on lines 43 and 48 and specifies the CHECK\_HALO directive as follows:

```
CSMS$CHECK_HALO(b, c, 'LOOP 52')
```

This form of the directive tells SMS to check the entire halo region. Since, for the lower and upper halo regions, only one of the halo points are up-to-date, the program will terminate with the same error message.

The directive can be added to the code on a permanent basis. When SMS\_CHECK\_HALO is "ON", CHECK\_HALO adds costly communication. However, if the SMS\_CHECK\_HALO environment variable is set to something other than "ON" then the halo checks are skipped and the CHECK\_HALO directives do not degrade performance. If, after a code change, the program generates the wrong answer, the halo checks can be turned back on to help isolate the problem.

## 14 Building a Parallel Program

### 14.1 Overview

This section describes how to use the PPP to translate Fortran code into SMS parallel source code. Output files, named automatically by PPP, will be introduced in Section 14.2. Several command line options to PPP are described in Section 14.3. In Section 14.4, a simple makefile is described which can be used to build a serial or SMS parallel code. In addition, various relevant compiler and linker options are discussed in this section. Building incorrect parallel source using PPP can result in both syntactic and semantic errors that must be corrected. Section 14.5 will discuss how to interpret these PPP generated messages. Finally, Section 14.6 will describe possible compiler errors due to namespace conflicts from PPP-generated source code.

### 14.2 PPP-Generated Output Files

Output files generated by PPP are named automatically. Include files will be named by appending “.SMS” to the original file name (e.g. *params.h* becomes *params.h.SMS*). All other source files will be named by appending “\_sms” to the body of the original filename (e.g. *main.f* becomes *main\_sms.f*). Intermediate files are also generated during the code translation process. These files, appended with the suffix “.tmp”, remain after PPP translation. When errors are detected in the code during code parallelization, PPP messages will be generated that reference these intermediate files (see Example 14-7). However, any corrections should still go into the original file from which translated code was generated by PPP.

### 14.3 Building SMS Parallel Source Code

The transformation of Fortran code into parallel SMS code requires the use of PPP. PPP translations are based on both its analysis of the original code and the SMS directives that were inserted into the code. This section describes how to use PPP to create parallel code at the command line, defines what code generation options are available, and gives some examples.

#### 14.3.1 PPP Command Line Options

All PPP code translations are managed through a command line script called **ppp**. A single file can be processed at a time and no inter-procedural analysis is done. PPP is invoked by: *ppp [options] filename*. Command line options currently available are:

<code>--checkfirst</code>	A useful optimization to avoid PPP processing of files that do not require translation. This option can be used to allow more flexible use of suffix rules (see Section 14.4). If no I/O statements or directives are found, no PPP
---------------------------	---

processing is done and the following message is output:

```

--comment      File has no directives - SKIPPING PPP PROCESSING
                Leave replaced lines in the code as Fortran
                comments. This can be useful for debugging the
                parallel code. Note: the string used to comment
                out the original code is "C-PPP".
--CompareOnly  Only translate COMPARE_VAR directives. Note
                that I/O statements will not be translated (see
                Section13.1).
--ExtendedSource Allow valid Fortran source to extend beyond 72
                characters.
--Fcommon file Name of an optional include file that is not
                part of the original source code. Typically it
                will contain data decomposition directives (see
                Example 14-4).
--Finclude file Name of an included file to be parallelized that
                is referenced in the source file being
                translated by PPP (see Example 14-2).
--Flookup file Name of file containing mapping of module names
                to source code files containing them (see
                Example 14-5).
--Fmodule file Name of module that is not part of original
                source code. Typically, it will contain data
                decomposition directives (see Example 14-5).
--Fvisible files Name(s) of file(s) to be made visible to PPP in
                order to correctly translate the current file.
                This option is only required for a series of
                interdependent include files (see Example 14-3).
--header       Indicates that the file that is about to be
                translated is a Fortran include file.
--help         Prints the command line options
--IncludePath path Include file search path. Similar to -I F77/F90
                compiler option
--r8           Indicates that an SMS program to be run on a
                machine whose normal default is 4 byte real
                numbers will be, instead, compiled so that the
                default is 8 byte real numbers.
--Verbose level Controls the output of PPP diagnostic and code
                analysis messages. Errors, Warnings and Notes
                are output based on the verbose level. (see
                Example 14-8).

```

### 14.3.2 Examples

Example 14-1 shows how to build a parallel version of an include file:

```

>> ppp --header params.h
[params.h]
    parameter(nx=50, ny=50)
CSMS$DECLARE_DECOMP(decomp, <nx, ny>)
C global variable declarations ...

```

**Example 14-1: Building any Fortran include file requires the --header option.**

Example 14-2 shows how to use the parallel version of an include file when translating an executable code file. Since the translation of *params.h* will result in an SMS parallel

version of this file (*params.h.SMS*), we use the `--Finclude` option to ensure this include file reference will be changed in the parallel version of *dynamics.f*.

```
>> ppp --Finclude=params.h --comment dynamics.f
[dynamics.f]
```

```
    program dynamics
    include 'params.h'
c   Fortran code ...
    end
```

GENERATED PARALLEL PSEUDO CODE

```
[dynamics_sms.f]
    program dynamics
C-PPP    include 'params.h'
    include 'params.h.SMS'
c   Fortran code
    end
```

**Example 14-2:** The `--Finclude` option is used to specify the Fortran include file *params.h* which is referenced in the file (*dynamics.f*) being translated. This ensures the parallel (translated) include file will be referenced in the translated output of *dynamics.f*.

Example 14-3 illustrates the use of the `--Fvisible` option. In this example, the file “*variables.h*” requires information about the data decompositions listed in “*params.h*” to correctly translate the declarations “*a*” and “*b*” enclosed within the `DISTRIBUTE` directive. In particular, the array dimensions *nx*, *ny* and *nz* must be translated to process local sizes using information provided by `DECLARE_DECOMP`. The `--Fvisible` option is used to make *params.h* “visible” to *variables.h*.

```
>> ppp --header params.h
>> ppp --Fvisible=params.h --header variables.h
>> ppp --Finclude=params.h --Finclude=variables.h main.f
```

```
[params.h]
```

```
    parameter(nx=50, ny=50)
CSMS$DECLARE_DECOMP(decomp, <nx, ny>)
```

```
C global variable declarations ...
```

```
[variables.h]
CSMS$DISTRIBUTE(decomp, nx, ny) BEGIN
    real a(nx, ny, nz)
    real b(nx, ny, nz)
CSMS$DISTRIBUTE END
```

```
[main.f]
    program main
    include 'params.h'
    include 'variables.h'
c other code
    end
```

**Example 14-3:** The `--Fvisible` option is used when inter-dependent include files must be translated.

In Example 14-1, the `DECLARE_DECOMP` was added to an include file that already existed (`params.h`). If the user prefers to insert the SMS directives into a separate “directives” file, the option `--Fcommon` is used instead of `--Finclude`. Example 14-4 illustrates the `--Fcommon` option.

```
>> ppp --header directives.inc
>> ppp --Fcommon=directives.inc dynamics.f
```

```
[directives.inc]
CSMS$DECLARE_DECOMP(decomp, 2)

[dynamics.f]
      program main
c more Fortran code
      end

                                GENERATED PARALLEL PSEUDO CODE
      program main
      include 'directives.inc.SMS'
      include 'params.h'
c more Fortran code
      end
```

**Example 14-4:** In this example `DECLARE_DECOMP`, defined in “directives.inc”, is required by “dynamics.f” when the parallel executable is built. It is not needed for a serial build.

Example 14-5 shows how to handle Fortran 90 modules in SMS. The `DECLARE_DECOMP` directive is inserted into a separate “directives” module in a way analogous to the include file in Example 14-4. In addition, the main program uses a module needed by the serial code.

```
[decomp.F]
1      module decomp
2      CSMS$DECLARE_DECOMP(dh, 2)
3      end module decomp

[m1_module.F]
1      module m1
2      CSMS$DISTRIBUTE(dh, 1, 2) BEGIN
3      real, allocatable :: u(:, :, :)
4      CSMS$DISTRIBUTE END
5      end module m1

[main.F]
1      program USE_MODULES
2      use m1
3      integer, parameter :: im = 10
4      integer, parameter :: jm = 20
5      integer, parameter :: km = 30
6
7      CSMS$CREATE_DECOMP(dh, <im,jm>, <1,1>)
8
9      allocate(u(im,jm,km))
10     CSMS$PRINT_MODE(ORDERED) BEGIN
11     print *, size(u,1), size(u,2), size(u,3)
12     CSMS$PRINT_MODE END
```

```
13      end
```

```
[my_lookup_file]
decomp          decomp.F
m1              m1_module_cpp.f
```

**Example 14-5: Code and module look-up file illustrating how SMS handles Fortran 90 modules. See text for explanation.**

To handle this code correctly, the decomposition module file must be translated first:

```
>> ppp decomp.f
```

Notice that the “--header” option is NOT used in this case in contrast to the case where an include file is translated.

Next, the module file is translated. Since the module file contains reference to the decomposition, “dh”, the decomposition module must be made visible to SMS during translation. This is done with the “-Fmodule=decomp” option. This option tells SMS to insert “use decomp” in the translated version of the module file. SMS must also know the mapping of modules to the files that contain them. In this case, module `decomp` is contained in file `decomp.F`. The lookup table specifies that mapping. The file containing this look-up table must also be specified when translating `m1_module.F`. The full command is :

```
>> ppp -Fmodule=decomp --Flookup=my_lookup_file m1_module.F
```

Finally, the command line for translating `main.F` is:

```
>> ppp -Fmodule=decomp --Flookup=my_lookup_file main.F
```

## 14.4 Building SMS Programs

A simple makefile is presented to aid the user in building an SMS program. The environment variable “SMS” must be set to the location where the SMS software has been installed. This can either be done explicitly in the Makefile (line 5) or via `setenv` prior to running `make` (e.g. `setenv SMS pathname`).

```
1  # standard make file used to build serial or SMS parallel
2  # executables
3  #
4  .SUFFIXES: .f .o
5  SMS = /usr/local/sms
6
7  # system-specific compilation flags (for a Compaq Alpha EV67)
8  FC = mpif90
9  FFLAGS = -O4 -arch host -tune host -fixed -I$(SMS)/include
10
11 # SMS link libraries
12 LIBS = -L$(SMS)/lib -lsms -lmpi
13
14 # PPP specific options set here
15 PPP = $(SMS)/bin/ppp
```

```

16  PPP_FLAGS = --Finclude=params.h --Finclude=variables.h --comment \
17  --checkfirst
18  PPP_HEADER_FLAGS = --header --comment
19
20  # include files
21  INCLUDES = params.h variables.h globals.h
22  PINCLUDES = ${INCLUDES:.h=.h.SMS}
23
24  # object files
25  OBJS = file1.o file2.o file3.o
26
27  TARGET = par_prog
28
29  # executable
30  $(TARGET): $(PINCLUDES) $(OBJS)
31             $(FC) -o $(TARGET) $(OBJS) $(FFLAGS) $(LIBS)
32
33  # suffix rules
34  .f.o: $(PINCLUDES)
35         $(PPP) $(PPP_FLAGS) $*.f
36         $(FC) -c $(FFLAGS) $*_sms.f
37         /bin/mv -f $*_sms.o $*.o
38
39  # include file translations
40  params.h.SMS:      params.h
41                 $(PPP) $(PPP_HEADER_FLAGS) params.h
42
43  variables.h.SMS:  variables.h params.h
44                 $(PPP) $(PPP_HEADER_FLAGS) --Fvisible=params.h variables.h
45
46  clean:
47         /bin/rm -f $(TARGET) *_sms.f *.SMS *.o *.tmp

```

**Example 14-6: A makefile for an SMS program.**

## 14.4.1 Makefile Compiler and Linker Options

The Fortran compiler flags (FLAGS on line 9) are set for a Compaq Alpha EV67. Other systems will require different options. A makefile provided in the SMS distribution (\$SMS/lib/makefile.header) gives recommended compilation flags (found in variable STD\_OPT\_FLAGS) that should be used when modifying FFLAGS for a different target machine.

## 14.4.2 Include File Handling

Include files are listed in the makefile variable INCLUDES. Parallel include files (line 22) translated using SMS are built using the explicit targets *params.h.SMS* and *variables.h.SMS* (lines 40-44). Notice the PPP command to build *variables.h.SMS* (line 44) contains the --Fvisible option in addition to the standard PPP flags defined by: PPP\_HEADER\_FLAGS at line 18. Since *variables.h.SMS* requires information from *params.h* for proper translation, this option is required (see Example 14-3).

PPP\_FLAGS (lines 16-17) lists the include files that are translated by PPP via the --Finclude option. This option is required to ensure any references to these files in Fortran source will be modified to their parallel filename (see Example 14-2).

### 14.4.3 Building the Executable

To build the SMS parallel executable “par\_prog” using this makefile, simply run make:

```
>> make
```

Translated source code is written to the file “*file1\_sms.f*”.

### 14.5 PPP Error Reporting

Two types of errors are reported by PPP: parsing errors and semantic errors. Parsing errors must be corrected before further translations of the input file are permitted. Semantic errors are reported as errors, warnings or notes. These messages can be controlled through the --verbose option of PPP (see Section 14.5.2).

#### 14.5.1 Parsing Errors

Parsing errors occur when PPP cannot resolve the Fortran source code to the grammar defined by the SMS directives (refer to the SMS Reference Manual), the Fortran 77 language, and currently supported Fortran 90 syntax. Further details about language extensions supported by SMS can be found at:

[http://www-ad.fsl.noaa.gov/ac/SMS\\_Supported\\_Fortran\\_Features.html](http://www-ad.fsl.noaa.gov/ac/SMS_Supported_Fortran_Features.html)

The parser currently supports statements or SMS directives that are up to 500 characters in length. Multiple statement lines are collapsed and white space is removed before statements are parsed. Statements longer than 500 characters will not be parsed correctly.

The form of a parsing error message is:

```
<filename> <line> <column> <error type> <message>
```

```
filename    - name of file being parsed
line        - line number
column      - column number in which error occurred
error type  - types are:
               ERROR, WARNING, NOTE
message     - diagnostic message
```

An example of a PPP-generated parsing error is shown in Example 14-7.

```
1 CSMS$DECLARE_DECOMP(spec_dh,<jtrun>)
2 CSMS$DISTRIBUTE(spec_dh, jtrun) BEGIN
3     real*8 cc(jtrun), bb(jtrun)
4 CSMS$DISTRIBUTE END
5
6 CSMS$PARALLEL(spec_dh, m) BEGIN
7     do 3 m=2, jtrun, 2
8         cc(m) = cc(m) + bb(m)
9     3 continue
10
```

```

11 C CSMS$PARALLEL END is missing
12
13         end

```

**Example 14-7: Code that generates a PPP parsing error.**

PPP generates the following error message:

```

"Loops_sms.f.tmp" 13 501 ERROR:  Syntax error
"Loops_sms.f.tmp" 13 501 NOTE  Parsing resumed here

```

This message indicates the parser failed in the file *Loops\_sms.f.tmp* at line 13 column 501. A parsing error occurring at column 501 indicates no resolution of the statement to the grammar by the end of the line. In the example, the parser expects a PARALLEL END directive before the end of the file. Naturally, the error should be corrected in the original file (*Loops.f*) rather than the PPP-generated file.

## 14.5.2 PPP Diagnostic Messages

Three levels of diagnostic messages are reported by PPP. A PPP ERROR is reported when a section of code targeted for translation contains a syntax or semantic error. A PPP WARNING is reported when PPP suspects that it may generate an incorrect translation. A PPP NOTE identifies a place where a particular type of transformation occurred or SMS limitation was detected. By default, all PPP ERROR messages will be output. Control of diagnostic messages is handled through the PPP command line option: “*--verbose = <value>*”. Three verbose options are supported:

value	message domain
1	PPP ERRORS only (DEFAULT)
2	PPP ERRORS and WARNINGS only
3	PPP ERRORS, WARNINGS and NOTES

While the error messages should always be addressed, warning messages may also be useful for detecting potential problems. For example, the code segment in Example 14-8 below causes PPP to generate the following important warning message:

```

./IO.f.tmp" 11 13 WARNING: This variable, decomposed by CSMS$DISTRIBUTE,
is being used outside of a parallel region.

```

This warning message indicates a problem on line 11, column 13 of the PPP-generated file *IO.f.tmp* (which is not shown). The variable *cc* was defined to be a distributed array (using DISTRIBUTE) but is being referenced outside a parallel region (PARALLEL). Further explanation on the use of these directives can be found in Section 2.3.

```

>> ppp --verbose=2 IO.p
1 CSMS$DISTRIBUTE(dh, m, n) BEGIN
2     real cc(m,n)
3 CSMS$DISTRIBUTE END

```

```
4      do i = 1, m
5          do j = 1, n
6              cc(i,j) = 0.0
7          enddo
8      enddo
9
10 c   more code ...
```

**Example 14-8:** Code that generates a **WARNING** because the decomposed variable “cc” is being used outside of a parallel region.

## 14.6 Compilation Errors

During the translation process PPP generates new variables for some translations. PPP variables are either automatically generated or defined explicitly by PPP. Explicitly defined names will always contain a double underscore in their name (e.g. `ppp__status`). To avoid compiler errors due to name space conflicts, avoid using variable names with double underscores in them. For example, the serial code cannot contain a declaration of a variable named `ppp__status` because PPP translation explicitly declares a variable named `ppp__status` for its own use. A compilation error would result because two variables would be declared with the same name.

## 15 Running an SMS Program

### 15.1 Introduction

Once a program has been translated into SMS parallel code (Section 14.3) and linked to the appropriate libraries (see Section 14.4), it can be run on one or more processes using the SMS program launcher *smsRun*. The standard syntax for *smsRun* is:

```
>> smsRun -np numprocs execname
```

By default, SMS uses an additional server process to perform I/O operations, and provide overall management and control services for the other processes (see Figure 11-5). For example, to run the executable *test* with two processes and one server process, the user would enter:

```
>> smsRun -np 2 test
```

It is possible to take advantage of the idle compute cycles available on the server process by setting SMS environment variable *SMS\_SERVER\_MODE* to *serverless*. This will permit computational and management functions to co-exist in a single process. This option is beneficial when only a small number of processes are available. However, as the numbers of processes grow, the cost of performing both server functions and computations will limit the performance of the other dependent processes.

Figure 11-5 assumes a single process is run on each processor. However, SMS permits the user to request more processes (using *smsRun*) than available processors. For example if *my\_program* was run with 20 processes

```
>> smsRun -np 20 my_program
```

on a system with only 16 processors, five processors would contain two processes, one would contain the server process, and the rest would each contain a single process designated to run the program. This is a bad idea because performance will suffer whenever multiple processes are scheduled on a single processor on most machines.

### 15.2 Optional Command Line Arguments

Several optional arguments to *smsRun* are permitted. One optional argument to control the number of I/O write-cache processes to be dedicated to the program's execution can be expressed by:

```
>> smsRun -np 20 my_program -wc 2
```

This command runs *my\_program* on 20 processes with one additional server process. Also, two additional, processes are used to cache data written to disk. The use of write-cache processes to improve output performance is discussed in Section 11.3.5.

The `-cf` option is used to gain more control over process layout as described in Section 3.3. In the example below, program `my_program` is run using the process layout specified in configuration file `my_config`:

```
>> smsRun -cf my_config my_program
```

The `-cv` option is used to assist debugging by turning on `COMPARE_VAR` directives as described in Section 13.1. In the first example below, program `my_program` is run simultaneously on 1 and 9 processes. Arrays specified in `COMPARE_VAR` directives are compared on-the-fly. If a difference is found, execution is halted and an error message is printed. In the second example, different executables are loaded (useful for comparing a serial run to a parallel run or for testing static memory codes that use minimum memory) and the second program uses a configuration file.

```
>> smsRun -np 1 my_program -np 9 my_program -cv
>> smsRun -np 1 my_program1 -cf my_config9 my_program9 -cv
```

Another option, `-sms-`, allows the user to specify machine-specific arguments to the underlying communication layer (e.g. MPI, SHMEM) directly. All arguments that follow this option will be ignored by SMS and passed directly to the communications software. The following command illustrates a way to pass the run-time option `-mpi_special` to the underlying program launcher (which is `mpirun` when SMS is built using MPI):

```
>> smsRun -np 3 my_program -sms- -mpi_special
```

Information about other machine-specific options for `smsRun` is available at the following SMS web site:

[http://www-ad.fsl.noaa.gov/ac/SMS\\_Run\\_Options.html](http://www-ad.fsl.noaa.gov/ac/SMS_Run_Options.html)

## 15.3 Run-time Environment Variables

Several environment variables can also be set to control the run-time behavior of SMS. The following environment variables are available:

<code>SMS_BITWISE</code>	- Set to "EXACT" to use bit-wise exact reductions (see Section 7.2)
<code>SMS_CHECK_HALO</code>	- Set to "ON" to execute checks of halo regions specified by <code>CHECK_HALO</code> directives (see Section 13.2).
<code>SMS_CLOSE_MODE</code>	- Use to improve output performance (see Section 11.3.4).
<code>SMS_IO_FORMAT</code>	- Use to specify file format for files that are read or written by SMS.
<code>SMS_IOC_SIZE</code>	- Use to improve output performance (see Section 11.3.5)
<code>SMS_RAN_RSTYLE</code>	- Use to improve input performance (see Section 11.3.6).
<code>SMS_RBC</code>	- Use to improve input performance (see Section 11.3.6).
<code>SMS_RBS</code>	- Use to improve input performance (see Section 11.3.6).

- SMS\_READ\_FORMAT - Use to specify file format for files that are read by SMS programs (see Section 11.1).
- SMS\_SERVER\_MODE - Set to "SERVERLESS" to avoid using a server process (see Section 11.3.2).
- SMS\_WBS - Use to improve output performance (see Section 11.3.5).
- SMS\_WRITE\_FORMAT - Use to specify file format for files that are written by SMS programs (see Section 11.1).
- SMS\_XFERMODE - Use to choose optimal communication patterns for TRANSFER. Options are: "LOGN" and "ORIGINAL".

## 15.4 Run-time Error Messages

When an error occurs in an SMS program, execution will usually terminate and SMS will generate an informational message describing the source file name, line number, and a brief summary of the problem. A complete set of SMS run-time error messages is available at the following SMS web site:

[http://www-ad.fsl.noaa.gov/ac/SMS\\_Messages.html](http://www-ad.fsl.noaa.gov/ac/SMS_Messages.html)

When the code in Example 3-3 is run with 2 processes, the following error message is generated (as seen in Section 3.4):

```

Process: 1 Error at: ./decomp_ex4_sms.f:10.1
Process: 1 Error status= -2202 : USER DECLARED STATIC ARRAY IS TOO
SMALL.
Process: 1 Aborting...
```

The first line of the error message indicates the file name and location within the file where the problem occurred. PPP-generated code frequently uses sub-numbering to handle multiply generated calls to SMS routines that stem from the same line of original code. In this example, a run-time error was detected by SMS at line 17 in code generated by the directive `CREATE_DECOMP` that can be found in temporary file: `decomp_ex4_sms.f.tmp` (not shown).

The second line gives the SMS error message. The error messages reflect the incorrect sizing of the decomposition *decomp*, declared by `DECLARE_DECOMP` and initialized by `CREATE_DECOMP`.

Once the problem is understood corrections to the code can be made. These corrections should go into the original file (in this case `decomp_ex4.f`) not in the temporary file where the problem was detected and probably diagnosed. Once changes are made, PPP can be executed to re-translate the input file from which a fresh executable can be built and tested.

---

**A**

adjacent dependence · 23  
ASYNC · 16, 97

---

**B**

bandwidth · 5  
Bit-wise Exact · 71  
broadcast · 93

---

**C**

cache · 30  
CHECK\_HALO · 112  
COMPARE\_VAR · 109  
Coupling · 83  
CREATE\_DECOMP · 12

---

**D**

DECLARE\_DECOMP · 12  
decomposition · 10  
dependence analysis · 5  
DISTRIBUTE · 12  
Distributed Memory · 5  
dynamic memory allocation · 10

---

**E**

environment variables · 123  
Exchange  
    Aggregating Exchanges · 60  
    Limited-Thickness Exchange · 59  
EXCHANGE · 52

---

**F**

FDA · 5  
finite difference approximation · 5  
FLUSH\_OUTPUT · 102, 103

---

**G**

gather · 68, 74, 91  
global array · 10  
global dependence · 17  
global indices · 10  
global sizes · 10, 53  
GLOBAL\_INDEX · 49

---

**H**

halo regions · 13, 24, 32  
HALO\_COMP · 62

---

**I**

I/O Performance Tuning · 99  
interior · 54

---

**L**

latency · 5  
*lbound* · 48  
local array · 10  
local indices · 10, 12, 14, 46  
local sizes · 10  
Lower Bounds · 44

---

**M**

module · 114, 116, 117  
MPI · 8

---

**N**

namelist · 96  
native · 95  
Nesting · 83  
numbers · 6, 18, 37, 38, 46, 71, 73

---

**O**

ORDERED · 17, 97

---

**P**

PARALLEL · 12

Periodic Boundary Conditions · 80

PPP · 8

PRINT\_MODE · 16, 18

process · 6

process configuration file · 32

---

**R**

recurrence relation · 27

REDUCE · 13, 18, 71

redundant computations · 62

---

**S**

scatter · 91

scope · 13

SERIAL · 46, 76

serverless · 99, 102, 122

SET\_NEST\_LEVEL · 89

---

SET\_TRANSFER\_INTERPOLATE · 85

shared memory · 5

SMS\_READ\_FORMAT · 95

SMS\_WRITE\_FORMAT · 95

smsRun · 8, 122

spectral transform method · 5

SPMD · 5

Standard Reductions · 71

stencil · 23

Stencil-Oriented Form · 85

---

**T**

TO\_LOCAL · 46

TRANSFER · 13, 29, 68

    Aggregation · 69

    Interpolating between grids · 88

Transfer-Interpolation · 83

---

**V**

vector · 30

---

**W**

write-cache · 104

---