

IBM Parallel Environment for Linux



Introduction

Version 4 Release 2

IBM Parallel Environment for Linux



Introduction

Version 4 Release 2

Note

Before using this information and the product it supports, read the information in “Notices” on page 63.

First Edition (April 2006)

This edition applies to version 4, release 2, modification 0 of IBM Parallel Environment for Linux (product number 5724-N05) and to all subsequent releases and modifications until otherwise indicated in new editions. Significant changes or additions to the text and illustrations are indicated by a vertical line (|) to the left of the change.

IBM welcomes your comments. A form for readers' comments may be provided at the back of this publication, or you may address your comments to the following address:

International Business Machines Corporation
Department 55JA, Mail Station P384
2455 South Road
Poughkeepsie, NY 12601-5400
United States of America

FAX (United States & Canada): 1+845+432-9405

FAX (Other Countries): Your International Access Code +1+845+432-9405

IBMLink™ (United States customers only): IBMUSM10(MHVRCFS)

Internet e-mail: mhvrdfs@us.ibm.com

If you would like a reply, be sure to include your name, address, telephone number, or FAX number.

Make sure to include the following in your comment or note:

Title and order number of this book

Page number or topic related to your comment

When you send information to IBM, you grant IBM a nonexclusive right to use or distribute the information in any way it believes appropriate without incurring any obligation to you.

© Copyright International Business Machines Corporation 2006. All rights reserved.

US Government Users Restricted Rights – Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

Contents

Figures	v
About this book	vii
Who should read this book	vii
How this book is organized	vii
Overview of contents	vii
Conventions and terminology used in this book	viii
Abbreviated names	viii
Prerequisite and related information	ix
Using LookAt to look up message explanations	ix
How to send your comments	x
National language support (NLS)	x
Functional restrictions for PE 4.2.0	x
Chapter 1. Understanding the environment	1
What is IBM Parallel Environment for Linux?	1
What is the Parallel Operating Environment?	1
Before you start	2
Running POE	3
Chapter 2. Message passing	21
The message passing model	21
Data decomposition	22
Functional decomposition	29
Duplication versus redundancy	31
Protocols supported	32
Thread debugging implications	34
Chapter 3. Diagnosing and correcting common problems	35
Messages	35
Message catalog errors	35
Finding PE messages	35
Logging POE errors to a file	36
Message format	36
Diagnosing problems using IVP	36
Cannot compile a parallel program	36
Cannot start a parallel job	37
Cannot execute a parallel program	38
The program runs but...	41
When a core dump is created	41
No output at all	41
The program hangs	42
Why did the program hang?	44
Other reasons for the program to hang	45
Bad output	46
Debugging and threads	46
Tuning the performance of a parallel application	47
Tuning the performance of threaded programs	47
Profile it	48
Chapter 4. Creating a safe program	49
What is a safe program?	49
Safety and threaded programs	49

Using threaded programs with non-thread-safe libraries	50
Message ordering	50
Program progress when two processes initiate two matching sends and receives	51
Communication fairness	51
Resource limitations	51
Appendix A. A sample program to illustrate messages	53
Figuring out what all of this means	57
Appendix B. Parallel Environment internals	59
What happens when I compile my applications?	59
How do my applications start?	59
How does POE talk to the nodes?	59
How are signals handled?	60
What happens when my application ends?	60
Appendix C. Accessibility	61
Using assistive technologies	61
Notices	63
Trademarks	65
Glossary	67
Index	73

Figures

1. Output from compiler scripts	8
---	---

About this book

This book provides suggestions and guidance for using the IBM Parallel Environment for Linux to develop and run Fortran, C, and C++ parallel applications. To make this book a little easier to read, the name *Parallel Environment for Linux* has been abbreviated to *PE* throughout.

Note that PE for Linux is based on its predecessor, PE for AIX, with which you may be familiar.

In this book, you will find information on basic parallel programming concepts and the Message Passing Interface (MPI) standard. You will also find information about the application development tools that are provided by PE such as the Parallel Operating Environment.

Note: Although this book includes information about using PE with the User Space protocol, User Space is not supported on PE for Linux at this time.

Who should read this book

This book is intended for application developers who are interested in creating and running parallel programs. To make the best use of this book, you should be familiar with the following:

- The Linux operating system
- One or more of the supported programming languages (Fortran, C, or C++)
- Basic parallel programming concepts.

This book is not intended to provide comprehensive coverage of the topics, nor is it intended to tell you everything there is to know about IBM Parallel Environment (PE). If you are new to either message passing parallel programming or to PE, you should find this book useful. For the latest information, always use the documents at:

<http://publib.boulder.ibm.com/infocenter/clresctr/index.jsp>

The purpose of this book is to get you started creating parallel programs with PE. Once you have mastered these initial concepts, you will need to know more about how PE works. For information on the Parallel Operating Environment (POE), see *IBM Parallel Environment: Operation and Use*.

How this book is organized

Overview of contents

This book contains the following information:

- **Chapter 1, “Understanding the environment,” on page 1** familiarizes you with the Parallel Operating Environment (POE).
- **Chapter 2, “Message passing,” on page 21** covers parallelization techniques and discusses their advantages and disadvantages. It discusses how you take a working serial program and create a parallel program that gives the same result.
- **Chapter 3, “Diagnosing and correcting common problems,” on page 35** outlines the possible causes for a parallel application to fail to execute correctly, and provides some tips on how to identify and correct problems.

- **Chapter 4, “Creating a safe program,” on page 49** provides you with some general guidelines for creating *safe* parallel MPI programs.
- **Appendix A, “A sample program to illustrate messages,” on page 53** provides a sample program, run with the maximum level of error messages. It points out the various types of messages you can expect, and tells you what they mean.
- **Appendix B, “Parallel Environment internals,” on page 59** provides some additional information about how the PE works with respect to your application.

Conventions and terminology used in this book

Note that in this document, LoadLeveler is also referred to as *Tivoli® Workload Scheduler LoadLeveler* and *TWS LoadLeveler*.

This book uses the following typographic conventions:

Convention	Usage
bold	Bold words or characters represent system elements that you must use literally, such as: command names, file names, flag names, path names, PE component names (poe , for example), and subroutines.
constant width	Examples and information that the system displays appear in constant-width typeface.
<i>italic</i>	<i>Italicized</i> words or characters represent variable values that you must supply. <i>Italics</i> are also used for book titles, for the first use of a glossary term, and for general emphasis in text.
[item]	Used to indicate optional items.
<Key>	Used to indicate keys you press.
\	The continuation character is used in coding examples in this book for formatting purposes.

In addition to the highlighting conventions, this manual uses the following conventions when describing how to perform tasks.

User actions appear in uppercase boldface type. For example, if the action is to enter the **tool** command, this manual presents the instruction as:

```
ENTER
    tool
```

Abbreviated names

Some of the abbreviated names used in this book follow.

CSM Clusters Systems Management
dsh distributed shell
GUI graphical user interface
IP Internet Protocol
LAPI Low-level Application Programming Interface
MPI Message Passing Interface

PE	IBM® Parallel Environment for Linux
PE MPI	IBM's implementation of the MPI standard for PE
PE MPI-IO	IBM's implementation of MPI I/O for PE
POE	parallel operating environment
RSCT	Reliable Scalable Cluster Technology
rsh	remote shell
STDERR	standard error
STDIN	standard input
STDOUT	standard output
xSeries	IBM @server xSeries®

Prerequisite and related information

The Parallel Environment library consists of:

- *IBM Parallel Environment: Introduction*, SA23-2218
- *IBM Parallel Environment: Installation*, SC23-5208
- *IBM Parallel Environment: Operation and Use*, SA23-2217
- *IBM Parallel Environment: MPI Programming Guide*, SA23-2219
- *IBM Parallel Environment: Messages*, SA38-0648
- *IBM Parallel Environment: MPI Subroutine Reference*, SA23-2220

To access the most recent Parallel Environment documentation in PDF and HTML format, refer to the IBM @server Cluster Information Center on the Web at:

<http://publib.boulder.ibm.com/infocenter/clresctr/index.jsp>

Both the current Parallel Environment books and earlier versions of the library are also available in PDF format from the IBM Publications Center Web site located at:

<http://www.ibm.com/shop/publications/order/>

It is easiest to locate a book in the IBM Publications Center by supplying the book's publication number. The publication number for each of the Parallel Environment books is listed after the book title in the preceding list.

Using LookAt to look up message explanations

LookAt is an online facility that lets you look up explanations for most of the IBM messages you encounter, as well as for some system abends and codes. You can use LookAt from the following locations to find IBM message explanations for Clusters for Linux:

- The Internet. You can access IBM message explanations directly from the LookAt Web site:
<http://www.ibm.com/eserver/zseries/zos/bkserv/lookat/>

- Your wireless handheld device. You can use the LookAt Mobile Edition with a handheld device that has wireless access and an Internet browser (for example, Internet Explorer for Pocket PCs, Blazer, or Eudora for Palm OS, or Opera for Linux® handheld devices). Link to the LookAt Mobile Edition from the LookAt Web site.

How to send your comments

Your feedback is important in helping to provide the most accurate and high-quality information. If you have comments about this book or other PE documentation:

- Send your comments by e-mail to: mhvrcfs@us.ibm.com
Be sure to include the name of the book, the part number of the book, the version of PE, and, if applicable, the specific location of the text you are commenting on (for example, a page number or table number).
- Fill out one of the forms at the back of this book and return it by mail, by fax, or by giving it to an IBM representative.

National language support (NLS)

For national language support (NLS), all PE components and tools display messages that are located in externalized message catalogs. English versions of the message catalogs are shipped with the PE licensed program, but your site may be using its own translated message catalogs. The PE components use the Linux environment variable **NLSPATH** to find the appropriate message catalog. **NLSPATH** specifies a list of directories to search for message catalogs. The directories are searched, in the order listed, to locate the message catalog. In resolving the path to the message catalog, **NLSPATH** is affected by the values of the environment variables **LC_MESSAGES** and **LANG**. If you get an error saying that a message catalog is not found and you want the default message catalog:

ENTER

```
export NLSPATH=/usr/share/locale/%L/%N
```

```
export LANG=en_US
```

The PE message catalogs are in English, and are located in the following directories:

```
/usr/share/locale/en_US/pempl.cat
```

```
/usr/share/locale/en_US/pepoe.cat
```

```
/usr/share/locale/en_US/liblapi.cat
```

If your site is using its own translations of the message catalogs, consult your system administrator for the appropriate value of **NLSPATH** or **LANG**.

Functional restrictions for PE 4.2.0

Although it is documented in this library, the use of the User Space protocol is restricted for Parallel Environment 4.2.0.

The following functions, although currently available with Parallel Environment for AIX, are not supported by Parallel Environment for Linux 4.2.0.

- Checkpoint/restart
- Lightweight corefiles
- pdbx parallel debugger

- Xprofiler profiling tool

Chapter 1. Understanding the environment

This section will help you understand the new environment, IBM Parallel Environment for Linux (PE). It covers:

- PE
- The Parallel Operating Environment (POE)
- Starting the POE
- Running simple commands
- Experimenting with parameters and environment variables
- Using a *host list* file versus a job management system (LoadLeveler) for requesting processor nodes
- Compiling and running a simple parallel application
- Some simple environment setup and debugging tips.

This book contains examples and illustrates various commands and programs as well as the output you get as a result of running them. When looking at these examples, keep in mind that the output you see on the system may not exactly match what is printed in the book. The included examples here give you a basic idea of what happens.

What is IBM Parallel Environment for Linux?

IBM Parallel Environment for Linux (PE) software lets you develop, analyze, and execute parallel applications written in Fortran, C, and C++. PE conforms to existing standards like UNIX[®] and MPI.

PE consists of the following:

- The Parallel Operating Environment (POE), for submitting and managing jobs.
- A message passing library (MPI), for communication among the tasks that make up a parallel program.
- Parallel utilities for easing file manipulation.

What is the Parallel Operating Environment?

The Parallel Operating Environment (POE) allows you to develop and execute the parallel applications across multiple operating system images, called **nodes**. When using POE, there is a single node (possibly a workstation) that is called the *home node* that manages interactions with users.

POE transparently manages the allocation of remote nodes where the parallel application actually runs. It also handles the various requests and communication between the home node and the remote nodes via the underlying network.

This approach eases the transition from serial to parallel programming by hiding the differences, and allowing you to continue using standard Linux tools and techniques. You have to tell POE what remote nodes to use, but once you have, POE does the rest.

The *processor node* is a physical entity or operating system image that is defined to the network. It can be a standalone machine, or a processor node within a cluster. From POE's point of view, a node is a single copy of the operating system.

If you are using a Symmetric Multiprocessor (SMP) system, it is important to know that, although an SMP node has more than one processing unit, it is still considered, and referred to as, a *processor node*.

Before you start

Before starting, check that you have addressed the items covered in this section.

Installation

Whoever installed POE should have verified that it was installed successfully by running the *Installation Verification Program (IVP)*. The *IBM Parallel Environment: Installation Guide* discusses the IVP.

The IVP tests to see if POE can do the following:

- Establish a remote execution environment
- Compile and execute the program
- Initialize the IP message passing environment
- Check that the MPI library is operable.

Access

Before running the job, you must first have access to computer resources in the system. Here are some things to consider:

- You must have the *same* user ID and group ID on the home node and each remote node on which you will be running the parallel application.
- POE will not allow you to run the application as root.

If you are using LoadLeveler to submit POE jobs, which includes all user space applications, then LoadLeveler is responsible for the security authentication. The security function in POE is not invoked when POE is run under LoadLeveler.

Security methods: PE uses an enhanced set of security methods, based on Cluster Security Services in RSCT (Reliable Scalable Cluster Technology). RSCT is a set of software components that provide a comprehensive clustering environment. RSCT is the infrastructure used by a variety of products to provide clusters with improved system availability, scalability, and ease of use.

Under Linux, PE supports a limited set of user authorization mechanisms. POE uses a configuration option for the system administrator to define the security mechanism which, on Linux, is limited to the compatibility method. When you set up a node, each user ID must be authorized to access that node or remote link from the initiating home node. To specify this user ID authorization, use the ***/etc/hosts.equiv*** or the ***.rhosts*** file, or both, as explained in “User authorization.”

User authorization

You must have remote execution authority on all the nodes in the system that you will use for parallel execution. The system administrator should:

- Authorize both the home node machine and the user name (or machine names) in the ***/etc/hosts.equiv*** file on each remote node, or
- Set up the ***.rhosts*** file in the home directory of the user ID for each node that you want to use. The contents of each ***.rhosts*** file can be either the explicit IP address of the home node, or the home node name. For more information about ***.rhosts*** files, see the *IBM Parallel Environment: Installation*.

/etc/hosts.equiv is checked first, and, if the home node and user/machine name do not appear there, it then looks to ***.rhosts***.

You can verify that you have remote execution authority by running a remote shell from the workstation where you intend to submit parallel jobs. For example, to test whether you have remote execution authority on node **202r1n10**, try the following command:

```
$ rsh 202r1n10 hostname
```

The response should be the remote host name. If it is not the remote host name, or the command cannot run, see the system administrator. Issue this command for every remote host on which you plan to have POE execute the job.

Refer to *IBM Parallel Environment: Installation* for more information on enabling **rsh**.

Host list file

One way to tell POE where to run the program is by using a *host list* file. The host list file is generally in the current working directory, but you can move it anywhere you like by specifying certain parameters. This file can be given any name, but the default name is *host.list*. Many people use *host.list* as the name to avoid having to specify another parameter. This file contains one of two different kinds of information; node names or pool numbers (a pool can also be designated by a string).

Node names refer to the hosts on which parallel jobs may be run. They may be specified as Domain Names (as long as those Domain Names can be resolved from the workstation where you submit the job) or as Internet addresses. Each host goes on a separate line in the host list file.

Here is an example of a host list file that specifies the node names on which four tasks will run:

```
202r1n09.hpssl.kgn.ibm.com
202r1n10.hpssl.kgn.ibm.com
202r1n11.hpssl.kgn.ibm.com
202r1n12.hpssl.kgn.ibm.com
```

Running POE

After you have checked all the items in “Before you start” on page 2, you are ready to run the POE. You can view POE as a way to run commands and programs on multiple nodes from a single point. Remember that these commands and programs are really running on the remote nodes. If you ask POE to perform some operation on a remote node, everything necessary to perform that operation must be available on the remote node.

There are two ways to influence the way the parallel program is executed; with environment variables or command-line option flags. You can set environment variables at the beginning of the session to influence each program that you execute. You also get the same effect by specifying the related command-line flag when you invoke POE, but its influence lasts only for that particular program execution. This book shows you how to use the command-line option flags to influence the way the program executes. “Running POE with environment variables” on page 5 gives you some high-level information, but you may also want to refer to *IBM Parallel Environment: Operation and Use* to learn more about using environment variables.

The following sections show you how to run a parallel job by requesting that POE use nodes in a host list file and tell you how to use a host list file to request nodes from LoadLeveler.

Some examples of running POE

The **poe** command enables you to load and execute programs on remote nodes. The syntax is:

```
poe [program] [options]
```

When you invoke **poe**, it allocates processor nodes for each task and initializes the local environment. It then loads the program and reproduces the local shell environment on each processor node. POE also passes the user program arguments to each remote node.

The simplest thing to do with POE is to run a command. When you try these examples on the system, use a host list file that contains the node names (as opposed to a pool number). These examples assume at least a four-node parallel environment. If you have more than four nodes, feel free to use more. If you have fewer than four nodes, you may duplicate lines. This example assumes that the file is called *host.list*, and is in the directory from which you are submitting the parallel job. If either of these conditions are not true, POE will not find the host list file unless you use the **-hostfile** option.

The **-procs 4** option tells POE to run this command on four nodes. It will use the first four in the host list file.

```
$ poe hostname -procs 4  
  
202r1n10.hpssl.kgn.ibm.com  
202r1n11.hpssl.kgn.ibm.com  
202r1n09.hpssl.kgn.ibm.com  
202r1n12.hpssl.kgn.ibm.com
```

What you see is the output from the **hostname** command run on each of the remote nodes. POE has taken care of submitting the command to each node, collecting the standard output and standard error from each remote node, and sending it back to the workstation. One thing that you do not see is an indication of which task is responsible for each line of output. In a simple example like this, it is not that important. If, however, you had many lines of output from each node, you would want to know which task was responsible for each line of output. To do that, you use the **-labelio** option:

```
$ poe hostname -procs 4 -labelio yes  
  
1:202r1n10.hpssl.kgn.ibm.com  
2:202r1n11.hpssl.kgn.ibm.com  
0:202r1n09.hpssl.kgn.ibm.com  
3:202r1n12.hpssl.kgn.ibm.com
```

Notice how each line starts with a number and a colon and that the numbering started at 0 (zero). The number is the task ID that the line of output came from (it is also the line number in the host list file that identifies the host which generated this output). Use this parameter to identify lines from a command that generates more output. Try this command:

```
$ poe perms -procs 2 -labelio yes
```

You should see something similar to this:

```
0:ppe_x86_base_32bit_sles900-4.2.1.0-0607a  
0:ppe_x86_64bit_sles900-4.2.1.0-0607a  
1:ppe_x86_base_32bit_sles900-4.2.1.0-0607a  
1:ppe_x86_64bit_sles900-4.2.1.0-0607a
```

```
0:lapi_x86_base_32bit_sles900-2.4.1.0-0607a
0:lapi_x86_64bit_sles900-2.4.1.0-0607a
1:lapi_x86_base_32bit_sles900-2.4.1.0-0607a
1:lapi_x86_64bit_sles900-2.4.1.0-0607a
```

The **perpms** command is displaying the installed RPMs on each node. Notice how the output from each of the remote nodes is intermingled. This is because as soon as a buffer is full on the remote node, POE sends it back to the workstation for display (in case you had any doubts that these commands were really being executed in parallel). The result is the jumbled mess that can be difficult to interpret. Fortunately, POE can clear things up with the **-stdoutmode** parameter.

Try this command:

```
$ poe perpms -procs 2 -labelio yes -stdoutmode ordered
```

You should see something similar to this:

```
0:ppe_x86_base_32bit_sles900-4.2.1.0-0607a
0:ppe_x86_64bit_sles900-4.2.1.0-0607a
0:lapi_x86_base_32bit_sles900-2.4.1.0-0607a
0:lapi_x86_64bit_sles900-2.4.1.0-0607a
1:ppe_x86_base_32bit_sles900-4.2.1.0-0607a
1:ppe_x86_64bit_sles900-4.2.1.0-0607a
1:lapi_x86_base_32bit_sles900-2.4.1.0-0607a
1:lapi_x86_64bit_sles900-2.4.1.0-0607a
```

POE holds all the output until the jobs either finish or POE itself runs out of space. If the jobs finish, POE displays the output from each remote node together. If POE runs out of space, it prints everything, and then starts a new page of output. You get less of a sense of the parallel nature of the program, but it is easier to understand.

Running POE with environment variables

If you are getting tired of typing the same command line options over and over again, you can set them as environment variables so that you do not have to put them on the command line. The environment variable names are the same as the command line option names (without the leading dash), but they start with **MP_**, all in upper case. For example, the environment variable name for the **-procs** option is **MP_PROCS**, and for the **-labelio** option it is **MP_LABELIO**. Setting these two variables like this:

```
$ export MP_PROCS=2
$ export MP_LABELIO=yes
```

allows you to run the **perpms** program with two processes and labeled output, without specifying either with the **poe** command.

Try this command:

```
poe perpms -stdoutmode ordered
```

In the previous example, the program ran with two processes, and the output was labeled.

Now, to see that the environment variable setting lasts for the duration of the session, try running the command below, without specifying the number of processes or labeled I/O.

```
$ poe hostname  
0:202r1n09.hpssl.kgn.ibm.com  
1:202r1n10.hpssl.kgn.ibm.com
```

Notice that the program still ran with two processes and you got labeled output.

Now try overriding the environment variables just set. To do this, use command line options when running POE. Try running the following command:

```
$ poe hostname -procs 4 -labelio no  
  
202r1n09.hpssl.kgn.ibm.com  
202r1n12.hpssl.kgn.ibm.com  
202r1n11.hpssl.kgn.ibm.com  
202r1n10.hpssl.kgn.ibm.com
```

This time, notice that the program ran with four processes and that the output was not labeled. No matter what the environment variables have been set to, you can always override them when you run POE.

To show that this was a temporary override of the environment variable settings, try running the following command again, without specifying any command line options.

```
$ poe hostname  
  
0:202r1n09.hpssl.kgn.ibm.com  
1:202r1n10.hpssl.kgn.ibm.com
```

Once again, the program ran with two processes, and the output was labeled.

Compiling

You probably have programs that you want to run in parallel. Chapter 2, “Message passing,” on page 21 talks about creating parallel programs in a more detail. Right now the topic is compiling a program for POE. You can compile almost any Fortran, C, or C++ program for execution under POE.

Before compiling, you should verify that the following has happened:

- POE is installed on the system
- You are authorized to use POE
- A Fortran, C Compiler, or C ++ compiler is installed on the system.

See *IBM Parallel Environment: MPI Programming Guide* for information on compilation restrictions for POE.

This example, showing how compiling works, uses the *Hello World* program. Here it is in C:

```

/*****
*
* Hello World C Example
*
* To compile:
* mpcc -o hello_world_c hello_world.c
*
*****/
#include<stdlib.h>
#include<stdio.h>
/* Basic program to demonstrate compilation and execution techniques */
int main()
{
printf("Hello, World!\n");
exit(0);
}

```

And here it is in Fortran:

```

C*****
C*
C* Hello World Fortran Example
C*
C* To compile:
C* mpxf1f -o hello_world_f hello_world.f
C*
C*****
C -----
C Basic program to demonstrate compilation and execution techniques
C -----
C program hello

implicit none
write(6,*)'Hello, World!'

stop
end

```

To compile these programs, you just invoke the appropriate compiler script:

```

$ mpcc -o hello_world_c hello_world.c

$ mpfort -o hello_world_f hello_world.f
** main === End of Compilation 1 ===
1501-510 Compilation successful for file hello_world.f.

```

POE scripts **mpcc**, **mpCC**, and **mpfort** link the parallel libraries that allow programs to run in parallel. Script **mpcc** generates thread-aware code by linking in the threaded version of MPI, including the threaded POE utility library.

All the compiler scripts accept all the same options that the non-parallel compilers do, as well as some options specific to POE. For a complete list of all parallel-specific compilation options, see *IBM Parallel Environment: Operation and Use*.

Running one of the POE compiler scripts creates an executable version of the source program that takes advantage of POE. However, before POE can run the program, you need to make sure that it is accessible on each remote node. You can do this by either copying it there, or by mounting the file system that the program resides in to each remote node.

Here is the output of the C program:

```
$ poe hello_world_c -procs 4
```

```
Hello, World!  
Hello, World!  
Hello, World!  
Hello, World!
```

And here is the output of the Fortran program:

```
$ poe hello_world_f -procs 4
```

```
Hello, World!  
Hello, World!  
Hello, World!  
Hello, World!
```

Figure 1. Output from compiler scripts

POE options

There are a number of options (command line flags) that you may want to specify when invoking POE. These options are covered in full detail in *IBM Parallel Environment: Operation and Use* but here are the ones you will most likely need to be familiar with at this stage.

-procs: When you set **-procs**, you are telling POE how many tasks the program will run. You can also set the **MP_PROCS** environment variable to do this (**-procs** can be used to temporarily override it).

-hostfile or -hfile: The default host list file used by POE to allocate nodes is called *host.list*. You can specify a file other than *host.list* by setting the **-hostfile** or **-hfile** options when invoking POE. You can also set the **MP_HOSTFILE** environment variable to do this (**-hostfile** and **-hfile** can be used to temporarily override it).

-labelio: You can set the **-labelio** option when invoking POE so that the output from the parallel tasks of the program are labeled by task id. This becomes especially useful when you are running a parallel program and the output is *unordered*. When you have output that is labeled output, you can easily determine which message the task returned.

You can also set the **MP_LABELIO** environment variable to do this (**-labelio** can be used to temporarily override it).

-infolevel or -ilevel: You can use the **-infolevel** or **-ilevel** options to specify the level of messages you want from POE. There are different levels of informational, warning, and error messages, plus several debugging levels. The **-infolevel** option generates large amounts of output. Use it with care. You can also set the **MP_INFOLEVEL** environment variable to do this (**-infolevel** and **-ilevel** can be used to temporarily override it).

-pmdlog: The **-pmdlog** option lets you specify that diagnostic messages should be logged to a file in **/tmp** on each of the remote nodes of the partition. These diagnostic logs are particularly useful for isolating the cause of abnormal termination. The **-pmdlog** option consumes a significant amount of system

resources. Use it with care. You can also set the **MP_PMDLOG** environment variable to do this (**-pmdlog** can be used to temporarily override it).

-stdoutmode: The **-stdoutmode** option lets you specify how you want the output data from each task in the program to be displayed. When you set this option to *ordered*, the output data from each parallel task is written to its own buffer, and later, all buffers are flushed, in task order, to STDOUT. The examples in this section show you how this works. Using the **-infolevel** option consumes a significant amount of system resources, which may affect performance. You can also set the **MP_STDOUTMODE** environment variable to do this (**-stdoutmode** can be used to temporarily override it).

Managing jobs

So far, you have explicitly specified to POE the set of nodes on which to run the parallel application. You did this by creating a list of hosts in a file called *host.list*, in the directory from which you submitted the parallel job. In the absence of any other instructions, POE selected host names out of this file until it had as many as the number of processes you told POE to use (with the **-procs** option).

Another way to tell POE which hosts to use is with LoadLeveler. LoadLeveler can manage jobs on a networked cluster of xSeries servers.

LoadLeveler is a job management system that allows users to run more jobs in less time by matching the jobs' processing needs with the available resources. LoadLeveler allocates nodes, one job at a time. This is necessary if a parallel application is communicating directly over a high performance interconnect. With the **-euilib** command line option (or the **MP_EUILIB** environment variable), you can specify how you want to do message passing. This option lets you specify the message passing subsystem library implementation, IP or User Space (US), that you wish to use. See *IBM Parallel Environment: Operation and Use* for more information. With LoadLeveler, you can also dedicate the parallel nodes to a single job, so there is no conflict or contention for resources. LoadLeveler allocates nodes from either the host list file, or from a predefined *pool*, which the system administrator usually sets up.

How the nodes are allocated: To know who is allocating the nodes and where they are being allocated from, you must always have a *host list* file or use the **MP_RMPOOL** environment variable or **-rmpool** command line option (unless you are using the **MP_LLFILE** environment variable or the **-llfile** command line option). See *IBM Parallel Environment: Operation and Use* for more information.

The default for the *host list* file is a file named *host.list* in the directory from which the job is submitted. This default may be overridden by the **-hostfile** command line option or the **MP_HOSTFILE** environment variable. For example, the following command:

```
$ poe hostname -procs 4 -hostfile $HOME/myHosts
```

uses a file called **myHosts**, located in the home directory. If the value of the **-hostfile** parameter does not start with a slash (/), it is taken as relative to the current directory. If the value starts with a slash (/), it is taken as a fully-qualified file name.

For specific examples of how a system administrator defines pools, see *Tivoli Workload Scheduler LoadLeveler: Using and Administering*. There is another way to designate the pool on which you want the program to run. If **myHosts** did not contain any pool numbers, you could use the:

- **MP_RMPOOL** environment variable which you can set to a number or string. This setting would last for the duration of the session.
- **-rmpool** command line option to specify a number or string when you invoke the program. This option would override the **MP_RMPOOL** environment variable.

If a host list file named **host.list** exists, or if a host list file is specified using **MP_HOSTFILE** or **-hostfile**, anything you specify with **MP_RMPOOL** or **-rmpool** will be ignored. If a file named **host.list** exists and you want to use **MP_RMPOOL** or **-rmpool** then **MP_HOSTFILE** or **-hostfile** must be set to NULL.

For more information about the **MP_RMPOOL** environment variable or the **-rmpool** command line option, see *IBM Parallel Environment: Operation and Use*.

You cannot have both host names and pool IDs in the same host list file.

The program executes exactly the same way, regardless of whether POE or LoadLeveler allocated the nodes. In the following example, the host list file contains a pool number which causes the job management system to allocate nodes. However, the output is identical to Figure 1 on page 8, where POE allocated the nodes from the host list file.

```
$ poe hello_world_c -procs 4 -hostfile pool.list
```

```
Hello, World!
Hello, World!
Hello, World!
Hello, World!
```

So, if the output looks the same, regardless of how the nodes are allocated, how do you know whether LoadLeveler was used? Well, POE knows a lot that it ordinarily does not tell you. If you coax it with the **-infolevel** option, POE will tell you more information than you ever wanted to know.

Getting a little more information

You can control the level of messages you get from POE as the program executes by using the **-infolevel** option of POE. The default setting is 1 (normal), which says that warning and error messages from POE will be written to STDERR. However, you can use this option to get more information about how the program executes. For example, with **-infolevel** set to 2, you see a couple of different things. First, you will see a message that says that POE has contacted LoadLeveler. Following that, you will see messages that indicate which nodes LoadLeveler passed back to POE for use.

For a description of the various **-infolevel** settings, see *IBM Parallel Environment: Operation and Use*.

Here is the *hello world* program again:

```
$poe ./hello_world -resd yes -procs 2 -labelio yes -infolevel 2
```

You should see output similar to the following:

```
INFO: 0031-364 Contacting LoadLeveler to set and query information for
interactive job
INFO: 0031-380 LoadLeveler step ID is c171f9sq01.ppd.pok.ibm.com.198.0
INFO: 0031-118 Host c171f9sq01.ppd.pok.ibm.com requested for task 0
INFO: 0031-118 Host c171f9sq02.ppd.pok.ibm.com requested for task 1
INFO: 0031-119 Host c171f9sq01.ppd.pok.ibm.com allocated for task 0
INFO: 0031-120 Host address 9.114.136.84 allocated for task 0
```

```

INFO: 0031-377 Using eth0 for mpi euidevice for task 0
INFO: 0031-119 Host c171f9sq02.ppd.pok.ibm.com allocated for task 1
INFO: 0031-120 Host address 9.114.136.85 allocated for task 1
INFO: 0031-377 Using eth0 for mpi euidevice for task 1
  0:INFO: 0031-724 Executing program: <./hello_world>
  1:INFO: 0031-724 Executing program: <./hello_world>
  0:Hello, world!
  0:INFO: 0031-306 pm_atexit: pm_exit_value is 0.
  1:Hello, world!
  1:INFO: 0031-306 pm_atexit: pm_exit_value is 0.
INFO: 0031-656 I/O file STDOUT closed by task 0
INFO: 0031-656 I/O file STDERR closed by task 0
INFO: 0031-656 I/O file STDOUT closed by task 1
INFO: 0031-656 I/O file STDERR closed by task 1
INFO: 0031-251 task 1 exited: rc=0
INFO: 0031-251 task 0 exited: rc=0
INFO: 0031-639 Exit status from pm_respond = 0

```

With **-infolevel** set to 2, you also see messages from each node that indicate the executable they are running and what the return code from the executable is. In the example above, you can differentiate between the **-infolevel** messages that come from POE itself and the messages that come from the remote nodes, because the remote nodes are prefixed with their task ID. If you did not set **-infolevel**, you would see only the output of the executable (Hello world!, in the previous example), interspersed with POE output from remote nodes.

With **-infolevel** set to 3, you get more information. In the following example, use the host list file that contains host names again (as opposed to a Pool ID), when you invoke POE.

Look at the following output. In this case, POE tells you that it is opening the host list file, the nodes it found in the file (along with their Internet addresses), the parameters to the executable being run, and the values of some of the POE parameters.

```
$poe ./hello_world_c -resd yes -procs 2 -labelio yes -infolevel 3
```

You should see output similar to the following:

```

INFO: DEBUG_LEVEL changed from 0 to 1
D1<L1>: Open of file ./host.list successful
D1<L1>: mp_euilib = ip
D1<L1>: 02/28 08:06:05.874667 task 0 c171f9sq01.ppd.pok.ibm.com 10
D1<L1>: 02/28 08:06:05.874802 task 1 c171f9sq02.ppd.pok.ibm.com 10
D1<L1>: node allocation strategy = 2
INFO: 0031-364 Contacting LoadLeveler to set and query information for
interactive job
D1<L1>: 02/28 08:06:05.896603 Calling ll_init_job.
D1<L1>: 02/28 08:06:06.054882 ll_init_job returned.
D1<L1>: Affinity is not requested; MP_TASK_AFFINITY: -1
D1<L1>: 02/28 08:06:06.055027 Job Command String:
#@ job_type = parallel
#@ environment = COPY_ALL
#@ node_usage = shared
#@ bulkxfer = NO
#@ class = No_Class
#@ queue
INFO: 0031-380 LoadLeveler step ID is c171f9sq01.ppd.pok.ibm.com.200.0
INFO: 0031-118 Host c171f9sq01.ppd.pok.ibm.com requested for task 0
INFO: 0031-118 Host c171f9sq02.ppd.pok.ibm.com requested for task 1
INFO: 0031-119 Host c171f9sq01.ppd.pok.ibm.com allocated for task 0
INFO: 0031-120 Host address 9.114.136.84 allocated for task 0
INFO: 0031-377 Using eth0 for mpi euidevice for task 0
INFO: 0031-119 Host c171f9sq02.ppd.pok.ibm.com allocated for task 1

```

```

INFO: 0031-120 Host address 9.114.136.85 allocated for task 1
INFO: 0031-377 Using eth0 for mpi euidevice for task 1
D1<L1>: Entering pm_contact, jobid is 0
D1<L1>: Jobid = 1141162358
D1<L1>: Spawning /etc/pmdv4 on all nodes
D1<L1>: 2 master nodes
D1<L1>: 02/28 08:06:06.340390 Calling ll_spawn_connect for node 0, host name
c171f9sq01.ppd.pok.ibm.com
D1<L1>: 02/28 08:06:06.340692 ll_spawn_connect returned for node 0, socket fd 5,
host name c171f9sq01.ppd.pok.ibm.com
D1<L1>: 02/28 08:06:06.340732 Calling ll_spawn_connect for node 1, host name
c171f9sq02.ppd.pok.ibm.com
D1<L1>: 02/28 08:06:06.341063 ll_spawn_connect returned for node 1, socket fd 6,
host name c171f9sq02.ppd.pok.ibm.com
D1<L1>: 02/28 08:06:06.341100 Calling pm_spawn_ready.
D1<L1>: 02/28 08:06:06.342803 returned from pm_spawn_ready. D1<L1>: Socket file
descriptor for master 0 (c171f9sq01.ppd.pok.ibm.com) is 5
D1<L1>: Socket file descriptor for master 1 (c171f9sq02.ppd.pok.ibm.com) is 6
D1<L1>: SSM_read on socket 5, source = 0, task id: 0, nread: 12, type:3.
D1<L1>: SSM_read on socket 6, source = 1, task id: 1, nread: 12, type:3.
D1<L1>: Leaving pm_contact, jobid is 1141162358
  0:INFO: 0031-724 Executing program: <./hello_world_c>
  1:INFO: 0031-724 Executing program: <./hello_world_c>
  1:INFO: DEBUG_LEVEL changed from 0 to 1
  0:INFO: DEBUG_LEVEL changed from 0 to 1
  1:D1<L1>: After bzero, threadsig.sa_handler = 0
  1:
  0:D1<L1>: After bzero, threadsig.sa_handler = 0
  0:
  0:D1<L1>: Before sigwait: threadsig.sa_handler = 0
  0:D1<L1>: mp_euilib is <ip>
  0:Hello, world!
  0:INFO: 0031-306 pm_atexit: pm_exit_value is 0.
  0:D1<L1>: Sending SSM_EXIT_REQ
  1:D1<L1>: Before sigwait: threadsig.sa_handler = 0
  1:D1<L1>: mp_euilib is <ip>
  1:Hello, world!
  1:INFO: 0031-306 pm_atexit: pm_exit_value is 0.
  1:D1<L1>: Sending SSM_EXIT_REQ
INFO: 0031-656 I/O file STDOUT closed by task 0
INFO: 0031-656 I/O file STDERR closed by task 0
D1<L1>: Accounting data from task 0 for source 0:
INFO: 0031-656 I/O file STDOUT closed by task 1
INFO: 0031-251 task 0 exited: rc=0
INFO: 0031-656 I/O file STDERR closed by task 1
D1<L1>: Accounting data from task 1 for source 1:
INFO: 0031-251 task 1 exited: rc=0
D1<L1>: All remote tasks have exited: maxx_errcode = 0
INFO: 0031-639 Exit status from pm_respond = 0
D1<L1>: Maximum return code from user = 0

```

The **-infolevel** messages give you more information about what is happening on the home node, but if you want to see what is happening on the remote nodes, you need to use the **-pmdlog** option. If you set **-pmdlog** to a value of **yes**, a log is written to each of the remote nodes that tells you what POE did while running each task.

If you issue the following command, a file is written in **/tmp**, of each remote node, called **mplog.jobid.taskid**,

```
$ poe hello_world -procs 4 -pmdlog yes
```

If **-infolevel** is set to 3 or higher, The job ID will be displayed in the output. If you do not know what the job ID is, it is probably the most recent log file. If you are

sharing the node with other POE users, the job ID will be *one* of the most recent log files (but you own the file, so you should be able to tell).

Here is a sample log file. In this example, all four tasks are running on the same node. For more information about how POE runs with multiple tasks on the same node, see Appendix A, “A sample program to illustrate messages,” on page 53.

```
Parallel Environment pmd4 version @(#) 2003/06/11 13:19:38
The ID of this process is 31326
The version of this pmd for version checking is 4100
The hostname of this node is c171f9sq01
The short hostname of this node is
The taskid of this task is 0
HOMENAME: c171f9sq01.ppd.pok.ibm.com
USERID: 1079
USERNAME: voe3
GROUPID: 1079
GROUPNAME: voe3
PWD: /u/voe3/pfc
PRIORITY: 0
NPROCS: 4
PMDLOG: 1
NEWJOB: 0
LIBPATH: /opt/ibmhpc/ppe.poe/lib/libmpi
LD_LIBRARY_PATH: _EMPTY_PATH
VERSION (of home node): 4100
JOBID: 1141163670
ENVC recv'd
envc: 70
envc is 70
env[0] = MODULE_VERSION_STACK=3.1.6
env[1] = LESSKEY=/etc/lesskey.bin
env[2] = NNTPSERVER=news
env[3] = INFODIR=/usr/local/info:/usr/share/info:/usr/info
env[4] =MANPATH=/opt/csm/man:/usr/local/man:/usr/share/man:/usr/X11R6/
man:/opt/gnome/share/man
env[5] = HOSTNAME=c171f9sq01
env[6] = GNOME2_PATH=/usr/local:/opt/gnome:/usr
env[7] = XKEYSYMDB=/usr/X11R6/lib/X11/XKeysymDB
env[8] = HOST=c171f9sq01
env[9] = SHELL=/bin/bash
env[10] = TERM= aixterm
env[11] = PROFILEREAD=true
env[12] = HISTSIZE=1000
env[13] = GNOME_PATH=/opt/gnome:/usr
env[14] = QTDIR=/usr/lib/qt3
env[15] = OLDPWD=/opt/ibmhpc/ppe.poe
env[16] = JRE_HOME=/usr/lib/java/jre
env[17] = USER=voe3
env[18] = LS_COLORS=
env[19] = XNLSPATH=/usr/X11R6/lib/X11/nls
env[20] = HOSTTYPE=ppc64
env[21] = PAGER=less
env[22] = MINICOM=-c on
env[23] = MODULE_VERSION=3.1.6
env[24] = GNOMEDIR=/opt/gnome
env[25] = PATH=/opt/csm/bin:/u/voe3/bin:/usr/local/bin:/usr/bin:/usr/X11R6/
bin:/bin:/usr/games:/opt/bin:/opt/gnome/bin:/opt/kde3/bin:/usr/lib/java/jre/bin:/
opt/bin:/u/voe3/bin
env[26] = CPU=ppc64
env[27] = JAVA_BINDIR=/usr/lib/java/jre/bin
env[28] = INPUTRC=/etc/inputrc
env[29] = PWD=/u/voe3/pfc
env[30] = JAVA_HOME=/usr/lib/java/jre
env[31] = LANG=en_US.UTF-8
env[32] = MODULEPATH=/usr/share/modules/versions:/usr/share/modules/
modulefiles
```

```

env[33] = LOADEDMODULES=
env[34] = TEXINPUTS=./u/voe3/.TeX:/usr/share/doc/.TeX:/usr/doc/.TeX
env[35] = SHLVL=1
env[36] = HOME=/u/voe3
env[37] = LESS_ADVANCED_PREPROCESSOR=no
env[38] = OSTYPE=linux
env[39] = LS_OPTIONS=-N --color=none -T 0
env[40] = XCURSOR_THEME=crystalwhite
env[41] = no_proxy=localhost
env[42] = WINDOWMANAGER=/usr/X11R6/bin/kde
env[43] = GTK_PATH=/usr/local/lib/gtk-2.0:/opt/gnome/lib/gtk-2.0:/usr/lib/gtk-2.0
env[44] = LESS=-M -I
env[45] = MACHTYPE=ppc64-suse-linux
env[46] = LOGNAME=voe3
env[47] = CVS_RSH=ssh
env[48] = MODULESHOME=/usr/share/modules
env[49] = PKG_CONFIG_PATH=/opt/gnome/lib/pkgconfig
env[50] = LESSOPEN=lessopen.sh %s
env[51] = ACLOCAL_PATH=/opt/gnome/share/aclocal
env[52] = INFOPATH=/usr/local/info:/usr/share/info:/usr/info:/opt/gnome/share/info
env[53] = LESSCLOSE=lessclose.sh %s %s
env[54] = G_BROKEN_FILENAMES=1
env[55] = JAVA_ROOT=/usr/lib/java
env[56] = COLORTERM=1
env[57] = _=/usr/bin/poe
env[58] = NLSPATH=/usr/share/locale/%L/%N
env[59] = MP_PROCS=4
env[60] = MP_PMDLOG=YES
env[61] = MP_EUIDEVICE=en0
env[62] = MP_PGMMODEL=SPMD
env[63] = MP_TASK_AFFINITY=-1
env[64] = MP_MSG_API=MPI
env[65] = MP_PRIORITY_LOG=YES
env[66] = MP_PRIORITY_NTP=NO
env[67] = MP_ISATTY_STDIN=1
env[68] = MP_ISATTY_STDOUT=1
env[69] = MP_ISATTY_STDERR=1
Couldn't open /etc/poe.limits
MASTERS: 1
TASKS: 4:0:1:2:3
Total number of tasks is 4
Task id for task 1 is 0
Task id for task 2 is 1
Task id for task 3 is 2
Task id for task 4 is 3
TASK_ENV: 0:1 MP_CHILD_INET_ADDR=@1:9.114.136.84,ip (1:1 MP_CHILD_INET_ADDR=\
@1:9.114.136.84,ip
(2:1 MP_CHILD_INET_ADDR=@1:9.114.136.84,ip (3:1 MP_CHILD_INET_ADDR=\
@1:9.114.136.84,ip
(Number of environment variables is 1
Environment specific data for task 1, task id 0 :
-- MP_CHILD_INET_ADDR=@1:9.114.136.84,ip
Number of environment variables is 1
Environment specific data for task 2, task id 1 :
-- MP_CHILD_INET_ADDR=@1:9.114.136.84,ip
Number of environment variables is 1
Environment specific data for task 3, task id 2 :
-- MP_CHILD_INET_ADDR=@1:9.114.136.84,ip
Number of environment variables is 1
Environment specific data for task 4, task id 3 :
-- MP_CHILD_INET_ADDR=@1:9.114.136.84,ip
Initial data msg received and parsed
Info level = 1
Doing ruserok() user validation
User validation complete
About to do user root chk
User root check complete

```

```

task information parsed
STDOUT socket SO_SNDBUF set to 4194048
STDOUT socket SO_RCVBUF set to 87552
newjob is 0.
msg read, type is 13
string = <./hello_world_c ^@>
SSM_CMD_STR recv'd
command_string is <./hello_world_c >
0: pm_putargs: argc = 1, k = 1
1: pm_putargs: argc = 1, k = 1
2: pm_putargs: argc = 1, k = 1
3: pm_putargs: argc = 1, k = 1
SSM_CMD_STR parsed
child pipes created
parent: task 0 forked, child pid is 31327
attach data sent for task 0
parent: task 1 forked, child pid is 31328
attach data sent for task 1
parent: task 2 forked, child pid is 31329
attach data sent for task 2
parent: task 3 forked, child pid is 31330
attach data sent for task 3
Entering child section...
child: pipes successfully duped for task 0
0: MP_COMMON_TASKS is <3:1:2:3>
0: partition id is <1141163670>
after initgroups (*group_struct).gr_gid = 1079
after initgroups (*group_struct).gr_name = voe3
Entering child section...
Entering child section...
child: pipes successfully duped for task 2
child: pipes successfully duped for task 1
2: MP_COMMON_TASKS is <3:0:1:3>
2: partition id is <1141163670>
pmd child: core limit is 1073741824, hard limit is -1
pmd child: rss limit is -1, hard limit is -1
pmd child: stack limit is -1, hard limit is -1
pmd child: data segment limit is -1, hard limit is -1
pmd child: cpu time limit is -1, hard limit is -1
pmd child: file size limit is -1, hard limit is -1
pmd child: process virtual memory limit is -1, hard limit is -1
pmd child: process virtual memory locked limit is -1, hard limit is -1
pmd child: process files locked limit is -1, hard limit is -1
pmd child: process opened files limit is 1024, hard limit is 1024
pmd child: process forked process limit is 253952, hard limit is 253952
0: (*group_struct).gr_gid = 1079
0: (*group_struct).gr_name = voe3
Entering child section...
1: MP_COMMON_TASKS is <3:0:2:3>
1: partition id is <1141163670>
child: pipes successfully duped for task 3
3: MP_COMMON_TASKS is <3:0:1:2>
1: (*group_struct).gr_gid = 1079
1: (*group_struct).gr_name = voe3
3: partition id is <1141163670>
2: (*group_struct).gr_gid = 1079
2: (*group_struct).gr_name = voe3
3: (*group_struct).gr_gid = 1079
3: (*group_struct).gr_name = voe3
0: userid, groupid and cwd set!
0: current directory is /u/voe3/pfc
3: userid, groupid and cwd set!
3: current directory is /u/voe3/pfc
0: about to start the user's program
3: about to start the user's program
3: argument list:
argv[0] for task 3 = ./hello_world_c

```

```

argv[1] (in hex) = 0
child: environment for task 3:
    task_env[0] = MP_CHILD_INET_ADDR=@1:9.114.136.84,ip

3: LIBPATH = /opt/ibmhcp/ppe.poe/lib/libmpi
0: argument list:
argv[0] for task 0 = ./hello_world_c

argv[1] (in hex) = 0
child: environment for task 0:
    task_env[0] = MP_CHILD_INET_ADDR=@1:9.114.136.84,ip

child: common environment data for all tasks:
    env[0] = MODULE_VERSION_STACK=3.1.6
    env[1] = LESSKEY=/etc/lesskey.bin
    env[2] = NNTPSERVER=news
    env[3] = INFODIR=/usr/local/info:/usr/share/info:/usr/info
    env[4] = MANPATH=/opt/csm/man:/usr/local/man:/usr/share/man:/usr/X11R6/man:\
/opt/gnome/share/man
    env[5] = HOSTNAME=c171f9sq01
    env[6] = GNOME2_PATH=/usr/local:/opt/gnome:/usr
    env[7] = XKEYSYMDB=/usr/X11R6/lib/X11/XKeysymDB
    env[8] = HOST=c171f9sq01
    env[9] = SHELL=/bin/bash
    env[10] = TERM=aixterm
    env[11] = PROFILEREAD=true
    env[12] = HISTSIZE=1000
    env[13] = GNOME_PATH=/opt/gnome:/usr
    env[14] = QTDIR=/usr/lib/qt3
    env[15] = OLDPWD=/opt/ibmhcp/ppe.poe
    env[16] = JRE_HOME=/usr/lib/java/jre
    env[17] = USER=voe3
    env[18] = LS_COLORS=
    env[19] = XNLSPATH=/usr/X11R6/lib/X11/nls
    env[20] = HOSTTYPE=ppc64
    env[21] = PAGER=less
    env[22] = MINICOM=-c on
env[23] = MODULE_VERSION=3.1.6
    env[24] = GNOMEDIR=/opt/gnome
    env[25] = PATH=/opt/csm/bin:/u/voe3/bin:/usr/local/bin:/usr/bin:\
/usr/X11R6/bin:/bin:/usr/games:\
/opt/bin:/opt/gnome/bin:/opt/kde3/bin:/usr/lib/java/jre/bin:/opt/bin:/u/voe3/bin
    env[26] = CPU=ppc64
    env[27] = JAVA_BINDIR=/usr/lib/java/jre/bin
    env[28] = INPUTRC=/etc/inputrc
    env[29] = PWD=/u/voe3/pfc
    env[30] = JAVA_HOME=/usr/lib/java/jre
    env[31] = LANG=en_US.UTF-8
    env[32] = MODULEPATH=/usr/share/modules/versions:/usr/share/modules/modulefiles
    env[33] = LOADEDMODULES=
    env[34] = TEXINPUTS=:/u/voe3/.TeX:/usr/share/doc/.TeX:/usr/doc/.TeX
    env[35] = SHLVL=1
    env[36] = HOME=/u/voe3
    env[37] = LESS_ADVANCED_PREPROCESSOR=no
    env[38] = OSTYPE=linux
    env[39] = LS_OPTIONS=-N --color=none -T 0
    env[40] = XCURSOR_THEME=crystalwhite
    env[41] = no_proxy=localhost
    env[42] = WINDOWMANAGER=/usr/X11R6/bin/kde
    env[43] = GTK_PATH=/usr/local/lib/gtk-2.0:/opt/gnome/lib/gtk-2.0:\
/usr/lib/gtk-2.0
    env[44] = LESS=-M -I
    env[45] = MACHTYPE=ppc64-suse-linux
    env[46] = LOGNAME=voe3
    env[47] = CVS_RSH=ssh
    env[48] = MODULESHOME=/usr/share/modules

```

```

env[49] = PKG_CONFIG_PATH=/opt/gnome/lib/pkgconfig
env[50] = LESSOPEN=lessopen.sh %s
env[51] = ACLOCAL_PATH=/opt/gnome/share/aclocal
env[52] = INFOPATH=/usr/local/info:/usr/share/info:/usr/info:\
/opt/gnome/share/info
env[53] = LESSCLOSE=lessclose.sh %s %s
env[54] = G_BROKEN_FILENAMES=1
env[55] = JAVA_ROOT=/usr/lib/java
env[56] = COLORTERM=1
env[57] = _=/usr/bin/poe
env[58] = NLSPATH=/usr/share/locale/%L/%N
env[59] = MP_PROCS=4
env[60] = MP_PMDLOG=YES
env[61] = MP_EUIDEVICE=en0
env[62] = MP_PGMMODEL=SPMD
env[63] = MP_TASK_AFFINITY=-1
env[64] = MP_MSG_API=MPI
env[65] = MP_PRIORITY_LOG=YES
env[66] = MP_PRIORITY_NTP=NO
env[67] = MP_ISATTY_STDIN=1
env[68] = MP_ISATTY_STDOUT=1
env[69] = MP_ISATTY_STDERR=1

0: LIBPATH = /opt/ibmhpc/ppe.poe/lib/libmpi
1: userid, groupid and cwd set!
1: current directory is /u/voe3/pfc
1: about to start the user's program
1: argument list:
argv[0] for task 1 = ./hello_world_c
argv[1] (in hex) = 0
child: environment for task 1:
    task_env[0] = MP_CHILD_INET_ADDR=@1:9.114.136.84,ip

1: LIBPATH = /opt/ibmhpc/ppe.poe/lib/libmpi
2: userid, groupid and cwd set!
2: current directory is /u/voe3/pfc
2: about to start the user's program
2: argument list:
argv[0] for task 2 = ./hello_world_c

argv[1] (in hex) = 0
child: environment for task 2:
    task_env[0] = MP_CHILD_INET_ADDR=@1:9.114.136.84,ip

2: LIBPATH = /opt/ibmhpc/ppe.poe/lib/libmpi
select: rc = 4
pulse is on, curr_time is 1141132345, send_time is 0, select time is 599
pulse sent at 1141132345 count is 0
pmd parent: STDOUT read OK for task 0
0: STDOUT: Hello, world!

0: pmd parent: cntl pipe read OK:
0: pmd parent: type: 26, srce: 0, dest: -2, bytes: 6
parent: SSM_CHILD_PID: 31327
pmd parent: STDOUT read OK for task 1
1: STDOUT: Hello, world!

1: pmd parent: cntl pipe read OK:
1: pmd parent: type: 26, srce: 1, dest: -2, bytes: 6
parent: SSM_CHILD_PID: 31328
select: rc = 4
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
0: pmd parent: cntl pipe read OK:
0: pmd parent: type: 47, srce: 0, dest: -2, bytes: 5
parent: child's version is 4100.
parent: home node version is 4100.
parent: this pmd version is 4100.

```

```

1: pmd parent: cntl pipe read OK:
1: pmd parent: type: 47, srce: 1, dest: -2, bytes: 5
parent: childs version is 4100.
parent: home node version is 4100.
parent: this pmd version is 4100.
pmd parent: STDOUT read OK for task 2
2: STDOUT: Hello, world!

2: pmd parent: cntl pipe read OK:
2: pmd parent: type: 26, srce: 2, dest: -2, bytes: 6
parent: SSM_CHILD_PID: 31329
select: rc = 3
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
0: pmd parent: cntl pipe read OK:
0: pmd parent: type: 17, srce: 0, dest: -1, bytes: 2
1: pmd parent: cntl pipe read OK:
1: pmd parent: type: 17, srce: 1, bytes: 2
2: pmd parent: cntl pipe read OK:
2: pmd parent: type: 47, srce: 2, dest: -2, bytes: 5
parent: childs version is 4100.
parent: home node version is 4100.
parent: this pmd version is 4100.
select: rc = 1
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
2: pmd parent: cntl pipe read OK:
2: pmd parent: type: 17, srce: 2, dest: -1, bytes: 2
select: rc = 1
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
3: pmd parent: cntl pipe read OK:
3: pmd parent: type: 26, srce: 3, dest: -2, bytes: 6
parent: SSM_CHILD_PID: 31330
select: rc = 1
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
3: pmd parent: cntl pipe read OK:
3: pmd parent: type: 47, srce: 3, dest: -2, bytes: 5
parent: childs version is 4100.
parent: home node version is 4100.
parent: this pmd version is 4100.
select: rc = 1
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
pmd parent: STDOUT read OK for task 3
3: STDOUT: Hello, world!

select: rc = 1
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
3: pmd parent: cntl pipe read OK:
3: pmd parent: type: 17, srce: 3, dest: -1, bytes: 2
select: rc = 1
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
in pmd select, SSM_read ok, SSM_type=34.
pulse received at 1141132345 received count is -5448
select: rc = 1
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
in pmd select, SSM_read ok, SSM_type=5.
select: rc = 3
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
pmd send SSM_IO_CLOSED to poe for stdout_open
0: count = 0 on stderr
pmd send SSM_IO_CLOSED to poe for stderr_open
in pmd signal handler for task 0, signal 17
0: wait status is 00000000
Exiting child for task 0, PID: 31327
err_data for task 0 is 0
select: rc = 3
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
pmd send SSM_IO_CLOSED to poe for stdout_open
1: count = 0 on stderr

```

```

pmd send SSM_IO_CLOSED to poe for stderr_open
in pmd signal handler for task 1, signal_17
1: wait status is 00000000
Exiting child for task 1, PID: 31328
err_data for task 1 is 0
2: wait status is 00000000
Exiting child for task 2, PID: 31329
err_data for task 2 is 0
in pmd signal handler, wait returned 0...
select: rc = 3
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
pmd send SSM_IO_CLOSED to poe for stdout_open
2: count = 0 on stderr
pmd send SSM_IO_CLOSED to poe for stderr_open
select: rc = 3
pulse is on, curr_time is 1141132345, send_time is 1141132345, select time is 599
pmd send SSM_IO_CLOSED to poe for stdout_open
3: count = 0 on stderr
pmd send SSM_IO_CLOSED to poe for stderr_open
in pmd signal handler for task 3, signal_17
3: wait status is 00000000
Exiting child for task 3, PID: 31330
err_data for task 3 is 0
parent: child exited and all pipes closed for all tasks
err_data for task 0 is 0
err_data for task 1 is 0
err_data for task 2 is 0
err_data for task 3 is 0
pmd_exit reached!, exit code is 0

```

Appendix A, “A sample program to illustrate messages,” on page 53 includes an example of setting **-infolevel** to 6, and explains the important lines of output.

Chapter 2. Message passing

If you are familiar with message passing parallel programming, and you are familiar with message passing protocols, you can skip ahead to Chapter 3, “Diagnosing and correcting common problems,” on page 35 for a discussion on using the PE tools. If you are familiar with message passing parallel programming, but you would like to know more about the PE message passing protocols, look at the information in “Protocols supported” on page 32.

This section discusses some of the techniques for creating a parallel program, using message passing, and the various advantages and pitfalls associated with each technique. It does not provide an in-depth tutorial on writing parallel programs. Instead, it is an introduction to basic message passing parallel concepts.

To create a successful parallel program start with a working sequential program. Complex sequential programs are difficult to get working correctly, without also having to worry about the additional complexity introduced by parallelism and message passing. It is easier to convert a working serial program to parallel, than it is to create a parallel program from scratch. As you become proficient at creating parallel programs, you will develop an awareness of which sequential techniques translate better into parallel implementations. Once aware, you can then make a point of using these techniques in your sequential programs. In this section, contains information on some of the fundamentals of creating parallel programs.

There are two common techniques for turning a sequential program into a parallel program; *data decomposition* and *functional decomposition*. Data decomposition means distributing the data that the program is processing among the parallel tasks. Each parallel task does approximately the same thing but on a different set of data. With functional decomposition, the function that the application is performing is distributed among the tasks. Each task operates on the same data, but does something different. Most parallel programs do not use data decomposition or functional decomposition exclusively. Rather, they use a mixture of the two, weighted more toward one type or the other. One way to implement either form of decomposition is through the use of message passing.

The message passing model

The message passing model of communication is typically used in distributed memory systems, where each processor node owns private memory, and is linked by an interconnection network. With message passing, each task operates exclusively in a private environment, but must cooperate with other tasks to interact. In this situation, tasks must exchange messages to interact with one another.

The challenge of the message passing model is in reducing message traffic over the interconnection network while ensuring that the correct and updated values of the passed data are promptly available to the tasks, when required. Optimizing message traffic boosts performance.

Synchronization is the act of forcing events to occur at the same time or in a certain order. Synchronization requires taking into account the logical dependence and the order of precedence among the tasks. You can describe the message passing model as self-synchronizing because the mechanism of sending and receiving messages involves implicit synchronization points. To put it another way, a message cannot be received if it has not already been sent.

Data decomposition

A good technique for making a sequential application parallel is to look for loops where each iteration does not depend on any prior iteration (this is also a prerequisite for either *unrolling* or eliminating loops). An example of a loop that has dependencies on prior iterations is the loop for computing the Factorial series. The value calculated by each iteration depends on the value resulting from the previous pass. If each iteration of a loop does not depend on a previous iteration, the data being processed can be processed in parallel, with two or more iterations being performed simultaneously.

The C program example below includes a loop with independent iterations. This example does not include the routines for computing the coefficient and determinant because they are not part of the parallelization at this point.

```
/******  
*  
* Matrix Inversion Program - serial version  
*  
* To compile:  
* cc -o inverse_serial inverse_serial.c  
*  
*****/  
  
#include<stdlib.h>  
#include<stdio.h>  
#include<assert.h>  
#include<errno.h>  
  
float determinant(float **matrix,  
                 int size,  
                 int * used_rows,  
                 int * used_cols,  
                 int depth);  
float coefficient(float **matrix,int size, int row, int col);  
void print_matrix(FILE * fptr,float ** mat,int rows, int cols);  
float test_data[8][8] = {  
    {4.0, 2.0, 4.0, 5.0, 4.0, -2.0, 4.0, 5.0},  
    {4.0, 2.0, 4.0, 5.0, 3.0, 9.0, 12.0, 1.0 },  
    {3.0, 9.0, -13.0, 15.0, 3.0, 9.0, 12.0, 15.0},  
    {3.0, 9.0, 12.0, 15.0, 4.0, 2.0, 7.0, 5.0 },  
    {2.0, 4.0, -11.0, 10.0, 2.0, 4.0, 11.0, 10.0 },  
    {2.0, 4.0, 11.0, 10.0, 3.0, -5.0, 12.0, 15.0 },  
    {1.0, -2.0, 4.0, 10.0, 3.0, 9.0, -12.0, 15.0 } ,  
    {1.0, 2.0, 4.0, 10.0, 2.0, -4.0, -11.0, 10.0 } ,  
};  
#define ROWS 8  
  
int main(int argc, char **argv)  
{  
  
    float **matrix;  
    float **inverse;  
    int rows,i,j;  
    float determ;  
    int * used_rows, * used_cols;  
  
    rows = ROWS;  
  
    /* Allocate markers to record rows and columns to be skipped */  
    /* during determinant calculation */  
    used_rows = (int *) malloc(rows*sizeof(*used_rows));  
    used_cols = (int *) malloc(rows*sizeof(*used_cols));  
  
    /* Allocate working copy of matrix and initialize it from static copy */
```

```

matrix = (float **) malloc(rows*sizeof(*matrix));
inverse = (float **) malloc(rows*sizeof(*inverse));
for(i=0;i<rows;i++)
{
    matrix[i] = (float *) malloc(rows*sizeof(**matrix));
    inverse[i] = (float *) malloc(rows*sizeof(**inverse));
    for(j=0;j<rows;j++)
        matrix[i][j] = test_data[i][j];
}

/* Compute and print determinant */
printf("The determinant of\n\n");
print_matrix(stdout,matrix,rows,rows);
determ=determinant(matrix,rows,used_rows,used_cols,0);
printf("\nis %f\n",determ);
fflush(stdout);
assert(determ!=0);

for(i=0;i<rows;i++)
{
    for(j=0;j<rows;j++)
    {
        inverse[j][i] = coefficient(matrix,rows,i,j)/determ;
    }
}

printf("The inverse is\n\n");
print_matrix(stdout,inverse,rows,rows);

return (0);
}

```

Before talking about making the algorithm parallel, look at what is necessary to create the program with PE. The example below shows the same program, but it is now aware of PE. You do this by using three calls in the beginning of the routine, and one at the end.

The first of these calls (**MPI_Init**) initializes the *MPI* environment, and the last call (**MPI_Finalize**) closes the environment. **MPI_Comm_size** sets the variable **tasks** to the total number of parallel tasks running this application, and **MPI_Comm_rank** sets **me** to the task ID of the particular instance of the parallel code that invoked it.

MPI_Comm_size actually gets the size of the communicator you pass in and **MPI_COMM_WORLD** is a predefined communicator that includes everybody. For more information about these calls, *IBM Parallel Environment: MPI Subroutine Reference* or other MPI publications may be of some help.

```

/*****
*
* Matrix Inversion Program - serial version enabled for parallel environment
*
* To compile:
* mpcc -g -o inverse_parallel_enabled inverse_parallel_enabled.c
*
*****/

#include<stdlib.h>
#include<stdio.h>
#include<assert.h>
#include<errno.h>
#include<mpi.h>

float determinant(float **matrix,int size, int * used_rows, int * used_cols,
                 int depth);
float coefficient(float **matrix,int size, int row, int col);

```

```

void print_matrix(FILE * fptr,float ** mat,int rows, int cols);
float test_data[8][8] = {
    {4.0, 2.0, 4.0, 5.0, 4.0, -2.0, 4.0, 5.0},
    {4.0, 2.0, 4.0, 5.0, 3.0, 9.0, 12.0, 1.0 },
    {3.0, 9.0, -13.0, 15.0, 3.0, 9.0, 12.0, 15.0},
    {3.0, 9.0, 12.0, 15.0, 4.0, 2.0, 7.0, 5.0 },
    {2.0, 4.0, -11.0, 10.0, 2.0, 4.0, 11.0, 10.0 },
    {2.0, 4.0, 11.0, 10.0, 3.0, -5.0, 12.0, 15.0 },
    {1.0, -2.0, 4.0, 10.0, 3.0, 9.0, -12.0, 15.0 },
    {1.0, 2.0, 4.0, 10.0, 2.0, -4.0, -11.0, 10.0 } ,
};
#define ROWS 8

int me, tasks, tag=0;

int main(int argc, char **argv)
{

    float **matrix;
    float **inverse;
    int rows,i,j;
    float determ;
    int * used_rows, * used_cols;

    MPI_Status status[ROWS]; /* Status of messages */
    MPI_Request req[ROWS]; /* Message IDs */

    MPI_Init(&argc,&argv); /* Initialize MPI */
    MPI_Comm_size(MPI_COMM_WORLD,&tasks); /* How many parallel tasks are there?*/
    MPI_Comm_rank(MPI_COMM_WORLD,&me); /* Who am I? */

    rows = ROWS;

    /* Allocate markers to record rows and columns to be skipped */
    /* during determinant calculation */
    used_rows = (int *) malloc(rows*sizeof(*used_rows));
    used_cols = (int *) malloc(rows*sizeof(*used_cols));

    /* Allocate working copy of matrix and initialize it from static copy */
    matrix = (float **) malloc(rows*sizeof(*matrix));
    inverse = (float **) malloc(rows*sizeof(*inverse));
    for(i=0;i<rows;i++)
    {
        matrix[i] = (float *) malloc(rows*sizeof(**matrix));
        inverse[i] = (float *) malloc(rows*sizeof(**inverse));
        for(j=0;j<rows;j++)
            matrix[i][j] = test_data[i][j];
    }

    /* Compute and print determinant */
    printf("The determinant of\n\n");
    print_matrix(stdout,matrix,rows,rows);
    determ=determinant(matrix,rows,used_rows,used_cols,0);
    printf("\nis %f\n",determ);
    fflush(stdout);

    for(i=0;i<rows;i++)
    {
        for(j=0;j<rows;j++)
        {
            inverse[j][i] = coefficient(matrix,rows,i,j)/determ;
        }
    }

    printf("The inverse is\n\n");
    print_matrix(stdout,inverse,rows,rows);

```

```

    /* Wait for all parallel tasks to get here, then quit */
    MPI_Barrier(MPI_COMM_WORLD);
    MPI_Finalize();

    exit(0);
}

float determinant(float **matrix, int size, int * used_rows, int * used_cols,
                 int depth)
{
    int col1, col2, row1, row2;
    int j, k;
    float total=0;
    int sign = 1;

    /* Find the first unused row */
    for(row1=0; row1<size; row1++)
    {
        for(k=0; k<depth; k++)
        {
            if(row1==used_rows[k]) break;
        }
        if(k>=depth) /* this row is not used */
            break;
    }
    assert(row1<size);

    if(depth==(size-2))
    {
        /* There are only 2 unused rows/columns left */

        /* Find the second unused row */
        for(row2=row1+1; row2<size; row2++)
        {
            for(k=0; k<depth; k++)
            {
                if(row2==used_rows[k]) break;
            }
            if(k>=depth) /* this row is not used */
                break;
        }
        assert(row2<size);

        /* Find the first unused column */
        for(col1=0; col1<size; col1++)
        {
            for(k=0; k<depth; k++)
            {
                if(col1==used_cols[k]) break;
            }
            if(k>=depth) /* this column is not used */
                break;
        }
        assert(col1<size);

        /* Find the second unused column */
        for(col2=col1+1; col2<size; col2++)
        {
            for(k=0; k<depth; k++)
            {
                if(col2==used_cols[k]) break;
            }
            if(k>=depth) /* this column is not used */
                break;
        }
        assert(col2<size);
    }
}

```

```

/* Determinant = m11*m22-m12*m21 */
return matrix[row1][col1]*matrix[row2][col2]
-matrix[row2][col1]*matrix[row1][col2];
}

/* There are more than 2 rows/columns in the matrix being processed */
/* Compute the determinant as the sum of the product of each element */
/* in the first row and the determinant of the matrix with its row */
/* and column removed */
total = 0;

used_rows[depth] = row1;
for(col1=0;col1<size;col1++)
{
    for(k=0;k<depth;k++)
    {
        if(col1==used_cols[k]) break;
    }
    if(k<depth) /* This column is used */
        continue;
    used_cols[depth] = col1;
    total += sign*matrix[row1][col1]*determinant(matrix,size,
used_rows,used_cols,depth+1);
    sign=(sign==1)?-1:1;
}
return total;
}

void print_matrix(FILE * fptr,float ** mat,int rows, int cols)
{
    int i,j;
    for(i=0;i<rows;i++)
    {
        for(j=0;j<cols;j++)
        {
            fprintf(fptr,"%10.4f ",mat[i][j]);
        }
        fprintf(fptr,"\n");
    }
    fflush(fptr);
}

float coefficient(float **matrix,int size, int row, int col)
{
    float coef;
    int * ur, *uc;

    ur = malloc(size*sizeof(matrix));
    uc = malloc(size*sizeof(matrix));
    ur[0]=row;
    uc[0]=col;
    coef = (((row+col)%2)?-1:1)*determinant(matrix,size,ur,uc,1);
    return coef;
}

```

In this particular example each parallel task is going to determine the entire inverse matrix, and they are all going to print it out. In the previous section, the output of all the tasks will be intermixed, so it will be difficult to figure out what the answer really is.

A better approach is to distribute the work among several parallel tasks and collect the results when they are done. In this example, the loop that computes the elements of the inverse matrix simply goes through the elements of the inverse

matrix, computes the coefficient, and divides it by the determinant of the matrix. Since there is no relationship between elements of the inverse matrix, they can all be computed in parallel.

Every communication call has an associated cost, so you need to balance the benefit of parallelism with the cost of communication. If you were to totally parallelize the inverse matrix element computation, each element would be derived by a separate task. The cost of collecting those individual values back into the inverse matrix would be significant. It might also outweigh the benefit of having reduced the computation cost and time by running the job in parallel. So, instead, you are going to compute the elements of each row in parallel, and send the values back, one row at a time. This way you spread some of the communication overhead over several data values. In this case, you will execute loop 1 in parallel in this next example.

```
*****
*
* Matrix Inversion Program - First parallel implementation
* To compile:
* mpicc -g -o inverse_parallel inverse_parallel.c
*
*****

#include<stdlib.h>
#include<stdio.h>
#include<assert.h>
#include<errno.h>
#include<mpi.h>
float determinant(float **matrix,int size, int * used_rows,
                  int * used_cols, int depth);
float coefficient(float **matrix,int size, int row, int col);
void print_matrix(FILE * fptr,float ** mat,int rows, int cols);

float test_data[8][8] = {
    {4.0, 2.0, 4.0, 5.0, 4.0, -2.0, 4.0, 5.0},
    {4.0, 2.0, 4.0, 5.0, 3.0, 9.0, 12.0, 1.0 },
    {3.0, 9.0, -13.0, 15.0, 3.0, 9.0, 12.0, 15.0},
    {3.0, 9.0, 12.0, 15.0, 4.0, 2.0, 7.0, 5.0 },
    {2.0, 4.0, -11.0, 10.0, 2.0, 4.0, 11.0, 10.0 },
    {2.0, 4.0, 11.0, 10.0, 3.0, -5.0, 12.0, 15.0 },
    {1.0, -2.0, 4.0, 10.0, 3.0, 9.0, -12.0, 15.0 },
    {1.0, 2.0, 4.0, 10.0, 2.0, -4.0, -11.0, 10.0 } ,
};

#define ROWS 8
int me, tasks, tag=0;

int main(int argc, char **argv)
{
    float **matrix;
    float **inverse;
    int rows,i,j;
    float determ;
    int * used_rows, * used_cols;

    MPI_Status status[ROWS]; /* Status of messages */
    MPI_Request req[ROWS]; /* Message IDs */

    MPI_Init(&argc,&argv); /* Initialize MPI */
    MPI_Comm_size(MPI_COMM_WORLD,&tasks); /* How many parallel tasks are there?*/
    MPI_Comm_rank(MPI_COMM_WORLD,&me); /* Who am I? */

    rows = ROWS;
```

```

/* You need exactly one task for each row of the matrix plus one task */
/* to act as coordinator. If you didn't have this, the last task */
/* reports the error (so everybody doesn't put out the same message */
if(tasks!=rows+1)
{
    if(me==tasks-1)
    fprintf(stderr,"%d tasks required for this demo"
"one more than the number of rows in matrix\n",rows+1);
    exit(-1);
}
/* Allocate markers to record rows and columns to be skipped */
/* during determinant calculation */
used_rows = (int *) malloc(rows*sizeof(*used_rows));
used_cols = (int *) malloc(rows*sizeof(*used_cols));

/* Allocate working copy of matrix and initialize it from static copy */
matrix = (float **) malloc(rows*sizeof(*matrix));
for(i=0;i<rows;i++)
{
    matrix[i] = (float *) malloc(rows*sizeof(**matrix));
    for(j=0;j<rows;j++)
        matrix[i][j] = test_data[i][j];
}

/* Everyone computes the determinant (to avoid message transmission) */
determ=determinant(matrix,rows,used_rows,used_cols,0);

if(me==tasks-1)
{
    /* The last task acts as coordinator */
    inverse = (float**) malloc(rows*sizeof(*inverse));
    for(i=0;i<rows;i++)
    {
        inverse[i] = (float *) malloc(rows*sizeof(**inverse));
    }
    /* Print the determinant */
    printf("The determinant of\n\n");
    print_matrix(stdout,matrix,rows,rows);
    printf("\nis %f\n",determ);
    /* Collect the rows of the inverse matrix from the other tasks */
    /* First, post a receive from each task into the appropriate row */
    for(i=0;i<rows;i++)
    {
        MPI_Irecv(inverse[i],rows,MPI_REAL,i,tag,MPI_COMM_WORLD,&(req[i]));
    }
    /* Then wait for all the receives to complete */
    MPI_Waitall(rows,req,status);
    printf("The inverse is\n\n");
    print_matrix(stdout,inverse,rows,rows);
}
else
{
    /* All the other tasks compute a row of the inverse matrix */
    int dest = tasks-1;
    float *one_row;
    int size = rows*sizeof(*one_row);

    one_row = (float*) malloc(size);
    for(j=0;j<rows;j++)
    {
        one_row[j] = coefficient(matrix,rows,j,me)/determ;
    }
    /* Send the row back to the coordinator */
    MPI_Send(one_row,rows,MPI_REAL,dest,tag,MPI_COMM_WORLD);
}
}
/* Wait for all parallel tasks to get here, then quit */

```

```

MPI_Barrier(MPI_COMM_WORLD);
MPI_Finalize();
}
exit(0);

```

Functional decomposition

Parallel servers and data mining applications are examples of functional decomposition. With functional decomposition, the function that the application is performing is distributed among the tasks. Each task operates on the same data, but does something different. The sine series algorithm is also an example of functional decomposition. With this algorithm, the work being done by each task is trivial. The cost of distributing data to the parallel tasks could outweigh the value of running the program in parallel, and parallelism would increase total time. Another approach to parallelism is to invoke different functions, each of which processes all of the data simultaneously. This is possible as long as the final or intermediate results of any function are not required by another function. For example, searching a matrix for the largest and smallest values as well as a specific value could be done in parallel.

This is a simple example, but suppose the elements of the matrix were arrays of polynomial coefficients. Further, suppose the search involved actually evaluating different polynomial equations using the same coefficients. In this case, it would make sense to evaluate each equation separately.

On a simpler scale, let us look at the series for the sine function:

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots - \frac{x^{2n+1}}{(2n+1)!}$$

The serial approach to solving this problem is to loop through the number of terms desired, accumulating the factorial value and the sine value. When the appropriate number of terms has been computed, the loop exits. The following example does exactly this. In this example, you have an array of values for which you want the sine, and an outer loop would repeat this process for each element of the array. Since you do not want to recompute the factorial each time, you need to allocate an array to hold the factorial values and compute them outside the main loop.

```

/*****
 *
 * Series Evaluation - serial version
 *
 * To compile:
 * cc -o series_serial series_serial.c -lm
 *
 *****/

#include<stdlib.h>
#include<stdio.h>
#include<math.h>

double angle[] = { 0.0, 0.1*M_PI, 0.2*M_PI, 0.3*M_PI, 0.4*M_PI,
                  0.5*M_PI, 0.6*M_PI, 0.7*M_PI, 0.8*M_PI, 0.9*M_PI, M_PI };

#define TERMS 8

int main(int argc, char **argv)
{
    double divisor[TERMS], sine;

```

```

int a, t, angles = sizeof(angle)/sizeof(angle[0]);

/* Initialize denominators of series terms */
divisor[0] = 1;
for(t=1;t<TERMS;t++)
{
    divisor[t] = -2*t*(2*t+1)*divisor[t-1];
}

/* Compute sine of each angle */
for(a=0;a<angles;a++)
{
    sine = 0;
    /* Sum the terms of the series */
    for(t=0;t<TERMS;t++)
    {
        sine += pow(angle[a],(2*t+1))/divisor[t];
    }
    printf("sin(%lf) + %lf\n",angle[a],sine);
}
}

```

In a parallel environment, you could assign each term to one task and just accumulate the results on a separate node. In fact, that is what the following example does.

```

/*****
*
* Series Evaluation - parallel version
*
* To compile:
* mpcc -g -o series_parallel series_parallel.c -lm
*
*****/

#include<stdlib.h>
#include<stdio.h>
#include<math.h>
#include<mpi.h>

double angle[] = { 0.0, 0.1*M_PI, 0.2*M_PI, 0.3*M_PI, 0.4*M_PI,
                  0.5*M_PI, 0.6*M_PI, 0.7*M_PI, 0.8*M_PI, 0.9*M_PI, M_PI };

int main(int argc, char **argv)
{
    double data, divisor, partial, sine;
    int a, t, angles = sizeof(angle)/sizeof(angle[0]);
    int me, tasks, term;

    MPI_Init(&argc,&argv); /* Initialize MPI */
    MPI_Comm_size(MPI_COMM_WORLD,&tasks); /* How many parallel tasks are there? */
    MPI_Comm_rank(MPI_COMM_WORLD,&me); /* Who am I? */

    term = 2*me+1; /* Each task computes a term */
    /* Scan the factorial terms through the group members */
    /* Each member will effectively multiply the product of */
    /* the result of all previous members by its factorial */
    /* term, resulting in the factorial up to that point */
    if(me==0)
        data = 1.0;
    else
        data = -(term-1)*term;
    MPI_Scan(&data,&divisor,1,MPI_DOUBLE,MPI_PROD,MPI_COMM_WORLD);

    /* Compute sine of each angle */
    for(a=0;a<angles;a++)
    {

```

```

    partial = pow(angle[a],term)/divisor;
    /* Pass all the partials back to task 0 and */
    /* accumulate them with the MPI_SUM operation */
    MPI_Reduce(&partial,&sine,1,MPI_DOUBLE,MPI_SUM,0,MPI_COMM_WORLD);
    /* The first task has the total value */
    if(me==0)
    {
        printf("sin(%lf) + %lf\n",angle[a],sine);
    }
}
MPI_Finalize();
}

```

With this approach, each task i uses its position in the **MPI_COMM_WORLD** communicator group to compute the value of one term. It first computes its working value as $2i+1$ and calculates the factorial of this value. Since $(2i+1)!$ is $(2i-1)! \times 2i \times (2i+1)$, if each task could get the factorial value computed by the previous task, all it would have to do is multiply it by $2i \times (2i+1)$. Fortunately, MPI provides the capability to do this with the **MPI_SCAN** function. When **MPI_SCAN** is invoked on the first task in a communication group, the result is the input data to **MPI_SCAN**. When **MPI_SCAN** is invoked on subsequent members of the group, the result is obtained by invoking a function on the result of the previous member of the group and its input data.

The MPI standard is documented in *MPI: A Message-Passing Interface Standard, Version 1.1* and is extended in *MPI: A Message-Passing Interface Standard, Version 2.0*, both of which are available from the University of Tennessee. The standard does not specify how to implement the scan function, so a particular implementation does not have to obtain the result from one task and pass it on to the next for processing. This is, however, a convenient way of visualizing the scan function, and the remainder of the discussion will assume that this is happening.

In the example, the function invoked is the built-in multiplication function, **MPI_PROD**. Task 0 (which is computing 1!) sets its result to 1. Task 2 is computing 3! which it obtains by multiplying 2 x 3 by 1! (the result of Task 0). Task 3 multiplies 3! (the result of Task 2) by 4 to get 4!. This continues until all the tasks have computed their factorial values. The input data to the **MPI_SCAN** calls is made negative so the signs of the divisors will alternate between plus and minus.

Once the divisor for a term has been computed, the loop through all the angles (θ) can be done. The partial term is computed as:

$$\pm \frac{\theta^n}{n!}$$

Then, **MPI_REDUCE** is called which is similar to **MPI_SCAN** except that instead of calling a function on each task, the tasks send their raw data to Task 0, which invokes the function on all data values. The function being invoked in the example is **MPI_SUM** which just adds the data values from all of the tasks. Then, Task 0 prints out the result.

Duplication versus redundancy

In the matrix inversion program, each task goes through the process of allocating the matrix and copying the initialization data into it. So why does not one task do this and send the result to all the other tasks? This example has a trivial

initialization process, but in a situation where initialization requires complex time-consuming calculations, this question is even more important.

To understand the answer to this question and, more importantly, be able to apply the understanding to answering the question for other applications, you need to stop and consider the application as a whole. If one task of a parallel application takes on the role of initializer, two things happen. First, all of the other tasks must wait for the initializer to complete (assuming that no work can be done until initialization is completed). Second, some sort of communication must occur to get the results of initialization distributed to all the other tasks. This not only means that there is nothing for the other tasks to do while one task is doing the initializing, there is also a cost associated with sending the results out. Although replicating the initialization process on each of the parallel tasks seems like unnecessary duplication, it allows the tasks to start processing more quickly because they do not have to wait to receive the data.

So, should all initialization be done in parallel? Not necessarily. If the initialization is just computation and setup based on input parameters, each parallel task can initialize independently. Although this seems counter-intuitive at first, because the effort is redundant, for the reasons given above, it is the right answer. Eventually you will get used to it. However, if initialization requires access to system resources that are shared by all the parallel tasks (such as file systems and networks), having each task attempt to obtain the resources will create contention in the system and hinder the initialization process. In this case, it makes sense for one task to access the system resources on behalf of the entire application. In fact, if multiple system resources are required, you could have multiple tasks access each of the resources in parallel. Once the data has been obtained from the resource, you need to decide whether to share the raw data among the tasks and have each task process it, or have one task perform the initialization processing and distribute the results to all the other tasks. You can base this decision on whether the amount of data increases or decreases during the initialization processing. Of course, you want to transmit the smaller amount.

Duplicating the same work on all the remote tasks (which is not the same as redundancy, which implies something can be eliminated) is not bad if:

- The work is inherently serial
- The work is parallel, but the cost of computation is less than the cost of communication
- The work must be completed before tasks can proceed
- Communication can be avoided by having each task perform the same work.

Protocols supported

IMPORTANT

Although this book includes information about using PE with the User Space protocol, User Space is not supported on PE for Linux at this time.

To perform data communication, PE interfaces with a low-level communication API (LAPI), which is a reliable transport provided with PE. LAPI interfaces with a lower level protocol, running in the user space (User Space protocol), which offers a low-latency and high-bandwidth communication path to user applications, running over a high performance switch. LAPI alternatively interfaces with the IP layer.

For optimal performance, PE uses the User Space (US) protocol as its communication path. However, PE also lets you run parallel applications that use the IP interface of LAPI.

The User Space interface allows user applications to take full advantage of the high speed interconnect, and you should use it whenever communication is a critical issue (for instance, when running a parallel application in a production environment). With LoadLeveler, you can use the User Space interface by more than one process per node at a given time.

Both the IP and User Space interfaces allow multiple tasks per job on a single node. As a result, you can use both interfaces in development or test environments, where more attention is paid to the correctness of the parallel program than to its speed-up, and therefore, more users can work on the same nodes at a given time. In both cases, data exchange always occurs between processes, without involving the POE Partition Manager daemon.

Shared memory message passing

For MPI programs in which multiple tasks run on the same computing node, using shared memory to send messages between tasks may be beneficial. This applies to programs running over either the IP or User Space protocol.

By setting the **MP_SHARED_MEMORY** environment variable to *YES*, you can select the shared memory protocol. If *all* the tasks of your program run on the same node, and you specify the shared memory protocol, shared memory is used exclusively for all MPI communications.

For more information on PE's shared memory support, see *IBM Parallel Environment: Operation and Use*.

To thread or not to thread - protocol implications

If you are unfamiliar with POSIX threads, do not try to learn both threads and MPI all at once. Get some experience writing and debugging single process multi-threaded programs first, then tackle multi-process multi-threaded programs.

While each threaded task has more than one independent instruction stream, all of a task's threads share the same address space, file system, and environment variables. In addition, all the threads in a threaded MPI task have the same MPI communicators, data types, ranks, and so on.

A parallel program using MPI normally depends on task parallelism with two or more tasks (or processes) that communicate by message passing. Each of these tasks, by default, has one user thread. An application may explicitly create additional threads within each task, resulting in thread level as well as task level parallelism. If thread creation is done, the application must manage both levels of parallelism properly.

In each threaded MPI task, the **MPI_INIT** routine must be called before any thread can make an MPI call, and all MPI calls must be completed before **MPI_FINALIZE** is called. The principal difference between a threaded task and a non-threaded task is that, in each threaded task, more than one blocking call may be in progress at any given time.

The underlying communication subsystem provides thread-dispatching, so that all blocking messages are given a chance to run when a message completes.

The MPI library creates the following service threads:

- A thread that periodically wakes up and calls the message passing dispatcher, and handles interrupts generated by arriving packets.
- Responder threads used to implement non-blocking collective communication calls and MPI I/O.

The service threads above are terminated when **MPI_FINALIZE** is called. These threads are not available to end users.

Thread debugging implications

To effectively debug a multi-threaded application, you must be aware of how the threads are dispatched. For more information about debugging a multi-threaded program, see the documentation for the GNU debugger (GDB).

Chapter 3. Diagnosing and correcting common problems

What do you do when something goes wrong with your parallel program? PE provides ways to identify and correct problems that arise when you are developing or executing your parallel program. This all depends on where in the process the problem occurred and what the symptoms are.

This section is probably more useful if you use it in conjunction with *IBM Parallel Environment: Operation and Use*. So, you might want to go find it, and keep it on hand for reference.

Here are the steps, greatly abbreviated, in the basic process of creating a parallel program:

1. Create and compile program
2. Start PE
3. Execute the program
4. Verify the output

The remainder of this section describes some of the common problems you might run into, and what to do when they occur. The sections in this section are labeled according to the *symptom* you might be experiencing.

Messages

Messages are an important part of diagnosing problems, so it is essential that you have access to them and that they are at the correct level.

Message catalog errors

You may get message catalog errors. This usually means that the message catalog could not be located or loaded. Check that your **NLSPATH** environment variable includes the path where the message catalog is located. The environment variable **NLSPATH** is used by the various PE components to find the appropriate message catalogs. If the message catalogs are not in the proper place, or your environment variables are not set properly, your system administrator can help. Refer your system administrator to “National language support (NLS)” on page x for more information.

The following are the PE message catalogs:

- pepoe.cat
- pempl.cat
- rsct.lapi.cat

Finding PE messages

There are a number of places that you can find PE messages:

- They are displayed on the home node when it is running POE (STDERR and STDOUT).
- If you set either the **MP_PMDLOG** environment variable or the **-pmdlog** command line option to yes, they are collected in the pmd log file of each task, in **/tmp** (STDERR and STDOUT).

You can also use LookAt to look up message explanations. For more information on how to do this see “Using LookAt to look up message explanations” on page ix

Logging POE errors to a file

You can also specify that diagnostic messages be logged to a file in **/tmp** on each of the remote nodes of your partition by using the **MP_PMDLOG** environment variable. The log file is called **/tmp/mplog.jobid.taskid**, where *jobid* is a unique identifier and *taskid* is the task number. The *jobid* is the same for all remote nodes. This file contains additional diagnostic information about why the user connection was not made. If the file is not there, then pmd did not start. Check the **/etc/xinetd.d/pmv4** and **/etc/services** entries and the executability of pmd for the root user ID again.

For more information about the **MP_PMDLOG** environment variable, see *IBM Parallel Environment: Operation and Use*.

Message format

Knowing which component a message is associated is helpful when trying to resolve a problem. PE messages include prefixes that identify the related component. The message identifiers for the PE components are as follows.

0031-nnnn

Parallel Operating Environment

0032-nnnn

Message Passing Interface

2660-nnnn

LAPI

where:

- The first four digits (such as 0031), identify the component that issued the message.
- *nnnn* identifies the sequence of the message in the group.

For more information about PE messages, see *IBM Parallel Environment: Messages*.

Note that you might find it helpful to run POE as you use this section.

Diagnosing problems using IVP

The *Installation Verification Program* (IVP) can be a useful tool for diagnosing problems. When you installed POE, you verified that everything turned out correctly by running the IVP. It verified that the:

- Location of the libraries was correct
- Binaries existed
- Partition Manager daemon was executable
- POE files were in order
- Sample IVP programs compiled correctly.

The IVP can provide some important first clues when you experience a problem, so you may want to rerun this program before you do anything else.

Cannot compile a parallel program

Programs for PE must be compiled with the current release of the compiler scripts you are using, such as **mpcc**, **mpCC**, **mpfort**. If the command you are trying to use cannot be found, make sure the installation was successful and that your *PATH* environment variable contains the path to the compiler scripts. These commands call the Fortran, C, and C++ compilers respectively, so you also need to make sure

that the underlying compiler is installed and accessible. Your system administrator should be able to assist you in verifying these things.

Cannot start a parallel job

Once you have successfully compiled your program, you either invoke it directly or start POE and then submit the program to it. In both cases, POE is started to establish communication with the parallel nodes. Problems that can occur at this point include: POE does not start, or cannot connect to the remote nodes.

These problems can be caused by other problems on the home node (where you are trying to submit the job), on the remote parallel nodes, or in the communication subsystem that connects them. You need to make sure that all the things POE expects to be set up really are set up. Here is what you do:

1. Make sure that you can execute POE. If you are a Korn or Bash shell user, type:

```
$ whence poe
```

If you are a C shell user, type:

```
$ which poe
```

If the result is just the shell prompt, you do not have POE in your path. It might mean that POE is not installed, or that your path does not point to it. Check that the file **/opt/ibmhpc/ppe.poe/bin/poe** exists and is executable, and that your PATH includes the directory **/opt/ibmhpc/ppe.poe/bin**.

2. Type:

```
$ env | grep MP_
```

Look at the settings of the environment variables beginning with **MP_**, (the POE environment variables). Check their values against what you expect, particularly **MP_HOSTFILE** (where the list of remote host names is to be found), **MP_RESO** (whether a job management system is to be used to allocate remote hosts) and **MP_RMPOOL** (the pool from which the job management system is to allocate remote hosts) values. If they are all not set, make sure that you have a file named **host.list** in your current directory. This file must include the names of all the remote parallel hosts that can be used. There must be at least as many hosts available as the number of parallel processes you specified with the **MP_PROCS** environment variable.

3. Type:

```
$ poe -procs 1
```

You should get the following message:

```
0031-503 Enter program name and flags for each node: _
```

If you do get this message, POE has successfully loaded and established communication with the first remote host in your host list file. It has also validated your use of that remote host, and is ready to go to work. If you type a command, for example, **date**, **hostname**, or **env**, you should get a response when the command executes on the remote host (like you would from **rsh**).

If you get some other set of messages, then the message text should give you some idea of where to look. Some common situations include:

- Cannot connect with the remote host

The path to the remote host is unavailable. Check to make sure that you are trying to connect to the host you think you are. If you are using LoadLeveler to allocate nodes from a pool, you may want to allocate nodes from a known list instead. **ping** the remote hosts in the list to see if a path can be established to them. If it can, run **rsh remote_host date** to verify that the remote host can be contacted and recognizes the host from which you submitted the job, so it can send results back to you.

Check the **/etc/services** file on your home node, to make sure that the Parallel Environment service is defined. Check the **/etc/services** and **/etc/xinetd.d/pmv4** files on the remote host to make sure that the PE service is defined, and that the Partition Manager Daemon (**pmd**) program invoked by **xinetd** on the remote node is executable.

For more information on configuring **rsh** and **inetd**, see *IBM Parallel Environment: Installation*.

- User not authorized on remote host

You need an ID on the remote host and your ID on the home host (the one from which you are submitting the job) must be authorized to run commands on the remote hosts. You do this by placing a **\$HOME/.rhosts** file on the remote hosts that identify your home host and ID. Brush up on “Access” on page 2 if you need to. Even if you have a **\$HOME/.rhosts** file, make sure that you are not denied access the **/etc/hosts.equiv** file on the remote hosts.

In some installations, your home directory is a mounted file system on both your home node and the remote host. In this case, check with your system administrator.

Even if the remote host is actually the same machine as your home node, you still need an entry in the **.rhosts** file.

- Other strangeness

On the home node, you can set or increase the **MP_INFOLEVEL** environment variable (or use the **-infolevel** command line option) to get more information out of POE while it is running. Although this does not give you any more information about the error, or prevent it, it gives you an idea of where POE was, and what it was trying to do when the error occurred. A value of 6 gives you more information than you could ever want. See Appendix A, “A sample program to illustrate messages,” on page 53 for an example of the output from this setting.

Cannot execute a parallel program

Once POE can be started, you need to consider the problems that can arise in running a parallel program, specifically initializing the message passing subsystem. The way to eliminate this initialization as the source of POE startup problems is to run a program that does not use message passing.

As discussed in “Running POE” on page 3, you can use POE to invoke a command or serial program on remote nodes. If you can get a command or simple program, like *Hello, World!*, to run under POE, but a parallel program does not, you can be pretty sure the problem is in the message passing subsystem. The message passing subsystem is the underlying implementation of the message passing calls used by a parallel program (in other words, an **MPI_SEND**). POE code that is linked into your executable by the compiler script (**mpcc**, **mpCC**, **mpfort**) initializes the message passing subsystem.

The Parallel Operating Environment (POE) supports two distinct communication subsystems, an IP-based system, and User Space optimized adapter support. The

subsystem choice is normally made at run time, by environment variables or command line options passed to POE. Use the IP subsystem for diagnosing initialization problems before worrying about the User Space (US) subsystem. Select the IP subsystem by setting the environment variable:

```
$ export MP_EUILIB=ip
```

Use specific remote hosts in your host list file and do not use LoadLeveler (set **MP_RESD=no**). If you do not have a small parallel program, compile the following sample program, **hello_parallel_world**.

Here is the **hello_parallel_world** program in C:

```

/*****
 *
 * Hello Parallel World C Example
 *
 * To compile:
 * mpcc -o hello_parallel_world_c hello_parallel_world.c
 *
 *****/
#include<stdlib.h>
#include<stdio.h>
#include<mpi.h>
/* Basic program to demonstrate compilation and execution techniques */
int main()
{
  MPI_Init(0,0);
  printf("Hello, Parallel World!\n");
  MPI_Finalize();
  exit(0);
}

```

You compile it in C like this:

```
$ mpcc -o hello_parallel_world_c hello_parallel_world.c
```

And here is the **hello_parallel_world** program in Fortran:

```

C*****
C*
C* Hello World Fortran Example
C*
C* To compile:
C* mpfort -o hello_parallel_world_f hello_parallel_world.f
C*
C*****
C -----
C Basic program to demonstrate compilation and execution techniques
C -----
C   program hello

implicit none
include mpif.h
INTEGER error
MPI_INIT(error)
write(6,*)'Hello, Parallel World!'
MPI_FINALIZE(error)

stop
end

```

You compile it in Fortran like this:

```
$ mpfort -o hello_parallel_world_f hello_parallel_world.f
```

Make sure that the executable can be loaded on the remote host that you are using.

Type the following command, and then look at the messages on the console. For C, type the command like this:

```
$ poe hello_parallel_world_c -procs 1 -infolevel 4
```

For Fortran, type the command like this:

```
$ poe hello_parallel_world_f -procs 1 -infolevel 4
```

If you get

```
Hello, Parallel World!
```

then the communication subsystem has been successfully initialized on the one node and things should be looking good. Just for kicks, make sure that there are two remote nodes in your host list file and try again with the following command. If you are using C, type the command like this:

```
$ poe hello_parallel_world_c -procs 2
```

If you are using Fortran, type the command like this:

```
$ poe hello_parallel_world_f -procs 2
```

If and when **hello_parallel_fortran** works with IP and device en0 (the Ethernet), try again with the high speed interconnect.

Each node has one name that it is known by on the external LAN to which it is connected, and another name that it is known by on the interconnect. If the node name you use is not the proper name for the network device you specify, the connection is not made. You can put the names in your host list file. Otherwise, use LoadLeveler to locate the nodes.

The following example assumes you are using C,

```
$ export MP_RESD=yes
$ export MP_EUILIB=ip
$ export MP_EUIDEVICE=sn_single
$ poe hello_parallel_world_c -procs 2 -ilevel 2
```

where **sn_single** is the switch device name. Look at the console lines containing the string **MPI euidevice**. These identify the device name that is actually being used for message passing (as opposed to the IP address that is used to connect the home node to the remote hosts.) If these are not device names, check the LoadLeveler configuration and the switch configuration.

Once IP works, and you are on a clustered server, you can try message passing using the User Space device support. Note that LoadLeveler allows you to run multiple tasks over the switch adapter while in User Space.

You can run **hello_parallel_world** with the User Space library by typing the following. This example assumes you are using C.

```
e$ export MP_RESD=yes
$ export MP_EUILIB=us
$ export MP_EUIDEVICE=sn_single
$ poe hello_parallel_world_c -procs 2 -ilevel 6
```

The console log should inform you that you are using User Space support, and that LoadLeveler is allocating the nodes for you. LoadLeveler tells you that it cannot

allocate the requested nodes if someone else is already running on them **and** has requested dedicated use of the switch, or if User Space capacity has been exceeded.

You can try for other specific nodes, or you can ask LoadLeveler for nonspecific nodes from a pool. You can refer to *IBM Parallel Environment: Operation and Use*.

The program runs but...

Occasionally, you may encounter problems running your program. This includes program errors resulting in abnormal termination and a core dump, lack of output, or program hangs. This section provides some basic suggestions for how to deal with those kinds of situations under Parallel Environment.

When a core dump is created

If your program creates a core dump (segmentation fault), POE saves a copy of the core file so you can debug it later. Unless you specify otherwise, POE saves the core file in the **coredir.taskid** directory under the current working directory, where *taskid* is the task number. For example, if your current working directory is **/u/mickey**, and your application creates a core dump while running on the node that is task 4, the core file will be located in **/u/mickey/coredir.4** on that node.

You can control where POE saves the core file by using the **-coredir** POE command line option or the **MP_COREDIR** environment variable.

No output at all

Should there be output?

If you are not getting output from your program and you think you ought to be, make sure you have enabled the program to send data back to you. If the **MP_STDOUTMODE** environment variable is set to a number, it is the number of the only task for which standard output will be displayed. If that task does not generate standard output, you will not see any.

There should be output

If **MP_STDOUTMODE** is set appropriately, the next step is to verify that the program is actually doing something. Start by observing how the program terminates (or fails to terminate). It will do one of the following things:

- Terminate without generating output other than POE messages.
- Fail to terminate after a **really** long time, still without generating output.

In the first case, you should examine any messages you receive. Since your program is not generating any output, all of the messages will be coming from POE.

In the second case, you will have to stop the program yourself (<Ctrl-c> should work).

One possible reason for lack of output could be that your program is terminating abnormally before it can generate any. POE will report abnormal termination conditions such as being killed, as well as non-zero return codes. Sometimes these messages are obscured in the blur of other errata, so it is important to check the messages carefully.

Figuring out return codes: It is important to understand POE's interpretation of return codes. If the exit code for a task is zero(0) or in the range of 2 to 127, then

POE will make that task wait until all tasks have exited. If the exit code is 1 or greater than 128 (or less than 0), then POE will terminate the entire parallel job abruptly (with a **SIGTERM** signal to each task). In normal program execution, one would expect to have each program go through **exit(0)** or **STOP**, and exit with an exit code of 0. However, if a task encounters an error condition (for example, a full file system), then it may exit unexpectedly. In these cases, the exit code is usually set to -1. If, however, you have written error handlers which produce exit codes other than 1 or -1, then POE's termination algorithm may cause your program to *hang* because one task has terminated abnormally, while the other tasks continue processing (expecting the terminated task to participate).

If the POE messages indicate the job was killed (either because of some external situation like low page space or because of POE's interpretation of the return codes), it may be enough information to fix the problem. Otherwise, you may have to do more analysis.

The program hangs

If you have gotten this far and the POE messages, and the additional checking by the message passing routines, have not shed any light on why your program is not generating output, the next step is to figure out whether your program is doing anything at all (besides not giving you output).

Let's look at the following example...it has a bug in it.

```

/*****
*
* Ray trace program with bug
*
* To compile:
* mpcc -g -o rtrace_bug rtrace_bug.c
*
* Description:
* This is a sample program that partitions N tasks into
* two groups, a collect node and N - 1 compute nodes.
* The responsibility of the collect node is to collect the data
* generated by the compute nodes. The compute nodes send the
* results of their work to the collect node for collection.
*
* There is a bug in this code. Please do not fix it in this file!
*
*****/

#include <mpi.h>

#define PIXEL_WIDTH 50
#define PIXEL_HEIGHT 50

int First_Line = 0;
int Last_Line = 0;

void main(int argc, char *argv[])
{
    int numtask;
    int taskid;

    /* Find out number of tasks/nodes. */
    MPI_Init( &argc, &argv);
    MPI_Comm_size( MPI_COMM_WORLD, &numtask);
    MPI_Comm_rank( MPI_COMM_WORLD, &taskid);

    /* Task 0 is the coordinator and collects the processed pixels */
    /* All the other tasks process the pixels */

```

```

if ( taskid == 0 )
    collect_pixels(taskid, numtask);
else
    compute_pixels(taskid, numtask);

printf("Task %d waiting to complete.\n", taskid);
/* Wait for everybody to complete */
MPI_Barrier(MPI_COMM_WORLD);
printf("Task %d complete.\n", taskid);
MPI_Finalize();
exit();
}

/* In a real implementation, this routine would process the pixel */
/* in some manner and send back the processed pixel along with its*/
/* location. Since you did process the pixel. all you do is      */
/* send back the location                                         */
compute_pixels(int taskid, int numtask)
{
    int section;
    int row, col;
    int pixel_data[2];
    MPI_Status stat;

    printf("Compute #%d: checking in\n", taskid);

    section = PIXEL_HEIGHT / (numtask - 1);

    First_Line = (taskid - 1) * section;
    Last_Line = taskid * section;

    for (row = First_Line; row < Last_Line; row ++)
        for ( col = 0; col < PIXEL_WIDTH; col ++)
            {
                pixel_data[0] = row;
                pixel_data[1] = col;
                MPI_Send(pixel_data, 2, MPI_INT, 0, 0, MPI_COMM_WORLD);
            }
    printf("Compute #%d: done sending. ", taskid);
    return;
}

/* This routine collects the pixels. In a real implementation, */
/* after receiving the pixel data, the routine would look at the*/
/* location information that came back with the pixel and move */
/* the pixel into the appropriate place in the working buffer */
/* Since you aren't doing anything with the pixel data, you don't */
/* bother and each message overwrites the previous one          */
collect_pixels(int taskid, int numtask)
{
    int pixel_data[2];
    MPI_Status stat;
    int mx = PIXEL_HEIGHT * PIXEL_WIDTH;

    printf("Control #%d: No. of nodes used is %d\n", taskid, numtask);
    printf("Control: expect to receive %d messages\n", mx);

    while (mx > 0)
        {
            MPI_Recv(pixel_data, 2, MPI_INT, MPI_ANY_SOURCE,
                    MPI_ANY_TAG, MPI_COMM_WORLD, &stat);
            mx--;
        }
    printf("Control node #%d: done receiving. ", taskid);
    return;
}

```

This example is from a ray tracing program that distributed a display buffer out to server nodes. The intent is that each task, other than Task 0, takes an equal number of full rows of the display buffer, processes the pixels in those rows, and then sends the updated pixel values back to the client. In the real application, the task would compute the new pixel value and send it as well, but in this example, you are just sending the row and column of the pixel. Because the client is getting the row and column location of each pixel in the message, it does not care which server each pixel comes from. The client is Task 0, and the servers are all the other tasks in the parallel job.

This example has a functional bug in it. With a little bit of analysis, the bug is probably easy to spot, and you may be tempted to fix it right away. PLEASE DO NOT!

When you run this program, you get the output shown below. Notice that the **-g** option is used when you compile the example. You are cheating a little because you know that there is going to be a problem, so you are compiling with debug information that is turned on right away.

```
$ mpcc -g -o rtrace_bug rtrace_bug.c
$ rtrace_bug -procs 4 -labelio yes
1:Compute #1: checking in
0:Control #0: No. of nodes used is 4
1:Compute #1: done sending. Task 1 waiting to complete.
2:Compute #2: checking in
3:Compute #3: checking in
0:Control: expect to receive 2500 messages
2:Compute #2: done sending. Task 2 waiting to complete.
3:Compute #3: done sending. Task 3 waiting to complete.
^C
ERROR: 0031-250 task 1: Interrupt
ERROR: 0031-250 task 2: Interrupt
ERROR: 0031-250 task 3: Interrupt
ERROR: 0031-250 task 0: Interrupt
```

No matter how long you wait, the program will not terminate until you press **<Ctrl-c>**.

So, you suspect the program is hanging somewhere. You know it starts executing because you get some messages from it. It could be a logical hang or it could be a communication hang.

Hangs and threaded programs

Coordinating the threads in a task requires careful locking and signaling. Deadlocks that occur because the program is waiting on locks that have not been released are common, in addition to the deadlock possibilities that arise from improper use of the MPI message passing calls.

If you have concluded that your program is hanging, attach your debugger to diagnose the problem.

Why did the program hang?

The symptom of the problem in the **rtrace_bug** program was a hang. Hangs can occur for the same reasons they occur in serial programs (in other words, loops without exit conditions). They may also occur because of message passing deadlocks or because of some subtle differences between the parallel and sequential environments.

Using the debugger to analyze sometimes indicates that the source of a hang is a message that was never received, even though it is a valid one, and even though it appears to have been sent. In these situations, the problem is probably due to lost messages in the communication subsystem. This is especially true if the lost message is intermittent or varies from run to run. This is either the program's fault or the environment's fault. Before investigating the environment, you should analyze the program's *safety* with respect to MPI. A *safe* MPI program is one that does not depend on a particular implementation of MPI. You should also examine the error logs for evidence of repeated message transmissions (which usually indicate a network failure).

Although MPI specifies many details about the interface and behavior of communication calls, it also leaves many implementation details unspecified (and it does not just omit them, it specifies that they are unspecified.) This means that certain uses of MPI may work correctly in one implementation and fail in another, particularly in the area of how messages are buffered. An application may even work with one set of data and fail with another in the same implementation of MPI. This is because, when the program works, it has stayed within the limits of the implementation. When it fails, it has exceeded the limits. Because the limits are unspecified by MPI, both implementations are valid. MPI *safety* is discussed further in Chapter 4, "Creating a safe program," on page 49.

Once you have verified that the application is *MPI-safe*, your only recourse is to blame lost messages on the environment. If the communication path is IP, use the standard network analysis tools to diagnose the problem. Look particularly at **mbuf** usage. You can examine **mbuf** usage with the **netstat** command. Note that the **netstat** command is not a distributed command, which means that it applies only to the node on which you execute it.

```
$ netstat -m
```

If the **mbuf** line shows any failed allocations, you should increase the **thewall** value of your network options. You can see your current setting with the **no** command. Note that the **no** command is not a distributed command which means that it applies only to the node on which you execute it.

```
$ no -a
```

The value presented for **thewall** is in KBytes. You can use the **no** command to change this value. You will have to have root access to do this. For example,

```
$ no -o thewall=16384
```

sets **thewall** to 16 MBytes.

Message passing between lots of remote hosts can tax the underlying IP system. Make sure that you look at all the remote nodes, not just the home node. Allow lots of buffers. If the communication path is user space (US), you will need to get your system support people involved to isolate the problem.

Other reasons for the program to hang

One final cause for no output is a problem on the home node (POE is hung). Normally, a *hang* is associated with the remote hosts waiting for each other, or for a termination signal. POE running on the home node is alive and well, waiting patiently for some action on the remote hosts. If you type **<Ctrl-c>** on the POE console, you will be able to successfully interrupt and terminate the set of remote hosts. See *IBM Parallel Environment: Operation and Use* for information on the **poekill** command.

There are situations where POE itself can hang. Usually these situations are associated with large volumes of input or output. Remember that POE normally gets standard output from each node. If each task writes a large amount of data to standard output, it may chew up the IP buffers on the machine running POE, causing it (and all the other processes on that machine) to block and hang. The only way to know that this is the problem is by seeing that the rest of the home node has hung. If you think that POE is hung on the home node, your only solution may be to kill POE there. Press <Ctrl-c> several times, or use the command **kill -9**. At present, there are only partial approaches to avoiding the problem. You can allocate lots of **mbufs** on the home node, and do not make the send and receive buffers too large.

Bad output

Bad output includes unexpected error messages. After all, who expects error messages or bad results (results that are not correct).

Error messages

You can track down the causes of error messages and correct them in parallel programs using techniques similar to those used for serial programs. One difference, however, is that you need to identify which task is producing the message, if it is not coming from all tasks. You can do this by setting the **MP_LABELIO** environment variable to **yes**, or using the **-labelio yes** command line parameter. Generally, the message will give you enough information to identify the location of the problem.

You may also want to generate *more* error and warning messages by setting the **MP_EUIDEVELOP** environment variable to **yes** when you first start running a new parallel application. This will give you more information about the things that the message passing library considers errors or unsafe practices.

Bad results

You can track down bad results and correct them in a parallel program in a fashion similar to that used for serial programs. The process in the previous debugging exercise can be more complicated because the processing and control flow on one task may be affected by other tasks. In a serial program, you can follow the exact sequence of instructions that were executed and observe the values of all variables that affect the control flow. However, in a parallel program, both the control flow and the data processing on a task may be affected by messages sent from other tasks. For one thing, you may not have been watching those other tasks. For another, the messages could have been sent a long time ago. Therefore, it is very difficult to correlate a message that you receive with a particular series of events.

Debugging and threads

So far, the discussion has been about debugging normal old serial or parallel programs, but you may want to debug a threaded program (or be aware of the threads used in the library). If this is the case, there are a few things you should consider.

Before you do anything else, you first need to understand the environment in which you are working. You have the potential to create a multi-threaded application, using a multi-threaded library, that consists of multiple distributed tasks. As a result, finding and diagnosing bugs in this environment may require a different set of debugging techniques that you are not used to using. Here are some things to remember.

When you attach to a running program, all the tasks you selected in your program will be stopped at their current points of execution. Typically, you want to see the current point of execution of your task. This stop point is the position of the program counter, and may be in any one of the many threads that your program may create OR any one of the threads that the MPI library creates. With non-threaded programs, it was adequate to just travel up the program stack until you reached your application code (assuming you compiled your program with the **-g** option). But with threaded programs, you now need to traverse across other threads to get to your thread(s) and then up the program stack to view the current point of execution of your code.

The MPI library itself will create a set of threads to process message requests. When you attach to a program that uses the MPI library, all of the threads associated with the POE job are stopped, including the ones created and used by MPI.

For more information on the threaded MPI library, see *IBM Parallel Environment: MPI Programming Guide*.

Tuning the performance of a parallel application

So far, the discussions have been about getting PE working, creating message passing parallel programs, debugging problems, and debugging parallel applications. When you get a parallel program running so that it gives us the correct answer, you are done. Not necessarily. In this area, parallel programs are just like sequential programs; just because they give you the correct answer does not mean they are doing it in the most efficient manner. For a program that is relatively short running or is run infrequently, it may not matter how efficient it is. For a program that consumes a significant portion of the system resources, you need to make the best use of those resources by tuning its performance.

The basic approach here is to tune a sequential program and then parallelize it. With this approach, the process is the same as for any sequential program, and you use one of the same tools; **gprof**. In this case, the parallelization process must take performance into account, and should avoid anything that adversely affects it.

With this approach, you use a standard sequential tool in the traditional manner. When you tune an application and then parallelize it, observe the communication performance, how it affects the performance of each of the individual tasks, and how the tasks affect each other. For example, does one task spend a lot of time waiting for messages from another? If so, perhaps you need to rebalance the workload. Or if a task starts waiting for a message long before it arrives, perhaps it could do more algorithmic processing before waiting for the message. When an application is made parallel and then tuned, you need a way to collect the performance data in a manner that includes both communication and algorithmic information. That way, if the performance of a task needs to be improved, you can decide between tuning the algorithm or tuning the communication.

Tuning the performance of threaded programs

There are some things you need to consider when you want to get the maximum performance out of the program.

Note: The PE implementation of MPI (PE MPI) is thread safe.

- Two environment variables affect the overhead of an MPI call in the threaded library:
 - **MP_SINGLE_THREAD=[nolyes]**
 - **MP_EUIDEVELOP=[nolyes|deblmin]**

A program that has only one MPI communication thread may set the environment variable **MP_SINGLE_THREAD=yes** before calling **MPI_INIT**. This will avoid some locking which is otherwise required to maintain consistent internal MPI state. The program may have other threads that do computation or other work, as long as they do not make MPI calls. Note that the implementation of MPI I/O and MPI one-sided communication is thread-based, and that these facilities may not be used when **MP_SINGLE_THREAD** is set to **yes**.

The **MP_EUIDEVELOP** environment variable lets you control how much checking is done when you run the program. Eliminating checking altogether (setting **MP_EUIDEVELOP** to **min**) provides performance (latency) benefits, but may cause critical information to be unavailable if the executable hangs due to message passing errors. For more information on **MP_EUIDEVELOP** and other POE environment variables, see *IBM Parallel Environment: Operation and Use*.

- Programs (threaded or non-threaded) that use the threaded MPI library can be profiled by using the **-pg** flag on the compilation and linking step of the program. The profile results (**gmon.out**) will contain only a summary of the information from all the threads per task together. Viewing the data using **prof** or **gprof** is limited to showing only this summarized data on a per task basis, not per thread. For more information on profiling programs, see “Profile it.”

Profile it

The first step in tuning a program is to find the areas within the program that execute most of the work. Locating these compute-intensive areas within the program lets you focus on the areas that give you the most benefit from tuning. The best way to find them is to *profile* the program.

To profile your program, compile it with the **-pg** flag to generate profiling data. The **-pg** flag compiles and links the executable so that when you run the program, the performance data gets written to output. After you have compiled your program with the **-pg** flag, run it again to see what you get. This generates a file called **gmon.out** for each task, which can then be analyzed with the GNU **gprof** command.

By default, POE stores the **gmon.out** files in a subdirectory called **profdir.task_id** in your current working directory. Note that you can change the name of this directory by using the **MP_PROFDIR** environment variable or the **-profdir** command line flag of the **poe** command.

Chapter 4. Creating a safe program

Going from serial to parallel programming means that you are on a different scale now. This section alerts you to some of the things you need to pay attention to as you create your parallel programs. In particular, the section provides information on creating a *safe* (mostly harmless) MPI program. *MPI: A Message-Passing Interface Standard, Version 1.1* which is available from the University of Tennessee (<http://www.mpi-forum.org/>) provides additional information. You may want to refer to that document.

What is a safe program?

Many people consider a program to be *safe* if message buffering is not required for the program to complete. In a program like this, you should be able to replace all standard sends with synchronous sends, and the program will still run correctly. This type of programming style is conservative; it provides good portability because program completion does not depend on the amount of available buffer space.

With PE, setting the **MP_EAGER_LIMIT** environment variable to **0** is equivalent to making all sends synchronous, including those used in collective communication. A good test of your program's safety is to set the **MP_EAGER_LIMIT** to **0**.

Some programmers prefer more flexibility and use an *unsafe* style that relies on buffering. In such cases, the use of standard send operations provides a compromise between performance and robustness. MPI attempts to supply sufficient buffering so that these programs will not result in deadlock. You can use the buffered send mode for programs that require more buffering, or in situations where you want more control. Since buffer overflow conditions are easier to diagnose than deadlock, you can also use this mode for debugging purposes.

You can use nonblocking message passing operations to avoid the need for buffering outgoing messages. This prevents deadlock situations due to a lack of buffer space, and improves performance by allowing computation and communication to overlap. It also avoids the overhead associated with allocating buffers and copying messages into buffers.

Safety and threaded programs

Sometimes message passing programs can hang or deadlock. This can occur when one task waits for a message that is never sent or when each task is waiting for the other task to send or receive a message. Within a task, a similar situation can occur when one thread is waiting for another thread to release a lock on a shared resource, such as a piece of memory. If thread *A*, which holds the lock, cannot run to the point at which it is ready to release it, the waiting thread *B* will never run. This may occur because thread *B* holds some other lock that thread *A* needs. Thread *A* cannot proceed until thread *B* does, and thread *B* cannot proceed until thread *A* does.

When programs are both multi-thread and multi-task, there is risk of *deadly embrace* involving both mutex and communication blocks. Say threads *A* and *B* are on task 0, and thread *A* holds a lock while waiting for a message from task 1. Thread *B* will send a message to task 1 only after it gets the lock that thread *A* holds. If task 1 will send the message that thread *A* is waiting for only after getting

the one that thread *B* cannot send, the job is in a 3-way deadly embrace, involving two threads at task 0 and one thread at task 1.

A problem that is more subtle occurs when two threads simultaneously access a shared resource without a lock protocol. The result may be incorrect without any obvious sign. For example, the following function is not thread-safe, because the thread may be preempted after the variable *c* is updated, but before it is stored.

```
int c; /* external, used by two threads */
void update_it()
{
    c++; /* this is not thread safe */
}
```

You probably should avoid writing threaded message passing programs until you are familiar with writing and debugging threaded, single-task programs.

Using threaded programs with non-thread-safe libraries

A threaded MPI program must meet the same criteria as any other threaded program; it must avoid using non-thread-safe functions in more than one thread (for example, **strtok**). In addition, it must use only thread-safe libraries, if library functions are called on more than one thread. All of the libraries may not be thread-safe, so you should carefully examine how they are used in your program.

Message ordering

With MPI, messages are *non-overtaking*. This means that the order of sends must match the order of receives. Assume a sender sends two messages (Message 1 and Message 2) in succession, to the same destination, and both match the same receive. The receive operation will receive Message 1 before Message 2. Likewise, if a receiver posts two receives (Receive 1 and Receive 2), in succession, and both are looking for the same message, Receive 1 will receive the message before Receive 2. Adhering to this rule ensures that sends are always matched with receives.

If a process in your program has a single thread of execution, then the sends and receives that occur follow a natural order. However, if a process has multiple threads, the various threads may not execute their relative send operations in any defined order. In this case, the messages can be received in any order.

Order rules apply within each communicator. Weakly synchronized threads can each use independent communicators to avoid many order problems.

The following is an example of using non-overtaking messages. The message sent by the first send must be received by the first receive, and the message sent by the second send must be received by the second receive.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_BSEND(buf1, count, MPI_REAL, 1, tag, comm, ierr)
    CALL MPI_BSEND(buf2, count, MPI_REAL, 1, tag, comm, ierr)
ELSE ! rank.EQ.1
    CALL MPI_RECV(buf1, count, MPI_REAL, 0, MPI_ANY_TAG, comm, status, ierr)
    CALL MPI_RECV(buf2, count, MPI_REAL, 0, tag, comm, status, ierr)
END IF
```

Program progress when two processes initiate two matching sends and receives

If two processes (or *tasks*) initiate two matching sends and receives, at least one of the operations (the send or the receive) will complete, regardless of other actions that occur in the system. The send operation will complete unless its matching receive operation has already been satisfied by another message, and has itself completed. Likewise, the receive operation will complete unless its matching send message is claimed by another matching receive that was posted at the same destination.

The following example shows two matching pairs that are intertwined in this manner. Here is what happens:

1. Both processes invoke their first calls.
2. *process 0*'s first send indicates buffered mode, which means it must complete, even if there is no matching receive. Since the first receive posted by *process 1* does not match, the send message gets copied into buffer space.
3. Next, *process 0* posts its second send operation, which matches *process 1*'s first receive, and both operations complete.
4. *process 1* then posts its second receive, which matches the buffered message, so both complete.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
  CALL MPI_BSEND(buf1, count, MPI_REAL, 1, tag1, comm, ierr)
  CALL MPI_SSEND(buf2, count, MPI_REAL, 1, tag2, comm, ierr)
ELSE ! rank.EQ.1
  CALL MPI_RECV(buf1, count, MPI_REAL, 0, tag2, comm, status, ierr)
  CALL MPI_RECV(buf2, count, MPI_REAL, 0, tag1, comm, status, ierr)
END IF
```

Communication fairness

MPI does not guarantee *fairness* in the way communications are handled. It is your responsibility to prevent starvation among the operations in your program.

One example of an *unfair* situation might be where a send, with a matching receive on another process, does not complete because another message, from a different process, overtakes the receive.

Resource limitations

If a lack of resources prevents an MPI call from executing, errors may result. Pending send and receive operations consume a portion of your system resources. MPI attempts to use a minimal amount of resource for each pending send and receive, but buffer space is required for storing messages sent in either standard or buffered mode when no matching receive is available.

When a buffered send operation cannot complete due to a lack of buffer space, the resulting error could cause your program to terminate abnormally. On the other hand, a standard send operation that cannot complete because of a lack of buffer space, will block and wait for buffer space to become available or for the matching receive to be posted. In some situations, this behavior is preferable because it avoids the error condition associated with buffer overflow.

Sometimes a lack of buffer space can lead to deadlock. The program in the following example will succeed even if no buffer space for data is available.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 1, tag, comm, ierr)
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 1, tag, comm, status, ierr)
ELSE ! rank.EQ.1
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 0, tag, comm, status, ierr)
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 0, tag, comm, ierr)
END IF
```

In this next example, neither process will send until the other process sends first. As a result, this program will always result in deadlock.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 1, tag, comm, status, ierr)
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 1, tag, comm, ierr)
ELSE ! rank.EQ.1
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 0, tag, comm, status, ierr)
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 0, tag, comm, ierr)
END IF
```

The example below shows how message exchange relies on buffer space. The message send by each process must be copied out before the send returns and the receive starts. Consequently, at least one of the two messages sent needs to be buffered for the program to complete. As a result, this program can execute successfully only if the communication system can buffer at least the words of data specified by *count*.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 1, tag, comm, ierr)
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 1, tag, comm, status, ierr)
ELSE ! rank.EQ.1
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 0, tag, comm, ierr)
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 0, tag, comm, status, ierr)
END IF
```

When standard send operations are used, deadlock can occur where both processes are blocked because buffer space is not available. This is also true for synchronous send operations. For buffered sends, if the required amount of buffer space is not available, the program will not complete either, and instead of deadlock, you will have buffer overflow.

Appendix A. A sample program to illustrate messages

This appendix provides sample output for a program run under POE with the maximum level of message reporting. It also points out the different types of messages you can expect, and explains what they mean.

To set the level of messages that get reported when you run your program, you can use the **-infolevel** (or **-ilevel**) option when you invoke POE. You can also use the **MP_INFOLEVEL** environment variable. Setting either of these to 6 gives you the maximum number of diagnostic messages when you run your program. For more information about setting the POE message level, see *IBM Parallel Environment: Operation and Use*.

Note that we are using numbered prefixes along the left-hand edge of the following output as a way to refer to particular lines. The prefixes are **not** part of the output you will see when you run your program. For an explanation of the messages denoted by these numbered prefixes, see “Figuring out what all of this means” on page 57.

This command produces output similar to the following:

```
> poe ./hello_world -procs 2 -rmpool 1 -infolevel 6
1 INFO: DEBUG_LEVEL changed from 0 to 4
2 D1<L4>: Open of file ./host.list successful
3 ATTENTION: 0031-379 Pool setting ignored when hostfile used
4 D1<L4>: mp_euilib = ip
5 D1<L4>: 02/28 10:47:09.464301 task 0 c171f9sq01.ppd.pok.ibm.com 10
6 D1<L4>: 02/28 10:47:09.464456 task 1 c171f9sq01.ppd.pok.ibm.com 10
7 D1<L4>: node allocation strategy = 2
8 INFO: 0031-364 Contacting LoadLeveler to set and query information for
  interactive job
9 D1<L4>: 02/28 10:47:09.484486 Calling ll_init_job.
10
11 D1<L4>: 02/28 10:47:09.626426 ll_init_job returned.
12
13 D1<L4>: Affinity is not requested; MP_TASK_AFFINITY: -1
14 D1<L4>: 02/28 10:47:09.626559 Job Command String:
15 #@ job_type = parallel
16 #@ environment = COPY_ALL
17 #@ node_usage = shared
18 #@ bulkxfer = NO
19 #@ class = No_Class
20 #@ queue
21
22 INFO: 0031-380 LoadLeveler step ID is c171f9sq01.ppd.pok.ibm.com.204.0
23 INFO: 0031-118 Host c171f9sq01.ppd.pok.ibm.com requested for task 0
24 INFO: 0031-118 Host c171f9sq01.ppd.pok.ibm.com requested for task 1
25 D4<L4>: 02/28 10:47:09.699382 Calling ll_event()...
26 D4<L4>: 02/28 10:47:09.774382 Return code from ll_event() was 0
27 D4<L4>: 02/28 10:47:09.774445 State returned for LL_StepState was 3
28 D4<L4>: 02/28 10:47:09.901215 LL_StepBulkXfer returned by LL was 0
29 D4<L4>: 02/28 10:47:09.901274 LL_StepNodeCount returned by LL was 1
30 D4<L4>: 02/28 10:47:09.901309 LL_NodeTaskCount returned by LL was 1
31 D4<L4>: 02/28 10:47:09.901336 LL_TaskTaskInstanceCount returned by LL was 2
32 INFO: 0031-119 Host c171f9sq01.ppd.pok.ibm.com allocated for task 0
33 INFO: 0031-120 Host address 9.114.136.84 allocated for task 0
34 D4<L4>: 02/28 10:47:09.902179 LL_TaskInstanceAdapterCount returned by
  LL was 1
35 D4<L4>: 02/28 10:47:09.902239 Device returned...
36 INFO: 0031-377 Using eth0 for mpi euidevice for task 0
37 D4<L4>: 02/28 10:47:09.902280 AdapterUsageRcxtBlocks returned value: 0
38 D4<L4>: 02/28 10:47:09.902305 Adapter info for task 0: proto=mpi, mode=ip,
```

```

address=9.114.136.84, name=eth0, tag=0
39 D4<L4>: 02/28 10:47:09.902344 MP_RCXT_BLKs set to: 0
40 INFO: 0031-119 Host c171f9sq01.ppd.pok.ibm.com allocated for task 1
41 INFO: 0031-120 Host address 9.114.136.84 allocated for task 1
42 D4<L4>: 02/28 10:47:09.902403 LL_TaskInstanceAdapterCount returned by LL was 1
43 D4<L4>: 02/28 10:47:09.902432 Device returned...
44 INFO: 0031-377 Using eth0 for mpi euidevice for task 1
45 D4<L4>: 02/28 10:47:09.902469 AdapterUsageRcxtBlocks returned value: 0
46 D4<L4>: 02/28 10:47:09.902494 Adapter info for task 1: proto=mpi, mode=ip,
address=9.114.136.84, name=eth0, tag=0
47 D4<L4>: 02/28 10:47:09.902523 MP_RCXT_BLKs set to: 0
48 D4<L4>: 02/28 10:47:09.902557 Return for get_job_info()
49 D1<L4>: Entering pm_contact, jobid is 0
50 D1<L4>: Jobid = 1141162211
51 D1<L4>: Spawning /etc/pmdv4 on all nodes
52 D1<L4>: 1 master nodes
53 D4<L4>: LoadLeveler Version 3 Release 3
54 D1<L4>: 02/28 10:47:09.903578 Calling ll_spawn_connect for node 0, host name
c171f9sq01.ppd.pok.ibm.com
55 D1<L4>: 02/28 10:47:09.903854 ll_spawn_connect returned for node 0, socket fd 5,
host name c171f9sq01.ppd.pok.ibm.com
56 D4<L4>: ll_spawn_connect successful with c171f9sq01.ppd.pok.ibm.com (task 0)
57 D4<L4>: Calling pm_spawn_ready (number of nodes 1)
58 D1<L4>: 02/28 10:47:09.903928 Calling pm_spawn_ready.
59
60 D4<L4>: Return from ll_spawn_write (socket file descriptor 5), rc = 0
61 D4<L4>: Return from ll_spawn_read (socket file descriptor 5), rc = 0
62 D4<L4>: Return from ll_spawn_read (socket file descriptor 5), rc = 1
63 D1<L4>: 02/28 10:47:09.905043 returned from pm_spawn_ready.
64
65 D1<L4>: Socket file descriptor for master 0 (c171f9sq01.ppd.pok.ibm.com) is 5
66 D1<L4>: SSM_read on socket 5, source = 0, task id: 0, nread: 12, type:3.
67 D1<L4>: Leaving pm_contact, jobid is 1141162211
68 D4<L4>: Command args:<>
69 D3<L4>: Message type 34 from source 0
70 D4<L4>: Task 0 pulse received, count is 0 curr_time is 1141141629
71 D4<L4>: Task 0 pulse acknowledged, count is 0 curr_time is 1141141629
72 D3<L4>: Message type 21 from source 0
73 INFO: 0031-724 Executing program: <./hw_mpi>
74 D3<L4>: Message type 21 from source 1
75 D3<L4>: Message type 21 from source 1
76 D3<L4>: Message type 21 from source 1
77 INFO: 0031-724 Executing program: <./hw_mpi>
78 D3<L4>: Message type 21 from source 0
79 INFO: DEBUG_LEVEL changed from 0 to 4
80 D3<L4>: Message type 21 from source 1
81 INFO: DEBUG_LEVEL changed from 0 to 4
82 D3<L4>: Message type 21 from source 1
83 D1<L4>: After bzero, threadsig.sa_handler = 0
84
85 D3<L4>: Message type 21 from source 1
86 D4<L4>: pm_async_thread sends cond sig
87 D3<L4>: Message type 21 from source 1
88 D4<L4>: pm_async_thread calls sigwait, in_async_thread=0
89 D3<L4>: Message type 21 from source 1
90 D1<L4>: Before sigwait: threadsig.sa_handler = 0
91 D3<L4>: Message type 21 from source 0
92 D3<L4>: Message type 21 from source 0
93 D1<L4>: After bzero, threadsig.sa_handler = 0
94 D3<L4>: Message type 21 from source 0
95
96 D3<L4>: Message type 21 from source 0
97 D3<L4>: Message type 21 from source 0
98 D3<L4>: Message type 21 from source 0
99 D4<L4>: pm_async_thread sends cond sig
100 D3<L4>: Message type 21 from source 0
101 D3<L4>: Message type 21 from source 0

```

```

102 D3<L4>: Message type 21 from source 0
103 D4<L4>: pm_async_thread calls sigwait, in_async_thread=0
104 D3<L4>: Message type 21 from source 0
105 D3<L4>: Message type 21 from source 0
106 D3<L4>: Message type 21 from source 0
107 D1<L4>: Before sigwait: threadsig.sa_handler = 0
108 D3<L4>: Message type 21 from source 1
109 D4<L4>: pm_main, wake up from timed cond wait
110 D3<L4>: Message type 21 from source 1
111 D1<L4>: mp_euilib is <ip>
112 D3<L4>: Message type 21 from source 1
113 D3<L4>: Message type 21 from source 1
114 D1<L4>: Executing _mp_init_msg_passing() from MPI_Init()...
115 D3<L4>: Message type 21 from source 0
116 D4<L4>: pm_main, wake up from timed cond wait
117 D3<L4>: Message type 21 from source 0
118 D1<L4>: mp_euilib is <ip>
119 D3<L4>: Message type 21 from source 1
120 D3<L4>: Message type 21 from source 0
121 D3<L4>: Message type 21 from source 1
122 D1<L4>: mp_css_interrupt is <0>
123 D3<L4>: Message type 21 from source 0
124 D1<L4>: Executing _mp_init_msg_passing() from MPI_Init()...
125 D3<L4>: Message type 21 from source 1
126 D3<L4>: Message type 21 from source 1
127 D1<L4>: About to call mpci_connect
128 D3<L4>: Message type 21 from source 0
129 D3<L4>: Message type 21 from source 0
130 D1<L4>: mp_css_interrupt is <0>
131 D3<L4>: Message type 21 from source 0
132 D3<L4>: Message type 21 from source 0
133 D1<L4>: About to call mpci_connect
134 D3<L4>: Message type 21 from source 0
135 D3<L4>: Message type 21 from source 0
136 INFO: 0031-619 32bit(ip) ppe_rsanlx MPCIE shared object was compiled on
Tue Feb 14 14:43:11 2006
137 D3<L4>: Message type 21 from source 0
138
139 D3<L4>: Message type 21 from source 1
140 2660-501 LAPI version #8.22 2006/01/24 1.168 src/rsct/lapi/lapi.c, lapi,
rsct_ros, ros0607a 32bit(ip) library compiled on Wed Feb 22 09:43:08 2006
141
142 D3<L4>: Message type 21 from source 1
143 LAPI is using lightweight lock.
144 D3<L4>: Message type 21 from source 0
145 2660-501 LAPI version #8.22 2006/01/24 1.168 src/rsct/lapi/lapi.c, lapi,
rsct_ros, ros0607a 32bit(ip) library compiled on Wed Feb 22 09:43:08 2006
146
147 D3<L4>: Message type 21 from source 0
148 LAPI is using lightweight lock.
149 D3<L4>: Message type 23 from source 0
150 D1<L4>: init_data for instance number 0, task 0: <158500948:32816>
151 D3<L4>: Message type 23 from source 1
152 D1<L4>: init_data for instance number 0, task 1: <158500948:32815>
153 D3<L4>: Message type 21 from source 0
154 LAPI_AM_LW_XFER is supported in LAPI library
155 D3<L4>: Message type 21 from source 1
156 LAPI_AM_LW_XFER is supported in LAPI library
157 D3<L4>: Message type 21 from source 0
158 The MPI shared memory protocol is used for the job
159 D3<L4>: Message type 21 from source 1
160 The MPI shared memory protocol is used for the job
161 D3<L4>: Message type 21 from source 0
162 D3<L4>: Message type 21 from source 1
163 D1<L4>: Elapsed time for mpci_connect: 0 seconds
164 D3<L4>: Message type 21 from source 0
165 D1<L4>: Elapsed time for mpci_connect: 0 seconds

```

```

166 D3<L4>: Message type 21 from source 0
167 D3<L4>: Message type 21 from source 1
168 D3<L4>: Message type 21 from source 0
169 D1<L4>: _css_init: rc from HPS0clk_init is 0
170
171 D3<L4>: Message type 21 from source 1
172 D1<L4>: _css_init: rc from HPS0clk_init is 0
173
174 D3<L4>: Message type 21 from source 0
175 D1<L4>: About to call _ccl_init
176 D3<L4>: Message type 21 from source 1
177 D1<L4>: About to call _ccl_init
178 D3<L4>: Message type 21 from source 0
179 D3<L4>: Message type 21 from source 1
180 D3<L4>: Message type 21 from source 0
181 D1<L4>: Elapsed time for _ccl_init: 0 seconds
182 D3<L4>: Message type 21 from source 1
183 D1<L4>: Elapsed time for _ccl_init: 0 seconds
184 D3<L4>: Message type 21 from source 0
185 D1<L4>: LAPI_AM_LW_XFER limit is 104 bytes
186
187 D3<L4>: Message type 21 from source 1
188 D1<L4>: LAPI_AM_LW_XFER limit is 104 bytes
189
190 D3<L4>: Message type 20 from source 0
191 Hello Penguin
192 D3<L4>: Message type 20 from source 1
193 Hello Penguin
194 D3<L4>: Message type 62 from source 0
195 D3<L4>: Message type 62 from source 1
196 D3<L4>: Message type 20 from source 0
197 Exiting...
198 D3<L4>: Message type 21 from source 0
199 D3<L4>: Message type 21 from source 0
200 INFO: 0031-306 pm_atexit: pm_exit_value is 0.
201 D3<L4>: Message type 21 from source 0
202 D1<L4>: Sending SSM_EXIT_REQ
203 D3<L4>: Message type 17 from source 0
204 D3<L4>: Message type 20 from source 1
205 Exiting...
206 D3<L4>: Message type 21 from source 1
207 D3<L4>: Message type 21 from source 1
208 D3<L4>: Message type 21 from source 1
209 INFO: 0031-306 pm_atexit: pm_exit_value is 0.
210 D3<L4>: Message type 21 from source 1
211 D3<L4>: Message type 21 from source 1
212 D1<L4>: Sending SSM_EXIT_REQ
213 D3<L4>: Message type 17 from source 1
214 D3<L4>: Message type 22 from source 0
215 INFO: 0031-656 I/O file STDOUT closed by task 0
216 D3<L4>: Message type 22 from source 0
217 INFO: 0031-656 I/O file STDERR closed by task 0
218 D3<L4>: Message type 15 from source 0
219 D1<L4>: Accounting data from task 0 for source 0:
220 D3<L4>: Message type 22 from source 1
221 INFO: 0031-656 I/O file STDOUT closed by task 1
222 D3<L4>: Message type 22 from source 1
223 INFO: 0031-656 I/O file STDERR closed by task 1
224 D3<L4>: Message type 15 from source 1
225 D1<L4>: Accounting data from task 1 for source 1:
226 D3<L4>: Message type 1 from source 0
227 INFO: 0031-251 task 0 exited: rc=0
228 D3<L4>: Message type 1 from source 1
229 INFO: 0031-251 task 1 exited: rc=0
230 D1<L4>: All remote tasks have exited: maxx_errcode = 0
231 INFO: 0031-639 Exit status from pm_respond = 0
232 D1<L4>: Maximum return code from user = 0

```

```

233 D2<L4>: In pm_exit... About to call pm_remote_shutdown
234 D2<L4>: Sending PMD_EXIT to task 0
235 D2<L4>: Elapsed time for pm_remote_shutdown: 0 seconds
236 D2<L4>: Return code from ll_close: 0
237 D2<L4>: In pm_exit... Calling exit with status = 0 at Tue Feb 28
10:47:15 2006

```

Figuring out what all of this means

When you set **-infolevel** to 6, you get the full complement of diagnostic messages, which we will explain here.

The previous example includes numbered prefixes along the left-hand edge of the output so that we can refer to particular lines, and then explain what they mean. Remember, that these prefixes are **not** part of your output. The table below points you to the line number of the messages that are of most interest, and provides a short explanation.

Lines	Message description
5-6	Names hosts that are used.
65	Indicates node with partition manager running.
69	Message type 34 indicates pulse activity (the pulse mechanism checked that each remote node was actively participating with the home node).
72	Message type 21 indicates a STDERR message.
111, 118	Indicates that the euilib message passing protocol was specified.
114, 124	Indicates message passing initialization has begun.
136	Timestamp of MPCPI shared object being executed.
140, 145	Timestamp of LAPI library being executed.
158, 160	Indicates MPI shared memory is being used.
163, 165	Indicates initialization of MPCPI has completed.
190, 191, 192, 193	Message type 20 shows STDOUT from your program..
200, 209	Indicates that the user's program has reached the exit handler. The exit code is 0.
203, 213	Message type 17 indicates the tasks have requested to exit.
215, 217, 221, 223	Indicates that the STDOUT and STDERR pipes have been closed.
218, 224	Message type 15 indicates accounting data.
234	Indicates that the home node is sending an exit.

Appendix B. Parallel Environment internals

This appendix provides some additional information about how the IBM Parallel Environment (PE) works with respect to your application. Much of this information is also explained in the *IBM Parallel Environment: MPI Programming Guide*.

What happens when I compile my applications?

In order to run your program in parallel, you first need to compile your application source code with one of the following scripts:

1. **mpcc**
2. **mpCC**
3. **mpfort**

To make sure the parallel execution works, these scripts add the following to your application executable:

- POE initialization module, so POE can determine that all nodes can communicate successfully, before giving control to the user application's main() routine.
- Signal handlers, for additional control in terminating the program during parallel tracing, and enabling the handling of the process termination signals. The *IBM Parallel Environment: MPI Programming Guide* explains the signals that are handled in this manner.

The compile scripts dynamically link the Message Passing library interfaces in such a way that the specific communication library that is used is determined when your application executes.

Applications created as static executables are not supported.

How do my applications start?

Because POE adds its entry point to each application executable, user applications do not need to be run under the **poe** command. When a parallel application is invoked directly, as opposed to under the control of the **poe** command, POE is started automatically. It then sets up the parallel execution environment and then re-invokes the application on each of the remote nodes.

Serial applications can be run in parallel only using the **poe** command. However, such applications cannot take advantage of the function and performance provided with the message passing libraries.

How does POE talk to the nodes?

A parallel job running under POE consists of a *home node* (where POE was started) and *n* tasks. Each task runs under the control of a Partition Manager daemon (pmd). There is one pmd for each job on each node on which the job's tasks run.

When you start a parallel job, POE contacts the nodes assigned to run the job (called *remote nodes*), and starts a pmd instance on each node. POE sends environment information to the pmd daemons for the parallel job (including the name of the executable) and the pmd daemons spawn processes to run the

executable. For tasks that run on the same node, the pmd daemon forks and manages all tasks for that job on that node. It routes messages to and from each remote task, and also coordinates with the home node to terminate each task.

The spawned processes have standard I/O redirected to socket connections back to the pmd daemons. Therefore, any output the application writes to STDOUT or STDERR is sent back to the pmd daemons. The pmd daemons, in turn, send the output back to POE via another socket connection, and POE writes the output to its STDOUT or STDERR. Any input that POE receives on STDIN is delivered to the remote tasks in a similar fashion.

The socket connections between POE and the pmd daemons are also used to exchange control messages for providing task synchronization, exit status, and signaling. These capabilities are available to control any parallel program run by POE, and they do not depend on the message passing library.

When POE executes without LoadLeveler[®], it is assumed that the Partition Manager Daemon (PMD) is started under **xinetd**. There is no consideration for running the PMD without **xinetd**.

When POE executes under LoadLeveler (including all User Space applications), the PMD is started by LoadLeveler.

How are signals handled?

POE installs signal handlers for most signals that cause program termination and interrupts, in order to control and notify all tasks of the signal. POE will exit the program normally with a code of (128 + signal). If the user program installs a signal handler for any of the signals POE supports, it should follow the guidelines presented in *IBM Parallel Environment: MPI Programming Guide*.

What happens when my application ends?

POE returns exit status (a return code value between 0 and 255) on the home node which reflects the composite exit status of the user application. The exit status can have various conditions and values and each can have a specific meaning. These are explained in *The IBM Parallel Environment: MPI Programming Guide*.

In addition, if the POE job-step function is used, the job control mechanism is the program's exit code. When the task exit code is 0 (zero), or in the range of 2 to 127, the job-step will be continued. If the task exit code is 1 or greater than 127, POE terminates the parallel job, as well as any remaining user programs in the job-step list. Also, any POE infrastructure failure detected (such as failure to open pipes to the child process) will terminate the parallel job as well as any remaining programs in the job-step list.

Appendix C. Accessibility

Accessibility features help a user who has a physical disability, such as restricted mobility or limited vision, to use software products successfully. The major accessibility features enable users to:

- Use assistive technologies such as screen readers and screen magnifier software
- Operate specific or equivalent features using only the keyboard
- Customize display attributes such as color, contrast, and font size.

Using assistive technologies

Assistive technology products, such as screen readers, function with user interfaces. Consult the assistive technology documentation for specific information when using such products to access interfaces.

Notices

This information was developed for products and services offered in the U.S.A.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing
IBM Corporation
North Castle Drive
Armonk, NY 10504-1785
U.S.A.

For license inquiries regarding double-byte (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

IBM World Trade Asia Corporation
Licensing
2-31 Roppongi 3-chome, Minato-ku
Tokyo 106-0032, Japan

The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law:

INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

This book refers to IBM's implementation of the Message Passing Interface (MPI) standard for Parallel Environment for Linux (PE). PE MPI intends to comply with the requirements of the Message Passing Interface Forum described below. PE MPI provides an implementation of MPI which is complete except for omitting the features described in the "Process Creation and Management" chapter of MPI-2.

Permission to copy without fee all or part of these Message Passing Interface Forum documents:

MPI: A Message Passing Interface Standard, Version 1.1

MPI-2: Extensions to the Message Passing Interface, Version 2.0

is granted, provided the University of Tennessee copyright notice and the title of the document appear, and notice is given that copying is by permission of the University of Tennessee. ©1993, 1997 University of Tennessee, Knoxville, Tennessee.

Any references in this information to non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this IBM product and use of those Web sites is at your own risk.

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Licensees of this program who wish to have information about it for the purpose of enabling: (i) the exchange of information between independently created programs and other programs (including this one) and (ii) the mutual use of the information which has been exchanged, should contact:

IBM Corporation
Department LJEB/P905
522 South Road
Poughkeepsie, NY 12601-5400
U.S.A

Such information may be available, subject to appropriate terms and conditions, including in some cases, payment of a fee.

The licensed program described in this document and all licensed material available for it are provided by IBM under terms of the IBM Customer Agreement, IBM International Program License Agreement or any equivalent agreement between us.

Any performance data contained herein was determined in a controlled environment. Therefore, the results obtained in other operating environments may vary significantly. Some measurements may have been made on development-level systems and there is no guarantee that these measurements will be the same on generally available systems. Furthermore, some measurements may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.

COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrates programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs.

If you are viewing this information softcopy, the photographs and color illustrations may not appear.

Trademarks

The following terms are trademarks of the International Business Machines Corporation in the United States or other countries or both:

- IBM
- IBMLink
- LoadLeveler
- Tivoli
- xSeries

UNIX is a registered trademark in the United States, other countries, or both and is licensed exclusively through X/Open Company Limited.

Linux is a trademark of Linus Torvalds in the United States, other countries, or both.

Other company, product, and service names may be the trademarks or service marks of others.

Glossary

A

AFS[®]. Andrew File System.

address. A value, possibly a character or group of characters that identifies a register, a device, a particular part of storage, or some other data source or destination.

API. Application programming interface.

application. The use to which a data processing system is put; for example, a payroll application, an airline reservation application.

argument. A parameter passed between a calling program and a called program or subprogram.

attribute. A named property of an entity.

Authentication. The process of validating the identity of a user or server.

Authorization. The process of obtaining permission to perform specific actions.

B

bandwidth. For a specific amount of time, the amount of data that can be transmitted. Bandwidth is expressed in bits or bytes per second (bps) for digital devices, and in cycles per second (Hz) for analog devices.

blocking operation. An operation that does not complete until the operation either succeeds or fails. For example, a blocking receive will not return until a message is received or until the channel is closed and no further messages can be received.

breakpoint. A place in a program, specified by a command or a condition, where the system halts execution and gives control to the workstation user or to a specified program.

broadcast operation. A communication operation where one processor sends (or broadcasts) a message to all other processors.

buffer. A portion of storage used to hold input or output data temporarily.

C

C. A general-purpose programming language. It was formalized by Uniforum in 1983 and the ANSI standards committee for the C language in 1984.

C++. A general-purpose programming language that is based on the C language. C++ includes extensions that support an object-oriented programming paradigm. Extensions include:

- strong typing
- data abstraction and encapsulation
- polymorphism through function overloading and templates
- class inheritance.

chaotic relaxation. An iterative relaxation method that uses a combination of the Gauss-Seidel and Jacobi-Seidel methods. The array of discrete values is divided into subregions that can be operated on in parallel. The subregion boundaries are calculated using the Jacobi-Seidel method, while the subregion interiors are calculated using the Gauss-Seidel method. See also *Gauss-Seidel*.

client. A function that requests services from a server and makes them available to the user.

cluster. A group of processors interconnected through a high-speed network that can be used for high-performance computing.

collective communication. A communication operation that involves more than two processes or tasks. Broadcasts, reductions, and the **MPI_Allreduce** subroutine are all examples of collective communication operations. All tasks in a communicator must participate.

communicator. An MPI object that describes the communication context and an associated group of processes.

compile. To translate a source program into an executable program.

condition. One of a set of specified values that a data item can assume.

core dump. A process by which the current state of a program is preserved in a file. Core dumps are usually associated with programs that have encountered an unexpected, system-detected fault, such as a Segmentation Fault or a severe user error. The current program state is needed for the programmer to diagnose and correct the problem.

core file. A file that preserves the state of a program, usually just before a program is terminated for an unexpected error. See also *core dump*.

D

data decomposition. A method of breaking up (or decomposing) a program into smaller parts to exploit

parallelism. One divides the program by dividing the data (usually arrays) into smaller parts and operating on each part independently.

data parallelism. Refers to situations where parallel tasks perform the same computation on different sets of data.

debugger. A debugger provides an environment in which you can manually control the execution of a program. It also provides the ability to display the program's data and operation.

domain name. The hierarchical identification of a host system (in a network), consisting of human-readable labels, separated by decimal points.

E

environment variable. (1) A variable that describes the operating environment of the process. Common environment variables describe the home directory, command search path, and the current time zone. (2) A variable that is included in the current software environment and is therefore available to any called program that requests it.

Ethernet. A baseband local area network (LAN) that allows multiple stations to access the transmission medium at will without prior coordination, avoids contention by using carrier sense and deference, and resolves contention by using collision detection and delayed retransmission. Ethernet uses carrier sense multiple access with collision detection (CSMA/CD).

event. An occurrence of significance to a task — the completion of an asynchronous operation such as an input/output operation, for example.

executable. A program that has been link-edited and therefore can be run in a processor.

execution. To perform the actions specified by a program or a portion of a program.

expression. In programming languages, a language construct for computing a value from one or more operands.

F

fairness. A policy in which tasks, threads, or processes must be allowed eventual access to a resource for which they are competing. For example, if multiple threads are simultaneously seeking a lock, no set of circumstances can cause any thread to wait indefinitely for access to the lock.

Fiber Distributed Data Interface (FDDI). An American National Standards Institute (ANSI) standard for a local area network (LAN) using optical fiber cables. An FDDI LAN can be up to 100 kilometers (62 miles) long, and

can include up to 500 system units. There can be up to 2 kilometers (1.24 miles) between system units and concentrators.

file system. The collection of files and file management structures on a physical or logical mass storage device, such as a diskette or minidisk.

foreign host. See *remote host*.

FORTRAN. One of the oldest of the modern programming languages, and the most popular language for scientific and engineering computations. Its name is a contraction of *FORmula TRANslation*. The two most common FORTRAN versions are FORTRAN 77, originally standardized in 1978, and FORTRAN 90. FORTRAN 77 is a proper subset of FORTRAN 90.

function cycle. A chain of calls in which the first caller is also the last to be called. A function that calls itself recursively is not considered a function cycle.

functional decomposition. A method of dividing the work in a program to exploit parallelism. The program is divided into independent pieces of functionality, which are distributed to independent processors. This method is in contrast to data decomposition, which distributes the same work over different data to independent processors.

functional parallelism. Refers to situations where parallel tasks specialize in particular work.

G

Gauss-Seidel. An iterative relaxation method for solving Laplace's equation. It calculates the general solution by finding particular solutions to a set of discrete points distributed throughout the area in question. The values of the individual points are obtained by averaging the values of nearby points. Gauss-Seidel differs from Jacobi-Seidel in that, for the $i+1$ st iteration, Jacobi-Seidel uses only values calculated in the i th iteration. Gauss-Seidel uses a mixture of values calculated in the i th and $i+1$ st iterations.

GDB. The GNU Project Debugger. GDB is an open-source tool, created and maintained by the GNU Project, that allows you to examine the source code of a program as it executes. GDB is a useful tool for determining why a program crashes and where, in the program, the problem occurs.

global max. The maximum value across all processors for a given variable. It is global in the sense that it is global to the available processors.

global variable. A variable defined in one portion of a computer program and used in at least one other portion of the computer program.

gprof. A UNIX command that produces an execution profile of C, COBOL, FORTRAN, or Pascal programs. The execution profile is in a textual and tabular format. It is useful for identifying which routines use the most CPU time. See the man page on **gprof**.

graphical user interface (GUI). A type of computer interface consisting of a visual metaphor of a real-world scene, often of a desktop. Within that scene are icons, which represent actual objects, that the user can access and manipulate with a pointing device.

GUI. Graphical user interface.

H

home node. The node from which an application developer compiles and runs his program. The home node can be any workstation on the LAN.

host. A computer connected to a network that provides an access method to that network. A host provides end-user services.

host list file. A file that contains a list of host names, and possibly other information, that was defined by the application that reads it.

host name. The name used to uniquely identify any computer on a network.

hot spot. A memory location or synchronization resource for which multiple processors compete excessively. This competition can cause a disproportionately large performance degradation when one processor that seeks the resource blocks, preventing many other processors from having it, thereby forcing them to become idle.

I

IBM Parallel Environment (PE) for Linux. A licensed program that provides an execution and development environment for parallel C, C++, and FORTRAN programs. It also includes tools for debugging, profiling, and tuning parallel programs.

installation image. A file or collection of files that are required in order to install a software product on system nodes. See also *licensed program*, and *package*.

Internet. The collection of worldwide networks and gateways that function as a single, cooperative virtual network.

Internet Protocol (IP). (1) The IP protocol lies beneath the UDP protocol, which provides packet delivery between user processes and the TCP protocol, which provides reliable message delivery between user processes.

IP. Internet Protocol.

J

Jacobi-Seidel. See *Gauss-Seidel*.

K

kernel. The core portion of the UNIX operating system that controls the resources of the CPU and allocates them to the users. The kernel is memory-resident, is said to run in *kernel mode* (in other words, at higher execution priority level than *user mode*), and is protected from user tampering by the hardware.

L

Laplace's equation. A homogeneous partial differential equation used to describe heat transfer, electric fields, and many other applications.

latency. The time interval between the initiation of a send by an origin task and the completion of the matching receive by the target task. More generally, latency is the time between a task initiating data transfer and the time that transfer is recognized as complete at the data destination.

licensed program. A collection of software packages sold as a product that customers pay for to license. A licensed program can consist of packages and file sets a customer would install. These packages and file sets bear a copyright and are offered under the terms and conditions of a licensing agreement. See also *filesset* and *package*.

LoadLeveler. A job management system that works with POE to let users run jobs and match processing needs with system resources, in order to make better use of the system.

local variable. A variable that is defined and used only in one specified portion of a computer program.

loop unrolling. A program transformation that makes multiple copies of the body of a loop, also placing the copies within the body of the loop. The loop trip count and index are adjusted appropriately so the new loop computes the same values as the original. This transformation makes it possible for a compiler to take additional advantage of instruction pipelining, data cache effects, and software pipelining.

See also *optimization*.

M

management domain . A set of nodes configured for manageability by the Clusters Systems Management (CSM) product. Such a domain has a management server that is used to administer a number of managed nodes. Only management servers have knowledge of the whole domain. Managed nodes only know about the

servers managing them; they know nothing of each other. Contrast with *peer domain*.

menu. A list of options displayed to the user by a data processing system, from which the user can select an action to be initiated.

message catalog. A file created from a message source file that contains application error and other messages, which can later be translated into other languages without having to recompile the application source code.

message passing. Refers to the process by which parallel tasks explicitly exchange program data.

Message Passing Interface (MPI). A standardized API for implementing the message-passing model.

MIMD. Multiple instruction stream, multiple data stream.

Multiple instruction stream, multiple data stream (MIMD). A parallel programming model in which different processors perform different instructions on different sets of data.

MPMD. Multiple program, multiple data.

Multiple program, multiple data (MPMD). A parallel programming model in which different, but related, programs are run on different sets of data.

MPI. Message Passing Interface.

N

network. An interconnected group of nodes, lines, and terminals. A network provides the ability to transmit data to and receive data from other systems and users.

Network Information Services. A set of network services (for example, a distributed service for retrieving information about the users, groups, network addresses, and gateways in a network) that resolve naming and addressing differences among computers in a network.

NIS. See *Network Information Services*.

node. (1) In a network, the point where one or more functional units interconnect transmission lines. A computer location defined in a network.

node ID. A string of unique characters that identifies the node on a network.

nonblocking operation. An operation, such as sending or receiving a message, that returns immediately whether or not the operation was completed. For example, a nonblocking receive will not wait until a message arrives. By contrast, a blocking receive will wait. A nonblocking receive must be completed by a later test or wait.

O

object code. The result of translating a computer program to a relocatable, low-level form. Object code contains machine instructions, but symbol names (such as array, scalar, and procedure names), are not yet given a location in memory. Contrast with *source code*.

optimization. A widely-used (though not strictly accurate) term for program performance improvement, especially for performance improvement done by a compiler or other program translation software. An optimizing compiler is one that performs extensive code transformations in order to obtain an executable that runs faster but gives the same answer as the original. Such code transformations, however, can make code debugging and performance analysis very difficult because complex code transformations obscure the correspondence between compiled and original source code.

option flag. Arguments or any other additional information that a user specifies with a program name. Also referred to as *parameters* or *command-line options*.

P

package. A collection of files, usually used to install a piece of software. The equivalent AIX term is fileset. One type of package is an RPM (RPM Package Manager) package.

parallelism. The degree to which parts of a program may be concurrently executed.

parallelize. To convert a serial program for parallel execution.

Parallel Operating Environment (POE). An execution environment that smooths the differences between serial and parallel execution. It lets you submit and manage parallel jobs. It is abbreviated and commonly known as POE.

parameter. (1) In FORTRAN, a symbol that is given a constant value for a specified application. (2) An item in a menu for which the operator specifies a value or for which the system provides a value when the menu is interpreted. (3) A name in a procedure that is used to refer to an argument that is passed to the procedure. (4) A particular piece of information that a system or application program needs to process a request.

partition. (1) A fixed-size division of storage.

Partition Manager. The component of the Parallel Operating Environment (POE) that allocates nodes, sets up the execution environment for remote tasks, and manages distribution or collection of standard input (STDIN), standard output (STDOUT), and standard error (STDERR).

PE. The Parallel Environment for Linux licensed program.

peer domain. A set of nodes configured for high availability by the RSCT configuration manager. Such a domain has no distinguished or master node. All nodes are aware of all other nodes, and administrative commands can be issued from any node in the domain. All nodes also have a consistent view of the domain membership. Contrast with *management domain*.

performance monitor. A utility that displays how effectively a system is being used by programs.

PID. Process identifier.

POE. Parallel Operating Environment.

pool. Groups of nodes on an SP system that are known to LoadLeveler, and are identified by a pool name or number.

point-to-point communication. A communication operation that involves exactly two processes or tasks. One process initiates the communication through a *send* operation. The partner process issues a *receive* operation to accept the data being sent.

procedure. (1) In a programming language, a block, with or without formal parameters, whose execution is invoked by means of a procedure call. (2) A set of related control statements that cause one or more programs to be performed.

process. A program or command that is actually running the computer. It consists of a loaded version of the executable file, its data, its stack, and its kernel data structures that represent the process's state within a multitasking environment. The executable file contains the machine instructions (and any calls to shared objects) that will be executed by the hardware. A process can contain multiple threads of execution.

The process is created with a **fork()** system call and ends using an **exit()** system call. Between **fork** and **exit**, the process is known to the system by a unique process identifier (PID).

Each process has its own virtual memory space and cannot access another process's memory directly. Communication methods across processes include pipes, sockets, shared memory, and message passing.

prof. A utility that produces an execution profile of an application or program. It is useful to identify which routines use the most CPU time. See the man page for **prof**.

profiling. The act of determining how much CPU time is used by each function or subroutine in a program. The histogram or table produced is called the execution profile.

pthread. A thread that conforms to the POSIX Threads Programming Model.

R

reduced instruction-set computer. A computer that uses a small, simplified set of frequently-used instructions for rapid execution.

reduction operation. An operation, usually mathematical, that reduces a collection of data by one or more dimensions. For example, the arithmetic SUM operation is a reduction operation that reduces an array to a scalar value. Other reduction operations include MAXVAL and MINVAL.

Reliable Scalable Cluster Technology. A set of software components that together provide a comprehensive clustering environment for AIX®. RSCT is the infrastructure used by a variety of IBM products to provide clusters with improved system availability, scalability, and ease of use.

remote host. Any host on a network except the one where a particular operator is working.

remote shell (rsh). A command that lets you issue commands on a remote host.

RISC. See *reduced instruction-set computer*.

RPM. RPM Package Manager (RPM) manages Linux software and update packages.

RSCT. See *Reliable Scalable Cluster Technology*.

RSCT peer domain. See *peer domain*.

S

secure shell (ssh). Secure Shell, sometimes known as Secure Socket Shell, is a Unix-based command interface and protocol for securely getting access to a remote computer.

shell script. A sequence of commands that are to be executed by a shell interpreter such as the Bourne shell (**sh**), the C shell (**csh**), or the Korn shell (**ksh**). Script commands are stored in a file in the same format as if they were typed at a terminal.

segmentation fault. A system-detected error, usually caused by referencing a non-valid memory address.

server. A functional unit that provides shared services to workstations over a network — a file server, a print server, or a mail server, for example.

signal handling. In the context of a message passing library (such as MPI), there is a need for asynchronous operations to manage packet flow and data delivery while the application is doing computation. This

asynchronous activity can be carried out either by a signal handler or by a service thread. The early IBM message passing libraries used a signal handler and the more recent libraries use service threads. The older libraries are often referred to as the *signal handling* versions.

Single program, multiple data (SPMD). A parallel programming model in which different processors execute the same program on different sets of data.

source code. The input to a compiler or assembler, written in a source language. Contrast with *object code*.

source line. A line of source code.

SPMD. Single program, multiple data.

standard error (STDERR). An output file intended to be used for error messages for C programs.

standard input (STDIN). The primary source of data entered into a command. Standard input comes from the keyboard unless redirection or piping is used, in which case standard input can be from a file or the output from another command.

standard output (STDOUT). The primary destination of data produced by a command. Standard output goes to the display unless redirection or piping is used, in which case standard output can go to a file or to another command.

STDERR. Standard error.

STDIN. Standard input.

STDOUT. Standard output.

stencil. A pattern of memory references used for averaging. A 4-point stencil in two dimensions for a given array cell, $x(i,j)$, uses the four adjacent cells, $x(i-1,j)$, $x(i+1,j)$, $x(i,j-1)$, and $x(i,j+1)$.

subroutine. (1) A sequence of instructions whose execution is invoked by a call. (2) A sequenced set of instructions or statements that can be used in one or more computer programs and at one or more points in a computer program. (3) A group of instructions that can be part of another routine or can be called by another program or routine.

synchronization. The action of forcing certain points in the execution sequences of two or more asynchronous procedures to coincide in time.

system administrator. (1) The person at a computer installation who designs, controls, and manages the use of the computer system. (2) The person who is responsible for setting up, modifying, and maintaining the Parallel Environment.

T

task. A unit of computation analogous to a process. In a parallel job, there are two or more concurrent tasks working together through message passing. Though it is common to allocate one task per processor, the terms *task* and *processor* are not interchangeable.

thread. A single, separately dispatchable, unit of execution. There can be one or more threads in a process, and each thread is executed by the operating system concurrently.

TPD. Third party debugger.

U

unrolling loops. See *loop unrolling*.

user. (1) A person who requires the services of a computing system. (2) Any person or any thing that can issue or receive commands and message to or from the information processing system.

User Space. A version of the message passing library that is optimized for direct access to the communication adapter. User Space maximizes performance by passing up all kernel involvement in sending or receiving a message.

utility program. A computer program in general support of computer processes; for example, a diagnostic program, a trace program, a sort program.

utility routine. A routine in general support of the processes of a computer; for example, an input routine.

V

variable. (1) In programming languages, a named object that may take different values, one at a time. The values of a variable are usually restricted to one data type. (2) A quantity that can assume any of a given set of values. (3) A name used to represent a data item whose value can be changed while the program is running. (4) A name used to represent data whose value can be changed, while the program is running, by referring to the name of the variable.

X

X Window System. The UNIX industry's graphics windowing standard that provides simultaneous views of several executing programs or processes on high resolution graphics displays.

Index

Special characters

- eulib 9
- hostfile option 8
- ilevel option 8
- infolevel option 8, 10, 38, 53
- labelio option 4, 8
- llfile 9
- pmdlog option 8
- procs option 4, 8
- rmpool 9
- stdoutmode option 5, 9

A

- abbreviated names viii
- access, to nodes 2
- Accessibility 61
- acronyms for product names viii
- allocation, node
 - high performance switch 9
 - host list file 9
- audience of this book vii

B

- bad output 46
 - bad results 46
 - error messages 46
- bad results 46

C

- common problems 35
 - bad output 46
 - cannot compile a parallel program 36
 - cannot connect with the remote host 37
 - cannot execute a parallel program 38
 - cannot start a parallel job 37
 - hangs 42
 - no output 41
 - user not authorized on remote host 38
- compiler scripts 7, 59
 - for threaded programs 7
- compiling 6
 - C example 6
 - examples 6
 - Fortran example 7
 - scripts 7, 59
- conventions viii

D

- data decomposition 21, 22
- debugging
 - threaded programs 46
- duplication 31

E

- efficiency 47
- environment variables 9
 - LANG 35
 - MP_EUIDEVELOP 46
 - MP_EULIB 9
 - MP_HOSTFILE 8, 37
 - MP_INFOLEVEL 8, 38, 53
 - MP_LABELIO 5, 8, 46
 - MP_LLFILE 9
 - MP_PMDLOG 9
 - MP_PROCS 5, 8, 37
 - MP_RESD 37, 39
 - MP_RMPOOL 9, 37
 - MP_STDOUTMODE 9, 41
 - NLSPATH 35
 - running POE with 5
- error messages 46
- errors
 - logging to a file 36

F

- functional decomposition 21, 29

H

- hangs 45
 - threaded programs 44
- high performance switch 39
 - and node allocation 9
- host list file 3, 9
- host list file, examples 3

I

- initialization, how implemented 32
- installation 2
- Installation Verification Program (IVP) 2, 36
- IP protocol 33

L

- LANG 35
- LAPI 32
- LoadLeveler 1, 3, 9, 10, 33, 40
 - and User Space support 40
- logging errors to a file 36
- LookAt message retrieval tool ix
- loops, unrolling 22
 - example 22

M

- message passing 21
 - definition 21

- message passing (*continued*)
 - synchronization 21
- message retrieval tool, LookAt ix
- messages
 - and problem determination 35
 - finding 35
 - format 36
 - interpreted 57
 - level reported 10, 53
 - PE message catalog components 36
 - PE message catalog errors 35
 - types 57
- MP_EUIDEVELOP 46
- MP_EUILIB 9
- MP_HOSTFILE 8, 9
- MP_INFOLEVEL 8, 38, 53
- MP_LABELIO 5, 8, 46
- MP_LLFILE 9
- MP_PMDLOG 9
- MP_PROCS 5, 8
- MP_RESD 39
- MP_RMPOOL 9
- MP_STDOUTMODE 41
- MP_STOUTMODE 9
- MPI_COMM_WORLD 31
- MPI_Comm_rank 23
- MPI_Comm_size 23
- MPI_Finalize 23
- MPI_Init 23
- MPI_PROD 31
- MPI_REDUCE 31
- MPI_Scan 31
- MPI_SCAN 31
- MPI_SUM 31
- myhosts file 9

N

- NLSPATH 35
- node allocation
 - high performance switch 9
 - host list file 9

O

- options
 - eulib 9
 - hostfile 8
 - ilevel 8
 - infolevel 8, 10, 53
 - labelio 8
 - llfile 9
 - pmdlog 8
 - procs 8
 - rmpool 9
 - stdoutmode 9

P

- Parallel Operating Environment
 - hostfile option 8

- Parallel Operating Environment (*continued*)
 - ilevel option 8
 - infolevel option 8, 10, 53
 - labelio option 8
 - pmdlog option 8
 - procs option 8
 - stdoutmode option 9
 - communication with nodes 59
 - compiling programs 59
 - description 1
 - exit status 60
 - how it works 59
 - internals 59
 - options 8
 - running 3
 - running, examples 4
 - signal handling 60
 - starting applications 59
- Parallel Operating Environment (POE), description 1
- parallel program 21
- parallel programs 35
 - profiling 48
 - safe 49
 - tuning 47
- parallel task 21
- Partition Manager Daemon 38
- performance 47
- POE
 - eulib 9
 - hostfile option 8
 - ilevel option 8
 - infolevel option 8, 10, 53
 - labelio option 8
 - llfile option 9
 - pmdlog option 8
 - proc option 8
 - rmpool option 9
 - stdoutmode option 9
 - communication with nodes 59
 - compiling programs 59
 - description 1
 - exit status 60
 - how it works 59
 - internals 59
 - options 8
 - running 3
 - running, examples 4
 - signal handling 60
 - starting applications 59
- POE options 8
- preface vii
- prerequisite knowledge for this book vii
- problem diagnosis 36
- problems, common
 - bad output 46
 - cannot compile a parallel program 36
 - cannot connect with the remote host 37
 - cannot execute a parallel program 38
 - cannot start a parallel job 37
 - hangs 42
 - no output 41

- problems, common (*continued*)
 - user not authorized on remote host 38
- processor node, defined 1
- profiling program 48
- protocol
 - IP 33
 - User Space (US) 32

R

- redundancy 31
- return code 41
- running POE 3
 - running 4
 - with environment variables 5

S

- safe coding practices 51
 - fairness 51
 - order 50
 - resource limitations 51
 - safe program, described 49
- safety
 - MPI programs 45
 - threaded programs 49
- sample program, to illustrate messages 53
- security, supported methods 2
- sequential program 21
- shared memory protocol 33
- sine series algorithm 29
- starting applications with POE 59
- startup problems 38
- synchronization 21

T

- threaded programs
 - debugging 34, 46
 - hangs 44
 - performance tuning 47
 - protocol implications 33
 - safety 49
- trademarks 65
- tuning 47
 - threaded programs 47

U

- unrolling loops 22
 - example 22
- User Space (US) protocol 32
- using assistive technologies 61

X

- xinetd 38

Reader's Comments– We'd like to hear from you

IBM Parallel Environment for Linux
Introduction
Version 4 Release 2

Publication No. SA23-2218-00

Overall, how satisfied are you with the information in this book?

	Very Satisfied	Satisfied	Neutral	Dissatisfied	Very Dissatisfied
Overall satisfaction	<input type="checkbox"/>				

How satisfied are you that the information in this book is:

	Very Satisfied	Satisfied	Neutral	Dissatisfied	Very Dissatisfied
Accurate	<input type="checkbox"/>				
Complete	<input type="checkbox"/>				
Easy to find	<input type="checkbox"/>				
Easy to understand	<input type="checkbox"/>				
Well organized	<input type="checkbox"/>				
Applicable to your tasks	<input type="checkbox"/>				

Please tell us how we can improve this book:

Thank you for your responses. May we contact you? Yes No

When you send comments to IBM, you grant IBM a nonexclusive right to use or distribute your comments in any way it believes appropriate without incurring any obligation to you. IBM or any other organizations will only use the personal information that you supply to contact you about the issues that you state on this form.

Name

Address

Company or Organization

Phone No.



Fold and Tape

Please do not staple

Fold and Tape



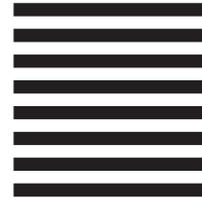
NO POSTAGE
NECESSARY
IF MAILED IN THE
UNITED STATES

BUSINESS REPLY MAIL

FIRST-CLASS MAIL PERMIT NO. 40 ARMONK, NEW YORK

POSTAGE WILL BE PAID BY ADDRESSEE

IBM Corporation
Department 55JA, Mail Station P384
2455 South Road
Poughkeepsie NY
12601-5400



Fold and Tape

Please do not staple

Fold and Tape



Program Number: 5724-N05

SA23-2218-00

