



Linux® Application Tuning Guide for SGI®
X86-64 Based Systems

007-5646-006

COPYRIGHT

© 2010–2013, SGI. All rights reserved; provided portions may be copyright in third parties, as indicated elsewhere herein. No permission is granted to copy, distribute, or create derivative works from the contents of this electronic documentation in any manner, in whole or in part, without the prior written permission of SGI.

LIMITED RIGHTS LEGEND

The software described in this document is "commercial computer software" provided with restricted rights (except as to included open/free source) as specified in the FAR 52.227-19 and/or the DFAR 227.7202, or successive sections. Use beyond license provisions is a violation of worldwide intellectual property laws, treaties and conventions. This document is provided with limited rights as defined in 52.227-14.

TRADEMARKS AND ATTRIBUTIONS

Altix, ICE, NUMalink, OpenMP, Performance Co-Pilot, SGI, the SGI logo, SHMEM, and UV are trademarks or registered trademarks of Silicon Graphics International Corp. or its subsidiaries in the United States and other countries.

Cray is a registered trademark of Cray, Inc. Dinkumware is a registered trademark of Dinkumware, Ltd. Intel, GuideView, Itanium, KAP/Pro Toolset, Phi, VTune, and Xeon are trademarks or registered trademarks of Intel Corporation, in the United States and other countries. Oracle and Java are registered trademarks of Oracle and/or its affiliates. Linux is a registered trademark of Linus Torvalds in several countries. Red Hat is a registered trademark of Red Hat, Inc. PostScript is a trademark of Adobe Systems Incorporated. TotalView and TotalView Technologies are registered trademarks and TVD is a trademark of Rogue Wave Software, Inc. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries. All other trademarks are the property of their respective owners.

New Features

Contains revised information that explains how to install SGI PerfSocket.

Record of Revision

Version	Description
001	November 2010 Original publication.
002	February 2011 Supports the SGI Performance Suite 1.1 release.
003	November 2011 Supports the SGI Performance Suite 1.3 release.
004	May 2012 Supports the SGI Performance Suite 1.4 release.
005	November 2013 Supports the SGI Performance Suite 1.7 release.
006	November 2013 Supports the SGI Performance Suite 1.7 release and includes a correction to the PerfSocket installation documentation.

Contents

About This Guide	xiii
Related SGI Publications	xiii
Related Publications From Other Sources	xv
Obtaining Publications	xvi
Conventions	xvi
Reader Comments	xvi
1. System Overview	1
Scalable Computing	1
An Overview of SGI System Architecture	2
The Basics of Memory Management	2
2. The SGI Compiling Environment	3
Compiler Overview	3
Environment Modules	4
Library Overview	6
Static Libraries	6
Dynamic Libraries	6
C/C++ Libraries	6
SHMEM Message Passing Libraries	7
Other Compiling Environment Features	8
3. Performance Analysis and Debugging	9
Determining System Configuration	9
Sources of Performance Problems	16
007-5646-006	vii

Profiling with <code>perf</code>	16
Profiling with <code>PerfSuite</code>	16
Using VTune for Remote Sampling	17
Other Performance Tools	17
Debugging Tools	18
Using the Intel Debugger <code>idb</code>	20
Using <code>ddd</code>	20
4. Monitoring Tools	23
System Monitoring Tools	23
Hardware Inventory and Usage Commands	23
<code>topology(1)</code> Command	23
<code>gtopology(1)</code> Command	24
Performance Co-Pilot Monitoring Tools	27
<code>hubstats(1)</code> Command	28
<code>linkstat-uv(1)</code> Command	28
Other Performance Co-Pilot Monitoring Tools	28
System Usage Commands	30
Memory Statistics and <code>nodeinfo</code> Command	33
5. Data Placement Tools	35
Data Placement Tools Overview	35
Distributed Shared Memory (DSM)	36
ccNUMA Architecture	36
Cache Coherency	37
Non-uniform Memory Access (NUMA)	37
Data Placement Practices	37
<code>taskset</code> Command	39

numactl Command	41
dplace Command	41
Using the dplace Command	42
dplace for Compute Thread Placement Troubleshooting Case Study	47
dlook Command	50
Using the dlook Command	50
omplace Command	57
Installing NUMA Tools	58
About cpusets	58
6. Performance Tuning	59
Single Processor Code Tuning	59
Getting the Correct Results	60
Managing Heap Corruption Problems	61
Using Tuned Code	61
Determining Tuning Needs	62
Using Compiler Options Where Possible	62
Tuning the Cache Performance	65
Managing Memory	67
Memory Use Strategies	67
Memory Hierarchy Latencies	68
Multiprocessor Code Tuning	68
Data Decomposition	69
Parallelizing Your Code	70
Use MPT	71
Use OpenMP	71
OpenMP Nested Parallelism	72

Use Compiler Options	72
Identifying Parallel Opportunities in Existing Code	73
Fixing False Sharing	73
Using <code>dplace</code> and <code>taskset</code>	74
Environment Variables for Performance Tuning	74
Understanding Parallel Speedup and Amdahl's Law	75
Adding CPUs to Shorten Execution Time	76
Understanding Parallel Speedup	76
Understanding Superlinear Speedup	77
Understanding Amdahl's Law	77
Calculating the Parallel Fraction of a Program	78
Predicting Execution Time with n CPUs	79
Gustafson's Law	80
Floating-point Program Performance	81
About MPI Application Tuning	81
MPI Application Communication on SGI Hardware	82
MPI Job Problems and Application Design	82
MPI Performance Tools	84
Using Transparent Huge Pages (THPs) in MPI, SHMEM, and UPC Applications	85
Enabling Huge Pages in MPI, SHMEM, and UPV Applications on Systems Without THP	86
7. Flexible File I/O	87
FFIO Operation	87
Environment Variables	88
Simple Examples	89
Multithreading Considerations	92
Application Examples	93

Event Tracing	94
System Information and Issues	94
8. I/O Tuning	95
Application Placement and I/O Resources	95
Layout of Filesystems and XVM for Multiple RAIDs	96
9. Suggested Shortcuts and Workarounds	97
Determining Process Placement	97
Example Using pthreads	98
Example Using OpenMP	100
Combination Example (MPI and OpenMP)	102
Resetting System Limits	104
Resetting the File Limit Resource Default	105
Resetting the Default Stack Size	107
Avoiding Segmentation Faults	107
Resetting Virtual Memory Size	109
Linux Shared Memory Accounting	110
OFED Tuning Requirements for UPC and SHMEM	111
Setting Java Environment Variables	112
10. Using PerfSocket	115
About SGI PerfSocket	115
Installing and Using PerfSocket	115
Installing PerfSocket (Administrator Procedure)	116
Running an Application With PerfSocket	117
About Security When Using PerfSocket	118
Troubleshooting	118

Index 119

About This Guide

This publication provides information about tuning application programs on the SGI® UV™ series systems running the Linux operating system. Application programs includes Fortran and C programs written with the Intel-provided compilers on SGI Linux systems. Some parts of this manual are also applicable to other SGI X86-64 based systems, such as the SGI® ICE™ X and SGI® Rackable™ systems.

This guide is written for experienced programmers, familiar with Linux commands and with either the C or Fortran programming languages. The focus in this document is on achieving the highest possible performance by exploiting the features of your SGI system. The material assumes that you know the basics of software engineering and that you are familiar with standard methods and data structures. If you are new to programming or software design, this guide will **not** be of use to you.

Related SGI Publications

The release notes for the SGI Foundation Suite and the SGI Performance Suite list SGI publications that pertain to the specific software packages in those products. The release notes reside in a text file in the `/docs` directory on the product media. For example, `SGI-MPI-1.x-readme.txt`. After installation, the release notes and other product documentation reside in the `/usr/share/doc/packages/product` directory.

All SGI publications are available on the Technical Publications Library at <http://docs.sgi.com>. The following publications provide information about Linux implementations on SGI systems:

- *SGI UV System Software Installation and Configuration Guide*

Explains how to install the operating system on an SGI UV system. This manual also includes information about basic configuration features such as CPU frequency scaling and partitioning.

- *SGI Cpuset Software Guide*

Explains how to use cpusets within your application program. Cpusets restrict processes within a program to specific processors or memory nodes.

- *Message Passing Toolkit (MPT) User Guide*

Describes industry-standard message passing protocol optimized for SGI computers. This manual describes how to tune the run-time environment to improve the performance of an MPI message passing application on SGI computers. None of these ways involve application code changes.

- *MPInside Reference Guide*

Documents the SGI MPInside MPI profiling tool.

- SGI hardware documentation.

SGI creates hardware manuals that are specific to each product line. The hardware documentation typically includes a system architecture overview and describes the major components. It also provides the standard procedures for powering on and powering off the system, basic troubleshooting information, and important safety and regulatory specifications.

The following procedure explains how to retrieve a list of hardware manuals for your system.

Procedure 0-1 To retrieve hardware documentation

1. Type the following URL into the address bar of your browser:

`docs.sgi.com`

2. In the search box on the Techpubs Library, narrow your search as follows:

- In the **search** field, type the model of your SGI system.

For example, type one of the following: "UV 2000", "ICE X", Rackable.

Remember to enclose hardware model names in quotation marks (" ") if the hardware model name includes a space character.

- Check **Search only titles**.
- Check **Show only 1 hit/book**.
- Click **search**.

Related Publications From Other Sources

Compilers and performance tool information for software that runs on SGI Linux systems is available from a variety of sources. The following additional documents might be useful to you:

- <http://sourceware.org/gdb/documentation/>

GDB: The GNU Project Debugger website with documentation, such as, *Debugging with GDB*, *GDB User Manual*, and so on.

- <http://www.intel.com/cd/software/products/asmo-na/eng/perflib/219780.htm>; documentation for Intel compiler products can be downloaded from this website.

Intel Software Network page with links to Intel documentation, such as, *Intel Professional Edition Compilers*, *Intel Thread Checker*, *Intel VTune Performance Analyzer*, and various Intel cluster software solutions.

- Intel provides detailed application tuning information including the Intel Xeon processor 5500 at [http://www.intel.com/Assets/en_US/PDF/manual/248966.pdf?wapkw= Intel Xeon processor 5500 Series tuning manual](http://www.intel.com/Assets/en_US/PDF/manual/248966.pdf?wapkw=Intel%20Xeon%20processor%205500%20Series%20tuning%20manual)
- Intel provides specific tuning information tutorial for Nehalem (Intel Xeon 5500) at <http://software.intel.com/sites/webinar/tuning-your-application-for-nehalem/>.
- Intel provides information for Westmere (Intel Xeon 5600) at <http://www.intel.com/itcenter/products/xeon/5600/index.htm>
- <http://software.intel.com/en-us/articles/intel-vtune-performance-analyzer-for-linux-documentation/>

Intel Software Network page with information specific to *Intel VTune Performance Analyzer* including links to documentation.

- Intel provides information about the Intel Performance Tuning Utility (PTU) at <http://software.intel.com/en-us/articles/intel-performance-tuning-utility/>.
- Information about the OpenMP Standard can be found at <http://openmp.org/wp/>.

The OpenMP API specification for parallel programming website is found here.

Obtaining Publications

You can obtain SGI documentation in the following ways:

- You can access the SGI Technical Publications Library at the following website:

<http://docs.sgi.com>

Various formats are available. This library contains the most recent and most comprehensive set of online books, release notes, man pages, and other information.

- You can view man pages by typing `man title` at a command line.

Conventions

The following conventions are used in this documentation:

[]	Brackets enclose optional portions of a command or directive line.
<code>command</code>	This fixed-space font denotes literal items such as commands, files, routines, path names, signals, messages, and programming language structures.
...	Ellipses indicate that a preceding element can be repeated.
user input	This bold, fixed-space font denotes literal items that the user enters in interactive sessions. (Output is shown in nonbold, fixed-space font.)
<i>variable</i>	Italic typeface denotes variable entries and words or concepts being defined.
<code>manpage(x)</code>	Man page section identifiers appear in parentheses after man page names.

Reader Comments

If you have comments about the technical accuracy, content, or organization of this publication, contact SGI. Be sure to include the title and document number of the publication with your comments. (Online, the document number is located in the

front matter of the publication. In printed publications, the document number is located at the bottom of each page.)

You can contact SGI in either of the following ways:

- Send e-mail to the following address:

techpubs@sgi.com

- Contact your customer service representative and ask that an incident be filed in the SGI incident tracking system:

<http://www.sgi.com/support/supportcenters.html>

SGI values your comments and will respond to them promptly.

System Overview

Tuning an application involves making your program run its fastest on the available hardware. The first step is to make your program run as efficiently as possible on a single processor system and then consider ways to use parallel processing.

This chapter provides an overview of concepts involved in working in parallel computing environments.

Scalable Computing

Scalability is computational power that can grow over a large number of CPUs. Scalability depends on the time between nodes on the system. *Latency* is the time to send the first byte between nodes.

A Symmetric Multiprocessor (SMP) is a parallel programming environment in which all processors have equally fast (symmetric) access to memory. These types of systems are easy to assemble and have limited scalability due to memory access times.

On a symmetric multiprocessor (SMP) machine, all data is visible from all processors. NonUniform Memory Access (NUMA) machines also have a shared address space. In both cases, there is a single shared memory space and a single operating system instance. However, in an SMP machine, each processor is functionally identical and has equal time access to every memory address. In contrast, a NUMA system has a shared address space, but the access time to memory varies over physical address ranges and between processing elements. The Intel® Xeon® 7500 series processor (Nehalem i7 architecture) and the Intel Xeon processor E5 family platform are examples of NUMA architecture. Each processor has its own memory and can address the memory attached to another processor through the Quick Path Interconnect (QPI). For more information, see the system architecture overview in "Data Placement Tools Overview" on page 35.

Another parallel environment is that of arrays, or clusters. Any networked computer can participate in a cluster. These are highly scalable, easy to assemble, but are often hard to use. There is no shared memory and there are frequently long latency times.

Massively Parallel Processors (MPPs) have a distributed memory and can scale to thousands of processors; they have large memories and large local memory bandwidth.

Scalable Symmetric Multiprocessors (S²MPs), as in the ccNUMA environment, combine qualities of SMPs and MPPs. They are logically programmable like an SMP and have MPP-like scalability.

An Overview of SGI System Architecture

For information about system architecture, see the hardware manuals that are available on the Tech Pubs Library at <http://docs.sgi.com>

The Basics of Memory Management

Virtual memory (VM), also known as virtual addressing, is used to divide a system's relatively small amount of physical memory among the potentially larger amount of logical processes in a program. It does this by dividing physical memory into *pages*, and then allocating pages to processes as the pages are needed.

A page is the smallest unit of system memory allocation. Pages are added to a process when either a page fault occurs or an allocation request is issued. Process size is measured in pages and two sizes are associated with every process: the total size and the resident set size (RSS). The number of pages being used in a process and the process size can be determined by using either the `ps(1)` or the `top(1)` command.

Swap space is used for temporarily saving parts of a program when there is not enough physical memory. The swap space may be on the system drive, on an optional drive, or allocated to a particular file in a filesystem. To avoid swapping, try not to overburden memory. Lack of adequate swap space limits the number and the size of applications that can run simultaneously on the system, and it can limit system performance. Access time to disk is orders of magnitude slower than access to random access memory (RAM). A system that runs out of memory and uses swap to disk while running a program will have its performance seriously affected, as swapping will become a major bottleneck. Be sure your system is configured with enough memory to run your applications.

Linux is a demand paging operating system, using a least-recently-used paging algorithm. Pages are mapped into physical memory when first referenced and pages are brought back into memory if swapped out. In a system that uses demand paging, the operating system copies a disk page into physical memory only if an attempt is made to access it, that is, a page fault occurs. A page fault handler algorithm does the necessary action. For more information, see the `mmap(2)` man page.

The SGI Compiling Environment

This chapter provides an overview of the SGI compiling environment on the SGI family of servers and covers the following topics:

- "Compiler Overview" on page 3
- "Environment Modules" on page 4
- "Library Overview" on page 6
- "Other Compiling Environment Features" on page 8

The remainder of this book provides more detailed examples of the use of the SGI compiling environment elements.

Compiler Overview

The Intel Fortran and C/C++ compilers are available from Intel Corporation or can be ordered from SGI on a reseller basis. The Fortran compiler supports OpenMP 3.0 and the C/C++ compiler is compatible with `gcc` and the C99 standard. Both the C/C++ and Fortran Intel compilers support OpenMP 3.0. It should be noted that Intel has renamed their compilers, as follows:

- Intel® C++ Compilers are now called Intel® C++ Composer XE.
- Intel® Fortran Compilers `gcc`, `g++`, and `g77` compilers are now called Intel® Visual Fortran Composer XE and Intel® Visual Fortran Composer XE with IMSL.

In addition, the GNU Fortran and C compilers are available on SGI systems.

The following is the general form of the compiler command line (note that the Fortran command is used in this example):

```
% ifort [options] filename.extension
```

An appropriate filename extension is required for each compiler, according to the programming language used (Fortran, C, C++, or FORTRAN 77).

Some common compiler options are:

- `-o filename`: renames the output to *filename*.

- `-g`: produces additional symbol information for debugging.
- `-O[level]`: invokes the compiler at different optimization *levels*, from 0 to 3.
- `-Idirectory_name`: looks for `include` files in *directory_name*.
- `-c`: compiles without invoking the linker; this options produces an `a.o` file only.

Many processors do not handle denormalized arithmetic (for gradual underflow) in hardware. The support of gradual underflow is implementation-dependent. Use the `-ftz` option with the Intel compilers to force the flushing of denormalized results to zero.

Note that frequent gradual underflow arithmetic in a program causes the program to run very slowly, consuming large amounts of system time (this can be determined with the `time` command). In this case, it is best to trace the source of the underflows and fix the code; gradual underflow is often a source of reduced accuracy anyway. `prctl(1)` allows you to query or control certain process behavior. In a program, `prctl` tracks where floating point errors occur.

Environment Modules

A *module* is a user interface that provides for the dynamic modification of a user's environment. By loading a module, a user does not have to change environment variables in order to access different versions of the compilers, loaders, libraries and utilities that are installed on the system.

Modules can be used in the SGI compiling environment to customize the environment. If the use of modules is not available on your system, its installation and use is highly recommended.

To view which modules are available on your system, use the following command (for any shell environment):

```
% module avail
```

To load modules into your environment (for any shell), use the following commands:

```
% module load intel-compilers-latest mpt/2.04
```

Note: The above commands are for example use only; the actual release numbers may vary depending on the version of the software you are using. See the release notes that are distributed with your system for the pertinent release version numbers.

The module help command provides a list of all arguments accepted, as follows:

```
sys:~> module help
```

```
Modules Release 3.1.6 (Copyright GNU GPL v2 1991):
Available Commands and Usage:
+ add|load          modulefile [modulefile ...]
+ rm|unload         modulefile [modulefile ...]
+ switch|swap       modulefile1 modulefile2
+ display|show      modulefile [modulefile ...]
+ avail             [modulefile [modulefile ...]]
+ use [-a|--append] dir [dir ...]
+ unuse             dir [dir ...]
+ update
+ purge
+ list
+ clear
+ help              [modulefile [modulefile ...]]
+ whatis            [modulefile [modulefile ...]]
+ apropos|keyword  string
+ initadd           modulefile [modulefile ...]
+ initprepend      modulefile [modulefile ...]
+ initrm           modulefile [modulefile ...]
+ initswitch       modulefile1 modulefile2
+ initlist
+ initclear
```

For details about using modules, see the `module(1)` man page.

Library Overview

Libraries are files that contain one or more object (.o) files. Libraries are used to simplify local software development by hiding compilation details. Libraries are sometimes also called *archives*.

The SGI compiling environment contains several types of libraries; an overview about each library is provided in this subsection.

Static Libraries

Static libraries are used when calls to the library components are satisfied at link time by copying text from the library into the executable. To create a static library, use the `ar(1)`, or an archiver command.

To use a static library, include the library name on the compiler's command line. If the library is not in a standard library directory, be sure to use the `-L` option to specify the directory and the `-l` option to specify the library filename.

To build an application to have all static versions of standard libraries in the application binary, use the `-static` option on the compiler command line.

Dynamic Libraries

Dynamic libraries are linked into the program at run time and when loaded into memory can be accessed by multiple programs. Dynamic libraries are formed by creating a Dynamic Shared Object (DSO).

Use the link editor command (`ld(1)`) to create a dynamic library from a series of object files or to create a DSO from an existing static library.

To use a dynamic library, include the library on the compiler's command line. If the dynamic library is not in one of the standard library directories, use the `-L path` and `-l library_shortname` compiler options during linking. You must also set the `LD_LIBRARY_PATH` environment variable to the directory where the library is stored before running the executable.

C/C++ Libraries

The following C/C++ libraries are provided with the Intel compiler:

- `libguide.a`, `libguide.so`: for support of OpenMP-based programs.
- `libsvml.a`: short vector math library
- `libirc.a`: Intel's support for Profile-Guided Optimizations (PGO) and CPU dispatch
- `libimf.a`, `libimf.so`: Intel's math library
- `libcprts.a`, `libcprts.so`: Dinkumware C++ library
- `libunwind.a`, `libunwind.so`: Unwinder library
- `libcxa.a`, `libcxa.so`: Intel's runtime support for C++ features

SHMEM Message Passing Libraries

The SHMEM application programming interface is implemented by the `libisma` library and is part of the Message Passing Toolkit (MPT) product on SGI systems. The SHMEM programming model consists of library routines that provide low-latency, high-bandwidth communication for use in highly parallelized, scalable programs. The routines in the SHMEM application programming interface (API) provide a programming model for exchanging data between cooperating parallel processes. The resulting programs are similar in style to Message Passing Interface (MPI) programs. The SHMEM API can be used either alone or in combination with MPI routines in the same parallel program.

A SHMEM program is SPMD (single program, multiple data) in style. The SHMEM processes, called processing elements or PEs, all start at the same time, and they all run the same program. Usually the PEs perform computation on their own subdomains of the larger problem, and periodically communicate with other PEs to exchange information on which the next computation phase depends.

The SHMEM routines minimize the overhead associated with data transfer requests, maximize bandwidth, and minimize data latency. Data latency is the period of time that starts when a PE initiates a transfer of data and ends when a PE can use the data.

SHMEM routines support remote data transfer through put operations, which transfer data to a different PE, get operations, which transfer data from a different PE, and remote pointers, which allow direct references to data objects owned by another PE. Other operations supported are collective broadcast and reduction, barrier synchronization, and atomic memory operations. An atomic memory operation is an

atomic read-and-update operation, such as a fetch-and-increment, on a remote or local data object.

For details about using the SHMEM routines, see the `intro_shmem(3)` man page or the *Message Passing Toolkit (MPT) User's Guide*.

Other Compiling Environment Features

The SGI compiling environment includes several other products as part of its distribution:

- `idb`: the Intel debugger (available if your system is licensed for the Intel compilers). This is a fully symbolic debugger and supports Fortran, C, and C++ debugging. The Intel® Debugger for Linux is a fully fledged Eclipse graphical user interface based debug solution. To use the command line debugger (`gdb` like), you need to use the `idbc` command. It works with the following compilers: Intel® C++ Compilers (now called Intel® C++ Composer XE) and Intel® Fortran Compilers `gcc`, `g++`, and `g77` compilers (now called Intel® Visual Fortran Composer XE and Intel® Visual Fortran Composer XE with IMSL). The Intel Debugger can debug both single and multi-threaded applications, serial and parallel code. For more information, see <http://software.intel.com/en-us/articles/idb-linux/>.
- `gdb`: the GNU project debugger, which supports C, C++ and Modula-2. It also supports Fortran 95 debugging when the `gdbf95` patch is installed.
- `ddd`: a graphical user interface to `gdb` and the other debuggers.
- TotalView: a licensed graphical debugger useful in an MPI environment (see <http://www.roguewave.com/products/totalview-family/totalview.aspx>)

These and other performance analysis tools are discussed in Chapter 3, "Performance Analysis and Debugging" on page 9.

Performance Analysis and Debugging

Tuning an application involves determining the source of performance problems and then rectifying those problems to make your programs run their fastest on the available hardware. Performance gains usually fall into one of three categories of measured time:

- User CPU time: time accumulated by a user process when it is attached to a CPU and is executing.
- Elapsed (wall-clock) time: the amount of time that passes between the start and the termination of a process.
- System time: the amount of time performing kernel functions like system calls, `sched_yield`, for example, or floating point errors.

Any application tuning process involves the following steps:

1. Analyzing and identifying a problem
2. Locating where in the code the problem is
3. Applying an optimization technique

This chapter describes the process of analyzing your code to determine performance bottlenecks. See Chapter 6, "Performance Tuning" on page 59, for details about tuning your application for a single processor system and then tuning it for parallel processing.

Determining System Configuration

One of the first steps in application tuning is to determine the details of the system that you are running. Depending on your system configuration, different options may or may not provide good results.

The `topology(1)` command displays general information about SGI systems, with a focus on node information. This includes node counts for blades, node IDs, NASIDs, memory per node, system serial number, partition number, UV Hub versions, CPU to node mappings, and general CPU information. The `topology` command is installed by the `pcp-sgi` RPM package.

The following is example output:

```
uv-sys:~ # topology
System type: UV2000
System name: harp34-sys
Serial number: UV2-00000034
Partition number: 0
    8 Blades
    256 CPUs
    16 Nodes
235.82 GB Memory Total
15.00 GB Max Memory on any Node
    1 BASE I/O Riser
    2 Network Controllers
    2 Storage Controllers
    2 USB Controllers
    1 VGA GPU
```

The `cpumap(1)` command displays logical CPUs and shows relationships between them in a human-readable format. Aspects displayed include hyperthread relationships, last level cache sharing, and topological placement. The `cpumap` command gets its information from `/proc/cpuinfo`, the `/sys/devices/system` directory structure, and `/proc/sgi_uv/topology`. When creating cpusets, the Socket numbers reported in the output section Processor Numbering on Socket(s) corresponds to the `mems` argument you would use in the definition of a cpuset. The cpuset `mems` argument is the list of memory nodes that tasks in the cpuset are allowed to use. For more information, see the *SGI Cpuset Software Guide* available at <http://docs.sgi.com>. The following is example output:

```
uv# cpumap
Thu Sep 19 10:17:21 CDT 2013
harp34-sys.americas.sgi.com

This is an SGI UV
model name      : Genuine Intel(R) CPU @ 2.60GHz
Architecture    : x86_64
cpu MHz        : 2599.946
cache size     : 20480 KB (Last Level)

Total Number of Sockets      : 16
Total Number of Cores       : 128 (8 per socket)
```

```

Hyperthreading           : ON
Total Number of Physical Processors : 128
Total Number of Logical Processors  : 256   (2 per Phys Processor)
    
```

UV Information

```

HUB Version:             UVHub  3.0
Number of Hubs:          16
Number of connected Hubs: 16
Number of connected NUMALink ports: 128
    
```

=====

Hub-Processor Mapping

Hub Location	Processor Numbers -- HyperThreads in ()							
-----	-----							
0 r001i11b00h0	0	1	2	3	4	5	6	7
	(128	129	130	131	132	133	134	135)
1 r001i11b00h1	8	9	10	11	12	13	14	15
	(136	137	138	139	140	141	142	143)
2 r001i11b01h0	16	17	18	19	20	21	22	23
	(144	145	146	147	148	149	150	151)
3 r001i11b01h1	24	25	26	27	28	29	30	31
	(152	153	154	155	156	157	158	159)
4 r001i11b02h0	32	33	34	35	36	37	38	39
	(160	161	162	163	164	165	166	167)
5 r001i11b02h1	40	41	42	43	44	45	46	47
	(168	169	170	171	172	173	174	175)
6 r001i11b03h0	48	49	50	51	52	53	54	55
	(176	177	178	179	180	181	182	183)
7 r001i11b03h1	56	57	58	59	60	61	62	63
	(184	185	186	187	188	189	190	191)
8 r001i11b04h0	64	65	66	67	68	69	70	71
	(192	193	194	195	196	197	198	199)
9 r001i11b04h1	72	73	74	75	76	77	78	79
	(200	201	202	203	204	205	206	207)
10 r001i11b05h0	80	81	82	83	84	85	86	87
	(208	209	210	211	212	213	214	215)
11 r001i11b05h1	88	89	90	91	92	93	94	95
	(216	217	218	219	220	221	222	223)
12 r001i11b06h0	96	97	98	99	100	101	102	103

3: Performance Analysis and Debugging

```

( 224 225 226 227 228 229 230 231 )
13 r001i11b06h1 104 105 106 107 108 109 110 111
( 232 233 234 235 236 237 238 239 )
14 r001i11b07h0 112 113 114 115 116 117 118 119
( 240 241 242 243 244 245 246 247 )
15 r001i11b07h1 120 121 122 123 124 125 126 127
( 248 249 250 251 252 253 254 255 )

```

=====

Processor Numbering on Node(s)

Node	(Logical) Processors															
	0	1	2	3	4	5	6	7	128	129	130	131	132	133	134	135
0	0	1	2	3	4	5	6	7	128	129	130	131	132	133	134	135
1	8	9	10	11	12	13	14	15	136	137	138	139	140	141	142	143
2	16	17	18	19	20	21	22	23	144	145	146	147	148	149	150	151
3	24	25	26	27	28	29	30	31	152	153	154	155	156	157	158	159
4	32	33	34	35	36	37	38	39	160	161	162	163	164	165	166	167
5	40	41	42	43	44	45	46	47	168	169	170	171	172	173	174	175
6	48	49	50	51	52	53	54	55	176	177	178	179	180	181	182	183
7	56	57	58	59	60	61	62	63	184	185	186	187	188	189	190	191
8	64	65	66	67	68	69	70	71	192	193	194	195	196	197	198	199
9	72	73	74	75	76	77	78	79	200	201	202	203	204	205	206	207
10	80	81	82	83	84	85	86	87	208	209	210	211	212	213	214	215
11	88	89	90	91	92	93	94	95	216	217	218	219	220	221	222	223
12	96	97	98	99	100	101	102	103	224	225	226	227	228	229	230	231
13	104	105	106	107	108	109	110	111	232	233	234	235	236	237	238	239
14	112	113	114	115	116	117	118	119	240	241	242	243	244	245	246	247
15	120	121	122	123	124	125	126	127	248	249	250	251	252	253	254	255

=====

Sharing of Last Level (3) Caches

Socket	(Logical) Processors															
	0	1	2	3	4	5	6	7	128	129	130	131	132	133	134	135
0	0	1	2	3	4	5	6	7	128	129	130	131	132	133	134	135
1	8	9	10	11	12	13	14	15	136	137	138	139	140	141	142	143
2	16	17	18	19	20	21	22	23	144	145	146	147	148	149	150	151
3	24	25	26	27	28	29	30	31	152	153	154	155	156	157	158	159

4	32	33	34	35	36	37	38	39	160	161	162	163	164	165	166	167
5	40	41	42	43	44	45	46	47	168	169	170	171	172	173	174	175
6	48	49	50	51	52	53	54	55	176	177	178	179	180	181	182	183
7	56	57	58	59	60	61	62	63	184	185	186	187	188	189	190	191
8	64	65	66	67	68	69	70	71	192	193	194	195	196	197	198	199
9	72	73	74	75	76	77	78	79	200	201	202	203	204	205	206	207
10	80	81	82	83	84	85	86	87	208	209	210	211	212	213	214	215
11	88	89	90	91	92	93	94	95	216	217	218	219	220	221	222	223
12	96	97	98	99	100	101	102	103	224	225	226	227	228	229	230	231
13	104	105	106	107	108	109	110	111	232	233	234	235	236	237	238	239
14	112	113	114	115	116	117	118	119	240	241	242	243	244	245	246	247
15	120	121	122	123	124	125	126	127	248	249	250	251	252	253	254	255

=====
HyperThreading

Shared Processors

(0, 128)	(1, 129)	(2, 130)	(3, 131)
(4, 132)	(5, 133)	(6, 134)	(7, 135)
(8, 136)	(9, 137)	(10, 138)	(11, 139)
(12, 140)	(13, 141)	(14, 142)	(15, 143)
(16, 144)	(17, 145)	(18, 146)	(19, 147)
(20, 148)	(21, 149)	(22, 150)	(23, 151)
(24, 152)	(25, 153)	(26, 154)	(27, 155)
(28, 156)	(29, 157)	(30, 158)	(31, 159)
(32, 160)	(33, 161)	(34, 162)	(35, 163)
(36, 164)	(37, 165)	(38, 166)	(39, 167)
(40, 168)	(41, 169)	(42, 170)	(43, 171)
(44, 172)	(45, 173)	(46, 174)	(47, 175)
(48, 176)	(49, 177)	(50, 178)	(51, 179)
(52, 180)	(53, 181)	(54, 182)	(55, 183)
(56, 184)	(57, 185)	(58, 186)	(59, 187)
(60, 188)	(61, 189)	(62, 190)	(63, 191)
(64, 192)	(65, 193)	(66, 194)	(67, 195)
(68, 196)	(69, 197)	(70, 198)	(71, 199)
(72, 200)	(73, 201)	(74, 202)	(75, 203)
(76, 204)	(77, 205)	(78, 206)	(79, 207)
(80, 208)	(81, 209)	(82, 210)	(83, 211)
(84, 212)	(85, 213)	(86, 214)	(87, 215)

```
( 88, 216) ( 89, 217) ( 90, 218) ( 91, 219)
( 92, 220) ( 93, 221) ( 94, 222) ( 95, 223)
( 96, 224) ( 97, 225) ( 98, 226) ( 99, 227)
( 100, 228) ( 101, 229) ( 102, 230) ( 103, 231)
( 104, 232) ( 105, 233) ( 106, 234) ( 107, 235)
( 108, 236) ( 109, 237) ( 110, 238) ( 111, 239)
( 112, 240) ( 113, 241) ( 114, 242) ( 115, 243)
( 116, 244) ( 117, 245) ( 118, 246) ( 119, 247)
( 120, 248) ( 121, 249) ( 122, 250) ( 123, 251)
( 124, 252) ( 125, 253) ( 126, 254) ( 127, 255)
```

The `x86info(1)` command displays x86 CPU diagnostics information. Type one of the following commands to load the `x86info(1)` command if the command is not already installed:

- On RHEL systems, type the following:

```
# yum install x86info.x86_64
```

- On SLES systems, type the following:

```
# zypper install x86info
```

The following is an example of `x86info(1)` command output:

```
uv44-sys:~ # x86info
x86info v1.25. Dave Jones 2001-2009
Feedback to .

Found 64 CPUs
-----
CPU #1
EFamily: 0 EModel: 2 Family: 6 Model: 46 Stepping: 6
CPU Model: Unknown model.
Processor name string: Intel(R) Xeon(R) CPU           E7520  @ 1.87GHz
Type: 0 (Original OEM) Brand: 0 (Unsupported)
Number of cores per physical package=16
Number of logical processors per socket=32
Number of logical processors per core=2
APIC ID: 0x0 Package: 0 Core: 0 SMT ID 0
-----
CPU #2
EFamily: 0 EModel: 2 Family: 6 Model: 46 Stepping: 6
```



```
CPU Model: Unknown model.
Processor name string: Intel(R) Xeon(R) CPU           E7520  @ 1.87GHz
Type: 0 (Original OEM)  Brand: 0 (Unsupported)
Number of cores per physical package=16
Number of logical processors per socket=32
Number of logical processors per core=2
APIC ID: 0x6    Package: 0  Core: 0  SMT ID 6
```

```
-----
CPU #3
EFamily: 0 EModel: 2 Family: 6 Model: 46 Stepping: 6
CPU Model: Unknown model.
Processor name string: Intel(R) Xeon(R) CPU           E7520  @ 1.87GHz
Type: 0 (Original OEM)  Brand: 0 (Unsupported)
Number of cores per physical package=16
Number of logical processors per socket=32
Number of logical processors per core=2
APIC ID: 0x10   Package: 0  Core: 0  SMT ID 16
```

...

You can also use the `uname` command, which returns the kernel version and other machine information. For example:

```
uv44-sys:~ # uname -a
Linux uv44-sys 2.6.32.13-0.4.1.1559.0.PTF-default #1 SMP 2010-06-15 12:47:25 +0200 x86_64 x86_64 x86_64
```

For more system information, change directory (`cd`) to the `/sys/devices/system/node/node0/cpu0/cache` directory.

For example:

```
uv44-sys:/sys/devices/system/node/node0/cpu0/cache # ls
index0 index1 index2 index3
```

Change directory to `index0` and list the contents, as follows:

```
uv44-sys:/sys/devices/system/node/node0/cpu0/cache/index0 # ls
coherency_line_size level number_of_sets physical_line_partition shared_cpu_list shared_cpu_map size type way
```

Sources of Performance Problems

There are usually three areas of program execution that can have performance slowdowns:

- CPU-bound processes: processes that are performing slow operations (such as `sqrt` or floating-point divides) or non-pipelined operations such as switching between add and multiply operations.
- Memory-bound processes: code which uses poor memory strides, occurrences of page thrashing or cache misses, or poor data placement in NUMA systems.
- I/O-bound processes: processes which are waiting on synchronous I/O, formatted I/O, or when there is library or system level buffering.

Several profiling tools can help pinpoint where performance slowdowns are occurring. The following sections describe some of these tools.

Profiling with `perf`

Linux Performance Events provides a performance analysis framework for systems that use Intel Xeon Phi technology. It includes hardware-level CPU performance monitoring unit (PMU) features, software counters, and tracepoints.

Before you use these profiling tools, make sure the `perf` RPM is installed. The `perf` RPM comes with the your operating system and is not an SGI product.

For more information, see the following man pages: `perf(1)`, `perf-stat(1)`, `perf-top(1)`, `perf-record(1)`, `perf-report(1)`, `perf-list(1)`. The `perf` RPM includes these man pages.

Profiling with `PerfSuite`

`PerfSuite` is a set of tools, utilities, and libraries that you can use to analyze application software performance Linux-based systems. You can use `PerfSuite` tools to perform performance-related activities, ranging from assistance with compiler optimization reports to hardware performance counting, profiling, and MPI usage summarization. `PerfSuite` is Open Source software. It is approved for licensing under the University of Illinois/NCSA Open Source License (OSI-approved).

For more information, see one of the following websites:

- <http://perfsuite.ncsa.uiuc.edu/>
- <http://perfsuite.sourceforge.net/>
- <http://www.ncsa.illinois.edu/UserInfo/Resources/Software/Tools/PerfSuite/>, which hosts NCSA-specific information about using PerfSuite tools

The `psrun` utility is a PerfSuite command line utility that gathers hardware performance information on an unmodified executable. For more information, see <http://perfsuite.ncsa.uiuc.edu/psrun/>.

Using VTune for Remote Sampling

The Intel[®] VTune[™] Performance Analyzer (now called Intel[®] VTune[™] Amplifier XE) does remote sampling experiments. The VTune data collector runs on the Linux system and an accompanying GUI runs on an IA-32 Windows machine, which is used for analyzing the results. VTune allows you to perform interactive experiments while connected to the host through its GUI. PTU (Performance Tuning Utility) is another tool which requires the Intel VTune license.

For details about using VTune, see the following URL:

<http://software.intel.com/en-us/articles/intel-vtune-performance-analyzer-for-linux-documentation/>

Other Performance Tools

The following performance tools also can be of benefit when you are trying to optimize your code:

- *Guide OpenMP Compiler* is an OpenMP implementation for C, C++, and Fortran from Intel.
- *Assure Thread Analyzer* from Intel locates programming errors in threaded applications with no recoding required.

For details about these products, see the following website:

<http://developer.intel.com/software/products/threading>

Note: These products have not been thoroughly tested on SGI systems. SGI takes no responsibility for the correct operation of third party products described or their suitability for any particular purpose.

Debugging Tools

Three debuggers are available to help you analyze your code:

- `gdb`: the GNU project debugger. This is useful for debugging programs written in C, C++, and Fortran 95. When compiling with C and C++, include the `-g` option on the compiler command line to produce the `dwarf2` symbols database used by `gdb`.

When using `gdb` for Fortran debugging, include the `-g` and `-O0` options. Do not use `gdb` for Fortran debugging when compiling with `-O1` or higher.

The debugger to be used for Fortran 95 codes can be downloaded from http://sourceforge.net/project/showfiles.php?group_id=56720 . (Note that the standard `gdb` compiler does not support Fortran 95 codes.) To verify that you have the correct version of `gdb` installed, use the `gdb -v` command. The output should appear similar to the following:

```
GNU gdb 5.1.1 FORTRAN95-20020628 (RC1)
Copyright 2002 Free Software Foundation, Inc.
```

For a complete list of `gdb` commands, see the `gdb` user guide online at http://sources.redhat.com/gdb/onlinedocs/gdb_toc.html or use the `help` option. Note that current instances of `gdb` do not report `ar.ec` registers correctly. If you are debugging rotating, register-based, software-pipelined loops at the assembly code level, try using `idb` instead.

- `idb`: the Intel debugger. This is a fully symbolic debugger for the Linux platform. The debugger provides extensive support for debugging programs written in C, C++, FORTRAN 77, and Fortran 90. `idb` includes a GUI and it supports both Intel and GNU compilers.

Running `idb` with the `-gdb` option on the shell command line provides `gdb(1)`-like user commands and debugger output.

- ddd: a GUI to a command line debugger. It supports `gdb` and `idb`. For details about usage, see the following subsection.
- TotalView: a licensed graphical debugger useful in an MPI environment (see <http://www.totalviewtech.com/>)

Figure 3-1 on page 19 shows a TotalView session.

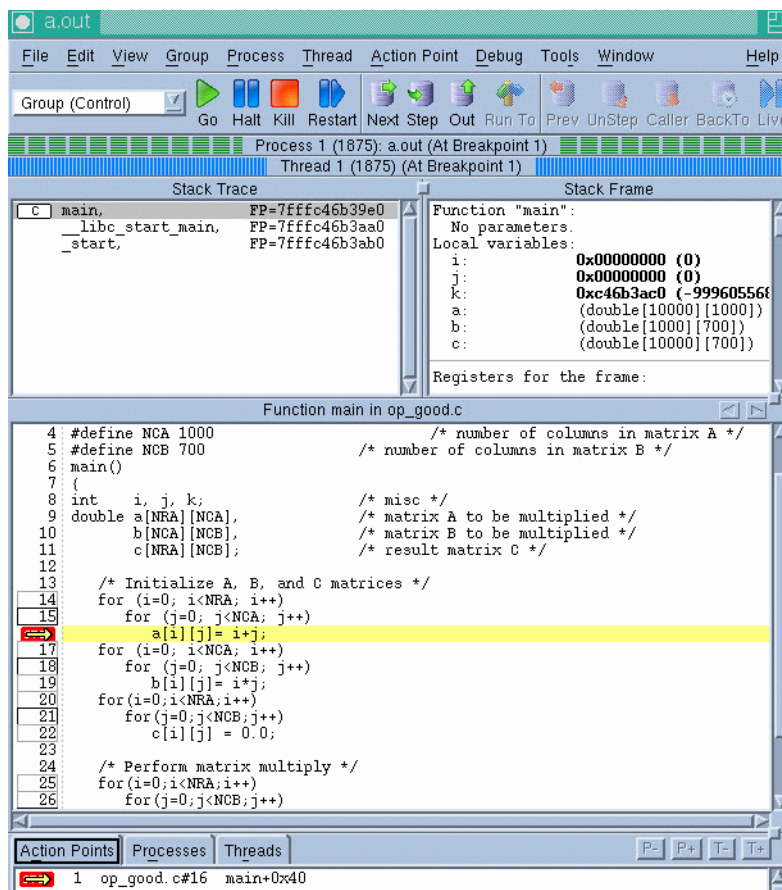


Figure 3-1 TotalView Session

Using the Intel Debugger `idb`

`idb` is part of the Intel® Compiler suite, both Fortran and C/C++ (now called Intel® Composer XE). You are asked during the installation if you want to install it or not. When running `idb` you get the GUI interface. When you invoke the `idbc` command, you get the command line interface.

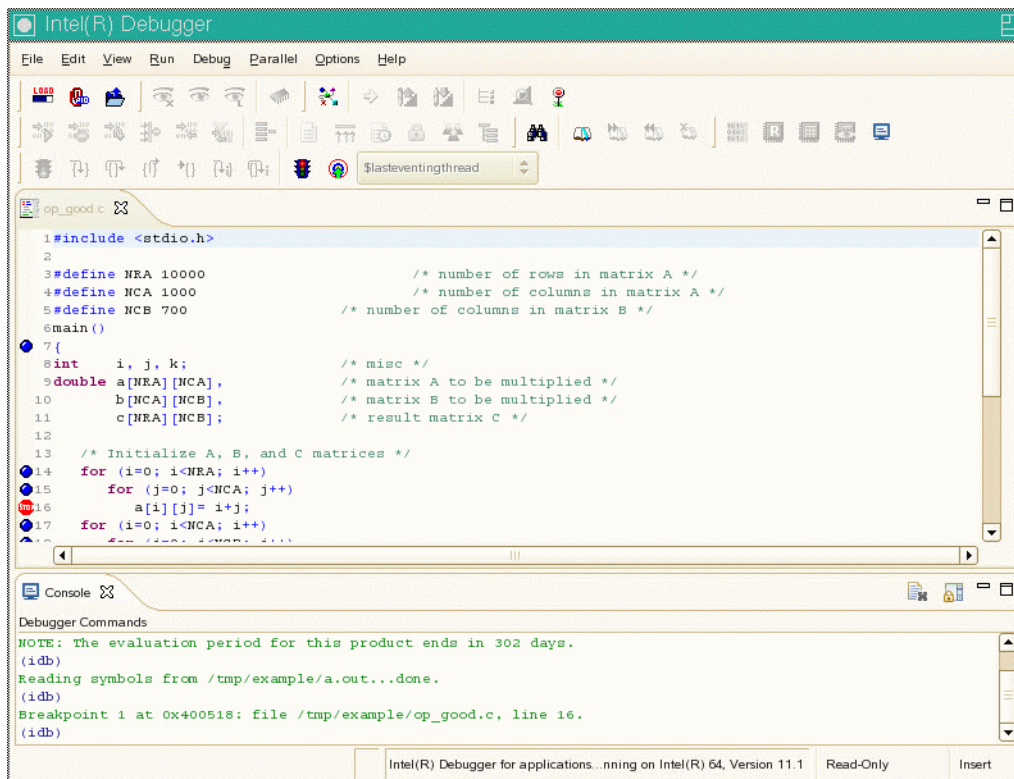


Figure 3-2 Intel® Debugger GUI

Using `ddd`

The `DataDisplayDebugger ddd(1)` tool is a GUI to an arbitrary command line debugger. When starting `ddd`, use the `--debugger` option to specify the debugger used (for example, `--debugger "idb"`). The default debugger used is `gdb`.

When the debugger is loaded the DataDisplayDebugger screen appears divided into panes that show the following information:

- Array inspection
- Source code
- Disassembled code
- A command line window to the debugger engine

These panes can be switched on and off from the **View** menu.

Some commonly used commands can be found on the menus. In addition, the following actions can be useful:

- Select an address in the assembly view, click the right mouse button, and select `lookup`. The `gdb` command is executed in the command pane and it shows the corresponding source line.
- Select a variable in the source pane and click the right mouse button. The current value is displayed. Arrays are displayed in the array inspection window. You can print these arrays to PostScript by using the **Menu>Print Graph** option.
- You can view the contents of the register file, including general, floating-point, NaT, predicate, and application registers by selecting **Registers** from the **Status** menu. The **Status** menu also allows you to view stack traces or to switch OpenMP threads.

Monitoring Tools

This chapter describes several tools that you can use to monitor system performance. The tools are divided into two general categories: system monitoring tools and nonuniform memory access (NUMA) tools.

System monitoring tools include the `topology(1)`, `top(1)` commands and the Performance Co-Pilot `pmchart(1)` command and other operating system commands such as the `vmstat(1)`, `iostat(1)` command and the `sar(1)` commands that can help you determine where system resources are being spent.

The `gtopology(1)` command displays a 3D scene of the system interconnect using the output from the `topology(1)` command.

System Monitoring Tools

You can use system utilities to better understand the usage and limits of your system. These utilities allow you to observe both overall system performance and single-performance execution characteristics. This section covers the following topics:

- "Hardware Inventory and Usage Commands" on page 23
- "Performance Co-Pilot Monitoring Tools" on page 27
- "System Usage Commands" on page 30
- "Memory Statistics and `nodeinfo` Command" on page 33

Hardware Inventory and Usage Commands

This section describes hardware inventory and usage commands and covers the following topics:

- "`topology(1)` Command" on page 23
- "`gtopology(1)` Command" on page 24

`topology(1)` Command

The `topology(1)` command provides topology information about your system.

Applications programmers can use the `topology` command to help optimize execution layout for their applications. For more information, see the `topology(1)` man page.

For an example of the `topology` command's output, see "Determining System Configuration" on page 9.

`gtopology(1)` Command

The `gtopology(1)` command is included as part of the `sgi-ppc` package of the SGI Accelerate, part of SGI Performance Suite software. It displays a 3D scene of the system interconnect using the output from the `topology(1)` command. See the man page for more details.

Figure 4-1 on page 25, shows the ring topology (the eight nodes are shown in pink, the NUMAlink connections in cyan) of an SGI system with 16 CPUs.

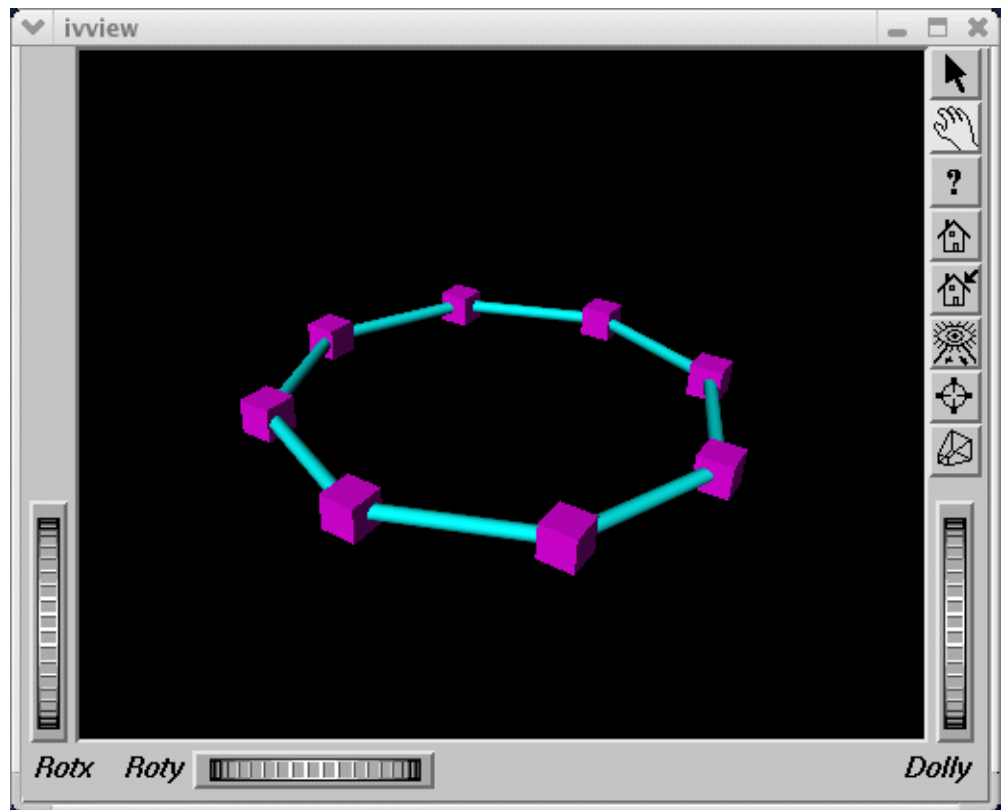


Figure 4-1 Ring Topology of an System with 16 CPUs

Figure 4-2 on page 26, shows the fat-tree topology of an SGI system with 32 CPUs. Again, nodes are the pink cubes. Routers are shown as blue spheres (if all ports are used) otherwise, yellow.

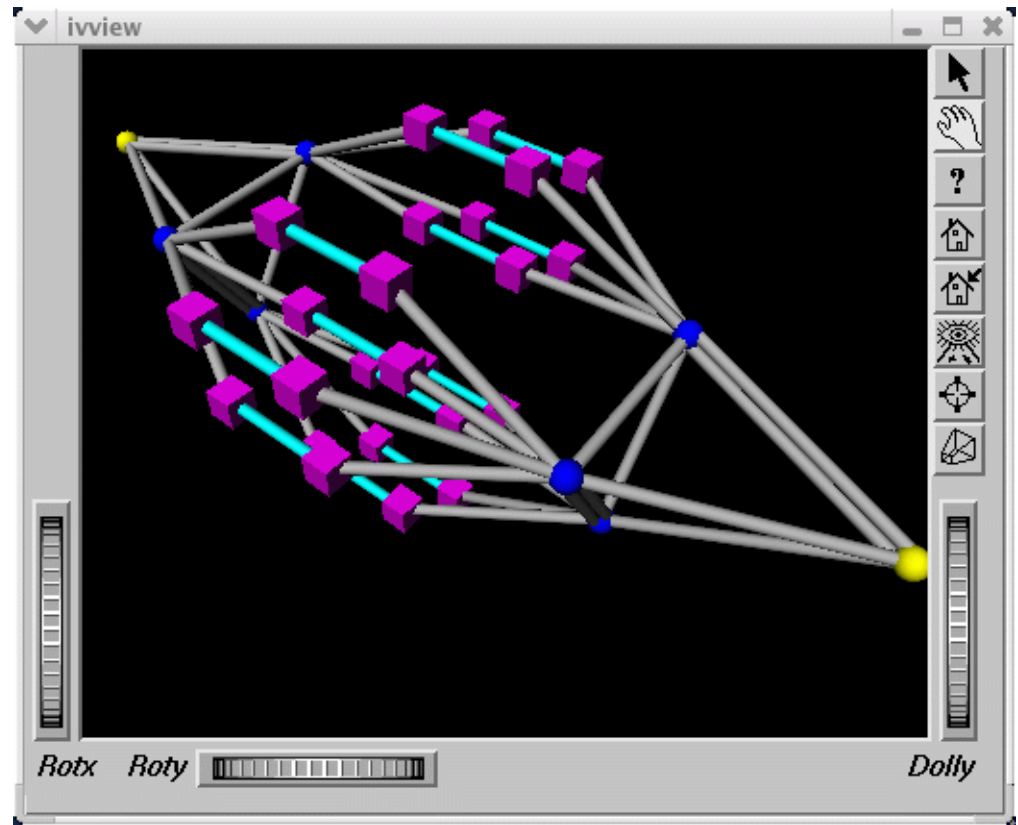


Figure 4-2 An SGI System with 32 CPUs Fat-tree Topology

Figure 4-3 on page 27, shows an SGI system with 512 CPUs. The dual planes of the fat-tree topology are clearly visible.

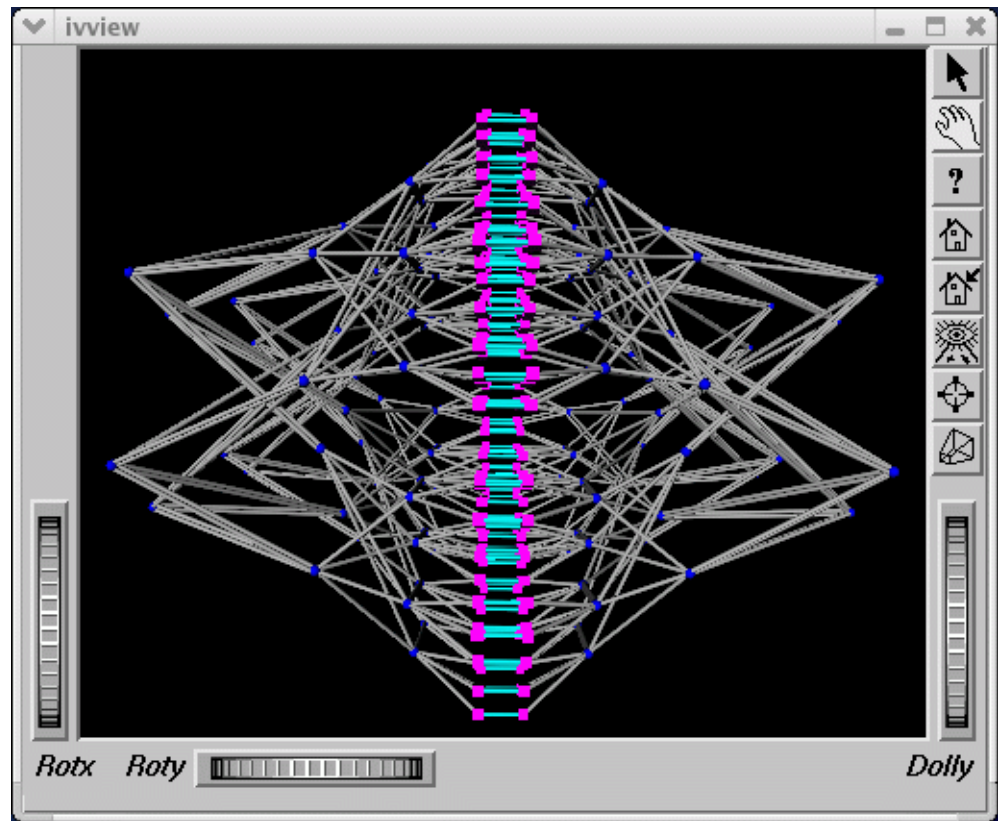


Figure 4-3 An SGI System with 512 CPUs

Performance Co-Pilot Monitoring Tools

This section describes Performance Co-Pilot monitoring tools and covers the following topics:

- "hubstats(1) Command" on page 28
- "linkstat-uv(1) Command" on page 28
- "Other Performance Co-Pilot Monitoring Tools" on page 28

hubstats(1) Command

The `hubstats(1)` command monitors NUMAlink traffic, directory cache operations, and global reference unit (GRU) traffic statistics on SGI UV systems. It does not work on any other platform. For more information, see the `hubstats(1)` man page.

linkstat-uv(1) Command

The `linkstat-uv(1)` command monitors NUMAlink traffic and error rates on SGI UV systems. It does not work on any other platform. This command returns information about packets and Mbytes sent/received on each NUMAlink in the system, as well as error rates. It is useful as a performance monitoring tool and as a tool for helping you to diagnose and identify faulty hardware. For more information, see the `linkstat-uv(1)` man page.

Note that this command is specific to SGI UV systems and does not return the same information as the `linkstat(1)` command.

Other Performance Co-Pilot Monitoring Tools

In addition to the UV specific tools described above, the `pcp` and `pcp-sgi` packages also provide numerous other performance monitoring tools, both graphical and text-based. It is important to remember that all of the performance metrics displayed by any of the tools described in this chapter can also be monitored with other tools such as `pmchart(1)`, `pmval(1)`, `pminfo(1)` and others. Additionally, the `pmlogger(1)` command can be used to capture Performance Co-Pilot archives, which can then be "replayed" during a retrospective performance analysis.

A very brief description of other Performance Co-Pilot monitoring tools follows. See the associated man page for each tool for more details.

- `pmchart(1)` — graphical stripchart tool, chiefly used for investigative performance analysis.
- `pmgsys(1)` — graphical tool showing miniature CPU, Disk, Network, LoadAvg and memory/swap in a miniature display, for example, useful for permanent residence on your desktop for the servers you care about.
- `pmgcluster(1)` — `pmgsys`, but for multiple hosts and thus useful for monitoring a cluster of hosts or servers.
- `clustervis(1)` — 3D display showing per-CPU and per-Network performance for multiple hosts.

- `nfsvis(1)` - 3D display showing NFS client/server traffic, grouped by NFS operation type
- `nodevis(1)` - 3D display showing per-node CPU and memory usage.
- `webvis(1)` - 3D display showing per-httpd traffic.
- `dkvis(1)` - 3D display showing per-disk traffic, grouped by controller.
- `diskstat(1)` - command line tool for monitoring disk traffic.
- `topdisk(1)` - command line, curses-based tool, for monitoring disk traffic.
- `topsys(1)` - command line, curses-based tool, for monitoring processes making a large numbers of system calls or spending a large percentage of their execution time in system mode using assorted system time measures.
- `pmgxvm(1)` - miniature graphical display showing XVM volume topology and performance statistics.
- `osvis(1)` - 3D display showing assorted kernel and system statistics.
- `pmdumptext(1)` - command line tool for monitoring multiple performance metrics with a highly configurable output format. Therefore, it is a useful tools for scripted monitoring tasks.
- `pmval(1)` - command line tool, similar to `pmdumptext(1)`, but less flexible.
- `pminfo(1)` - command line tool, useful for printing raw performance metric values and associated help text.
- `pmprobe(1)` - command line tool useful for scripted monitoring tasks.
- `pmie(1)` - a performance monitoring inference engine. This is a command line tool with an extraordinarily powerful underlying language. It can also be used as a system service for monitoring and reporting on all sorts of performance issues of interest.
- `pmieconf(1)` - command line tool for creating and customizing "canned" `pmie(1)` configurations.
- `pmlogger(1)` - command line tool for capturing Performance Co-Pilot performance metrics archives for replay with other tools.
- `pmlogger_daily(1)` and `pmlogger_check(1)` - cron driven infrastructure for automated logging with `pmlogger(1)`.

- `pmcd(1)` – the Performance Co-Pilot metrics collector daemon
- `PCPIntro(1)` - introduction to Performance Co-Pilot monitoring tools, generic command line usage and environment variables
- `PMAPI(3)` - introduction to the Performance Co-Pilot API libraries for developing new performance monitoring tools
- `PMDA(3)` - introduction to the Performance Co-Pilot Metrics Domain Agent API, for developing new Performance Co-Pilot agents
- `topology(1)` – displays general information about SGI UV systems, with a focus on node information. This includes counts of the CPUs, nodes, routers and memory, as well as, various I/O devices. More detailed information is available for node IDs, NASIDs, memory per node, system serial number, partition number, UV Hub versions, CPU to node mappings, I/O device descriptions, and general CPU and I/O information.

System Usage Commands

Several commands can be used to determine user load, system usage, and active processes.

To determine the system load, use the `uptime(1)` command, as follows:

```
uv44-sys:~ # uptime
 3:48pm up 2:50, 5 users, load average: 0.12, 0.25, 0.40
```

The output displays the current time, the length of time the system has been up, the number of users on the system, and the average number of jobs in the run queue over the last one, five, and 15 minutes.

To determine who is using the system and for what purpose, use the `w(1)` command, as follows:

```
uv44-sys:~ # w
 15:47:48 up 2:49, 5 users, load average: 0.04, 0.27, 0.42
USER      TTY      LOGIN@  IDLE   JCPU   PCPU   WHAT
root      pts/0    13:10   1:41m  0.07s  0.07s  -bash
root      pts/2    13:31   0.00s  0.14s  0.02s  w
boetcher pts/4    14:30   2:13   0.73s  0.73s  -csh
root      pts/5    14:32   1:14m  0.04s  0.04s  -bash
root      pts/6    15:09   31:25  0.08s  0.08s  -bash
```


The output from this command shows who is on the system, the duration of user sessions, processor usage by user, and currently executing user commands.

To determine active processes, use the `ps(1)` command, which displays a snapshot of the process table. The `ps -A` command selects all the processes currently running on a system as follows:

```
[user@profit user]# ps -A
  PID TTY          TIME CMD
    1 ?            00:00:06 init
    2 ?            00:00:00 migration/0
    3 ?            00:00:00 migration/1
    4 ?            00:00:00 migration/2
    5 ?            00:00:00 migration/3
    6 ?            00:00:00 migration/4
      .
      .
      .
 1086 ?            00:00:00 sshd
  1120 ?            00:00:00 xinetd
  1138 ?            00:00:05 ntpd
  1171 ?            00:00:00 arrayd
  1363 ?            00:00:01 amd
  1420 ?            00:00:00 crond
  1490 ?            00:00:00 xfs
  1505 ?            00:00:00 sesdaemon
  1535 ?            00:00:01 sesdaemon
  1536 ?            00:00:00 sesdaemon
  1538 ?            00:00:00 sesdaemon
```

To monitor running processes, use the `top(1)` command. This command displays a sorted list of top CPU utilization processes.

The `vmstat(8)` command reports virtual memory statistics. It reports information about processes, memory, paging, block IO, traps, and CPU activity. For more information, see the `vmstat(8)` man page.

In the following `vmstat(8)` command, the `10` specifies a 10-second delay between updates.

```
uv44-sys:~ # vmstat 10
procs -----memory----- ---swap-- -----io----- --system-- -----cpu-----
 r  b   swpd   free   buff   cache   si   so    bi   bo    in   cs us sy id wa st
 2  0     0 235984032 418748 8649568    0    0    0    0    0    0  0  0  0 100  0  0
```

```
1 0      0 236054400 418748 8645216    0    0    0 4809 256729 3401 0 0 100 0 0
1 0      0 236188016 418748 8649904    0    0    0 448 256200 631 0 0 100 0 0
2 0      0 236202976 418748 8645104    0    0    0 341 256201 1117 0 0 100 0 0
1 0      0 236088720 418748 8592616    0    0    0 847 257104 6152 0 0 100 0 0
1 0      0 235990944 418748 8648460    0    0    0 240 257085 5960 0 0 100 0 0
1 0      0 236049568 418748 8645100    0    0    0 4849 256749 3604 0 0 100 0 0
```

Without the *delay* parameter, the output returns averages since the last reboot. Additional reports give information on a sampling period of length *delay*. The process and memory reports are instantaneous in either case.

The `iostat(1)` command monitors system input/output device loading by observing the time the devices are active, relative to their average transfer rates. You can use information from the `iostat` command to change system configuration information to better balance the input/output load between physical disks. For more information, see the `iostat(1)` man page.

In the following `iostat(1)` command, the `10` specifies a 10-second *interval* between updates

```
uv44-sys:~ # iostat
Linux 2.6.32.13-0.4.1.1559.0.PTF-default (uv44-sys) 10/18/2010 _x86_64_

avg-cpu:  %user   %nice %system %iowait  %steal   %idle
           0.18    0.00    0.04    0.02    0.00   99.77

Device:            tps    Blk_read/s    Blk_wrtn/s    Blk_read    Blk_wrtn
sda                 3.02         72.80         16.28       722432     161576
sdb                 0.01          0.34          0.00         3419         0
```

The `sar(1)` command returns the content of selected, cumulative activity counters in the operating system. Based on the values in the *count* and *interval* parameters, the command writes information *count* times spaced at the specified *interval*, which is in seconds. For more information, see the `sar(1)` man page. The following example shows the `sar(1)` command with a request for information about CPU 1, a count of 10, and an interval of 10:

```
uv44-sys:~ # sar -P 1 10 10
Linux 2.6.32-416.el6.x86_64 (harp34-sys)                09/19/2013                _x86_64_                (2

11:24:54 AM      CPU      %user      %nice      %system      %iowait      %steal      %idle
```

11:25:04 AM	1	0.20	0.00	0.10	0.00	0.00	99.70
11:25:14 AM	1	10.10	0.00	0.30	0.00	0.00	89.60
11:25:24 AM	1	99.70	0.00	0.30	0.00	0.00	0.00
11:25:34 AM	1	99.70	0.00	0.30	0.00	0.00	0.00
11:25:44 AM	1	8.99	0.00	0.60	0.00	0.00	90.41
11:25:54 AM	1	0.10	0.00	0.20	0.00	0.00	99.70
11:26:04 AM	1	38.70	0.00	0.10	0.00	0.00	61.20
11:26:14 AM	1	99.80	0.00	0.10	0.00	0.00	0.10
11:26:24 AM	1	80.42	0.00	0.70	0.00	0.00	18.88
11:26:34 AM	1	0.10	0.00	0.20	0.00	0.00	99.70
Average:	1	43.78	0.00	0.29	0.00	0.00	55.93

Memory Statistics and nodeinfo Command

nodeinfo(1) is a tool for monitoring per-node NUMA memory statistics on SGI UV systems. The nodeinfo tool reads /sys/devices/system/node/*/meminfo and /sys/devices/system/node/*/numastat on the local system to gather NUMA memory statistics.

Sample memory statistic from the nodeinfo(1) command are, as follows:

```
uv44-sys:~ # nodeinfo
Memory Statistics Tue Oct 26 12:01:58 2010
uv44-sys
----- Per Node KB ----- Preferred Alloc ----- -- Loc/Rem ---
node      Total      Free      Used      Dirty      Anon      Slab      hit      miss foreign interlv      local      remote
0      16757488      16277084      480404      52      34284      36288      20724      0      0      0      20720      4
1      16777216      16433988      343228      68      6772      17708      4477      0      0      0      3381      1096
2      16777216      16438568      338648      76      6908      12620      1804      0      0      0      709      1095
3      16760832      16429844      330988      56      2820      16836      1802      0      0      0      708      1094
4      16777216      16444408      332808      88      10124      13588      1517      0      0      0      417      1100
5      16760832      16430300      330532      72      1956      17304      4546      0      0      0      3453      1093
6      16777216      16430788      346428      36      3236      15292      3961      0      0      0      2864      1097
7      16760832      16435532      325300      44      1220      14800      3971      0      0      0      2877      1094
TOT      134148848      131320512      2828336      492      67320      144436      42802      0      0      0      35129      7673
Press "h" for help
```

From an interactive nodeinfo session, enter h for a help statement:

```
Display memory statistics by node.
q      quit
+      Increase starting node number. Used only if more nodes than will
```

4: Monitoring Tools

fit in the current window.
- Decrease starting node number. Used only if more nodes than will fit in the current window.
b Start output with node 0.
e Show highest node number.
k show sizes in KB.
m show sizes in MB.
p show sizes in pages.
t Change refresh rate.
A Show/Hide memory policy stats.
H Show/Hide hugepage info.
L Show/Hide LRU Queue stats.

Field definitions:

hit - page was allocated on the preferred node
miss - preferred node was full. Allocation occurred on THIS node by a process running on another node that was full

foreign - Preferred node was full. Had to allocate somewhere else.

interlv - allocation was for interleaved policy

local - page allocated on THIS node by a process running on THIS node
remote - page allocated on THIS node by a process running on ANOTHER node

(press any key to exit from help screen)

Data Placement Tools

This chapter describes data placement tools you can use on an SGI UV system. It covers the following topics:

- "Data Placement Tools Overview" on page 35
- "taskset Command" on page 39
- "dplace Command" on page 41
- "dlook Command" on page 50
- "omplace Command" on page 57
- "numactl Command" on page 41
- "Installing NUMA Tools" on page 58

Data Placement Tools Overview

On an symmetric multiprocessor (SMP) machine, all data is visible from all processors. Nonuniform Memory Access (NUMA) machines also have a shared address space. In both cases, there is a single shared memory space and a single operating system instance. However, in an SMP machine, each processor is functionally identical and has equal time access to every memory address. In contrast, a NUMA system has a shared address space, but the access time to memory vary over physical address ranges and between processing elements. The Intel Xeon 7500 series processor (Nehalem i7 architecture) is an example of NUMA architecture. Each processor has its own memory and can address the memory attached to another processor through the Quick Path Interconnect (QPI).

The SGI UV 1000 series is a family of multiprocessor distributed shared memory (DSM) computer systems that initially scale from 32 to 2560 Intel processor cores, as a cache-coherent single system image (SSI). The SGI UV 100 series is a family of multiprocessor distributed shared memory (DSM) computer systems that initially scale from 16 to 768 Intel processor cores as a cache-coherent SSI. The SGI UV 2000 series scales from 32 to 4,096 Intel processor cores, as a cache-coherent SSI.

In both UV series systems, there are two levels of NUMA: intranode managed by the QPI and internode managed through the HUB ASIC and NUMALink 5 or NUMALink 6. This section covers the following topics:

- "Distributed Shared Memory (DSM)" on page 36
- "ccNUMA Architecture" on page 36
- "Cache Coherency" on page 37
- "Non-uniform Memory Access (NUMA)" on page 37

Distributed Shared Memory (DSM)

In the SGI UV 2000, SGI UV 1000, or SGI UV 100 series systems, memory is physically distributed both within and among the IRU enclosures (compute/memory/I/O blades); however, it is accessible to and shared by all devices connected by NUMALink within the single-system image (SSI). This is to say that all components connected by NUMALink sharing a single Linux operating system, operate and share the memory fabric of the system. Memory latency is the amount of time required for a processor to retrieve data from memory. Memory latency is lowest when a processor accesses local memory. Note the following sub-types of memory within a system:

- If a processor accesses memory that is on a compute node blade, the memory is referred to as the node's local memory.
- If processors access memory located in other blade nodes within the IRU, (or other NUMALink IRUs) the memory is referred to as remote memory.
- The total memory within the NUMALink system is referred to as global memory.

ccNUMA Architecture

As the name implies, the cache-coherent non-uniform memory access (ccNUMA) architecture has two parts, cache coherency and nonuniform memory access, which are discussed in the sections that follow.

Cache Coherency

The SGI UV systems use caches to reduce memory latency. Although data exists in local or remote memory, copies of the data can exist in various processor caches throughout the system. Cache coherency keeps the cached copies consistent.

To keep the copies consistent, the ccNUMA architecture uses directory-based coherence protocol. In directory-based coherence protocol, each block of memory (64 bytes) has an entry in a table that is referred to as a directory. Like the blocks of memory that they represent, the directories are distributed among the compute/memory blade nodes. A block of memory is also referred to as a cache line.

Each directory entry indicates the state of the memory block that it represents. For example, when the block is not cached, it is in an unowned state. When only one processor has a copy of the memory block, it is in an exclusive state. And when more than one processor has a copy of the block, it is in a shared state; a bit vector indicates which caches may contain a copy.

When a processor modifies a block of data, the processors that have the same block of data in their caches must be notified of the modification. The SGI UV systems use an invalidation method to maintain cache coherence. The invalidation method purges all unmodified copies of the block of data, and the processor that wants to modify the block receives exclusive ownership of the block.

Non-uniform Memory Access (NUMA)

In DSM systems, memory is physically located at various distances from the processors. As a result, memory access times (latencies) are different or "non-uniform." For example, it takes less time for a processor blade to reference its locally installed memory than to reference remote memory.

Data Placement Practices

For cc-NUMA systems like the SGI UV systems, there is a performance penalty to access remote memory versus local memory. Because the Linux operating system has a tendency to migrate processes, the importance of using a placement tool becomes more apparent. Various data placement tools are described in this section.

Special optimization applies to SGI UV systems to exploit multiple paths to memory, as follows:

- By default, all pages are allocated with a “first touch” policy.
- The initialization loop, if executed serially, will get pages from single node.
- In the parallel loop, multiple processors will access that one memory.

Perform initialization in parallel, such that each processor initializes data that it is likely to access later for calculation.

Figure 5-1 on page 38, shows how to code to get good data placement.

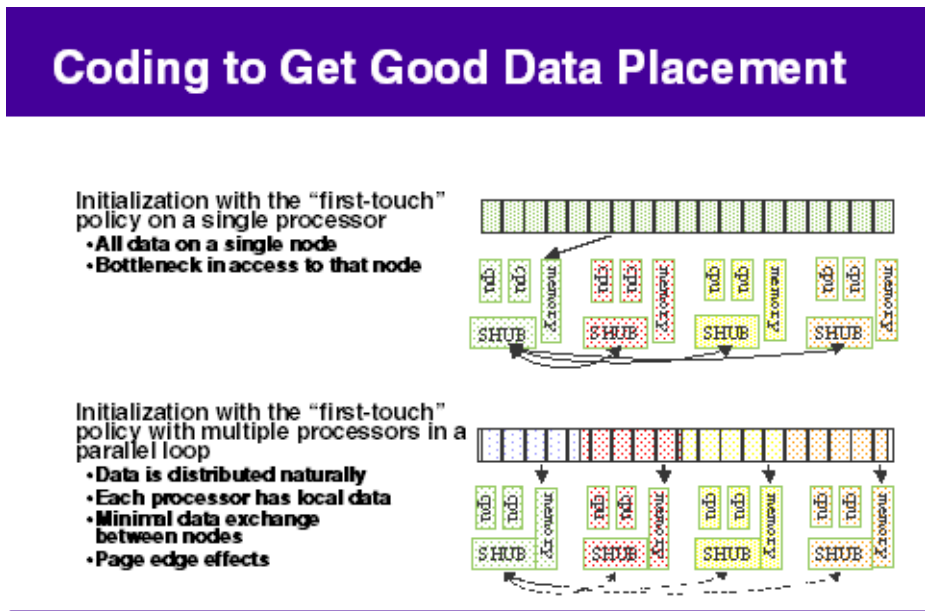


Figure 5-1 Coding to Get Good Data Placement

The data placement facilities include `taskset(1)`, `dplace(1)`, and `cpuset`. These tools are all built upon the `CpuMemSets` API. The `taskset` and `dplace` utilities enable you to avoid poor data locality caused by process or thread drift from CPU to CPU. `taskset` restricts execution to the listed set of CPUs (see the `taskset -c --cpu-list` option); however, processes are still free to move among listed CPUs.

`dplace` binds processes to specified CPUs in round-robin fashion; once pinned, they do not migrate, so you can use this for high performance and reproducibility of parallel codes. Cpusets are named subsets of system cpus/memories and are used extensively in batch environments. For more information about cpusets, see the *SGI Cpuset Software Guide*.

The following topics describe the data placement utilities:

- "taskset Command" on page 39
- "numactl Command" on page 41
- "dplace Command" on page 41
- "dlook Command" on page 50
- "omplace Command" on page 57
- "Installing NUMA Tools" on page 58
- "About cpusets" on page 58

taskset Command

The `taskset(1)` command retrieves or sets a CPU affinity of a process, as follows:

```
taskset [options] mask command [arg]...
taskset [options] -p [mask] pid
```

The `taskset` command sets or retrieves the CPU affinity of a running process or launches a new command with a given CPU affinity. CPU affinity is a scheduler property that "bonds" a process to a given set of CPUs on the system. The Linux scheduler will honor the given CPU affinity and the process will not run on any other CPUs. Note that the Linux scheduler also supports natural CPU affinity; the scheduler attempts to keep processes on the same CPU as long as practical for performance reasons. Therefore, forcing a specific CPU affinity is useful only in certain applications.

The CPU affinity is represented as a bitmask, with the lowest order bit corresponding to the first logical CPU and the highest order bit corresponding to the last logical CPU. Not all CPUs may exist on a given system but a mask may specify more CPUs than are present. A retrieved mask will reflect only the bits that correspond to CPUs physically on the system. If an invalid mask is given (that is, one that corresponds to

no valid CPUs on the current system) an error is returned. The masks are typically given in hexadecimal. For example:

```
0x00000001          is processor #0
0x00000003          is processors #0 and #1
0xFFFFFFFF          is all processors (#0 through #31)
```

When `taskset` returns, it is guaranteed that the given program has been scheduled to a legal CPU.

The `taskset` command does not pin a task to a specific CPU. It only restricts a task so that it does not run on any CPU that is not in the `cpulist`. For example, if you use `taskset` to launch an application that forks multiple tasks, it is possible that multiple tasks will initially be assigned to the same CPU even though there are idle CPUs that are in the `cpulist`. Scheduler load balancing software will eventually distribute the tasks so that CPU bound tasks run on different CPUs. However, the exact placement is not predictable and can vary from run-to-run. After the tasks are evenly distributed (assuming that happens), nothing prevents tasks from jumping to different CPUs. This can affect memory latency since pages that were node-local before the jump may be remote after the jump.

If you are running an MPI application, SGI recommends that you do not use the `taskset` command. The `taskset` command can pin the MPI shepherd process (which is a waste of a CPU) and then putting the remaining working MPI rank on one of the CPUs that already had some other rank running on it. Instead of `taskset`, SGI recommends using the `dplace(1)` (see "dplace Command" on page 41) or the environment variable `MPI_DSM_CPULIST`. The following example assumes a job running on eight CPUs. For example:

```
# mpirun -np 8 dplace -s1 -c10,11,16-21 myMPIapplication ...
To set MPI_DSM_CPULIST variable, perform a command similar to the following:
```

```
setenv MPI_DSM_CPULIST 10,11,16-21 mpirun -np 8 myMPIapplication ...
```

If they are using a batch scheduler that creates and destroys `cpusets` dynamically, you should use `MPI_DSM_DISTRIBUTE` environment variable instead of either `MPI_DSM_CPULIST` environment variable or the `dplace` command.

For more detailed information, see the `taskset(1)` man page.

To run an executable on CPU 1 (the `cpumask` for CPU 1 is `0x2`), perform the following:

```
# taskset 0x2 executable name
```

To move pid 14057 to CPU 0 (the cpumask for cpu 0 is 0x1), perform the following:

```
# taskset -p 0x1 14057
```

To run an MPI Abaqus/Std job on a UV system with eight CPUs, perform the following:

```
# taskset -c 8-15 ./runme < /dev/null &
```

The stdin is redirected to /dev/null to avoid a SIGTTIN signal for MPT applications.

The following example uses the `taskset` command to lock a given process to a particular CPU (CPU5) and then uses the `profile(1)` command to profile it. It then shows how to use `taskset` to move the process to another CPU (CPU3).

```
# taskset -p -c 5 16269
pid 16269's current affinity list: 0-15
pid 16269's new affinity list: 5
```

```
# taskset -p 16269 -c 3
pid 16269's current affinity list: 5
pid 16269's new affinity list: 3
```

numactl Command

The `numactl(8)` command runs processes with a specific NUMA scheduling or memory placement policy. The policy is set for command and inherited by all of its children. In addition it can set persistent policy for shared memory segments or files. For more information, see the `numactl(8)` man page.

dplace Command

You can use the `dlook(1)` and `dplace(1)` commands to improve the performance of processes running on your SGI nonuniform memory access (NUMA) machine. The following topics explain how to use the `dplace(1)` command to bind a related set of processes to specific CPUs or nodes to prevent process migration. This can improve the performance of your application since it increases the percentage of memory accesses that are local.

Using the `dplace` Command

The `dplace` command allows you to control the placement of a process onto specified CPUs. Its format is as follows:

```
dplace [-e] [-c cpu_numbers] [-s skip_count] [-n process_name] \  
      [-x skip_mask] [-r [l|b|t]] [-o log_file] [-v 1|2] \  
      command [command-args]  
dplace [-p placement_file] [-o log_file] command [command-args]  
dplace [-q] [-qq] [-qqq]
```

Scheduling and memory placement policies for the process are set up according to `dplace` command line arguments.

By default, memory is allocated to a process on the node on which the process is executing. If a process moves from node to node while it running, a higher percentage of memory references are made to remote nodes. Because remote accesses typically have higher access times, process performance can be diminished. CPU instruction pipelines also have to be reloaded.

You can use the `dplace` command to bind a related set of processes to specific CPUs or nodes to prevent process migrations. In some cases, this improves performance since a higher percentage of memory accesses are made to local nodes.

Processes always execute within a `CpuMemSet`. The `CpuMemSet` specifies the CPUs on which a process can execute. By default, processes usually execute in a `CpuMemSet` that contains all the CPUs in the system. For information on `CpusMemSets`, see the *SGI Cpuset Software Guide*.

The `dplace` command invokes an SGI kernel hook (module called `numatools`) to create a placement container consisting of all the CPUs (or a or a subset of CPUs) of a `cpuset`. The `dplace` process is placed in this container and by default is bound to the first CPU of the `cpuset` associated with the container. Then `dplace` invokes `exec` to execute the command.

The command executes within this placement container and remains bound to the first CPU of the container. As the command forks child processes, they inherit the container and are bound to the next available CPU of the container.

If you do not specify a placement file, `dplace` binds processes sequentially in a round-robin fashion to CPUs of the placement container. For example, if the current `cpuset` consists of physical CPUs 2, 3, 8, and 9, the first process launched by `dplace` is bound to CPU 2. The first child process forked by this process is bound to CPU 3, the next process (regardless of whether it is forked by parent or child) to 8, and so on.

If more processes are forked than there are CPUs in the cpuset, binding starts over with the first CPU in the cpuset.

For more information on `dplace(1)` and examples of how to use the command, see the `dplace(1)` man page.

The `dplace(1)` command accepts the following options:

- `-c cpu_numbers`: Specified as a list of cpus, optionally strided cpu ranges, or a striding pattern. Example: "`-c 1`", "`-c 2-4`", "`-c 1,4-8,3`", "`-c 2-8:3`", "`-c CS`", "`-c BT`". The specification "`-c 2-4`" is equivalent to "`-c 2,3,4`" and "`-c 2-8:3`" is equivalent to `2,5,8`. Ranges may also be specified in reverse order: "`-c 12-8`" is equivalent to `12,11,10,9,8`. CPU numbers are **NOT** physical cpu numbers. They are logical cpu number that are relative to the cpus that are in the set of allowed cpus as specified by the current cpuset. A cpu value of "x" (or "*"), in the argument list for `-c` option, indicates that binding should not be done for that process. "x" should be used only if the `-e` option is also used. Cpu numbers start at 0. For striding patterns any subset of the characters (B)lade, (S)ocket, (C)ore, (T)hread may be used and their ordering specifies the nesting of the iteration. For example "SC" means to iterate all the cores in a socket before moving to the next CPU socket, while "CB" means to pin to the first core of each blade, then the second core of every blade, etc. For best results, use the `-e` option when using stride patterns. If the `-c` option is not specified, all cpus of the current cpuset are available. The command itself (which is `exec'd` by `dplace`) is the first process to be placed by the `-c cpu_numbers`.
- `-e`: Exact placement. As processes are created, they are bound to cpus in the exact order that the cpus are specified in the cpu list. Cpu numbers may appear multiple times in the list. A cpu value of "x" indicates that binding should not be done for that process. If the end of the list is reached, binding starts over at the beginning of the list.
- `-o`: Write a trace file to `<log file>` that describes the placement actions that were made for each `fork`, `exec`, and so on. Each line contains a timestamp, process id:thread number, cpu that task was executing on, taskname | placement action. Works with version 2 only.
- `-s skip_count`: Skips the first `skip_count` processes before starting to place processes onto CPUs. This option is useful if the first `skip_count` processes are "shepherd" processes that are used only for launching the application. If `skip_count` is not specified, a default value of 0 is used.

- `-n process_name`: Only processes named `process_name` are placed. Other processes are ignored and are not explicitly bound to CPUs.

The `process_name` argument is the basename of the executable.

- `-r`: Specifies that text should be replicated on the node or nodes where the application is running. In some cases, replication will improve performance by reducing the need to make offnode memory references for code. The replication option applies to all programs placed by the `dplace` command. See the `dplace(5)` man page for additional information on text replication. The replication options are a string of one or more of the following characters:

l	Replicate library text
b	Replicate binary (<code>a.out</code>) text
t	Thread round-robin option

- `-x skip_mask`: Provides the ability to skip placement of processes. The `skip_mask` argument is a bitmask. If bit N of `skip_mask` is set, then the N+1th process that is forked is not placed. For example, setting the mask to 6 prevents the second and third processes from being placed. The first process (the process named by the command) will be assigned to the first CPU. The second and third processes are not placed. The fourth process is assigned to the second CPU, and so on. This option is useful for certain classes of threaded applications that spawn a few helper processes that typically do not use much CPU time.

Note: OpenMP with Intel applications should be placed using the `-x` option with a `skip_mask` of 2 (`-x2`). This could change in future versions of OpenMP. For applications compiled using the Native Posix Thread Library (NPTL), use the `-x2` option.

- `-v`: Provides the ability to run in version 1 or version 2 compatibility mode if the kernel support is available. If not specified, version 2 compatibility is selected. See COMPATIBILITY section of the `dplace(1)` man page for more details. Note: version 1 requires kernel support for PAGG.
- `-p placement_file`: Specifies a placement file that contains additional directives that are used to control process placement.
- `command [command-args]`: Specifies the command you want to place and its arguments.

- `-q`: Lists the global count of the number of active processes that have been placed (by `dplace`) on each CPU in the current cpuset. Note that CPU numbers are logical CPU numbers within the cpuset, **not** physical CPU numbers. If specified twice, lists the current `dplace` jobs that are running. If specified three times, lists the current `dplace` jobs and the tasks that are in each job.

Example 5-1 Using the `dplace` command with MPI Programs

You can use the `dplace` command to improve placement of MPI programs on NUMA systems and verify placement of certain data structures of a long running MPI program by running a command such as the following:

```
mpirun -np 64 /usr/bin/dplace -s1 -c 0-63 ./a.out
```

You can then use the `dlook(1)` command to verify placement of certain data structures of a long running MPI program by using the `dlook` command in another window on one of the slave thread PIDs to verify placement. For more information on using the `dlook` command, see "dlook Command" on page 50 and the `dlook(1)` man page.

Example 5-2 Using `dplace` command with OpenMP Programs

To run an OpenMP program on logical CPUs 4 through 7 within the current cpuset, perform the following:

```
%efc -o prog -openmp -O3 program.f
%setenv OMP_NUM_THREADS 4
%dplace -x6 -c4-7 ./prog
```

The `dplace(1)` command has a static load balancing feature so that you do not necessarily have to supply a CPU list. To place `prog1` on logical CPUs 0 through 3 and `prog2` on logical CPUs 4 through 7, perform the following:

```
%setenv OMP_NUM_THREADS 4
%dplace -x6 ./prog1 &
%dplace -x6 ./prog2 &
```

You can use the `dplace -q` command to display the static load information.

Example 5-3 Using the `dplace` command with Linux commands

The following examples assume that the command is executed from a shell running in a cpuset consisting of physical CPUs 8 through 15.

Command	Run Location
<code>dplace -c2 date</code>	Runs the <code>date</code> command on physical CPU 10.
<code>dplace make linux</code>	Runs <code>gcc</code> and related processes on physical CPUs 8 through 15.
<code>dplace -c0-4,6 make linux</code>	Runs <code>gcc</code> and related processes on physical CPUs 8 through 12 or 14.
<code>taskset 4,5,6,7 dplace app</code>	The <code>taskset</code> command restricts execution to physical CPUs 12 through 15. The <code>dplace</code> command sequentially binds processes to CPUs 12 through 15.

To use the `dplace` command accurately, you should know how your placed tasks are being created in terms of the `fork`, `exec`, and `pthread_create` calls. Determine whether each of these worker calls are an MPI rank task or are they groups of pthreads created by rank tasks? Here is an example of two MPI ranks, each creating three threads:

```
cat <<EOF > placefile
firsttask cpu=0
exec name=mpiapp cpu=1
fork name=mpiapp cpu=4-8:4 exact
thread name=mpiapp oncpu=4 cpu=5-7 exact thread name=mpiapp oncpu=8
cpu=9-11 exact EOF

# mpirun is placed on cpu 0 in this example # the root mpiapp is
placed on cpu 1 in this example

# or, if your version of dplace supports the "cpurel=" option:
# firsttask cpu=0
# fork name=mpiapp cpu=4-8:4 exact
# thread name=mpiapp oncpu=4 cpurel=1-3 exact

# create 2 rank tasks, each will pthread_create 3 more # ranks will be
on 4 and 8
# thread children on 5,6,7 9,10,11
dplace -p placefile mpirun -np 2 ~cpw/bin/mpiapp -P 3 -l

exit
```


You can use the debugger to determine if it is working. It should show two MPI rank applications, each with three pthreads, as follows:

```
>> pthreads | grep mpiapp
px *(task_struct *)e00002343c528000 17769 17769 17763 0 mpiapp
    member task: e000013817540000 17795 17769 17763 0 5 mpiapp
    member task: e000013473aa8000 17796 17769 17763 0 6 mpiapp
    member task: e000013817c68000 17798 17769 17763 0 mpiapp
px *(task_struct *)e0000234704f0000 17770 17770 17763 0 mpiapp
    member task: e000023466ed8000 17794 17770 17763 0 9 mpiapp
    member task: e00002384cce0000 17797 17770 17763 0 mpiapp
    member task: e00002342c448000 17799 17770 17763 0 mpiapp
```

And you can use the debugger, to see a root application, the parent of the two MPI rank applications, as follows:

```
>> ps | grep mpiapp
0xe00000340b300000 1139 17763 17729 1 0xc8000000 - mpiapp
0xe00002343c528000 1139 17769 17763 0 0xc8000040 - mpiapp
0xe0000234704f0000 1139 17770 17763 0 0xc8000040 8 mpiapp
```

Placed as specified:

```
>> oncpus e00002343c528000 e000013817540000 e000013473aa8000
>> e000013817c68000 e0
000234704f0000 e000023466ed8000 e00002384cce0000 e00002342c448000
task: 0xe00002343c528000 mpiapp cpus_allowed: 4
task: 0xe000013817540000 mpiapp cpus_allowed: 5
task: 0xe000013473aa8000 mpiapp cpus_allowed: 6
task: 0xe000013817c68000 mpiapp cpus_allowed: 7
task: 0xe0000234704f0000 mpiapp cpus_allowed: 8
task: 0xe000023466ed8000 mpiapp cpus_allowed: 9
task: 0xe00002384cce0000 mpiapp cpus_allowed: 10
task: 0xe00002342c448000 mpiapp cpus_allowed: 11
```

dplace for Compute Thread Placement Troubleshooting Case Study

This section describes common reasons why compute threads do not end up on unique processors when using commands such as `dplace(1)` or `profile.pl` (see "Profiling with PerfSuite" on page 16).

In the example that follows, a user used the `dplace -s1 -c0-15` command to bind 16 processes to run on 0-15 CPUs. However, output from the `top(1)` command shows only 13 CPUs running with CPUs 13, 14, and 15 still idle and CPUs 0, 1 and 2 are shared with 6 processes.

```
263 processes: 225 sleeping, 18 running, 3 zombie, 17 stopped
CPU states:  cpu    user   nice  system   irq  softirq  iowait   idle
              total 1265.6%   0.0%  28.8%   0.0%   11.2%   0.0%  291.2%

cpu00  100.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu01   90.1%   0.0%   0.0%   0.0%   0.0%   9.7%   0.0%   0.0%
cpu02   99.9%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu03   99.9%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu04  100.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu05  100.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu06  100.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu07   88.4%   0.0%  10.6%   0.0%   0.0%   0.8%   0.0%   0.0%
cpu08  100.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu09   99.9%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu10   99.9%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu11   88.1%   0.0%  11.2%   0.0%   0.0%   0.6%   0.0%   0.0%
cpu12   99.7%   0.0%   0.2%   0.0%   0.0%   0.0%   0.0%   0.0%
cpu13   0.0%   0.0%   2.5%   0.0%   0.0%   0.0%   0.0%  97.4%
cpu14   0.8%   0.0%   1.6%   0.0%   0.0%   0.0%   0.0%  97.5%
cpu15   0.0%   0.0%   2.4%   0.0%   0.0%   0.0%   0.0%  97.5%
Mem: 60134432k av, 15746912k used, 44387520k free, 0k shrd,
```

672k buff

351024k active, 13594288k inactive

Swap: 2559968k av, 0k used, 2559968k free
2652128k cached

PID	USER	PRI	NI	SIZE	RSS	SHARE	STAT	%CPU	%MEM	TIME	CPU	COMMAND
7653	ccao	25	0	115G	586M	114G	R	99.9	0.9	0:08	3	mocassin
7656	ccao	25	0	115G	586M	114G	R	99.9	0.9	0:08	6	mocassin
7654	ccao	25	0	115G	586M	114G	R	99.8	0.9	0:08	4	mocassin
7655	ccao	25	0	115G	586M	114G	R	99.8	0.9	0:08	5	mocassin
7658	ccao	25	0	115G	586M	114G	R	99.8	0.9	0:08	8	mocassin
7659	ccao	25	0	115G	586M	114G	R	99.8	0.9	0:08	9	mocassin
7660	ccao	25	0	115G	586M	114G	R	99.8	0.9	0:08	10	mocassin
7662	ccao	25	0	115G	586M	114G	R	99.7	0.9	0:08	12	mocassin
7657	ccao	25	0	115G	586M	114G	R	88.5	0.9	0:07	7	mocassin
7661	ccao	25	0	115G	586M	114G	R	88.3	0.9	0:07	11	mocassin
7649	ccao	25	0	115G	586M	114G	R	55.2	0.9	0:04	2	mocassin
7651	ccao	25	0	115G	586M	114G	R	54.1	0.9	0:03	1	mocassin
7650	ccao	25	0	115G	586M	114G	R	50.0	0.9	0:04	0	mocassin
7647	ccao	25	0	115G	586M	114G	R	49.8	0.9	0:03	0	mocassin
7652	ccao	25	0	115G	586M	114G	R	44.7	0.9	0:04	2	mocassin
7648	ccao	25	0	115G	586M	114G	R	35.9	0.9	0:03	1	mocassin

An application can start some threads executing for a very short time yet the threads still have taken a token in the CPU list. Then, when the compute threads are finally started, the list is exhausted and restarts from the beginning. Consequently, some threads end up sharing the same CPU. To bypass this, try to eliminate the "ghost" thread creation, as follows:

- Check for a call to the "system" function. This is often responsible for the placement failure due to unexpected thread creation.
- When all the compute processes have the same name, you can do this by issuing a command, such as the following:

```
dplace -c0-15 -n compute-process-name ...
```

- You can also run `dplace -e -c0-32` on 16 CPUs to understand the pattern of the thread creation. If by chance, this pattern is the same from one run to the other (unfortunately race between thread creation often occurs), you can find the right flag to `dplace`. For example, if you want to run on CPU 0-3, with `dplace -e -c0-16` and you see that threads are always placed on CPU 0, 1, 5, and 6, then `dplace -e -c0,1,x,x,x,2,3` or `dplace -x24 -c0-3` (24 = 11000, place the 2 first and skip 3 before placing) should place your threads correctly.

dlook Command

You can use the `dlook(1)` and `dplace(1)` commands to improve the performance of processes running on your SGI nonuniform memory access (NUMA) machine. The following topics explain how to use the `dlook(1)` command to find out where in memory the operating system is placing your application's pages and how much system and user CPU time it is consuming.

Using the dlook Command

The `dlook(1)` command allows you to display the memory map and CPU usage for a specified process as follows:

```
dlook [-a] [-p] [-h] [-l] [-n] [-o outfile] [-s secs] command [command-args]  
dlook [-a] [-p] [-h] [-l] [-n] [-o outfile] [-s secs] pid
```

For each page in the virtual address space of the process, `dlook(1)` prints the following information:

- The object that owns the page, such as a file, SYSV shared memory, a device driver, and so on.
- The type of page, such as random access memory (RAM), FETCHOP, IOSPACE, and so on.
- If the page type is RAM memory, the following information is displayed:
 - Memory attributes, such as, SHARED, DIRTY, and so on
 - The node on which the page is located
 - The physical address of the page (optional)

Two forms of the `dlook(1)` command are provided. In one form, `dlook` prints information about an existing process that is identified by a process ID (PID). To use this form of the command, you must be the owner of the process or be running with root privilege. In the other form, you use `dlook` on a command you are launching and thus are the owner.

The `dlook(1)` command accepts the following options:

- `-a`: Shows the physical addresses of each page in the address space.
- `-h`: Explicitly lists holes in the address space.
- `-l`: Shows libraries.
- `-p`: Show raw hardware page table entries.
- `-o`: Outputs to file name (*outfile*). If not specified, output is written to stdout.
- `-s`: Specifies a sample interval in seconds. Information about the process is displayed every second (*secs*) of CPU usage by the process.

An example for the `sleep` process with a PID of 191155 is as follows:

```
$ dlook 191155
```

```
Peek:  sleep
Pid: 191155 Fri Sep 27 17:14:01 2013
```

```
Process memory map:
 00400000-00406000 r-xp 00000000 08:08 262250 /bin/sleep
 [0000000000400000-0000000000401000]      1 page  on node
```

5: Data Placement Tools

```
4 MEMORY|SHARED
[000000000401000-0000000000402000]      1 page  on node
5 MEMORY|SHARED
[000000000403000-0000000000404000]      1 page  on node
7 MEMORY|SHARED
[000000000404000-0000000000405000]      1 page  on node
8 MEMORY|SHARED

00605000-00606000 rw-p 00005000 08:08 262250 /bin/sleep
[0000000000605000-0000000000606000]      1 page  on node
2 MEMORY|RW|DIRTY

00606000-00627000 rw-p 00000000 00:00 0 [heap]
[0000000000606000-0000000000608000]      2 pages  on node
2 MEMORY|RW|DIRTY

7ffff7dd8000-7ffff7ddd000 rw-p 00000000 00:00 0
[00007ffff7dd8000-00007ffff7dda000]      2 pages  on node
2 MEMORY|RW|DIRTY
[00007ffff7ddc000-00007ffff7ddd000]      1 page  on node
2 MEMORY|RW|DIRTY

7ffff7fde000-7ffff7fe1000 rw-p 00000000 00:00 0
[00007ffff7fde000-00007ffff7fe1000]      3 pages  on node
2 MEMORY|RW|DIRTY

7ffff7ffa000-7ffff7ffb000 rw-p 00000000 00:00 0
[00007ffff7ffa000-00007ffff7ffb000]      1 page  on node
2 MEMORY|RW|DIRTY

7ffff7ffb000-7ffff7ffc000 r-xp 00000000 00:00 0 [vdso]
[00007ffff7ffb000-00007ffff7ffc000]      1 page  on node
7 MEMORY|SHARED

7ffff7ffe000-7ffff7fff000 rw-p 00000000 00:00 0
[00007ffff7ffe000-00007ffff7fff000]      1 page  on node
2 MEMORY|RW|DIRTY

7fffffff0000-7fffffff0000 rw-p 00000000 00:00 0 [stack]
[00007fffffff0000-00007fffffff0000]      2 pages  on node
2 MEMORY|RW|DIRTY
```

```
ffffffff600000-ffffffff601000 r-xp 00000000 00:00 0 [vsyscall]
[ffffffff600000-ffffffff601000]      1 page   on node
0 MEMORY|DIRTY|RESERVED
```

The `dlook` command gives the name of the process (peek: `sleep`), the process ID, and time and date it was invoked. It provides total user and system CPU time in seconds for the process.

Under the heading **Process memory map**, the `dlook` command prints information about a process from the `/proc/pid/cpu` and `/proc/pid/maps` files. On the left, it shows the memory segment with the offsets below in decimal. In the middle of the output page, it shows the type of access, time of execution, the PID, and the object that owns the memory (in this case, `/lib/ld-2.2.4.so`). The characters `s` or `p` indicate whether the page is mapped as sharable (`s`) with other processes or is private (`p`). The right side of the output page shows the number of pages of memory consumed and on which nodes the pages reside. A page is 16,384 bytes.

The node numbers reported by the `dlook` command correspond to `Socket` numbers reported by the `cpumap(1)` command under the section `Processor Numbering` on `Socket(s)`. For more information, see the `cpumap(1)` command description in "Determining System Configuration" on page 9.

Dirty memory means that the memory has been modified by a user.

In the second form of the `dlook` command, you specify a command and optional command arguments. The `dlook` command issues an `exec` call on the command and passes the command arguments. When the process terminates, `dlook` prints information about the process, as shown in the following example:

dlook date

```
Thu Aug 22 10:39:20 CDT 2002
```

```
Exit: date
```

```
Pid: 4680      Thu Aug 22 10:39:20 2002
```

Process memory map:

```
2000000000030000-200000000003c000 rw-p 0000000000000000 00:00 0
```

5: Data Placement Tools

```

                [200000000030000-200000000003c000]          3 pages on node  3 MEMORY|DIRTY

20000000002dc000-20000000002e4000 rw-p 0000000000000000 00:00 0
                [20000000002dc000-20000000002e4000]          2 pages on node  3 MEMORY|DIRTY

2000000000324000-2000000000334000 rw-p 0000000000000000 00:00 0
                [2000000000324000-2000000000328000]          1 page  on node  3 MEMORY|DIRTY

4000000000000000-400000000000c000 r-xp 0000000000000000 04:03 9657220  /bin/date
                [4000000000000000-400000000000c000]          3 pages on node  1 MEMORY|SHARED

6000000000008000-6000000000010000 rw-p 0000000000008000 04:03 9657220  /bin/date
                [600000000000c000-6000000000010000]          1 page  on node  3 MEMORY|DIRTY

6000000000010000-6000000000014000 rwxp 0000000000000000 00:00 0
                [6000000000010000-6000000000014000]          1 page  on node  3 MEMORY|DIRTY

60000fff80000000-60000fff80004000 rw-p 0000000000000000 00:00 0
                [60000fff80000000-60000fff80004000]          1 page  on node  3 MEMORY|DIRTY

60000fffffff4000-60000ffffffc000 rwxp ffffffffcccc000 00:00 0
                [60000fffffff4000-60000ffffffc000]          2 pages on node  3 MEMORY|DIRTY

```

If you use the `dlook` command with the `-s secs` option, the information is sampled at regular intervals. The output for the command `dlook -s 5 sleep 50` is as follows:

```

Exit:  sleep
Pid:  5617      Thu Aug 22 11:16:05 2002

```

```

Process memory map:
2000000000030000-200000000003c000 rw-p 0000000000000000 00:00 0
                [2000000000030000-200000000003c000]          3 pages on node  3 MEMORY|DIRTY

2000000000134000-2000000000140000 rw-p 0000000000000000 00:00 0

20000000003a4000-20000000003a8000 rw-p 0000000000000000 00:00 0
                [20000000003a4000-20000000003a8000]          1 page  on node  3 MEMORY|DIRTY

20000000003e0000-20000000003ec000 rw-p 0000000000000000 00:00 0
                [20000000003e0000-20000000003ec000]          3 pages on node  3 MEMORY|DIRTY

```



```

4000000000000000-40000000000008000 r-xp 0000000000000000 04:03 9657225 /bin/sleep
    [4000000000000000-40000000000008000]                2 pages on node 3 MEMORY|SHARED

60000000000004000-60000000000008000 rw-p 00000000000004000 04:03 9657225 /bin/sleep
    [60000000000004000-60000000000008000]                1 page on node 3 MEMORY|DIRTY

60000000000008000-6000000000000c000 rwxp 0000000000000000 00:00 0
    [60000000000008000-6000000000000c000]                1 page on node 3 MEMORY|DIRTY

60000fff80000000-60000fff80004000 rw-p 0000000000000000 00:00 0
    [60000fff80000000-60000fff80004000]                1 page on node 3 MEMORY|DIRTY

60000fffffff4000-60000fffffff0000 rwxp ffffffff00000000 00:00 0
    [60000fffffff4000-60000fffffff0000]                2 pages on node 3 MEMORY|DIRTY
    
```

You can run a Message Passing Interface (MPI) job using the `mpirun` command and print the memory map for each thread, or redirect the output to a file, as follows:

Note: The output has been abbreviated to shorten the example and bold headings added for easier reading.

```
mpirun -np 8 dlook -o dlook.out ft.C.8
```

Contents of dlook.out:

Exit: ft.C.8

Pid: 2306 **Fri Aug 30 14:33:37 2002**

Process memory map:

```

2000000000030000-200000000003c000 rw-p 0000000000000000 00:00 0
    [2000000000030000-2000000000034000]                1 page on node 21 MEMORY|DIRTY
    [2000000000034000-200000000003c000]                2 pages on node 12 MEMORY|DIRTY|SHARED

2000000000044000-2000000000060000 rw-p 0000000000000000 00:00 0
    [2000000000044000-2000000000050000]                3 pages on node 12 MEMORY|DIRTY|SHARED
    
```

...

Exit: ft.C.8

Pid: 2310 **Fri Aug 30 14:33:37 2002**

5: Data Placement Tools

Process memory map:

```
2000000000030000-200000000003c000 rw-p 0000000000000000 00:00 0
  [2000000000030000-2000000000034000]          1 page on node 25 MEMORY|DIRTY
  [2000000000034000-200000000003c000]          2 pages on node 12 MEMORY|DIRTY|SHARED

2000000000044000-2000000000060000 rw-p 0000000000000000 00:00 0
  [2000000000044000-2000000000050000]          3 pages on node 12 MEMORY|DIRTY|SHARED
  [2000000000050000-2000000000054000]          1 page on node 25 MEMORY|DIRTY

...
```

Exit: ft.C.8

Pid: 2307 Fri Aug 30 14:33:37 2002

Process memory map:

```
2000000000030000-200000000003c000 rw-p 0000000000000000 00:00 0
  [2000000000030000-2000000000034000]          1 page on node 30 MEMORY|DIRTY
  [2000000000034000-200000000003c000]          2 pages on node 12 MEMORY|DIRTY|SHARED

2000000000044000-2000000000060000 rw-p 0000000000000000 00:00 0
  [2000000000044000-2000000000050000]          3 pages on node 12 MEMORY|DIRTY|SHARED
  [2000000000050000-2000000000054000]          1 page on node 30 MEMORY|DIRTY

...
```

Exit: ft.C.8

Pid: 2308 Fri Aug 30 14:33:37 2002

Process memory map:

```
2000000000030000-200000000003c000 rw-p 0000000000000000 00:00 0
  [2000000000030000-2000000000034000]          1 page on node 0  MEMORY|DIRTY
  [2000000000034000-200000000003c000]          2 pages on node 12 MEMORY|DIRTY|SHARED

2000000000044000-2000000000060000 rw-p 0000000000000000 00:00 0
  [2000000000044000-2000000000050000]          3 pages on node 12 MEMORY|DIRTY|SHARED
  [2000000000050000-2000000000054000]          1 page on node 0  MEMORY|DIRTY
```

...

For more information on the `dlook` command, see the `dlook` man page.

omplace Command

The `omplace(1)` command is a tool for controlling the placement of MPI processes and OpenMP threads.

The `omplace` command causes the successive threads in a hybrid MPI/OpenMP job to be placed on unique CPUs. The CPUs are assigned in order from the effective CPU list within the containing `cpuset`.

This command is a wrapper script for `dplace(1)` that can be used with MPI, OpenMP, pthreads, and hybrid MPI/OpenMP and MPI/pthreads codes. It generates the proper `dplace` placement file syntax automatically. It also supports some unique options like block-strided CPU lists.

The CPU placement is performed by dynamically generating a placement file and invoking `dplace` with the MPI job launch. For example, To run two MPI processes with four threads per process, and to display the generated placement file run a command similar to the following:

```
# mpirun -np 2 omplace -nt 4 -vv ./a.out
```

The threads would be place as follows:

```
rank 0 thread 0 on CPU 0
rank 0 thread 1 on CPU 1
rank 0 thread 2 on CPU 2
rank 0 thread 3 on CPU 3
rank 1 thread 0 on CPU 4
rank 1 thread 1 on CPU 5
rank 1 thread 2 on CPU 6
rank 1 thread 3 on CPU 7
```

For more information, see the `omplace(1)` man page and “Run-Time Tuning” chapter in the *Message Passing Toolkit (MPT) User’s Guide*.

Installing NUMA Tools

To use the `dlook(1)`, `dplace(1)`, and `topology(1)` commands, you must load the `numatools` kernel module. Perform the following steps:

1. To configure `numatools` kernel module to be started automatically during system startup, use the `chkconfig(8)` command as follows:

```
chkconfig --add numatools
```

2. To turn on `numatools`, enter the following command:

```
/etc/rc.d/init.d/numatools start
```

This step will be done automatically for subsequent system reboots when `numatools` are configured on by using the `chkconfig(8)` utility.

The following steps are required to disable `numatools`:

1. To turn off `numatools`, enter the following:

```
/etc/rc.d/init.d/numatools stop
```

2. To stop `numatools` from initiating after a system reboot, use the `chkconfig(8)` command as follows:

```
chkconfig --del numatools
```

About cpusets

The `cpuset` facility is a workload manager tool that permits a system administrator to restrict the number of processor and memory resources that a process or set of processes may use. A *cpuset* defines a list of CPUs and memory nodes. A process contained in a `cpuset` can execute only on the CPUs in that `cpuset` and can only allocate memory on the memory nodes in that `cpuset`. Essentially, `cpusets` provide you with a CPU and memory containers or *soft partitions* within which you can run sets of related tasks. Using `cpusets` on an SGI UV system improves cache locality and memory access times and can substantially improve an application's performance and runtime repeatability.

For information about `cpusets`, see *SGI Cpuset Software Guide*.

Performance Tuning

After analyzing your code to determine where performance bottlenecks are occurring, you can turn your attention to making your programs run their fastest. One way to do this is to use multiple CPUs in parallel processing mode. However, this should be the last step. The first step is to make your program run as efficiently as possible on a single processor system and then consider ways to use parallel processing.

Intel provides detailed application tuning information including information about the Intel Xeon processor 5500 at this location

<http://developer.intel.com/Assets/PDF/manual/248966.pdf> and specific tuning information tutorial for Nehalem (Intel Xeon 5500) at

<http://software.intel.com/sites/webinar/tuning-your-application-for-nehalem/>.

This chapter describes the process of tuning your application for a single processor system, and then tuning it for parallel processing. It also addresses how to improve the performance of floating-point programs and MPI applications. This chapter includes the following topics:

- "Single Processor Code Tuning" on page 59
- "Multiprocessor Code Tuning" on page 68
- "Understanding Parallel Speedup and Amdahl's Law" on page 75
- "Gustafson's Law" on page 80
- "Floating-point Program Performance" on page 81
- "About MPI Application Tuning" on page 81
- "Using Transparent Huge Pages (THPs) in MPI, SHMEM, and UPC Applications" on page 85
- "Enabling Huge Pages in MPI, SHMEM, and UPV Applications on Systems Without THP" on page 86

Single Processor Code Tuning

Several basic steps are used to tune performance of single-processor code:

- Get the expected answers and then tune performance. For details, see "Getting the Correct Results" on page 60.
- Use existing tuned code, such as that found in math libraries and scientific library packages. For details, see "Using Tuned Code" on page 61.
- Determine what needs tuning. For details, see "Determining Tuning Needs" on page 62.
- Use the compiler to do the work. For details, see "Using Compiler Options Where Possible" on page 62.
- Consider tuning cache performance. For details, see "Tuning the Cache Performance" on page 65.
- Set environment variables to enable higher-performance memory management mode. For details, see "Managing Memory" on page 67.

Getting the Correct Results

One of the first steps in performance tuning is to verify that the correct answers are being obtained. Once the correct answers are obtained, tuning can be done. You can verify answers by initially disabling specific optimizations and limiting default optimizations. This can be accomplished by using specific compiler options and by using debugging tools.

The following compiler options emphasize tracing and porting over performance:

- `-O`: the `-O0` option disables all optimization. The default is `-O2`.
- `-g`: the `-g` option preserves symbols for debugging. In the past, using `-g` automatically put down the optimization level. In Intel compiler today, you can use `-O3` with `-g`.
- `-fp-model`: the `-fp-model` option lets you specify the compiler rules for:
 - Value safety
 - Floating-point (FP) expression evaluation
 - FPU environment access
 - Precise FP exceptions
 - FP contractions

Default is `-fp-model fast=1`. Note that `-mp` option is an old flag replaced by `-fp-model`.

- `-r`; `-i`: the `-r8` and `-i8` options set default real, integer, and logical sizes to 8 bytes, which are useful for porting codes from Cray, Inc. systems. **This explicitly declares intrinsic and external library functions.**

Some debugging tools can also be used to verify that correct answers are being obtained. See "Debugging Tools" on page 18 for more details.

Managing Heap Corruption Problems

You can use environment variables to check for heap corruption problems in programs that use `glibc malloc/free` dynamic memory management routines.

Set the `MALLOC_CHECK_` environment variable to 1 to print diagnostic messages or to 2 to abort immediately when heap corruption is detected.

Overruns and underruns are circumstances where an access to an array is outside the declared boundary of the array. Underruns and overruns cannot be simultaneously detected. The default behavior is to place inaccessible pages immediately after allocated memory.

Using Tuned Code

Where possible, use code that has already been tuned for optimum hardware performance.

The following mathematical functions should be used where possible to help obtain best results:

- MKL: Intel's Math Kernel Library. This library includes BLAS, LAPACK, and FFT routines.
- VML: the Vector Math Library, available as part of the MKL package (`libmkl_vml_itp.so`).
- Standard Math library

Standard math library functions are provided with the Intel compiler's `libimf.a` file. If the `-lm` option is specified, `glibc libm` routines are linked in first.

Documentation is available for MKL and VML at the following website:

http://intel.com/software/products/perflib/index.htm?iid=ipp_home+software_libraries&

Determining Tuning Needs

Use the following tools to determine what points in your code might benefit from tuning:

- `time`: Use this command to obtain an overview of user, system, and elapsed time.
- `gprof`: Use this tool to obtain an execution profile of your program (a `pcsamp` profile). Use the `-p` compiler option to enable `gprof` use.
- `VTune`: This is an Intel performance monitoring tool. You can run it directly on your SGI UV system. The Linux server/Windows client is useful when you are working on a remote system.
- `psrun` is a `PerfSuite` (see <http://perfsuite.ncsa.uiuc.edu/>) command-line utility that allows you to take performance measurements of unmodified executables. `psrun` takes as input a configuration XML document that describes the desired measurement.

For information about other performance analysis tools, see Chapter 3, "Performance Analysis and Debugging" on page 9.

Using Compiler Options Where Possible

Several compiler options can be used to optimize performance. For a short summary of `ifort` or `icc` options, use the `-help` option on the compiler command line. Use the `-dryrun` option to show the driver tool commands that `ifort` or `icc` generate. This option does not actually compile.

Use the following options to help tune performance:

- `-ftz`: Flushes underflow to zero to avoid kernel traps. Enabled by default at `-O3` optimization.
- `-fno-alias`: Assumes no pointer aliasing. Pointer aliasing can create uncertainty about the possibility that two unrelated names might refer to the identical memory; because of this uncertainty, the compiler will assume that any two pointers can point to the same location in memory. This can remove optimization opportunities, particularly for loops.

Other aliasing options include `-ansi_alias` and `-fno_fnalias`. Note that incorrect alias assertions may generate incorrect code.

- `-ip`: Generates single file, interprocedural optimization; `-ipo` generates multifile, interprocedural optimization.

Most compiler optimizations work within a single procedure (like a function or a subroutine) at a time. This **intra**-procedural focus restricts optimization possibilities because a compiler is forced to make worst-case assumptions about the possible effects of a procedure. By using **inter**-procedural analysis, more than a single procedure is analyzed at once and code is optimized. It performs two passes through the code and requires more compile time.

- `-O3`: Enables `-O2` optimizations plus more aggressive optimizations, including loop transformation and prefetching. *Loop transformation* are found in a transformation file created by the compiler; you can examine this file to see what suggested changes have been made to loops. *Prefetch instructions* allow data to be moved into the cache before their use. A prefetch instruction is similar to a load instruction.

Note that Level 3 optimization may not improve performance for all programs.

- `-opt_report`: Generates an optimization report and places it in the file specified in `-opt_report_file`.
- `-override_limits`: This is an undocumented option that sometimes allows the compiler to continue optimizing when it has hit an internal limit.
- `-prof_gen` and `-prof_use`: Generates and uses profiling information. These options require a three-step compilation process:
 1. Compile with proper instrumentation using `-prof_gen`.
 2. Run the program on one or more training datasets.
 3. Compile with `-prof_use`, which uses the profile information from the training run.
- `-S`: Compiles and generates an assembly listing in the `.s` files and does not link. The assembly listing can be used in conjunction with the output generated by the `-opt_report` option to try to determine how well the compiler is optimizing loops.

- `-vec-report`: For information specific to the vectorizer. Intel Xeon 7500 series processors can perform short vector operations which provides a powerful performance boost.
- `-fast`: equivalent to writing: `-ipo -O3 -no-prec-div -static -xHost`
- `-xHost`: Can generate instructions for the highest instruction set and processor available on the compilation host.
- Specific processor architecture to compile for: `-xSSE4.2` for Nehalem EP/EX, for example. Useful if compiling in a different system than an SGI UV.
- `-xSSE4.2`: Can generate Intel® SSE4 Efficient Accelerated String and Text Processing instructions supported by Intel® Core i7 processors. Can generate Intel® SSE4 Vectorizing Compiler and Media Accelerator, Intel® SSSE3, SSE3, SSE2, and SSE instructions and it can optimize for the Intel® Core™ processor family.

Another important feature of new Intel compilers is the Source Checker, which is enabled using the flag `-diag-enable + options`. The source checker is a compiler feature that provides advanced diagnostics based on detailed analysis of source code. It performs static global analysis to find errors in software that go undetected by the compiler itself. It is a general source code analysis tool that provides an additional diagnostic capability to help you debug your programs. You can use source code analysis options to detect potential errors in your compiled code.

Specific processor architecture to compile for: `-xSSE4.2` for Nehalem EP/EX Useful if compiling in a different system than an SGI UV.

`-xSSE4.2`: Can generate Intel® SSE4 Efficient Accelerated String and Text Processing instructions supported by Intel® Core i7 processors. Can generate Intel® SSE4 Vectorizing Compiler and Media Accelerator, Intel® SSSE3, SSE3, SSE2, and SSE instructions and it can optimize for the Intel® Core™ processor family. Another important feature of new Intel compilers is the Source Checker, which is enabled using the flag `-diag-enable + options`. The source checker is a compiler feature that provides advanced diagnostics based on detailed analysis of source code. It performs static global analysis to find errors in software that go undetected by the compiler itself. general source code analysis tool that provides an additional diagnostic capability to help you debug your programs. You can use source code analysis options to detect potential errors in your compiled code including the following:

- Incorrect usage of OpenMP directives
- Inconsistent object declarations in different program units

- Boundary violations
- Uninitialized memory
- Memory corruptions
- Memory Leaks
- Incorrect usage of pointers and allocatable arrays
- Dead code and redundant executions
- Typographical errors or uninitialized variables
- Dangerous usage of unchecked input

Source checker analysis performs a general overview check of a program for all possible values simultaneously. This is in contrast to run-time checking tools that execute a program with a fixed set of values for input variables; such checking tools cannot easily check all edge effects. By not using a fixed set of input values, the source checker analysis can check for all possible corner cases. In fact, you do not need to run the program for Source Checker, the analysis is performed at compilation time. Only requirement is a successful compilation. Important caveat: Limitations of Source Checker Analysis: Since the source checker does not perform full interpretation of analyzed programs, it can generate so called false-positive messages. This is a fundamental difference between the compiler and source checker generated errors; in the case of the source checker, you decide whether the generated error is legitimate and needs to be fixed.

Tuning the Cache Performance

The processor cache stores recently-used information in a place where it can be accessed quickly. This discussion uses the following terms:

- A *cache line* is the minimum unit of transfer from next-higher cache into this one.
- A *cache hit* is reference to a cache line which is present in the cache.
- A *cache miss* is reference to a cache line which is not present in this cache level and must be retrieved from a higher cache (or memory or swap space).
- The *hit time* is the time to access the upper level of the memory hierarchy, which includes the time needed to determine whether the access is a hit or a miss.

- A *miss penalty* is the time to replace a block in the upper level with the corresponding block from the lower level, plus the time to deliver this block to the processor. The time to access the next level in the hierarchy is the major component of the miss penalty.

There are several actions you can take to help tune cache performance:

- Avoid large power-of-2 (and multiples thereof) strides and dimensions that cause *cache thrashing*. Cache thrashing occurs when multiple memory accesses require use of the same cache line. This can lead to an unnecessary number of cache misses.

To prevent cache thrashing, redimension your vectors so that the size is not a power of two. Space the vectors out in memory so that concurrently accessed elements map to different locations in the cache. When working with two-dimensional arrays, make the leading dimension an odd number; for multidimensional arrays, change two or more dimensions to an odd number.

Consider the following example: a cache in the hierarchy has a size of 256 KB (or 65536 4-byte words). A Fortran program contains the following loop:

```
real data(655360,24)
...
do i=1,23
  do j=1,655360
    diff=diff+data(j,i)-data(j,i+1)
  enddo
enddo
```

The two accesses to `data` are separated in memory by 655360×4 bytes, which is a simple multiple of the cache size; they consequently load to the same location in the cache. Because both data items cannot simultaneously coexist in that cache location, a pattern of replace on reload occurs that considerably reduces performance.

- Use a memory stride of 1 wherever possible. A loop over an array should access array elements from adjacent memory addresses. When the loop iterates through memory by consecutive word addresses, it uses every word of every cache line in sequence and does not return to a cache line after finishing it.

If memory strides other than 1 are used, cache lines could be loaded multiple times if an array is too large to be held in memory at one time.

- Cache bank conflicts can occur if there are two accesses to the same 16-byte-wide bank at the same time.

A maximum of four performance monitoring events can be counted simultaneously.

- Group together data that is used at the same time and do not use vectors in your code, if possible. If elements that are used in one loop iteration are contiguous in memory, it can reduce traffic to the cache and fewer cache lines will be fetched for each iteration of the loop.
- Try to avoid the use of temporary arrays and minimize data copies.

Managing Memory

Nonuniform memory access (NUMA) uses hardware with memory and peripherals distributed among many CPUs. This allows scalability for a shared memory system but a side effect is the time it takes for a CPU to access a memory location. Because memory access times are nonuniform, program optimization is not always straightforward.

Codes that frequently allocate and deallocate memory through `glibc malloc/free` calls may accrue significant system time due to memory management overhead. By default, `glibc` strives for system-wide memory efficiency at the expense of performance.

In compilers up to and including version 7.1.x, to enable the higher-performance memory management mode, set the following environment variables:

```
% setenv MALLOC_TRIM_THRESHOLD_ -1
% setenv MALLOC_MMAP_MAX_ 0
```

Because allocations in `ifort` using the `malloc` intrinsic use the `glibc malloc` internally, these environment variables are also applicable in Fortran codes using, for example, Cray pointers with `malloc/free`. But they do not work for Fortran 90 allocatable arrays, which are managed directly through Fortran library calls and placed in the stack instead of the heap. The example, above, applies only to the `cs` shell and the `tcsh` shell.

Memory Use Strategies

This section describes some general memory use strategies, as follows:

- Register reuse: do a lot of work on the same data before working on new data
- Cache reuse: the program is much more efficient if all of the data and instructions fit in cache; if not, try to use what is in cache a lot before using anything that is not in cache.
- Data locality: try to access data that is near each other in memory before data that is far.
- I/O efficiency: do a bunch of I/O all at once rather than a little bit at a time; do not mix calculations and I/O.

Memory Hierarchy Latencies

Programmers tend to think of memory as a flat, random access storage device. It is critical to understand that memory is a hierarchy to get good performance. Memory latency differs within the hierarchy. Performance is affected by where the data resides. Registers: 0 cycles latency (cycle = 1/freq) L1 cache: 1 cycle L2 cache: 5-6 cycles L3 cache: 12-17 cycles Main memory: 130-1000+ cycles. CPUs which are waiting for memory are not doing useful work. Software should be "hierarchy-aware" to achieve best performance:

- Perform as many operations as possible on data in registers
- Perform as many operations as possible on data in the cache(s)
- Keep data uses spatially and temporally local
- Consider temporal locality and spatial locality

Memory hierarchies take advantage of temporal locality by keeping more recently accessed data items closer to the processor. Memory hierarchies take advantage of spatial locality by moving contiguous words in memory to upper levels of the hierarchy.

Multiprocessor Code Tuning

Before beginning any multiprocessor tuning, first perform single processor tuning. This can often obtain good results in multiprocessor codes also. For details, see "Single Processor Code Tuning" on page 59.

Multiprocessor tuning consists of the following major steps:

- Determine what parts of your code can be parallelized. For background information, see "Data Decomposition" on page 69.
- Choose the parallelization methodology for your code. For details, see "Parallelizing Your Code" on page 70.
- Analyze your code to make sure it is parallelizing properly. For details, see Chapter 3, "Performance Analysis and Debugging" on page 9.
- Check to determine if false sharing exists. False sharing refers to OpenMP, not MPI. For details, see "Fixing False Sharing" on page 73.
- Tune for data placement. For details, see "Using `dplace` and `taskset`" on page 74.
- Use environment variables to assist with tuning. For details, see "Environment Variables for Performance Tuning" on page 74.

Data Decomposition

In order to efficiently use multiple processors on a system, tasks have to be found that can be performed at the same time. There are two basic methods of defining these tasks:

- Functional parallelism

Functional parallelism is achieved when different processors perform different functions. This is a known approach for programmers trained in modular programming. Disadvantages to this approach include the difficulties of defining functions as the number of processors grow and finding functions that use an equivalent amount of CPU power. This approach may also require large amounts of synchronization and data movement.

- Data parallelism

Data parallelism is achieved when different processors perform the same function on different parts of the data. This approach takes advantage of the large cumulative memory. One requirement of this approach, though, is that the problem domain be *decomposed*. There are two steps in data parallelism:

1. Data decomposition

Data decomposition is breaking up the data and mapping data to processors. Data can be broken up explicitly by the programmer by using message

passing (with MPI) and data passing (using the SHMEM library routines) or can be done implicitly using compiler-based MP directives to find parallelism in implicitly decomposed data.

There are advantages and disadvantages to implicit and explicit data decomposition:

- **Implicit decomposition advantages:** No data resizing is needed; all synchronization is handled by the compiler; the source code is easier to develop and is portable to other systems with OpenMP or High Performance Fortran (HPF) support.
- **Implicit decomposition disadvantages:** The data communication is hidden by the user
- **Explicit decomposition advantages:** The programmer has full control over insertion of communication and synchronization calls; the source code is portable to other systems; code performance can be better than implicitly parallelized codes.
- **Explicit decomposition disadvantages:** Harder to program; the source code is harder to read and the code is longer (typically 40% more).

2. The final step is to divide the work among processors.

Parallelizing Your Code

The first step in multiprocessor performance tuning is to choose the parallelization methodology that you want to use for tuning. This section discusses those options in more detail.

You should first determine the amount of code that is parallelized. Use the following formula to calculate the amount of code that is parallelized:

$$p = N(T(1) - T(N)) / T(1)(N - 1)$$

In this equation, $T(1)$ is the time the code runs on a single CPU and $T(N)$ is the time it runs on N CPUs. Speedup is defined as $T(1)/T(N)$.

If $speedup/N$ is less than 50% (that is, $N > (2-p)/(1-p)$), stop using more CPUs and tune for better scalability.

CPU activity can be displayed with the `top` or `vmstat` commands or accessed by using the Performance Co-Pilot tools (for example, `pmval`)

`kernel.percpu.cpu.user`) or by using the Performance Co-Pilot visualization tools `pmchart`.

Next you should focus on a parallelization methodology, as discussed in the following subsections.

Use MPT

You can use the Message Passing Interface (MPI) from the SGI Message Passing Toolkit (MPT). MPI is optimized and more scalable for SGI UV series systems than generic MPI libraries. It takes advantage of the SGI UV architecture and SGI Linux NUMA features. MPT is included with the SGI MPI, part of the SGI Performance Suite software.

Use the `-lmpi` compiler option to use MPI. For a list of environment variables that are supported, see the `mpi` man page.

`MPIO_DIRECT_READ` and `MPIO_DIRECT_WRITE` are supported under Linux for local XFS filesystems in SGI MPT version 1.6.1 and beyond.

MPI provides the MPI-2 standard MPI I/O functions that provide file read and write capabilities. A number of environment variables are available to tune MPI I/O performance. See the `mpi_io(3)` man page for a description of these environment variables.

Performance tuning for MPI applications is described in more detail in Chapter 6 of the *Message Passing Toolkit (MPT) User's Guide*.

Use OpenMP

OpenMP is a shared memory multiprocessing API, which standardizes existing practice. It is scalable for fine or coarse grain parallelism with an emphasis on performance. It exploits the strengths of shared memory and is directive-based. The OpenMP implementation also contains library calls and environment variables.

To use OpenMP directives with C, C++, or Fortran codes, you can use the following compiler options:

- `ifort -openmp` or `icc -openmp`: These options use the OpenMP front-end that is built into the Intel compilers. The latest Intel compiler OpenMP runtime name is `libiomp5.so`. The latest Intel compiler also supports the GNU OpenMP OpenMP library as an either/or option (not to be mixed-and-matched with the Intel version).

For details about OpenMP usage see the OpenMP standard, available at <http://www.openmp.org/specs>.

OpenMP Nested Parallelism

This section describes OpenMP nested parallelism. For additional information, see the `dplace(1)` man page.

The following Open MP nested parallelism output shows 2 primary threads and 4 secondary threads, called master/nested:

```
% cat place_nested
firsttask cpu=0
thread name=a.out oncpu=0 cpu=4 noplac=1 exact onetime thread name=a.out oncpu=0
cpu=1-3 exact thread name=a.out oncpu=4 cpu=5-7 exact

% dplace -p place_nested a.out
Master thread 0 running on cpu 0
Master thread 1 running on cpu 4
Nested thread 0 of master 0 gets task 0 on cpu 0 Nested thread 1 of master 0 gets task 1 on cpu 1
Nested thread 2 of master 0 gets task 2 on cpu 2 Nested thread 3 of master 0 gets task 3 on cpu 3
Nested thread 0 of master 1 gets task 0 on cpu 4 Nested thread 1 of master 1 gets task 1 on cpu 5
Nested thread 2 of master 1 gets task 2 on cpu 6 Nested thread 3 of master 1 gets task 3 on cpu 7
```

Use Compiler Options

Use the compiler to invoke automatic parallelization. Use the `-parallel` and `-par_report` option to the `ifort` or `icc` compiler. These options show which loops were parallelized and the reasons why some loops were not parallelized. If a source file contains many loops, it might be necessary to add the `-override_limits` flag to enable automatic parallelization. The code generated by `-parallel` is based on the OpenMP API; the standard OpenMP environment variables and Intel extensions apply.

There are some limitations to automatic parallelization:

- For Fortran codes, only `DO` loops are analyzed
- For C/C++ codes, only `for` loops using explicit array notation or those using pointer increment notation are analyzed. In addition, `for` loops using pointer arithmetic notation are not analyzed nor are `while` or `do/while` loops. The compiler also does not check for blocks of code that can be run in parallel.

Identifying Parallel Opportunities in Existing Code

Another parallelization optimization technique is to identify loops that have a potential for parallelism, such as the following:

- Loops without data dependencies; a *data dependency conflict* occurs when a loop has results from one loop pass that are needed in future passes of the same loop.
- Loops with data dependencies because of temporary variables, reductions, nested loops, or function calls or subroutines.

Loops that do not have a potential for parallelism are those with premature exits, too few iterations, or those where the programming effort to avoid data dependencies is too great.

Fixing False Sharing

If the parallel version of your program is slower than the serial version, false sharing might be occurring. False sharing occurs when two or more data items that appear not to be accessed by different threads in a shared memory application correspond to the same cache line in the processor data caches. If two threads executing on different CPUs modify the same cache line, the cache line cannot remain resident and correct in both CPUs, and the hardware must move the cache line through the memory subsystem to retain coherency. This causes performance degradation and reduction in the scalability of the application. If the data items are only read, not written, the cache line remains in a shared state on all of the CPUs concerned. False sharing can occur when different threads modify adjacent elements in a shared array. When two CPUs share the same cache line of an array and the cache is decomposed, the boundaries of the chunks split at the cache line.

You can use the following methods to verify that false sharing is happening:

- Use the performance monitor to look at output from `pfmon` and the `BUS_MEM_READ_BRIL_SELF` and `BUS_RD_INVALID_ALL_HITM` events.
- Use `pfmon` to check `DEAR` events to track common cache lines.
- Use the Performance Co-Pilot `pmshub` utility to monitor cache traffic and CPU utilization.

If false sharing is a problem, try the following solutions:

- Use the hardware counter to run a profile that monitors storage to shared cache lines. This will show the location of the problem.

- Revise data structures or algorithms.
- Check shared data, static variables, common blocks, and private and public variables in shared objects.
- Use critical regions to identify the part of the code that has the problem.

Using `dplace` and `taskset`

The `dplace` command binds processes to specified CPUs in a round-robin fashion. Once bound to a process, they do not migrate. `dplace` numbering is done in the context of the current CPU memory set. See Chapter 4, "Monitoring Tools" on page 23 for details about `dplace`.

The `taskset(1)` command and the `numactl(8)` command restrict execution to the listed set of CPUs. Processes are still free to move among listed CPUs.

Environment Variables for Performance Tuning

You can use several different environment variables to assist in performance tuning. For details about environment variables used to control the behavior of MPI, see the `mpi(1)` man page.

Several OpenMP environment variables can affect the actions of the OpenMP library. For example, some environment variables control the behavior of threads in the application when they have no work to perform or are waiting for other threads to arrive at a synchronization semantic; other variables can specify how the OpenMP library schedules iterations of a loop across threads. The following environment variables are part of the OpenMP standard:

- `OMP_NUM_THREADS` (The default is the number of CPUs in the system.)
- `OMP_SCHEDULE` (The default is `static`.)
- `OMP_DYNAMIC` (The default is `false`.)
- `OMP_NESTED` (The default is `false`.)

In addition to the preceding environment variables, Intel provides several OpenMP extensions, two of which are provided through the use of the `KMP_LIBRARY` variable.

The `KMP_LIBRARY` variable sets the run-time execution mode, as follows:

- If set to `serial`, single-processor execution is used.
- If set to `throughput`, CPUs yield to other processes when waiting for work. This is the default and is intended to provide good overall system performance in a multiuser environment.
- If set to `turnaround`, worker threads do not yield while waiting for work. Setting `KMP_LIBRARY` to `turnaround` may improve the performance of benchmarks run on dedicated systems, where multiple users are not contending for CPU resources.

If your program gets a segmentation fault immediately upon execution, you may need to increase `KMP_STACKSIZE`. This is the private stack size for threads. The default is 4 MB. You may also need to increase your shell `stacksize` limit.

Understanding Parallel Speedup and Amdahl's Law

There are two ways to obtain the use of multiple CPUs. You can take a conventional program in C, C++, or Fortran, and have the compiler find the parallelism that is implicit in the code.

You can write your source code to use explicit parallelism, stating in the source code which parts of the program are to execute asynchronously, and how the parts are to coordinate with each other.

When your program runs on more than one CPU, its total run time should be less. But how much less? What are the limits on the speedup? That is, if you apply 16 CPUs to the program, should it finish in 1/16th the elapsed time?

This section covers the following topics:

- "Adding CPUs to Shorten Execution Time" on page 76
- "Understanding Parallel Speedup" on page 76
- "Understanding Amdahl's Law" on page 77
- "Calculating the Parallel Fraction of a Program" on page 78
- "Predicting Execution Time with n CPUs" on page 79

Adding CPUs to Shorten Execution Time

You can distribute the work your program does over multiple CPUs. However, there is always some part of the program's logic that has to be executed serially, by a single CPU. This sets the lower limit on program run time.

Suppose there is one loop in which the program spends 50% of the execution time. If you can divide the iterations of this loop so that half of them are done in one CPU while the other half are done at the same time in a different CPU, the whole loop can be finished in half the time. The result: a 25% reduction in program execution time.

The mathematical treatment of these ideas is called Amdahl's law, for computer pioneer Gene Amdahl, who formalized it. There are two basic limits to the speedup you can achieve by parallel execution:

- The fraction of the program that can be run in parallel, p , is never 100%.
- Because of hardware constraints, after a certain point, there is less and less benefit from each added CPU.

Tuning for parallel execution comes down to doing the best that you are able to do within these two limits. You strive to increase the parallel fraction, p , because in some cases even a small change in p (from 0.8 to 0.85, for example) makes a dramatic change in the effectiveness of added CPUs.

Then you work to ensure that each added CPU does a full CPU's work and does not interfere with the work of other CPUs. In the SGI UV architectures this means:

- Spreading the workload equally among the CPUs
- Eliminating false sharing and other types of memory contention between CPUs
- Making sure that the data used by each CPU are located in a memory near that CPU's node

Understanding Parallel Speedup

If half the iterations of a loop are performed on one CPU, and the other half run at the same time on a second CPU, the whole loop should complete in half the time. For example, consider the typical C loop in Example 6-1.

Example 6-1 Typical C Loop

```
for (j=0; j<MAX; ++j) {  
    z[j] = a[j]*b[j];  
}
```

The compiler can automatically distribute such a loop over n CPUs (with n decided at run time based on the available hardware), so that each CPU performs MAX/n iterations.

The speedup gained from applying n CPUs, $Speedup(n)$, is the ratio of the one-CPU execution time to the n -CPU execution time: $Speedup(n) = T(1) \div T(n)$. If you measure the one-CPU execution time of a program at 100 seconds, and the program runs in 60 seconds with two CPUs, $Speedup(2) = 100 \div 60 = 1.67$.

This number captures the improvement from adding hardware. $T(n)$ ought to be less than $T(1)$; if it is not, adding CPUs has made the program slower, and something is wrong. So $Speedup(n)$ should be a number greater than 1.0, and the greater it is, the better. Intuitively you might hope that the speedup would be equal to the number of CPUs (twice as many CPUs, half the time) but this ideal can seldom be achieved.

Understanding Superlinear Speedup

You expect $Speedup(n)$ to be less than n , reflecting the fact that not all parts of a program benefit from parallel execution. However, it is possible, in rare situations, for $Speedup(n)$ to be larger than n . When the program has been sped up by more than the increase of CPUs it is known as *superlinear speedup*.

A superlinear speedup does not really result from parallel execution. It comes about because each CPU is now working on a smaller set of memory. The problem data handled by any one CPU fits better in cache, so each CPU executes faster than the single CPU could do. A superlinear speedup is welcome, but it indicates that the sequential program was being held back by cache effects.

Understanding Amdahl's Law

There are always parts of a program that you cannot make parallel, where code must run serially. For example, consider the loop. Some amount of code is devoted to setting up the loop, allocating the work between CPUs. This housekeeping must be done serially. Then comes parallel execution of the loop body, with all CPUs running concurrently. At the end of the loop comes more housekeeping that must be done

serially; for example, if n does not divide `MAX` evenly, one CPU must execute the few iterations that are left over.

The serial parts of the program cannot be speeded up by concurrency. Let p be the fraction of the program's code that can be made parallel (p is always a fraction less than 1.0.) The remaining fraction $(1-p)$ of the code must run serially. In practical cases, p ranges from 0.2 to 0.99.

The potential speedup for a program is proportional to p divided by the CPUs you can apply, plus the remaining serial part, $1-p$. As an equation, this appears as Example 6-2.

Example 6-2 Amdahl's law: $Speedup(n)$ Given p

$$Speedup(n) = \frac{1}{(p/n) + (1-p)}$$

Suppose $p = 0.8$; then $Speedup(2) = 1 / (0.4 + 0.2) = 1.67$, and $Speedup(4) = 1 / (0.2 + 0.2) = 2.5$. The maximum possible speedup (if you could apply an infinite number of CPUs) would be $1 / (1-p)$. The fraction p has a strong effect on the possible speedup.

The reward for parallelization is small unless p is substantial (at least 0.8); or to put the point another way, the reward for increasing p is great no matter how many CPUs you have. The more CPUs you have, the more benefit you get from increasing p . Using only four CPUs, you need only $p = 0.75$ to get half the ideal speedup. With eight CPUs, you need $p = 0.85$ to get half the ideal speedup.

There is a slightly more sophisticated version of Amdahl's law which includes communication overhead, showing also that if the program has no serial part that as we increase the number of cores the amount of computation per core diminishes and the communication overhead (unless there is not communication and we have trivial parallelization) increases, also diminishing the efficiency of the code and the speedup. The equation is: $Speedup(n) = n / (1 + a*(n-1) + n*(tc/ts))$ Where: n : number of processes a : the fraction of the given task not dividable into concurrent subtasks ts : time to execute the task in a single processor tc : communication overhead If $a=0$ and $tc=0$ (no serial part and no communications) like in a trivial parallelization program, you will get linear speedup.

Calculating the Parallel Fraction of a Program

You do not have to guess at the value of p for a given program. Measure the execution times $T(1)$ and $T(2)$ to calculate a measured $Speedup(2) = T(1) / T(2)$. The

Amdahl's law equation can be rearranged to yield p when $Speedup(2)$ is known, as in Example 6-3.

Example 6-3 Amdahl's law: p Given $Speedup(2)$

$$p = \frac{2}{1} * \frac{SpeedUp(2) - 1}{SpeedUp(2)}$$

Suppose you measure $T(1) = 188$ seconds and $T(2) = 104$ seconds.

$$SpeedUp(2) = 188/104 = 1.81$$

$$p = 2 * ((1.81-1)/1.81) = 2*(0.81/1.81) = 0.895$$

In some cases, the $Speedup(2) = T(1)/T(2)$ is a value greater than 2; in other words, a superlinear speedup ("Understanding Superlinear Speedup" on page 77). When this occurs, the formula in Example 6-3 returns a value of p greater than 1.0, which is clearly not useful. In this case you need to calculate p from two other more realistic timings, for example $T(2)$ and $T(3)$. The general formula for p is shown in Example 6-4, where n and m are the two CPU counts whose speedups are known, $n > m$.

Example 6-4 Amdahl's Law: p Given $Speedup(n)$ and $Speedup(m)$

$$p = \frac{Speedup(n) - Speedup(m)}{(1 - 1/n)*Speedup(n) - (1 - 1/m)*Speedup(m)}$$

Predicting Execution Time with n CPUs

You can use the calculated value of p to extrapolate the potential speedup with higher numbers of CPUs. The following example shows the expected time with four CPUs, if $p=0.895$ and $T(1)=188$ seconds:

$$Speedup(4) = 1 / ((0.895/4) + (1-0.895)) = 3.04$$

$$T(4) = T(1)/Speedup(4) = 188/3.04 = 61.8$$

The calculation can be made routine using the computer by creating a script that automates the calculations and extrapolates run times.

These calculations are independent of most programming issues such as language, library, or programming model. They are not independent of hardware issues, because Amdahl's law assumes that all CPUs are equal. At some level of parallelism, adding a CPU no longer affects run time in a linear way. For example, on some

architectures, cache-friendly codes scale closely with Amdahl's law up to the maximum number of CPUs, but scaling of memory intensive applications slows as the system bus approaches saturation. When the bus bandwidth limit is reached, the actual speedup is less than predicted.

Gustafson's Law

Gustafson's law proposes that programmers set the size of problems to use the available equipment to solve problems within a practical fixed time. Therefore, if faster (more parallel) equipment is available, larger problems can be solved in the same time. Amdahl's law is based on fixed workload or fixed problem size. It implies that the sequential part of a program does not change with respect to machine size (for example, the number of processors). However, the parallel part is evenly distributed by n processors. The effect of Gustafson's law was to shift research goals to select or reformulate problems so that solving a larger problem in the same amount of time would be possible. In particular, the law redefines efficiency as a need to minimize the sequential part of a program, even if it increases the total amount of computation. The bottom line is that by running larger problems, it is hoped that the bulk of the calculation will increase faster than the serial part of the program, allowing for better scaling. There is a slightly more sophisticated version of Amdahl's law which includes communication overhead, showing also that if the program has no serial part that as we increase the number of cores the amount of computation per core diminishes and the communication overhead (unless there is not communication and you have trivial parallelization) increases, also diminishing the efficiency of the code and the speedup. The equation is:

$$\text{Speedup}(n) = n / (1 + a * (n-1) + n * (tc/ts))$$

The preceding equation uses the following variables:

- n : number of processes
- a : the fraction of the given task not dividable into concurrent subtasks
- ts : time to execute the task in a single processor
- tc : communication overhead

If $a=0$ and $tc=0$, which indicated no serial part and no communications, like in a trivial parallelization program, the result is linear speedup.

Floating-point Program Performance

Certain floating-point programs experience slowdowns due to excessive floating point traps called Floating-Point Software Assist (FPSWA).

This happens when the hardware cannot complete a floating point operation and requests help (emulation) from software. This happens, for instance, with denormals numbers.

The symptoms are a slower than normal execution, FPSWA message in the system log (run `dmesg`). The average cost of a FPSWA fault is quite high around 1000 cycles/fault.

By default, the kernel prints a message similar to the following in the system log:

```
foo(7716): floating-point assist fault at ip 4000000000200e1
        isr 0000020000000008
```

The kernel throttles the message in order to avoid flooding the console.

It is possible to control the behavior of the kernel on FPSWA faults using the `prctl(1)` command. In particular, it is possible to get a signal delivered at the first FPSWA. It is also possible to silence the console message.

About MPI Application Tuning

When you design your MPI application, make sure to include the following in your design:

- The pinning of MPI processes to CPUs
- The isolating of multiple MPI jobs onto different sets of sockets and Hubs

You can achieve this design by configuring a batch scheduler to create a `cpuset` for every MPI job. MPI pins its processes to the sequential list of logical processors within the containing `cpuset` by default, but you can control and alter the pinning pattern using `MPI_DSM_CPULIST`. For more information about these programming practices, see the following:

- The `MPI_DSM_CPULIST` discussion in the *Message Passing Toolkit (MPT) User's Guide*.
- The `omplace(1)` and `dplace(1)` man pages.

- The *SGI Cpuset Software Guide*.

MPI Application Communication on SGI Hardware

On an SGI UV system, the following two transfer methods facilitate MPI communication between processes:

- Shared memory
- The global reference unit (GRU), which is part of the SGI UV Hub ASIC

The SGI UV series systems use a scalable nonuniform memory access (NUMA) architecture to allow the use of thousands of processors and terabytes of RAM in a single Linux operating system instance. As in other large shared memory systems, memory is distributed to processor sockets and accesses to memory are cache coherent. Unlike other systems, SGI UV systems use a network of Hub ASICs connected over NUMALink to scale to more sockets than any other x86-64 system, with excellent performance out of the box for most applications.

When running on SGI UV systems with SGI's Message Passing Toolkit (MPT), applications can attain higher bandwidth and lower latency for MPI calls than when running on more conventional distributed memory clusters. However, knowing your SGI UV system's NUMA topology and the performance constraints that it imposes can still help you extract peak performance. For more information about the SGI UV hub, SGI UV compute blades, Intel QPI and SGI NUMALink, see your SGI UV hardware system user guide.

The MPI library chooses the transfer method depending on internal heuristics, the type of MPI communication that is involved, and some user-tunable variables. When using the GRU to transfer data and messages, the MPI library uses the GRU resources it allocates via the GRU resource allocator, which divides up the available GRU resources. It allocates buffer space and control blocks between the logical processors being used by the MPI job.

MPI Job Problems and Application Design

The MPI library chooses buffer sizes and communication algorithms in an attempt to deliver the best performance automatically to a wide variety of MPI applications, but user tuning might be needed to improve performance. The following are some application performance problems and some ways that you might be able to improve MPI performance:

- Primary HyperThreads are idle.

Most high performance computing MPI programs run best when they use only one HyperThread per core. When an SGI UV system has multiple HyperThreads per core, logical CPUs are numbered such that primary HyperThreads are the high half of the logical CPU numbers. Therefore, the task of scheduling only on the additional HyperThreads may be accomplished by scheduling MPI jobs as if only half the full number exists, leaving the high logical CPUs idle.

You can use the `cpumap(1)` command to determine if cores have multiple HyperThreads on your SGI UV system. The command's output includes the following:

- The number of physical and logical processors
- Whether Hyperthreading is enabled
- How shared processors are paired

If an MPI job uses only half of the available logical CPUs, set `GRU_RESOURCE_FACTOR` to 2 so that the MPI processes can use all the available GRU resources on a hub rather than reserving some of them for the idle HyperThreads. For more information about GRU resource tuning, see the `gru_resource(3)` man page.

- Message bandwidth is inadequate.

Use either huge pages or transparent huge pages (THP) to ensure that your application obtains optimal message bandwidth.

To specify the use of hugepages, use the `MPI_HUGEPAGE_HEAP_SPACE` environment variable. The `MPI_HUGEPAGE_HEAP_SPACE` environment variable defines the minimum amount of heap space that each MPI process can allocate using huge pages. For information about this environment variable, see the `MPI(1)` man page.

To use THPs, see "Using Transparent Huge Pages (THPs) in MPI, SHMEM, and UPC Applications" on page 85.

Some programs transfer large messages via the `MPI_Send` function. To enable unbuffered, single-copy transport in these cases, you can set `MPI_BUFFER_MAX` to 0. For information about the `MPI_BUFFER_MAX` environment variable, see the `MPI(1)` man page.

- MPI small or near messages are very frequent.

For small fabric hop counts, shared memory message delivery is faster than GRU messages. To deliver all messages within an SGI UV host via shared memory, set `MPI_SHARED_NEIGHBORHOOD` to `host`. For more information, see the `MPI(1)` man page.

- Memory allocations are nonlocal.

MPI application processes normally perform best if their local memory is allocated near the socket assigned to use it. This cannot happen if memory on that socket is exhausted by the application or by other system consumption, for example, file buffer cache. Use the `nodeinfo(1)` command to view memory consumption on the nodes assigned to your job, and use `bcfree(1)` to clear out excessive file buffer cache. PBS Professional batch scheduler installations can be configured to issue `bcfree` commands in the job prologue.

For information about PBS Professional, including the availability of scripts, see the PBS Professional documentation and the `bcfree(1)` man page.

MPI Performance Tools

SGI supports several MPI performance tools. You can use the following tools to enhance or troubleshoot MPI program performance:

- **MPInside.** MPInside is an MPI profiling tool that can help you optimize your MPI application. The tool provides information about data transferred between ranks, both in terms of speed and quantity. For information about MPInside, see one of the following:
 - The `MPInside(1)` man page.
 - The *MPInside Reference Guide*.
- **SGI Perfboost.** SGI PerfBoost uses a wrapper library to run applications compiled against other MPI implementations under the SGI Message Passing Toolkit (MPT) product on SGI platforms. The PerfBoost software allows you to run SGI MPT which is a version of MPI optimized for SGI large, shared-memory systems and can take advantage of the UV Hub.

For more information, see the *Message Passing Toolkit (MPT) User's Guide* available at <http://docs.sgi.com>.

- **SGI PerfCatcher.** SGI PerfCatcher uses a wrapper library to return MPI and SHMEM function profiling information. The information returned includes

percent CPU time, total time spent per function, message sizes, and load imbalances. For more information, see the following:

- The `perfcatch(1)` man page
- The *Message Passing Toolkit (MPT) User Guide*

Using Transparent Huge Pages (THPs) in MPI, SHMEM, and UPC Applications

On SGI UV systems, THP is important because it contributes to attaining the best GRU-based data transfer bandwidth in Message Passing Interface (MPI), SHMEM, and Unified Parallel C (UPC) programs. On newer kernels, the THP feature is enabled by default. If THP is disabled on your SGI UV system, see "Enabling Huge Pages in MPI, SHMEM, and UPV Applications on Systems Without THP" on page 86.

The THP feature, however, can affect the performance of some OpenMP threaded applications in a negative way. For certain OpenMP applications, some threads in some shared data structures might be forced to make more nonlocal references because the application assumes a smaller, 4-KB page size.

The THP feature affects users in the following ways:

- Administrators:

To activate the THP feature on a system-wide basis, write the keyword `always` to the following file:

```
/sys/kernel/mm/transparent_hugepage/enabled
```

To disable THP, type the following command:

```
echo never > /sys/kernel/mm/transparent_hugepage/enabled
```

- MPI application programmers:

To determine whether THP is enabled on your system, type the following command and note the output:

```
cat /sys/kernel/transparent_hugepage/enabled
```

The output is as follows on a system for which THP is enabled:

```
[always] never
```

In the output, the bracket characters ([]) appear around the keyword that is in effect.

Enabling Huge Pages in MPI, SHMEM, and UPV Applications on Systems Without THP

If the THP capability is disabled on your SGI UV system, you can use the `MPI_HUGEPAGE_HEAP_SPACE` environment variable and the `MPT_HUGEPAGE_CONFIG` command to create huge pages.

The `MPT_HUGEPAGE_CONFIG` command configures the system to allow huge pages to be available to MPT's memory allocation interceptors. The `MPI_HUGEPAGE_HEAP_SPACE` environment variable enables an application to use the huge pages reserved by the `MPT_HUGEPAGE_CONFIG` command.

For more information, see the `MPI_HUGEPAGE_HEAP_SPACE` environment variable on the `MPI(1)` man page, or see the `mpt_hugepage_config(1)` man page.

Flexible File I/O

Flexible File I/O (FFIO) provides a mechanism for improving the file I/O performance of existing applications without having to resort to source code changes, that is, the current executable remains unchanged. Knowledge of source code is not required, but some knowledge of how the source and the application software work can help you better interpret and optimize FFIO results. To take advantage of FFIO, all you need to do is to set some environment variables before running your application. This chapter covers the following topics:

- "FFIO Operation" on page 87
- "Environment Variables" on page 88
- "Simple Examples" on page 89
- "Multithreading Considerations" on page 92
- "Application Examples " on page 93
- "Event Tracing " on page 94
- "System Information and Issues " on page 94

FFIO Operation

The FFIO subsystem allows you to define one or more additional I/O buffer caches for specific files to augment the Linux kernel I/O buffer cache. The FFIO subsystem then manages this buffer cache for you. In order to accomplish this, FFIO intercepts standard I/O calls like open, read, and write, and replaces them with FFIO equivalent routines. These routines route I/O requests through the FFIO subsystem which utilizes the user defined FFIO buffer cache. FFIO can bypass the Linux kernel I/O buffer cache by communicating with the disk subsystem via direct I/O. This gives you precise control over cache I/O characteristics and allows for more efficient I/O requests. For example, doing direct I/O in large chunks (say 16 megabytes) allows the FFIO cache to amortize disk access. All file buffering occurs in user space when FFIO is used with direct I/O enabled. This differs from the Linux buffer cache mechanism which requires a context switch in order to buffer data in kernel memory. Avoiding this kind of overhead, helps FFIO to scale efficiently. Another important distinction is that FFIO allows you to create an I/O buffer cache dedicated to a specific application.

The Linux kernel, on the other hand, has to manage all the jobs on the entire system with a single I/O buffer cache. As a result, FFIO typically outperforms the Linux kernel buffer cache when it comes to I/O intensive throughput.

Environment Variables

There are only two environment variables that you need to set in order to use FFIO. They are `LD_PRELOAD` and `FF_IO_OPTS`.

In order to enable FFIO to trap standard I/O calls, you must set the `LD_PRELOAD` environment variable.

For SGI systems, perform the following:

```
# export LD_PRELOAD="/usr/lib64/libFFIO.so"
```

The `LD_PRELOAD` software is a Linux feature that instructs the linker to preload the indicated shared libraries. In this case, `libFFIO.so` is preloaded and provides the routines which replace the standard I/O calls. An application that is not dynamically linked with the `glibc` library will not work with FFIO, since the standard I/O calls will not be intercepted. To disable FFIO, perform the following:

```
unset LD_PRELOAD
```

The FFIO buffer cache is managed by the `FF_IO_OPTS` environment variable. The syntax for setting this variable can be quite complex. A simple method for defining this variable is as follows:

```
# export FF_IO_OPTS 'string(eie.direct.mbytes:size:num:lead:share:stride:0)'
```

You can use the following parameters with the `FF_IO_OPTS` environment variable:

<i>string</i>	Matches the names of files that can use the buffer cache.
<i>size</i>	Number of 4k blocks in each page of the I/O buffer cache.
<i>num</i>	Number of pages in the I/O buffer cache.
<i>lead</i>	The maximum number of read-ahead pages.
<i>share</i>	A value of 1 means a shared cache, 0 means private.

stride

Note that the number after the `stride` parameter is always 0.

The following example shows a command that creates a shared buffer cache of 128 pages where each page is 16 megabytes (that is, 4096*4k). The cache has a lead of six pages and uses a stride of one, as follows:

```
setenv FF_IO_OPTS 'test*(eie.direct.mbytes:4096:128:6:1:1:0)'
```

Each time the application opens a file, the FFIO code checks the file name to see if it matches the string supplied by `FF_IO_OPTS`. The file's path name is not considered when checking for a match against the string. So in the example supplied above, file names like `/tmp/test16` and `/var/tmp/testit` would both be a match.

More complicated usages of `FF_IO_OPTS` are built upon this simpler version. For example, multiple types of file names can share the same cache, as follows:

```
setenv FF_IO_OPTS 'output* test*(eie.direct.mbytes:4096:128:6:1:1:0)'
```

Multiple caches may also be specified with `FF_IO_OPTS`. In the example that follows, files of the form `output*` and `test*` share a 128 page cache of 16 megabyte pages. The file `special42` has a 256 page private cache of 32 megabyte pages, as follows:

```
setenv FF_IO_OPTS 'output* test*(eie.direct.mbytes:4096:128:6:1:1:0) special42(eie.direct.mbytes:8192:256:6:0:1:0)'
```

Additional parameters can be added to `FF_IO_OPTS` to create feedback that is sent to standard output. Examples of doing this diagnostic output will be presented in the following section.

Simple Examples

This section walks you through some simple examples using FFIO.

Assume that `LD_PRELOAD` is set for the correct library and `FF_IO_OPTS` is defined, as follows:

```
setenv FF_IO_OPTS 'test*(eie.direct.mbytes:4096:128:6:1:1:0)'
```

This example uses a small C program called `fio` that reads four megabyte chunks from a file for 100 iterations. When the program runs it produces output, as follows:

```
./fio -n 100 /build/testit
Reading 4194304 bytes 100 times to /build/testit
Total time = 7.383761
```

Throughput = 56.804439 MB/sec

It can be difficult to tell what FFIIO may or may not be doing even with a simple program such as shown above. A summary of the FFIIO operations that occurred can be directed to standard output by making a simple addition to FF_IO_OPTS, as follows:

```
setenv FF_IO_OPTS 'test*(eie.direct.mbytes:4096:128:6:1:1:0, event.summary.mbytes.notrace )'
```

This new setting for FF_IO_OPTS generates the following summary on standard output when the program is run:

```
./fio -n 100 /build/testit
Reading 4194304 bytes 100 times to /build/testit
Total time = 7.383761
Throughput = 56.804439 MB/sec
```

```
event_close(testit)  eie <-->syscall  (496 mbytes)/( 8.72 s)= 56.85 mbytes/s
oflags=0x0000000000004042=RDWR+CREAT+DIRECT
sector size =4096(bytes)
cblks =0  cbits =0x0000000000000000
current file size =512 mbytes  high water file size =512 mbytes
```

function	times called	wall time	all hidden	mbytes requested	mbytes delivered	min request	max request	avg request
open	1	0.00						
read	2	0.61		32	32	16	16	16
reada	29	0.01	0	464	464	16	16	16
fcntl								
recall								
reada	29	8.11						
other	5	0.00						
flush	1	0.00						
close	1	0.00						

Two synchronous reads of 16 megabytes each were issued (for a total of 32 megabytes) and 29 asynchronous reads (reada) were also issued (for a total of 464 megabytes). Additional diagnostic information can be generated by specifying the .diag modifier, as follows:

```
setenv FF_IO_OPTS 'test*(eie.direct.diag.mbytes:4096:128:6:1:1:0 )'
```

The `.diag` modifier may also be used in conjunction with `.event.summary`, the two operate independently from one another, as follows:

```
setenv FF_IO_OPTS 'test*(eie.diag.direct.mbytes:4096:128:6:1:1:0, event.summary.mbytes.notrace )'
```

An example of the diagnostic output generated when just the `.diag` modifier is used is, as follows:

```
./fio -n 100 /build/testit
Reading 4194304 bytes 100 times to /build/testit
Total time = 7.383761
Throughput = 56.804439 MB/sec

eie_close EIE final stats for file /build/testit
eie_close Used shared eie cache 1
eie_close 128 mem pages of 4096 blocks (4096 sectors), max_lead = 6 pages
eie_close advance reads used/started :      23/29      79.31%   (1.78 seconds wasted)
eie_close write hits/total           :           0/0      0.00%
eie_close read hits/total            :           98/100  98.00%
eie_close mbytes transferred      parent --> eie --> child      sync      async
eie_close                          0              0              0          0
eie_close                          400            496            2          29 (0,0)
eie_close                          parent <-- eie <-- child

eie_close EIE stats for Shared cache 1
eie_close 128 mem pages of 4096 blocks
eie_close advance reads used/started :      23/29      79.31%   (0.00 seconds wasted)
eie_close write hits/total           :           0/0      0.00%
eie_close read hits/total            :           98/100  98.00%
eie_close mbytes transferred      parent --> eie --> child      sync      async
eie_close                          0              0              0          0
eie_close                          400            496            2          29 (0,0)
```

Information is listed for both the file and the cache. An mbytes transferred example is shown below:

```
eie_close mbytes transferred      parent --> eie --> child      sync      async
eie_close                          0              0              0          0
eie_close                          400            496            2          29 (0,0)
```

The last two lines are for write and read operations, respectively. Only for very simple I/O patterns, the difference between (parent -> eie) and (eie -> child) read statistics

can be explained by the number of read aheads. For random reads of a large file over a long period of time, this is not the case. All write operations count as `async`.

Multithreading Considerations

FFIO will work with applications that use MPI for parallel processing. An MPI job assigns each thread a number or rank. The master thread has rank 0, while the remaining threads (called slave threads) have ranks from 1 to N-1 where N is the total number of threads in the MPI job. It is important to consider that the threads comprising an MPI job do not (necessarily) have access to each others address space. As a result, there is no way for the different MPI threads to share the same FFIO cache. By default, each thread defines a separate FFIO cache based on the parameters defined by `FF_IO_OPTS`.

Having each MPI thread define a separate FFIO cache based on a single environment variable (`FF_IO_OPTS`) can waste a lot of memory. Fortunately, FFIO provides a mechanism that allows the user to specify a different FFIO cache for each MPI thread via the following environment variables:

```
setenv FF_IO_OPTS_RANK0 'result*(eie.direct.mbytes:4096:512:6:1:1:0)'  
setenv FF_IO_OPTS_RANK1 'output*(eie.direct.mbytes:1024:128:6:1:1:0)'  
setenv FF_IO_OPTS_RANK2 'input*(eie.direct.mbytes:2048:64:6:1:1:0)'  
.  
.  
.  
setenv FF_IO_OPTS_RANKN-1 ... (N = number of threads).
```

Each rank environment variable is set using the exact same syntax as `FF_IO_OPTS` and each defines a distinct cache for the corresponding MPI rank. If the cache is designated shared, all files within the same ranking thread will use the same cache. FFIO works with SGI MPI, HP MPI, and LAM MPI. In order to work with MPI applications, FFIO needs to determine the rank of callers by invoking the `mpi_comm_rank()` MPI library routine. Therefore, FFIO needs to determine the location of the MPI library used by the application. This is accomplished by having the user set one, and only one, of the following environment variables:

```
setenv SGI_MPI /usr/lib  
or  
setenv LAM_MPI *see below  
or  
setenv HP_MPI *see below
```

*LAM and HP MPIs are usually distributed via a third party application. The precise paths to the LAM and the HP MPI libraries are application dependent. Please refer to the application installation guide to find the correct path.

In order to use the rank functionality, both the MPI and `FF_IO_OPTS_RANK0` environment variables must be set. If either variable is not set, then the MPI threads all use `FF_IO_OPTS`. If both the MPI and the `FF_IO_OPTS_RANK0` variables are defined but, for example, `FF_IO_OPTS_RANK2` is undefined, all rank 2 files would generate a no match with FFIO. This means that none of the rank 2 files would be cached by FFIO. In this case, the software does not default to `FF_IO_OPTS`.

Fortran and C/C++ applications that use the `pthread`s interface will create threads that share the same address space. These threads can all make use of the single FFIO cache defined by `FF_IO_OPTS`.

Application Examples

FFIO has been deployed successfully with several HPC applications such as Nastran and Abaqus. In a recent customer benchmark, an eight-way Abaqus throughput job ran approximately twice as fast when FFIO was used. The FFIO cache used 16 megabyte pages (that is, `page_size = 4096`) and the cache size was 8.0 gigabytes. As a rule of thumb, it was determined that setting the FFIO cache size to roughly 10-15% of the disk space required by Abaqus yielded reasonable I/O performance. For this benchmark, the `FF_IO_OPTS` environment variable was defined by:

```
setenv FF_IO_OPTS '*.fct *.opr* *.ord *.fil *.mdl* *.stt* *.res *.sst *.hdx *.odb* *.023
*.nck* *.sct *.lop *.ngr *.elm *.ptn* *.stp* *.eig *.lnz* *.mass *.inp* *.scn* *.ddm
*.dat* fort*(eie.direct.nodiag.mbytes:4096:512:6:1:1:0,event.summary.mbytes.notrace)'
```

For the MPI version of Abaqus, different caches were specified for each MPI rank, as follows:

```
setenv FF_IO_OPTS_RANK0 '*.fct *.opr* *.ord *.fil *.mdl* *.stt* *.res *.sst *.hdx *.odb* *.023
*.nck* *.sct *.lop *.ngr *.ptn* *.stp* *.elm *.eig *.lnz* *.mass *.inp *.scn* *.ddm
*.dat* fort*(eie.direct.nodiag.mbytes:4096:512:6:1:1:0,event.summary.mbytes.notrace)'
```

```
setenv FF_IO_OPTS_RANK1 '*.fct *.opr* *.ord *.fil *.mdl* *.stt* *.res *.sst *.hdx *.odb* *.023
*.nck* *.sct *.lop *.ngr *.ptn* *.stp* *.elm *.eig *.lnz* *.mass *.inp *.scn* *.ddm
*.dat* fort*(eie.direct.nodiag.mbytes:4096:16:6:1:1:0,event.summary.mbytes.notrace)'
```

```
setenv FF_IO_OPTS_RANK2 '*.fct *.opr* *.ord *.fil *.mdl* *.stt* *.res *.sst *.hdx *.odb* *.023
*.nck* *.sct *.lop *.ngr *.ptn* *.stp* *.elm *.eig *.lnz* *.mass *.inp *.scn* *.ddm
*.dat* fort*(eie.direct.nodiag.mbytes:4096:16:6:1:1:0,event.summary.mbytes.notrace)'
```

```
setenv FF_IO_OPTS_RANK3 '*.fct *.opr* *.ord *.fil *.mdl* *.stt* *.res *.sst *.hdx *.odb* *.023
*.nck* *.sct *.lop *.ngr *.ptn* *.stp* *.elm *.eig *.lnz* *.mass *.inp *.scn* *.ddm
*.dat* fort*(eie.direct.nodiag.mbytes:4096:16:6:1:1:0,event.summary.mbytes.notrace)'
```

Event Tracing

By specifying the `.trace` option as part of the event parameter the user can enable the event tracing feature in FFIO, as follows:

```
setenv FF_IO_OPTS 'test*(eie.direct.mbytes:4096:128:6:1:1:0, event.summary.mbytes.trace )'
```

This option generates files of the form `ffio.events.pid` for each process that is part of the application. By default, event files are placed in `/tmp` but this destination can be changed by setting the `FFIO_TMPDIR` environment variable. These files contain time stamped events for files using the FFIO cache and can be used to trace I/O activity (for example, I/O sizes and offsets).

System Information and Issues

Applications written in C, C++, and Fortran are supported. C and C++ applications can be built with either the Intel or gcc compiler. Only Fortran codes built with the Intel compiler will work.

The following restrictions on FFIO must also be observed:

- The FFIO implementation of `pread/pwrite` is not correct (the file offset advances).
- Do not use FFIO to do I/O on a socket.
- Do not link your application with the `librt` asynchronous I/O library.
- Calls that operate on files in `/proc`, `/etc`, and `/dev` are not intercepted by FFIO.
- Calls that operate on `stdin`, `stdout`, and `stderr` are not intercepted by FFIO.
- FFIO is not intended for generic I/O applications such as `vi`, `cp`, or `mv`, and so on.

I/O Tuning

This chapter describes tuning information that you can use to improve I/O throughput and latency.

Application Placement and I/O Resources

It is useful to place an application on the same node as its I/O resource. For graphics applications, for example, this can improve performance up to 30 percent.

For example, for an SGI UV system with the following devices:

```
# gfxtopology
```

```
Serial number: UV-00000021
Partition number: 0
8 Blades
248 CPUs
283.70 Gb Memory Total
5 I/O Risers
```

Blade Location	NASID	PCI Address	X Server Display	Device
0 r001i01b08	0	0000:05:00.0	-	Matrox Pilot
4 r001i01b12	8	0001:02:01.0	-	SGI Scalable Graphics Capture
6 r001i01b14	12	0003:07:00.0	Layout0.0	nVidia Quadro FX 5800
		0003:08:00.0	Layout0.1	nVidia Quadro FX 5800
7 r001i01b15	14	0004:03:00.0	Layout0.2	nVidia Quadro FX 5800

For example, to run an OpenGL graphics program, such as `glxgears(1)`, on the third graphics processing unit using `numactl(8)`, type the following command:

```
% numactl -N 14 -m 14 /usr/bin/glxgears -display :0.2
```

This example assumes the X server was started with `:0 == Layout0`.

The `--N` parameter specifies to run the command on node 14. The `-m` parameter specifies to allocate memory only from node 14.

You could also use the `dplace(1)` command to place the application, see "`dplace Command`" on page 41.

Layout of Filesystems and XVM for Multiple RAIDs

There can be latency spikes in response from a RAID and such a spikes can in effect slow down all of the RAIDs as one I/O completion waits for all of the striped pieces to complete.

These latency spikes impact on throughput may be to stall all the I/O or to delay a few I/Os while others continue. It depends on how the I/O is striped across the devices. If the volumes are constructed as stripes to span all devices, and the I/Os are sized to be full stripes, the I/Os will stall, since every I/O has to touch every device. If the I/Os can be completed by touching a subset of the devices, then those that do not touch a high latency device can continue at full speed, while the stalled I/Os can complete and catch up later.

In large storage configurations, it is possible to lay out the volumes to maximize the opportunity for the I/Os to proceed in parallel, masking most of the effect of a few instances of high latency.

There are at least three classes of events that cause high latency I/O operations, as follows:

- Transient disk delays - one disk pauses
- Slow disks
- Transient RAID controller delays

The first two events affect a single logical unit number (LUN). The third event affects all the LUNs on a controller. The first and third events appear to happen at random. The second event is repeatable.

Suggested Shortcuts and Workarounds

This chapter contains suggested workarounds and shortcuts that you can use on your SGI system. It covers the following topics:

- "Determining Process Placement" on page 97
- "Resetting System Limits" on page 104
- "Linux Shared Memory Accounting" on page 110
- "OFED Tuning Requirements for UPC and SHMEM" on page 111
- "Setting Java Environment Variables" on page 112

Determining Process Placement

This section describes methods that can be used to determine where different processes are running. This can help you understand your application structure and help you decide if there are obvious placement issues.

There are some set-up steps to follow before determining process placement (note that all examples use the C shell):

1. Set up an alias as in this example, changing *guest* to your username:

```
% alias pu "ps -edaf|grep guest"
% pu
```

The `pu` command shows current processes.

2. Create the `.toprc` preferences file in your login directory to set the appropriate `top` options. If you prefer to use the `top` defaults, delete the `.toprc` file.

```
% cat <<EOF>> $HOME/.toprc

YEAbcDgHIjklMnoTP|qrsuzV{FWX
2mlt
EOF
```

3. Inspect all processes and determine which CPU is in use and create an alias file for this procedure. The CPU number is shown in the first column of the `top` output:

```
% top -b -n 1 | sort -n | more
% alias top1 "top -b -n 1 | sort -n "
```

Use the following variation to produce output with column headings:

```
% alias top1 "top -b -n 1 | head -4 | tail -1;top -b -n 1 | sort -n"
```

4. View your files (replacing *guest* with your username):

```
% top -b -n 1 | sort -n | grep guest
```

Use the following variation to produce output with column headings:

```
% top -b -n 1 | head -4 | tail -1;top -b -n 1 | sort -n grep guest
```

Example Using pthreads

The following example demonstrates a simple usage with a program name of *th*. It sets the number of desired OpenMP threads and runs the program. Notice the process hierarchy as shown by the PID and the PPID columns. The command usage is the following, where *n* is the number of threads:

```
% th n
```

```
% th 4
```

```
% pu
```

```
UID      PID    PPID    C  STIME TTY          TIME CMD
root     13784  13779   0  12:41 pts/3        00:00:00 login --
guest1
guest1   13785  13784   0  12:41 pts/3        00:00:00 -csh
guest1   15062  13785   0  15:23 pts/3        00:00:00 th 4  <-- Main thread
guest1   15063  15062   0  15:23 pts/3        00:00:00 th 4  <-- daemon thread
guest1   15064  15063   99 15:23 pts/3        00:00:10 th 4  <-- worker thread 1
guest1   15065  15063   99 15:23 pts/3        00:00:10 th 4  <-- worker thread 2
guest1   15066  15063   99 15:23 pts/3        00:00:10 th 4  <-- worker thread 3
guest1   15067  15063   99 15:23 pts/3        00:00:10 th 4  <-- worker thread 4
guest1   15068  13857   0  15:23 pts/5        00:00:00 ps -aef
guest1   15069  13857   0  15:23 pts/5        00:00:00 grep guest1
```

```
% top -b -n 1 | sort -n | grep guest1
```

```
LC %CPU  PID USER      PRI  NI  SIZE  RSS SHARE STAT %MEM  TIME COMMAND
 3  0.0 15072 guest1    16   0  3488 1536  3328 S    0.0  0:00 grep
```

```

 5  0.0 13785 guest1    15  0  5872 3664  4592 S    0.0  0:00 csh
 5  0.0 15062 guest1    16  0 15824 2080  4384 S    0.0  0:00 th
 5  0.0 15063 guest1    15  0 15824 2080  4384 S    0.0  0:00 th
 5 99.8 15064 guest1    25  0 15824 2080  4384 R    0.0  0:14 th
 7  0.0 13826 guest1    18  0  5824 3552  5632 S    0.0  0:00 csh
10 99.9 15066 guest1    25  0 15824 2080  4384 R    0.0  0:14 th
11 99.9 15067 guest1    25  0 15824 2080  4384 R    0.0  0:14 th
13 99.9 15065 guest1    25  0 15824 2080  4384 R    0.0  0:14 th
15  0.0 13857 guest1    15  0  5840 3584  5648 S    0.0  0:00 csh
15  0.0 15071 guest1    16  0 70048 1600 69840 S    0.0  0:00 ort
15  1.5 15070 guest1    15  0  5056 2832  4288 R    0.0  0:00top

```

Now skip the Main and daemon processes and place the rest:

```
% /usr/bin/dplace -s 2 -c 4-7 th 4
```

```
% pu
```

```

UID          PID  PPID  C  STIME TTY          TIME CMD
root         13784 13779  0 12:41 pts/3        00:00:00 login --
guest1
guest1      13785 13784  0 12:41 pts/3        00:00:00 -csh
guest1      15083 13785  0 15:25 pts/3        00:00:00 th 4
guest1      15084 15083  0 15:25 pts/3        00:00:00 th 4
guest1      15085 15084 99 15:25 pts/3        00:00:19 th 4
guest1      15086 15084 99 15:25 pts/3        00:00:19 th 4
guest1      15087 15084 99 15:25 pts/3        00:00:19 th 4
guest1      15088 15084 99 15:25 pts/3        00:00:19 th 4
guest1      15091 13857  0 15:25 pts/5        00:00:00 ps -aef
guest1      15092 13857  0 15:25 pts/5        00:00:00 grep guest1

```

```
% top -b -n 1 | sort -n | grep guest1
```

```

LC %CPU  PID USER      PRI  NI  SIZE  RSS SHARE STAT %MEM  TIME COMMAND
 4 99.9 15085 guest1    25   0 15856 2096 6496 R    0.0  0:24 th
 5 99.8 15086 guest1    25   0 15856 2096 6496 R    0.0  0:24 th
 6 99.9 15087 guest1    25   0 15856 2096 6496 R    0.0  0:24 th
 7 99.9 15088 guest1    25   0 15856 2096 6496 R    0.0  0:24 th
 8  0.0 15095 guest1    16   0  3488 1536 3328 S    0.0  0:00 grep
12  0.0 13785 guest1    15   0  5872 3664 4592 S    0.0  0:00 csh
12  0.0 15083 guest1    16   0 15856 2096 6496 S    0.0  0:00 th
12  0.0 15084 guest1    15   0 15856 2096 6496 S    0.0  0:00 th

```

```
15  0.0 15094 guest1    16   0 70048 1600 69840 S    0.0  0:00 sort
15  1.6 15093 guest1    15   0  5056 2832  4288 R    0.0  0:00 top
```

Example Using OpenMP

The following example demonstrates a simple OpenMP usage with a program name of md. Set the desired number of OpenMP threads and run the program, as shown below:

```
% alias pu "ps -edaf | grep guest1"
% setenv OMP_NUM_THREADS 4
% md
```

The following output is created:

```
% pu
```

```
UID      PID  PPID  C  STIME TTY          TIME CMD
root     21550 21535  0  21:48 pts/0        00:00:00 login -- guest1
guest1   21551 21550  0  21:48 pts/0        00:00:00 -csh
guest1   22183 21551 77  22:39 pts/0        00:00:03 md    <-- parent / main
guest1   22184 22183  0  22:39 pts/0        00:00:00 md    <-- daemon
guest1   22185 22184  0  22:39 pts/0        00:00:00 md    <-- daemon helper
guest1   22186 22184 99  22:39 pts/0        00:00:03 md    <-- thread 1
guest1   22187 22184 94  22:39 pts/0        00:00:03 md    <-- thread 2
guest1   22188 22184 85  22:39 pts/0        00:00:03 md    <-- thread 3
guest1   22189 21956  0  22:39 pts/1        00:00:00 ps -aef
guest1   22190 21956  0  22:39 pts/1        00:00:00 grep guest1
```

```
% top -b -n 1 | sort -n | grep guest1
```

```
LC %CPU  PID USER      PRI  NI  SIZE  RSS  SHARE STAT %MEM  TIME COMMAND
 2  0.0 22192 guest1    16   0 70048 1600 69840 S    0.0  0:00 sort
 2  0.0 22193 guest1    16   0  3488 1536  3328 S    0.0  0:00 grep
 2  1.6 22191 guest1    15   0  5056 2832  4288 R    0.0  0:00 top
 4 98.0 22186 guest1    26   0 26432 2704  4272 R    0.0  0:11 md
 8  0.0 22185 guest1    15   0 26432 2704  4272 S    0.0  0:00 md
 8 87.6 22188 guest1    25   0 26432 2704  4272 R    0.0  0:10 md
 9  0.0 21551 guest1    15   0  5872 3648  4560 S    0.0  0:00 csh
 9  0.0 22184 guest1    15   0 26432 2704  4272 S    0.0  0:00 md
 9 99.9 22183 guest1    39   0 26432 2704  4272 R    0.0  0:11 md
```

```
14 98.7 22187 guest1      39   0 26432 2704 4272 R    0.0  0:11 md
```

From the notation on the right of the `pu` list, you can see the `-x 6` pattern.

```
place 1, skip 2 of them, place 3 more [ 0 1 1 0 0 0 ]
now, reverse the bit order and create the dplace -x mask
[ 0 0 0 1 1 0 ] --> [ 0x06 ] --> decimal 6
dplace does not currently process hex notation for this bit mask)
```

The following example confirms that a simple `dplace` placement works correctly:

```
% setenv OMP_NUM_THREADS 4
% /usr/bin/dplace -x 6 -c 4-7 md
% pu
UID      PID  PPID  C STIME TTY      TIME CMD
root     21550 21535  0 21:48 pts/0    00:00:00 login -- guest1
guest1   21551 21550  0 21:48 pts/0    00:00:00 -csh
guest1   22219 21551 93 22:45 pts/0    00:00:05 md
guest1   22220 22219  0 22:45 pts/0    00:00:00 md
guest1   22221 22220  0 22:45 pts/0    00:00:00 md
guest1   22222 22220 93 22:45 pts/0    00:00:05 md
guest1   22223 22220 93 22:45 pts/0    00:00:05 md
guest1   22224 22220 90 22:45 pts/0    00:00:05 md
guest1   22225 21956  0 22:45 pts/1    00:00:00 ps -aef
guest1   22226 21956  0 22:45 pts/1    00:00:00 grep guest1
```

```
% top -b -n 1 | sort -n | grep guest1
```

LC	%CPU	PID	USER	PRI	NI	SIZE	RSS	SHARE	STAT	%MEM	TIME	COMMAND
2	0.0	22228	guest1	16	0	70048	1600	69840	S	0.0	0:00	sort
2	0.0	22229	guest1	16	0	3488	1536	3328	S	0.0	0:00	grep
2	1.6	22227	guest1	15	0	5056	2832	4288	R	0.0	0:00	top
4	0.0	22220	guest1	15	0	28496	2736	21728	S	0.0	0:00	md
4	99.9	22219	guest1	39	0	28496	2736	21728	R	0.0	0:12	md
5	99.9	22222	guest1	25	0	28496	2736	21728	R	0.0	0:11	md
6	99.9	22223	guest1	39	0	28496	2736	21728	R	0.0	0:11	md
7	99.9	22224	guest1	39	0	28496	2736	21728	R	0.0	0:11	md
9	0.0	21551	guest1	15	0	5872	3648	4560	S	0.0	0:00	csh
15	0.0	22221	guest1	15	0	28496	2736	21728	S	0.0	0:00	md

Combination Example (MPI and OpenMP)

For this example, explicit placement using the `dplace -e -c` command is used to achieve the desired placement. If an `x` is used in one of the CPU positions, `dplace` does not explicitly place that process.

If running without a `cpuset`, the `x` processes run on any available CPU.

If running with a `cpuset`, you have to renumber the CPU numbers to refer to “logical” CPUs (0 . . . n) within the `cpuset`, regardless of which physical CPUs are in the `cpuset`. When running in a `cpuset`, the unplaced processes are constrained to the set of CPUs within the `cpuset`.

For information about `cpusets`, see the *SGI Cpuset Software Guide*.

The following example shows a hybrid MPI and OpenMP job with two MPI processes, each with two OpenMP threads and no `cpusets`:

```
% setenv OMP_NUM_THREADS 2
% efc -O2 -o hybrid hybrid.f -lmpi -openmp

% mpirun -v -np 2 /usr/bin/dplace -e -c x,8,9,x,x,x,x,10,11 hybrid

-----
# if using cpusets ...
-----
# we need to reorder cpus to logical within the 8-15 set [0-7]

% cpuset -q omp -A mpirun -v -np 2 /usr/bin/dplace -e -c x,0,1,x,x,x,x,2,3,4,5,6,7 hybrid

# We need a table of options for these pairs. "x" means don't
# care. See the dplace man page for more info about the -e option.
# examples at end

-np  OMP_NUM_THREADS  /usr/bin/dplace -e -c <as shown> a.out
---  -----
  2         2         x,0,1,x,x,x,x,2,3
  2         3         x,0,1,x,x,x,x,2,3,4,5
  2         4         x,0,1,x,x,x,x,2,3,4,5,6,7

  4         2         x,0,1,2,3,x,x,x,x,x,x,x,x,4,5,6,7
  4         3
x,0,1,2,3,x,x,x,x,x,x,x,x,4,5,6,7,8,9,10,11
```



```

      4          4
x,0,1,2,3,x,x,x,x,x,x,x,x,4,5,6,7,8,9,10,11,12,13,14,15
Notes:          0 <- 1 -> <- 2 -> <- 3 -> <----- 4
----->

```

Notes:

0. mpi daemon process
1. mpi child procs, one per np
2. omp daemon procs, one per np
3. omp daemon helper procs, one per np
4. omp thread procs, (OMP_NUM_THREADS - 1) per np

```

-----
# Example - -np 2 and OMP_NUM_THREADS 2
-----

```

```

% setenv OMP_NUM_THREADS 2
% efc -O2 -o hybrid hybrid.f -lmpi -openmp

% mpirun -v -np 2 /usr/bin/dplace -e -c x,8,9,x,x,x,x,10,11 hybrid

% pu

```

```

UID          PID  PPID  C STIME TTY          TIME CMD
root    21550 21535  0 Mar17 pts/0 00:00:00 login -- guest1
guest1  21551 21550  0 Mar17 pts/0 00:00:00 -csh
guest1  23391 21551  0 00:32 pts/0 00:00:00 mpirun -v -np 2

/usr/bin/dplace
guest1  23394 23391  2 00:32 pts/0 00:00:00 hybrid <-- mpi daemon
guest1  23401 23394 99 00:32 pts/0 00:00:03 hybrid <-- mpi child 1
guest1  23402 23394 99 00:32 pts/0 00:00:03 hybrid <-- mpi child 2
guest1  23403 23402  0 00:32 pts/0 00:00:00 hybrid <-- omp daemon 2
guest1  23404 23401  0 00:32 pts/0 00:00:00 hybrid <-- omp daemon 1
guest1  23405 23404  0 00:32 pts/0 00:00:00 hybrid <-- omp daemon hlpr 1
guest1  23406 23403  0 00:32 pts/0 00:00:00 hybrid <-- omp daemon hlpr 2
guest1  23407 23403 99 00:32 pts/0 00:00:03 hybrid <-- omp thread 2-1
guest1  23408 23404 99 00:32 pts/0 00:00:03 hybrid <-- omp thread 1-1
guest1  23409 21956  0 00:32 pts/1 00:00:00 ps -aef
guest1  23410 21956  0 00:32 pts/1 00:00:00 grep guest1

```

```
% top -b -n 1 | sort -n | grep guest1
```

LC	%CPU	PID	USER	PRI	NI	SIZE	RSS	SHARE	STAT	%MEM	TIME	COMMAND
0	0.0	21551	guest1	15	0	5904	3712	4592	S	0.0	0:00	csch
0	0.0	23394	guest1	15	0	883M	9456	882M	S	0.1	0:00	hybrid
4	0.0	21956	guest1	15	0	5856	3616	5664	S	0.0	0:00	csch
4	0.0	23412	guest1	16	0	70048	1600	69840	S	0.0	0:00	sort
4	1.6	23411	guest1	15	0	5056	2832	4288	R	0.0	0:00	top
5	0.0	23413	guest1	16	0	3488	1536	3328	S	0.0	0:00	grep
8	0.0	22005	guest1	15	0	5840	3584	5648	S	0.0	0:00	csch
8	0.0	23404	guest1	15	0	894M	10M	889M	S	0.1	0:00	hybrid
8	99.9	23401	guest1	39	0	894M	10M	889M	R	0.1	0:09	hybrid
9	0.0	23403	guest1	15	0	894M	10M	894M	S	0.1	0:00	hybrid
9	99.9	23402	guest1	25	0	894M	10M	894M	R	0.1	0:09	hybrid
10	99.9	23407	guest1	25	0	894M	10M	894M	R	0.1	0:09	hybrid
11	99.9	23408	guest1	25	0	894M	10M	889M	R	0.1	0:09	hybrid
12	0.0	23391	guest1	15	0	5072	2928	4400	S	0.0	0:00	mpirun
12	0.0	23406	guest1	15	0	894M	10M	894M	S	0.1	0:00	hybrid
14	0.0	23405	guest1	15	0	894M	10M	889M	S	0.1	0:00	hybrid

Resetting System Limits

To regulate these limits on a per-user basis for applications that do not rely on `limit.h`, the `limits.conf` file can be modified. System limits that can be modified include maximum file size, maximum number of open files, maximum stack size, and so on. You can view this file is, as follows:

```
[user@machine user]# cat /etc/security/limits.conf
# /etc/security/limits.conf
#
#Each line describes a limit for a user in the form:
#
#          #
#Where:
# can be:
#          - an user name
#          - a group name, with @group syntax
#          - the wildcard *, for default entry
#
# can have the two values:
```

```

#         - "soft" for enforcing the soft limits
#         - "hard" for enforcing hard limits
#
# can be one of the following:
#         - core - limits the core file size (KB)
#         - data - max data size (KB)
#         - fsize - maximum filesize (KB)
#         - memlock - max locked-in-memory address space (KB)
#         - nofile - max number of open files
#         - rss - max resident set size (KB)
#         - stack - max stack size (KB)
#         - cpu - max CPU time (MIN)
#         - nproc - max number of processes
#         - as - address space limit
#         - maxlogins - max number of logins for this user
#         - priority - the priority to run user process with
#         - locks - max number of file locks the user can hold
#
#           #
#*          soft   core           0
#*          hard   rss           10000
#@student   hard   nproc          20
#@faculty   soft   nproc          20
#@faculty   hard   nproc          50
#ftp        hard   nproc          0
#@student   -      maxlogins      4

# End of file

```

For instructions on how to change these limits, see "Resetting the File Limit Resource Default" on page 105.

Resetting the File Limit Resource Default

Several large user applications use the value set in the `limit.h` file as a hard limit on file descriptors and that value is noted at compile time. Therefore, some applications may need to be recompiled in order to take advantage of the SGI system hardware.

To regulate these limits on a per-user basis for applications that do not rely on `limit.h`, the `limits.conf` file can be modified. This allows the administrator to

set the allowed number of open files per user and per group. This also requires a one-line change to the `/etc/pam.d/login` file.

Follow this procedure to execute these changes:

1. Add the following line to `/etc/pam.d/login`:

```
session required /lib/security/pam_limits.so
```

2. Add the following line to `/etc/security/limits.conf`, where *username* is the user's login and *limit* is the new value for the file limit resource:

```
[username] hard nofile [limit]
```

The following command shows the new limit:

```
ulimit -H -n
```

Because of the large number of file descriptors that some applications require, such as MPI jobs, you might need to increase the system-wide limit on the number of open files on your SGI system. The default value for the file limit resource is 1024. The default 1024 file descriptors allows for approximately 199 MPI processes per host. You can increase the file descriptor value to 8196 to allow for more than 512 MPI processes per host by adding the following lines to the `/etc/security/limits.conf` file:

```
* soft nofile 8196
* hard nofile 8196
```

The `ulimit -a` command displays all limits, as follows:

```
sys:~ # ulimit -a
core file size          (blocks, -c) 1
data seg size           (kbytes, -d) unlimited
scheduling priority     (-e) 0
file size               (blocks, -f) unlimited
pending signals         (-i) 511876
max locked memory       (kbytes, -l) 64
max memory size         (kbytes, -m) 55709764
open files              (-n) 1024
pipe size               (512 bytes, -p) 8
POSIX message queues    (bytes, -q) 819200
real-time priority      (-r) 0
stack size              (kbytes, -s) 8192
```

```
cpu time          (seconds, -t) unlimited
max user processes (-u) 511876
virtual memory    (kbytes, -v) 68057680
file locks        (-x) unlimited
```

Resetting the Default Stack Size

Some applications will not run well on an SGI system with a small stack size. To set a higher stack limit, follow the instructions in "Resetting the File Limit Resource Default" on page 105 and add the following lines to the `/etc/security/limits.conf` file:

```
* soft stack 300000
* hard stack unlimited
```

This sets a soft stack size limit of 300000 KB and an unlimited hard stack size for all users (and all processes).

Another method that does not require root privilege relies on the fact that many MPI implementation use `ssh`, `rsh`, or some sort of login shell to start the MPI rank processes. If you merely need to bump up the soft limit, you can modify your shell's startup script. For example, if your login shell is `bash` then add something like the following to your `.bashrc` file:

```
% ulimit -s 300000
```

Note that SGI MPT MPI allows you to set your stack size limit larger with the `ulimit` or `limit` shell command before launching an MPI program with `mpirun(1)` or `mpiexec_mpt(1)`. MPT will propagate the stack limit setting to all MPI processes in the job.

For more information on default settings, also see "Resetting the File Limit Resource Default" on page 105.

Avoiding Segmentation Faults

The default stack size in the Linux operating system is 8MB (8192 kbytes). This value need to be increased to avoid "Segmentation Faults" errors.

You can use the `ulimit -a` command to view the stack size, as follows:

```
uv44-sys:~ # ulimit -a
core file size          (blocks, -c) unlimited
```

```
data seg size          (kbytes, -d) unlimited
file size              (blocks, -f) unlimited
pending signals        (-i) 204800
max locked memory      (kbytes, -l) unlimited
max memory size        (kbytes, -m) unlimited
open files             (-n) 16384
pipe size              (512 bytes, -p) 8
POSIX message queues   (bytes, -q) 819200
stack size            (kbytes, -s) 8192
cpu time               (seconds, -t) unlimited
max user processes     (-u) 204800
virtual memory         (kbytes, -v) unlimited
file locks             (-x) unlimited
```

To change the value, perform a command similar to the following:

```
uv44-sys:~ # ulimit -s 300000
```

There is a similar variable for OpenMP programs. If you get a segmentation fault right away while running a program parallelized with OpenMP, a good idea is to increase the `KMP_STACKSIZE` to a larger size. The default size in Intel Compilers is 4MB.

For example, to increase it to 64MB in `cs` shell perform the following:

```
setenv KMP_STACKSIZE 64M
```

in `bash`:

```
export KMP_STACKSIZE=64M
```

Resetting Virtual Memory Size

The virtual memory parameter `vmemoryuse` determines the amount of virtual memory available to your application. If you are running with `csch`, use `csch` commands, such as, the following:

```
limit
limit vmemoryuse 7128960
limit vmemoryuse unlimited
```

The following MPI program fails with a memory-mapping error because of a virtual memory parameter `vmemoryuse` value set too low:

```
% limit vmemoryuse 7128960

% mpirun -v -np 4 ./program
MPI: libxmpi.so 'SGI MPI 4.9 MPT 1.14 07/18/06 08:43:15'
MPI: libmpi.so 'SGI MPI 4.9 MPT 1.14 07/18/06 08:41:05'
MPI: MPI_MSGS_MAX = 524288
MPI: MPI_BUFS_PER_PROC= 32
mmap failed (mmap_base) for 504972 pages (8273461248
bytes) Killed n
```

The program now succeeds when virtual memory is unlimited:

```
% limit vmemoryuse unlimited

% mpirun -v -np 4 ./program
MPI: libxmpi.so 'SGI MPI 4.9 MPT 1.14 07/18/06 08:43:15'
MPI: libmpi.so 'SGI MPI 4.9 MPT 1.14 07/18/06 08:41:05'
MPI: MPI_MSGS_MAX = 524288
MPI: MPI_BUFS_PER_PROC= 32

HELLO WORLD from Processor 0

HELLO WORLD from Processor 2

HELLO WORLD from Processor 1

HELLO WORLD from Processor 3
```

If you are running with `bash`, use `bash` commands, such as, the following:

```
ulimit -a
ulimit -v 7128960
ulimit -v unlimited
```

Linux Shared Memory Accounting

The Linux operating system does not calculate memory utilization in a manner that is useful for certain applications in situations where regions are shared among multiple processes. This can lead to over-reporting of memory and to processes being killed by schedulers erroneously detecting memory quota violation.

The `get_weighted_memory_size` function weighs shared memory regions by the number of processes using the regions. Thus, if 100 processes are each sharing a total of 10GB of memory, the weighted memory calculation shows 100MB of memory shared per process, rather than 10GB for each process.

Because this function applies mostly to applications with large shared-memory requirements, it is located in the SGI NUMA tools package and made available in the `libmemacct` library available from a new package called `memacct`. The library function makes a call to the `numatools` kernel module, which returns the weighted sum back to the library, and then returns back to the application.

The usage statement for the `memacct` call is, as follows:

```
cc ... -lmemacct
#include <sys/types.h>
extern int get_weighted_memory_size(pid_t pid);
```

The syntax of the `memacct` call is, as follows:

```
int *get_weighted_memory_size(pid_t pid);
```

Returns the weighted memory (RSS) size for a `pid`, in bytes. This weights the size of shared regions by the number of processes accessing it. Return -1 when an error occurs and set `errno`, as follows:

ESRCH	Process <code>pid</code> was not found.
ENOSYS	The function is not implemented. Check if <code>numatools</code> kernel package is up-to-date.

Normally, the following errors should not occur:

ENOENT	Can not open <code>/proc/numatools</code> device file.
EPERM	No read permission on <code>/proc/numatools</code> device file.
ENOTTY	Inappropriate <code>ioctl</code> operation on <code>/proc/numatools</code> device file.
EFAULT	Invalid arguments. The <code>ioctl()</code> operation performed by the function failed with invalid arguments.

For more information, see the `memacct(3)` man page.

OFED Tuning Requirements for UPC and SHMEM

You can specify the maximum number of queue pairs (QPs) for SHMEM and UPC applications when run on large clusters over an OFED fabric, such as InfiniBand. If the `log_num_qp` parameter is set to a number that is too low, the system generates the following message:

```
MPT Warning: IB failed to create a QP
```

SHMEM and UPC codes use the InfiniBand RC protocol for communication between all pairs of processes in the parallel job, which requires a large number of QPs. The `log_num_qp` parameter defines the \log_2 of the number of QPs. The following procedure explains how to specify the `log_num_qp` parameter.

Procedure 9-1 To specify the `log_num_qp` parameter

1. Log into one of the hosts upon which you installed the MPT software as the root user.
2. Use a text editor to open file `/etc/modprobe.d/libmlx4.conf`.
3. Add a line similar to the following to file `/etc/modprobe.d/libmlx4.conf`:

```
options mlx4_core log_num_qp=21
```

By default, the maximum number of queue pairs on RHEL platforms is 2^{18} (262144).

By default, the maximum number of queue pairs on SLES platforms is 2^{17} (131072).

4. Save and close the file.

5. Repeat the preceding steps on other hosts.

Setting Java Environment Variables

When Java software starts, it checks the environment in which it is running and configures itself to fit, assuming that it owns the entire environment. The default for some Java implementations (for example, IBM J9 1.4.2) is to start a garbage collection (GC) thread for every CPU it sees. Other Java implementations use other algorithms to decide the number of GCs to start, but the number will generally be 0.5 to 1 times the number of CPUs. On a 1 or 2 socket system, that works out well enough. This strategy does not scale well to large core-count systems, however.

Java command line options allow you to control the number of GC threads that the Java virtual machine (JVM) will use. In many cases, a single GC thread is sufficient, as set in the examples above. In other cases, a larger number may be appropriate and can be set with the applicable environment variable or command line option. Properly tuning the number of GC threads for an application is an exercise in performance optimization, but a reasonable starting point is to use one GC thread per active worker thread.

For Sun Java (now Oracle Java):

```
-XX:ParallelGCThreads
```

For IBM Java:

```
-Xgcthreads
```

An example command line option:

```
java -XX:+UseParallelGC -XX:ParallelGCThreads=1
```

The system administrator may choose to implement the option to limit the number of GC threads to a reasonable value with an environment variable set in the global profile, for example the `/etc/profile.local` file so casual Java users can avoid difficulties.

Environment variable settings, are as follows:

For Sun Java (now Oracle Java):

```
JAVA_OPTIONS="-XX:ParallelGCThreads=1"
```

For IBM Java:

```
IBM_JAVA_OPTIONS="-Xgcthreads1"
```


Using PerfSocket

This chapter includes the following topics:

- "About SGI PerfSocket" on page 115
- "Installing and Using PerfSocket" on page 115
- "About Security When Using PerfSocket" on page 118
- "Troubleshooting" on page 118

About SGI PerfSocket

The SGI PerfSocket feature improves an application's TCP/IP communication within a host. The PerfSocket software intercepts local TCP/IP communication and routes the communication through shared memory, which eliminates much of the overhead incurred when communication data passes through the operating system kernel.

PerfSocket includes a system library, a daemon, a kernel module, and a wrapper command that enables PerfSocket functionality within specified processes. Only processes specifically run with the PerfSocket wrapper command are run with this feature.

SGI includes the PerfSocket technology in the SGI Accelerate™ product within the SGI Performance Suite. After you install SGI Accelerate, you can use the procedure in the following topic to install PerfSocket as an optional feature:

"Installing PerfSocket (Administrator Procedure)" on page 116.

Installing and Using PerfSocket

The SGI Accelerate installation process does not install PerfSocket. The following procedures explain how to install and use PerfSocket:

- "Installing PerfSocket (Administrator Procedure)" on page 116
- "Running an Application With PerfSocket" on page 117

Installing PerfSocket (Administrator Procedure)

When you install the SGI Performance Suite software, the installer does not install PerfSocket. You need to complete the procedure in this topic to install PerfSocket.

The PerfSocket RPM contains the PerfSocket libraries, applications, kernel modules, and man(1) pages. The installer writes the majority of PerfSocket's files to the `/opt/sgi/perfsocket` directory.

The following procedure explains how to install PerfSocket and how to start the PerfSocket daemon on an SGI UV computer system.

Procedure 10-1 To install Perfsocket on an SGI UV system

1. Log in as root.
2. Install the PerfSocket software.

This command sequence differs, depending on your platform, as follows:

- On RHEL platforms, type the following command:

```
# yum install perfsocket
```

When RHEL displays the download size and displays the `Is this ok [y/N]:` prompt, type `y` and press Enter.

- On SLES platforms, type the following commands:

```
# zypper refresh
# zypper install perfsocket
```

When SLES displays the download size and displays the `Continue? [y/n/?] (y):` prompt, type `y` and press Enter.

3. Type the following command to turn on the PerfSocket service:

```
# chkconfig perfsocket on
```

4. Type the following command to verify that the PerfSocket service is on:

```
# chkconfig --list | grep perfsocket
perfsocket      0:off  1:off  2:on   3:on   4:on   5:on   6:off
```

5. Type the following command to start the PerfSocket daemon without a reboot:

```
# service perfsocket start
Starting the PerfSocket daemon
```

6. (Optional) Type the following command to verify that the kernel module is loaded:

```
# lsmod | grep perfsock
perfsock          71297  10
```

7. (Optional) Type the following command to display the PerfSocket processes:

```
# ps ax | grep perfsocketd
10308 ?        Ss      0:00 /opt/sgi/perfsocket/sbin/perfsocketd
10319 pts/0    S+      0:00 grep perfsocketd
```

Running an Application With PerfSocket

You do not need to recompile an application in order to use PerfSocket. An application that runs with PerfSocket automatically detects whether the endpoints of its communication are also using PerfSocket. All applications that use any particular socket endpoint must be run with PerfSocket in order for PerfSocket to accelerate the communication.

The following procedure explains how to invoke PerfSocket to run with your application programs.

Procedure 10-2 To run an application with PerfSocket

1. Type the following command to load the PerfSocket environment module:

```
$ module load perfsocket
```

The previous command adds the PerfSocket wrapper command to your `PATH` variable and adds the PerfSocket libraries to your `LD_LIBRARY_PATH`.

2. For each command that you want to run with PerfSocket, prefix the command with the `perfsocket(1)` command.

For example, if applications `a.out` and `b.out` communicate with TCP/IP, type the following commands to enable them to communicate through PerfSocket:

```
$ perfsocket ./a.out &
$ perfsocket ./b.out
```

For more information, see the `perfsocket(1)` and `perfsocketd(1)` man pages.

About Security When Using PerfSocket

The PerfSocket daemon facilitates application communication under the following conditions:

- An application that uses PerfSocket connects to another application that also uses PerfSocket
- The applications that connect run on the same host
- The user ID is identical for both the connecting process and the receiving process

Only the user who owns the applications can read the shared memory structure used for communication. No additional copies of the data are made. If a process that uses PerfSocket calls `exec()`, all PerfSocket-enabled sockets are duplicated to `/dev/null`.

Troubleshooting

If PerfSocket detects an unsupported condition, stop using PerfSocket.

SGI implemented and tested PerfSocket with the majority of socket and file I/O APIs. It is possible that you might attempt to use PerfSocket with a rarely used API or an unsupported usage pattern. If you encounter an unsupported condition, PerfSocket logs an error message and aborts the application.

PerfSocket writes its log messages to the system log, `/var/log/messages`.

Index

A

- Amdahl's law, 75
 - execution time given n and p , 79
 - parallel fraction p , 78
 - parallel fraction p given speedup(n), 78
 - speedup(n) given p , 78
 - superlinear speedup, 77
- application placement and I/O resources, 95
- application tuning process, 9
- automatic parallelization
 - limitations, 72
- avoiding segmentation faults, 107

C

- cache bank conflicts, 67
- cache coherency, 37
- Cache coherent non-uniform memory access (ccNUMA) systems, 84
- cache performance, 66
- ccNUMA
 - See also "cache coherent non-uniform memory access", 84
- ccNUMA architecture, 36
- cluster environment, 1
- commands
 - dlook, 50
 - dplace, 42
 - topology, 24
- common compiler options, 3
- compiler command line, 3
- compiler libraries
 - C/C++, 6
 - dynamic libraries, 6
 - message passing, 7

- overview, 6
- compiler libraries
 - static libraries, 6
- compiler options
 - tracing and porting, 60
- compiler options for tuning, 62
- compiling environment, 3
 - compiler overview, 3
 - debugger overview, 8
 - libraries, 6
 - modules, 4
- Configuring MPT
 - OFED, 111
- CPU-bound processes, 16
- Cpuset Facility
 - advantages, 58
 - overview, 58

D

- data decomposition, 69
- data dependency, 73
- data parallelism, 69
- data placement tools, 35
 - cpusets, 38
 - dplace, 38
 - overview, 35
 - taskset, 38
- data Pplacement practices, 37
- debugger overview, 8
- debuggers, 18
 - gdb, 8
 - idb, 8
 - TotalView, 8
- denormalized arithmetic, 4
- determining parallel code amount, 70

- determining tuning needs
 - tools used, 62
- distributed shared memory (DSM), 36
- dlook command, 50
- dplace command, 42

E

- Environment variables, 74
- explicit data decomposition, 70

F

- False sharing, 73
- file limit resources
 - resetting, 105
- Flexible File I/O (FFIO), 92
 - environment variables to set, 88
 - operation, 87
 - overview, 87
 - simple examples, 89
- floating-point programs, 81
- Floating-Point Software Assist, 81
- FPSWA
 - See "Floating-Point Software Assist", 81
- functional parallelism, 69

G

- gdb tool, 18
- Global reference unit (GRU), 82
- GNU debugger, 18
- gtopology command, 24
- Gustafson's law, 80

I

- I/O tuning

120

- application placement, 95
- layout of filesystems, 96
- I/O-bound processes, 16
- idb tool, 18
- implicit data decomposition, 70
- iostat command, 32

J

- Java environment variables
 - setting, 112

L

- latency, 1
- layout of filesystems, 96
- limits
 - system, 104
- linkstat command, 28
- Linux shared memory accounting, 110

M

- memory
 - cache coherency, 37
 - ccNUMA architecture, 36
 - distributed shared memory (DSM), 36
 - non-uniform memory access (NUMA), 37
- memory accounting, 110
- memory management, 2, 67
- memory page, 2
- memory strides, 66
- memory-bound processes, 16
- Message Passing Toolkit
 - for parallelization, 71
- modules, 4
 - command examples, 4
- MPI on SGI UV systems

- general considerations, 81
 - job performance types, 82
 - other ccNUMA performance issues, 84
 - MPI on UV systems, 82
 - MPI profiling, 84
 - MPInside profiling tool, 84
 - MPP definition, 2
- N**
- non-uniform memory access (NUMA), 37
 - NUMA Tools
 - command
 - dlook, 50
 - dplace, 42
 - installing, 58
- O**
- OFED configuration for MPT, 111
 - OpenMP, 71
 - environment variables, 74
 - Guide OpenMP Compiler, 17
- P**
- parallel execution
 - Amdahl's law, 75
 - parallel fraction p, 78
 - parallel speedup, 76
 - parallelization
 - automatic, 72
 - using MPI, 71
 - using OpenMP, 71
 - perf tool, 16
 - performance
 - Assure Thread Analyzer, 17
 - Guide OpenMP Compiler, 17
 - VTune, 17
 - performance analysis, 9
 - Performance Co-Pilot monitoring tools, 27
 - hubstats, 28
 - linkstat, 28
 - Other Performance Co-Pilot monitoring tools, 28
 - performance gains
 - types of, 9
 - performance problems
 - sources, 16
 - PerfSuite script, 16
 - process placement
 - determining, 97
 - MPI and OpenMP, 102
 - set-up, 97
 - using OpenMP, 100
 - using pthreads, 98
 - profiling
 - MPI, 84
 - perf, 16
 - PerfSuite, 16
 - ps command, 31
- R**
- resetting default system stack size, 107
 - resetting file limit resources, 105
 - resetting system limit resources, 104
 - resetting virtual memory size, 109
 - resident set size, 2
- S**
- sar command, 32
 - scalable computing, 1
 - segmentation faults, 107
 - setting Java environment variables, 112
 - SGI PerfBoost, 84
 - SGI PerfCatcher, 84
 - SHMEM, 7

- shortening execution time, 76
- shubstats command, 28
- SMP definition, 1
- stack size
 - resetting, 107
- suggested shortcuts and workarounds, 97
- superlinear speedup, 77
- swap space, 2
- system
 - overview, 1
- system configuration, 9
- system limit resources
 - resetting, 104
- system limits
 - address space limit, 105
 - core file siz, 105
 - CPU time, 105
 - data size, 105
 - file locks, 105
 - file size, 105
 - locked-in-memory address space, 105
 - number of logins, 105
 - number of open files, 105
 - number of processes, 105
 - priority of user process, 105
 - resetting, 104
 - resident set size, 105
 - stack size, 105
- system monitoring tools, 23
 - command
 - topology, 23
- system usage commands, 30
 - iostat, 32
 - ps, 31
 - sar, 32
 - uptime, 30
 - vmstat, 31
 - w, 31

T

- taskset command, 39
- tools
 - Assure Thread Analyzer, 17
 - Guide OpenMP Compiler, 17
 - perf, 16
 - PerfSuite, 16
 - VTune, 17
- topology command, 23, 24
- tuning
 - cache performance, 66
 - debugging tools
 - idb, 18
 - dplace, 74
 - environment variables, 74
 - false sharing, 73
 - heap corruption, 61
 - managing memory, 67
 - multiprocessor code, 68
 - parallelization, 70
 - profiling
 - perf, 16
 - PerfSuite script, 16
 - VTune analyzer, 17
 - single processor code, 59
 - using compiler options, 62
 - using dplace, 74
 - using math functions, 61
 - using taskset, 74
 - verifying correct results, 60

U

- uname command, 15
- unflow arithmetic
 - effects of, 4
- uptime command, 30
- UV Hub, 82

V

virtual addressing, 2
virtual memory, 2
vmstat command, 31
VTune performance analyzer, 17

W

w command, 31